

# Comprehensive Review of K-Means Clustering Algorithms

**Eric U. Oti<sup>1\*</sup>, Michael O. Olusola<sup>2</sup>, Francis C. Eze<sup>3</sup>, and Samuel U. Enogwe<sup>4</sup>**

<sup>1\*</sup>Department of Statistics, Federal Polytechnic, Ekowe Bayelsa State

<sup>2,3</sup>Department of Statistics, Nnamdi Azikiwe University, Awka Anambra State

<sup>4</sup>Department of Statistics, Michael Okpara University of Agriculture, Umudike, Nigeria

## ABSTRACT

*This paper presents a comprehensive review of existing techniques of k-means clustering algorithms made at various times. The k-means algorithm is aimed at partitioning objects or points to be analyzed into well-separated clusters. There are different algorithms for k-means clustering of objects such as traditional k-means algorithm, standard k-means algorithm, basic k-means algorithm, and the conventional k-means algorithm, this is perhaps the most widely used version of the k-means algorithms. These algorithms use the Euclidean distance as their metric and minimum distance rule approach by assigning each data point (objects) to its closest centroids.*

**Keywords:** Centroid, Cluster Analysis, Euclidean Distance, K-Means, Unsupervised Classification.

## 1. INTRODUCTION

Data clustering is an important topic of research and it has its applications in various field like statistics, data mining, computer science, pattern recognition, image processing, marketing, psychiatry, etc. [1-3]. Clustering is seen as purely a multivariate technique but it can also be applied to univariate and bivariate data. Clustering or grouping is done on the basis of similarities or distances [4] and it is one of the best approach of multivariate analysis and a common methodology for statistical data analysis.

Clustering (cluster analysis) was originated in anthropology by Driver and Kroeber [5] and was introduced to psychology in 1938 (Zubin [6]). Cattell [7] introduced mathematical procedures for organizing objects based on observed similarity. It was not until Sokal and Sneath's [8] publication of Principles of Numerical Taxonomy; that clustering methods gained widespread acceptances in the sciences, and motivated world-wide research on clustering methods and thereby initiated the publication of a broad range of books such as those of Fisher [9], Tryon and Bailey [10], Jardine and Sibson [11], Anderberg [1], Hartigan [12], Spáth [13-14], Aldenderfer and Blashfield [15], Romesburg [16], Fukunaga [17], Kaufman and Rousseeuw [18], Berkhin [19], Mirkin [20], etc.

K-means is the most popular clustering formulation in which the goal is to maximize the expected similarity between data items and their associated cluster centroids [21].

The purpose of this paper is to present a comprehensive review of existing techniques of k-means clustering algorithm made at various times.

The rest of this paper is organized as follows: section 2 discussed hierarchical and partitioning clustering methods as main group of cluster analysis. Section 3 discussed k-means clustering. Furthermore, section 4 discussed about related literature of k-means clustering. Finally, section 5 is the conclusion of the paper.

## 2. CLUSTER ANALYSIS

Cluster analysis or clustering is an unsupervised classification mechanism where a set of data, usually multidimensional is classified into groups (clusters) such that members of one cluster are similar to one another with respect to some predetermined criterion [12, 22-23].

Cluster analysis can be divided into two main groups which are based on the structure of their output namely: hierarchical non-hierarchical (Partitioning) clustering methods. Hierarchical clustering also known as hierarchical cluster analysis is an algorithm

that groups similar objects into groups called clusters. The clusters are merged (agglomerative methods) or split (divisive methods) step-by-step based on the similarity measure. The results of a hierarchical clustering method entails that agglomerative and divisive methods can be displayed graphically using a tree diagram known as dendrogram. The dendrogram shows all the steps in the hierarchical procedure which includes the similarities or distances at which clusters are merged. While partitioning clustering methods partition the data object set into clusters where every pair of object clusters is either distinct or has some members in common. Partitioning clustering begins with a starting cluster partition which is iteratively improved until a locally optimal partition is reached [12].

### 3. K-MEANS CLUSTERING

K-means is an iterative procedure that partition N objects into K disjoint clusters. K-means is perhaps the most widely used clustering method, and especially the best-known of the partitioning-based clustering methods that uses centroids for cluster presentation [23]. The quality of k-means clustering is measured through the within-cluster squared error criterion [24,25,26].

