

# HANDLING BIGDATA USING MAPREDUCE CLUSTERING AND MACHINE INTELLIGENCE TECHNIQUES

<sup>1</sup>Shyam Mohan J S, <sup>2</sup>P.Shanmugapriya

<sup>1</sup>Assistant Professor, <sup>2</sup>Associate Professor,  
Department of CSE,

SCSVMV University, Enathur , Kanchipuram , Tamilnadu – 631 561. India

## **Abstract :**

Cluster Analysis or Bunch examination is one the most flexible strategies in measurable science. It is especially utilized in investigating datasets of high multifaceted nature like sub-atomic science or spatial information. Huge Data, Machine Learning, Computer Vision and Computational Biology are other run of the mill fields where grouping is fundamentally engaged. Conventional calculations neglect to deal with tremendous and high dimensional information as the datasets are of high volume, high speed and diverse assortments. Viable bunching calculations give advantages to some continuous logical utilizations of huge high dimensional datasets. Greater parts of the Companies have begun chipping away at Hadoop MapReduce Algorithms for bunching information. Grouping procedures are connected on various datasets. Huge Data is prevalent for handling, putting away and overseeing tremendous volumes of information. Bunching of such gigantic and complex datasets has turned into a testing task in the zone of huge information examination. Customary grouping calculations are not adaptable for overseeing substantial datasets. For little datasets, K-Means calculation is most appropriate for discovering likenesses between elements dependent on separation measures. For colossal and complex datasets, machine insight calculations are executed on hadoop MapReduce to shape bunches. MapReduce-Machine Intelligence Clustering (MMC) are structured and actualized in Hadoop and Amazon Elastic MapReduce (EMR) for various datasets (colossal and high dimensional) to create groups with most extreme intra-bunch and least between bunch separations. The after effects of the MMC grouping calculations show critical upgrades in execution time contrasted and customary bunching calculations and it is both viable and proficient.

**Keywords - MapReduce-Machine Intelligence Clustering (MMC) , Hadoop.**

## **I. INTRODUCTION**

Preparing colossal datasets is a troublesome undertaking that includes the utilization of many complex apparatuses. For powerful preparing of Big Data, Massively Parallel Processing (MPP) and MapReduce are utilized. Routinely, MapReduce functionalities are given by NoSQL frameworks and information can be duplicated from NoSQL frameworks into scientific frameworks, for example, Hadoop for MapReduce.

Operational frameworks NoSQL Database is joined with Hadoop (MongoDB with Hadoop). Association is set up by utilizing API's. By and large, activities performed by operational frameworks and scientific advances are delegated:

NoSQL – Provide Big Data and give insight to the clients.

MPP and Hadoop – Performing bits of knowledge utilizing examination.

Group examination is one the most adaptable strategies in measurable science. Group investigation is especially utilized in investigating datasets of high multifaceted nature like sub-atomic science or spatial information. Enormous Data, Machine Learning, Computer Vision and Computational Biology are other run of the mill fields where grouping is principally engaged. Bunching gives outline of datasets by gathering of comparative articles that encourages the translation for further investigation of information.

