

Report

by Deeja Chhabra

Submission date: 17- May- 2018 10:24AM (UTC+0530)

Submission ID: 964924812

File name: Download_File.pdf (1.12M)

Word count: 5639

Character count: 29450

Movie Review Analysis using Hive

Submitted in partial fulfilment of the requirements of the degree of

Bachelor of Technology

by

Soinaya Rao (140240)

Deeja Chhabra(140078)

under the supervision of

Dr. Anil Kumar Dahiya



A Leading Women's University

COMPUTER SCIENCE AND ENGINEERING
SCHOOL OF ENGINEERING AND TECHNOLOGY

MODY UNIVERSITY OF SCIENCE AND TECHNOLOGY, LAKSHMANGARH

May.2018



Certificate

This is to certify that the report entitled Mot ie Rei ies Analysis using Hit e has been undertaken and » ritten under my superb ision and it describes the original research iv carried out by Deeja Chhabra (140078). Somaja Rao(140240). in Computer Science and Engineering for the dcgrcc of B.Tcch. To the best of my knowledge and belief, this › ork is original and has not been submitted elsewhere for any degree from any other institution in India or abroad.

Dr. Anil Kumar Dahiya
Assistant Profccsor

Dr. Anil Kumar Dahija
Hcad of Dcpartmcnt

Approx al Certil'icate

This report entitled Movie Reii» Analysis using hive by Soinaja Rao, Deeja Chhabra and is approved for the degree of Bachelor of Technology.

Examiner(s)

Super isor(s)

Date _____

Declaration

I declare that this written submission represents my ideas in my own words and where others' ideas or words have been included, I have adequately cited and referenced the original sources. I also declare that I have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea / data / fact / source in my submission. I understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or to whom proper permission has been taken when needed.

I also attest that this Business Director report is free of plagiarism of any kind and that I am fully aware of the Plagiarism Prevention Policy of Mody University of Science and Technology.

Somaya

Rao(140240)

Deeja

Chhabr(14007

8)

Date:

07.05.2018

ACKNOWLEDGEMENT

We would like to express special thanks of gratitude to the Dr. V.K Jain (**Dean CET**) and **Dr.Anil Kumar Dahiy'a** (HOD CSE) for introducing Major Project in our curriculum which helped us in learning hands on certain good technologies.

We would also like to thank our project Supervisor Dr. Anil Kumar Dahiy'a for his support and guidance throughout the schedule of the project. This project could not have been completed without his help.

Thank you all.

Deeja Chhabra(140078)

Somaya Rao(140240)

ABSTRACT

Big data usually includes data sets with sizes beyond the ability of commonly used software tools to capture, curate, manage, and process data within a tolerable elapsed time

Hadoop is an open source, Java-based programming framework that supports the processing and storage of extremely large data sets in a distributed computing environment. It is part of the Apache project sponsored by the Apache Software Foundation.

Hadoop quickly emerged as a foundation for big data processing tasks, such as scientific analytics, business and sales planning, and processing enormous volumes of sensor data, including from internet of things sensors.

Hive is a data warehouse infrastructure tool to process structured data in Hadoop. It resides on top of Hadoop to summarize Big Data, and makes querying and analyzing easy.

It provides a simple query language called HiveQL, which is based on SQL and which enables users familiar with SQL to easily perform ad-hoc querying, summarization and data analysis

With every releasing model whether it is hit or flop, by generating revenues to the audience we are here to transform the voluminous raw data into some meaningful information by applying Big Data Analytics

CONTENTS

Chapter i: Introduction	1
Chapter 2: Literature Review	3
Chapter 3: Research Gap	5
3.1 Map Reducing	5
3.2 Hive	5
Chapter 4 Design and implementation	7
4.1 Hive	7
4.2 Major components of Hive	8
4.3 Operation	9
4.4 Security	10
4.5 HSQL	11
4.6 Research Questions	11
Chapter 5 : System Requirements	12
5.1 Hadoop	12
5.1.2 Hive	13
5.2 Hive Job Execution Flow'	14

Chapter 6: Experimental Results	17
6.1 Movie dataset.	17
6.2 Research Questions .	18
Chapter 7 : Conclusion and Future Work	22
7.1 Conclusion	22
7.2 Future Work.	22
References,.	25

LIST OF FIGURES

FIG NO.	TITLE OF FIGURE	PAGE NO.
4.1	Components of Hit e	7
5.1	Hit c Job Exccutioii Floor	14
6.1	Total number of Mo ics	1 X
6.2	Maxim um rating of mo> ics	18
6.3	Count number of inaximu in ratings	19
6.4. 1	Mot ies having rating 1 & 2	20
6.4.2	Mo res ha ing rating I & 2	20
6.5.1	List of mod ies	21
6.5.2	List of mod ies	21
6.6	Mot ies u ith duration of 2 hours	22

LIST OF TABLE

TABLE NO,	TITLE OF TABLE	PAGE NO.
4.I	Components of Hit e	7

CHAPTER 1

Introduction

1.1 Introduction to bigdata

Big Data is a term that refers to dataset » hose size or i olume. complexity.i arietj ,rate of growth or eracity of data » hich organisations handled hat e reached such unbelief able hi'cl thnt traditional processing and analytical tools failed to process. A data ii hich is beyond to the storage capacity and »'liich is bejoiid to the processing po» er.that data » e can call as Big Data.

Or er the last fe» years, there has been an incredible explosion of data IBM reported that 2.5 billion gigab; tes of data ii as gncrated ci en' day in 2012 Data is grow ing faster than ci'cr before then one solution » as » e can use iiiiultiple computers called distributed systems. but in this there are high chances of sj stein failure. programming complexity is also high as it is difficult to synchronizc daia and process . then comes the another solution known as Hadoop.

1.2 Hadoop:Solution for Big data

Hadoop is n fraincii ork that allow s for distribntcd processing of large daia scts across clusters of commodity computers using simple programming models It is inspired bj technical document published by Google. The » ord 'hadoop does not liar e aiij meaning long Cutting discoi'crd Hadoop and nained it after his son's yelloii-colored toy clephnnt

In recent years, a great proliferation has been witnessed in the amount of data generated. This rote is growing os datn is continuously growing amassing such huge omount of dota is a tough job end thus it needs some out of the box thinking os it cannot be tackled with traditional tools and tcchniques . This inadequacy led to the birth of the term Big Datn and olong with it the challenges such as storage. processing. visualization and privacy. Big Data is not just about being big in size.

