

A Multi-Objective Particle Swarm Optimization Approach for Optimizing K-Means Clustering Centroids

Aina Latifa Riyana Putri^{1*}, Joko Riyono², Christina Eni Pujiastuti³, Supriyadi⁴ ¹Data Science, Telkom University, Purwokerto, Indonesia ^{2,3,4} Universitas Trisakti, Jakarta, Indonesia

¹ainaqp@telkomuniversity.ac.id, ²jokoriyono@trisakti.ac.id, ³christina.eni@trisakti.ac.id, ⁴supri@trisakti.ac.id

Abstract

The K-Means algorithm is a popular unsupervised learning method used for data clustering. However, its performance heavily depends on centroid initialization and the distribution shape of the data, making it less effective for datasets with complex or non-linear cluster structures. This study evaluates the performance of the standard K-Means algorithm and proposes a Multiobjective Particle Swarm Optimization K-Means (MOPSO+K-Means) approach to improve clustering accuracy. The evaluation was conducted on five benchmark datasets: Atom, Chainlink, EngyTime, Target, and TwoDiamonds. Experimental results show that K-Means is effective only on datasets with clearly separated clusters, such as EngyTime and TwoDiamonds, achieving accuracies of 95.6% and 100%, respectively. In contrast, MOPSO+K-Means achieved a substantial accuracy improvement on the complex Target dataset, increasing from 0.26% to 59.2%. The TwoDiamonds dataset achieved the most desirable trade-off: it had the lowest SSW (1323.32), relatively high SSB (2863.34), and lowest standard deviation values, indicating compact clusters, good separation, and high consistency across runs. These findings highlight the potential of swarm-based optimization to achieve consistent and accurate clustering results on datasets with varying structural complexity.

Keywords: centroid; k-means; multiobjective particle swarm optimization; the sum of square within; the sum of square between

How to Cite: A. Latifa Riyana Putri, J. Riyono, and C. Eni Pujiastuti, "A Multi-Objective Particle Swarm Optimization Approach for Optimizing K-Means Clustering Centroids", J. RESTI (Rekayasa Sist. Teknol. Inf.), vol. 9, no. 3, pp. 542 - 550, Jun. 2025. Permalink/DOI: https://doi.org/10.29207/resti.v9i3.6533

Received: April 11, 2025 Accepted: June 4, 2025 Available Online: June 21, 2025

This is an open-access article under the CC BY 4.0 License Published by Ikatan Ahli Informatika Indonesia

1. Introduction

Clustering is one of the data mining techniques aimed at grouping data into several clusters based on similarities in their characteristics. One of the most popular clustering methods is k-Means clustering [1], which is a non-hierarchical clustering algorithm that works by initializing centroids randomly. The objects are then grouped into k clusters based on their distances to the centroids, and the centroids' positions are iteratively updated until convergence is reached.

The advantages and disadvantages of this algorithm have been widely discussed in various studies. The k-Means algorithm is well-known for being efficient and scalable in processing large datasets [2]. Previous research, such as that conducted by [3] and [4], has shown that k-Means can produce more compact clusters compared to hierarchical clustering methods. However, k-Means has some limitations, particularly related to the random initialization of centroids, which can cause the clustering results to vary each time the algorithm is run [5]. Additionally, k-Means tends to get stuck in local optima, leading to suboptimal cluster assignments [6]. Its sensitivity to outliers is also a major concern, as extreme values can significantly shift the centroids' positions [7]. Furthermore, the assumption that clusters are spherical [8] and of uniform size makes k-Means less effective when dealing with datasets that have complex cluster shapes or varying densities.

To understand the quality of clustering results generated by k-Means, it is important to review the objectives and evaluation metrics. In k-Means clustering, the main goal is to form optimal clusters, where members of each cluster are highly like one another but significantly different from members of other clusters [9]. To achieve this, two primary metrics commonly used are the Sum of Squares Within-cluster (SSW) and the Sum of Squares Between-cluster (SSB). SSW measures the density of the cluster, indicating how tightly the data points within a cluster are grouped around the centroid [10]. Smaller SSW values indicate greater similarity between data points within the same cluster. On the other hand, SSB measures the distance between centroids, reflecting the separation between clusters [11]. Larger SSB values indicate greater distance between clusters, making the clustering more effective at distinguishing between different groups of data.

