

Big Data Framework for storage Extraction and Identification of Data using Hadoop Distributed File system

B. Suvarnamukhi, M. Seshashayee

Abstract: Big data is all about the developing challenge that associations face in today's world, As they manage enormous and quickly developing wellsprings of information or data, with the complex range of analysis and the problem includes computing infrastructure, accessing mixed data both structured and unstructured data from various sources such as networking, Recording and stored images. Hadoop is the open source software framework includes no of compartments that are specifically designed for solving large-scale distributed data storage. MapReduce is a parallel programming design for processing

Keywords: Big data, Hadoop, MapReduce, Parallel Programming.

I. INTRODUCTION

Hadoop Distributed File System: HDFS can be used to work on large volumes of data which is available on the web like data node and name node monitors where the information is put away. The methods of establishment are independent mode, pseudo circulated mode, and completely conveyed mode.

Independent mode: As the name suggests everything runs in single Java virtual machine (JVM). These modes are good in development and testing only with small data.

Pseudo circulated mode: In this mode Multiple JVMs are used has a single node in each single machine clusters are simulated. Fully Distributed mode: All the components run has a separate nodes, these are very good for staging and production.

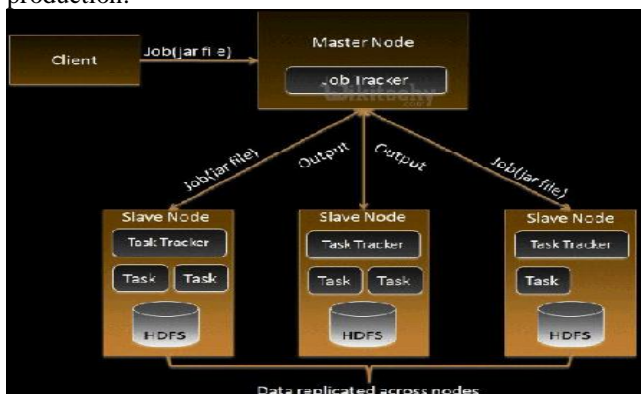


Fig1. Hadoop distributed File system:

II. BACKGROUND

The pattern extraction from the data is existing since centuries. Bayes theorem and regression analysis were used for classifying patterns in datasets. The advancement of technology has increased the capability of computer automation in data collection. Thus data mining is the process of discovering several hidden patterns from raw data to make useful information. Data sets are specified by Input formats, which is defined by Input data and identifies the partition of the data. The input splits the data, records the data of reader objects and extracts the key value pairs.

The input files will read all the files from the specified directory and send them to the mapper. The mapper delegates the files list to a method by overriding the subclasses.

Example: create your own file with file input format as ABC, which can read as * ABC from directory list.

Record Readers object: Each input format provides its own record reader implementation.

```
public class objectpositioninputformat extends
Fileinput Format<t,point3d>
{
Public RecReader<t,point3d>getRecReader
{
Input.split input, JobConf job,Reporter reporter)throws
IOException
{
Rep.setstatus(input .toString());
Returnnew objposRecReader(job,(filesplit)
input);
}
inputsplit[]getsplits(jobConfjob,int numsplits)throws
IOException;
}
```

Revised Manuscript Received on December 13, 2019.

* Correspondence Author

B. Suvarnamukhi*, Dept. of CSE, Assistant Professor in St Mary's Group of Institutions India., mukhi.suvarna@gmail.com

M. Seshashayee, Department of Computer Science, GITAM Deemed to be University, Visakhapatnam, India., mshashayee@gmail.com

III. ANALYTICS IN BIG DATA CONCEPTS

Enormous Data examination contrasts from the conventional information investigation in volume, speed and assortment which are the attributes of the information being. To distinguish the necessities for the investigation of Big Data. This is a bit by bit procedure is expected to compose every one of the exercises and procedure the undertakings associated with obtaining, handling and dissecting the information. In Data examination process are partitioned as accompanying Nine (09) phases.