The definition is broadened using five characteristics or

—V'sl, These ore:

- Volume : This characteristic signifies huge voluminous data: it is in orders of terabytes and even petabytes.
- Velocity . This characteristic signifies the high velocity with which the data is generated.
- Variety : This characteristic refers to the huge variety in the big data.
- Value: This characteristic refers to the intrinsic value contained in big data.
- Veracity This characteristic refers to uncertainties in big data such as missing, duplicate and incomplete entries,

Hadoop is an excellent and robust analytics platform for Big Data which can process huge data sets at a really quick speed by providing scalability. It can manage all aspects of Big Data such as volume, velocity and variety by storing and processing the data over a cluster of nodes . Hadoop has two major components in its architecture that is MapReduce and the Hadoop Distributed File System (HDFS). With the introduction of YARN (Yet Another Resource Negotiator) in later releases, Hadoop was integrated with a number of wonderful components which can be used for storing, processing and analyzing data a lot more efficiently, thus aiding in exploration of data for undiscovered facts at a smooth pace. Some of these components that work on top of Hadoop are Hite.Flume, Sqoop, Hbase and Oozie.

1.3 Motivation

In today's world, every establishment is facing ever growing challenges which need to be coped up quickly and efficiently. With continually increasing Mobile industry and entertainment industry it wants to enjoy all the perks of facility. So if they invest on mobile resources that they can have received with rating and all the information related to that mobile.

The best place to look up to find room for improvement is the voluminous raw data that is generated on a regular basis from various sources by applying Big Data Analytics (BDA)[2]. BDA refers to the tools and practices that can be used for transforming the raw data into meaningful and crucial information. The main goal with every releasing mobile is whether it is

hit or flop. by generating reviews to the audience i.e. are here to transform the ioluminous ran
data into sortie meaningful information by applying Big Data Analytics

CHAPTER 2

Literature Review

Vidyasagar S.D[2017] did a survey on Big Data and Hadoop system and found that organizations need to process and handle petabytes of Data in an efficient and inexpensive manner. According to him if there is any node failure then we can lose some information. Hadoop is an Efficient, reliable. Open Source Apache License. Hadoop is used to deal with large data sets. Author explained its need, uses and application. Now days, Hadoop is playing an important role in Big Data. Vidyasagar S.D concluded that Hadoop is designed to run on cheap commodity hardware, it automatically handles data replication and node failure, it does the hard work — you can focus on processing data, Cost Saving and efficient and reliable data processing.

Sujatha .Va , Prasanna Devi Sb , Vinu Kiran Sb ,ManivannanS[2016] had analyzed a large scale Diabetic data sets for several patients to find the length of time taken for treatment for each class of Diabetes and the risk of re-admission of diabetic patients performing Bigdata analytics, the type of diabetes and its outcome is highly acted as a high risk sample of patient data sets. They have collected and integrated different sources of diabetic information for several patients, from primary and secondary treatment information to administrative information, to analyze the effect of patient care processes such as type of treatments and each patient's behavior on which results multifaceted nature of chronic care that they take into their account to predict the significant factors and length of stay. Nowadays by using electronic medical equipments with high quality and high degree calibrations. they are able to gather large amounts of real-time diabetic data sets That requires the usage of distributed platforms for making Big Data analysis that results on making decisions based on available data and its trends This type of Bigdata analysis will's geographical and environmental information of patients' enables the capability of

interpreting the ethnicity of data gathered and extract new analysis to identify survival options and treatment timelines from them.

Sergey V. Kovalchuk , **Artem V. Itharchuk** , **Jiaqi Liao** , **Sergey V. Ivanov** , **Alexander V. Boukhanovsky** [2016] presents a technology for dynamic knowledge-based building of Domain-Specific Languages (DSL) to describe data-intensive scientific discovery tasks using BigData technology. Their proposed technology supports high level abstract definition of analytic and simulation parts of the task as well as integration into the composite scientific solutions. Automatic translation of the abstract task definition enables seamless integration of various data sources within single solution.

Vennila S and **Priyadarshini I** [2015] Cloud computing provides flexible infrastructure and high storage capacity for BigData applications. The MapReduce framework is most preferable for processing huge volume of unstructured data set in BigData. Increase in data volume leads to flexible and scalable privacy preservation of such dataset over the MapReduce framework in BigData applications. A survey has been taken for the MapReduce framework based big data privacy preservation in Cloud environment.

Jeff Sedayaoet. AI [2014] suggested to use Hadoop to analyze the data and obtain useful results for the Human Factors analysts. At the same time, the requirements of analysing were learned and anonymized data sets need to be carefully analyzed to determine whether they are vulnerable to attack.

Meiko Jensen et. AI [2014] explained that the field of privacy in big data contexts contains a bunch of key challenges that must be addressed by research. Many of these challenges do not stem from technical issues, but merely are based on legislation and organizational matters. Nevertheless, it can be anticipated that it is as feasible to meet each of the challenges discussed here by means of appropriate technical measures.

B. Saralailevia , **N. Pazzh rajaa** , **P. Victor Paula** , **M.S. Saleem Bashab** , **P. Dhavachelvanc** [2013] Their paper shows the big data information and characteristics used in world wide. The issues are also mentioned to give idea about the big data issues in real time

The security issue is pointed more in order to increase the security in big data. We can improve security in big data by using any one of the approach or by combining these three approaches in Hadoop Distributed File System which is the base layer in Hadoop, where it contains large number of blocks

Mohd Rehan Ghazia and **Durgaprasad** Gangodkara[2018] discuss Hadoop and its components in detail which comprise of MapReduce and Hadoop Distributed File System (HDFS). MapReduce engine uses JobTracker and TaskTracker that handle monitoring and execution of job. HDFS a distributed file- system which comprise of NameNode, DataNode and Secondary NameNode for efficient handling of distributed storage purpose. The details provided can be used for developing large scale distributed applications that can exploit computational power of multiple nodes for data and compute intensive applications.

CHAPTER 3

Research Gap

20

3.1 Map Reduce

MapReduce works by breaking the processing into two phases: the map phase and the reduce phase. Each phase has key-value pairs as input and **output**, the types of which may be chosen by the programmer. The programmer also specifies two functions: the map function and the reduce function.