To address these issues, this study proposes an optimization approach based on Multi-Objective Particle Swarm Optimization (MOPSO) to determine more optimal centroids. MOPSO is a variant of Particle Swarm Optimization (PSO) [12] developed to solve problems with multiple objectives, making it suitable for clustering tasks that involve balancing two criteria simultaneously: minimizing the Sum of Squares Within-cluster (SSW) and maximizing the Sum of Squares Between-cluster (SSB), thereby producing clusters that are balanced in terms of both homogeneity and separation. The main advantage of MOPSO over other optimization algorithms lies in its convergence speed, efficient global exploration capabilities, and ease of implementation. Unlike evolutionary algorithms that use complex selection and mutation processes, MOPSO relies on a simple particle interaction mechanism in the search space. Additionally, MOPSO uses an external archive to store the best non-dominated solutions (Pareto optimal), facilitating decision-making based on trade-offs between cluster homogeneity and separation.

This approach will be evaluated using several benchmark datasets commonly used in clustering studies [13], namely Atom, Chainlink, Engytime, Target, and Two Diamonds. Unlike previous studies that often use classic datasets such as Iris, this study specifically selects these five datasets because they present more complex challenges in the data clustering process. The Atom dataset challenges the algorithm in separating very close clusters, while Chainlink has a topological structure that is interconnected, which is difficult for centroid-based methods to handle. Engytime has an uneven density distribution, which may complicate the identification of proper cluster boundaries. The Target dataset presents a non-linear pattern that is hard for standard k-Means to capture, while Two Diamonds involves clusters that are very close together, making optimal separation difficult.

The performance evaluation was conducted by comparing the clustering results against the ground truth, which refers to the true labels that are known beforehand and used as a reference to assess the accuracy of the clustering performed by the algorithm. This comparison will be made using accuracy metrics and by testing the MOPSO-based approach against conventional clustering methods such as standard k-Means. With MOPSO-based optimization, this method is expected to be able to produce more separated and uniform clusters, thus outperforming data complexity. The findings of this study are anticipated to contribute

significantly to the development of more optimal clustering methods for various applications in data mining and machine learning.

2. Methods

The methodology section will sequentially present the analytical methods employed in this study.

2.1 Experimental Testing and Simulation of the Proposed Muti-Objective Particle Swarm Optimization (PSO) Algorithm

As part of the empirical analysis, this study also evaluates the effectiveness of the proposed method using a set of benchmark datasets. These datasets are employed in the testing process to assess the capability of MOPSO in determining the optimal centroids for the K-Means algorithm, which is a critical step in enhancing clustering quality across various data scenarios. The first step in this study is the selection of datasets to be used for testing the effectiveness of the proposed method. The datasets used in this research are Atom, Chainlink, Engytime, Target, and Two Diamonds.

The Atom dataset [14] is a commonly used benchmark for evaluating clustering algorithms under complex spatial conditions. It exists in a three-dimensional space (\mathbb{R}^3) and consists of two main clusters: a dense core cluster with 100 data points located at the center, and a larger, more dispersed outer hull cluster with 400 data points that geometrically encloses the core. This structure creates what is known as an overlapping convex hull, where the outer cluster fully surrounds the inner one. The significant difference in density and spatial arrangement poses challenges for traditional clustering algorithms, especially those that rely on distance measures such as K-Means.

The Chainlink dataset [15], [16] is a benchmark designed to evaluate the ability of clustering algorithms to handle complex and interrelated data structures. It consists of two clusters, each containing 300 data points, forming an interlocked chain-like structure in three-dimensional space (\mathbb{R}^3). Each cluster is shaped like a ring, and the two rings are intertwined, creating a configuration known as linear nonseparable entanglement. This refers to a condition where the clusters cannot be linearly separated due to their intertwined spatial arrangement. Although the clusters are globally distinct, many points from one cluster are locally closer to points from the other, which introduces a conflict between global separability and local proximity. Additionally, both clusters have nearly identical densities and inter-point distances, making it difficult to distinguish them based solely on size or distribution. This makes Chainlink particularly challenging for distance-based algorithms such as K-Means.

The EngyTime dataset is a benchmark used to evaluate the effectiveness of clustering algorithms in handling overlapping clusters with varying densities [17]. It contains 2,000 data points grouped into two clusters in a two-dimensional space (\mathbb{R}^2), based on two variables: "Engy" and "Time". This dataset represents a simplified yet realistic density-based clustering problem, like those encountered in applications such as flow cytometry and sonar signal processing. EngyTime is generated from a mixture of two 2D Gaussian distributions, making it a suitable test case for evaluating how well clustering algorithms can distinguish overlapping groups. The clusters in this dataset are not separated by empty space, and they differ in density, which creates a significant challenge for traditional centroid-based methods like K-Means. These algorithms rely primarily on distance metrics and often ignore local density, which may lead to incorrect groupings when faced with overlapping, unevenly distributed data.

