

Analysis of Data Using Hadoop and Mapreduce

^[1] Bharti Kalra, ^[2] Dr. Anuranjan Misra, ^[3] Dr. D. K. Chauhan

^[1] P.h.D. Scholar(CSE) ^[2] Prof. & Head, Dept of CSE & IT ^[3] Director Technical
Noida International University, Greater Noida(India)

^[1]kalra.bharti1@gmail.com^[2] amc290@gmail.com, ^[3]prof.dkchauhan@gmail.com

Abstract: — We are in the world of technology, where data is moving around all the time. This data is becoming bigdata when it comes in huge volume and data can be structured, unstructured or semi-structured. hadoop is the technology to analysis the big data. The objective of this paper is to analysis the data that is collected from the open source using the hadoop and mapreduce programming model.

I.INTRODUCTION

Data is the most important assets for any organization. This Data can be collected by individual, or by any team or by any organization for any objective. So Big data, when this term comes into mind, so many definition moving around. Many people and organizations describe bigdata as its own way, based on the parameters (5vs)i.e. volume, variety, velocity ,variability and complexity. According to [2] Big Data is defined as the very large amount of datasets that hold variety of data. Big data is data that exceeds the processing capacity of conventional database systems. The data is too big, moves too fast, or doesn't fit the strictures of your database architectures. To gain value from this data, you must choose an alternative way to process it [1].

II.BIG DATA ANALYTICS

According to a report [3], Co-sponsored by IBM, Big data analytics is where advanced analytic techniques operate on big data sets. Hence, big data analytics is really about two things—big data and analytics—plus how the two have teamed up to create one of the most profound trends in business intelligence (BI) today. Let's start by defining advanced analytics, then move on to big data and the combination of the two.

Big data analytics is the process of examining big data to uncover hidden patterns, unknown correlations and other useful information that can be used to make better decisions. With big data analytics, data scientists and others can analyse huge volumes of data that conventional analytics and business intelligence solutions can't touch.[4]

III.HADOOP

Apache Hadoop is an open source software framework for storage and large scale processing of data-sets on clusters of commodity hardware. Hadoop is an Apache top-level project being built and used by a global community of

contributors and users. It is licensed under the Apache License 2.0[5].

The Apache Hadoop framework is composed of the following modules

A.Hadoop Common: Contains libraries and utilities needed by other Hadoop modules

B.Hadoop Distributed File System (HDFS): A distributed file-system that stores data on the commodity machines, providing very high aggregate bandwidth across the cluster

C.Hadoop YARN: A resource-management platform responsible for managing compute resources in clusters and using them for scheduling of users' applications

D.Hadoop MapReduce: A programming model for large scale data processing in addition to the previous overall configurations, the following individual components are also included on HDInsight clusters.

E.Ambari: Cluster provisioning, management, and monitoring.

F.Avro (Microsoft .NET Library for Avro): Data serialization for the Microsoft .NET environment.

G.Hive & HCatalog: Structured Query Language (SQL)-like querying, and a table and storage management layer.

H.Mahout: Machine learning.

I.Oozie: Workflow management.

J.Phoenix: Relational database layer over HBase.

K.Pig: Simpler scripting for MapReduce transformations.

L.Sqoop: Data import and export.

M.Tez: Allows data-intensive processes to run efficiently at scale.

O.ZooKeeper: Coordination of processes in distributed systems.

This Paper proposes an implementation of Analytics of Big Data using Hadoop and Mapreduce.

IV. SIMULATION BACKGROUND

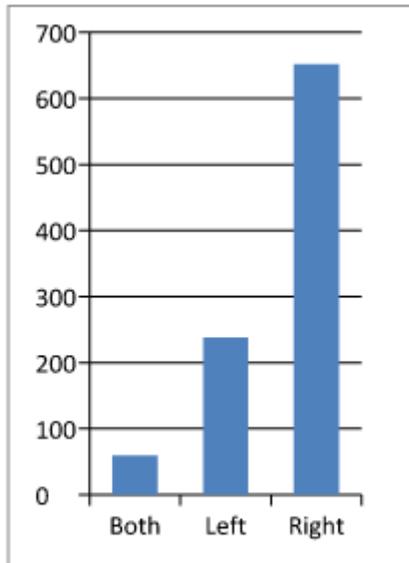


Fig C

VI. CONCLUSION

Our aim in this paper to analysis the research data through the mapreduce Programming model and summarizes the result of some aspects of multiples columns.

REFERENCES

- [1] <https://www.oreilly.com/ideas/what-is-big-data>
- [2] using *privacy by design* to achieve big data innovation without compromising privacy, by deloitte, june 10, 2014.
- [3] big data analytics fourth quarter 2011 by philip russomftp://ftp.software.ibm.com/software/tw/defining_big_data_through_3v_v.pdf
- [4]https://www.umb.edu/academics/caps/corporate/big_data_analytics
- [5] <http://opensource.com/life/14/8/intro-apache-hadoop-big-data>
- [6]<https://azure.microsoft.com/en-in/documentation/articles/hdinsight-hadoop-introduction/>
- [7]<http://www.stat.ncsu.edu/people/monahan/courses/st590g/baseball/lahman591-csv/readme59.txt>