# Clustering of Multidimensional Big Data using Enhanced K-Mean Algorithm

**Jagdish Kushwah,  Shailesh Jaloree, R.S.Thakur**

*Abstract: One of the basic issues with K-means clustering is that it just merges to nearby optimum which is simpler than comprehending for worldwide optima however can prompt less ideal union. This is especially valid for enormous information as the underlying focuses assume a significant job on the exhibition of this calculation. The paper proposes a novel K-means clustering algorithm which presents a technique to discover advanced area of beginning focuses and introductory number of bunches. This outcome in getting last arrangement of bunches to meet internationally, encouraging quick and exact grouping over enormous datasets. Distributed computing executes huge scope and complex processing. A lot of information are economically and proficiently broke down by utilizing parallelism method. To get parallelism and versatile registering, using Amazon Web Services with R Studio Flexible Process Cloud occasion which partitions the activity among different hubs. The proposed system presents an exceptionally serious exhibition taking significant less calculation time and financially savvy. It very well may be contrasted with complex Hadoop Disseminated Record Framework and MapReduce A significant disadvantage with Apache Hadoop is its MapReduce worldview that is exceptionally open when a procedure emphasizes number of times. R performs execution inside memory which is quicker and less mind boggling when contrasted with Read/Keep in touch with the circle over and again in MapReduce. The examination work is mimicked on some well known genuine datasets from UCI AI storehouse. The outcomes affirm that the proposed work models a vigorous and adaptable procedure for grouping huge datasets.*

*Index Terms: Artificial Intelligence, Big Data, Cloud Computing, K-means, MapReduce; R.*

## I.  INTRODUCTION

Prior, data innovation didn't have present day's immense skyline and was restricted to colleges, government foundations, innovation associations and huge business houses before the development of distributed computing. The advancement of distributed computing served data innovation to the majority with much decreased expense and huge decisions out of the applications accessible [1]. With this promotion of data innovation, it is presently being seen that gigantic volume of data are created through online life, messages, Web of Things (IoT), web crawlers, exchange records, money related markets, interactive media and more in either organized or unstructured arrangement. This made

  **Jagdish Kushwaha\***,  is currently pursuing Ph.D.  degree program in computer applications in BU Bhopal, ,+919630831212. E-mail: mca.jagdish@gmail.com
  **Dr.(Prof) Shailesh Jaloree** ,Professor department of computer science & applied                     mathematics                     SATI VIDISHA(MP),9424418356.shailesh_jaloree@rediffmail.com.
  **Dr (Prof) R.S.Thakur** ,Professor and HOD department of mathematics, Bioinformatics & Computer Application MANIT Bhopal(MP) 9826241996.ramthakur2000@gmail.com.

another worldwide universe of information known as Large Information [1]. Today business world basically requires Huge Information and Distributed computing [2]. Huge Information gives wanted business experiences and cloud makes it conceivable to store and investigate this information by giving computational powers for all intents and purposes to the clients. Additionally, this innovation guarantees that such information can be effectively available. Contemporary advancements have made it conceivable to dissect very huge and shifted datasets which are perplexing also, for instance Apache Hadoop. This is making a lot of significant worth for associations and aiding business development. Openness accommodation and cost viability offered by cloud innovation in dealing with these gigantic measures of information is making a lot of significant worth for associations and encouraging in business development [2]. Grouping is a major information investigation method to discover valuable examples from a huge database [3]. These examples are valuable for the scientists and information laborers, for example, money related examiner and director to take right choices. Applying grouping to large information is a troublesome assignment since it forces certain difficulties. High

Computational expenses bring about when enormous sizes of information for example terabytes to petabytes of information are mulled over for grouping. To accomplish quality outcomes as speedy as would be prudent, the matter of significance is the best approach to deal with this issue. K-implies grouping is a solo AI calculation. It is favored as the engaging quality lies in its effectiveness with $O(n*K*i*a)$, where n, K, I and an equivalents number of information focuses, groups, emphasess and properties individually. In any case, it will in general be non-definable even with littlest significantly number of groups [4]. Hence, adaptability is a significant test in breaking down huge information applying K-implies.

## II.  LITERATURE REVIEW

The exploration for grouping large information is a ceaseless wonder and much has been done right now the ongoing hardly any years. Scientists have been attempting to improve further and advance huge information examination. Analysts in [5] dealt with multidimensional enormous dataset to talk about issues while bunching huge datasets with MapReduce. They proposed an equal bunching technique through Hadoop MapReduce system which concentrated on a key factor of diminishing I/O and system costs. Quicker execution and scaling were two pivotal focuses to manage. Getting from the conduct of winged animals in a herd, one of the arrangements has been found.

In the ongoing past, the designs preparing unit (GPU) picked up fascination.

It can take care of rapidly issues of parallelism. Analysts in [6] applied this idea utilizing CUDA stage from NVIDIA GPU. So also, in [7] utilized DBSCAN calculation on GPU to increase superior.