The Target dataset [18] is a benchmark designed to test the robustness of clustering algorithms when faced with overlapping clusters and the presence of outliers. This dataset exists in two-dimensional space (\mathbb{R}^2) and comprises 743 data points, divided into two main clusters and four small outlier groups. The first cluster is a dense spherical structure with 365 data points, while the second cluster forms an enclosing ring with 395 data points. This circular arrangement results in overlapping convex hulls, a geometric configuration that is particularly difficult to resolve using centroid-based algorithms like K-Means, which assume wellseparated, linearly distinguishable clusters. The dataset also includes four corner-located outlier groups, each containing four data points. These outliers introduce additional complexity by potentially skewing the centroid calculations or being misclassified as independent clusters. The combination of dense corering overlap and peripheral noise makes the Target dataset a comprehensive benchmark for evaluating both the accuracy and stability of clustering methods.

The TwoDiamonds dataset [19], [20] is a benchmark commonly used to evaluate the ability of clustering algorithms to distinguish between weakly connected yet distinct cluster structures. It consists of 400 data points distributed evenly across two clusters, each shaped like a diamond, and located in a two-dimensional space (\mathbb{R}^2). The two clusters are positioned in adjacent square regions that almost touch at one side, forming a configuration that resembles two diamonds placed side by side. The primary challenge of this dataset lies in the weak connection between the two clusters. While the clusters are globally distinct, the narrow gap separating them can mislead clustering algorithms-particularly those based on local distance metrics like K-Meansinto interpreting them as a single elongated cluster. Successfully separating the clusters in this dataset requires the algorithm to capture the overall geometric structure rather than relying purely on inter-point distances.

2.2 Standard K-Means Implementation as a baseline Comparison

As a baseline for performance comparison, the next step involves applying the standard K-Means clustering algorithm to each benchmark dataset. K-Means begins by randomly initializing centroids, followed by an iterative process of assigning data points to the nearest centroid based on Euclidean distance and updating the centroids until convergence. The number of clusters (k)is predetermined based on the ground truth of each dataset.

Due to the random nature of centroid initialization, K-Means may produce different clustering results in different runs. To address this, multiple independent runs are performed to assess consistency. The resulting cluster assignments are then evaluated using a confusion matrix, which enables the calculation of clustering accuracy by comparing the predicted clusters to the true class labels.

This baseline evaluation provides a reference point for comparing the clustering performance of the proposed MOPSO-KMeans method. By analyzing both standard K-Means and the optimized version under the same conditions and metrics, a more comprehensive assessment of the benefits and improvements introduced by the proposed approach can be achieved.

3. Results and Discussions

3.1 Design and Workflow of the Proposed MOPSO Method

To improve the quality of clustering results, this study implements the Multi-Objective Particle Swarm Optimization (MOPSO) algorithm to optimize the selection of centroids in the K-Means algorithm. This approach simultaneously considers two objectives: minimizing the Sum of Squared Within-Cluster (SSW) and maximizing the Sum of Squared Between-Cluster (SSB).

The first objective function aims to minimize the Sum of Squared Within-Cluster (SSW) shown in Equation 1.

$$f_{1} = \min\left(\sum_{j=1}^{k} \sum_{x_{i} \in C_{j}} \|x_{i} - \mu_{j}\|^{2}\right)$$
(1)

k is the number of clusters; x_i is the i-th data point; μ_j is the centroid of cluster C_j ; $||x_i - \mu_j||^2$ is the squared Euclidean distance between the data point and the cluster centroid.

The second objective function shown in Equation 2 aims to maximize the Sum of Squared Between-Cluster (SSB), which is expressed as the minimization of its negative:

$$f_{2} = -\min\left(-\sum_{j=1}^{k} n_{j} \|\mu_{j} - \mu\|^{2}\right)$$
(2)

 n_j is the number of data points in cluster j; μ is the global centroid of the entire dataset; $\|\mu_j - \mu\|^2$ is the squared

distance between the cluster centroid and the global centroid.

The complete procedure of the proposed Multi-Objective Particle Swarm Optimization (PSO) algorithm can be summarized as follows, with its pseudocode presented on Figure 1.

