Review the Clustering Algorithm in Big Data

Daw Thet Thet Khaing

Lecturer, University Of Computer Studies, Yangon

ABSTRACT

As today's organizations are capturing exponentially larger amounts of data than ever, now is the time for organizations to rethink how they digest that data. Through advanced algorithms and analytics techniques, organizations can harness this data, discover hidden patterns, and use the newly acquired knowledge to achieve competitive advantages. Big Data: Algorithms, Analytics, and Applications bridges the gap between the vastness of Big Data and the appropriate computational methods for scientific and social discovery. It covers fundamental issues about Big Data, including efficient algorithmic methods to process data, better analytical strategies to digest data, and representative applications in diverse fields, such as medicine, science, and engineering. Clustering is an essential data mining tool that plays an important role for analyzing big data. MapReduce is one of the most famous frameworks, and it has attracted great attention because of its flexibility, ease of programming, and fault tolerance. However, the framework has evident performance limitations, especially for iterative algorithms. We summarize these techniques, discuss their uniqueness and limitations, and explain how they address the challenging issues of iterative programs. We also perform an in-depth review to understand the problems and the solving techniques for parallel clustering algorithms. Hence, we believe that no well-rounded review provides a significant comparison among parallel clustering algorithms using MapReduce.

Keyword : *bigdata, mapreduce, cluster*

1.INTRODUCTION

In today's era of big data, we are dealing with a massive amount of data that are increasing significantly in different ways. For example, Yahoo! Web graph is reported to reach 1 billion nodes and 7 billion edges in 2002 [1]; social networking websites such as Twitter and Facebook data span several terabyte; Wikipedia and YouTube data are of similar size, which produces hundreds of gigabytes per minute [2]. To deal with this huge amount of data, clustering has become an essential tool that empowers data scientists for analyzing and discovering distribution patterns in a broad range of data. Accordingly, numerous data mining algorithms have been studied to deal with large-scale data sets. However, traditional data mining techniques are no longer able to capture the large amount of data. Hence, many researchers have shifted their focus to parallel clustering algorithms that would improve the bottleneck of traditional cluttering methods on a single machine [3]. In this way, parallel processing models have recently attracted considerable attention, which are used to process large-scale data sets. MapReduce has become the most popular parallel processing model that is used in many companies, including Facebook, mainly because of its simplicity [4]. Despite its advantages, MapReduce has to deal with numerous critical issues and challenges as explained in [5] such as wasting bandwidth, I/O, and CPU cycles for iterative programs, including PageRank [6], social network analysis, neural-network analysis, and clustering.

This survey aims to provide a valuable guidebook of problems and solving techniques in the context of iterative applications. Furthermore, we aim to provide a point of reference for future works that aim to deal with big data sets using the parallel clustering algorithms. To summarize, the main contributions of this work are as follows:

The most significant past literature related to MapReduce processing model was reviewed with an emphasis on iterative processing improvement techniques.

A significant number of research papers related to parallel clustering algorithms based on MapReduce were discussed.

The advantages and disadvantages of current works are discussed, which reveal the differentiating factors among the existing parallel algorithms using MapReduce.

Observations were made based on the problem and solutions to provide opportunities for future work.

The most relevant work to the first contribution of this paper is the recent survey of large-scale analytical query processing in MapReduce [7]. However, our work narrows down this paper to focus more on iterative processing techniques. In addition, our study provides an in-depth analysis of various clustering algorithms that can be executed in parallel using MapReduce. To the best of our knowledge, no relevant survey with an emphasis on parallel clustering algorithms based on MapReduce has been conducted.

The rest of this paper is organized as follows: Section 2 provides an overview of big data with focus on MapReduce/Hadoop. Section 3 presents the existing approaches to improve the performance of MapReduce based on iterative processing techniques and reviews current parallel clustering algorithms based on MapReduce. Section 4 analyzes the current parallel clustering algorithms. Section 5 concludes this work.

2. ISSUE IN BIGDATA

Important issues have been reviewed and discussed in this section. We can describe the characteristics of big data using three Vs, also issues in big data.Big data requires a revolutionary step forward from traditional data analysis, characterized by its three main components through which it is formed: variety, velocity and volume as shown in Figure 1[3, 8, 13, and 17].



Fig -1 the three Vs of Big data

1) Variety: Variety makes big data really big. Big data comes from a great variety of sources and generally has in three types: structured, semi structured and unstructured.Structured data inserts a data warehouse already tagged and easily sorted but unstructured data is random and difficult to analyze. Semi-structured data does not conform to fixed fields but contains tags to separate data elements [4, 17].

2) Volume: Volume or the size of data now is larger than terabytes and petabytes. The grand scale and rise of data outstrips traditional store and analysis techniques [4, 16].