K-means algorithm is used to minimize the problem of k-means, and it has many variants which will be discussed next but to be able to use any of the k-means algorithm, the number of clusters present in the data need to be known; multiple runs or trials will be necessary to find the best number of clusters. There is no best k-means algorithm, as the tendency of generating global optimum depends on the characteristics of the data set, the size and also the number of variables in the cases. The k-means clustering methods have two phases of iteration namely: the assignment or initialization phase which involves an iterative process where each data point is assigned to its nearest centroid using Euclidean metric; the next is the centroid update phase, where clusters centroids are updated given the partition obtained by the previous phase. The iterative process stops when no data point change clusters or some maximum number of iterations is reached [21].

Forgy [27] proposed a batch algorithm called the traditional k-means algorithm; the algorithm is based on the minimization of the average squared Euclidean distance between the data points and the cluster's center known as centroid, where centroid is the center of a geometric object and it is seen as a generalization of the mean. The Forgys algorithm start by choosing the number of cluster k representing the cluster centers, it then assigns data point of the data set to the cluster having the closest centroid, update new centroids for each cluster by averaging the data points or objects belonging to the cluster, if there is no change in the cluster center, then the iteration stops.

Lloyd [28] proposed the standard k-means algorithm which is also a batch algorithm, the difference between Forgys algorithm and Lloyds algorithm is that Forgys algorithm treats data distribution as continuous while Lloyds algorithm treats data distribution as discrete case.

MacQueen [24] proposed the basic k-means algorithm which is an online algorithm, the algorithm is similar to Forgys and Lloyds algorithm when it comes to the initialization process but differs from the two algorithms when it comes to the update process. During the update of MacQueens algorithm, the centroids are updated by re-calculating the points any time there is a change in the centroid, and when each points is currently assigned to the cluster with the nearest centroid, the process stops.

Hartigan and Wong [29] proposed a conventional k-means algorithm which is a non-Forgy (or non-Lloyd) heuristic that updates cluster centers considering each points, rather than after each pass over the entire data set. This algorithm searches for the partition of data space with locally optimal within-cluster sum of squares of errors (SSE); which means that it may assign a point to another subspace, even if it currently belongs to the subspace of the closest centroid. If the centroid has been updated for each data point included, the within-cluster sum of squares for each data point if included in another cluster is calculated. If one of the cluster sum of squares (SSE 2 in the equation below, for all  $i \neq 1$ ), the point is assigned to this new cluster

$$SSE\ 2 = \frac{N_i \sum_j^k \|x_{ij} - c_i\|^2}{N_i - 1} < SSE\ 1 = \frac{N_1 \sum_j^k \|x_{ij} - c_i\|^2}{N_1 - 1}$$

Where  $N_i$  is the number of points included in cluster  $k$ ,  $x_{ij}$  is the  $j$ th point in the  $i$ th cluster and  $c_i$  is the  $i$ th point in the cluster center. The iteration continues until no point changes cluster.

### 4. RELATED LITERATURE

Jancey [30] proposed a variant which is a modification for the Forgy's k-means algorithm (cf. Anderberg, [1]) which is expected to accelerate convergence and inferior local minima. In this variant, the new cluster center is not the mean of the old and added points, but the new center is updated by reflecting the old center through the mean of the new cluster.

In order to avoid poor local solutions, a number of genetic algorithm based methods have been developed [31, 32]. Likas et al. [33] developed the global k-means clustering algorithm which is a deterministic and incremental global optimization method. It is also independent on any initial parameters and employs k-means procedure as a local search procedure, since the exhaustive global k-means method is computationally expensive.

Faber [34] proposed a variant of the Lloyd's k-means algorithm called the continuous k-means algorithm. The reference points in the continuous k-means algorithm are chosen as a random sample from the whole population of the data point while in the standard k-means algorithm, the initial reference points are chosen more or less arbitrarily. During the update process, the continuous k-means algorithm examines only a random sample of the data points while the standard k-means algorithm examines all of the data set in sequence. If the data set is very large and the sample is a representative of the data set, then the continuous k-means algorithm should converge much faster than the algorithm that examines every point in sequence.

