# Mining of Datasets using Big Data Technique: Hadoop Platform

Rohit Kumar[1], Daulat Sihag[2], Khushboo Sukhija[3]

[1]M.Tech Scholar, Department of Computer Science Engineering, JCDM College of Engineering, Sirsa
[2]Assistant Professor Department of Computer Science Engineering, JCDM College of Engineering, Sirsa
[3]Assistant Professor, School of Computer Science, Lingayas University, Faridabad

*Abstract*— BIG DATA IS THE FUTURE OF IT INDUSTRY. *Here see the methodology i.e. ETL process used for analysis of big data by using Hadoop ecosystem. The analysis of big data extracts business values from the raw data and helps in gaining competitive advantage by different organisations. There is a drastic growth of data in the web applications and social networking and such data are said be as Big Data. It requires huge amount of time consumption to retrieve those datasets. It lacks in performance analysis. To overcome this problem the Hive queries with the integration of Hadoop are used to generate the report analysis for thousands of datasets. The objective is to store the data persistently along with the past history of the data set and performing the report analysis of that data set. The main aim of this system is to improve performance through parallelization of various operations such as loading the data, index building and evaluating the queries. Thus the performance analysis is done with parallelization. HDFS file system is used to store the data after performing the MapReduce operations and the execution time is decreased when the number of nodes gets increased. The performance analysis is tuned with the parameters such as the execution time and number of nodes..*

*Keywords*— *Big Data, Hadoop, HDFS.*

## I. INTRODUCTION

To generate information it requires massive collection of data. The data can be simple numerical figures and text documents, to more complex information such as spatial data, multimedia data, and hypertext documents. To take complete advantage of data; the data retrieval is simply not enough, it requires a tool for automatic summarization of data, extraction of the essence of information stored, and the discovery of patterns in raw data. With the enormous amount of data stored in files, databases, and other repositories, it is increasingly important, to develop powerful tool for analysis and interpretation of such data and for the extraction of interesting knowledge that could help in decision-making. The only answer to all above is 'Data Mining'.

Data mining is the extraction of hidden predictive information from large databases; it is a powerful technology with great potential to help organizations focus on the most important. information in their data warehouses [1][2][3][4].Data mining tools predict future trends and behaviours, helps organizations to make proactive knowledge-driven decisions[2]. The automated, prospective analyses offered by data mining move beyond the analyses of past events provided by retrospective tools typical of decision support systems. Data mining tools can answer the questions that traditionally were too time consuming to resolve. They prepare databases for finding hidden patterns, finding predictive information that experts may miss because it lies outside their expectations.

Data mining, popularly known as Knowledge Discovery in Databases (KDD), it is the nontrivial extraction of implicit, previously unknown and potentially useful information from data in databases [3][5]. Though, data mining and knowledge discovery in databases (or KDD) are frequently treated as synonyms, data mining is actually part of the knowledge discovery process.

### 1.1 Knowledge Discovery Process

Data mining is one of the tasks in the process of knowledge discovery from the database. The steps in the KDD process contain:

- **Data cleaning:** It is also known as data cleansing; in this phase noise data and irrelevant data are removed from the collection.
- **Data integration:** In this stage, multiple data sources, often heterogeneous, are combined in a common source.
- **Data selection:** The data relevant to the analysis is decided on and retrieved from the data collection.
- **Data transformation:** It is also known as data consolidation; in this phase the selected data is transformed into forms appropriate for the mining procedure.
- **Data mining:** It is the crucial step in which clever techniques are applied to extract potentially useful patterns [1][3].
- **Pattern evaluation:** In this step, interesting patterns representing knowledge are identified based on given measures.
- **Knowledge representation:** It is the final phase in which the discovered knowledge is visually presented to the user. This essential step uses visualization techniques to help users understand and interpret the data mining results.

