

Data Mining Algorithm in Cloud Computing Using Map Reduce Framework

Achi Sandeep, K. Rammohan Goud

Assistant Professor, CSE Department, Sri Indu College of Engineering and Technology, JNTU Hyderabad, Hyderabad, India

ABSTRACT

Today's Cloud computing technology has been emerged to manage large data sets efficiently and due to rapid growth of data, large-scale data processing is becoming a major point of information technique. The Hadoop Distributed File System (HDFS) is designed for reliable storage of very large data sets and to stream those data sets at high bandwidth to user applications. In a large cluster, hundreds of servers both host directly attached storage and execute user application tasks. By distributing storage and computation across many servers, the resource can grow on demand while remaining economical at every size. Map Reduce has been widely used for large-scale data analysis in the Cloud. Hadoop is an open source implementation of Map Reduce which can achieve better performance with the allocation of more compute nodes from the cloud to speed up computation; however, this approach of "renting more nodes" isn't cost effective in a pay-as-you-go environment.

Keywords: Cloud Computing, Distributed Data Mining, Hadoop, Hadoop Distributed File System, Map Reduce.

I. INTRODUCTION

These days large amount of data is created every day so with this rapid explosion of data we are moving towards the terabytes to petabytes. This trend creates the demand for the advancement in data collection and storing technology. Hence, there is a growing need to run data mining algorithm on very large data sets. Cloud computing is a new business model containing pool of resources constituting large number of computers. It distributes the computation task to its pool of resources so that applications can obtain variety of software services on demand. Another feature of cloud computing is that it provides unlimited storage and computing power which leads us to mine mass amount of data.

Hadoop is the software framework for writing applications that rapidly process large amount of data in parallel on large clusters of compute nodes. It provides a distributed file system and a framework for the analysis and transformation of very large data sets using the Map Reduce paradigm. The volume of data,

collectively called data sets, generated by the application is very large. Therefore, there is a need of processing large data sets efficiently.

Map Reduce is a generic execution engine that parallelizes computation over a large cluster of machines. An important characteristic of Hadoop is the partitioning of data and computation across many hosts, and executing application computations in parallel close to their data. A Hadoop cluster scales Computation capacity, storage capacity and IO bandwidth by simply adding commodity servers.

Big data has been used to convey the all sorts of concepts, including huge quantities of data (with respect to volume, velocity, and variety), social media analytics, next generation data management capabilities, real-time data, and much more. Now organizations are starting to understand and explore how to process and analyze a vast array of information in new ways.

Data mining is the process of finding correlations or patterns among fields in large data sets and building up the knowledge base, based on the given constraints.

The overall goal of data mining is to extract knowledge from an existing data set and transform it into a human understandable structure for further use. This process is often referred to as Knowledge Discovery in data sets (KDD). It encompasses data storage and access, scaling algorithms to very large data sets and interpreting results. The data cleansing and data access process included in data warehousing facilitate the KDD process.

Based on the increasing demand for parallel computing environment of cloud and parallel mining algorithm, we study different mining algorithms. Association rule based algorithm, Priority algorithm, is improved in order to combine it with the Map Reduce programming model of cloud and mine large amount of data. With emerging trends in Cloud-computing, data mining enters a new era, which can have a new implementation. We can use cloud-computing techniques with Data mining to reach high capacity and efficiency by using parallel computational nature of the cloud. As Map Reduce provides good parallelism for the computation, it's very suitable for us to implement data mining system based on Map Reduce.

In a distributed computing environment bunch of loosely coupled processing nodes are connected by the network. Each node contributes into the execution or distribution / replication of data. It is referred as a cluster of nodes. There are various methods of setting up a cluster, one of which is usually referred to as cluster framework. Such frameworks enforce the setting up processing and replication nodes for data. Examples are Aneka and Apache Hadoop (called Map / Reduce). The other methods involve setting up of cluster nodes on ad-hoc basis and not being bound by a rigid framework. Such methods just involve a set of API calls for remote method invocation (RMI) as a part of inter-process communication.

