

A Specialized Log Analysis Engine in Distributed Environment

Abhiruchi Shinde¹, Neha Vautre², Prajakta Yadav³, Sapna Kumari⁴

^{1,2,3,4}Dept of Computer Engineering, SITS, Maharashtra, India

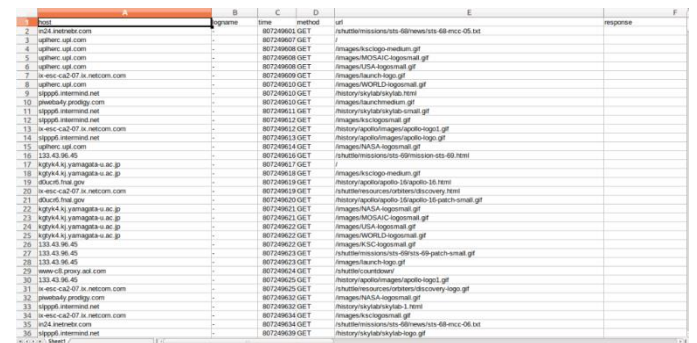
Abstract - Log file or logs in computing are the files for keeping record of the events that occur in the operating system or communication between the users or operating systems. Log files contains large amount of valuable information about the system operation status, usage, user behavior analysis etc. Due to extensive use of digital appliances in today's modern era log file analysis has become a necessary task to track system operation or user behavior and acquire important knowledge based on it. These kinds of files are generated at stupendous rate and to analyze them is tedious task and a burden to corporations and various organizations. In order to analyze large dataset, and to store it efficiently, economically and effectively we need to have an effective solution which needs not only the massive and stable data processing ability but also the adaptation to a variety of scenarios under the requirement of efficiency. Such capabilities can't be achieved from standalone analysis tools or even single cloud computing framework. The main objective of the proposed system is to design an application for log analysis and applying the data mining algorithm to get the results which will be useful for system administrator to take proper decisions. The combination of Hadoop, Spark and the data warehouse and analysis tools of Hive and Shark makes it possible to provide a unified platform with batch analysis and in-memory computing capacity in order to process log in a high available, stable and efficient way. Statistics based on customer feedback data from the system will help in greater expansion of business and a company that will have such data to its disposal and ready to use in the distributed environment for log analysis

Key Words: Log, Weblog, Hadoop, Spark, Log analysis.

1. INTRODUCTION

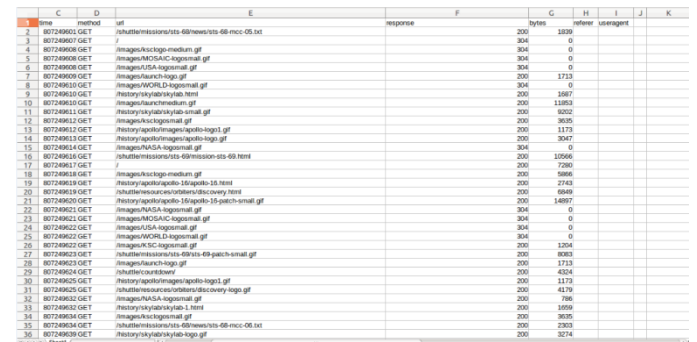
Big data analytics is the process of examining huge amount of data present in structured or unstructured form generated at a high speed to uncover hidden patterns, unknown correlations, market trends, customer preferences and other useful business information. The analytical findings can lead to more effective marketing, new revenue opportunities, better customer service, improved operational efficiency, competitive advantages over rival organizations and other business benefits. . With big data analytics, data scientists and others can analyze huge volumes of data that conventional analytics and business intelligence solutions can't touch. Log analysis is the way to gather information of the number of access users, user behavior, operation status etc. There have been some free powerful log analysis tools like Awstats ,

Webalizer , and Google Analytics. But they are either standalone or have the limitation of data scale. Consequently the higher number of log producing on a daily basis has led to the analysis task to be more hectic and tedious, in large databases, so we here propose a system for log analysis. In this system the web server logs are analyzed in distributed environment. Due to digitization in today's modern era log file analysis has become a necessary task to track system operation or user behavior. Log analysis is necessary for any organization to determine how well their website is performing as marketing tool. Log analysis is a tedious task. There is a need of an effective solution for integration and parallel processing of data. We will present an engine in distributed environment for log analysis using the integration of hadoop and hive and also spark and shark for the purpose. The main objective of the project is to produce such a engine that reduces the tedious and alone machine introducing the distributed environment also reducing the overall time required for the computation to take place.