## II.   RELATED WORK

D. Abadi analyse the large scale data analysis with the traditional DBMS. The data management is scalable but there is replication of data. Replication of data leads to the fault tolerance [1]. J.ABABI, AVI SILBERCHATZ analyse massive datasets on very large clusters is done within the HadoopDB architecture for the real world application like business data warehousing. It approaches for parallel databases in performance. Still there is no scalability. It consumes huge amount of time for execution [4]. Farah Habib Chan chary analyse large datasets among the clusters of machines are efficiently stored in the cloud storage systems. So that the same information on more than one system could operate the datasets even if any one of the system's power fails [3]. According to IBM the amount of unstructured and multi-structured data within an average organization is about 80% (Savvas, 2011). Taking account the average data growth, annually by 59% (Pettey & Goasduff, 2011), this percentage will likely be much higher in a few years. Not only the volume of data is becoming a problem, also the variety and velocity are issues we need to look at (Russom, 2011). This phenomenon is called "big data" and is identified as one of the biggest IT trends for 2012 (Pettey, 2012) [5][6][7].

Time to market and innovation of new products are nowadays the key factors for enterprises. Data warehouses and BI support the process of making business decisions. These instruments allow the discovery of hidden pattern like asking unknown relations between certain facts or entities. This causes an important change: in the past, the question which has been run against the system was already known at the time of collecting the data, today it is common practice to catch all the data to hold it for questions which will be asked in the future [9]. That is the reason why Big Data is a hot growing topic in information science. To tackle the challenges of Big Data, a new type of technologies has emerged. Most of these technologies are distributed across many machines and have been grouped under the term "NoSQL". NoSQL is actually an acronym that expands to "Not Only SQL". In some ways, these new technologies are more complex than traditional databases, and in other ways they are simpler. There isn't any solution that fits all situations. You must do some compromises. This point will also be analyzed in this thesis. These new systems can scale to vastly larger sets of data, but using these systems require also new sets of technologies. Many of these technologies were first pioneered by two big companies: Google and Amazon. The most popular is probably the MapReduce computation framework introduced by Google in 2004. Amazon created an innovative distributed key-value store called Dynamo. The open source community responded in the year following with Hadoop (free implementation of MapReduce), HBase, MongoDB, Cassandra, RabbitMQ and countless other projects [Mar12]. The heterogeneous mixture learning technology is an advanced technology used in big data analysis. In the above, we introduced difficulties that are inherent in heterogeneous mixture data analysis, the basic concept of heterogeneous mixture learning and the results of a demonstration experiment that dealt with electricity demand predictions. As the big data analysis increases its importance, heterogeneous mixture data mining technology is also expected to play a significant role in the market. The range of application of heterogeneous mixture learning will be expanded broader than ever in the future [8].

## III.   METHODOLOGY USED

### 3.1   The Hadoop Platform

Hadoop is an open-source ecosystem used for storing, managing and analysing a large volume of data, designed to allow distributed processing of data sets across thousands of machines.
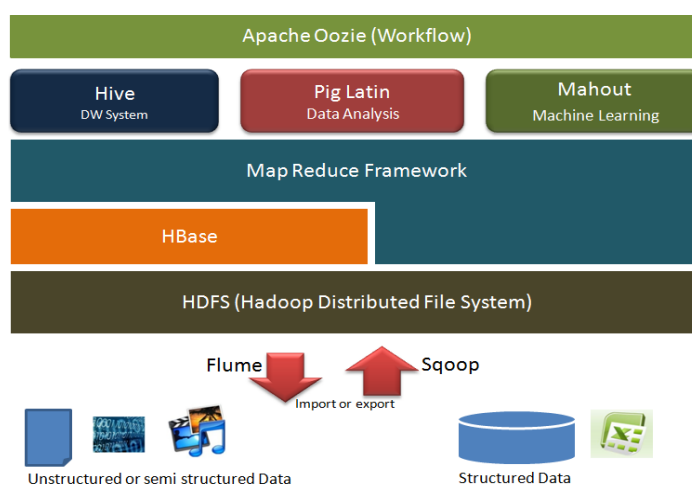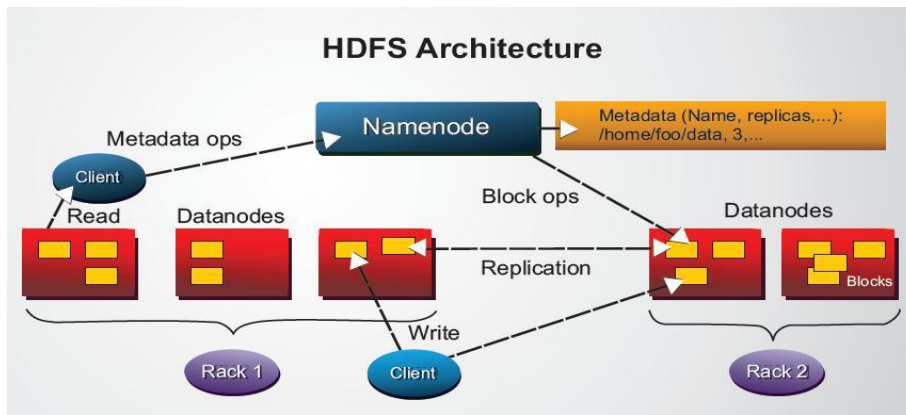


