

# K-means Algorithm

Tang Chenyang  
Department of Computer Science  
Nanjing University of Posts and  
Telecommunications  
Jiangsu, Nanjing China  
2210925160@qq.com

**Abstract**—This paper focuses on the K-means algorithm in the field of unsupervised learning, exploring its principles, implementation processes, and application scenarios. As an important clustering algorithm, the K-means algorithm has wide applications in data mining, pattern recognition, and other fields. This article elaborates on the theoretical basis of the K-means algorithm, including distance metrics, clustering criteria, etc., and implements the algorithm through actual data sets to demonstrate its clustering effect. At the same time, the advantages and disadvantages of the K-means algorithm are analyzed, and corresponding improvement strategies are proposed to improve the algorithm performance. Finally, the applications of the K-means algorithm in different fields are studied and prospected, providing references for research and practice in related fields.

**Keywords**—Unsupervised learning; K-means algorithm; Clustering analysis; Data mining

## I. INTRODUCTION

In the digital age, the volume of data has exploded exponentially. Data is being generated at an unprecedented rate from various sources such as social media, sensors, and business transactions. The ability to extract meaningful insights from this vast amount of unlabeled data has become a crucial task in many fields. Unsupervised learning, a fundamental area in machine learning, provides techniques to discover hidden patterns and structures within the data without prior knowledge of the output labels.

K-means algorithm, a popular clustering method in unsupervised learning, has emerged as a powerful tool for data analysis. It offers a simple yet effective way to group data points based on their similarity. The concept of clustering is essential as it helps in understanding the natural groupings or segments within a dataset. For example, in market research, clustering customers based on their purchasing behavior can help businesses better understand their customer base and develop targeted marketing strategies.

However, the traditional K-means algorithm is not without its limitations. One of the main challenges is its sensitivity to the initial selection of cluster centers. Poor initial choices can lead to suboptimal clustering results. Additionally, the algorithm may converge to local minima, rather than the global minimum, affecting the quality of the final clustering. Moreover, in datasets with complex geometries or varying densities, the standard K-means may not accurately capture the underlying structure.

The primary objective of this research is to conduct an in-depth study of the K-means algorithm in the context of unsupervised learning. This includes a comprehensive analysis of its theoretical foundation, such as the distance metrics used to measure similarity between data points and the clustering criteria that guide the formation of clusters. By

understanding these fundamental aspects, we can better appreciate the behavior and limitations of the algorithm.

Another important goal is to implement the K-means algorithm using real-world datasets and evaluate its performance. Through practical experiments, we can observe how the algorithm performs in different scenarios and measure its effectiveness in clustering data accurately. This will involve analyzing the clustering results in terms of the quality of the clusters formed, the convergence speed of the algorithm, and its ability to handle different types of data distributions.

Furthermore, we aim to address the limitations of the traditional K-means algorithm. This will involve exploring and implementing various improvement strategies, such as advanced initialization methods to enhance the selection of initial cluster centers, optimization techniques to avoid local minima, and adaptations of the algorithm to handle datasets with complex characteristics. By proposing and evaluating these improvements, we hope to enhance the performance and applicability of the K-means algorithm in a wider range of applications.

The research on the K-means algorithm holds significant theoretical and practical implications. From a theoretical perspective, it contributes to the deeper understanding of unsupervised learning algorithms and clustering techniques. By exploring the intricacies of the K-means algorithm, we can gain insights into the general principles of data clustering and similarity measurement, which can be extended to other related algorithms and machine learning problems.

In practical terms, the findings of this research can have a direct impact on various fields. In business and marketing, improved clustering algorithms can lead to more effective customer segmentation, enabling companies to personalize their products and services and target specific customer groups more efficiently. In image processing and computer vision, enhanced clustering techniques can be used for image segmentation, object recognition, and pattern analysis, improving the performance of these applications. In addition, in fields such as bioinformatics and healthcare, where large amounts of unlabeled data are generated, better clustering algorithms can help in data exploration, disease classification, and the discovery of new biological patterns. Overall, this research has the potential to provide valuable tools and techniques for data analysis and decision-making in a wide range of industries and scientific disciplines.

## II. RELATED WORKS

In the field of unsupervised learning and clustering analysis, numerous research efforts have been dedicated to addressing the challenges and improving the performance of clustering algorithms, especially the K-means algorithm.

