

Market Basket Analysis Algorithm with Map/Reduce of Cloud Computing

Jongwook Woo

Computer Information Systems Department
California State University
Los Angeles, CA

Yuhang Xu

Computer Information Systems Department
California State University
Los Angeles, CA

Abstract – *Map/Reduce approach has been popular in order to compute huge volumes of data since Google implemented its platform on Google Distributed File Systems (GFS) and then Amazon Web Service (AWS) provides its services with Apache Hadoop platform. Map/Reduce motivates to redesign and convert the existing sequential algorithms to Map/Reduce algorithms for big data so that the paper presents Market Basket Analysis algorithm with Map/Reduce, one of popular data mining algorithms. The algorithm is to sort data set and to convert it to (key, value) pair to fit with Map/Reduce. It is executed on Amazon EC2 Map/Reduce platform. The experimental results show that the code with Map/Reduce increases the performance as adding more nodes but at a certain point, there is a bottle-neck that does not allow the performance gain. It is believed that the operations of distributing, aggregating, and reducing data in Map/Reduce should cause the bottle-neck.*

Keywords: Map/Reduce, Market Basket Analysis, Data Mining, Association Rule, Hadoop, Cloud Computing

1 Introduction

Before Internet and Web did not exist, we did not have enough data so that it was not easy to analyze people, society, and science etc with the limited volumes of data. Contradicting to the past, after Internet and web, it has been more difficult to analyze data because of its huge volumes, that is, tera- or peta-bytes of data. Google faced to the issue as she collected big data and the existing file systems were not sufficient to handle the data efficiently. Besides, the legacy computing power and platforms were not useful for the big data. Thus, she implemented Google File Systems (GFS) and Map/Reduce parallel computing platform, which Apache Hadoop project is motivated from.

Hadoop is the parallel programming platform built on Hadoop Distributed File Systems (HDFS) for Map/Reduce computation that processes data as (key, value) pairs. Hadoop has been receiving highlights for the enterprise computing because business world always has the big data such as log files for web transactions. Hadoop is useful to process such big data for business intelligence so that it has been used in data mining for past few years. The era of

Hadoop means that the legacy algorithms for sequential computing need to be redesigned or converted to Map/Reduce algorithms. Therefore, in this paper, a Market Basket Analysis algorithm in data mining with Map/Reduce is proposed with its experimental result in Elastic Compute Cloud (EC2) and (Simple Storage Service) S3 of Amazon Web Service (AWS).

People have talked about Cloud Computing that is nothing else but the services we have used for several years: hosting service, web email service, document sharing service, and map API service etc. It is categorized into Software as a Service (SaaS), Platform as a Service (PaaS), and Infrastructure as a Service (IaaS). SaaS is to use a service via Internet without installing or maintaining the software, for example, web email services. PaaS is to have a computing or storage service without purchasing hardware or software, for example, hosting services. IaaS is to have utility computing service that is similar to SaaS but to purchase only the amount of time to use the service like AWS [6, 7]. AWS provides S3, EC2, and Elastic MapReduce services for Map/Reduce computation as IaaS and SaaS in cloud computing.

In this paper, section 2 is related work. Section 3 describes Map/Reduce and Hadoop as well as other related projects. Section 4 presents the proposed Map/Reduce algorithm for Market Basket Analysis. Section 5 shows the experimental result. Finally, section 6 is conclusion.

2 Related Work

Association Rule or Affinity Analysis is the fundamental data mining analysis to find the co-occurrence relationships like purchase behavior of customers. The analysis is legacy in sequential computation and many data mining books illustrate it.

Aster Data has SQL MapReduce framework as a product [9]. Aster provides *nPath SQL* to process big data stored in the DB. Market Basket Analysis is executed on the framework but it is based on its SQL API with MapReduce Database.

As far as we understand, there is not any other to present Market Basket Analysis algorithms with Map/Reduce. The approach in the paper is to propose the

algorithm and to convert data to (key, value) pair and execute the code on Map/Reduce platform.