Algorithm 1: Pseudocode of the proposed MOPSO+K-Means Algorithm

Input:

Dataset D with n data points; Number of particles N; Maximum number of iterations T; Objective functions: fl = intra-cluster distance (minimize) and f2 = inter-cluster separation (maximize)

Process:

1. Initialization Phase:

2. For i = 1 to N Do:

3. Randomly initialize centroids i with K positions in the data space.

4. Initialize velocity i = 0.

5. Assign each data point in D to the nearest centroid in centroids i.

6. Evaluate objective functions fl and f2.

7. Set best_objective i = [fl, f2].

9. End For

10. Evaluate dominance among all particles.

11. Save non-dominated solutions into repository.

12. Search Phase:

13. For i = 1 to N Do:

14. For i = 1 to N Do:

15. Select global best from repository using crowding distance.

16. Update velocity i using PSO velocity update rule.

17. Update centroids i using velocity_i.

18. Assign data points to the nearest centroid in centroids_i.

19. Evaluate objective functions fl and f2.

10. Apply mutation (optional).

21. If centroids i dominates best position_i Then:

22. Apply mutation (optional).

23. End

Figure 1. Pseudocode of the Proposed MOPSO+K-Means Algorithm

The proposed MOPSO+K-Means algorithm integrates Multi-Objective Particle Swarm Optimization (MOPSO) with the traditional K-Means clustering method to overcome the limitations of random centroid initialization. As outlined in Figure 1, the process begins with the initialization phase, where each particle represents a potential solution in the form of a set of cluster centroids. Objective functions are defined to minimize intra-cluster distance (SSW) and maximize inter-cluster separation (SSB), enabling a balanced evaluation of cluster compactness and separation.

The fitness of each particle is assessed based on these two objectives, and the non-dominated solutions are stored in a Pareto-based repository. In the search phase, the particle positions are updated using both personal and global bests selected from the repository using the crowding distance. This iterative process continues until the stopping criteria are met, gradually refining the solutions toward the optimal trade-offs between the two objectives.

At the end of the MOPSO optimization, the repository contains a set of Pareto-optimal centroid configurations. From these, one or more candidate solutions can be selected to initialize the K-Means algorithm. This hybrid approach allows K-Means to begin with welloptimized centroid positions, potentially resulting in more stable and accurate clustering outcomes compared to traditional random initialization. In addition to the Particle Swarm Optimization (PSO) algorithm, previous studies have also explored other metaheuristic approaches to improve clustering quality. Among the most widely used are Genetic Algorithm [21] (GA) and Differential Evolution [22] (DE). However, findings from several prior works suggest that PSO tends to offer advantages in terms of convergence speed and simplicity of parameters. PSO only requires a few parameters to be configured, such as inertia weight and learning coefficients. In contrast, GA and DE involve more complex parameter settings, including crossover rate, mutation rate, and selection strategies [23] which can significantly influence performance if not properly adjusted. Furthermore, PSO is known for its ability to maintain a good balance between exploration and exploitation during the search process [24], making it especially suitable for clustering problems with high complexity.

3.2 Experimental Results and Analysis

To evaluate the performance of the proposed Multi-Objective PSO, a series of experiments were conducted on a set of well-known benchmark datasets. These five datasets have been described in the Research Method.

The experimental setup in this study is defined as follows: the swarm size is set to N=40, and each benchmark dataset is tested independently 30 times, with each execution consisting of 100 iterations. All PSO-based algorithms are terminated upon reaching this maximum number of iterations. The performance of the proposed MOPSO+K-Means method is evaluated using standard clustering metrics, namely the best value, average value, and standard deviation of SSW, SSB, and accuracy. These metrics are used to assess the effectiveness and stability of the proposed method in comparison to standard K-Means across various benchmark datasets.

The performance of MOPSO+K-Means was evaluated using commonly used optimization metrics, namely the average solution and standard deviation. These metrics were used to assess the effectiveness of MOPSO+K-Means in solving the benchmark clustering tasks.

Table 1. Clustering With MOPSO+K-Means

Dataset	Item	SSW	SSB
Atom	Avg.	1191.22	1414.52
	Std.	230.85	238.22
ChainLink	Avg.	1531.74	1711.26
	Std.	167.04	168.519
EngyTime	Avg.	49122.71	85124.33
	Std.	2951.153	3913.554
Target	Avg.	4126.614	7658.074
•	Std.	892.147	1006.303
TwoDiamonds	Avg.	1323.322	2863.340
	Std.	42.822	44.191

Although the numerical results demonstrate that the proposed MOPSO+K-Means algorithm performs competitively across various datasets, a deeper analysis provides further insights into its behavior and performance dynamics. From a clustering quality perspective on Table 1, the ideal objective is to minimize the Sum of Squared Within-cluster distances (SSW) while maximizing the Sum of Squared Betweencluster distances (SSB). In this context, the EngyTime dataset exhibits the highest SSW (49122.71) and SSB (85124.33), indicating that the dataset likely has a large or widely spread structure. Despite this complexity, the algorithm successfully maintains strong inter-cluster separation. In contrast, the TwoDiamonds dataset achieves the most desirable trade-off: it has the lowest SSW (1323.32), relatively high SSB (2863.34), and the lowest standard deviation values, indicating compact clusters, good separation, and high consistency across runs.

