# Mining Analysis on Customers Data using Big Data tools

**P.Sai Prasad**

*Ph.D. Research Scholar*
*Dept. of Computer Science*
*SSSUTMS*
*Sehore, M.P.*

**K.Rajesh**

*Assistant Professor*
*Dept.of Computer Science*
*BVSR Engineering College*
*Chimakurthy, A.P.*

**Abstract: In this paper the authors presented that how to analyse the Customers data when it is huge and various raw format data. So, for this we have used the Hadoop Big Data tools map-reduce Pig and Hive which makes easier to analyse the various raw formats of data.**

**Keyword: Hadoop, Big Data, Map- reduce, Hive**

## INTRODUCTION

Mostly the universe runs on the Business. So they get customers raw data of various formats and they don't know what is the data is there, for that we can use technique of Hadoop and Big data Tools for analysing the data of raw formats. First we should know what is Hadoop , Map reduce and hive.

Hadoop follows Master-Slave architecture. Hadoop uses HDFS to store files efficiently in the cluster. When a file is placed in HDFS it is broken down into blocks, 64 MB block size by default. These blocks are then replicated across the different nodes (*DataNodes*) in the cluster. The default replication value is 3, i.e. there will be 3 copies of the same block in the cluster. We will see later on why we maintain replicas of the blocks in the cluster.

A Hadoop cluster can comprise of a single node (single node cluster) or thousands of nodes.

**Hadoop is made up of 2 parts:**
1. HDFS – Hadoop Distributed File System
2. MapReduce – The programming model that is used to work on the data present in HDFS.

**Hadoop Distributed File System (HDFS) Building blocks of Hadoop:**
A. Namenode
B. Datanode
C. Secondary Name node
D. JobTracker
E. TaskTracker
**NameNode**
The NameNode in Hadoop is the node where Hadoop stores all the location information of the files in HDFS. In other words, it holds the metadata for HDFS.

**Secondary NameNode**
**IMPORTANT** – The *Secondary NameNode* is not a failover node for the *NameNode*. The secondary name node is responsible for performing periodic housekeeping functions for the *NameNode*. It only creates checkpoints of the file system present in the *NameNode*.

**DataNode**
The *DataNode* is responsible for storing the files in HDFS. It manages the file blocks within the node. It sends information to the *NameNode* about the files and blocks stored in that node and responds to the *NameNode* for all filesystem operations.

**JobTracker**
*JobTracker* is responsible for taking in requests from a client and assigning *TaskTrackers* with tasks to be performed. The *JobTracker* tries to assign tasks to the *TaskTracker* on the *DataNode* where the data is locally present (Data Locality). If that is not possible it will at least try to assign tasks to *TaskTrackers* within the same rack. If for some reason the node fails the *JobTracker* assigns the task to another *TaskTracker* where the replica of the data exists since the data blocks are replicated across the *DataNodes*. This ensures that the job does not fail even if a node fails within the cluster.

**TaskTracker**
*TaskTracker* is a daemon that accepts tasks (Map, Reduce and Shuffle) from the *JobTracker*.
The *TaskTracker* keeps sending a heart beat message to the *JobTracker* to notify that it is alive. Along with the heartbeat it also sends the free slots available within it to process tasks. *TaskTracker* starts and monitors the Map & Reduce Tasks and sends progress/status information back to the *JobTracker*.

**MapReduce**
MapReduce is the programming model that uses Java as the programming language to retrieve data from files stored in the HDFS. All data in HDFS is stored as files. Even MapReduce was built in-line with another paper by Google.

Google, apart from their papers did not release their implementations of GFS and MapReduce. However, the Open Source Community built Hadoop and MapReduce based on those papers. The initial adoption of Hadoop was at Yahoo Inc., where it gained good momentum and went onto be a part of their production systems. After Yahoo, many organizations like LinkedIn,
Facebook, Snapdeal Netflix and many more have successfully implemented Hadoop within their organizations.
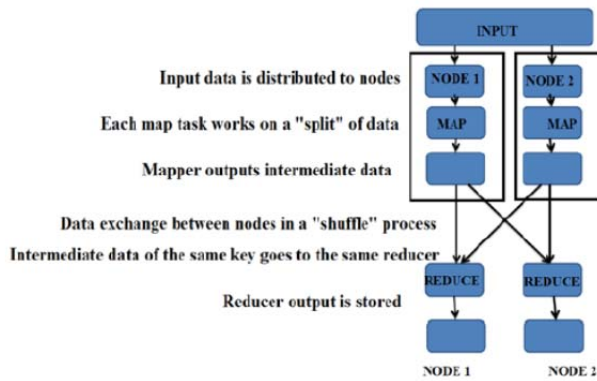
Fig.1.MapReduce Job submission flow mechanism

In this paper author's analyzed online data with the help of Big Data technology and visualized

**Analysis of Online Shopping Data**
Big Data technology analyzes the any type of data whether it is structured or unstructured. In Big Data technology everyone can analysis the data in a few second that takes more time in traditional one.