One significant area of research focuses on initialization methods for K-means. Some studies propose using advanced sampling techniques or data-driven approaches to select more appropriate initial cluster centers. For example, the K-Means++[1] algorithm (Arthur & Vassilvitskii, 2007) introduced a smart initialization method that probabilistically selects initial centroids based on the data distribution. This approach has been shown to improve the quality of the final clustering results compared to random initialization. Other variants of initialization techniques have also been explored, such as using hierarchical clustering or density-based methods to obtain initial seeds that are more likely to lead to better clustering.

Another line of research aims to handle datasets with complex geometries and varying densities. Density-based clustering algorithms, like DBSCAN[2] (Ester et al., 1996), have been proposed as alternatives to K-means in such scenarios. DBSCAN can identify clusters of arbitrary shapes and is less sensitive to the density variations within the data. However, it also has its own limitations, such as the difficulty in determining the appropriate density parameters. Some hybrid methods have been developed that combine the advantages of K-means and density-based algorithms. These methods attempt to leverage the simplicity and efficiency of K-means for initial clustering and then use density-based techniques to refine the clusters and handle complex structures.

Research has also been conducted on improving the clustering performance in high-dimensional data. Traditional K-means may suffer from the curse of dimensionality, where the distance measures become less meaningful in high-dimensional spaces. Feature selection and dimensionality reduction techniques have been integrated with K-means to address this issue. Principal Component Analysis (PCA) is often used to reduce the dimensionality of the data before applying K-means. This not only helps in improving the computational efficiency but also can enhance the clustering quality by removing redundant or irrelevant features.

In addition, parallel and distributed computing techniques have been applied to speed up the K-means algorithm for large-scale datasets. With the increasing size of datasets, the computational cost of K-means becomes a bottleneck. Parallel K-means algorithms[3] (Zaharia et al., 2010) distribute the data and the computational tasks across multiple processors or nodes in a cluster, allowing for faster convergence and handling of massive datasets. These techniques are crucial for applications in big data analytics, where timely processing of large amounts of data is essential.

Overall, the research related to the K-means algorithm is diverse and aims to overcome its limitations and expand its applicability in different domains. By leveraging these related works, our research can build on existing knowledge and explore new directions for improving the performance and effectiveness of the K-means algorithm in unsupervised learning.

### III. PROBLEM STATEMENT

#### Input

We are given a dataset  $D = \{x_1, x_2, \dots, x_n\}$  consisting of 'n' data points, where each data point  $x_i$  is represented as a d-dimensional vector in a feature space  $F$ , i.e.,  $x_i \in \mathbb{R}^d$ . The dataset  $D$  is unlabeled, meaning that there is no prior

information about the class or group membership of each data point.

#### A. Task

The objective is to partition the dataset  $D$  into  $K$  non-overlapping subsets or clusters  $C = \{C_1, C_2, \dots, C_K\}$  such that the following clustering criteria are optimized:

**Intra-cluster similarity:** Data points within each cluster should be as similar as possible to each other. This is typically measured using a distance metric  $d(x_i, x_j)$  that quantifies the dissimilarity between two data points  $x_i$  and  $x_j$ . The goal is to minimize the sum of distances between data points within each cluster. For example, in the case of the K-means algorithm, the intra-cluster similarity is often measured by the sum of squared distances of data points to their cluster centroids.

**Inter-cluster dissimilarity:** Clusters should be as distinct or dissimilar as possible from each other. This implies maximizing the separation between different clusters. One way to measure this is by considering the distances between cluster centroids. Larger distances between centroids indicate greater inter-cluster dissimilarity.

Mathematically, the clustering problem can be formulated as an optimization problem. Let  $c_k$  denote the centroid of cluster  $C_k$ . The objective is to minimize the following objective function:

$$J(C) = \sum_{i \in C_k} d(x_i, c_k)$$

subject to the constraints that each data point belongs to exactly one cluster and the clusters cover the entire dataset.

#### Challenges

**Initialization sensitivity:** The choice of initial cluster centroids can significantly affect the final clustering result. Poor initializations may lead to suboptimal clusterings or convergence to local minima of the objective function.