For other datasets, such as Target and Atom, the SSW and SSB values fall in a mid-range category, but the relatively high standard deviations, particularly in Atom, highlight variability in the clustering results, suggesting sensitivity to initial conditions or swarm dynamics. Similarly, ChainLink yields results comparable to Atom but also shows considerable variability (Std SSW: 167.04, Std SSB: 168.52), likely reflecting the complexity or overlap in its cluster structure. These findings suggest that while MOPSO+K-Means can handle both simple and complex datasets, its stability may be challenged under certain data conditions.

Overall, TwoDiamonds emerges as the most stable dataset for this method, whereas EngyTime, despite strong average performance, reveals higher variability possibly due to noise or dispersed data points. These behavioral differences point to both strengths and limitations of the method. The ability of MOPSO+K-Means to find competitive clustering solutions is evident, but further enhancements could improve its robustness. Future research should explore adaptive parameter strategies, improved initialization techniques, or hybrid models to increase the method's reliability across diverse datasets. Additionally, expanding testing to high-dimensional or real-world datasets would help evaluate its scalability and broader applicability.

In this section, an analysis is conducted on the clustering outcomes derived from the implementation of the standard K-Means algorithm and the MOPSO-enhanced K-Means across a range of benchmark datasets. These results highlight the capabilities and limitations of both approaches when confronted with datasets exhibiting varying structural complexities.



Figure 2. Clustering Result (a) Atom Dataset, (b) K-Means for Atom Dataset, (c) MOPSO+K-Means for Atom Dataset

The clustering results on the Atom dataset highlight the limitations of The Atom dataset illustrates the challenge of clustering data with a concentric structure, comprising a dense core surrounded by a shell. As depicted in Figure 2(a), the ground truth clearly reflects this core-shell configuration. However, the standard K-Means algorithm on Figure 2(b) produces a vertical partition, disregarding the radial nature of the data. This misalignment stems from K-Means assumption of spherical and convex clusters, which proves inadequate for capturing non-linear distributions. The integration of MOPSO+K-Means, as shown in Figure 2(c), results in a more accurate division that successfully distinguishes between the core and the surrounding shell. This demonstrates the ability of MOPSO to guide K-Means toward more structure-aware clustering outcomes in complex spatial configurations.

A similar pattern is observed in the ChainLink dataset, which features two intertwined, non-convex clusters on Figure 3(a). The standard K-Means algorithm on Figure 3(b) again fails to separate the data meaningfully, as it assigns points based on straight-line distances, ignoring the dataset's intricate shape. Conversely, the MOPSO+K-means on Figure 3(c) more effectively untangles the two chains, maintaining their topological distinction. This improved result underscores the role of MOPSO in adapting centroid placement to fit non-linear geometries that would otherwise confound traditional methods.

The EngyTime dataset, in contrast, presents a linearly separable structure, offering a more favorable scenario for K-Means. In Figure 4(a), the data exhibits two welldefined, adjacent clusters. The standard K-Means output on Figure 4(b) broadly captures the separation but misclassifies several points near the decision boundary. With optimized centroids on Figure 4(c), the clustering becomes cleaner, with reduced boundary ambiguity. This shows that even in simpler datasets,

MOPSO+K-Means contributes by refining the clustering precision and minimizing convergence to suboptimal solutions.



Figure 3. Clustering Result (a) Chainlink Dataset, (b) K-Means for Chainlink Dataset, (c) MOPSO+K-Means for Chainlink Dataset



Figure 4. Clustering Result (a) EngyTime Dataset, (b) K-Means for EngyTime Dataset, (c) MOPSO+K-Means for EngyTime Dataset



Figure 5. Clustering Results (a) Target Dataset, (b) K-Means for Target Dataset, (c) MOPSO+K-Means for Target Dataset

The analysis of the Target dataset provides further insight into the impact of complex geometries on clustering performance. This dataset contains concentric circular clusters and scattered peripheral groups on Figure 5(a). The clustering produced by the standard K-Means algorithm on Figure 5(b) results in fragmented groupings that poorly reflect the actual layout. Its tendency to impose spherical boundaries leads to significant structural mismatches. The MOPSO+K-Means on Figure 5(c) successfully aligns with the circular patterns and isolates the outlying groups more accurately, reinforcing the value of optimization in adapting to irregular spatial distributions.