IV. MAPREDUCE

MapReduce: MapReduce is a programming model used in processing, huge amounts of data in parallel. It contains two methods namely Map () and Reduce (). In Mapper it process the several chunks of the data and in Reduce stage it process the data received from the Mapper stage.

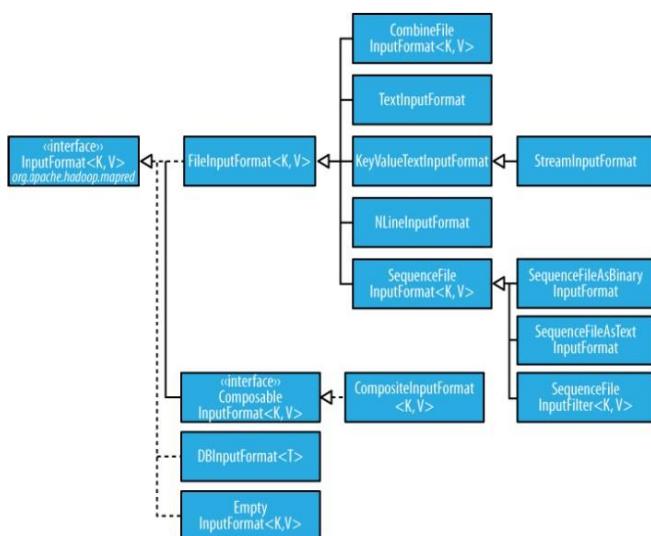


Fig 2. Input split file

Table 1 Lifecycle with stages

S NO	LIFE CYCLE PHASES	PHASE
1	Evaluation method	phase 1
2	Identification process	phase 2
3	Acquisition and process in Filtering	phase 3
4	Data Extraction	phase 4
5	Validating and Cleansing the Data	phase 5
6	Aggregation methods and Representation	phase 6
7	Phases in Analysis process	phase 7
8	Visualization process in Data	phase 8
9	Data Utilization and Analysis Results.	phase 9

Evaluation method

In Analytics the lifecycle phase should always begin with a well-structured evaluation method that represents a complete detail of the justification, motivation which carry out the analysis. This helps in decision makers to understand the resources. The goals of the analysis should be specific and relevant.

Identification process

The second stage is utilized to recognize the datasets which are vital for the procedure of information gathered from different sources identifying an enormous assortment of information sources may build the likelihood of finding concealed examples.

Acquisition and process in Filtering

On account of inward datasets, a rundown of accessible datasets from inside sources, for example, information shops, are normally executed and coordinated against a predefined dataset.

In outside datasets, information can be gathered from the outsider information suppliers, for example, information from business sectors and openly accessible information some types of outer information might be inserted inside web journals or different sorts of substance based sites by means of mechanized instruments.

Data from all of the data sources that were identified in the identification stage. The collected data can be self-operating can be done for the removal of bad data.

Extraction phase or pulling out Data

The basic step from the data identification is considering the input in the analysis phase and pulls out the data in a specified format, with the incompatible data applies the solution, this basic step is to dissimilar categories of data are collected from the various external sources.

Process of Validating and Cleansing the Data

Weak data can crooked and shows false analysis results. As in existing data, the structure of data is established in advance and data is pre-validated.

Consider two datasets like Dataset A1 and Dataset A2. If the value in Dataset 2 is valid data and the respective value in Dataset 1 should be valid.