Kanungo et al. [35] presented a simple and efficient implementation of Lloyd's k-means clustering algorithm which they called the filtering algorithm. The filtering algorithm is easy to implement which requires a kd-tree (cf. Bentley [36]) as the only major data structure. A kd-tree is a binary tree, which represents a hierarchical sub-division of the point set's bounding box using axis aligned splitting hyperplanes. Each node of the kd-tree is associated with a closed box, called a cell. The root's cell is the bounding box of the point set. If the cell contains at most one point (or, more generally, fewer than some small constant), then it is declared to be a leaf. Otherwise, the root's cell is splitting into two hyper-rectangles by an axis orthogonal hyperplane. The points of the cell are then partitioned to one side or the other of this hyperplane. The resulting sub-cells are the children of the original cell, thus leading to a binary tree structure.

Bagirov and Mardaneh [37] proposed a new variant of the global k-means algorithm which is known as the modified global k-means (MGKM) algorithm because it is said to be effective for solving clustering problems in gene expression data sets. In their algorithm, a starting point for the  $k^{th}$  cluster center is computed by minimizing the so-called auxiliary cluster function. The effectiveness of this algorithm highly depends on its starting point. The algorithm computes clusters incrementally and to compute k-partition of a data set, it uses  $k - 1$  cluster centers.

Nazeer and Sebastian [38] discussed in their paper about one major drawback of k-means algorithm, they proposed an enhanced method that deals with improving the accuracy and efficiency of k-means algorithm. Both the phases of the original k-means algorithm were modified. The initial centroids are determined systematically so as to produce clusters with better accuracy in the first phase. The second phase makes use of a variant of the clustering method discussed in Fahim et al. [39]. It starts by forming the initial clusters based on the relative distance of each data point from the initial centroids. These clusters are subsequently fine-tuned by using a heuristic approach thereby improving the efficiency.

Huang et al. [40] proposed a k-means type clustering algorithm that can automatically calculate variable weights and it is referred to as weighted-k-means (W-K-Means). The weighted-k-means adds a new step to the basic k-means algorithm to iteratively update the variable weights based on the current partition of the data and also a formula for weights calculation was proposed as well. The variable weights produced by the proposed weighted-k-means algorithm measured the importance of variables in clustering and can be used in variable selection in data mining applications where large and complex real data are often involved.

Amorim [41] proposed the constrained Minkowski Weighted K-Means algorithm which calculates cluster specific feature weights that can be interpreted as features rescaling factors. Naturally, the Minkowski weighted k-means (MWK-Means) algorithm requires a Minkowski exponent,  $p$ , which can be approximated via semi-supervised learning [42]. Weight  $w_{kw}$  was introduced, which depends on both cluster  $k$ , and feature  $v$ , allowing a given feature  $v$ , to have different weights at different cluster  $k$ ; also, the use of the Minkowski distance to the power of  $p$  was introduced, analogous to the Euclidean squared distance  $d_p(y_i, c_k) = \sum_{v=1}^V W_{kv}^p |y_{iv} - c_{kv}|^p$  where  $v$  represents the features and  $p$  is the Minkowski exponent,  $w_{kv}^p$  is the weight variable to take into account the Minkowski exponent  $p$ . Wagstaff et al. [43] introduced constrained clustering k-means which makes use of limited amount of background knowledge by applying pairwise must-link and cannot-link rules to entities and likewise is the Minkowski weighted k-means.

Oti et al. [44] proposed a new k-means clustering method that adds cluster centers one by one as clusters are being formed, such that when a point is moved from the initial configuration, the cluster centroids will be updated or recalculated before computing

the squared distance. The  $i$ th coordinate, where  $i = 1, 2, \dots, k$  of the centroid is updated using  $\bar{c}_i$ , new =  $\frac{N_k \bar{c}_i + \bar{c}_{ij}}{N_{k+1}}$  if the  $j$ th point is added to the cluster and  $\bar{c}_i$ , new =  $\frac{N_k \bar{c}_i - \bar{c}_{ij}}{N_{k-1}}$  if the  $j$ th point is removed from the cluster. Here  $N_k$  is said to be the number of points (cases) in the old cluster with centroid  $\bar{c}^l = (\bar{c}_1, \bar{c}_2, \dots, \bar{c}_k)$  or perhaps the cluster size and centroid  $\bar{c}_k$  is a multidimensional vector which minimizes the sum of squared distance to clusters elements.