The method of setting up a cluster depends upon the data densities and up on the scenarios listed below:

1. The data is generated at various locations and needs to be accessed locally most of the time for processing.
2. The data and processing is distributed to the machines in the cluster to reduce the impact of

any particular machine being overloaded that damages its processing.

II. RELATED WORK

Distributed Data Mining in Peer-to-Peer Networks (P2P) ^[1] offers an overview of the distributed data-mining applications and algorithms for peer-to-peer environments. It describes both exact and approximate distributed data-mining algorithms that work in a decentralized manner. It illustrates these approaches for the problem of computing and monitoring clusters in the data residing at the different nodes of a peer-to-peer network. This paper focuses on an emerging branch of distributed data mining called peer-to-peer data mining. It also offers a sample of exact and approximate P2P algorithms for clustering in such distributed environments

Architecture for data mining in distributed environments ^[2] describes system architecture for scalable and portable distributed data mining applications. This approach presents a document metaphor called *Living Documents* for accessing and searching for digital documents in modern distributed information systems. The paper describes a corpus linguistic analysis of large text corpora based on collocations with the aim of extracting semantic relations from unstructured text.

Distributed Data Mining of Large Classifier Ensembles ^[3] presents a new classifier combination strategy that scales up efficiently and achieves both high predictive accuracy and tractability of problems with high complexity. It induces a global model by learning from the averages of the local classifiers output. The effective combination of large number of classifiers is achieved this way.

Map-Reduce for Machine Learning on Multi core ^[4] discuss the ways to develop a broadly applicable parallel programming paradigm that is applicable to different learning algorithms. By taking advantage of the summation form in a map-reduce framework, this paper tries to parallelize a wide range of machine learning algorithms and achieve a significant speedup on a dual processor cores.

III. STUDY OF DATA MINING ALGORITHMS

A. K-Means Clustering

The K-mean clustering algorithm [7] is used to cluster the huge data set into smaller cluster.

In data mining, k-means clustering is a method of cluster analysis, which aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean. The problem is computationally difficult (NP-hard), however there are efficient heuristic algorithms that are commonly employed and converge fast to a local optimum. These are usually similar to the expectation-maximization algorithm for mixtures of Gaussian distributions via an iterative refinement approach employed by both algorithms. Additionally, they both use cluster centers to model the data, however k-means clustering tends to find clusters of comparable spatial extent, while the expectation-maximization mechanism allows clusters to have different shapes. The most common algorithm uses an iterative refinement technique. Due to its ubiquity, it is often called the k-means algorithm; which is also referred to as Lloyd's algorithm, particularly in the computer science community.

1. Algorithm:

Given an initial set of k means $m_1(1) \dots m_k(1)$, The algorithm proceeds by alternating between two steps:

1. **Assignment step:** Assign each observation to the cluster with the closest mean.
2. **Update step:** Calculate the new means to be the Centroid of the observations in the cluster.

In the beginning we determine number of cluster K and we assume the centroid or center of these clusters. We can take any random objects as the initial centroid or the first K objects in sequence can also serve as the initial centroid.

Then the K means algorithm will do the three steps below until convergence. Iterate until stable:

- ✓ Determine the centroid coordinate
- ✓ Determine the distance of each object to the centroid
- ✓ Group the object based on minimum distance

2. Euclidean Distance

In mathematics, the Euclidean distance or Euclidean metric is the "ordinary" distance between two points that one would measure with a ruler, and is given by the Pythagorean formula. By using this formula as distance, Euclidean space (or even any inner product space) becomes a metric space. The Euclidean distance between point's p and q is the length of the line segment connecting them. In Cartesian coordinates, if $p = (p_1, p_2, \dots, p_n)$ and $q = (q_1, q_2, \dots, q_n)$ are two points in Euclidean n - space, then the distance from p to q or from q to p is given by:

Finding Frequent Item set by using Apriori data mining algorithm:

Require Items $I = \{i_1, i_2, i_n\}$, data set D , user-defined support threshold

Ensure $F(D, _):=$ Frequent sets from D w R T. that particular threshold

1. $C_1 := \{\{i\} | i \in I\}$ //Start with singleton sets
2. $k := 1$
3. while $C_k \neq \{\}$ do
4. //Pruning Part
5. for all transactions $(tid, I) \in D$ do
6. for all candidate sets $X \in C_k$ do
7. if $X \subseteq I$ then
8. $support(X)++$
9. end if
10. end for
11. end for //Computes the supports of all candidate sets
12. $F_k := \{X | support(X) \geq _ \}$ //Extracts all frequent sets
13. //Generating Part
14. for all $X, Y \in F_k, X[j] = Y[j]$ for $1 \leq j \leq k-1$, and $X[k] < Y[k]$ do
15. $I = X \cup \{Y[k]\}$ //Join step
16. if $\exists J \subseteq I, |J| = k : J \in F_k$ then
17. $C_{k+1} := C_{k+1} \cup I$ //Prune step
18. end if
19. end for
20. $k++$
21. end while

In short we are trying to perform following steps:

1. Generate C_{k+1} , candidates of frequent itemsets of size $k + 1$, from the frequent item sets of size k .
2. Scan the database and calculate the support of each candidate of frequent item sets.
3. Add those item sets that satisfies the minimum $d(p,q) = d(q,p) = \sqrt{(q_1-p_1)^2 + (q_2-p_2)^2 + \dots + (q_n-p_n)^2}$

B. Apriori

Apriori [7] is one of the key algorithms to generate frequent item sets. Analyzing frequent item set is a crucial step in Analyzing structured data and in finding association relationship between items. This stands as an elementary foundation to supervised learning.

Association – It aims to extract interesting correlations, frequent patterns associations or casual structures among sets of items in the transaction databases or other data repositories and describes association relationship among different attributes.

support requirement to F_{k+1} .

The Apriori algorithm is shown above in line 13 generates C_{k+1} from F_k in the following two step process:

a. Join step: Generate R_{k+1} , the initial candidates of frequent item sets of size $k + 1$ by taking the union of the two frequent item sets of size k , P_k and Q_k that have the first $k-1$ elements in common.

$$R_{k+1} = P_k \cup Q_k = \{i_{tem1}, \dots, i_{temk-1}, i_{temk}, i_{temk_}\}$$

$$P_k = \{i_{tem1}, i_{tem2}, \dots, i_{temk-1}, i_{temk}\} \quad Q_k = \{i_{tem1}, i_{tem2}, \dots, i_{temk-1}, i_{temk_}\}$$

where, $i_{tem1} < i_{tem2} < \dots < i_{temk} < i_{temk_}$.

b. Prune step: Check if all the item sets of size k in R_{k+1} are frequent and generate C_{k+1} by removing those that do not pass this requirement from R_{k+1} . This is because any subset of size k of C_{k+1} that is not frequent cannot be a subset of a frequent item set of size $k + 1$.

Function subset in line 5 finds all the candidates of the frequent item sets included in transaction t . Apriori, then, calculates frequency only for those candidates generated this way by scanning the database. It is

evident that Apriori scans the database at most $k_{max}+1$ times when the maximum size of frequent item sets is set at k_{max} .

The Apriori achieves good performance by reducing the size of candidate sets. However, in situations with very many frequent item sets, large item sets, or very low minimum support, it still suffers from the cost of generating a huge number of candidate sets.

IV. RESEARCH METHODOLOGY

A. Cloud computing:

It consists of shared computing resources which are opposed to local servers or devices. Users [6] can pay on the basis of resource usage as timely basis. The major goal of cloud computing is to provide easily scalable access to computing resources and IT (Information Technology) services for achieving better performance. Cloud computing basically provides three different types of service based architectures are SaaS, PaaS, and IaaS.

SaaS (Software as-a-service): It offers application as a service on the internet.