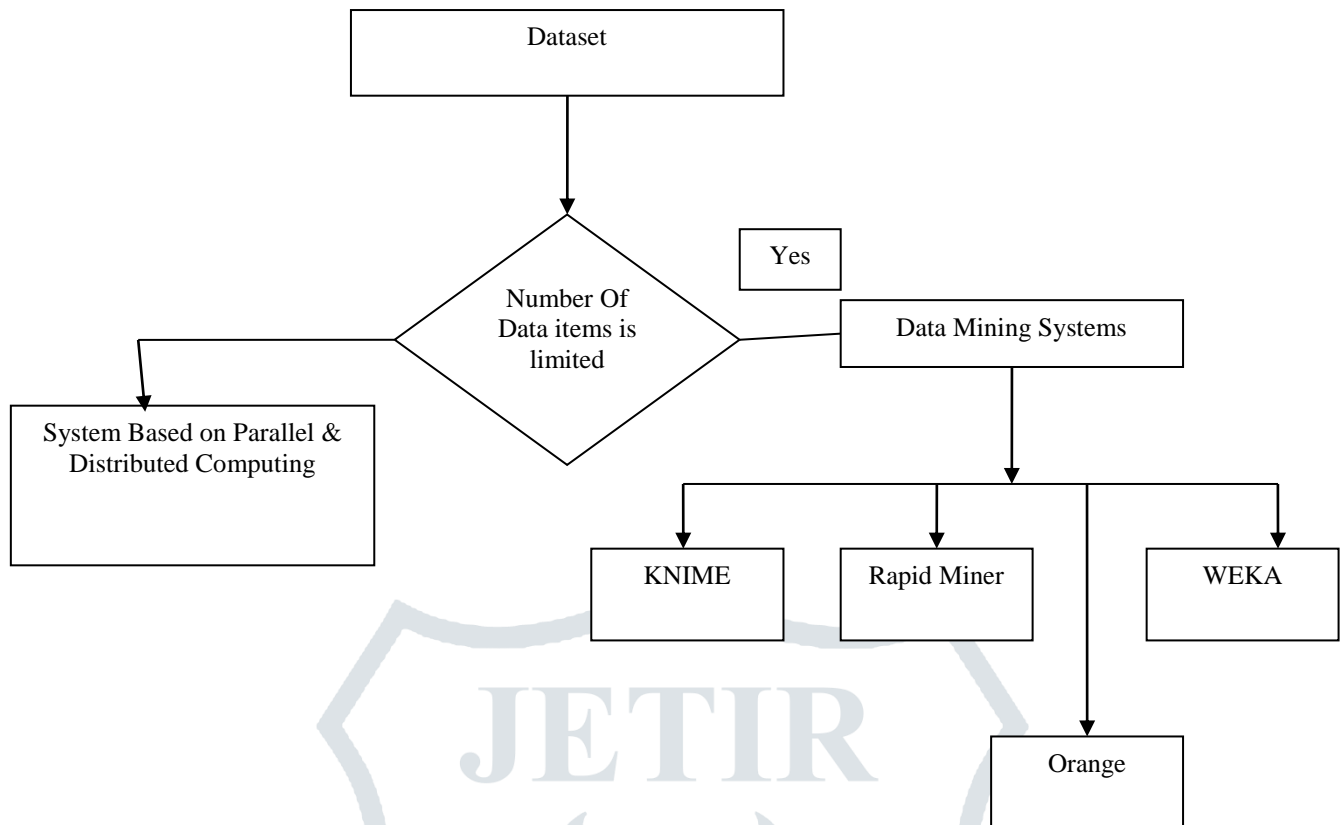
Grouping is isolated into two subgroups: Hard and Soft Clustering. In hard bunching, every datum point either has a place with a group totally or may not. In delicate grouping, rather than putting every datum point into a different bunch, a likelihood or probability of that information point to be in those groups is relegated.

A machine connecting with a domain in a smart way is called as machine insight. Creating calculations that empower a framework to interface with condition is viewed as Machine Learning. Numerous associations utilize man-made reasoning for settling on powerful basic leadership from the information gathered and to design its strategy for future.

Machine knowledge is utilized in assortment of utilizations:

1. Optimize procedure and computerize the procedure.
2. Extraction and grouping of information.
3. Detecting, examining and foreseeing the patterns or examples.

Figure 1: Existing Methods for Clustering



## II. MOTIVATION & PROBLEM STATEMENT

In the ongoing years, because of the expanding patterns in innovation has prompted the ascent in enormous volumes of information coming about to more affordable stockpiling gadget space. Preparing such enormous datasets is a troublesome errand as the datasets are spread over numerous areas. For instance web server log document that records the client's exercises in a specific site. The activity of information examiners is getting to be hard to viably break down the data from such gigantic voluminous information. Conventional instruments utilized in Data Mining neglect to process datasets of gigantic size and framing bunches is an intricate assignment as the information estimate is constantly developing and it is tedious procedure. Existing bunching calculations likewise endure wastefulness because of voluminous information.

Some of the time, the PC preparing worldview additionally changes when the datasets are developing in exponential size. In this way, bunching calculations should give successful groups that are versatile and this is accomplished by running calculations in parallel.