Lastly, the TwoDiamonds dataset poses a moderate challenge due to its diamond-shaped, linearly separated clusters on Figure 6(a). While the standard K-Means on Figure 6(b) performs reasonably well, some inconsistencies are evident along the boundary, suggesting less-than-optimal centroid positioning. By contrast, the optimized version on Figure 6(c) yields a more symmetric and faithful clustering result, demonstrating how MOPSO can enhance K-Means even in datasets that are linearly separable but geometrically unconventional.

Table 2 presents the performance evaluation results of the MOPSO-K-Means algorithm on five benchmark datasets: Atom, ChainLink, EngyTime, Target, and TwoDiamonds. The evaluation was carried out using commonly used optimization metrics, namely the average solution and standard deviation of the SSW (Sum of Squares Within) and SSB (Sum of Squares Between), along with the best accuracy achieved for each dataset. The objective of this evaluation is to assess the effectiveness of the MOPSO-K-Means algorithm in producing optimal cluster partitions.



Figure 6. Clustering Result (a) TwoDiamonds Dataset, (b) K-Means Clustering Result TwoDiamonds Dataset, (c) MOPSO+K-Means Clustering Result TwoDiamonds Dataset

Table 2. Comparing Accuracy and Time Computational

Dataset	Accuracy K-means	Best Accuracy MOPSO+K- Means	Average Time Computational MOPSO+K- Means
Atom	54.4%	52.8%	5.1137 seconds
ChainLink	50%	50.2%	6.6889 seconds
EngyTime	95.6%	95.7%	7.617 seconds
Target	0.2692%	59.2%	8.539 seconds
TwoDiamonds	100%	100%	0.844 seconds

Compared to the conventional K-Means algorithm, the results indicate that MOPSO-K-Means generally performs better on most datasets. On the Atom dataset, K-Means achieved an accuracy of 54.4%, while MOPSO-K-Means recorded an accuracy of 52.8%. Although there was a slight decrease, the SSW and SSB values obtained by MOPSO-K-Means still reflect a good and stable cluster distribution, with relatively low standard deviations. For the ChainLink dataset, K-Means achieved 50% accuracy, while MOPSO-K-Means achieved 50.2%, suggesting a slightly better performance in separating the clusters.

Next, on the EngyTime dataset, K-Means reached an accuracy of 95.60%, while MOPSO-K-Means achieved 95.7%. The difference is very small, indicating that both algorithms are equally effective in clustering data with clear cluster structures. However, the most significant improvement was observed on the Target dataset. K Means achieved only 0.26% accuracy, while MOPSO-K-Means improved the accuracy to 59.2%. This demonstrates that MOPSO-K-Means is more capable of handling datasets with complex or non-linearly separable cluster structures. Lastly, on the TwoDiamonds dataset, both K-Means and MOPSO-K-Means achieved perfect accuracy (100%), indicating that this dataset has a very clear structure that can be easily separated by both algorithms.

The Average Time Computational *MOPSO+K-Means* column in Table 2 presents the average computational

time required by the algorithm to complete one clustering process for each dataset. This time is measured from the beginning of the centroid optimization using the MOPSO algorithm to the final clustering result produced by K-Means. The values represent the average time taken across 30 independent trials. They indicate the efficiency of the algorithm in solving clustering tasks, which is influenced by the complexity of the data patterns and the algorithm's ability to converge toward optimal solutions.

For example, the TwoDiamonds dataset demonstrates the lowest average computational time (0.844 seconds), which can be attributed to its well-separated and clearly structured clusters. This allows MOPSO to converge quickly, requiring minimal exploration. In contrast, the Target dataset shows the highest computational time (8.539 seconds), suggesting that the algorithm needed significantly more effort to explore the solution space due to the highly irregular and overlapping nature of the data clusters. Similarly, EngyTime and Chainlink also required more computational time, which reflects the greater complexity in their data distributions and the increased difficulty in identifying distinct clusters. Meanwhile, Atom shows a moderate computational time, indicating that while the data is not entirely straightforward, it is still manageable for the algorithm to optimize efficiently.

Overall, the evaluation results show that MOPSO+K-Means has advantages in terms of flexibility and effectiveness in identifying complex cluster structures that conventional K-Means struggles to handle. The relatively small standard deviations across most datasets also indicate that this algorithm can produce stable and consistent solutions in each optimization run. Therefore, MOPSO+K-Means can be considered a more reliable alternative for clustering tasks involving datasets with diverse characteristics.

However, one key limitation of the current study is the assumption that the number of clusters (k) is known beforehand. Although this facilitates benchmarking against ground truth labels, it does not reflect the realities of unsupervised learning tasks, where k must be inferred from the data. In practice, k is a hyperparameter that must be estimated carefully. Common strategies include the elbow method, which identifies diminishing returns in intra-cluster variance as k increases, the silhouette score, which quantifies cluster cohesion and separation, and the gap statistic, which compares clustering performance against that of random reference distributions.