time	method	url	response	bytes	referer	useragent
01/01/2017 10:00:00	GET	/shuttle/mission/ats-68/news/ats-68-mcc-05.txt	200	1889		
01/01/2017 10:00:01	GET	/	304	0		
01/01/2017 10:00:02	GET	/images/ckc/logo-medium.gif	304	0		
01/01/2017 10:00:03	GET	/images/ckc/cac-logomail.gif	304	0		
01/01/2017 10:00:04	GET	/images/usa-logomail.gif	304	0		
01/01/2017 10:00:05	GET	/images/launch-logo.gif	200	1728		
01/01/2017 10:00:06	GET	/images/WORLD-logomail.gif	200	1267		
01/01/2017 10:00:07	GET	/shuttle/mission/ats-68/news/ats-68.html	200	14857		
01/01/2017 10:00:08	GET	/images/ckc/logo-medium.gif	200	1032		
01/01/2017 10:00:09	GET	/images/usa-logomail.gif	200	1173		
01/01/2017 10:00:10	GET	/images/ckc/cac-logomail.gif	200	1095		
01/01/2017 10:00:11	GET	/images/launch-logo.gif	200	1173		
01/01/2017 10:00:12	GET	/images/WORLD-logomail.gif	200	1267		
01/01/2017 10:00:13	GET	/shuttle/mission/ats-68/news/ats-68.html	200	14857		
01/01/2017 10:00:14	GET	/images/ckc/logo-medium.gif	200	1032		
01/01/2017 10:00:15	GET	/images/usa-logomail.gif	200	1173		
01/01/2017 10:00:16	GET	/images/ckc/cac-logomail.gif	200	1095		
01/01/2017 10:00:17	GET	/images/launch-logo.gif	200	1173		
01/01/2017 10:00:18	GET	/images/WORLD-logomail.gif	200	1267		
01/01/2017 10:00:19	GET	/shuttle/mission/ats-68/news/ats-68.html	200	14857		
01/01/2017 10:00:20	GET	/images/ckc/logo-medium.gif	200	1032		
01/01/2017 10:00:21	GET	/images/usa-logomail.gif	200	1173		
01/01/2017 10:00:22	GET	/images/ckc/cac-logomail.gif	200	1095		
01/01/2017 10:00:23	GET	/images/launch-logo.gif	200	1173		
01/01/2017 10:00:24	GET	/images/WORLD-logomail.gif	200	1267		
01/01/2017 10:00:25	GET	/shuttle/mission/ats-68/news/ats-68.html	200	14857		
01/01/2017 10:00:26	GET	/images/ckc/logo-medium.gif	200	1032		
01/01/2017 10:00:27	GET	/images/usa-logomail.gif	200	1173		
01/01/2017 10:00:28	GET	/images/ckc/cac-logomail.gif	200	1095		
01/01/2017 10:00:29	GET	/images/launch-logo.gif	200	1173		
01/01/2017 10:00:30	GET	/images/WORLD-logomail.gif	200	1267		
01/01/2017 10:00:31	GET	/shuttle/mission/ats-68/news/ats-68.html	200	14857		
01/01/2017 10:00:32	GET	/images/ckc/logo-medium.gif	200	1032		
01/01/2017 10:00:33	GET	/images/usa-logomail.gif	200	1173		
01/01/2017 10:00:34	GET	/images/ckc/cac-logomail.gif	200	1095		
01/01/2017 10:00:35	GET	/images/launch-logo.gif	200	1173		
01/01/2017 10:00:36	GET	/images/WORLD-logomail.gif	200	1267		