In 2013, Xiao Cai and co-creators [8] compactly pondered assortment of information from different sources, each source speaking to an alternate part of the information. Since each source has its own individual angle, accordingly, bunching of large information here gets troublesome. They introduced a novel strategy to consolidate amalgamated portrayal of enormous datasets. Moreover, analysts in [9] introduced a calculation to improve K-implies grouping by acquiring introductory centroids. Kodinariya1and Makwana looked into K-implies grouping on deciding number of bunches [10]. Kim et al., in 2014 [11] proposed a bunching calculation which was thickness based. It had the option to investigate datasets of various densities.

In [12] Cui et al., the analysts referenced MapReduce to be temperamental because of emphasess which include restarting occupations over and over. They proposed another model of handling enormous information utilizing K-implies calculation without emphasess. Utilization of K-implies bunching techniques in identifying similitude between records or copyright infringement is another subject of significance which has been taken up by K.Vani and D. Gupta in their paper [13] distributed in the year 2015. For contrasting records the major prerequisites are assortment of relatively huge number of information, classification and developing programmed techniques for fast examination just as investigation. They used various varieties of K-implies calculation to contrast and N-gram strategy and vector space model technique to assess at long last execution and investigation by profiting dataset from Dish 2013. When taken in execution of the calculation effectiveness and exactness are resolved in their procedure. Besides, Await and Shedge in [14] additionally chipped away at grouping content archives for comparability check. Pondering that yield is completely subject to contribution of number of bunches, they insisted catchphrase as information. They made subset of reports for accomplishing wanted outcomes.

Tsai with co-researchers in [15] write that one in all the most reasons of failure in bunch giant datasets by ancient clustering strategies is that the majority of them are designed for centralized systems. In their paper they planned to unravel this drawback by associate algorithmic rule that they termed as MapReduce region (MRBH) which accelerated bunch.

Problem of initial cluster centers in K-means algorithmic rule was tried to be resolved by researchers in Shanghai dialect et al. [16] by sampling the big dataset and used biconvex hull and opposite Chung points. They applied MapReduce framework for parallel execution of the algorithmic rule.

X. Cui in next analysis work [17] once more improved K-means. in contrast to standard K-means algorithmic rule the new algorithm used each inside cluster concentration and between cluster segregation.

Comparing K-means and K-medoids bunch algorithms J. Kaur and H. Singh bestowed a brand new hybrid approach furthering birch and K-means. Firstly, a tree of hierarchal bunch is formed. Then K-means partitioning is applied to cut back the quantity of clusters for economical performance of it. Utility of this algorithmic rule has been found blue banking sector [18].

Clustering of pictures additionally has become a distinguished field of analysis attributable to its business and social importance presently. Dhanachandra and co-authors in [19] worked on image partitioning victimisation K-means. In their method, image is partly elongated to boost its quality within the beginning. within the next step, supported the potential worth of the information points, centers of the clusters are generated. Finally those parts that don't seem to be needed are filtered from the image.

Authors of [20] compared K-means and K-medoids. They applied nearly 10 thousand transactions of KEEL dataset repository. Results delineate that K-medoids is best than K-means in terms of execution time, noise decrease and selecting initial center.

Automating the quantity of clusters was evolved by authors in [21]. it had been once more associate algorithmic rule {for big|for giant|for giant} information analytics dividing large datasets into K partitions. This discovery is simulated victimisation Spark platform.

Further, researchers studied K-means, Fuzzy C-Means (FCM), hierarchal bunch algorithmic rule like Balanced unvarying Reducing and bunch victimisation Hierarchies (BIRCH) and grid based mostly clustering algorithm bunch In QUEst (CLIQUE) that are the prevailing effective algorithms. Ajin and Lekshmy [22] compared them in context of huge information. a brand new inhume and intra K-means bunch (KM-I2C) algorithmic rule was developed in [23] by dynamical clustering distance metric that used parallelization tools through Hadoop.

Rehioui et al. [24] in their analysis bestowed a brand new version of DENCLUE referred to as DENCLUE-IM that avoids complexities of different DENCLUE algorithms for quick calculation of huge information. They compared the planned approach with DENCLUE, DENCLUE-SA and DENCLUE-GA.

A concept of huge information in kind of streams and the way to method them was bestowed by Giacomo Aletti and also the author in their analysis in [25]. Datasets having parts of characteristics that notice correlation were thought-about. They used Mahalanobis distances for assignment of information to clusters estimating total range of clusters Scientists in [26] took a shot at bunches where limits are not firm with sureness. There are lower and upper approximations. The calculation proposed by them depended on weighted separation measure with Gaussian capacity for registering the new community for each group. Without considering choice of beginning focuses of bunches, Vijay et al. in [27] introduced a calculation which they named as Change Based Moving K-implies (VBKM). They applied another separation measure. Additionally, the exploration work utilized an alternate methodology for moving information focuses between groups to limit inside bunch separation.