To overcome this limitation, future work could explore the extension of MOPSO to simultaneously optimize both the number of clusters and the centroid positions. This could be framed as a multi-objective optimization problem—balancing cluster compactness, separation, and model complexity—or as a constrained optimization task where k is bounded within a reasonable range. Such an approach would enhance the applicability of the method in real-world scenarios, allowing it to autonomously discover both the optimal clustering structure and its corresponding parameters without relying on prior knowledge.

4. Conclusions

Based on the analysis and evaluation of five benchmark datasets, it can be concluded that the performance of the K-Means algorithm is highly dependent on the shape and structural characteristics of the clusters in the data. On datasets with simple and linearly separable structures, such as TwoDiamonds and EngyTime, K-Means performs very well, achieving high accuracy up to 100%. However, on datasets with non-linear or complex structures, such as Atom, ChainLink, and Target, the algorithm fails to properly separate clusters, resulting in low accuracy and poor alignment with the ground truth.

To address these limitations, the MOPSO+K-Means approach was introduced as an alternative solution. Based on the experimental results, this algorithm shows significant performance improvement on datasets with complex structures—most notably on the Target dataset, where the accuracy increased from 26% (K-Means) to 59.2% (MOPSO+K-Means). In addition, the obtained SSW and SSB values, along with relatively low standard deviations, indicate that MOPSO-K-Means can produce stable and consistent clustering solutions.

Overall, MOPSO+K-Means has proven to be more flexible and reliable in handling various types of cluster structures, making it a more suitable choice for clustering tasks involving non-convex or non-linearly separable data distributions. Such characteristics are often found in real-world applications, including biomedical data analysis where patient subgroups may form irregular patterns, sensor grouping in IoT

environments with overlapping signal zones, and document clustering tasks where semantic relationships are not linearly separable. These application domains can benefit from algorithms that offer both structural flexibility and solution stability, as demonstrated by MOPSO+K-Means in this study.

Acknowledgements

This work was supported by Telkom University.

References

- T. M. Ghazal *et al.*, "Performances of K-Means Clustering Algorithm with Different Distance Metrics," *Intelligent Automation & Soft Computing*, vol. 30, no. 2, pp. 735–742, Aug. 2021, doi: 10.32604/IASC.2021.019067.
- [2] M. Ahmed, R. Seraj, and S. M. S. Islam, "The k-means Algorithm: A Comprehensive Survey and Performance Evaluation," *Electronics 2020, Vol. 9, Page 1295*, vol. 9, no. 8, p. 1295, Aug. 2020, doi: 10.3390/ELECTRONICS9081295.
- [3] A. Abdulhafedh, "Incorporating K-means, Hierarchical Clustering and PCA in Customer Segmentation," *Journal* of City and Development, vol. 3, no. 1, 2021.
- [4] P. Patel, B. Sivaiah, and R. Patel, "Approaches for finding Optimal Number of Clusters using K-Means and Agglomerative Hierarchical Clustering Techniques," in 2022 International Conference on Intelligent Controller and Computing for Smart Power, ICICCSP 2022, 2022. doi: 10.1109/ICICCSP53532.2022.9862439.
- [5] A. Vouros, S. Langdell, M. Croucher, and E. Vasilaki, "An empirical comparison between stochastic and deterministic centroid initialisation for K-means variations," *Mach Learn*, vol. 110, no. 8, pp. 1975–2003, Aug. 2021, doi: 10.1007/S10994-021-06021-7/FIGURES/8.
- [6] A. Qtaish, M. Braik, D. Albashish, M. T. Alshammari, A. Alreshidi, and E. J. Alreshidi, "Optimization of K-means clustering method using hybrid capuchin search algorithm," *Journal of Supercomputing*, vol. 80, no. 2, 2024, doi: 10.1007/s11227-023-05540-5.
- [7] Z. Zhang, Q. Feng, J. Huang, Y. Guo, J. Xu, and J. Wang, "A local search algorithm for k-means with outliers," *Neurocomputing*, vol. 450, pp. 230–241, Aug. 2021, doi: 10.1016/J.NEUCOM.2021.04.028.
- [8] A. A. Wani, "Comprehensive analysis of clustering algorithms: exploring limitations and innovative solutions," *PeerJ Comput Sci*, vol. 10, pp. 1–45, Aug. 2024, doi: 10.7717/PEERJ-CS.2286/FIG-14.
- [9] R. Gustriansyah, N. Suhandi, and F. Antony, "Clustering optimization in RFM analysis Based on k-Means," *IJEECS*, vol. 18, no. 1, pp. 470–477, Apr. 2020, doi: 10.11591/ijeecs.v18.i1.pp470-477.
- [10] R. Richard, H. Cao, and M. Wachowicz, "An Automated Clustering Process for Helping Practitioners to Identify Similar EV Charging Patterns across Multiple Temporal Granularities," *International Conference on Smart Cities* and Green ICT Systems, pp. 67–77, 2021, doi: 10.5220/0010485000670077.
- [11] M. T. Guerreiro *et al.*, "Anomaly Detection in Automotive Industry Using Clustering Methods—A Case Study," *Applied Sciences 2021, Vol. 11, Page 9868*, vol. 11, no. 21, p. 9868, Oct. 2021, doi: 10.3390/APP11219868.
- [12] M. Jain, V. Saihjpal, N. Singh, and S. B. Singh, "An Overview of Variants and Advancements of PSO Algorithm," *MDPI Applied Sciences*, 2022, doi: 10.3390/app12178392.
- [13] M. C. Thrun and A. Ultsch, "Clustering benchmark datasets exploiting the fundamental clustering problems," *Data Brief*, vol. 30, p. 105501, Jun. 2020, doi: 10.1016/J.DIB.2020.105501.
- [14] A. Ultsch, "Strategies for an Artificial Life System to cluster high dimensional Data," 2004, Accessed: Apr. 14,