Fig1: Screenshot of web server logs



time	method	url	response	bytes	referer	useragent
01/01/2017 10:00:00	GET	/shuttle/mission/ats-68/news/ats-68-mcc-05.txt	200	1889		
01/01/2017 10:00:01	GET	/	304	0		
01/01/2017 10:00:02	GET	/images/ckc/logo-medium.gif	304	0		
01/01/2017 10:00:03	GET	/images/ckc/cac-logomail.gif	304	0		
01/01/2017 10:00:04	GET	/images/usa-logomail.gif	304	0		
01/01/2017 10:00:05	GET	/images/launch-logo.gif	200	1728		
01/01/2017 10:00:06	GET	/images/WORLD-logomail.gif	200	1267		
01/01/2017 10:00:07	GET	/shuttle/mission/ats-68/news/ats-68.html	200	14857		
01/01/2017 10:00:08	GET	/images/ckc/logo-medium.gif	200	1032		
01/01/2017 10:00:09	GET	/images/usa-logomail.gif	200	1173		
01/01/2017 10:00:10	GET	/images/ckc/cac-logomail.gif	200	1095		
01/01/2017 10:00:11	GET	/images/launch-logo.gif	200	1173		
01/01/2017 10:00:12	GET	/images/WORLD-logomail.gif	200	1267		
01/01/2017 10:00:13	GET	/shuttle/mission/ats-68/news/ats-68.html	200	14857		
01/01/2017 10:00:14	GET	/images/ckc/logo-medium.gif	200	1032		
01/01/2017 10:00:15	GET	/images/usa-logomail.gif	200	1173		
01/01/2017 10:00:16	GET	/images/ckc/cac-logomail.gif	200	1095		
01/01/2017 10:00:17	GET	/images/launch-logo.gif	200	1173		
01/01/2017 10:00:18	GET	/images/WORLD-logomail.gif	200	1267		
01/01/2017 10:00:19	GET	/shuttle/mission/ats-68/news/ats-68.html	200	14857		
01/01/2017 10:00:20	GET	/images/ckc/logo-medium.gif	200	1032		
01/01/2017 10:00:21	GET	/images/usa-logomail.gif	200	1173		
01/01/2017 10:00:22	GET	/images/ckc/cac-logomail.gif	200	1095		
01/01/2017 10:00:23	GET	/images/launch-logo.gif	200	1173		
01/01/2017 10:00:24	GET	/images/WORLD-logomail.gif	200	1267		
01/01/2017 10:00:25	GET	/shuttle/mission/ats-68/news/ats-68.html	200	14857		
01/01/2017 10:00:26	GET	/images/ckc/logo-medium.gif	200	1032		
01/01/2017 10:00:27	GET	/images/usa-logomail.gif	200	1173		
01/01/2017 10:00:28	GET	/images/ckc/cac-logomail.gif	200	1095		
01/01/2017 10:00:29	GET	/images/launch-logo.gif	200	1173		
01/01/2017 10:00:30	GET	/images/WORLD-logomail.gif	200	1267		
01/01/2017 10:00:31	GET	/shuttle/mission/ats-68/news/ats-68.html	200	14857		
01/01/2017 10:00:32	GET	/images/ckc/logo-medium.gif	200	1032		
01/01/2017 10:00:33	GET	/images/usa-logomail.gif	200	1173		
01/01/2017 10:00:34	GET	/images/ckc/cac-logomail.gif	200	1095		
01/01/2017 10:00:35	GET	/images/launch-logo.gif	200	1173		
01/01/2017 10:00:36	GET	/images/WORLD-logomail.gif	200	1267		

Fig2: Screenshot of web server logs

2. EXISTING SYSTEM

There are many systems already available satisfying similar purposes. They are Webalizer and Awstats. We briefly understand the two existing technologies.

Webalizer:- The Webalizer is an application tool that generates web pages by analyzing logs, from access and usage logs, i.e. it is a software for analysis of web server logs. It is one of the most commonly used web server administration tools. It was initiated by Bradford Barret in 1997. Statistics commonly reported by Webalizer include hits, visits, referrers, the visitors' countries, and the amount of data downloaded. These statistics can be viewed graphically and presented by different time frames, such as by day, hour, or month.

3. PROPOSED SYSTEM

The previous systems as known earlier consists of the standalone machines, owing to which, there occurs a huge problem of data scalability. To address this issue and various other limitations of the previous systems, we propose to add distributed environment to the existing system. Adding this feature we try to remove the bugs of the previously proposed systems and make them work more efficiently as well as effectively, to produce desired results to the client to maintain a specific position in the market. We use map reduce operations along with the Hadoop system in HDFS to perform operations like knowing the number of hits for a particular URL, error detection, number of bytes transmitted, number of hits for each IP address. Descriptions of the entire map reduce phases along with the system architecture and results of the same are given below.