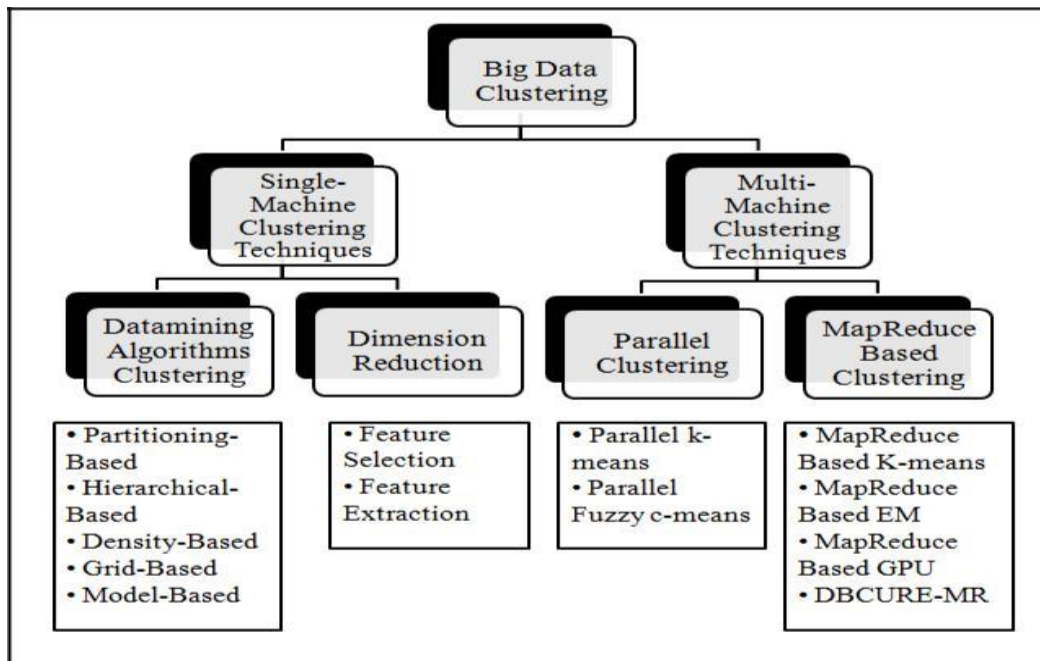
Difficulties looked in grouping of vast datasets:

1. Retrieving information from voluminous datasets represent a testing errand to the information experts.
2. For parallelizing conventional grouping calculations on MapReduce for substantial datasets is a troublesome assignment.

The above issues are tended to and settled in this paper:

1. Parallel bunching calculations have been executed on MapReduce to discover ideal groups for gigantic datasets.
2. Clustering calculations were executed on MapReduce in a solitary and multi machines.
3. Quality bunches are acquired from assorted information datasets.
4. Quality groups are created from low or high dimensional datasets from missing , off base or wrong information.
5. Results got are straightforward and unambiguous.

Figure 2: Big Data Clustering



### III. CLUSTERING TECHNIQUES FOR LARGE AND HIGH-DIMENSIONAL DATASETS

#### 3.1 Data Reduction

In a solitary machine, grouping for substantial datasets is finished by Balanced Iterative Reducing and Clustering utilizing Hierarchies (BIRCH), that limits the information and yield costs. MMC gives decrease of calculation time without influencing the execution.

#### 3.2 Sampling of Datasets

For malignancy datasets, Clustering Using Representatives (CURE) calculation is utilized for adaptability by examining the information.

#### 3.3 Decomposing and Space parceling of datasets

For spatial datasets, the normal groups are substantial; consequently calculation time is diminished by breaking down the information and after that Clustering in Quest (CLIQUE) calculation is kept running on the deteriorated information.

#### 3.4 Finding Projected Clusters

For quality articulation datasets, finding anticipated bunches is finished by utilizing Proclus iterative calculation. In the event that bunches must be developed one by one, Document Clustering (DOC) calculation is utilized.

#### 3. Probabilistic Approach for discovering Clusters

For high dimensional information, Statistical Subspace Clustering (SSC) technique is utilized.