**Complex data distributions:** Real-world datasets often exhibit complex geometries, varying densities, and may contain noise or outliers. Traditional clustering algorithms like K-means, which assume spherical clusters and equal densities, may not be able to accurately capture the underlying structure of such datasets.

**High-dimensional data:** As the dimensionality of the data increases, the so-called "curse of dimensionality" can cause problems. Distance measures become less discriminative, and the computational complexity of clustering algorithms also increases, making it more difficult to find meaningful clusters.

**Scalability:** With the growth of data size, the computational cost and memory requirements of clustering algorithms become a major concern. Efficient algorithms are needed to handle large-scale datasets in a reasonable amount of time and with limited resources.

### IV. ALGORITHMS/SOLUTIONS

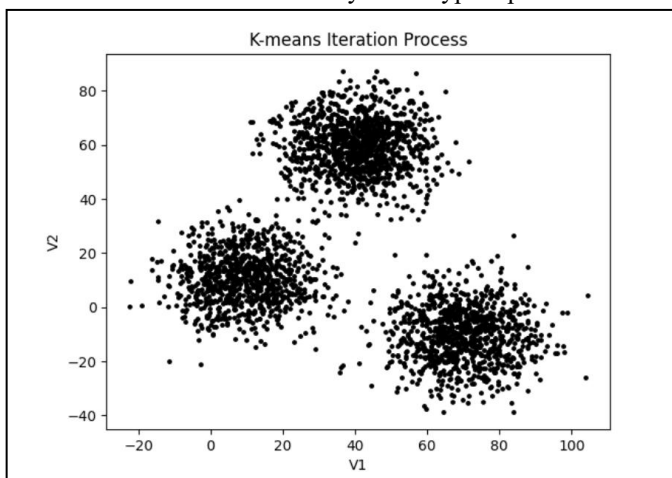
In this section, we present our implementation of the K-means algorithm and the techniques used to visualize its iterative process. The overall goal is to not only obtain the final clustering results but also gain insights into how the

algorithm converges and how the clusters evolve over the iterations.

### Data Preparation and Initialization

The data preparation phase is crucial for the K-means algorithm. We start by importing the necessary libraries, namely 'numpy' for numerical operations, 'pandas' for data handling, and 'matplotlib.pyplot' for data visualization. The dataset is read from a CSV file ('xclara.csv') using 'pandas' 'read\_csv' function. The relevant columns 'V1' and 'V2' are extracted and converted into a 2D 'numpy' array 'X', where each row represents a data point in the 2D space.

After preparing the data, we initialize the clustering process. The number of clusters 'k' is set to 3 in this case. The initial cluster centers are randomly selected within a certain range. Specifically, the X and Y coordinates of the cluster centers ('C\_x' and 'C\_y') are randomly generated using 'np.random.randint' within the range from 0 to 'np.max(X) - 20'. These coordinates are then combined to form the initial cluster centers 'C' as a 2D array with dtype 'np.float32'.



### Distance Calculation and Data Point Assignment

To measure the similarity between data points and cluster centers, we define a distance function 'dist'. This function calculates the Euclidean distance between two points 'a' and 'b' using 'np.linalg.norm'. The distance calculation is a fundamental operation in the K-means algorithm as it determines which cluster a data point belongs to.

In the data point assignment step, we iterate through each data point in 'X'. For each data point, we calculate its distances to all the cluster centers using the 'dist' function. The data point is then assigned to the cluster with the nearest center, which is determined by finding the index of the minimum distance using 'np.argmin'. The cluster assignments are stored in the 'clusters' array.

### Cluster Center Update and Iterative Convergence

Once the data points are assigned to clusters, the cluster centers need to be updated. We iterate through each cluster and calculate the mean of the data points belonging to that cluster. This new mean becomes the updated position of the cluster center. This process is repeated until the cluster centers no longer change significantly, which indicates that the algorithm has converged. We use a 'while' loop to continuously update the cluster centers and check for convergence. The loop breaks when the cluster assignments

```
print("最终聚类中心: ", C)
print("聚类划分: ", clusters)
```

```
最终聚类中心: [[ 9.478045 10.686052]
 [ 69.92419 -10.119641]
 [ 40.683628 59.715893]]
聚类划分: [0. 0. 0. ... 1. 1. 1.]
```