Hadoop MapReduce is a software framework for easily writing applications which process vast amounts of data (multi-terabyte data-sets) in-parallel on large clusters (thousands of nodes) of commodity hardware in a reliable, fault-tolerant manner.

A MapReduce job usually splits the input data-set into independent chunks which are processed by the map tasks in a completely parallel manner. The framework sorts the outputs of the maps, which are then input to the reduce tasks. Typically both the input and the output of the job are stored in a file-system. The framework takes care of scheduling tasks, monitoring them and re-executes the failed tasks.

Typically the compute nodes and the storage nodes are the same, that is, the MapReduce framework and the Hadoop Distributed File System are running on the same set of nodes. This configuration allows the framework to effectively schedule tasks on the nodes where data is already present, resulting in very high aggregate bandwidth across the cluster.

The MapReduce framework consists of a single master JobTracker and one slave TaskTracker per cluster-node. The master is responsible for scheduling the jobs' component tasks on the slaves, monitoring them and re-executing the failed tasks. The slaves execute the tasks as directed by the master.

3.2 Hive

One criticism of MapReduce is that the development cycle is very long. MapReduce allows you, as the programmer, to specify a map function followed by a reduce function, but working out how to fit your data processing into this pattern, which often requires multiple MapReduce stages, can be a challenge. Writing the mappers and reducers, compiling and packaging the code, submitting the job(s), and retrieving the results is a time-consuming business, and even with Streaming, which removes the compile and package step, the experience is still involved.

Hive raises the level of abstraction for processing large datasets. With Pig, the data structures are much richer, typically being multivalued and nested, and the set of transformations you can apply to the data are much more powerful. They include joins, for example, which are not for the faint of heart in MapReduce. Hive is made up of two pieces:

- The language used to express data flows, called HiveQL.
- The execution environment to run HiveQL programs. There are currently two environments: local execution in a single JVM and distributed execution on a Hadoop cluster.

CHAPTER 4

Design and Implementation

4.1 Hive

Apache Hive is a data warehouse software project built on top of Apache Hadoop for processing data summarization, query, and analysis. Hive gives an SQL-like interface to query data stored in various databases and file systems that integrate with Hadoop. Traditional SQL queries must be implemented in the Map reduce low API to execute SQL applications and queries over distributed data. Hive provides the necessary SQL abstraction to integrate SQL-like queries (HiveQL) into the underlying low without the need to implement queries in the low-level Java API. Since most data warehousing applications work with SQL-based querying languages, Hive aids portability of SQL-based applications to Hadoop

4.2 Major components of hive

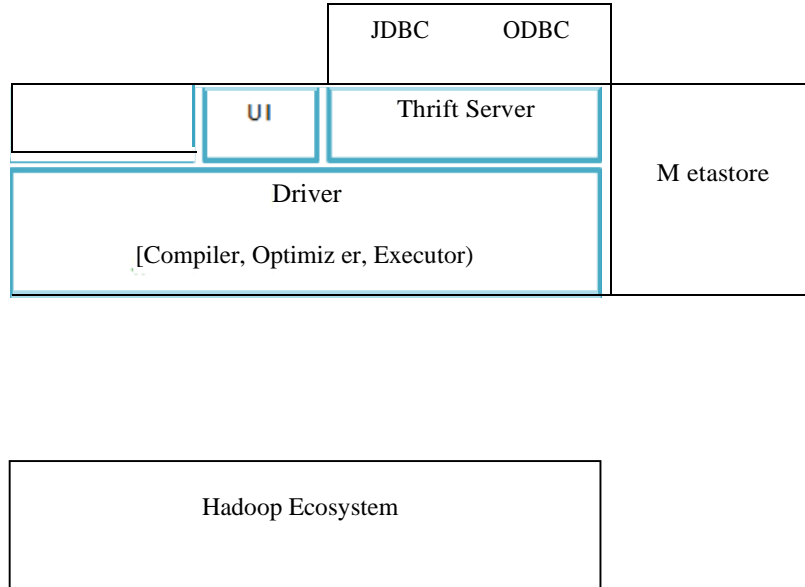


Fig.4.1 Components of Hive

- **Metastore:** Stores metadata for each of the tables such as their schema and location. It also includes the partition metadata which helps the driver to track the progress of various data sets distributed over the cluster. The data is stored in a traditional RDBMS format. The metadata helps the driver to keep a track of the data and it is highly crucial. Hence, a backup server regularly replicates the data which can be retrieved in case of data loss.
- **Driver:** Acts like a controller which receives the HiveQL statements. It starts the execution of statement by creating sessions and monitors the life cycle and progress of the execution. It stores the necessary metadata generated during the execution of an HiveQL statement. The driver also acts as a collection point of data or query result obtained after the Reduce operation.
- **Compiler:** Performs compilation of the HiveQL query, which converts the query to an execution plan. This plan contains the tasks and steps needed to be performed by the Hadoop MapReduce to get the output as translated by the query. The compiler converts the query to an abstract syntax tree (AST). After checking for compatibility and compile time errors, it converts the AST to a directed acyclic graph (DAG). The DAG divides operators to MapReduce stages and tasks based on the input query and data.
- **Optimizer:** Performs various transformations on the execution plan to get an optimized DAG. Transformations can be aggregated together, such as converting a pipeline of joins to a single join, for better performance. It can also split the tasks, such as applying a transformation on data before a reduce operation, to provide better performance and scalability. However, the logic of transformation used for optimization can be modified or pipelined using another optimizer.
- **Executor:** After compilation and optimization, the executor executes the tasks. It interacts with the job tracker of Hadoop to schedule tasks to be run. It takes care of pipelining the tasks by making sure that a task with dependency gets executed only if all other prerequisites are
- **CLI, UI, and Thrift Server:** A command-line interface (CLI) provides a user interface for an external user to interact with Hive by submitting queries, instructions and monitoring the

process status. Thrift server allows external clients to interact with Hive server framework. similar to the JDBC or ODBC protocols.