### IV. MACHINE INTELLIGENCE ALGORITHMS

Numerous viewpoints in bunching of enormous datasets are actualized; K-Means calculation is taken for handling. K-Means calculation is connected for US wrongdoing datasets that is accessible for nothing. Without loss of liberality, we take the datasets with various fields arbitrarily. Broadened variants of K-Means calculation is connected for cross or multidimensional datasets and joined to give the general aftereffects of the two datasets and the outcomes acquired are connected for MI frameworks for powerful basic leadership. To diminish commotion from insignificant traits, Principal Component Analysis (PCA) calculation is utilized. For expansive datasets, we center around diminishing memory and processor loads at normal interims. Results acquired incorporate non-recursive execution of the calculation on different datasets for compelling basic leadership and picturing bunches in various viewpoints. In the event that the quantity of bunches is constrained, group part depends on the nature of the split.

#### 4.1 K-Means MapReduce calculation (KM-MR)

Input Information :

O : {o1,o2,o3,... ..on};/number of items to be bunched

X : X number of groups

Mi : Maximum number of emphases

Yield :

Wanted yield with number of groups

K-Means – MR(values or information)

i ← 0

```

For each datapoint  $d \in D$  do
IC ← SELECT(X,d)
INPUT(d)
WRITE(IC)
PC ← IC
while (genuine)
call to job.mapper()
call to job.reducer()
NC = READ ()
In the event that refresh ((NC,PC)>0)
PC=NC
else
refresh NC to result
i++
Result=READ ()

```

#### 4.2 Modified K-Means Clustering Algorithm (M - KM)

Map Phase Algorithm:

Info:

M dimensional information objects( $m_1, m_2, m_3, \dots, m_n$ ) for every mapper

X: number of bunches

Peruse beginning bunch centroids as  $i_1, i_2, i_3, \dots, i_k$

Yield:

yield list<a,b>

list\_new : new centroid list

set  $k=0$

list\_new=0

for all  $d \in D$

for all  $ij \in T$  do

$b_i \leftarrow \emptyset$  where  $b_i$  speaks to centroid nearest to the information object

InC ← ∞

ItC ← ∞

For all  $o_i \in O$  do

$i \leftarrow 0$

$l(o_i) \leftarrow \text{Euclidean Distance}(o_i, o_j), j \in \{1, 2, 3, \dots, k\}$

$i \leftarrow 0$

$b \leftarrow 0$

rehash

for each  $e_i \in E$  do

$\text{minDist} \leftarrow \text{Euclidean Distance}(o_i, c_j), j \in \{1, 2, 3, \dots, k\}$

if( $\text{curr\_centroid}=0$  or  $l(o_i) < \text{minDist}$ ) at that point

refresh InC

else

refresh ItC

$b_i \leftarrow b_i + 1$

$i \leftarrow i + 1$

make a yield list<a,b> with each item and the group centroid that it has a place with

rehash until assembly

#### 4.3 Reduce Phase Algorithm:

Information:

Let  $(a, b) \rightarrow \text{key, esteem}$  where  $a = l(o_i)$

value = objects allocated centroids by mappers

$O_i$  speaks to mapper yields

Yield:

list\_new : new centroid list(NC)

list\_new=0

NC ← ∅

for all  $x \in O_1$

centroid ← x.key

information object ← x.value

NC ← dataobject

for all  $c_i \in M$  do

NC ← ∅

sum\_objects ← ∅

num\_objects ← ∅

for all  $o_i \in O$  do

sum\_objects += object

num\_object++

NC ← (sum\_objects/num\_objects)

outputlist ← NC list

Return NC

Formulas to calculate inter and intra clusters:

$$InC = \frac{1}{2} \left( \frac{\sum_{i=1}^{O1} \sum_{j=1}^{O2} (Ai - Bj)^2}{O1 * O2} \right)$$

$$ItC = \frac{1}{2} \left( \frac{\sum_{i,j=1}^{O1+O2} (Ai - Bj)^2}{(O1 + O2) * (O1 + O2 - 1)} \right)$$

Where InC is inter cluster distance and O1, O2,.....are data points in clusters 1 , 2 and so on.