## 5. CONCLUSION

In this paper, we have reviewed existing techniques in k-means clustering. This work shows that there are several variants of k-means clustering algorithms from sixties to recent times which have addressed some drawback of k-means algorithms.

In the future, we will look at the computational time complexity of some of the variants of k-means clustering algorithms and its analysis in terms of relative accuracy and efficiency.

## ACKNOWLEDGMENT

The authors wish to thank the referees for their worthwhile comments and suggestions.

## REFERENCES

1. M. R. Anderberg, "Cluster Analysis for Applications", New York: Academic Press, (1973).
2. S. Brohee and J. V. Helden, "Evaluation of clustering algorithms for protein-protein interaction networks", BMC Bioinformatics, (2006), 7(1): 488.
3. B. Everitt, S. Landau, M Leese. and D. Stahl, Cluster Analysis, 5<sup>th</sup>ed., John Wiley and Sons, (2011).
4. R. A. Johnson and D. W. Wichern, Applied Multivariate Statistical Analysis, 5<sup>th</sup> ed., Eaglewood Cliffs, NJ: Prentice-Hall, (2002).
5. H. E. Driver and A. L. Kroeber , "Quantitative Expression of Cultural Relationships," University of California Publications of American Archaeology and Ethnology, 1932, vol. 4, pp.211-256.
6. J. Zubin, "A technique for measuring like-mindedness", The Journal of Abnormal and Social Psychology. 1938, vol. 4. pp508-516.
7. R. Cattell, " $r_p$  and other coefficients of pattern similarity", Psychometrika, 1949, vol. 4. pp.279-298.
8. R. R. Sokal and P. H. A.. Sneath, Principles of Numerical Taxonomy. San Francisco: California, (1963).
9. W. D. Fisher, Clustering and Aggregation in Economics, Johns Hopkins Press, Baltimore, Maryland., (1968)
10. R. C. Tryon and D. E. Bailey, Cluster Analysis., McGraw Hill, New York, (1970)
11. N. Jardine and R Sibson, Mathematical Taxonomy, John Wiley and Sons, Ltd, Chichester, (1871)
12. J. A. Hartigan ,Clustering Algorithms, John Wiley & Sons. Inc., New York, (1075).
13. H. Spáth, Cluster Analysis Algorithms. West Sussex, UK: Ellis Horwood Limited, (1980).
14. H. Spáth, Cluster Dissection and Analysis: theory, FORTRAN programs, examples. (Translator: Johannes Goldschmidt.) Ellis Horwood Ltd Wiley, Chichester, (1985).
15. M. S. Aldenderfer, and, R. K.. Blashfield, Cluster Analysis. Beverly Hills, CA: Sage Publications, (1984).
16. C. Romesburg, Cluster Analysis for Researchers, London, Wadesworth, (1984).
17. K. Fukunaga, Introduction to Statistical Pattern Recognition, 2<sup>nd</sup> ed., Academic Press, (1990).
18. L. Kaufman and P. J. Rousseeuw, Finding Groups in Data: An Introduction to Cluster Analysis, John Wiley & Sons, Inc. New York, (1990).
19. P. Berkhin, A Survey of Clustering Data Mining Techniques. In: Kogan J., Nicholas C., Teboulle M. (eds) Grouping Multidimensional Data, Springer, Berlin, Heidelberg, (2006).
20. B. Mirkin, Clustering: A Data Recovery Approach, 2<sup>nd</sup> ed., Chapman and Hall/CRC, (2013).
21. N. Slonim, E. Aharoni and K. Crammer, "Hartigan's K-Means Versus Lloyd's K-Means-Is It Time for a Change? Proceedings of the Twenty-Third International Joint Conference on Artificial Intelligence, 2013, pp. 1677-1684.
22. A. K. Jain and R. Dubes, Algorithms for Clustering Data. Englewood Cliffs, NJ: Prentice-Hall, (1988).
23. V. Estivill-Castro Why so many clustering algorithms: a position paper. ACM SIGKDD Explorations Newsletter, 2002, 4(1), 65-75.
24. J. MacQueen, "Some methods for classification and analysis of multivariate observations". In: Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, 1967, vol 1, 281-297.
25. C. Yuan and H. Yang, "Research on k-value selection method of k-means clustering algorithm", Multidisciplinary Scientific Journal, 2019, 2(2), 226-235.