4.3 Operations

Table 4.1 Operations in Hive

Step No.	Operation
	Execute Query
	The Hive interface such as Command Line or Web UI sends query to Driver (any database driver such as JDBC, ODBC, etc.) to execute
2	Get Plan
	The driver takes the help of query compiler that parses the query to check the syntax and query plan or the requirement of query.
3	Get Metadata
	The compiler sends metadata request to Metastore (any database).
4	Send Metadata
	Metastore sends metadata as a response to the compiler.
	Send Plan
	The compiler checks the requirement and resends the plan to the driver. Up to

here, the parsing and compiling of a query is complete

6 Execute Plan

The driver sends the execute plan to the execution engine.

7 **Execute Job**

Internally, the process of execution job is a MapReduce job. The execution engine sends the job to JobTracker, which is in Name node and it assigns this job to TaskTracker, which is in Data node. Here, the query executes MapReduce job,

7.1 Metadata Ops

Meanwhile in execution, the execution engine can execute metadata operations with Metastore.

8 Pseudo Result

The execution engine receives the results from Data nodes.

9 Scud Results

The execution engine sends those resultant values to the driver.

10 Scud Results

The driver sends the results to Hadoop Interfaces

4.4 Security

Hiveⁿ 0.12.0 added integration with Hadoop security. Hadoop began rising Kerberos authorization support to provide security. Kerberos allows for mutual authentication between client and server. In this instance, the client's request for a ticket is passed along with the request. The previous versions of Hadoop had several issues such as users being able to spoof their username by setting the `hadoop.job.ugi` property and also MapReduce operations being run under the same user: `hadoop` or `mapred`. With Hive 0.12.0's integration with Hadoop security, these issues have largely been fixed. Tez jobs are run by the user who launched it and the username can no longer be spoofed by setting the `hadoop.job.ugi` property. Permissions for newly created files in Hive are dictated by the HDFS. The Hadoop distributed filesystem authorization model uses three entities: user, group and others with three permissions: read, write and execute. The default permissions for newly created files can be set by changing the `umask` value for the Hive configuration variable `hive.files.umask.value`.

4.5 HIVEQL

While based on SQL, HiveQL does not strictly follow the full SQL-92 standard. HiveQL offers extensions not in SQL, including multi-table inserts and create table as select, but only offers basic support for indexes. HiveQL lacked support for transactions and materialized views, and only limited subquery support. Support for insert, update, and delete with full ANSI functionality was made available with release 0.14.

Internally, a compiler translates HiveQL statements into a directed acyclic graph of MapReduce, Tez, or Spark jobs, which are submitted to Hadoop for execution.

4.6 Research Questions

1. Total Number of movies in 2017

It will give the total number of count of movies released in the year 2017

2. Finding maximum rating of the movie

It will show the movie having maximum rating.

3. Count the number of movies having the maximum rating.

It will display the count .

4. Count number of movies between rating 1 and 2

It will display the list of movies having least ratings

5. Find the list of stars and number of movies released each year.

It will show the details of movies.

6 Find number of movies with duration of 2 hours.

It will show the duration of all the movies.

Chapter 5

Technology Used

5.1 System Requirements: Software used: Hadoop, Hive

5.1.1 Hadoop:

It is a Java based programming framework which stores data and processes on it and distributes it in columnar form and runs on the basis of clusters with the commodity hardware.

The first component to provide online access is HBase, a key-value store that relies on HDFS for its underlying storage. HBase provides both online read/write access of individual rows and batch operations for reading and writing data in bulk, making it a good solution for building applications on. The real enabler for new processing models in Hadoop was the introduction of YARN (which stands for Yet Another Resource Negotiator) in Hadoop 2.

YARN is a cluster resource management system, which allows any distributed program (not just MapReduce) to run on data in a Hadoop cluster. In the last few years, there has been a flowering of different processing patterns that work with Hadoop. Here is a sample: Interactive SQL. By dispensing with MapReduce and using a distributed query engine that uses dedicated query engines (like Impala) or container reuse (like Hive on Tez), it's possible to achieve low-latency responses for SQL queries on Hadoop while still scaling up to large dataset sizes. Iterative processing. Many workflows—such as those in machine learning—are iterative in nature, so it's much more efficient to hold each intermediate working set in memory, compared to loading from disk on each iteration. The architecture of MapReduce does not allow this, but it's straightforward with Spark, for example, and it enables a highly exploratory style of working with datasets. Stream processing. Streaming systems like Storm, Spark Streaming, and Flink make it possible to run real-time, distributed computations on unbounded streams of data and commit results to Hadoop storage or

external systems. Search The Solr search platform can run on a Hadoop cluster, indexing documents as they are added to HDFS, and serving search queries from indexes stored in HDFS

5.1.2 Hive

Hive is a data warehousing infrastructure tool to process structured data in Hadoop. It resides on top of Hadoop to summarize Big Data, and makes query writing and analyzing easy.

Initially Hive was developed by Facebook, later the Apache Software Foundation took it up and developed it further as an open source under the name Apache Hite. It is used by different companies. For example, **Amazon uses it in Amazon Elastic MapReduce.**

Hive is not

- A relational database
- A design for OnLine Transaction Processing (OLTP)
- A language for real-time queries and row-level updates

Features of Hive

- It stores schema in a database and processed data into HDFS.
- It is designed for OLAP.
- It provides SQL type language for querying called HiteQL or HiteQL.
- It is familiar, fast, scalable, and extensible.