The greater part of previously mentioned inquires about depend on Hadoop worldview. A significant downside with Apache Hadoop is its MapReduce worldview that is profoundly responsive when a procedure repeats number of times. Cycle requires composing back information to the record framework and subsequently cost of contribution just as yield increments. Subsequently MapReduce based K-implies turns out to be a lot of expensive. Besides, K-implies calculation is should have been executed on various occasions to get the ideal number of groups. Multifaceted nature of these calculation further increments as and when it is applied to enormous information. This exploration work recommends to process tremendous measure of information by embracing parallelism through AWS distributed computing condition. R performs execution inside memory which is quicker and less mind boggling when contrasted with Read/Keep in touch with the plate over and again in MapReduce. Hubs of R Studio Server occasion are utilized for separating calculation occupations to be executed through web dissemination. Rather than randomizing starting seed setup, the investigation proposed to acquire the advanced area of beginning seeds and introductory number of bunches which is eminent for huge information grouping. Also, the examination finds the last arrangement of bunches by blending the underlying groups utilizing a benchmark.

## III. EXPLORING STANDARD K-MEANS ALGORITHM

Let us initially consider the conventional K-implies calculation for producing groups of a dataset. It is regular information that K-implies isolates the informational collection into K parts where K indicates a positive whole number and stands as a client contribution to the calculation. Each group has a centroid. The calculation checks the places of these centroids as the calculation emphasizes. Arbitrary qualities are put to introduce the centroids before the main cycle. The calculation stops when centroids areas become static during emphasis [28]. This notable calculation performs two stages for every emphasis:-
1. Allocate each item xi to a nearest bunch centroid cj. This portion is acquired by Euclidean measure between the article and the bunch's centroid (got at a past cycle).
2. Update the centroids of the bunches dependent on new groups individuals.
The over two stages proceed to a limited number of redundancies, until there is no change or modification of the group habitats. Assume {x1,....,xn} are the given perceptions of a multidimensional gigantic dataset D with d measurements. The goal is to decide a lot of K or group centroids C = {c1,… .,cK} that limits the inside group contortion given in condition 1:

$$W(C) = \sum_{i=1}^{n} \|X_i - c\|^2 \qquad (1)$$

The above equation defines a cluster assignment rule as

$$f(x) = \arg\min_{f \in \{1,\dots,K\}} X - a_f^2 \qquad (2)$$

The problem of choosing K, the number of clusters and cluster centers, can be considered to be a model determination problem.
Algorithm: Standard K-means
Input: Dataset D={x1,….,xn} and

K (number of clusters to be generated)
Step 1: Select centers C randomly from D
Step 2: If u > iter then /* loop repeats until convergence
Step 3: for each xi compute distance from all centers Cj dist(xi, Cj)
Step 4: Assign xi to closest Cj (min(dist(xi, Cj))
Step 5: Calculate new centroid

$$n_i = \frac{1}{\phantom{x}} \sum X_i$$

Step 6: end for
Step 7: End if
Step 8: Exit

This is a standard method of grouping datasets utilizing K-implies calculation. However, the calculation is powerless against specific issues:

1. The strategy to introduce the centroids isn't determined. Arbitrarily picking K of the examples is predominant.

2. Nature of grouping relies upon the underlying qualities for the centroids and it happens that problematic allotments are discovered number of times. Attempting various distinctive beginning stages is just an answer which is utilized for the most part.

3. Estimation of K influences results.

4. The standard calculation is straightforward yet it has high time multifaceted nature when the datasets are huge multidimensional. Under this condition the memory of a solitary machine couldn't be adequate.

## IV. SCALING UP K-MEANS VIA MAPREDUCE

Standard information warehousing, the board and examination frameworks need instruments to break down Enormous Information. Huge Information is put away in dispersed document framework models because of its particular attributes. Apache's Hadoop is generally utilized for putting away and overseeing Huge Information. Investigation of Huge Information is a significant issue as it includes enormous circulated document frameworks which ought to be adaptable and deficiency tolerant. MapReduce is generally applied for the proficient examination of Huge Information [29].

These days, K-implies bunching is being utilized for Large Information investigation in MapReduce system by numerous information researchers. The pivotal part of usage of K-implies calculation is the plan of Mapper and Reducer capacities. Mapper plays out the activity of allotting perceptions to nearest bunch focus and Reducer overhauls group focuses as mean of relegated perceptions.

For a given xi in the dataset, the Guide stage processes the squared separation among xi and each arbitrarily chosen group places and gets the mean μi which limits this separation. A key-esteem pair is transmitted with this present mean's file I and the information point xi as the worth.

$$Z_i \leftarrow \arg\min_j \|\mu_j - x_i\|_2^2 \; \forall \; emit(z_i, x_i)$$

The Lessen stage is re-focusing step. The Reducer aggregates all the doled out perceptions and partitions by the absolute number of perceptions appointed to a similar group. This gives normal of doled out perceptions therefore [29].

$$\mu_j = \frac{1}{n_j} \sum_{i, z=k} x_i$$

Sum=0
Count=0
For x in x –in-cluster j
Sum+=x

Count+=1

emit (j, sum/count)

