

EPH-Enhancement of Parallel Mining using Hadoop

Neha Mangla

Associate Professor, A.I.T.,Bangalore
Apj.neha@gmail.com

Sushma K

MTech Scholar, A.I.T., Bangalore
sushmakrishna2@gmail.com

ABSTRACT

Data in this era is generating at tremendous rate so now it is need of today to handle the data to gain useful insight, this data can be useful for researcher and accommodation to do analysis. As we know traditional system cannot handle more than terabytes of data since it affects performance and also storage is very costly. Bigdata is a innovative technique analyze, store, manage, distribute and capture datasets. To achieve compressed storage in this implement a parallel mining algorithm called as enhancement of parallel mining using Hadoop. Hadoop is a platform which enables the distributing processing using mapreduce programming. This help in getting result at very fast rate as result in less time help in competing for growth of business. For the analysis in this paper unstructured datasets from real-time is taken and converted to structured format and process in mapreduce. It is found in literature existing mining algorithm for real time datasets lacks in fault tolerance, load balancing, data distribution and automatic parallelization. To overcome these disadvantages we implement mapreduce for association analysis. In EPH we improve performance by distributing load across the computing nodes .In our proposed solution we use real-world celestial spectral data .The graphical representation of traditional system comparison with Hadoop is shown in this paper.

Keywords: Bigdata, Hadoop, Mapreduce, Parallel Mining, Association analysis, Enhancement of parallel mining using Hadoop(EPH).

I. INTRODUCTION

Information Technology is growing rapidly and volume of data is increasing such as social media, balck box so on and it is reaching petabytes of data threshold and as increase in data also

increases computational requirements which includes fault tolerance, load balancing, data distribution and automatic parallelization. In terms of academics and business bigdata is become the key role. Here efficient parallel mining algorithm techniques are used to easy and fast processing of data. In this we consider a data mining tool called as R tool to compare with the proposed system where we process our unstructured data perform the association analysis on the datasets and represent using graphical presentation in R tool the time required to process the datasets is more it recursively process the datasets and it cant process large amount of datasets which is great drawback.

Association analysis cannot be done on the real time datasets so this one of drawback .The real time data can be processed in Mapreduce. Firstly we generate realistic data by creating developer account and by streaming the data in the flume and the unstructured data is taken to hue and particular data can be searched using the solaris. We can also change colour, bolds, and highlights and so on in solaris only. Now the data as to be converted to structured data using Hive and this structured data is fed into Mapreduce is a programming model which consists of two phase:-

1. Firstly, Mapper phase which is each separate line and produces a key value pair.
2. Secondly, to do the association analysis and it is represented using graphical representation. After the all these task performed we compare the two system speed, processing speed, availability, an time taken for the execution of data and many more criteria.

The contributions of this paper are:-

I made complete overview about parallel mining of realistic datasets and converted to structured datasets.

I developed a parallel mining method using Map reduce programming model I also gave complete overview about existing traditional system R tool and did the association analysis using the datasets.I

also show the load balancing and how data is being distributed in the clustering nodes and processing is done.

The comparison of both the system is showed and it is measured in turn of processing speed, scalability ,availability, performance and what kind of synthetic and real world datasets that can be processed in these tools.

The paper contains as follows in Section II it describes background knowledge how the sentiment analysis of twitter data Section III The Overview about the existing system and how they are processing the Twitter data Section IV Overview about EPH Section V Implementation EPH Algorithm Section VI EPH over existing system Section VII Conclusion.

II. PRELIMINARY

Nowadays Wide area network plays a important role in the world this tend me to work on Twitter data. In this section we are going to learn about how we do Sentimental analysis on Twitter data from the mapreduce and R Tool. A sentimental analysis is nothing but retrieving the tweets from the Twitter API and calculating the score for that particular data. Firstly, we explain the sentimental analysis for here in existing system we are using the corpus based approach for the sentimental analysis which also include Natural language processing and machine learning processing.

There are some key characteristics of this tweets data:-

- *Message Length:* The maximum length of a Twitter message is 140 characters. This is different from previous sentiment classification research that focused on classifying longer texts, such as product and movie reviews.
- *Writing technique:* The occurrence of incorrect spellings and cyber slang in tweets is more often in comparison with other domains. As the messages are quick and short, people use acronyms, misspell, and use emoticons and other characters that convey special meanings.