**FIG 1. HADOOP ECOSYSTEM**

**FIG 2. HDFS ARCHITECTURE**

A cluster running Hadoop means running a set of daemons on the different servers of the network. The daemons include:

- **Name Node:** it is the master of HDFS that directs the slave Data Nodes daemons. It is recommended to run the Name Node on its own machine. The negative aspect to the Name Node is that it is a single point of failure.
- **Data Node:** It is responsible of writing and reading HDFS blocks to actual files on the local file system. A Data Node may communicate with other Data Nodes to replicate its data for redundancy.
- **Secondary Name Node:** it is an assistant daemon for monitoring purposes.
- **Job Tracker:** Manages Map Reduce job execution. The Job Tracker daemon determines the execution plan, assigns nodes to different tasks. There is only one Job Tracker per Hadoop cluster.
- **Task Tracker:** Manages the execution of individual tasks on the machine. Communicates to Job Tracker in order to obtain task requests.

**3.2    FileZilla**

FileZilla is free, cross-platform FTP application software.  Consisting of FileZilla Client and FileZilla Server. Binaries are available for *Windows*, *Linux*, and *Mac OS X*. It supports *FTP*, *SFTP*, and *FTPS*.

Filezilla is used for moving files from local systems to Cloudera directory using FTP client from your Windows.



**FIG 3. TRANSFERRING FILES FROM WINDOW SERVER TO CLOUDERA**

**IV.    PRESENT WORK**

**4.1    Hive Component of Hadoop**

It Support structured data, e.g., creating tables, as well as extensibility for unstructured data.

- Hive Query for Creating Tables

Create table user (UserId int, Age int, Gender char) row format delimited fields;
- Hive Query for inserting/loading data into table
Load data Local inpath '/users/local/users.txt' into Table user;

**4.2    Flow Chart**



**FIG 4. OUTLINE OF PRESENT WORK**

## V.   RESULTS



**FIG 5. CREATE TABLE LOAN_STATS_YEAR07_YEAR11.**



**FIG 6. CREATE TABLE LOAN_STATS_YEAR07_YEAR11**



**FIG 7. CREATE TABLE LOAN_STATS_YEAR112_YEAR13**

**FIG 8. CREATE TABLE LOAN_STATS_YEAR12_YEAR13**



**FIG 9. LOADING DATA INTO TABLE "LOAN_STATS_YEAR07_YEAR11**



**FIG 10. LOADING DATA INTO TABLE "LOAN_STATS_YEAR12_YEAR13"**



**FIG 11. EXECUTION**

**FIG 12. EXECUTION**



**FIG 13. EXECUTION**

## VI.　CONCLUSION AND FUTURE SCOPE

### 6.1　Conclusion

In this present work, Data mining is performed by using the Hadoop Ecosystem Approach. The presented work is about to performed data mining on large loan data sets using Hive component of Hadoop Ecosystem. Here, the hive queries are performed for mining the useful data like "we are pulling those Bank Customers whose Loans were processed successfully starting from year 2007 till 2013.They were proven as the best customers for banks as their payment schedule and other incomes were verified and made on timely basis. They were provided Loan at ROI = 7% and for 3years duration." This information can be mined in less time because of parallelization feature of Hadoop ecosystem.