In the past for analysis no one used the Big Data technology to analyze and visualize the large scale online data like gigabytes or terabytes. Therefore we need parallel techniques. With the help of Big Data technology, we can easily tackle the huge amount of online data in parallel.

Hadoop distribution platform for implementing Big Data technologies i.e. Hadoop, Pig and Hive etc. It provides virtual Linux environment for map reducing. Apache

Hadoop is a framework running on the top of the cloudera .Apache Hadoop is used for processing the huge data across the cluster of commodity computers using simple programming model. Map Reduce is generally used for splitting the task across processor. Map Reduce is written in Java programming language. For running ClouderaCDH5, VM Player has been used to execute the ClouderaCDH5 machine as a virtual framework. Local machine must have required 4GB RAM for successfully configure the ClouderaCDH5 machine on VM Player.

**Algorithm for analysis of Online Data**
An algorithm used for analysis of online data has two phases i.e. mapper phase and reducer phase.

**Mapper**
The Mapper code reads the input files as <Key,Value> pairs and emits key value pairs.

**Sort and shuffling is done Reducer**
The Reducer code reads the outputs generated by the different mappers as <Key,Value> pairs and emits key value pairs.

**SIMULATION WORK**
Simulation work has been initialized by installation of clouderaCDH3 on windows.
Take some file named as retail.txt


Fig 2.Raw Data

**First create DataBase**

Hive>create database onlinetransaction;

**Create table in Database**

Hive(onlinetransaction)>create table transaction(txnno INT,txndate STRING,custno INT,amount DOUBLE,category STRING,product STRING,city STRING,state STRING,spendby STRING) row format delimited fields terminated by ',' stored as textfile;

**Loading data from local**

Hive>load data local in path 'home/cloudera/desktop/retail.txt' overwrite into table transaction;

**Operations on Online Raw Data**

If we want category wise total amount

**Hive>Select category, Sum(amount) from transaction group by category;**

After that, these analyzed data presented on Hadoop and connected directly to Big Data Visualization tool i.e. Tableau tool. All the data have collected in the Big Data tool.

```
state    string
spendby string
Time taken: 0.055 seconds
hive> select category , sum(amount) from txnrecords group by category;
Total MapReduce 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 mapred.reduce.tasks=<number>
Starting Job = job_201409131257_0081, Tracking URL = http://localhost:50030/jobdetails.jsp?jobid=job_201409131257_0081
Kill Command = /usr/lib/hadoop/bin/hadoop job  -Dmapred.job.tracker=localhost:8021 -kill job_201409131257_0081
2014-10-18 12:04:40,339 Stage-1 map = 0%,   reduce = 0%
2014-10-18 12:04:45,362 Stage-1 map = 100%,   reduce = 0%
2014-10-18 12:04:54,408 Stage-1 map = 100%,   reduce = 100%
Ended Job = job_201409131257_0081
OK
Air Sports      99316.90000000001
Combat Sports   164730.66999999995
Dancing 42603.71000000001
Exercise & Fitness      766463.8699999992
Games   374932.69999999995
Gymnastics      327225.3400000001
Indoor Games    288506.0399999998
Jumping 205842.66999999995
Outdoor Play Equipment  294728.28000000014
Outdoor Recreation      846678.6400000013
Puzzles 61564.75
Racquet Sports  166976.05999999988
Team Sports     617461.3799999999
Water Sports    531815.9700000004
Winter Sports   321973.5599999998
Time taken: 17.581 seconds
hive>
```

Fig 3: showing the amount by category wise

## CONCLUSION

In this study, we have analyzed and visualized the online data of India using Big Data technology for a given data in fraction of second which was not possible by traditional techniques. In this paper authors are thus using Big Data technology to better understand the data and minimize the potential damage that an online data can cause.

## REFERENCES

[1] ApacheHive .https://hive.apache.org/index.html.
[2] Prajapati V (2014) Big Data Analytics with R and Hadoop. Shroff Publishers & Distributors Pvt Ltd.
[3] Cloudera.Http://www.cloudera.com/content/cloudera/en/products-and-Services/cdh.html.
[4] Hadoop: The Definitive Guide by Tom White, 3rd Edition, O'reilly
[5] Hadoop in Action by Chuck Lam, MANNING Publ.
[6] Hadoop for Dummies by Dirk deRoos, Paul C.Zikopoulos, Roman B.Melnyk,Bruce Brown, Rafael Coss
[7] Hadoop in Practice by Alex Holmes, MANNING Publ.
[8] Hadoop MapReduce Cookbook,Srinath Perera, Thilina Gunarathne