3) Velocity : Velocity is required not only for big data, but also all processes. For time limited processes, big data should be used as it streams into the organization in order

to maximize its value [4,16].

3. CLUSTERING ALGORITHMS

This paper presents various clustering algorithms with by considering the properties of Big Data characteristics such as size, noise, dimensionality, computations of algorithms, shape of cluster, etc [10] [11]. The overview of clustering algorithms is depicted in Table 1.

	No	Big Data Clustering Technique		
		Algorithm Type	Algorithm	
	1	Partition Based	K-means	
			K-medoids	
			K-modes	
			РАМ	
			CLARA	
			CLARANS	
			FCM	
	2	Hierarchical	BIRCH	
		Based	CURE	
		ALC: NOT THE OWNER	ROCK	
		di la constante di la constant	Chameleon	
			ECHIDNA	
R	a and a second		Wards	100
	1	10	SNN	
	1.10		CACTUS	
	6		GRIDCLUST	
	3	Density Based	DSCAN	
	2		OPTICS	
			DBCLASD	
			GDBSCAN	
			DENCL	
			SUBCLU	
	4	Grid Based	STING	
			Wave Cluster	
		1.0	BANG	
		1000	CLIOUE	
			OptiGrid	
			MAFIA	
			ENCLUS	
			PROCLUS	
		10 MI 100 M	ORCLUS	
			FC	
			STIRR	
	5	Model Based	EM	
		Justa Bused	COBWEB	
			CLASSIT	
	100		SOM	
			SLINK	

Table 1. An Overview of Algorithms for |BigData Mining

2.1 Partition Based Clustering Algorithm

All objects are considered initially as a single cluster. The objects are divided into no of partitions by iteratively locating the points between the partitions. The partitioning algorithms like K-means, K-medoids (PAM, CLARA, CLARANS, and FCM) and K-modes. Partition based algorithms can found clusters of Non convex shapes.

1)FCM - Fuzzy CMEANS algorithm: [9] the algorithm is based on the K-means concept to partition dataset into Clusters.

The algorithm is as follows:

- Calculate the cluster centroids and the objective value and initialize fuzzy matrix.
- Computer the membership values stored in the matrix.

The paper presents list of all algorithms and their efficiency based on the input parameter to mine the Big Data as described below:

- If the value of objective is between consecutive iterations is less than the stopping condition then stop.
- This process is continuous until a partition matrix and clusters are formed.

2.2 Hierarchical Clustering Algorithms

There are two approaches to perform Hierarchical clustering techniques Agglomerative (top-bottom) and Divisive (bottom- top). In Agglomerative approach, initially one object is selected and successively merges the neighbor objects based on the distance as minimum, maximum and average. The process is continuous until a desired cluster is formed. The Divisive approach deals with set of objects as single cluster and divides the cluster into further clusters until desired no of clusters are formed. BIRCH, CURE, ROCK, Chameleon, Echidna, Wards, SNN, GRIDCLUST, CACTUS are some of Hierarchical clustering algorithms in which clusters of Non convex, Arbitrary Hyper rectangular are formed.

1) BIRCH - Balanced Iterative Reducing and Clustering using Hierarchies: It is an agglomerative hierarchical algorithm which uses a Clustering Feature (CF-Tree) and incrementally adjusts the quality of sub clusters.

The algorithm is as follows:

- Load data into memory: CF Tree is constructed with one scan of the data. Subsequent phases become fast, accurate and less order sensitive.
- Condense data: Rebuilt the CF tree with larger T.
- Global Clustering: Use the existing clustering algorithm on CF leaves.
- Cluster refining-Do additional passes over the dataset and reassign data points to the closest centroids from above step.
- The process continuous until to form k no of clusters.

2) *CURE- Clustering Using REpresentatives*: A hierarchy of Divisive approach is used and it selects well scattered points from the cluster and then shrinks towards the center of the cluster by a specified function. Adjacent clusters are merged successively until the no of clusters reduces to desired no of clusters.

The algorithm is as follows:

- Initially all points are in separate clusters, each cluster is defined by the point in the cluster.
- The Representative points of a cluster are generated by first selecting well scattered objects for the cluster and then perform shrinking or moving towards the cluster by a specified factor.
- At each step of the algorithm, two clusters with closest pair of representative point are chosen and merged together to form cluster.

3) *ROCK - Robust Clustering algorithm for Categorical attributes*: It is a hierarchical clustering algorithm in which to form clusters it uses a link strategy. From bottom to top links are merging together to form a cluster.

The algorithm is as follows:

Initially consider set of points in which every point is a cluster and compute the links between each pair of points.