Ai is ith data point in cluster 1 and jth data point in clusters A and B.

### V. ANALYSIS OF MMC ALGORITHMS

The tasks that are performed using MMC algorithms are shown in the given below table1.

**Table 1: Tasks performed by MMC**

Objective	Tasks implemented using MMC	Datasets used
Finding the top rated URL for a particular product	Data is stored in HDFS in such a format which consists of unique id for each URL with date and date after map operation. This output results are used by the users to check the top rated URL for a particular search.	Amazon customer review datasets
Identification of microRNA Sequences	First, a file containing all microRNA sequences or combinations were generated. The next tasks are taken as the input for the Hadoop job and each line is taken as a individual Map task.	Cancer Datasets available in Broad institute
Detecting hidden patterns from biomedical datasets	Hierarchical clustering with K-means algorithm	Biomedical repository available in PubMed
Finding user session from log file	Based on unique id and timeout, all the pages visited are splitted if the user is in idle mode for 30 minutes.	Web log dataset available in Kaggle.
Objective	Tasks implemented using MMC	Datasets used
Finding the particular image from a huge set of complex image database	K-Means algorithm implemented on MapReduce. Map: Allocates the data objects to centroids with minimum distance centroid. Reduce: recalculates the value of the centroids when all objects have been assigned.	Coco datasets

Analyzing the data points in twitter	Subspace clustering Outputs from the each mapper are combined using a combiner function and local centroid values are calculated and the reducer calculates the global centroid from the combiner output.	Twitter datasets
--------------------------------------	--	------------------

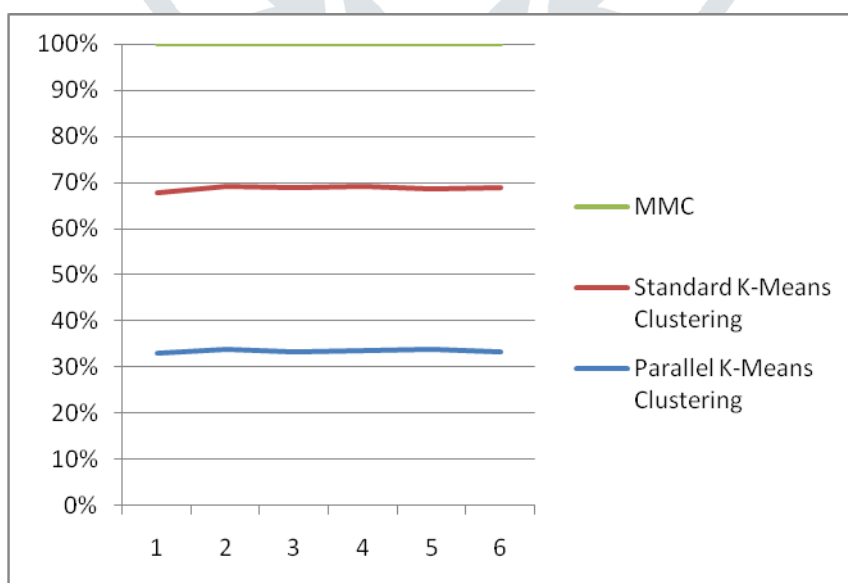
**VI. TESTED ENVIRONMENT**

The proposed bunching calculations are executed in R with different info and yield parameters. Information parameters are the dataset to be apportioned, greatest number of bunches, group measure and so forth. The yield acquired from the given info parameters are the groups files for each article. The yield is imagined in diagrams for different bunching strategies.

**Table 2: Comparison of total execution time**

	Parallel K-Means Clustering	Standard K-Means Clustering	MMC
Average	4324.16	4503.19	4215.18
	3889.32	4029.23	3520.52
	3694.19	3909.18	3415.32
	3963.24	4145.32	3624.23
	3123.26	3200.12	2900.16
	3693.28	3923.18	3423.44