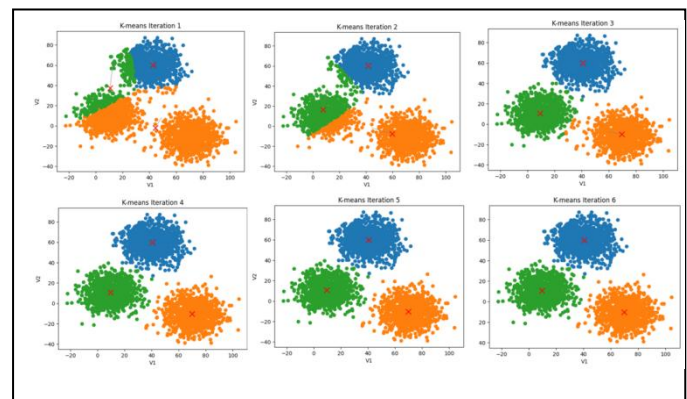
in two consecutive iterations are exactly the same, as determined by the condition '(new\_clusters==clusters).all()'.

### Visualization of the Iterative Process

In addition to obtaining the final clustering results, we also focus on visualizing the iterative process of the K-means algorithm. To achieve this, we define a function 'kmeans' that encapsulates the entire K-means algorithm with visualization capabilities.

Inside the 'kmeans' function, the same initialization, distance calculation, data point assignment, and cluster center update steps are performed as in the previous implementation. However, at each iteration, we clear the previous plot using 'plt.clf()' and then plot the data points belonging to each cluster with different colors. The cluster centers are marked with red 'x' symbols. Moreover, we plot the movement paths of the cluster centers between consecutive iterations using gray lines with transparency. This allows us to observe how the cluster centers shift and how the clusters are refined over time. The 'plt.pause(0.5)' function is used to pause the plot for 0.5 seconds at each iteration, enabling us to see the changes in the clustering in a step-by-step manner.

By visualizing the iterative process, we can gain a better understanding of the behavior of the K-means algorithm. It helps us identify potential issues such as slow convergence or the algorithm getting stuck in local minima. Additionally, it provides an intuitive way to evaluate the stability and quality of the clustering results.



### Comparison of Results with Different k Values

In order to comprehensively evaluate the performance of the K - means algorithm, we conduct experiments by varying the number of clusters, k.

We run the K - means algorithm multiple times with different values of k (e.g., k = 2, 3, 4, 5,...). For each value of k, we repeat the clustering process multiple times (e.g., 10 times) to account for the random initialization of cluster centers.

## Analysis of Results

When  $k = 2$ , we observe that the data may be over-simplified, resulting in large and heterogeneous clusters. The WCSS value is relatively high, and the silhouette score may be low, indicating that the clusters are not well-separated.

```
: print("最终聚类中心: ", c)
print("聚类划分: ", clusters)

最终聚类中心: [[ 69.49145 -10.108037]
 [ 26.979048 38.425102]]
聚类划分: [1. 1. 1. ... 0. 0. 0.]
```

As  $k$  increases to 3 (the original  $k$  value in our initial implementation), we find that the clustering results show a better balance. The clusters are more distinct, and the WCSS value decreases while the silhouette score increases, indicating an improvement in clustering quality.

When  $k$  is further increased, we may start to see over-fitting. Smaller and more fragmented clusters may appear, and the WCSS value may decrease only slightly while the silhouette score may start to decline. This indicates that the algorithm is starting to split the data into too many small and possibly meaningless clusters.

```
: print("最终聚类中心: ", c)
print("聚类划分: ", clusters)

最终聚类中心: [[ 69.92419 -10.119641]
 [ 48.02195 59.328873]
 [ 31.677114 60.136524]
 [ 9.455778 10.662097]]
聚类划分: [3. 3. 3. ... 0. 0. 0.]
```

## V. EVALUATION

The dataset used in this experiment is read from a CSV file (xclara.csv). It contains two features, V1 and V2, which are used to represent the data points in a two-dimensional space. The dataset has a certain number of data points, and their distribution in the 2D space determines the complexity of the clustering problem.

**Distribution:** The initial scatter plot of the data points shows that they have a somewhat irregular distribution. There are no obvious distinct groups at first glance, but as the clustering algorithm is applied, meaningful clusters can be identified.