2025. [Online]. Available: https://www.researchgate.net/publication/228932819

- [15] A. Ultsch, G. Guimaraes, D. Korus, and H. Li, "Knowledge Extraction from Artificial Neural Networks and Applications," *Parallele Datenverarbeitung mit dem Transputer*, pp. 148–162, 1994, doi: 10.1007/978-3-642-78901-4_11.
- [16] P. Mangiameli, S. K. Chen, and D. West, "A comparison of SOM neural network and hierarchical clustering methods," *Eur J Oper Res*, vol. 93, no. 2, pp. 402–417, Sep. 1996, doi: 10.1016/0377-2217(96)00038-0.
- [17] S. P. Chatzis and D. I. Kosmopoulos, "A variational Bayesian methodology for hidden Markov models utilizing Student's-t mixtures," *Pattern Recognit*, vol. 44, no. 2, pp. 295–306, Feb. 2011, doi: 10.1016/J.PATCOG.2010.09.001.
- [18] J. Poelmans, M. M. Van Hulle, S. Viaene, P. Elzinga, and G. Dedene, "Text mining with emergent self organizing maps and multi-dimensional scaling: A comparative study on domestic violence," *Appl Soft Comput*, vol. 11, no. 4, pp. 3870–3876, Jun. 2011, doi: 10.1016/J.ASOC.2011.02.026.
- [19] A. Ultsch, "U*-Matrix : a Tool to visualize Clusters in high dimensional Data," 2004.
- [20] A. Ultsch, "Density Estimation and Visualization for Data Containing Clusters of Unknown Structure," *Studies in*

Classification, Data Analysis, and Knowledge Organization, pp. 232–239, 2005, doi: 10.1007/3-540-28084-7_25.

- [21] B. Khusul Khotimah, F. Irhamni, and T. Sundarwati, "A GENETIC ALGORITHM FOR OPTIMIZED INITIAL CENTERS K-MEANS CLUSTERING IN SMEs," J Theor Appl Inf Technol, vol. 15, no. 1, 2016, Accessed: Jun. 13, 2025. [Online]. Available: www.jatit.org
- [22] H. He, B. Sun, Y. Yang, and S. Liu, "An improved K-Means clustering based on differential evolution," *J Phys Conf Ser*, vol. 2595, no. 1, p. 012010, Sep. 2023, doi: 10.1088/1742-6596/2595/1/012010.
- [23] A. Hassanat, K. Almohammadi, E. Alkafaween, E. Abunawas, A. Hammouri, and V. B. S. Prasath, "Choosing Mutation and Crossover Ratios for Genetic Algorithms—A Review with a New Dynamic Approach," *Information 2019, Vol. 10, Page 390*, vol. 10, no. 12, p. 390, Dec. 2019, doi: 10.3390/INFO10120390.
- [24] M. Jain, V. Saihjpal, N. Singh, and S. B. Singh, "An Overview of Variants and Advancements of PSO Algorithm," *Applied Sciences 2022, Vol. 12, Page 8392*, vol. 12, no. 17, p. 8392, Aug. 2022, doi: 10.3390/APP12178392.