3 Map/Reduce in Hadoop

Map/Reduce is an algorithm used in Artificial Intelligence as functional programming. It has been received the highlight since re-introduced by Google to solve the problems to analyze huge volumes of data set in distributed computing environment. It is composed of two functions to specify, “Map” and “Reduce”. They are both defined to process data structured in (key, value) pairs.

3.1 Map/Reduce in parallel computing

Map/Reduce programming platform is implemented in the Apache Hadoop project that develops open-source software for reliable, scalable, and distributed computing. Hadoop can compose hundreds of nodes that process and compute peta- or tera-bytes of data working together. Hadoop was inspired by Google's MapReduce and GFS as Google has had needs to process huge data set for information retrieval and analysis [1]. It is used by a global community of contributors such as Yahoo, Facebook, and Twitters. Hadoop's subprojects include Hadoop Common, HDFS, MapReduce, Avro, Chukwa, HBase, Hive, Mahout, Pig, and ZooKeeper etc [2].

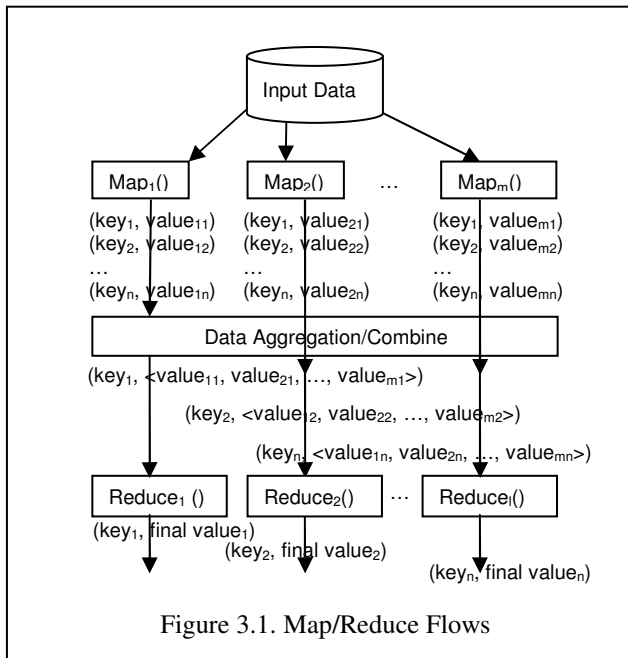


Figure 3.1. Map/Reduce Flows

The map and reduce functions run on distributed nodes in parallel. Each map operation can be processed independently on each node and all the operations can be performed in parallel. But in practice, it is limited by the data source and/or the number of CPUs near that data. The reduce functions are in the similar situation because they are from all the output of the map operations. However, Map/Reduce can handle significantly huge data sets since

data are distributed on HDFS and operations move close to data for better performance [5].

Hadoop is restricted or partial parallel programming platform because it needs to collect data of (key, value) pairs as input and parallelly computes and generates the list of (key, value) as output on map/reduce functions. In map function, the master node parts the input into smaller sub-problems, and distributes those to worker nodes. Those worker nodes process smaller problems, and pass the answers back to their master node. That is, map function takes inputs (k1, v1) and generates <k2, v2> where < > represents list or set. Between map and reduce, there is a combiner that resides on map node, which takes inputs (k2, <v2>) and generates <k2, v2>.

In reduce function, the master node takes the answers to all the sub-problems and combines them in some way to get the output, the answer to the problem [1, 2]. That is, reduce function takes inputs (k2, <v2>) and generates <k3, v3>. Figure 3.1 illustrates Map/Reduce control flow where each value_{mn} is simply 1 and gets accumulated for the occurrence of items together in the proposed Market Basket Analysis Algorithm.

3.2 Database for Big Data

Input/Output files are processed on HDFS instead of using HBase DB in the paper. However, as HBase is interesting and will be integrated with the algorithm in the future, the section briefly introduces HBase.

There are some drawbacks when we use RDBMS to handle huge volumes of data, like impossible deleting, slow inserting, and random failing. HBase on HDFS is distributed database that supports structured data storage for horizontally scalable tables. It is column oriented semi-structured data store.