Build a heap and maintain heap for each cluster.

A goodness measure based on the criterion function will be calculated between pairs of clusters.

Merge the clusters which have maximum value of criteria function.

4) Chameleon - : It is an agglomerative hierarchical clustering algorithm of dynamic modeling which deals with two phase approach of clustering

The algorithm is as follows:

A two phase approach of partition and merge is used to form a cluster.

• During Partition phase-Initially consider all data points as a single cluster.

- Using a graph partitioning algorithm divide the cluster into a relatively large no of small clusters using hMETIS method.
- The process terminates when a large sub cluster contains slightly more than a specified no of vertices.
- In merge phase using agglomerative hierarchical approach select pairs of clusters whose inter connectivity and relative closeness are reaches the threshold value.
- Merge the clusters which are having the highest inter connectivity and closeness.

The algorithm is repeated until none of the adjacent clusters satisfy the two conditions.

5) ECHIDNA: It is an agglomerative hierarchical approach for clustering the network traffic data.

The steps of algorithm are given below:

- The input data is extracted from network traffic consists of a 6 Tuple value of numerical and categorical attributes.
- Each record iteratively builds a hierarchical tree of clusters called CF-Tree.
- Insert each record into the closest cluster using a combined distance function for all attributes into CF-Tree.
- The radius of a cluster determines if a record should be absorbed into the cluster or if the cluster should be split.
- Once the cluster is created and all the significant nodes are to form a Cluster Tree.

The Cluster Tree is further compressed to create a concise and meaningful report.

6) SNN - Shared Nearest Neighbors : A hierarchy of top to bottom approach is used for grouping the objects.

The steps of algorithm are given below:

- A proximity matrix should be maintained for the distances of set of points.
- Objects are clustered together based on the nearest neighbor and the object with maximum distance can be avoided.

7) CACTUS – Clustering CaTegorical Data Using Summaries: It is a very fast and scalable algorithm for finding the clusters. A hierarchy structure is used to generate maximum segments or clusters. A two step procedure deals with the description of algorithm as follows:

- Attributes are strongly connected if the data points are having larger frequency.
- Clusters are formed based on the co-occurrences of attribute value pairs.
- A cluster is formed if any segment is having no of elements α times greater than elements of other.

8) *GRIDCLUST - GRID based hierarchical CLUSTering algorithm :* A clustering algorithm of hierarchical method based on grid structure.

The algorithm is as follows:

- Initially partition the data set into data space to form grid structure and the topological distributions are maintained.
- Once data is assigned to the blocks of cells or grids density values are calculated and sorted according to their values.
- The largest dense block was considered as cluster center.
- Using the Neighbor search algorithm a cluster can be formed with the remaining blocks.

A. Density Based Clustering Algorithms

Data objects are categorized into core points, border points and noise points. All the core points are connected together based on the densities to form cluster. Arbitrary shaped clusters are formed by various clustering algorithms such as DBSCAN, OPTICS, DBCLASD, GDBSCAN, DENCLU and SUBCLU.

1) DBSCAN – Density Based SCAN clustering algorithm :

It is a connectivity based algorithm which consists of 3 points namely core, border and noise.

The algorithm is as follows:

- Set of points to be considered to form a graph.
- Create an edge from each point c to the other point in the neighborhood of c.
- If set of nodes N not contain any core points then terminate N.
- Select a node X that must be reached form c.

Repeat the procedure until all core points forms a cluster.

2) *OPTICS – Ordering Points To Identify the Clustering Structure* : It is also a connectivity based density algorithm. OPTICS is an extension of DBSCAN algorithm which is also based on the same parameters as DBSCAN algorithm. The run time of OPTICS is 1.6 times greater than DBSCAN algorithm.

The algorithm is as follows:

- Among the set of points select a point is a core point if at least Minpts are found in the core distance.
- For each point c create an edge from c to other point with a core distance of c.
- Select set of nodes which contain core points as a cluster that reaches from c.

3) DBCLASD – Distribution Based Clustering of Large Spatial Databases : It is Connectivity based and application based clustering algorithm for mining of large spatial data bases.

The algorithm is as follows:

- Construct set of candidates C based on the query.
- The point will be remains within the cluster if the distance between set of C has expected distribution.
- Otherwise the point will be considered as unsuccessful candidate.
- The process is continuous until all points with expected distribution form cluster.

4) GDBSCAN – Generalized Density Based Spatial Clustering of ApplicatioN : A connectivity based density algorithm in which it form clusters with point objects and as well as spatial attributes.