- The other value in Dataset A2 is not true against its corresponding value in Dataset A1. Suppose any value is not found, it is inserted from Dataset A1

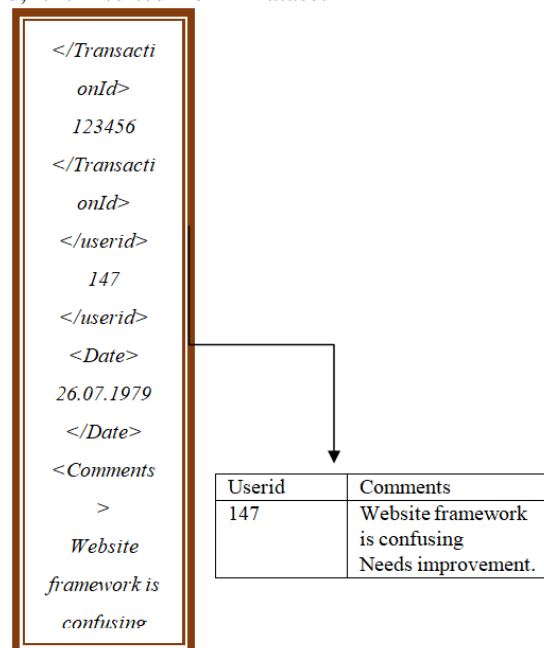


Fig 3. Data Extraction

Example

1	-	-	-
2	-	-	-
3	-	-	-----
4	-	-	-----
5	-	-	-----

= √ ●
≠ X ●

1	-	-	-
2	-	-	-
3	-	-	-
4	-	-	-
5	-	-	-

Dataset A1

Dataset B1

V. CONCLUSION

This paper is all about Data, Storage of Data using HDFs and various types of input data .It also tells methods of retrieving hidden format data. We also discussed MapReduce techniques for processing of data.

REFERENCES

1. 'Kunal Jain Amrit pal', Pinki Aggrawal,Sanjay Aggrawal "A Performance Analysis of Mapreduce Task with Large Number of Files Dataset in Big Data using hadoop" Fourth International Conference on Communication Systems and Network Technologies, 2014. [4] Kaur, Anureet.
2. J.Q.Simon Woolf, S.H Purnell , S.M.Zimmerman, E.BCamberos, G.J.Haley R.P. 2015.'Translating Evidence into Population Health Improvement: Strategies and Barriers, 'Annual review of public health (36), pp. 463-482.
3. Suvarnamukhi, B., and M. Seshashayee. "Big Data Concepts and Techniques in Data Processing." (2018).
4. N.Baro, E.D.SamuelBeuscart, RegisChazard, Emmanuel. 2015. "Toward a Literature-Driven Definition of BigData in Healthcare," *BioMed research international* (2015), 2015, pp. 639021-639021.
5. BigData Survey, Technologies, Opportunities, and Challenges <http://www.hindawi.com/journals/tswj/2014/712826>.
6. Ghose, Anindya, and Arun Sundararajan. "Evaluating pricing strategy using e-commerce data: Evidence and estimation challenges." *Statistical Science* 21.2 (2006): 131-142.

AUTHORS PROFILE



B.SuvarnaMukhi pursuing Ph.D. in Department of Computer Science at GITAM (Deemed-to-be University) Visakhapatnam, Andhra Pradesh and Working as an Assistant Professor in St Mary's Group of Institutions Hyderabad. Having 9 Years of Teaching Experience. Published Research paper in International Journal of Computer Science and Engineering (IJCSE) volume 6, Issue 10 – October 2018, published Research paper in International Journal of Innovative Technology and Exploring Engineering (IJITEE) volume 8 Issue 12 – October 2019. Member of IAENG Interested to do research area in BIG DATA.



M.Seshashayee, has been awarded Ph.D. in Computer Science and Technology, and presently working as Assistant Professor in the Department of Computer Science at GITAM (Deemed to be University), Visakhapatnam, and Andhra Pradesh. Her Research specialization is in Image Segmentation Methods using Data Mining Techniques. She has published more than 12 publications in International Journals. Editorial Member for Honorary Editorial Board of International Journal of Computational Mathematical Ideas (IJCMI), Reviewer in International Journal of Innovations in Computer Science and Engineering, a member of CSI and IAENG.