It is relatively easy to integrate with Hadoop Map/Reduce because HBase consists of a core map that is composed of keys and values - each key is associated with a value. Users store data rows in labeled tables. A data row has a sortable key and an arbitrary number of columns. The table is stored sparsely, so that rows in the same table can have different columns.

Using the legacy programming languages such as Java, PHP, and Ruby, we can put data in the map as Java JDBC does for RDBMS. The file storage of HBase can be distributed over an array of independent nodes because it is built on HDFS. Data is replicated across some participating nodes. When the table is created, the table's column families are generated at the same time. We can retrieve data from HBase with the full column name in a certain form. And then HBase returns the result according to the given queries as SQL does in RDBMS [10].

3.3 The Issues of Map/Reduce

Although there are advantages of Map/Reduce, for some researchers and educators, it is:

1. A giant step backward in the programming paradigm for large-scale data intensive applications
2. Not new at all - it represents a specific implementation of well known techniques developed tens of years ago, especially in Artificial Intelligence
4. Data should be converted to the format of (key, value) pair for Map/Reduce, which misses most of the features that are routinely included in current DBMS
5. Incompatible with all of the tools or algorithms that have been built [4].

However, the issues clearly show us not only the problems but also the opportunity where we can implement algorithms with Map/Reduce approach, especially for big data set. It will give us the chance to develop new systems and evolve IT in parallel computing environment. It started a few years ago and many IT departments of companies have been moving to Map/Reduce approach in the states.

4 Market Basket Analysis Algorithm

Market Basket Analysis is one of the Data Mining approaches to analyze the association of data set. The basic idea is to find the associated pairs of items in a store when there are transaction data sets as in Figure 4.1.

If store owners list a pair of items that are frequently occurred, s/he could control the stocks more intelligently, to arrange items on shelves and to promote items together etc. Thus, s/he should have much better opportunity to make a profit by controlling the order of products and marketing.

Transaction 1: cracker, icecream, beer
 Transaction 2: chicken, pizza, coke, bread
 Transaction 3: baguette, soda, hering, cracker, beer
 Transaction 4: bourbon, coke, turkey
 Transaction 5: sardines, beer, chicken, coke
 Transaction 6: apples, peppers, avocado, steak
 Transaction 7: sardines, apples, peppers, avocado, steak
 ...

Figure 4.1 Transaction data at a store

Total number of Items: 322,322
 Ten most frequent Items:

cracker, beer	6,836
artichok, avocado	5,624
avocado, baguette	5,337
bourbon, cracker	5,299
baguette, beer	5,003
corned, hering	4,664
beer, hering	4,566
...	

Figure 4.2 Top 10 pair of items frequently occurred at store

For example, people have built and run Market Basket Analysis codes – sequential codes - that compute the top 10 frequently occurred pair of transactions as in Figure 4.2. At the store, when customers buy a cracker, they purchase a beer as well, which happens 6,836 times and bourbon as well in 5,299 times. Thus, the owner can refer to the data to run the store.

4.1 Data Structure and Conversion

The data in Figure 4.1 is composed of the list of transactions with its transaction number and the list of products. For Map/Reduce operation, the data set should be structured with (key, value) pairs. The simplest way used in the paper is to pair the items as a key and the number of key occurrences as its value in the basket, especially for all transactions, without the transaction numbers. Thus, Figure 4.1 can be restructured as Figure 4.3 assuming collecting a pairs of items in order 2 – two items as a key.

```
< (cracker, icecream), (cracker, beer) >
< (chicken, pizza), (chicken, coke), (chicken, bread) >
< (baguette, soda), (baguette, hering), (baguette,
cracker), (baguette, beer) >
< (bourbon, coke), (bourbon, turkey) >
< (sardines, beer), (sardines, chicken), (sardines, coke)
>
...
```

Figure 4.3 Data Set restructured for Map/Reduce

However, if we select the two items in a basket as a key, there should be incorrect counting for the occurrence of the items in the pairs. As shown in Figure 4.4, transactions *n* and *m* have the items (cracker, icecream, beer) and (icecream, beer, cracker), which have the same items but in different order.