The algorithm is as follows:

- An attribute object P is selected and retrieves all objects densities whether they are reachable from P with respect to neighborhood of the object (NPred) and minimum weighted cardinality (Min weight).
- If P is a core object this procedure yields a density connected set Ci with respect to NPred and Min weight.
- Otherwise it does not belong to any density connected set Ci.

This procedure is iteratively applied to each object P which has not yet been classified.

5) DENCLUE – DENsity based CLUstEring : Among all algorithms of density based clustering approach DENCLUE is the algorithm which is based on the density function. Arbitrary shape of good quality of clusters can be formed with large amount of data set.

The algorithm is as follows:

- Consider the data set in the grid structure and find the high density cells based on mean value (highest).
- If d (mean (c1), mean (c2)) < 4a then connectc1, c2.
- Find the density attractors using Hill-Climbing approach and they should be local maxima of overall density function.
- Merge the attractors and they can be identified as clusters.

6) SUBCLU – SUBspace CLUstering : It is an efficient approach to the subspace clustering and which is based on the formal clustering notion. It can detect clusters of arbitrary shape.

The algorithm is as follows:

• Initially generate all 1-D subspace clusters in which at least one cluster in the subspace found.

- Generate (k+1) dimensional candidate subspaces. Test candidates and generate (k+1) dimensional clusters.
- All the clusters in the higher dimensional subspace will be the subsets of clusters which are detected in the first clustering.

The process continues until (k+1)–D clusters are formed from k-D clusters.

B. Grid Based Clustering Algorithms

Data objects are categorized into core points, border points and noise points. All the core points are connected together based on the densities to form cluster. Arbitrary shaped clusters are formed by various clustering algorithms such as DBSCAN, OPTICS, DBCLASD, GDBSCAN, DENCLU and SUBCLU.

1) STING – STatistical Information Grid based method : It is similar to BIRCH hierarchical algorithm to form a cluster with spatial data bases.

The algorithm is as follows:

- Initially the spatial data stored into rectangular cells using a hierarchical grid structure.
- Partition each cell into 4 child cells at the next level with each child corresponding to a quadrant of the • parent cell.
- Calculate probability of each cell whether it is relevant or not. If the cell is relevant then apply same calculations on each cell one by one.
- Find the regions of relevant cells in order to form cluster. •

2) Wave Cluster - Among all the clustering algorithms, this is based on signal processing : The algorithm works with numerical attributes and has multi-resolution. Outliers can be detected easily.

The algorithm is as follows:

- Fit all the data points into a cell. Apply wavelet transform to filter the data points. •
- Apply discrete Wavelet transform to accumulate data points. •
- High amplitude signals are applied to the corresponding cluster interiors and high frequency is applied to • find boundary of cluster.
- Signals are applied to the attribute space in order to form cluster with more sharp and eliminates outliers • easily.

3) BANG - Grid based clustering algorithm : It is an extension of GRIDCLUST algorithm which initially considers all data points as blocks but it uses BANG structure to maintain blocks.

The algorithm is as follows:

- Divide the feature space into rectangular blocks which contains up to a maximum of P max data points. •
- Build a binary tree to maintain the density indices of all blocks are calculated and sorted in decreasing or-• der.
- Starting with the highest density index, all neighbor blocks are determined and classified in decreasing order to form a cluster.
- The process is repeated for the remaining blocks.

4) CLIOUE – CLustering In OUEst : A subspace clustering algorithm for numerical attributes in which bottom top approach is used to form clusters.

The algorithm is as follows:

- Consider set of data points, at one pass apply equal width to the set of points to form grid cells. •
- Let the rectangular cells into subspace whose density exceed τ are placed into equal grids.
- The process is continuous recursively to form (q-1) dimensional units to q dimensional units.
- The subspaces are connected to each other to form cluster with equal width.

5) *OPTI GRID* – Optimal Grid : *The algorithm is designed to cluster large spatial data bases*. The algorithm is as follows:

- Define the data set with best cutting hyper planes through a set of selected projections.
- Select best local optima cutting plane.
- Insert all the cutting planes with a score greater than or equal to minimum cut score into a BEST CUT.
- Select q cutting planes of the highest score form BEST CUT and construct a multi dimensional grid G using the q cutting planes.
- Insert all data points in D into G and determine the highly populated grid cells in G and form a set of clusters C.

6) *MAFIA – Merging of Adaptive Finite IntervAls :* It is descendant of CLIQUE algorithm in which instead of using a fixed size cell grid structure with an equal number of bins in each dimension, it constructs an adaptive grid to improve the quality of clustering.

The algorithm is as follows:

- In a single pass an adaptive grid structure was constructed by considering set of all points.
- Compute the histogram by reading blocks of data into memory using bins.
- Bins are grouped together based on the dominance factor α.
- Select the bins that are α times dense greater than average as p candidate dense units (CDU).
- Recursively the process continuous to form new p-CDU's and merge adjacent CDU's into clusters.