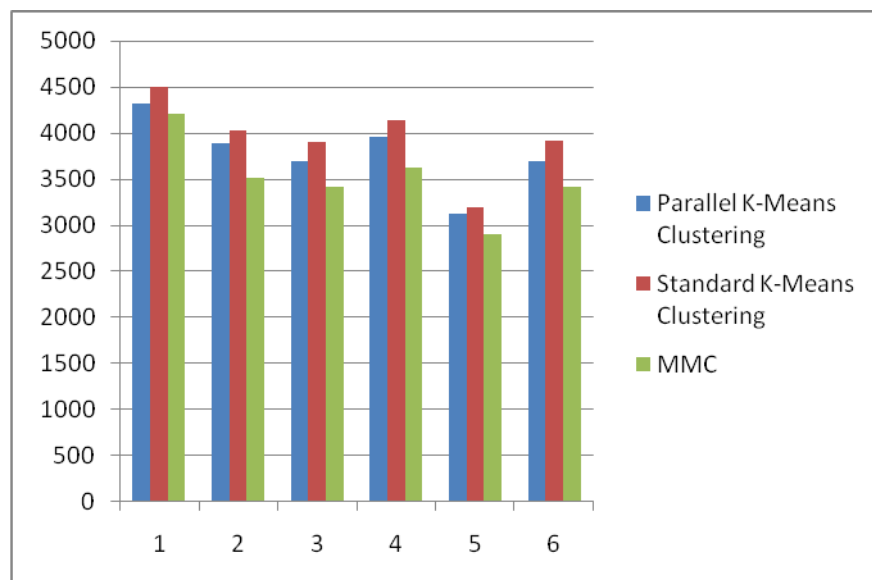
**Figure 2: Comparison of total execution time**



X-Axis – Time period  
Y-Axis – Percentatge of execution times



Figure 3: Comparison of total execution time



X-Axis – Time period

Y-Axis – Percentage of execution times

## VII. CONCLUSION

The proposed research method is useful for cluster identification and improves the performance of the process and provides fault – tolerance of the machines during the entire process. At each step of process, backup of the entire process till the last step is taken so as to keep the datasets available for future process. Simulated results show that the proposed method can be adopted for clustering datasets in optimal time.

## REFERENCES

- [1]. Robson L. F. Cordeiro et.al,(2011),KDD –ACM ,”Clustering Very Large Multi-dimensional Datasets with MapReduce”, 978-1-4503-0813-7/11/08.
- [2].Sridhar Ramaswamy et.al, “Efficient Algorithms for Mining Outliers from Large Data Sets” Korea Advanced Institute of Science and Technology, Advanced Information Technology Research Center at KAIST.
- [3]. Dongkuan Xu et.al,” A Comprehensive Survey of Clustering Algorithms”, Ann. Data. Sci. – Springer (2015)- DOI 10.1007/s40745-015-0040-1.
- [4]. Nivranshu Hans et.al, ,” Big Data Clustering Using Genetic Algorithm On Hadoop MapReduce”, INTERNATIONAL JOURNAL OF SCIENTIFIC & TECHNOLOGY RESEARCH, VOLUME 4, ISSUE 04, APRIL 2015, ISSN 2277-8616.
- [5]. Jeffrey Dean et.al,” MapReduce: Simplified Data Processing on Large Clusters”, Google, Inc., USENIX Association OSDI ’14: 6th Symposium on Operating Systems Design and Implementation.
- [6]. Weizhong Zhao et.al, “Parallel K-Means Clustering Based on MapReduce”, M.G. Jaatun, G. Zhao, and C. Rong (Eds.): CloudCom 2014, LNCS 5931, pp. 674–679, 2014. Springer-Verlag Berlin Heidelberg 2014.
- [7]. Max Bodoia,” MapReduce Algorithms for k-means Clustering”, 2014.
- [8]Y. unliang Chen et.al,” Mining association rules in big data with NGEP”, Cluster Comput. (2015) –Springer 18:577–585. DOI 10.1007/s10586-014-0419-3.
- [9]. Iulian V. Iliş , Thesis (2010),” Cluster Analysis for Large, High-Dimensional Datasets: Methodology and Applications”.
- [10]. Tanvir Habib Sardar, Zahid Ansari ,”Partition based clustering of large datasets using MapReduce framework: An analysis of recent themes and directions”, Future Computing and Informatics Journal xx (2018),pp:1-15.