## V. PROPOSED WORK

To accomplish worldwide optima in dissecting large information, it is important to acquire improved beginning communities. For this, the examination work presents a few enhancements in the grouping procedure of conventional K-implies calculation. The work comprises of two procedures Algorithm1 2. Algorithm 1 figures and creates introductory focuses. It takes a broad estimation of K. For group 1, the Algorithm takes a perception arbitrarily from the preparation dataset and the point comparing to this perception is the primary beginning seed. At that point, Algorithm 1 applies probability capacity to discover balance K-1 focuses which is given by condition 3 underneath:

$$L(x) = \frac{\sum_{i=1}^{n} \sum_{j=1}^{K} |x_i - c_j|}{\sum_{i=1}^{n} \sum_{j=1}^{K} \max |x_i - c_j|}$$

(3)

To locate the underlying focuses, following calculation is applied:

Algorithm 1: Finding Ideal Area of Beginning Seeds

Stage 1: Select c1 haphazardly from D
Stage 2: On the off chance that |center|< K, at that point Next inspecting of residual focuses to be taken up by applying condition (3)
Stage 3: End If
Stage 4: Leave

When the underlying focuses are chosen information focuses are relegated to them as it is acted in the essential K-Means Algorithm. Algorithm 2 thinks about the separation between these enormous quantities of beginning centroids with edge esteem and unions these focuses so as to diminish number of cycles to unite ideally and locate the last arrangement of all inclusive upgraded great quality bunches.

$$\omega = \frac{1}{K} \sum_{i=1}^{K} c_i$$

(4)

where K is the quantity of groups created by algorithm 1 and ci speaks to the underlying clusteroids. The mean separation between any two sets of focuses gives the edge esteem. The bunches are converged by the benchmark as given. The separation between all the centroids is figured. These separations are contrasted and the edge esteem given in condition 4. In the event that it is seen that edge esteem is more prominent than the separation between centroids then they are joined to shape a unit group. Normal of all such joined focuses speaks to the new group centroid.

Algorithm 2: Consolidating Centroids dependent on an Edge Worth
Info: D and beginning K communities
Yield: C= {c1,c2,… ..,ck} (set of group centroids) L= {l(d) | d= 1,2,… ..,n} (set of bunch names of D)

Stage 1: Set K, focuses from D executing Calculation 1
Stage 2: centroid = focus
Stage 3: on the off chance that u > iter, at that point/* circle rehashes until combination
Stage 4: for x=1 to n
Stage 5: for y=1 to k
Stage 6: distance[y] = ||d[x] – center[y]||
Stage 7: end for
Stage 8: distancemax=max separation y
Stage 9: centroidx= d∈ clusterx clusterx
Stage 10: end for
Stage 11: end if
Stage 12: while for all cx
Stage 13: while for all cy
Stage 14: in the event that x = y, at that point add x to merge[x] proceed
Stage 15: end if
Stage 16: in the event that dist x,y ≤ω, at that point add y to combine x erase y from centroid
Stage 17: end if
Stage 18: end while
Stage 19: register new centroid centroidfinal= mergex
Stage 20: end while Step
21: Exit

## VI. EXPLORATORY STRUCTURE

### A. Machines

The exploration work is reproduced on Amazon Elastic Compute Cloud (EC2). It grants versatile processing power in the Amazon Web Services (AWS) cloud. A virtual processing condition format is chosen from EC2. This is c3.xlarge case of R Studio Server. It has following qualities: group of 14 ECUs, 4 vCPUs having 2.8 GHz Intel Xeon E5-2680v2, RAM allocated is 7.5 GB memory with 2x40 GB Stockpiling Limit. Ubuntu-16.04-LTS-64 piece working framework was introduced on every hub. I/O execution is high. Coding of the proposed calculation is written in R content. All the cases are propelled utilizing a similar Amazon Machine Image (AMI).

### B. Datasets

Utilized in the investigation are genuine datasets imported from UCI Machine Learning Archive. Fig. 1 shows the info dataset. The datasets are YouTube Multiview Computer games Dataset and Every day and Sports Exercises Datasets. After information cleaning, these crude datasets are changed into reasonable arrangement. YouTube Multiview Computer games dataset is isolated into Data1, Data2 and Data3. Day by day and Sports Exercises datasets are parceled into DSA1, DSA2, DSA3 and DSA4 datasets. The datasets change in size from 128 MB to 934MB. Every day and Sports Exercises Dataset is gathered by giving sensors to screen movement of sports exercises at the rate 19 every day. Information esteems are genuine as recorded by the authorities. The insights concerning the datasets are given in Table 1. YouTube Multiview Computer games Dataset has almost 120 thousands records with 1 million properties. Qualities of this dataset are multivariate. Traits contain whole number and genuine qualities.