7) ENCLUS – Entropy based CLUStering : The algorithm is entropy based algorithm for clustering large data sets. ENCLUS is an adaptation of CLIQUE algorithm.

The algorithm is as follows:

• The objects whose subspaces are spanned by attribute A1....AP with an entropy criteria H (A1....AP) < ω (a threshold) are selected for clustering.

8) *PROCLUS – PROjected CLUStering algorithm* : The algorithm also uses medoids which is same as K – medoids clustering criteria.

The algorithm is defined in three step procedure as follows:

- Initialization: Consider the set of all points and select data points randomly.
- Iteration phase: select medoids of the clusters as data point and define a subspace to each medoids.
- Refinement phase: select best medoids form set of medoids which has all dimensions. Select another medoids which is nearest to best medoids.

All the data points within this distance will be formed as a cluster.

9) ORCLUS- ORiented projected CLUStering generation algorithm : It is similar to PROCLUS clustering algorithm but it focuses on non-axis parallel subspace.

The algorithm is defined by three strategies - assignment, subspace determination and merge as follows:

- Assignment: During this phase the algorithm iteratively assigns all the data points to the nearest cluster centers.
- Sub space determination: To determine sub space calculate co-variance matrix for each cluster and Eigen vectors with the least Eigen values.
- Merge: Clusters which are near to each other and have similar directions are merged.

10) FC – Fractal Clustering algorithm : The algorithm deals with hierarchy approach works with several layers of grids for numeric attributes and identifies clusters of irregular shapes.

The algorithm is as follows:

- Start with a data sample and a threshold value is considered for a given set of points.
- Initialize threshold value, scan full data incrementally.

- Using HFD-Hausdorff Fractal Dimension (HFD) method adds an incoming point to each cluster.
- If the smallest increase exceeds a threshold τ value, a point is declared an outlier and shape of the cluster is declared as irregular.
- Otherwise a point is assigned to cluster.

11) STIRR – Sieving Through Iterated ReinfiRcement : This algorithm deals with spectral partitioning using dynamic system as follows:

- Set of attributes are considered and weights W= W v are assigned to each attribute.
- Weights are assigned to set of attributes using combining operator ϕ defined as
- ϕ (W1...Wn-1) = W1+.....+Wn-1.
- At a particular point the process is stopped to achieve dynamic system.

C. Model Based Clustering Algorithms

Set of data points are connected together based on various strategies like statistical methods, conceptual methods, and robust clustering methods. There are two approaches for model based algorithms one is neural network approach and another one is statistical approach. Algorithms such as EM, COBWEB, CLASSIT, SOM, and SLINK are well known Model based clustering algorithms.

1) *EM* – *Expectation and Maximization* : This algorithm is based on two parameters- expectation (E) and maximization (M).

- E: The current model parameter values are used to evaluate the posterior distribution of the latent variables. Then the objects are fractionally assigned to each cluster based on this posterior distribution as Q(θ, θT) = E[log p(x g, x m | θ) x g, θT]
- M: The fractional assignment is given by re-estimating the model parameters with the maximum likelihood rule as

$\theta t + 1 = \max Q (\theta, \theta T)$

The process is repeated until the convergence condition is satisfied.

2) COBWEB – Model based clustering algorithm : It is an Incremental clustering algorithm, which builds taxonomy of clusters without having a predefined number of clusters. The clusters are represented probabilistically by conditional probability P(A = v | C) with which attribute A has value v, given that the instance belongs to class C.

The algorithm is as follows:

- The algorithm starts with an empty root node.
- Instances are added one by one.
- For each instance, the following options (operators) are considered:
- -classifying the instance into an existing class;
- -creating a new class and place the instance into it.
- -combining two classes into a single class (merging) and placing the new instance in the resulting hierarchy;
- -split the class into two classes (splitting) and placing the new instance in the resulting hierarchy.
- The algorithm searches the space of possible hierarchies by applying the above operators and an evaluation function based on the category utility.

3) SOM- Self Organized Map algorithm : A Model based clustering incremental clustering algorithm, which is based on the grid structure.

The algorithm is defined by a two step process:

- Place the grid of nodes along a plane where data points are distributed.
- Sample the data point and subject the closest node and neighboring node to its influence. Sampling another point and so on.
- The procedure is repeated until all data points have been sampled several times.
- Each cluster is defined with reference to a node specifically comprised by those data points for which it represents the closest node.

4) *SLINK – Single LINK clustering algorithm* : A Model based clustering algorithm in which a hierarchy approach is used to form clusters.