So from our data analytics we conclude that those customers are best market for Banks in future and can be given priority over other customers. Companies can create separate operational data store (ODS) to make inventory of those customers for faster search & processing of loan.

### 6.2　Future Work

In this present work, we have performed data mining on large structured data (big data) by executing hive queries on hive component of Hadoop Ecosystem.

Our future work will focus on Analysis or mining of unstructured data like images, audios, videos, graphs using map reduce component of Hadoop Ecosystem.

Obtained results from the system shows that the effective Data mining is been performed..

### REFERENCES

[1]　D.Abadi,"Data management in the cloud: Limitations and opportunities,"*IEEE Transactions on Data Engineering, Vol.32, No.1, March 2009.*

[2]　Y. Xu, P. Kostamaa, and L. Gao. "Integrating Hadoop and Parallel DBMS", *Proceedings of ACM SIGMOD, International conference on Data Management, New york, NY, USA 2010.*

[3]　Farah Habib Chan chary,"Data Migration: Connecting databases in the cloud" *IEEE JOURNAL ON COMPUTER SCIENCE, vol no-40, page no-450-455, MARCH 2012.*

[4]　Kamil Bajda at el."Efficient processing of data warehousing queries in a split execution environment", *JOURNAL ON DATA WAREHOUSING, vol no-35, ACM, JUNE 2011.*

[5]　Savvas, A. (2011, October 25). IBM: Businesses unable to analyze 90 percent of their data. ComputerworldUK. Retrieved August 9, 2012.

[6]   Pettey, C., & Goasduff, L. (2011, June 27). Gartner Says Solving "Big Data" Challenge Involves More Than Just Managing Volumes of Data. Stamford: Gartner.

[7]   Russom, P. (2011). Big Data Analytics. TDWI Research.

[8]   FUJIMAKI Ryohei, MORINAGA Satoshi**,"** The Most Advanced Data Mining of the Big Data Era", *NEC TECHNICAL JOURNAL Vol.7 No.2/2012FLEXChip Signal Processor (MC68175/D)*, Motorola, 1996.

[9]   Pavlo, A., Paulson, E., Rasin, A., Abadi, D.J., DeWitt, D.J., Madden, S., Stonebraker, "A Comparison of Approaches to Large-Scale Data Analysis", In: SIGMOD '09.

[10]  Das et al., 2010, Das S., Sismanis Y., Beyer K.S., Gemulla R., Haas P.J., McPherson J., Ricardo: Integrating R and Hadoop, *In: Proceedings of the 2010 ACM SIGMOD International Conference on Management of data (SIGMOD '10), 2010, pp. 987-998*

[11]  IBM 2012, what is big data: Bring big data to the enterprise, http://www-01.ibm.com/software/data/bigdata/, IBM.

[12]  Labrinidis and Jagadish 2012, A. Labrinidis and H. Jagadish, Challenges and Opportunities with Big Data, In *Proc. of the VLDB Endowment, 5(12):2032-2033, 2012*

[13]  Wu X. and Zhu X. 2008, Mining with Noise Knowledge: Error-Aware Data Mining, *IEEE Transactions on Systems, Man and Cybernetics, Part A, vol.38, no.4, pp.917-932.*

[14]  Badrul Chowdhury, Tilmann Rabl, Pooya Saadatpanah," A BigBench Implementation in the Hadoop Ecosystem" available at http://msrg.org

[15]  Wei Fan, Albert Bifet," Mining Big Data: Current Status, and Forecast to the Future" *SIGKDD Explorations Volume 14, Issue 2*

[16]  Kalyani M Raval**, "** Data Mining Techniques", *International Journal of Advanced Research in Computer Science and Software Engineering  Volume 2, Issue 10, October 2012*

[17]  Jim (Zhanwen) Li1, Siyuan He2, Liming Zhu1,3,Xiwei Xu," Challenges to Error Diagnosis in Hadoop Ecosystems"**,**School of Computer Science and Engineering, University of New South Wales, Sydney, Australia