Availability: The amount of data available is immense. More people tweet in the public domain as compared to Facebook (as Facebook has many privacy settings) thus making data more readily available. The Twitter API facilitates collection of tweets for training.

Topics: Twitter users post messages about a range of topics unlike other sites which are designed for a

specific topic. This differs from a large fraction of past research, which focused on specific domains such as movie reviews.

Real time: Blogs are updated at longer intervals of time as blogs characteristically are longer in nature and writing them takes time. Tweets on the other hand being limited to 140 letters and are updated very often. This gives a more real time feel and represents the first reactions to events

Secondly to process the data in the map reduce the sentimental analysis is done in the Flume.

Flume is a reliable, distributed and available service for efficiently aggregating, collecting and moving the large amount of celestial data and log data.

Flume components interact in the following way:

Client starts the flow of the flume

The transmission of the Event is done by the client to a source operating with the Agent

This Event is received by the Source and delivered to one or more channels.

Channels can be drained due the one more sink with the same Agent.

The ingestion rate from the drain rate can be calculated by Channels using producer-consumer model of data exchange

When spikes in client side activity cause data to be generated faster than can be handled by the provisioned destination capacity can handle, the **Channel** size increases. This allows sources to continue normal operation for the duration of the spike.

The **Sink** of one **Agent** can be chained to the **Source** of another **Agent**. This chaining enables the creation of complex data flow topologies.

Because Flume's distributed architecture requires no central coordination point. Each agent runs independently of others with no inherent single point of failure, and Flume can easily scale horizontally

Initially I create a gmail account and followed by a Twitter account in this and start with our developer account and go for the option manage your app and generate the keys and access token.



Figure 1. Snapshot of creation of developer account

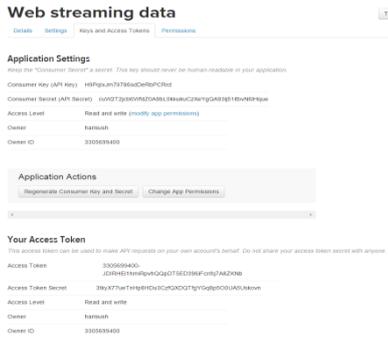
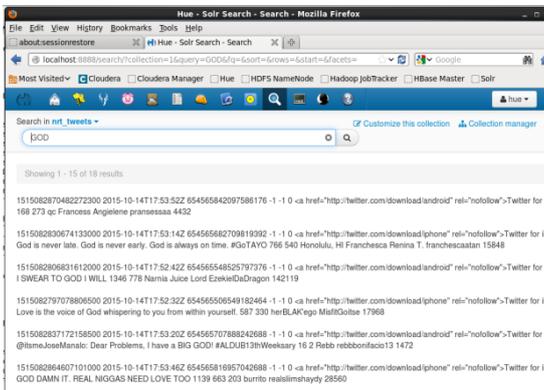


Figure 2. I have developer account with name WebStreaming data with keys and access tokens.



Figure 3. Coding in Hadoop to connect twitter with Flume



flume.

Figure 5 Query done in Hive platform

The data which is present now is unstructured data shown in Figure 5 which has to be converted to the structured format where this job is done by the hive platform only after conversion it can be processed to the mapreduce and it is processed for the further analysis.

III. EXISTING SYSTEM

Here I am going to learn about the how the tweets can be extracted using the code in the R tool with the called as twitter R[1][2], we need to install this package first. Later the tweets are converted to dataframe then corpus. The code is given below:

```
#
convert
tweets
to a
data
frame
df <-
do.call
("rbind", lapply(rdmTweets, as.data.frame))
dim(df)
library(tm)
# build a corpus, and specify the source to be
character vectors
myCorpus <- Corpus(VectorSource(df$text))
```

Next step is used for stemmed words to retrieve their radicals from the twitter data. I need to install few packages they are Snowball, RWeka, Rjava and RWekajars. Now we have to convert this unstructured to structure in terms of matrix. Where here row means terms column means entity we build a term-document matrix from the above processed corpus with function TermDocumentMatrix()[7] After the completion of this stage a row will be generated.

It is now list Frequent term and association. findFreqTerms() this function is list the number of frequent less than 10(10 is just a example).Now I will plot the bar graph for the twitter data which is having the maximum tweets and retweets show in Figure 7.