- Starts with set of points, let each point as a singleton cluster.
- Using Euclidean distance determine the distance between the two points.
- Merge the links between all points' shortest links first.
- Combine the single links to form a cluster.

Ν	Clustering Algorithms				
0	Algorithm Name	Data Size	Dataset Type	Time Complexit y	
1	K-means[12]	Large	Numeri cal	O(n k d)	
2	K- medoid [14]	Small	Categor ical	O(n ² dt)	
3	k-modes[13]	Large	Categor ical	O(n)	
4	PAM[15]	Small	Numeri cal	$O(k(n-k)^2)$	
5	CLARA[16]	Large	Numeri cal	$O(k(40+k)^{2}+k(n-k))$	
6	CLARANS[17]	Large	Numeri cal	O(kn2)	
7	FCM[9]	Large	Numeri cal	O(n)	
8	BIRCH[19]	Large	Numeri cal	O(n)	
9	CURE[20]	Large	Numeri cal and Categor ical	O(n2logn)	
1 0	ROCK[21]	Large	Numeri cal and Categor	O(n2+nm m- ma+n2logn	

Table 2 Various Clustering for BigDAta

				ical)
	1 1	Chameleon[22]	Large	All types data	O(n ²)
	1 2	ECHIDNA[23]	Large	Multi- variant	$O(N*B(1+\log_B m))$
	1 3	Wards[39]	Small	Numeri cal	no
	1 4	SNN [41]	Small	Categor ical	O(n ²)
1	1 5	CACTUS[45]	Small	Categor ical	O(c N)
1	1 6	GRIDCLUST[46]	Small	Numeri cal	O(n)
	1 7	DBSCAN[24]	Large	Numeri cal	O(n log n)for spatial data
	1 8	OPTICS[25]	Large	Numeri cal	O(n log n)
	1 9	DBCLASD[26]	Large	Numeri cal	O(3n ²)
	2 0	GDBSCAN[43]	Large	Numeri cal	no
	2 1	DENCLUE[27]	Large	Numeri cal	O(log D)
100	2 2	SUBCLU[42]	Large	Numeri cal	no
	2 3	STING[29]	Large	Spatial	O(k)
	2 4	Wave Cluster[28]	Large	Numeri cal	O(n)
	2 5	BANG[18]	Large	Numeri cal	O(n)
	2 6	CLIQUE[30]	Large	Numeri cal	O(C k + m k)
	2 7	OptiGrid[31]	Large	Spatial	O(n d) to

					O(ndlogn)	
	2 8	MAFIA[44]	Large	Numeri cal	$O(cp + p^{N})$	
	2	ENCLUS[36]	Large	Numeri	O(ND+ m	
	3	PROCLUS[37]	Large	Spatial	O(n)	
	0 3 1	ORCLUS[38]	Large	Spatial	O(d ³)	
ĺ	3 2	FC[11]	Large	Numeri cal	O(n)	
/	3 3	STIRR[11]	Large	Categor ical	O(n)	
6	3 4	EM[32]	Large	Spatial	O(knp)	
	3 5	COBWEB[33]	Small	Numeri cal	O(n ²)	
	3 6	CLASSIT[34]	Small	Numeri cal	O(n ²)	
3	3 7	SOM's[35]	Small	variant	O(n ² m)	
	3 8	SLINK[40]	Large	Numeri cal	O(n ²)	

4. CONCLUSIONS

This paper analyzed different clustering algorithms required for processing Big Data. The study revealed that to identify the outliers in large data sets, the algorithms that should be used are BIRCH, CLIQUE, and ORCLUS. To perform clustering, various algorithms can be used but to get appropriate results the present study suggests that – by using CURE and ROCK algorithms on categorical data, arbitrary shaped clusters will be created. By using COBWEB and CLASSIT algorithms on numerical data with model based, non-convex shape clusters can be formed. For spatial data STING, OPTIGRID, PROCLUS and ORCLUS algorithms when applied yield arbitrary shaped clusters.

REFERENCES

- [1] Aggarwal C, Zhai C. A survey of text clustering algorithms. Mining Text Data. New York, NY, USA. Springer-Verlag: 2012. p. 77–128.
- [2] Brank J, Grobelnik M, Mladenic D. A survey of ontology evaluation techniques. Proceedings Conf Data Mining and Data Warehouses (SiKDD); 2005. p. 166–9.