4. SYSTEM ARCHITECTURE

In any analytical tool, pre-processing is necessary, because Log file may contain noisy & ambiguous data which may affect result of analysis process. Log pre-processing is an important step to filter and organize only appropriate information before applying Map Reduce algorithm. Preprocessing reduces size of log file also it increases quality of available data. The purpose of log pre-processing is to improve log quality and increase accuracy of results.

4.1. MAP REDUCE FRAMEWORK

Map Reduce is a simple programming model for parallel processing of large volume of data.. Fundamental concept of Map Reduce is to transform lists of input data to lists of output data. Map Reduce does the conversion twice for the two major tasks: Map and Reduce just by dividing whole workload into number of tasks and distributing them over different machines in the Hadoop cluster. For the application which require tedious task of log data analysis, Map Reduce implementation in Hadoop is one of the best

solutions. Map Reduce is divided into two phases: Map phase and Reduce phase.

Map phase:-

Input to the Map Reduce is log file, each record in log file is considered as an input to a Map task. Map function takes a key-value pair as an input thus producing intermediate result in terms of key-value pair. It takes each attribute in the record as a key and Maps each value in a record to its key generating intermediate output as key-value pair.

Reduce phase:-

Reduce task takes key and its list of associated values as an input. It combines values for input key by reducing list of values as single value which is the count of occurrences of each key in the log file, thus generating output in the form of key-value pair

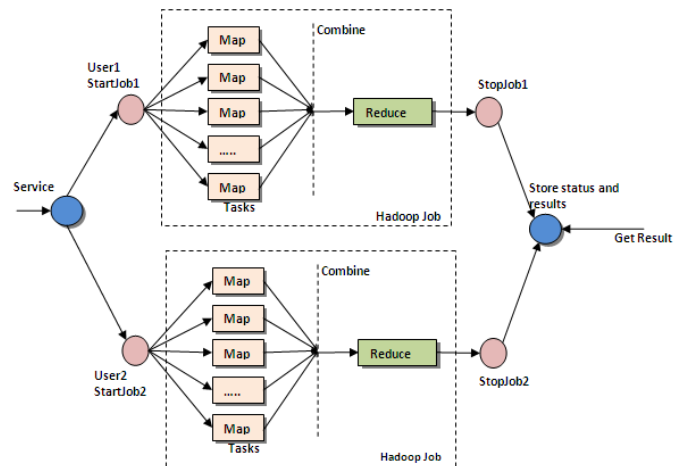


Fig.3: Map Reduce operation

The proposed system consists of major components like Web servers, implementing Hadoop storage and Map Reduce programming model and user interface.

The basic Log analysis flow in our system starts with the log data sets collection, and then ETL processing is done on log data to make it suitable for processing. Log analysis is done on hadoop as well as spark using the analytical tools hive and shark respectively. Spark does not have it's own file system so it will store it's data in HDFS. Parallel processing will be done on hadoop and for iterative query processing on shark. The result of log analysis can be given as an input to the tool named zeppelin which will show the statistical analysis of result in the form of chart.

Data cleaning is the first phase carried out in the proposed work as a pre-processing step in NASA web server log files. The NASA log file contains a number of records that corresponds to automatic requests originated by web robots, that includes a large amount of erroneous,

Counting the number of times a particular [P or host requesting the server is one of the crucial tasks in web log analysis and could be significantly useful in host centric recommendation systems and in analyzing and identifying potential security attackers.

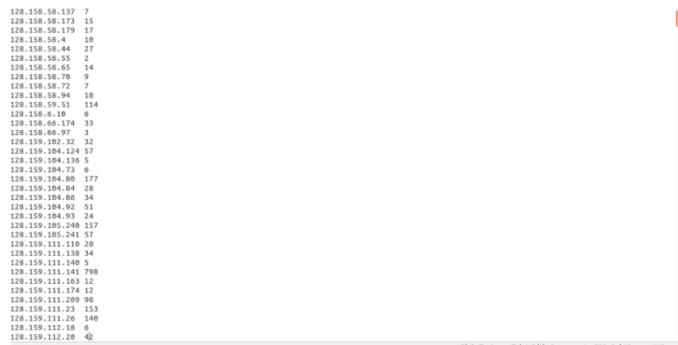


Fig.9: Number of hits for each IP address

6. DISCUSSIONS

We here discuss about the various methods about the flow of log analysis and the phases for completion of the proposed system. Following is the state diagram for the proposed system.