Figure 4. Twitter data retrieved in

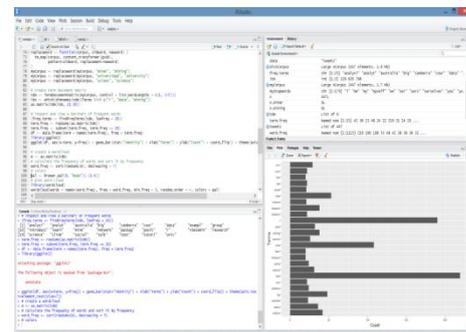


Figure 6. It shows the bar graph the maximum tweets that is happened.

I can also generate the word cloud for the particular datasets. I can calculate the maximum frequency and I can generate the word cloud for those maximum words. I can generate the networks of the terms for the particular data. a term-document matrix, term DocMatrix, is loaded into R, later it is converted to adjacency matrix. This matrix contains the maximum frequencies and it is in the form of table. Figure 7 shows the graph of the maximum number of words repeated in

the large number of datasets. To improve the storage efficiency of the system also.

Here the aggregation stage where in this process the given dataset is sent to mapper stage for processing and the association analysis is taken where it accepts the data key value pair in the given system it count the maximum of all tweets present in the table and later it is processed in intelligent graph to obtain the expected result

Now I will the processing in the mapreduce before I start with mapreduce I should all the services using the command called as start-all.sh. I can check in the weather services are started or not using jps command and I can also check in hadoop administration shown in Figure 10 and Figure 11. I can create a new file in the HDFS here I have created a file called as fi. Using the command show in the Figure 12 I will copy the twitter data in the file fi and I can weather the data is copied or not in hadoop adminstartion as show in Figure 13.

Now it is time to start the mapreduce process and in the mapreduce process the for each instace of how much the mapreduce shows the data their only and wait the reduce to finish 100% as shown in the Figure 14. If the mapreduce successful then the fie called success will be created in fi inside the output file as shown in Figure 15. If the process is unsuccessful then it gives error message then and their only.

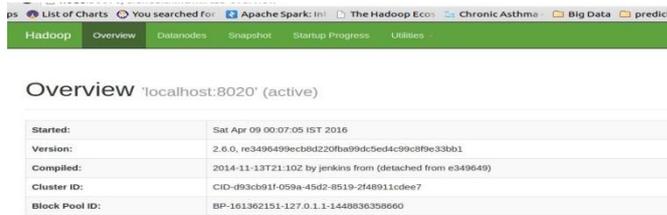
The obtained output must be put into the Intelligent graph which is written in php it is processed as shown in the Figure 16 The result obtained by the mapreduce is shown in term of line graph, bar graph, pie graph as shown in the Figure 17,18 and 19 .Even other types can also be implemented.

```

user@node:~$ jps
2830 SecondaryNameNode
3567 Jps
2551 NameNode
3030 ResourceManager
2670 DataNode
3155 NodeManager
user@node:~$

```

Figure 10. Start all the services



Summary

Figure 11. Checking in the hadoop administration services are started or not.

```

DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.
16/04/09 00:16:49 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
put: -: No such file or directory
user@node:~$ cd /home/user/workspace/FlDooop/twtterdata/
user@node:~/workspace/FlDooop/twtterdata$ hadoop dfs -put twitter_data.txt /fi
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.
16/04/09 00:25:58 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
put: -: No such file or directory
user@node:~/workspace/FlDooop/twtterdata$ hadoop dfs -put twitter_data.txt /fi
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.
16/04/09 00:27:34 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
user@node:~/workspace/FlDooop/twtterdata$

```

Figure 12 The data is copied to the fi file.