- [3] B.Gerhardt, K. Griffin and R. Klemann, "Unlocking Value in the Fragmented World of Big Data Analytics", Cisco Internet Business Solutions Group, June 2012, http://www.cisco.com/web/about/ac79/docs/sp/InformationInfomediaries.pdf.
- [4] C. Eaton, D. Deroos, T. Deutsch, G. Lapis and P.C. Zikopoulos, Understanding Big Data: Analytics for Enterprise Class Hadoop and Streaming Data, Mc Graw-Hill Companies, 978-0-07-179053-6,2012.
- [5]
- [6] Xu R, Wunsch D. Survey of clustering algorithms. IEEE Transactions on Neural Networks. 2005 May; 16(3):645–78.
- [7] J. Dean and S. Ghemawat. Mapreduce: Simplified data processing on large clusters. OSDI '04, pages 137–150, 2008.
- [8] Intel IT Center, "Planning Guide: Getting Started with Hadoop", Steps IT Managers Can Take to Move Forward with Big Data Analytics, June 2012.
- [9] <u>http://www.intel.com/content/dam/www/public/us/en/documents/g</u> uides/getting-started-with-hadoop-planning-guide.pdf.
- [10] Fahad A, Alshatri N, Tari Z, Alamri A. A survey of clustering algorithms for Big Data: Taxonomy and empirical analysis. IEEE Transactions on Emerging Topics in Computing. 2014 Sep; 2(3):267–79
- [11] Berkhin P. Survey of clustering data mining techniques in grouping multidimensional data. Springer. 2006; 25–71.
- [12] Yadav C, Wang S, Kumar M. Algorithms and approaches to handle large data sets A survey. International Journal of Computer Science and Network. 2013; 2(3):1–5.
- [13] R.D. Schneider, Hadoop for Dummies Special Edition, John Wiley&Sons Canada, 978-1-118-25051-8, 2012.
- [14] Bezdek JC, Ehrlich R, Full W. FCM: The Fuzzy C-Means Clustering algorithm. Computers and Geosciences. 1984; 10(2-3):191–203.
- [15] Fahad A, Alshatri N, Tari Z, Alamri A. A survey of clustering algorithms for Big Data: Taxonomy and empirical analysis. IEEE Transactions on Emerging Topics in Computing. 2014 Sep; 2(3):267–79.
- [16] S. Madden, "From Databases to Big Data", IEEE Internet Computing, v. 16, pp. 4-6, June 2012.
- [17] Berkhin P. Survey of clustering data mining techniques in grouping multidimensional data. Springer. 2006; 25–71.
- [18] Macqueen J. Some methods for classification and analysis of multivariate observations. Proceedings 5th Berkeley Symposium on Mathematical Statistics Probability; Berkeley, CA, USA. 1967. p. 281–97.
- [19] Huang Z. A fast clustering algorithm to cluster very large categorical data sets in data mining. Proceedings SIGMOD Workshop Res Issues Data Mining Knowl Discovery; 1997. p. 1–8.
- [20] Park HS, Jun CH. A simple and fast algorithm for K-medoids clustering. Expert Systems Applications. 2009 Mar; 36(2.2):3336–41.
- [21] Ng RT, Han J. Efficient and effective clustering methods for spatial data mining. Proceedings Int Conf Very Large Data Bases (VLDB); 1994. p. 144–55.
 [22] Kaufman L, Rousseau PJ. Finding groups in data: An introduction to cluster analysis. USA, Johns and Sons
- [22] Kaufman L, Rousseau PJ. Finding groups in data: An introduction to cluster analysis. USA, Johns and Sons Wiley; 2008.
- [23] Ng RT, Han J. CLARANS: A method for clustering objects for spatial data mining. IEEE Transactions on Knowledge Data Engineering (TKDE). 2002 Sep/Oct; 14(5):1003–16.
- [24] Schikuta E, Erchart M. The BANG Clustering system: Grid– based data analysis. Lecture Notes in Computer Science. 1997; 1280:513–24.
- [25] Zhang T, Ramakrishna R, Livny M. BIRCH: An efficient data clustering method for very large databases. Proceedings of the ACM SIGMOD International Conference on Management of Data. 1996 Jun; 25(2):103– 14.
- [26] Guha S, Rastogi R, Shim K. Cure: An efficient clustering algorithm for large data bases. Proceedings of the ACM SICMOID international Conference on Management of Data. 1998 Jun; 27(2):73–84.
- [27] Guha S, Rastogi R, Shim K. Rock: A robust clustering algorithm for categorical attributes. 15th International Conference on Data Engineering; 1999. p. 512–21.
- [28] Karypis G, Han EH, Kumar V. Chameleon: Hierarchical clustering using dynamic modeling. IEEE Computer. 1999 Aug; 32(8): 68–75.
- [29] Mahmood AN, Leckie C, Udaya P. An efficient clustering scheme to exploit hierarchical data in network traffic analysis. IEEE Transactions on Knowledge. Data Engineering. 2008 Jun; 20(6):752–67.
- [30] Ester M, Kriegel HP, Sander J, Xu X. A density-based algorithm for discovering clusters in large spatial databases with noise. Proceedings ACM SIGKDD Conf Knowl Discovery Ad Data Mining (KDD); 1996. pp. 226–31.