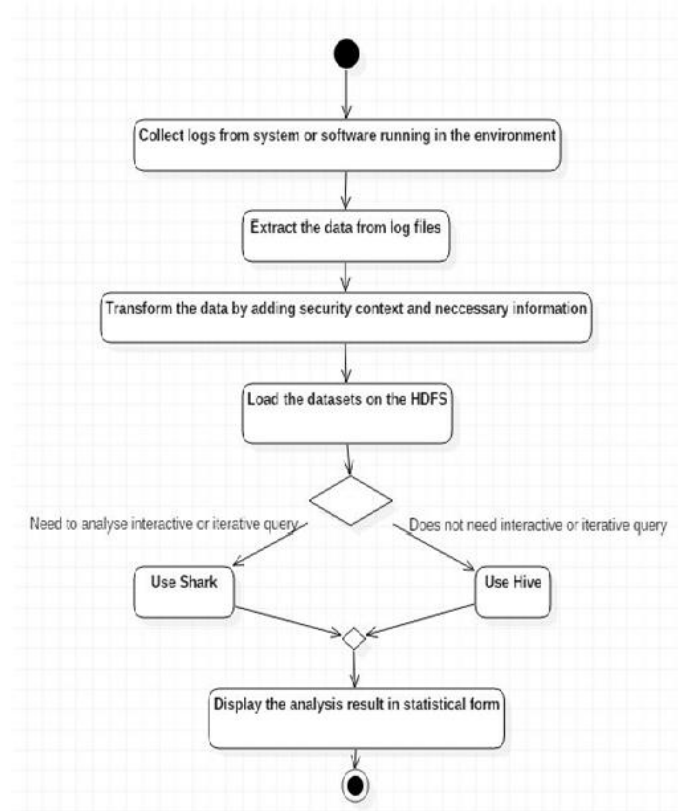


Fig.5: State diagram

7. CONCLUSION

Log analysis is the process to gather information of the number of access users, user behavior, and system operation status, etc. We design and implement a lightweight distributed framework, consisting of a minimized set of components. The framework is different from the general ones and is specially designed for log analysis of web server logs. This paper analyzes and compares the respective characteristics of Hadoop and Spark framework and Hive/ Shark. Combining the characteristics that are useful to us of Hadoop as well as Spark we propose a cloud platform analysis model with high stability, availability and efficiency for batch data analysis in comparison with standalone log analysis tools and system simply based on Hadoop or the combination of Hive.

REFERENCES

- [1] Y. Q. Wei, G. G. Zhou, D. Xu, Y. Chen, "Design of the Web Log Analysis System Based on Hadoop", Advanced Materials Research, Vols. 926-930, pp. 2474-2477, 2014
- [2] Xiuqin LIN, Peng WANG, Bin WU "LOG ANALYSIS IN CLOUD COMPUTING ENVIRONMENT WITH HADOOP AND SPARK", Beijing University of Posts and Telecommunications, Beijing 100876.
- [3] Xiaokui Shu, John Smiy 2013.Massive Distributed and Parallel Log Analysis For Organizational Security Industry/University Cooperative Research Center (I/UCRC), and NSF grant CAREER CNS0953638
- [4] Yanish Pradhananga, Shridevi Karande, Chandraprakash Karande, "CBA: Cloud-based Bigdata Analytics",978-1-4799-6892-3/15 \$31.00 © 2015 IEEE DOI 10.1109/ICCUBEA.2015.18
- [5] Markku Hinkka, Teemu Leht, Keijo Heljanko, "Assessing Big Data SQL Frameworks for Analyzing Event Logs",978-1-4673-8776-7/16 \$31.00 © 2016 IEEE DOI 10.1109/PDP.2016.26
- [6] Jaladi Udaya Santhi, Sivaiah Bellamkonda, Dr.N.Gnaneswara Rao, "Analysis of web server log files using Hadoop MapReduce to preprocess the log files and to explore the session identification and network anomalies", 3rd International Conference on Electrical, Electronics, Engineering Trends, Communication, Optimization and Sciences (EEECOS)-2016, 2016.