5.2 Hive Job Execution Flow

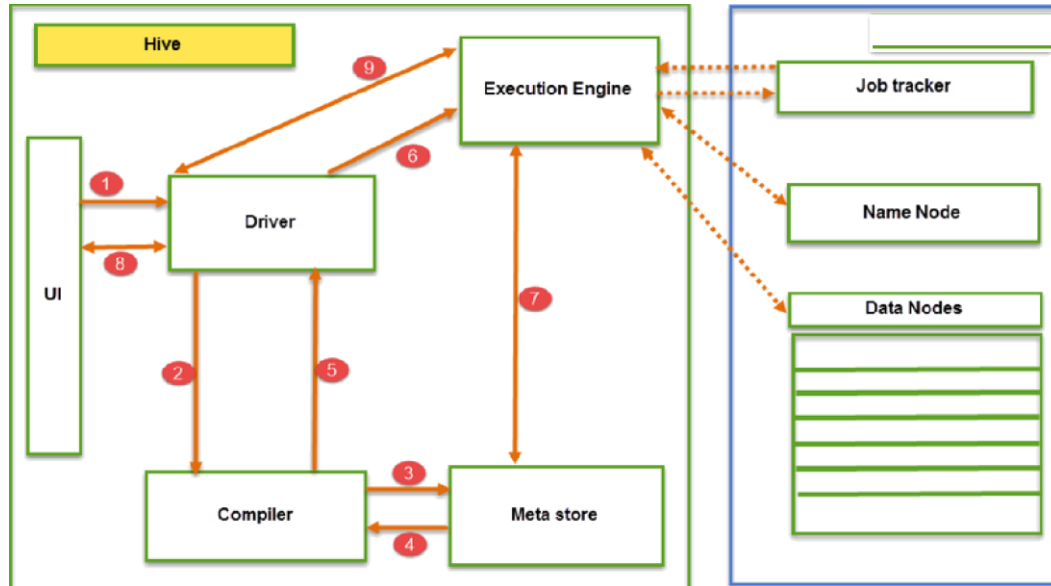


Fig 5.1 Hive Job Execution Flow

From the above screenshot we can understand the Job execution flow in Hive with Hadoop

The data flow in Hive behaves in the following pattern:

1. Executing Query from the UI (User Interface)
2. The driver is interacting with Compiler for getting the plan. (Here plan refers to query execution) process and its related metadata information gathering
3. The compiler creates the plan for a job to be executed. Compiler communicating with Meta store for getting metadata request
4. Meta store sends metadata information back to compiler
5. Compiler communicating with Driver with the proposed plan to execute the query
6. Driver Sending execution plans to Execution engine
7. Execution Engine (EE) acts as a bridge between Hive and Hadoop to process the query. For DFS operations.

- EE should first contacts Name Node .and then to Data nodes to ,get the values stored in tables.
- EE is going to fetch desired records from Data Nodes. The .actual.data of tables resides iii.data node only. While from Name Node it only fetches the metadata information for the query.
- It collects aotual .data from data nodes related to mentioned .query
- Execution Engine (EE) communicates bi-directioniudly with Meta stnre present in Hive tn perform DDL (Data Definition Language) operafions. Here DDL operafions like KREATE, DROP atid ALTERING tables and databases are!dons. Meta store will store inforñiation about .database name, table. names and column names only. It will fetch data related to query mendoned.
- Execution Engine (EE) in turn .communicates with Haitoop <daemons such as Name node,, Data nodes, and job oacker tn execute the query on top of Hadoop file sysUni

8. Fetching suits from driver

9. sending re.suits to Executi.on engin.e., Once the results fetched from .data n.odes to the EE, it will seiid results.back to dfiver arid to UI (froiit eiid)

! i@ Contioufiusly inc.oust with itado.op ifile sytem and its daemons v|a Execution engine. The .dotte<i artow in tire Job flow <diagram sliows. the Execution eiigine coinmuiiicationii with Hadoop.daemons.

- Different modes tif Hivñ

Hive can .opeiate iii two modes. <iependiig oii tire size of<iata iiodes iii Hadoop.

These inodes. are,

- Local mode
- Map. reñuce mode

When to use! Local mode:

- If the Hadoop installed under pseudo mode with having one data node we use Hive in this mode
- If the data size is smaller in term of limited to single local machine, we can use this mode
- Processing will be very fast on smaller data sets present in the local machine

When to use Map reduce mode:

- If Hadoop is having multiple data nodes and data is distributed across different nodes we use Hive in this mode
- It will help on large amount of data sets and query going to execute in parallel way
- Processing of large data sets with better performance can be achieved through this

In Hive, we can set this property to mention which mode Hive can work? By default, it works on Map Reduce mode and for local mode you can have the following setting.

Hive to work in local mode set

SET mapred.job.tracker=local;

From the Hive version 0.7 it supports a note to run map reduce jobs in local mode automatically

Chapter 6

Experimental Results

6.1 Movie Data Set

1,The Nightmare Before Christmas,1993,3.9,4568
2,The Mummy,1932,3.5,4388
3,Orphans of the Storm,1921,3.2,9062
4,The Object of Beauty,1991,2.8,6150
5,Night Tide,1963,2.8,5126
6,One Magic Christmas,1985,3.8,5333
7,Muriel's Wedding,1994,3.5,6323
8,Mother's Boys,1994,3.4,5733
9,Nosferatu: Original Version,1929,3.5,5651
10,Nick of Time,1995,3.4,5333
11,Broken Blossoms,1919,3.3,5367
12,Big Night,1996,3.6,6561
13,The Birth of a Nation,1915,2.9,12118
14,The Boys from Brazil,1978,3.6,7417
15,BLJ DO ! NOT SP,1971,).9,FA05
16,The Breakfast Club,1985,4.0,5823
17,The Bride of Frankenstein,1935,3.7,4485
18,Beautiful Girls,1996,3.5,6755
19,Bustin' Loose,1981,3.7,5598
20,The Beguiled,1971,3.4,6307
21,Born on the Fourth of July,1989,3.4,8646
22,Broadcast News,1987,3.4,7940
23,Swimming with Sharks,1994,3.3,5586
24,Beavis and Butt-head Do America,1996,3.4,4852
25,Brighton Beach Memoirs,1986,3.4,6564
26,The Best of Times,1986,3.4,6247
27,Brassed Off,1996,3.5,6040
28,Last Tango in Paris,1972,3.1,7732
29,Leprechaun 2,1994,3.2,5125
30,Incident at Oglala: The Leonard Peltier Story,1992,3.7,5487
31,Kalifornia,1993,3.4,7095