**Table I: Statistics of 7 large-scale real datasets**

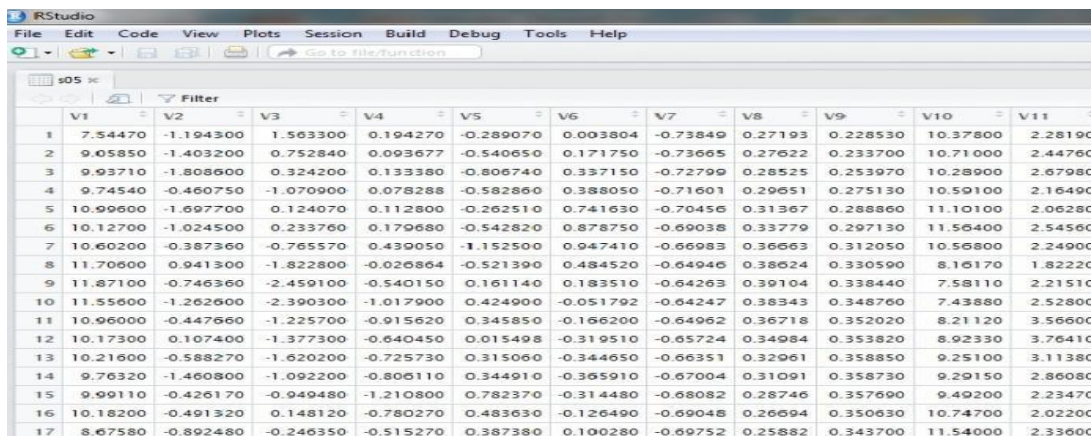| Dataset | No. of Samples | No. of Dimensions | Size |
|---------|----------------|-------------------|------|
| Data1 | 97,935 | 838 | 706MB |
| Data2 | 97,935 | 838 | 706MB |
| Data3 | 97,935 | 838 | 706MB |
| DSA1 | 2,85,000 | 45 | 128MB |
| DSA2 | 5,70,000 | 45 | 257MB |
| DSA3 | 1,140,000 | 45 | 467MB |
| DSA4 | 2,280,000 | 45 | 934MB |
| | | | |



**Fig. 1 Input dataset to the proposed system**

## VII.    RESULTS AND DISCUSSION

Comparing the presented work with standard K-means algorithm and K-means MapReduce (MR) algorithm, we sum up the experimental results run on the setup.
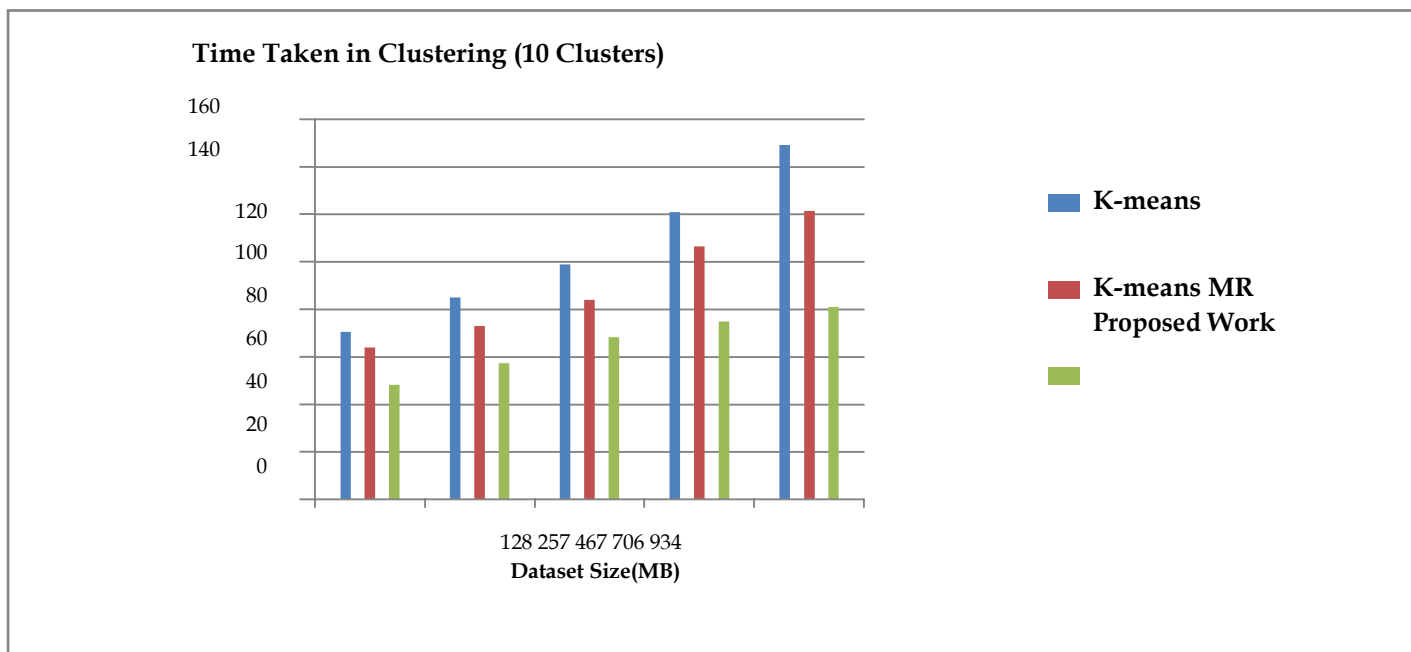
### A.   Performance Assessment



