# Clustering of Datasets Using K-Means Algorithm in SPARK

Subiksha N[1], Pallavi R Reddy[1], Mounica B[2]

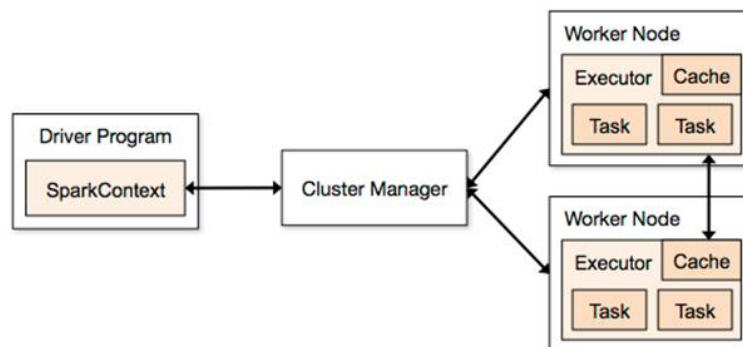B. E Student, Department of ISE, New Horizon College of Engineering, Bengaluru, Karnataka, India[1]

Assistant Professor, Department of ISE, New Horizon College of Engineering, Bengaluru, Karnataka, India[2]

**ABSTRACT**: With The growing trend of large volume of data generated every day, storing and processing the data is a great challenge. The existing solutions like SAS, R, Excel and MapReduce prove to be inefficient. MapReduce is a part of Apache Hadoop that allows distributive processing of unstructured data, where each distributive node has its own storage. Most of the algorithms are iterative in nature. But MapReduce is bane as it involves undesirable amount of read and writes to process such iterative algorithms. As an advancement, in this paper, we use the SPARK to implement algorithms like K-means, linear regression etc on real-time data. Spark is an unified engine that can run Hadoop, Mesos, cloud or standalone. It stores the intermediate result in RAM, thus avoiding read/write from/to disk. Like the MapReduce, Spark is used for batch processing. In addition, Spark can be used to handle streaming data, queries and machine learning. In this paper, we prefer to use the cloud services to access the large datasets because cloud is flexible, scalable and involves less hardware cost than using our own infrastructure which requires many software's (Hadoop, Ubuntu, Java etc) to be installed. Once the data is obtained, we run the K-means algorithm on Spark using EMR (Elastic MapReduce).

**KEYWORDS**: MapReduce, iterative algorithms, Spark, K-means, cloud services

## I. INTRODUCTION

Hadoop is used extensively in industries to analyze data sets, as this framework is based on simple programming models. The solutions are scalable,flexible,cost-effective and fault-tolerant. The main concern is to maintain speed in processing large datasets. To overcome this Spark was introduced for improving the computational speed. In memory cluster computing is the main feature that increases the processing speed. Spark and Hadoop are different, Spark is not a modified version of Hadoop. Spark has its own cluster management. Storage and processing is the two ways in which Spark uses Hadoop. It is used only for storage as Spark has its own cluster management.Spark is free and an open source framework, written in Java.

DATAFRAME

Collections of data that is organized into columns is a dataframe. It is equal to relational tables with some good optimization techniques. This can be constructed from an array of different sources. This was mainly designed for modern Big Data and data science applications

Some of the features of dataframes are:

- It can process the data in the form of Kilobytes till petabytes.

- This Supports different data forms and storage systems like HDFS, HIVE tables, mysql, etc.

- This can be easily integrated with all Big Data tools and frameworks using the Spark-Core component.

- This also provides an API for Python, Java, Scala, and R.

RDD (Resilient Distributed Dataset)

RDD's are the fundamental data structures of spark. RDD is Resilient Distributed Datasets which is an immutable distributed collection of objects. RDD's can have any type objects including Python, Java or Scala. It is a partitioned collection of records in the read-only mode. These are fault tolerant collection of elements which can be operated in parallel. There are two different ways to create RDD's: to parallelize an existing collection or for referencing a dataset which is stored in an external storage like HDFS or HBase. To achieve faster and efficient MapReduce operations Spark uses the RDD concept.

COMPONENTS OF SPARK

Apache Spark Core- This is the engine for Spark execution in which all the functionality is built upon. The main function is it provides in-memory computing and datasets which are stored on external storage systems can be accessed, like the data can be there on cloud or in local like HDFS.

Spark SQL- This is the component which provides support for structured and semi-structured data. It is placed above the Spark Core.

Spark Streaming- The data is taken in the form of mini-batches and RDD is performed on mini-batches.

MLib(Machine Learning library)- This is placed above Spark because of distributed memory based Spark architecture. This is nine times faster than the Hadoop disk based version of Mahout.

GraphX- It is a graph processing framework. It provides an API for graph computation. Provides an optimized run time for this abstraction.

## II. RELATED WORK

MapReduce

BigData can be processed in parallel on model nodes using MapReduce programming model[1].As data is accumulating in large quantities everyday it is not possible to process BigData using the traditional methods. Traditional systems use the centralized server to process and store the data. But this is not suitable for large data. So, Google developed the MapReduce algorithm. In MapReduce the task is divided into many small parts and each part is given to a computer. This algorithm performs two tasks: Map- in this one dataset is converted to another dataset. The individual elements are broken down into tuples like the key-value pairs. Reduce- The input to the reduce is the output
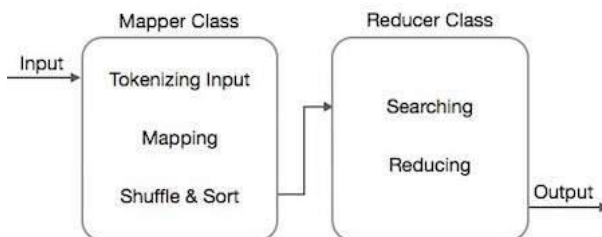
of the Map function. This combines the key value pairs into smaller set of tuples. Reduce is performed after the Map task.