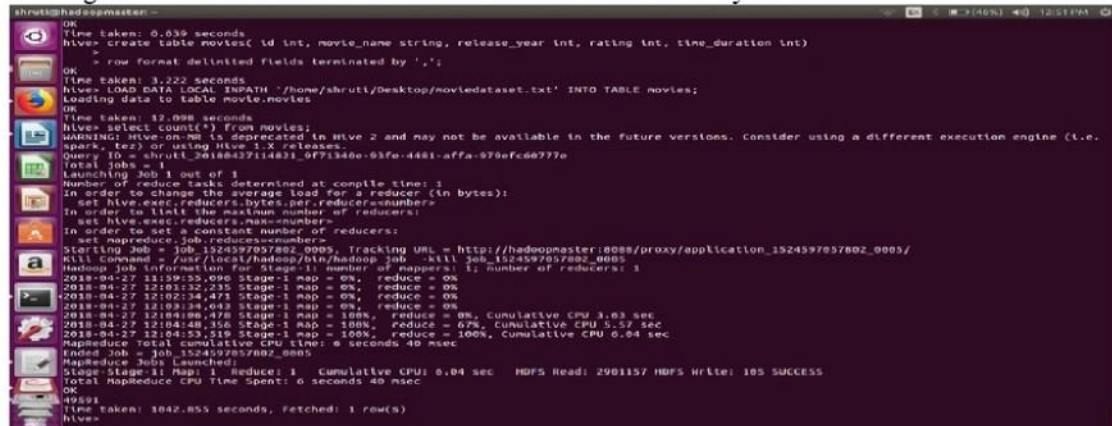
33,Jingle All the Way,1996,3.6,5371
34,Killing Zoe,1993,3.4,5773
35,King of Beggars,1992,3.6,6025
36,Into the Woods,1990,4.0,9077

38jO TOO D(ep,t9%,).9j8i)
39,CtA)XOO,t070,3,2,BO44
40 Internal Affairs 1990 3.5 6885

6.2 Research Questions

1. Total Number of movies in 2017

It will give the total number of count of movies released in the year 2017.

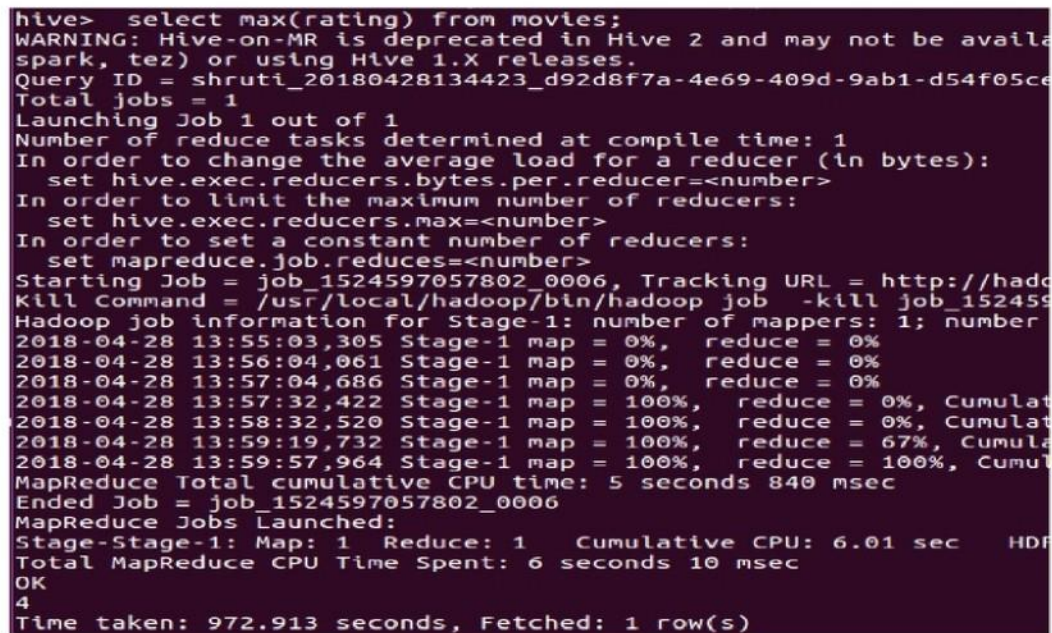


```
shrut@hadoopmaster:~$
OK
Time taken: 0.039 seconds
hive> create table movies( id int, movie_name string, release_year int, rating int, time_duration int)
> row format delimited fields terminated by ',';
OK
Time taken: 3.222 seconds
hive> load data local INPATH '/home/shrut/Desktop/moviedataset.txt' INTO TABLE movies;
Loading data to table movie.movies
OK
Time taken: 12.098 seconds
hive> select count(*) from movies;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
Query ID = shruti_20180427114831_9f71340e-93fe-4481-af6a-979efcd0777e
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1524597057802_0005, Tracking URL = http://hadoopmaster:8088/proxy/application_1524597057802_0005/
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_1524597057802_0005
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2018-04-27 12:00:35,090 Stage-1 map = 0%, reduce = 0%
2018-04-27 12:01:32,235 Stage-1 map = 0%, reduce = 0%
2018-04-27 12:02:24,673 Stage-1 map = 0%, reduce = 0%
2018-04-27 12:03:24,043 Stage-1 map = 0%, reduce = 0%
2018-04-27 12:04:08,476 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 3.81 sec
2018-04-27 12:04:49,316 Stage-1 map = 100%, reduce = 67%, Cumulative CPU 5.37 sec
2018-04-27 12:04:53,519 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 6.04 sec
MapReduce total cumulative CPU time: 6 seconds 40 msec
Ended Job = job_1524597057802_0005
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 6.04 sec HDFS Read: 290137 HDFS Write: 105 SUCCESS
Total MapReduce CPU Time Spent: 6 seconds 40 msec
OK
40591
Time taken: 1042.855 seconds, Fetched: 1 row(s)
hive>
```

Fig 6.1total number of movies

2. Finding maximum rating of the movie

It will show the movie having maximum rating.



```
hive> select max(rating) from movies;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
Query ID = shruti_20180428134423_d92d8f7a-4e69-409d-9ab1-d54f05ce
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1524597057802_0006, Tracking URL = http://hadoopmaster:8088/proxy/application_1524597057802_0006/
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_1524597057802_0006
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2018-04-28 13:55:03,305 Stage-1 map = 0%, reduce = 0%
2018-04-28 13:56:04,061 Stage-1 map = 0%, reduce = 0%
2018-04-28 13:57:04,686 Stage-1 map = 0%, reduce = 0%
2018-04-28 13:57:32,422 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 3.81 sec
2018-04-28 13:58:32,520 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 5.37 sec
2018-04-28 13:59:19,732 Stage-1 map = 100%, reduce = 67%, Cumulative CPU 6.04 sec
2018-04-28 13:59:57,964 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 6.04 sec
MapReduce total cumulative CPU time: 5 seconds 840 msec
Ended Job = job_1524597057802_0006
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 6.01 sec HDFS Read: 290137 HDFS Write: 105 SUCCESS
Total MapReduce CPU Time Spent: 6 seconds 10 msec
OK
4
Time taken: 972.913 seconds, Fetched: 1 row(s)
```