26. T. Hastie, R. Tibshirani and J. Friedman., The Elements of Statistical Learning: Data Mining, Inference and Prediction, 2<sup>nd</sup> ed., Springer-Verlag, (2009).
27. E. W. Forgy, "Cluster analysis of multivariate data: efficiency versus interpretability of classifications", Biometrics, 1965, vol. 21, 768-769.
28. S. Lloyd, "Least squares quantization in PCM. IEEE Transaction on Information Theory", 1982, 28(2), 129-137.
29. J. A. Hartigan and M. A. Wong, "Algorithm AS 136: A k-means clustering algorithm", Journal of the Royal Statistical Society. Series C (Applied Statistics), 1979, 28(1):100-108.
30. R. C. Jancey, "Multidimensional group analysis", Australian Journal of Botany, 1966, 14(1), 127-130.
31. K. Krishna and M. Murty, "Genetic K-Means algorithm", IEEE Transactions on Systems, Man, and Cybernetics-Part B: Cybernetics, 1999, 29(3): 433-439.
32. S. Bandyopadhyay and U. Maulik, "An evolutionary technique based on k-means algorithm for optimal clustering in  $R^N$ ", Information Science, 2002, vol. 146, pp. 221-237.
33. A. Likas, N. Vlassis and J. Verbeek, "The global k-means clustering algorithm", Pattern Recognition 2003, 36(2), 451-461.
34. V. Faber , "Clustering and the continuous k-means algorithm", Los Alamos Science, 1994, vol. 22, 138-144.
35. T. Kanungo, D. Mount, N. Netanyahu, C. Piatko, R. Silverman and A. Wu, "An efficient k-means clustering algorithm: Analysis and Implementation", IEEE Transactions on Pattern Analysis and Machine Intelligence, 2002, 24(7), 881-892.
36. J. L. Bentley, "Multidimensional binary search trees used for associative searching: Communication of the ACM, 1975, 18(9), 509-517.
37. A. M. Bagirov and K. Mardaneh, "Modified global k-means algorithm for clustering in gene expression datasets Conference Proceedings Workshop on Intelligent Systems for Bioinformatics, 2006, Vol. 73: 23-28.
38. [38] K. A. A Nazeer. and M. P. Sebastian, "Improving the accuracy and efficiency of the k-means clustering algorithm, Proceedings of the World Congress on Engineering, 2009, vol. 1, 1-3.
39. [39] A. M. Fahim A. M. Salem, F. A. Torkey and M. A. Ramadan, "An Efficient enhanced k-means clustering algorithm," Journal of Zhejiang University, 2006, 10(7): 1626-1633.
40. J. Z. Huang M. K. Ng, H. Rong and Z. Li, "Automated variable weighting in k-means type clustering". In: IEEE Transactions on Pattern Analysis and Machine Intelligence, 2005, 27(5), 657-668.
41. R. C. Amorim, "Constrained clustering with Minkowski weighted k-means", In: Proceedings of the 13<sup>th</sup> IEEE International Symposium on Computational Intelligence and Informatics, 2012, pp.13-17.
42. R. C. Amorim and B. Mirkin, "Minkowski metric, feature weighting and anomalous cluster initializing in k-means clustering", Pattern Recognition, 2012, vol. 45, pp.1061-1075.
43. K. Wagstaff, C. Cardie, S. Rogers and S. Schroedl, "Constrained k-means clustering with background knowledge", In Proceedings of the 8<sup>th</sup> International Conference on Machine Learning. 2001, pp. 577-584.
44. E. U. Oti, S. I. Onyeagu and R. A. Slink, "A modified k-means clustering method for effective updating of cluster centroid", Journal of Basic Physical Research, 2019, 9(2), 123-137.

\*Corresponding Author: Oti, Eric U. [eluchcollections@gmail.com](mailto:eluchcollections@gmail.com);