Figure 13 Shows the copied file in the hadoop adminstartion

```

0140671592_0001
6/04/09 00:38:31 INFO impl.YarnClientImpl: Submitted application application_140671592_0001
0140671592_0001
6/04/09 00:38:32 INFO mapreduce.Job: The url to track the job: http://node:8088/proxy/application_1460140671592_0001/
6/04/09 00:38:32 INFO mapreduce.Job: Running job: job_1460140671592_0001
6/04/09 00:38:59 INFO mapreduce.Job: Job job_1460140671592_0001 running in uber mode : false
6/04/09 00:38:59 INFO mapreduce.Job: map 0% reduce 0%
6/04/09 00:39:54 INFO mapreduce.Job: map 100% reduce 0%
6/04/09 00:40:27 INFO mapreduce.Job: map 100% reduce 100%
6/04/09 00:40:29 INFO mapreduce.Job: Job job_1460140671592_0001 completed successfully
6/04/09 00:40:29 INFO mapreduce.Job: Counters: 49
File System Counters
FILE: Number of bytes read=108
FILE: Number of bytes written=213083
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=129776
HDFS: Number of bytes written=90
HDFS: Number of read operations=6
HDFS: Number of large read operations=0

```

Figure 14. The mapreduce process for the file

Browse Directory

Permission	Owner	Group	Size	Replication	Block Size	Name
-rw-r--	user	supergroup	0 B	1	128 MB	._SUCCESS
-rw-r--	user	supergroup	90 B	1	128 MB	part-00000

Hadoop, 2014.

Figure 15 The success file created in fi file if mapreduce process complete successfully

```

user@node: /var/www/html
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
Bytes Read=229670
File Output Format Counters
Bytes Written=90
user@node:~$ cd /var/www/html
user@node: /var/www/html$ cp -R /home/user/workspace/FtDooop/Parllel\ Mini
erGraph /var/www/html
user@node: /var/www/html$ chmod -R 777 TwitterGraph
user@node: /var/www/html$ ls -lrt
total 40
-rwxr-xr-x 1 user user 11510 Nov 20 17:57 index.html
-rwxr-xr-x 1 user user 55 Nov 20 18:47 deploy.sh
-rwxr-xr-x 1 user user 537 Nov 20 19:05 hello.php
-rwxr-xr-x 1 user user 764 Nov 20 20:09 hello.php
drwxr-xr-x 6 user user 4096 Nov 21 13:22 BIG_DATA_PROJECT_GRAPH
drwxrwxr-x 8 user user 4096 Nov 26 14:51 graph

```

Figure 16 Initializing to the twitter graph and processing in it

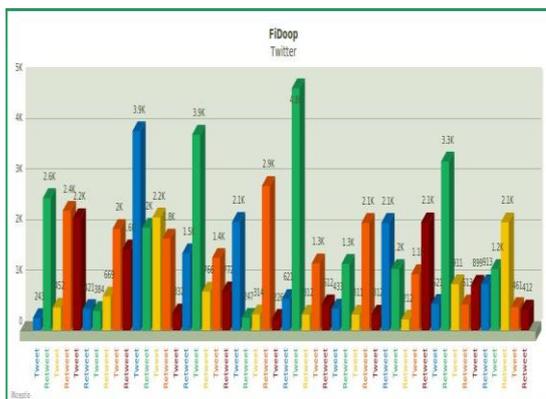


Figure 17 The expected result obtained in bar graph

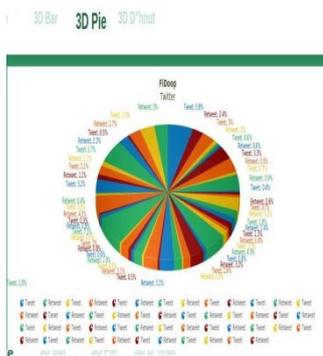


Figure 18 The result obtained in pie graph

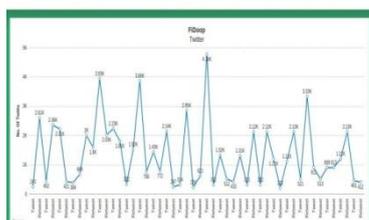


Figure 19 The result obtained in line graph

VI. EPH OVER EXSISTING TOOL

In this thesis I have work on the traditional system that is R tool now in this section. According to the programs done by me I have taken few factors for the comparison with the system.Few factors are:-