Fig 6.2 Maximum rating of movies

3. Count the number of movies having the maximum rating.

It will display the count.

```
hive>
> select count(*) from movies where rating = 4;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available
spark, tez) or using Hive 1.X releases.
Query ID = shruti_20180428140428_2e799b0e-5073-4cd6-a61e-967abf4eda87
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1524597057802_0007, Tracking URL = http://hadoopma
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_1524597057
Hadoop job information for Stage-1: number of mappers: 1; number of r
2018-04-28 14:07:06,674 Stage-1 map = 0%, reduce = 0%
2018-04-28 14:07:55,470 Stage-1 map = 100%, reduce = 0%, Cumulative
2018-04-28 14:08:36,318 Stage-1 map = 100%, reduce = 100%, Cumulative
MapReduce Total cumulative CPU time: 6 seconds 920 msec
Ended Job = job_1524597057802_0007
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 6.92 sec HDFS Re
Total MapReduce CPU Time Spent: 6 seconds 920 msec
OK
1477
```

Fig 6.3 Count number of maximum rating

4. Count number of movies between rating 1 and 2

It will display the list of movies having least ratings

```
hive> select distinct movie_name from movies where rating between 1 and 2;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions of Hive (you can still use Hive-on-MapReduce with spark, tez) or using Hive 1.X releases.
Query ID = shrutl_20180505120705_8bb2b2fb-6753-4110-acfe-c6706ec0db2f
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1525499465151_0001, Tracking URL = http://hadoopmaster:8088/proxy/aj
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_1525499465151_0001
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2018-05-05 12:07:59,049 Stage-1 map = 0%, reduce = 0%
2018-05-05 12:08:59,120 Stage-1 map = 0%, reduce = 0%, Cumulative CPU 2.91 sec
2018-05-05 12:09:03,854 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 4.43 sec
2018-05-05 12:09:26,030 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 8.96 sec
MapReduce Total cumulative CPU time: 8 seconds 960 msec
Ended Job = job_1525499465151_0001
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 8.96 sec HDFS Read: 2902193 HDFS W
Total MapReduce CPU Time Spent: 8 seconds 960 msec
OK
...At First Sight
```

Fig 6.4.1 Movies having rating 1 and 2

```
.Com for Murder
10 Items or Less
10 Items or Less: Season 1
10 Items or Less: Season 2
10 Items or Less: Season 3
100 Below Zero
100 Years Of Evil
11 11 11
12 12 12
13 13 13
1313: Actor Slash Model
1313: Bermuda Triangle
1313: Bigfoot Island
1313: Billy the Kid
1313: Cougar Cult
1313: Frankenqueen
```

Fig 6A.2 Movies having rating 1 and 2

1. Find the list of years and number of movies released each year.

It will show the details of movies,

```
hive> select release_year, count(distinct movie_name) from movies group by release_year
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions of Hive. You should use Tez or the new Hive release (Hive 3) instead. Please use spark, tez) or using Hive 1.X releases.
Query ID = shruti_20180505121555_48bf360a-0f4b-4957-9d4e-c465376fb454
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1525499465151_0002, Tracking URL = http://hadoopmaster:8080/
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_1525499465151_0002
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2018-05-05 12:16:12,857 Stage-1 map = 0%, reduce = 0%
2018-05-05 12:16:50,548 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 6.8 sec
2018-05-05 12:17:28,422 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 12.5 sec
MapReduce Total cumulative CPU time: 12 seconds 500 msec
Ended Job = job_1525499465151_0002
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 12.5 sec HDFS Read: 29020 B HDFS Write: 0 B
Total MapReduce CPU Time Spent: 12 seconds 500 msec
OK
NULL 0
1913 3
1914 20
1915 1
1916 1
1918 1
1919 3
1920 6
1921 2
1922 2
1923 4
1924 5
1925 5
1926 2
1927 4
1928 2
```

Fig 6.5.1 List of movies

```
1974 178
1975 134
1976 118
1977 136
1978 230
1979 140
1980 106
1981 112
1982 153
1983 267
1984 302
1985 333
1986 287
1987 280
1988 333
1989 419
1990 467
1991 364
1992 479
1993 561
1994 515
1995 590
1996 685
1997 783
1998 838
1999 1173
2000 900
2001 1167
2002 1109
2003 1389
2004 1368
```

Fig 6.5.2 List of movies

6. Find number of movies with duration of 2 hours.

It will show the duration of all the movies.

```
Time taken: 97.7748 seconds, Fetched: 102 row(s)
hive> select count(movie_name) from movies where time_duration = (2*60*60);
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions of Hive. You should use Tez or MapReduce with Hive. You may also use the Hive-on-Presto engine or the Hive-on-Spark engine, or using Hive 1.X releases.
Query ID = shruti_20180505122108_0e0131f3-2d09-4bb5-8d58-116c4e9bb397
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1525499465151_0003, Tracking URL = http://hadoopmaster:8088/proxy
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_1525499465151_0003
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2018-05-05 12:21:25,705 Stage-1 map = 0%, reduce = 0%
2018-05-05 12:21:55,871 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 5.63 sec
2018-05-05 12:22:33,432 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 9.06 sec
MapReduce Total cumulative CPU time: 9 seconds 60 msec
Ended Job = job_1525499465151_0003
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 9.06 sec HDFS Read: 2902201 HDFS Write: 102400
Total MapReduce CPU Time Spent: 9 seconds 60 msec
OK
2
Time taken: 87.463 seconds, Fetched: 1 row(s)
```

Fig 6.6 Movies with Duration of 2 hour

Chapter 7

Conclusion and Future Work

7.1 Conclusion

This is a novel approach where the user rating decisions has been taken to purview along with inherent movie attributes to model the classification approach. An experimental insight has also been provided for the post release aspect of the movie that relates initial budget with each of the financial returns.

Big Data Analytics refers to the tools and practices that can be used for transforming this raw data into meaningful and crucial information which helps in forming a decision support system to make it easier for people to decide for a movie. One should know which movie to watch. Also, with the review it provides us with the duration of the movie and also the release date of the movie.