PaaS (Platform as-a-service): This is to be used by developers for developing new applications.

IaaS (Infrastructure as-a-service): It is basically deals by providers to provide features on-demand Utility.

Table 1. Feature Comparison Of Commercial Offerings For Cloud Computing.

Properties	Amazon	Google	Microsoft	Manjras
	EC2	App Engine	Azure	oft. Aneka
Service Type	IaaS	IaaS - PaaS	IaaS - PaaS	PaaS
Support for (value offer)	Compute/ Storage	Compute (web applications)	Compute	Compute
Value added service provider	Yes	Yes	Yes	Yes
User Access Interface	Web API Command Line Tool	Web API Command Line Tool	Azure Web Portal	Web APIs Custom GUI
Virtualization	OS on Xen Hypervisor	Application Container	Service Container	Service Container
Platform (OS & runtime)	Linux, Windows	Linux	.NET on Windows	.NET on Windows, Mono, Linux
Deployment model	Customizable VM	Web apps (Python, Java, Ruby)	Azure Services	Applications (C#, C++, VB)
If PaaS, ability to deploy on 3rd party IaaS	N.A.	No	No	Yes

B. Map Reduce:

Map Reduce [5] is a programming model for processing large data sets, and the name of an implementation of the model by Google. Map Reduce is typically used to perform distributed computing on clusters of computers. The model is inspired by map and reduces functions commonly used in functional programming, although their purpose in the Map Reduce framework is not the same as their original forms. Map Reduce libraries have been written in many programming languages. A popular free implementation is Apache Hadoop.

Map Reduce is a framework for processing the parallelizable problems across huge data sets using a large number of computers (nodes), collectively referred to as a cluster (if all nodes are on the same local network and use similar hardware) or a grid (if the nodes are shared across geographically and administratively distributed systems, and use more heterogeneous hardware). Computational processing can occur on data stored either in a file system (unstructured) or in a database (structured). Map Reduce can take advantage of locality of data, processing data on or near the storage assets to decrease transmission of data.

1) **“Map” step:** The master node takes the input, divides it into smaller sub-problems, and

distribute them to worker nodes. A worker node may do the again in t , leading to a multilevel tree structure. The work r node processes the smaller problem, and passes the answer back to its master node.

1. **“Reduce” step:** The master node then collects the answers to all the sub-problems and combines them in some way to form the output -the answer to the problem it was originally trying to solve.

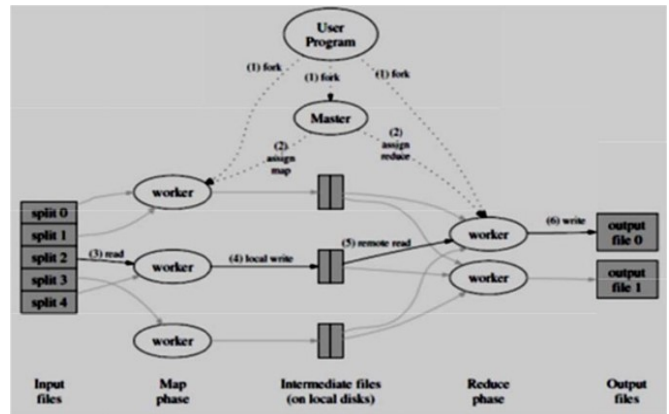


Figure 1. execution over view

Data flow of the system is given below; the frozen part of the Map Reduce framework is a large distributed sort. The above figure consists of following parts:

- i. **Input reader:** It divides the input into appropriate (in practice typically 64 MB to 512 MB as PE HDFS) and the framework assigns one split to one Map function. The input reader reads the data from stable storage (typically a in our case Hadoop distributed file system) and generates key/value pairs
- ii. **Map function:** Each Map function takes a series of key/value pairs, processes each, and generates zero or more output key/value pairs. The input and output types of the map can be and often are different from each other.
- iii. **Partition function:** Each Map function output is allocated to a particular reducer by the application's partition function for sharing purposes. The partition function is given the key and the number of reducers and returns the index of the desired reducer.

d) Comparison function:

The input for every Reduce is fetched from the machine where the Map run and sorted using the application's compare function.

e) Reduce function:

The framework calls the application's Reduce function for each unique key in the sorted order. It also iterates through the values that are associated with that key and produces zero or more outputs.

f) Output writer:

It writes the output of the Reduce function to stable storage, usually a Hadoop distributed file system.