**Fig. 2 Comparison of Execution Time for 10 Clusters**

**Fig. 2 shows execution times of the datasets in generation of 10 clusters by the three algorithms.**



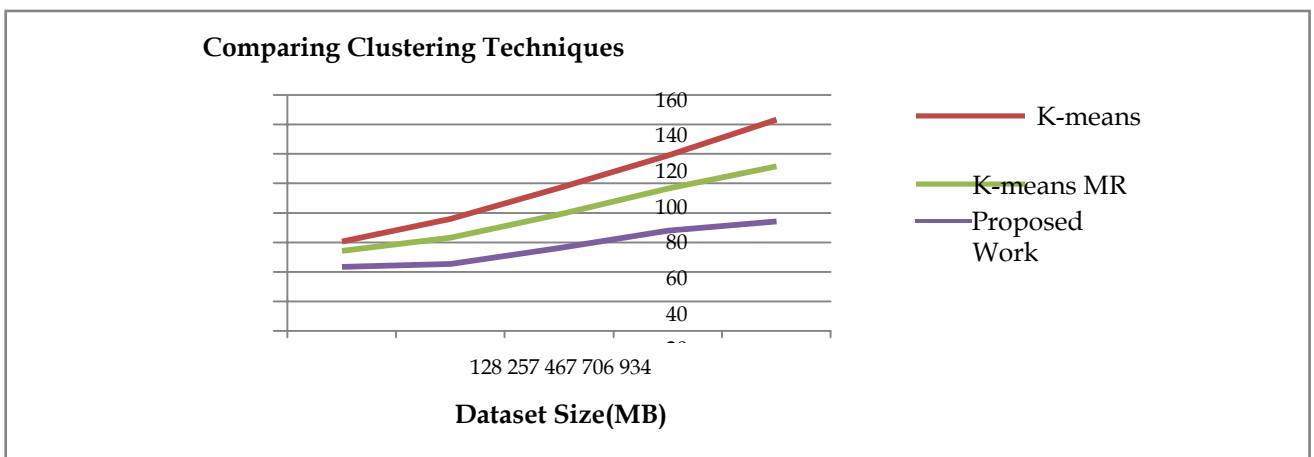**Fig. 3 Comparison of Execution Time for 15 Clusters**



**Fig. 4 Comparing Execution Time with Dataset size**
**Table2: Execution time in seconds**

| Dataset | Standard K-means | K-means MapReduce | Proposed Work |
|---------|------------------|-------------------|---------------|
| 128 | 60.57 | 54.13 | 43.37 |
| 257 | 75.88 | 63.04 | 45.42 |
| 467 | 97.00 | 78.92 | 56.25 |
| 706 | 118.97 | 96.54 | 67.94 |
| 934 | 143.25 | 111.49 | 74.16 |

Table 2 presents execution times of the three contrasting calculations for 15 bunches and dataset size going from 128 MB to 934 MB. In Fig. 4 the x-hub speaks to dataset estimates in super bytes and y-hub speaks to execution time like a flash. The red bend in the diagram speaks to execution time of standard K-implies which changes between 60.57 seconds to 143.25 seconds. The green bend in the diagram speaks to execution time of K-implies MapReduce calculation which differs between 54.13 seconds to 111.49 seconds. The blue bend in the diagram speaks to execution time of our proposed work which shifts between 43.37seconds to 74.16 seconds. Clearly the proposed algorithm reduces astoundingly the hour of execution contrasted with different calculations. Here, littlest size dataset is 128 MB and the biggest size dataset is 934 MB. The execution time is decreased by utilization of the proposed alculation by 28% and 20% separately contrasted with standard K-means and K-implies MR in execution of littlest dataset. Though the decreases in execution times are 48% and 33% separately contrasted with standard K-means and K-implies MR in execution of the biggest dataset. Proposed calculation outflanks different calculations. It is unmistakably progressively invaluable to apply for the bigger datasets.

### A. Cluster Legitimacy

Group legitimacy is a term broadly alluded when evaluation of the consequences of a bunching calculation is performed. For estimating "goodness" of a bunching result, there are a few legitimacy lists which are applied.

S_Dbw is one of the well known legitimacy files. It has been proposed in [30]. Notwithstanding group smallness and division, S_Dbw thinks about thickness of the bunches too. Lower S_Dbw esteem demonstrates better bunching procedure. Applying S_Dbw varieties inside and between groups are estimated. Between groups fluctuation gauges the normal partition of bunches signified by condition 5.

$$Sep = \frac{1}{n_c} \sum_{i=1}^{n_c} \frac{\|\sigma(v_i)\|}{\|\cdots\|} \qquad (5)$$

Here nc = count of clusters
vi = centroid of the ith cluster
σ = variance
Within cluster density is defined by equation 6 as given below.

$$\qquad (6)$$

Here mij = midpoint of the separation among vi and vj group centroids.