7.2 Future Work

- The abundance of movie data in terms of review, rating or even detail information (for example the information maintained by (IMDB) in the internet has encouraged many researches to formulate techniques to analyze the pattern in movie data. Most of the researches are devoted to develop recommendation systems of movie according to user review's.
- This analysis can be further carried out on Fully Distributed cluster mode that is hadoop daemons run on a cluster of machine.
- Similar analysis can be further carried out in different sector

REFERENCES

- [1] Apache Hadoop. <https://hadoop.apache.org/>
 - [2] Riyaz P.A., Sureliia Mariain Varghese, Lei'eraging MapReduce ii itli Hadoop for Weather Data Analytics . IOSR - Journal for Computer Science, 2015. i olume 17, issue 3. pp 6-12. [3] <http://spideropsiiet.co.in/site1/blog/2014/09/18/apache-pig-tool-big-data-analytics/>
 - [4] Andrcii Pai'lo. —A Comparison of Approaches to Large-Scale Data Analvsisl. SIGMOD. 200'3.
 - [5] Dean. J, and Gheinau at. S..—MapRediice: a flexible data processing toolll. ACM 2010.
 - [6] Greenpliin Analytics Workbench .i'sit us at u'u'u.greenpluin.co.in
 - [7] HadoopTutorial: <http://ldecs.elper.yahoo.co.in/liadoop/tutorial/module1.html>
 - [8] Harslia»ardhan S Bliosale. Prof. De endra P Gadekar —A Revie» Paper on Big Data and Hadoopll HUB,i » » .ibinbigdatahub.com/infograpliic/four- **s-big-datn**.
 - [10] Sanjay Rathee. —Big Data and Hadoop ii ith components like Fliime. Pig. Hii'c and Jaqll,Inteniatinal Conference on Cloud,BigData and Tnist2013.Vol. 1.5,Noi ember 20
 - [11]Narkhede. S..& Baraskar. T. —HMR log anal zer: anale ze Webapplication logs os'er Hadoop M;ipRcdiicc11. International Journal ofUbiComp. pp 41-51, 2013.
 - {12}Hadoop — Adi'antages and Disadi'antages <http://uuu.j2eebrain.com/jawa-J2ee-hadoop-advantages- and-disadi outages.html> Accessed on: 13-04-2017.
 - [13]Shilpa. Manjit Kaur. —BIG Data and Methodology-A rei iewl. International Journal of Adi'anccd Research inComputcr Science and Softii'arc Enginccring. Volume 3, Issue 1G. October 2013 [14] <http://ii'ii'ii hitehoiise got'/blog/2012/03/29/big-data-big-deal/Machine learning>.
 - [15] Ahmed Elda» y. Mohamed F. Mokbel — A Demonstration of Spatial Hadoop An Efficient MapReduce Frames ork for Spatial DatnIProceedings of the VLDB Endou'ment. Vol 6. No. 12
- Copyright 2013 VLDB Endou itient 215050
- [1 fi] MrigaiA Mridul . Akaslideep Khajuria. Sneliasish Dutta. Kumar N — Anal'sis of Bidgata using
 - Apache Hadoop and Map Reducel Volume 4. Issue 5. Maj 201411.27
 - {17} P D Londhe. S.S Kiimbhar. R.S Sul. and A J. Khadse. —Processing big data using liadoop fraitiieu orkll. in Proceedings of 4tli SARC-IRF International Conference. Neu Delhi. India. Apr. 27. 2014. pp 72-75

11 T. White. —MapReduce and the hadoop distributed file system, in Hadoop: The definitive guide. 1st edition. O'Reilly Media. Inc . Yahoo press. 2012

11 D. Borthakui. "The hadoop distributed file system architecture and design," Hadoop Project Website [online] Available: <http://hadoop.apache.org/core/docs/current/hdfs/design.pdf>

12 J. L. & C. Djer. "Data-intensive text processing with mapreduce". in Synthesis Lectures on Human Language Technologies. vol. 3. no. 1. pp. **1-177**. 2010.

Report

ORIGINALITY REPORT

69%

SIMILARITY INDEX

58%

INTERNET SOURCES

32%

PUBLICATIONS

35%

STUDENT PAPERS

PRIMARY SOURCES

1

en.wikipedia.org

Internet Source

14%

2

www.guru99.com

Internet Source

10%

3

Jain, Arushi, and Vishal Bhatnagar. "Crime Data Analysis Using Pig with Hadoop",
Procedia Computer Science, 2016.

Publication

7%

4

www.tutorialspoint.com

Internet Source

4%

5

hadoop.apache.org

Internet Source

4%

6

V. Sujatha, S. Prasanna Devi, S. Vinu Kiran, S. Manivannan. "Bigdata Analytics on Diabetic Retinopathy Study (DRS) on Real-time Data Set Identifying Survival Time and Length of Stay", Procedia Computer Science, 2016

Publication

4%

7

docs.com

8

www.ijcsit-apm.com

Internet Source

2%

9

www.uvpce.ac.in

Internet Source

2%

10

S. Vennila, J. Priyadarshini. "Scalable Privacy Preservation in Big Data a Survey", Procedia Computer Science, 2015

Publication

2%

11

arxiv.org

Internet Source

2%

12

www.safaribooksonline.com

Internet Source

2%

13

B. Saraladevi, N. Pazhaniraja, P. Victor Paul, M.S. Saleem Basha, P. Dhavachelvan. "Big Data and Hadoop-a Study in Security Perspective", Procedia Computer Science, 2015

Publication

2%

14

research2.fit.edu

Internet Source

1%

15

www.oalib.com

Internet Source

1%

16

Student Paper

1 %

17

research.ijcaonline.org

Internet Source

1 %

18

Submitted to Higher Education Commission
Pakistan

Student Paper

1 %

19

Arushi Jain, Vishal Bhatnagar. "Crime Data
Analysis Using Pig with Hadoop", Procedia
Computer Science, 2016

Publication

1 %

20

Submitted to Institute of Technology, Nirma
University

Student Paper

1 %

21

rku.ac.in

Internet Source

1 %

22

Submitted to Charotar University of Science
And Technology

Student Paper

1 %

23

Submitted to University of Strathclyde

Student Paper

1 %

24

www.ijcaonline.org

Internet Source

1 %

Exclude quotes On

Exclude bibliography On

Exclude matches < 1%