Transaction *n*: cracker, icecream, beer
 Transaction *m*: icecream, beer, cracker

Convert to (key, value): cross operation

Transaction *n*: ((cracker, icecream), 1),
 ((cracker, beer), 1), ((icecream, beer), 1)
 Transaction *m*: ((icecream, beer), 1),
 ((icecream, cracker), 1), ((beer, cracker), 1)

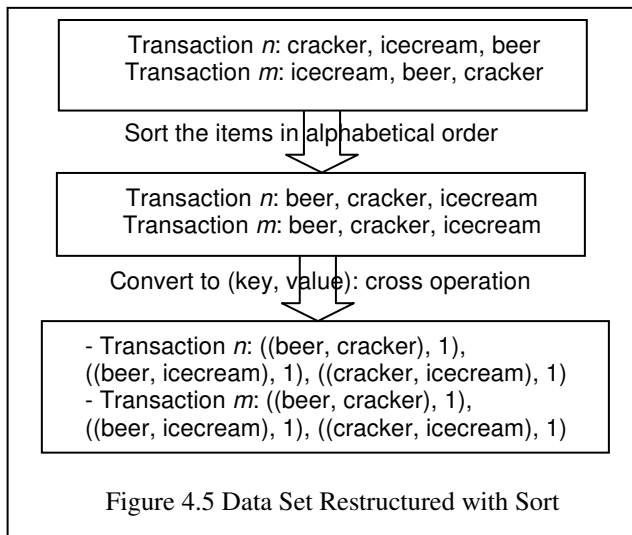
Figure 4.4 Data Set Restructured for the same list

That is, for (cracker, icecream, beer), the possible pair of items in (key, value) are ((cracker, icecream), 1), ((cracker, beer), 1), ((icecream, beer), 1). And, for (icecream, beer, cracker), the possible pair of items are

((icecream, beer), 1), ((icecream, cracker), 1), ((beer, cracker), 1).

Therefore, we have total SIX different pair of items that occurs only once respectively, which should be THREE different pairs. That is, keys (cracker, icecream) and (icecream, cracker) are not same even though they are, which is not correct.

We can avoid this issue if we sort the transaction in alphabetical order before generating (key, value) as shown in Figure 4.5. Now each transaction have the following THREE pair of items ((beer, cracker), 1), ((beer, icecream), 1), ((cracker, icecream), 1). That is TWO different pair of items that occurs twice respectively so that we accumulate the value of the occurrence for these two transactions as follows: ((beer, cracker), 2), ((beer, icecream), 2), ((cracker, icecream), 2), which is correct to count the total number of occurrences.



4.2 The algorithm

The Market Basket Analysis (MBA) Algorithms for Mapper and Reducer are illustrated in Figures 4.6 and 4.7 respectively. Mapper reads the input data and creates a list of items for each transaction. As a mapper of a node reads each transaction on Hadoop, it assigns mappers to number of nodes, where the assigning operation in Hadoop is hidden to us. For each transaction, its time complexity is $O(n)$ where n is the number of items for a transaction.

Then, the items in the list are sorted to avoid the duplicated keys as shown in Figures 4.4 and 4.5. Its time complexity is $O(n \log n)$ on merge sort. Then, the sorted items should be converted to pairs of items as keys, which is a cross operation in order to generate cross pairs of the items in the list as shown in Figures 4.4 and 4.5. Its time complexity is $O(n \times m)$ where m is the number of pairs that occurs together in the transaction. Thus, the time complexity of each mapper is $O(n + n \log n + n \times m)$.

```

1: Reads each transaction of input file and generates
the data set of the items:
(<V1>, <V2>, ..., <Vn>) where <Vn>: (vn1, vn2,... vnm)

2: Sort all data set <Vn> and generates sorted data set
<Un>:
(<U1>, <U2>, ..., <Un>) where <Un>: (un1, un2,... unm)

3: Loop While <Un> has the next element;
  note: each list Un is handled individually
  3.1: Loop For each item from un1 to unm of <Un> with
  NUM_OF_PAIRS
    3.a: generate the data set <Yn>: (yn1, yn2,... ynl);
    ynl: (unx, uny) is the list of self-crossed pairs of
    (un1, un2,... unm) where unx ≠ uny
    3.b: increment the occurrence of ynl;
    note: (key, value) = (ynl, number of occurrences)
  3.2: End Loop For
4. End Loop While

5. Data set is created as input of Reducer: (key,
<value>) = (ynl, <number of occurrences>)
  