- [31] Ankerst M, Breunig M, Kriegel HP, Sander J. Optics: Ordering points to identify the clustering structure. Proceedings of the ACM SIGMOD International Conference on Management of Data. 1999 Jun; 28(2):49–60.
- [32] Xu X, Ester M, Krieger HP, Sander J. A distribution-based clustering algorithm for mining in large spatial databases. Proceedings 14th IEEE International Conference on Data Engineering (ICDE); Orlando, FL. 1998 Feb 23-27. p. 324–31.
- [33] Hinneburg A, Keim DA. An efficient approach to clustering in large multimedia databases with noise. Proceedings ACM SIGKDD Conf Knowl Discovery Ad Data Mining (KDD); 1998. p. 58–65.
- [34] Sheikholeslami G, Chatterjee S, Zhang A. Wave cluster: A multi resolution clustering approach for very large spatial databases. Proceedings Int Conf Very Large Data Bases (VLDB); 1998. p. 428–39.
- [35] Wang W, Yang J, Muntz R. Sting: A statistical information grid approach to spatial data mining. Proceedings 23rd Int Conf Very Large Data Bases (VLDB); 1997. p. 186–95.
- [36] Jain AK, Dubes RC. Algorithms for Clustering Data. Upper Saddle River, NJ, USA, Prentice-Hall; 1988.
- [37] Hinneburg A, Keim DA. Optimal grid-clustering: Towards breaking the curse of dimensionality in highdimensional clustering. Proceedings 25th Int Conf Very Large Data Bases (VLDB); 1999. p. 506–17.
- [38] Dempster AP, Laird NM, Rubin DB. Maximum likelihood from incomplete data via thee algorithm. Journal of the Royal Statistical Society. 1977; 39(1):1–38.
- [39] Fisher DH. Knowledge acquisition via incremental conceptual clustering. Machine Learning. 1987 Sep; 2(2):139–72.
- [40] Gennari JH, Langley P, Fisher D. Models of incremental concept formation. Artificial Intelligence. 1989 Sep; 40(1-3):11-61.
- [41] Kohonen T. The self-organizing map. Neurocomputing. 1998 Nov; 21(1-3):1-6.
- [42] Cheng CH, Fu AW, Zhang Y. Entropy based sub space clustering for mining numerical data. Proceedings of the fifth ACM SIGMOID International Conference on Knowledge discovery and Data Mining; 1999. p. 84–93.
- [43] Milenova BL, Campos M. Clustering large databases with numeric and nominal values using orthogonal projections. O Cluster; 2006. p. 1–11.
- [44] Aggarwal CC, Yu PS. Finding generalized projected clusters in high dimensional spaces. Proceedings of the 2000 ACM SIGMOID International Conference on Management of Data. 2000 Jun; 29(2):70–81.
- [45] Xu R, Wunsch D. Survey of clustering algorithms. IEEE Transactions on Neural Networks. 2005 May; 16(3):645–78.
- [46] Han J, Kamber M. Data Mining: Concepts and Techniques. 12 A Survey on Clustering Techniques for Big Data Mining
- [47] 2nd edition. San Mateo, CA, USA, Morgan Kaufmann; 2006. 41. Hubert L, Arabie P. Comparing partitions. Journal of Classification. 1985 Dec; 2(1):193–218.
- [48] Kailing K, Kriegel HP, Kroger P. Density-connected subspace clustering for high- dimensionality data. Proceedings of the 2004 SIAM International Conference on Data Mining; 2010. p. 246–57.
- [49] Varghese BM, Unnikrishanan A, Paulose Jacob K. Spatial clustering algorithms An overview. Asian Journal of Computer Science and Information Technology. 2014; 3(1):1–8.
- [50] Cheng W, Wang W, Batista S. Grid Based Clustering. 2009. p. 12–24.
- [51] Ganti V, Gehrke J, Ramakrishna R. CACTUS- Clustering Categorical Data Using Summaries. Proceeding of the fifth ACM SIGMOID International Conference on Knowledge Discovery and Datamining; 1999. p. 73–83.
- [52] Cao Q, Bouqata B, Mackenzie PD, Messiar D, Salvo J. A grid-based clustering method for mining frequent trips from large-scale, event-based telemetries datasets. The 2009 IEEE International Conference on Systems, Man and Cybernetics; San Anonio, TX, USA. 2009 Oct. p. 2996–3001.