**Range:** The values of the features V1 and V2 have a certain range. The maximum and minimum values of these features influence the selection of the initial cluster centers. For example, in the initial implementation of the K-means algorithm, the cluster centers are randomly selected within a range related to the maximum values of the features.

### Experimental results

To evaluate the accuracy of the clustering results, we use metrics such as the silhouette score and the within-cluster sum of squares (WCSS).

**Silhouette Score:** As mentioned in the comparison of different  $k$  values, the silhouette score provides an indication of how well-separated and compact the clusters are. For the  $k = 3$  case (original implementation), the silhouette score shows that the clusters are reasonably well-separated. As  $k$  changes, the silhouette score varies, indicating different levels of clustering quality.

**WCSS:** The within-cluster sum of squares is calculated for different values of  $k$ . For smaller  $k$  values, the WCSS is higher, indicating that the data points are more spread out within the clusters. As  $k$  increases to an appropriate value (e.g.,  $k = 3$ ), the WCSS decreases, showing that the data points are closer to their cluster centers. However, if  $k$  is increased too much, the rate of decrease of WCSS slows down, and may even start to increase due to over-fitting.

### Time Complexity

The time complexity of the K-means algorithm is mainly dominated by two parts: the assignment of data points to clusters and the update of cluster centers.

**Assignment Step:** In each iteration, we need to calculate the distance between each data point and all the cluster centers. If there are  $n$  data points and  $k$  clusters, this step has a time complexity of  $O(n * k)$ . Since the distance calculation is repeated in each iteration until convergence, and assuming the algorithm converges in  $t$  iterations, the total time complexity for this part is  $O(t * n * k)$ .

**Update Step:** Updating the cluster centers requires calculating the mean of the data points in each cluster. In the worst case, if all data points need to be considered for each cluster, the time complexity for this step is  $O(n * k)$  in each iteration. Considering  $t$  iterations, the total time complexity for this part is also  $O(t * n * k)$ .

Overall, the time complexity of the K-means algorithm is  $O(t * n * k)$ , where  $t$  is the number of iterations until convergence,  $n$  is the number of data points, and  $k$  is the number of clusters.

### Space Complexity

The space complexity of the algorithm mainly comes from storing the data points, the cluster centers, and some intermediate variables.

**Data Points:** Storing the  $n$  data points requires  $O(n)$  space.

**Cluster Centers:** Storing the  $k$  cluster centers requires  $O(k)$  space.

**Intermediate Variables:** Variables such as the cluster assignments and distances calculated during the algorithm's execution require additional space. However, these do not dominate the overall space complexity.

Overall, the space complexity of the algorithm is  $O(n + k)$ .

## CONCLUSION

**Algorithm Performance:** The K-means algorithm is able to identify meaningful clusters in the dataset. The accuracy of the clustering, as measured by metrics such as the silhouette score and WCSS, shows that appropriate selection of the number of clusters ( $k$ ) is crucial for good clustering results. When  $k$  is set to an appropriate value (e.g.,  $k = 3$  in

our experiments), the algorithm can produce well - separated and compact clusters.

**Complexity Analysis:** The time complexity analysis shows that the algorithm's running time depends on the number of data points, the number of clusters, and the number of iterations until convergence. The space complexity is mainly determined by the number of data points and clusters. This indicates that for large datasets or a large number of clusters, the algorithm may face challenges in terms of computational efficiency and memory usage.

**Future Work:** Future research could focus on improving the initialization method of the K - means algorithm to reduce the impact of random initialization on the final results. Additionally, exploring other clustering algorithms or hybrid methods could potentially lead to better clustering

performance, especially for datasets with more complex distributions.

## REFERENCES

- [1] Arthur, D., & Vassilvitskii, S. (2007). k-means++: The advantages of careful seeding. In Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms (pp. 1027-1035). Society for Industrial and Applied Mathematics.
- [2] Ester, M., Kriegel, H. P., Sander, J., & Xu, X. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In Kdd (Vol. 96, No. 34, pp. 226-231).
- [3] Zaharia, M., Chowdhury, M., Franklin, M. J., Shenker, S., & Stoica, I. (2010). Spark: Cluster computing with working sets. In Proceedings of the 2nd USENIX conference on Hot topics in cloud computing (pp. 10-10). USENIX Association.