```

Figure 4.6. MBA Algorithm for Mapper

The reducer is to accumulate the number of values per key. Thus, its time complexity is $O(v)$ where v is the number of values per key.

```

1: Read (ynl, <number of occurrences>) data from
multiple nodes

2. Add the values for ynl to have (ynl, total number of
occurrences)
  
```

Figure 4.7. MBA Algorithm for Reducer

4.3 The code

The *ItemCount.java* code is implemented on Hadoop 0.20.2 and 0.21.0 and executable on stand-alone and clustered modes. The code generates the top 10 associated items that customers purchased together as shown in Figure 4.2. Anyone can download the files to test it, which takes the sample input "AssociationsSP.txt" as introduced in the blog [8]. The sample input has 1,000 transactions with data as shown in Figure 4.1.

5 Experimental Result

We have 5 transaction files for the experiment: 400 MB (6.7M transactions), 800MB (13M transactions), 1.6 GB (26M transactions). Those are run on small instances of AWS EC2 which allows to instantiate number of nodes requested, where each node is of 1.0-1.2 GHz 2007 Opteron or Xeon Processor, 1.7GB memory, 160GB storage on 32 bits platform. The data are executed on 2, 5, 10, 15, and 20

nodes respectively and its execution times are shown in Table 5.1. For 13 and 26 Mega transactions of 2 nodes, it took too long to measure the execution times so that we do not execute them and its times are Not Applicable (NA) in the Table 5.1.

Trax Nodes	6.7M (400MB)	13M (800MB)	26M (1.6GB)
2	9,133	NA	NA
5	5,544	8,717	15,963
10	2,910	5,998	8,845
15	2,792	2,917	5,898
20	2,868	2,911	5,671

Table 5.1. Execution time (sec) at Map Task:

The output of the computation in Table 5.2 presents the number of items (total: 1.3G) and keys (total: 212) that are associated pair of items in order 2, especially for data of 26M transactions in file size 1.6GB. And, the 10 most frequently occurred items and its frequency, which is (key, value), are shown.

Total number of keys in order 2: 212
Total number of items: 1,255,922,927

Items Paired (key)	Frequency (value)
cracker, heineken	208,816,643
artichok, avocado	171,794,426
avocado, baguette	163,027,463
bourbon, cracker	161,866,763
baguette, heineken	152,824,775
corned_b, hering	142,469,636
heineken, hering	139,475,906
bourbon, heineken	126,310,383
baguette, cracker	125,699,308
artichok, heineken	125,180,072

Table 5.2. 10 most frequently associated items on 1.6GB 26M transactions

Figure 5.1 based on Table 5.1 is the chart of the experimental result with 400 MB (6.7M transactions), 800MB (13M transactions), 1.6 GB (1,600 MB 26M transactions). The more the nodes are, the faster the computation times are. Since the algorithm is simply to sort the data set and then convert it to (key, value) pairs, the linear result is expected. The performance is linearly increased by some nodes for some transaction data sets but it has the limitation. For 400MB file, there is not much difference among nodes 10, 15 and 20. Similarly, for 800MB and 1.6GB files, there are not many differences between nodes 15 and 20. There is bottleneck in EC2 small instance, which shows that there is a trade-off between the number of nodes and the operations of distributing

transactions data to nodes, aggregating the data, and reducing the output data for each key so that it should not have much performance gain even though adding more nodes for faster parallel computation.