- **Parallelism:** The mapreduce supports the parallelism processing data parallel we can process even 1TB of data in just minute (minimum of 4 nodes) but when consider in the R tool it doesn't provide any parallelism in its structure.
- **Datasets:** As the data increases we can also include nodes in the cluster which increase the performance of system and so large datasets are processed in the system but In the increase in data the R tool performance is drained.
- **Fault tolerance:** In mapreduce iif some task is being interrupted then task tracker sends the heartbeat signal to the job tracker and the re-execution is done automatically but in R tool no such facilities.
- **CPU Time:** This is a little harder to understand, because this time will not change if a job does not change. This is the time consumed by CPU, which is used for processing instructions of a computer program. It is only related to the number of instructions in a program, which means if we parallelize a program, the total execution time might be the third of the original, but the CPU time would not change because the amount of instructions does not change but are executed in parallel.
- **Availability :** when we store the data in HDFS if there is 1 cluster present but by default 3 replication is created in the HDFS.So if one datasets is duplicated or deleted we can retrieve the data but when compared with the traditional system there is no such advantage is provided.
- **Scalable:** Hadoop is highly storage platform distributes the data across the many nodes it is not necessary that we install expensive servers but any kind of servers can be installed by this way the storage can be expanded .This is not provided in datamining tool
- **Load balancing:** The datasets on the cluster is automatically scattered between the node in the cluster without any user instruction if there is more load on single load then it is automatically transferred to the other node where this cannot be done in the traditional tool.

- Data distribution :By default the namenode distribute the data between the cluster in the node so user did not any kind instruction for the distribution of the data and in the data mining tool it is not possible.

VII. CONCLUSION

In the paper I applied and worked both on data mining and hadoop as per the result obtained Datamining tool which applied could not manage loadbalancing,datadistribution ,fault tolerance and many more problems is faced in the datamining tool so this solved by hadoop platform in the mapreduce when compared with the data mining tool mapreduce is advantageous in all the terms. EPH incorporates the parallel mechanism for the twitter data where I can achieve compressed storage and it is necessity not to build any conditional pattern. Firstly in this processes I retrieved the data using the flume which is in unstructured format and later stored in the HDFS and from the HDFS it is then processed to the hive platform and converted to the structured data.

I performed the analysis in the data mining tool I retrieved the data from the twitter data obtained in flume I processed the data in the r tool and we converted the data unstructured to structured data and processed the data to obtain the bar graph where the maximum frequency of words where it was obtained after using a package called as igraph and later the same data was used to obtain the tree graph.

In the hadoop system all services are started the data present in hdfs which is obtained from flume.The important stage of the paper is the mapreduce algorithm which was obtained where their three algorithm implemented mapper,reducer and driver config code was implemented where the data format specified in the for inputting the data is taken and the splitting of the data is done and in map stage it is converted to key value pair and it is sent reduce stage in the reduce stage the aggregation is done calculating the frequency of the tweets data present then later it is processed to the Driver function and aggregation is done. Later I run the mapreduce process as the process gets completed a success file gets created in the HDFS and later we connect the result with the intelligent graph system and result is obtained in term of bar graph,line graph and 3D graph and so on having the maimum tweet and retweet in the table.

VIII. REFERENCES

1. Mangla N and Khola R.K ” Optimization of IP Routing with Content Delivery Network” published in “Networking and Information Technology (ICNIT), 2010 InternationalConference “of IEEE Explore,June11-122010, E-ISBN : 978-1-4244-7578-0 pp.424-428 .
2. Apoorv Agarwal, Fadi Biadisy, and Kathleen Mckeown. 2009. Contextual phrase-level polarity analysis using lexical affect scoring and syntactic n-grams. Proceedings of the 12th Conference of the European Chapter of the ACL (EACL 2009), pages 24–32, March.
3. Luciano Barbosa and Junlan Feng. 2010. Robust sentiment detection on twitter from biased and noisy data. Proceedings of the 23rd International Conference on Computational Linguistics: Posters, pages 36–44.
4. Neha Mangla and Dr. R.K.Khola “A Way to Implement BGP with Geographic Information”published in International Journal of Electronics Engineering, 2 (2), 2010, pp.349–353.
5. NehaMangla and Dr.R.K.Khola,” Application Based Route Optimization”, IOSR Journal of Engineering (IOSRJEN) ISSN: 2250-3021 Volume 2, Issue 8 (August 2012), PP 78-82.
6. Adam Bermingham and Alan Smeaton. 2010. Classifying sentiment in microblogs: is brevity an advantage is brevity an advantage? ACM, pages 1833–1836. C. Fellbaum. 1998. Wordnet, an electronic lexicaldatabase. MIT Press.
7. Lammel, R.: Google’s MapReduce Programming Model - Revisited. Science of Computer Programming 70, 1–30 (2008)