As an example, the illustrative problem of counting the average word length of every word occurrences in

- a large collection of documents in Map Reduce is represented as following: The input key/value to the
- b Map function is a document name, and its contents. The function scans through the document and emits each word plus the associated word length of the occurrences of that word in the document. Shuffling groups together occurrences of the same word in all documents, and passes them to the Reduce function. The Reduce function sums up all the word length for all occurrences. Then divide it by the count of that word and emits the word and its overall average word length of every word occurrences.

Example: Consider the problem of counting the average word length in a large collection of documents. The user would write code similar to the following pseudo-code:

function map(String key, String value):

key: document name //value: document contents

for each word w in value:

Emit Intermediate(w, word length);

function reduce(String key, Iterator values)

key: word

values: list of counts

double sum = 0, count = 0, result = 0; for each v in values:

sum += Parse Int(v); count++;

result = sum / count;

Emit(w, As Double(result));

Here, each document is split into words, and each word length is counted by the map function, using the word as the result key. The framework puts together all the pairs with the same key and feeds them to the same call to reduce, thus this function just needs to sum all of its input values to find the total appearances of that word. Then for finding average word length we divide the sum by the count of that word.

V. CONCLUSION

There are many new technologies emerging at a rapid rate, each with technological advancements and with the potential of making ease in use of technology. However one must be very careful to understand the limitations and security risks posed in utilizing these Technologies. Neither Map Reduce-like software, nor Parallel databases are ideal solutions for data analysis in the cloud. Hybrid solution that combines the fault tolerance, heterogeneous cluster, and ease of use out-of-the-box capabilities of Map Reduce with the efficiency, performance, and tool plug ability of shared-nothing parallel systems could have a significant impact on the cloud market. We will work on bringing together ideas from Map Reduce and data mining algorithms, also to combine the advantages of Map Reduce-like software with the efficiency and shared work advantages that come with loading data and creating performance enhancing data structures.

VI. REFERENCES

- [1]. Souptik Datta, Kanishka Bhaduri, Chris Giannella, Ran Wolff, and Hillol Kargupta, Distributed Data Mining in Peer-to-Peer Networks, University of Maryland, Baltimore County, Baltimore, MD, USA, Journal IEEE Internet Computing archive Volume 10 Issue 4, Pages 18-26, July 2006.
- [2]. Mafruz Zaman Ashrafi, David Taniar, and Kate A. Smith, A Data Mining Architecture for Distributed Environments, pages 27-34, Springer-Verlag London, UK, 2007.
- [3]. Grigorios Tsoumakas and Ioannis Vlahavas, Distributed Data Mining of Large Classifier Ensembles, SETN-2008, Thessaloniki, Greece, Proceedings, Companion Volume, pp. 249-256, 11-12 April 2008.

- [4]. Cheng-Tao Chu et.al.,Map-Reduce for Machine Learning on Multicore,CS Department,Stanford University,Stanford,CA,2006.
- [5]. Jeffrey Dean and Sanjay Ghemawat,Map Reduce:Simplified data processing on large clusters.InOSDI,pages 137-150,2004.
- [6]. Daniel J.Abadi, Yale University, DataManagement in the Cloud: Limitations and Opportunities, Bulletin of the IEEE ComputerSociety Technical Committee on Data Engineering 2009
- [7]. "Top 10 algorithms in data mining", Springer-Verlag London Limited 2007
- [8]. James I.Johnson,SQL in the Clouds,IEEE journal Cloud Computing,2009.