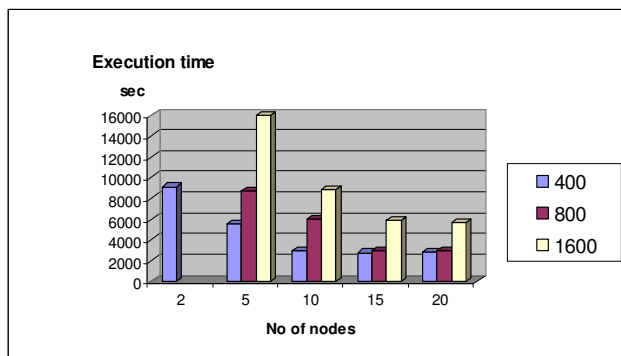


Figure 5.1. Chart for Execution Time

In summary, the experimental data illustrates that even though the number of nodes are added, at a certain point, there is a bottleneck that cannot increase the performance because of the time to distribute, aggregate, and reduce the data to Map/Reduce nodes.

6 Conclusion

Hadoop with Map/Reduce motivates the needs to propose new algorithms for the existing applications that have had algorithms for sequential computation. Besides, it is (key, value) based restricted parallel computing so that the legacy parallel algorithms need to be redesigned with Map/Reduce.

In the paper, the Market Basket Analysis Algorithm on Map/Reduce is presented, which is association based data mining analysis to find the most frequently occurred pair of products in baskets at a store. The data set shows that associated items can be paired with Map/Reduce approach. Once we have the paired items, it can be used for more studies by statically analyzing them even sequentially, which is beyond this paper.

The algorithm has been executed on EC2 small instances of AWS with nodes 2, 5, 10, 15, and 20. The execution times of the experiments show that the proposed algorithm gets better performance while running on large number of nodes to a certain point. However, from a certain point, Map/Reduce does not guarantee to increase the performance even though we add more nodes because there is a bottle-neck for distributing, aggregating, and reducing the data set among nodes against computing powers of additional nodes.

7 Reference

- [1] "MapReduce: Simplified Data Processing on Large Clusters", Jeffrey Dean and Sanjay Ghemawa, Google Labs, pp. 137-150, OSDI 2004

- [2] Apache Hadoop Project, <http://hadoop.apache.org/>,
- [3] “Building a business on an open source distributed computing”, Bradford Stephens , O'Reilly Open Source Convention (OSCON) 2009, July 20-24, 2009, San Jose, CA
- [4] “MapReduce Debates and Schema-Free”, Woohyun Kim, Coord, March 3 2010
- [5] “Data-Intensive Text Processing with MapReduce”, Jimmy Lin and Chris Dyer, Tutorial at the 11th Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL HLT 2010), June 2010, Los Angeles, California
- [6] “ Introduction to Cloud Computing”, Jongwook Woo, the 10th KOCSEA 2009 Symposium, UNLV, Dec 18-19, 2009
- [7] “The Technical Demand of Cloud Computing”, Jongwook Woo, Korean Technical Report of KISTI (Korea Institute of Science and Technical Information), Feb 2011
- [8] “Market Basket Analysis Example in Hadoop, <http://dal-cloudcomputing.blogspot.com/2011/03/market-basket-analysis-example-in.html>”, Jongwook Woo, March 2011
- [9] “SQL MapReduce framework ”, Aster Data, <http://www.asterdata.com/product/advanced-analytics.php>
- [10] Apache HBase, “<http://hbase.apache.org/>”
- [11] “Data-Intensive Text Processing with MapReduce”, Jimmy Lin and Chris Dyer, Morgan & Claypool Publishers, 2010.
- [12] GNU Coord, <http://www.coordguru.com/>
- [13] “Integrated Information Systems Architecture in e-Business”, Jongwook Woo, Dong-Yon Kim, Wonhong Cho, MinSeok Jang, The 2007 international Conference on e-Learning, e-Business, Enterprise Information Systems, e-Government, and Outsourcing, Las Vegas (June 26-29, 2007)