The MapReduce works well only on batch processing[2]. This is basically designed only for one specific problem domain. Often this programming paradigm is hard to understand.



### III. PROPOSED WORK

To overcome the time overhead involved in executing iterative algorithms like K-Means using MapReduce , in this paper we prefer to use a different framework called Spark. Spark is an unified engine that is 100 times faster than MapReduce.

*A) Case Study*

We consider the Uber trip data for clustering based on location and time. The input data has three features

- Data and time
- Latitude
- Longitude

The input data is in the form of CSV (Comma Separated) file.  Initially, we consider five clusters as, B02617, B02598, B02682, B02764 and B02512.

![IJIRCCE logo]

ISSN(Online): 2320-9801
ISSN (Print): 2320-9798

# International Journal of Innovative Research in Computer and Communication Engineering

*(An ISO 3297: 2007 Certified Organization)*

*Website: www.ijircce.com*

### Vol. 5, Issue 5, May 2017

B) *Implementation*

Step 1: As mentioned earlier, we prefer to use the cloud services (Amazon web services) to obtain the desired infrastructure. We create a cluster consisting of say, three "m4.large" machines, each machine having the following configuration:

- vCPU Count                - 2
- RAM                    -8Gib
- Instance Storage          -EBS (Elastic Block Storage)
- Network Performance    -Moderate
- EBS-Optimized            -450Mbps

In this step, we also choose the required software configuration (software's) like Spark, Hadoop, Hive etc.



We require a key pair with name of our choice to access these services from our local machine.

Step 2: Next, we upload the data on which we work to the S3 (Simple Storage Service). S3 is a storage service provided by Amazon web services. It can handle up to many terabytes of data. Once the data is uploaded, we can use the S3 for further references to it.

Step 3: Next, we install PuTTy to access the cloud services from our local machine using key value pair created in step 1. This instantiates EMR (Elastic Map Reduce) which has all the required software's like Spark, Mahout, Hive pre installed.



Step 4: Using vi commands, we type the code for the K-means algorithm[3] in any programming language. We consider, Python.

| Methods/functions used | Meaning |
|---|---|
| 1. Dataframe | Dataframe consists of the Dataset in named columns. |
| 2. Dataframe.show() | Displays the first 20 rows of the dataset. |
| 3. VectorAssembler (*InputCols, OutputCol*) | It combines the various columns specified in *InputCols* parameter into a single column in the form of list with the name given in *OutputCol* parameter. |
| 4. StringIndexer (*InputCos, OutputCol*) | It converts the string values into integers. |
| 5. Randomsplit([n1,n2]) | Splits the input data into two parts as specified in the parameter.Ex- n1=0.8 , n2=0.2 |
| 6. Kmeans (*k, seed, featuresCol*) | Performs the clustering into 'k' clusters on the parameter *featuresCol* with any intial seed value. |
| 7. saveAsTextFile (output_path) | Saves the output to the specified location in the S3. |
| 8. Dataframe.cache() | Stores the dataframe in the cache memory. |

Step 5: Next, we run the code using the following command
*spark-submit codefile.py  inputfile_path outputfile_path*

*codefile.py* is the file containing python code
*inputfile_path* is the S3 location of the raw input data
*outputfile_path* is the S3 location of the final output.

## IV.  RESULTS

The result of the clustering is stored in S3. It consists of a column named "prediction" which displays the cluster to which the algorithm has assigned the data.



In the implementation, we split the input data into two parts – training data and test data, thus giving a touch to supervised machine learning. We then allow the algorithm to predict the clusters for the test data.  The accuracy of the predictions is checked using a function available in Spark-

*evaluator=MulticlassClassificationEvaluator(labelCol="label",predictionCol="predicted_cluster", metricName="accuracy")*

*accuracy_lg = evaluator.evaluate(prediction)*

This gives the accuracy of the algorithm.

## V.  CONCLUSION AND FUTURE WORK

In this paper, we discussed how to apply the clustering algorithm to large datasets using Spark in order to improve the time of execution by caching the intermediate results in RAM. We have worked on structured and semi-structured data sets so far. This algorithm can be extended to take inputs in the form of unstructured data sets like images and video, which requires the use of other software's integrated with Spark like Tensorflow.

## REFERENCES

1. Madhavi Vaidya, 'Parallel Processing of cluster by Map Reduce', International Journal of Distributed and Parallel Systems (IJDPS) Vol.3, No.1, January 2012.
2. Jeffrey Dean and Sanjay Ghemawat,  'MapReduce: Simplified Data Processing on Large Clusters', Google, Inc.
3. Jyoti Yadav, Monika Sharma, 'A Review of K-mean Algorithm', International Journal of Engineering Trends and Technology (IJETT) – Volume 4 Issue 7- July 2013.
4. V Srinivas Jonnalagadda, P Srikanth, Krishnamachari Thumati, Sri Hari Nallamala, 'A Review Study of Apache Spark in Big Data Processing', International Journal of Computer Science Trends and Technology (IJCST) – Volume 4 Issue 3, May - Jun 2016
5. Abhishek Bhattachary, Shefali Bhatnagar, 'Big Data and Apache Spark: A Review', International Journal of Engineering Research & Science (IJOER), Vol-2, Issue-5 May- 2016