As given in [30], the thickness work is characterized by the quantity of focuses in a hyper-circle whose range is equivalent to the normal standard deviation of groups. Definitely the normal standard deviation of groups inferred is as alluded beneath:

$$Density_{inter} = \frac{1}{n_c(n_c-1)} \sum_{i=1}^{n_c} \sum_{\substack{i=1 \\ i \neq j}} \frac{density\ (m_{ij})}{\{density\ (v_i), density\ (v_j)\}}$$

**Table 3: S_Dbw for distinct K**

| K | Proposed Work | Traditional K-means | K-means MR |
|---|---|---|---|
| 5 | 0.21513 | 0.37164 | 0.34412 |
| 10 | 0.22181 | 0.39089 | 0.32735 |
| 15 | 0.24472 | 0.43276 | 0.35195 |
| 20 | 0.21366 | 0.39595 | 0.33981 |
| 25 | 0.26735 | 0.48357 | 0.35422 |

Table 3 shows estimations of S_Dbw for proposed work, standard K-means and K-implies MapReduce calculations for number of bunches going from 5 to 25. It is seen that the work under reference has lower estimations of S_Dbw contrasted with the other two calculations displaying particularly great quality bunches age

## VIII. CONCLUSION

Right now, proposed upgrades in customary K-implies calculation to defeat its restrictions in grouping enormous information. The introduced work incorporates two calculations. Calculation 1 registers area of introductory seeds which are huge in number. Calculation 2 consolidations these seeds dependent on an edge Worth. We checked the exhibition of our work on UCI genuine datasets. The outcomes mirror that these calculations accomplish the reason. Exploratory outcomes additionally demonstrate that the proposed work exceeds expectations customary K-means and K-implies MapReduce calculations. Proposed calculation works skillfully on large datasets and is savvy on the system conveyed. Time decrease is of indispensable significance in any calculation. Any analysis which accomplishes extensive time decrease in calculation and conveys quality outcomes finds a legitimate spot in inquire about. Here time decrease accomplished is significant.

## REFERENCES

1. V. K. Jain and S. Kumar, "Big Data Analytic Using Cloud Computing," 2015 Second International Conference on Advances in Computing and Communication Engineering, Dehradun, 2015, pp. 667-672. doi: 10.1109/ICACCE.2015.112
2. A. K. Manekar and G. Pradeepini, "Cloud Based Big Data Analytics a Review," 2015 International Conference on Computational Intelligence and Communication Networks (CICN), Jabalpur, 2015, pp.785-788.doi: 10.1109/CICN.2015.160
3. D. Pandove and S. Goel, "A comprehensive study on clustering approaches for big data mining," 2015 2nd International Conference on Electronics and Communication Systems (ICECS), Coimbatore, 2015, pp.1333-1338. doi: 10.1109/ECS.2015.7124801
4. Jia Qiao and Y. Zhang, "Study on K-means method based on Data-Mining," 2015 Chinese Automation Congress (CAC), Wuhan, 2015, pp. 51-54. doi: 10.1109/CAC.2015.7382468
5. R.L.F. Cordeiro, C. Junior Traina, A.J.M. Traina, J. López, U.Kang and C.Faloutsos. (2011). "Clustering very large multidimensional datasets with MapReduce", In: Proceedings of KDD'11, ACM, California, August 21–24. 2011.
6. X. Cui, J.S. Charles and T. Potok. (2013). "GPU enhanced parallel computing for large scale data clustering". FutureGeneration Computer Systems, 29(7), 1736-1741, (2013).
7. G.Andrade, G. Ramos, D. Madeira, R. Sachetto, R. Ferreira and L. Rocha. (2013). "G-DBSCAN: A gpu accelerated algorithm for densitybased clustering". Procedia Computer Science, 18, 369-378
8. X. Cai, F. Nie and H. Huang. (2013). Proceedings of the Twenty-Third International Conference on Artificial Intelligence, Pages 2598-2604, Beijing, China, August 03-09, 2013, ISBN:978-1-57735-6332-2
9. U. Ghosia, U. Ahmad and M. Ahmad. (2013). "Improved K-Means Clustering Algorithm by Getting Initial Cenroids", World Applied Sciences Journal 27 (4): 543-551, 2013, ISSN 1818-4952, © IDOSI Publications, 2013, DOI: 10.5829/idosi.wasj.2013.27.04.1142.
10. T.M. Kodinariya1 and P.R. Makwana. (2013). "Review on determining number of Cluster in K-Means Clustering", International Journal of Advance Research in Computer Science and Management Studies, ISSN: 2321-7782 (Online), Volume 1, Issue 6, November 2013.
11. Y. Kim, K. Shim, M.S. Kim and J.S. Lee. "DBCURE-MR: an efficient density-based clustering algorithm for large data using MapReduce". Information Systems 42, 15-35 (2014).
12. X. Cui, P. Zhu and X. Yang. J Supercomput (2014) 70: 1249. https://doi.org/10.1007/s11227-014-1225-7
13. D. Gupta and K. Vani. (2015). "Using K-means cluster based techniques in external plagiarism detection," 2014 International Conference on Contemporary Computing and Informatics (IC3I), Mysore, 2014, pp. 1268-1273. doi: 10.1109/IC3I.2014.7019659
14. P. Bide and R Shedge, "Improved Document Clustering using k-means algorithm," 2015 IEEE International Conference on Electrical, Computer and Communication Technologies (ICECCT), Coimbatore, 2015, pp. 1-5.
15. C.W. Tsai, C.H. Hsieh and M.C. Chiang. (2015). "Parallel black hole clustering based on MapReduce", In: Proceedings of IEEE International Conference on Systems, Man and Cybernetics, 2015. DOI: 10.1109/SMC.2015.445
16. K. Wu, W. Zeng, T. Wu, and Y. An. (2015). "Research and improve on K-means based on hadoop". Software Engineering and Service Science (ICSESS). 2015 6th IEEE conference, 23-15 september, 2015. DOI: 10.1109/ICSESS.2015.7339068.
17. X. Cui and F. Wang. (2015), "An Improved Method for K-Means Clustering," 2015 International Conference on Computational Intelligence and Communication Networks (CICN), Jabalpur, 2015, pp. 756-759. doi: 10.1109/CICN.2015.154
18. J. Kaur and H. Singh, "Performance evaluation of a novel hybrid clustering algorithm using birch and K-means," 2015 Annual IEEE India Conference (INDICON), New Delhi, 2015, pp. 1-6. doi: 10.1109/INDICON.2015.7443414
19. N. Dhanachandra, K.M. Yambem and J. Chanu. "Image Segmentation Using K - means Clustering Algorithm and Subtractive Clustering Algorithm", Procedia Computer Science, Volume 54, 2015, Pages 764-771,https://doi.org/10.1016/j.procs.2015.06.090
20. P. Arora, Deepali and S. Varshney. "Analysis of K-Means and K-Medoids Algorithm For Big Data", Volume 78, 2016, Pages 507-512. https://doi.org/10.1016/j.procs.2016.02.095
21. A. Sinha and P.K. Jana. "A novel K-Means based clustering algorithm for big Data", IEEE International Conference on Advances in Computing, Communications and Informatics (ICACCI), 21-24Sept. 2016, Electronic ISBN: 978-1-5090-2029-4.
22. V.M. Ajin and L.D. Kumar. (2016). "Big data and clustering algorithms". 2016 International Conference on Research Advances in Integrated Navigation Systems (RAINS). DOI: 10.1109/RAINS.2016.7764405. 6-7 May 2016. IEEE
23. C. Shreedhar, N. Kasiviswanath and P.C. Reddy. "Clustering large datasets using K-means modified inter and intra clustering (KM-I2C) in Hadoop". Journal of Big Data (2017). DOI: 10.1186/s40537-017-0087-2. Springer..

24. H. Rehioui, A. Idrissi, Abourezq, M. and Zegrari, F. (2016). DENCLUE-IM: A New Approach for Big Data Clustering. Procedia Computer Science, Volume 83, 2016, pages 560-567, DOI: 10.1016/j.procs.2016.04.265. ELSEVIER.
25. G. Aletti and A. Micheletti. (2017). "A clustering algorithm formultivariate data streams with correlated components". Journal of Big Data 2017:48. DOI:10.1186/s40537-017-0109-0
26. T. Zhang and F. Ma Fumin. (2017). "Improved rough k-means clustering algorithm based on weighted distance measure with Gaussian function", InternationalJournal of Computer Mathematics, 94:4, 663-675,DOI: 10.1080/00207160.2015.1124099
27. V. Vijay, V.P. Raghunath, A. Singh and S. N. Omkar, "Variance Based Moving K-Means Algorithm," 2017IEEE 7th International Advance Computing Conference (IACC), Hyderabad, 2017, pp. 841-847. doi: 10.1109/IACC.2017.0173
28. A. Saini, J. Minocha, J. Ubriani and D. Sharma, "New approach for clustering of big data: DisK-means," 2016 International Conference on Computing, Communication and Automation (ICCCA), Noida, 2016,pp.122-126. doi: 10.1109/CCAA.2016.7813702
29. S. A. Ghamdi and G. D. Fatta, "Efficient Parallel K-Means on MapReduce Using Triangle Inequality," 2017 IEEE 15th Intl Conf on Dependable, Autonomic and Secure Computing, 15th Intl Conf on Pervasive Intelligence and Computing, 3rd Intl Conf on Big Data Intelligence and Computing and Cyber Science and Technology Congress(DASC/PiCom/DataCom/CyberSciTech), Orlando, FL, 2017, pp. 985-992. doi: 10.1109/DASC-PICom DataCom-CyberSciTec.2017.163Halkidi, M., Batistakis, Y. & Vazirgiannis, M. Journal of Intelligent Information Systems (2001) 17: 107.https://doi.org/10.1023/A:1012801612483
30. UCI Machine Learning Repository: https://archive.ics.uci.edu (2015)

## AUTHORS PROFILE

**Jagdish Kushwaha** is currently pursuing Ph.D. degree program in computer applications in BU Bhopal, ,+919630831212. E-mail: mca.jagdish@gmail.com

**Dr.(Prof) Shailesh Jaloree** ,Professor & HOD department of computer science & applied mathematics SATI Vidisha(MP),9424418356.shailesh_jaloree@rediffmail.com.

**Dr (Prof) R.S.THAKUR,** Professor Department of mathematics, Bioinformatics & Computer Application MANITBhopal (MP) 41996.ramthakur2000@gmail.com.

*Retrieval Number: F4126049620/2020©BEIESP*
*DOI: 10.35940/ijitee.F4126.049620*
*Journal Website: www.ijitee.org*

1937

*Published By:*
*Blue Eyes Intelligence Engineering & Sciences Publication*