# IMPLEMENTATION OF AM_EBIZ ALGORITHM FOR FINDING FREQUENT ITEMSETS

## Dr.Shikha Dubey[*]

**Abstract**

In recent years, knowledge discovery has been focused on advancements with the vision of benefits to the corporate sector. This research work explicates the concepts of data mining, especially the importance of Frequent Pattern Mining (FPM) and Hadoop MapReduce programming model**.** The organization sectors such as social networks, web browsing, telecommunication, utilities, transportation, insurance, banking and many others generate vast amount of data every day. The proposed work for frequent item set mining of large volume of data is different from all the previous works, since it uses AM e_biz Algorithm with map reduce programming model. It provides compressed storage facility, reduced execution time. It is applicable in many of data mining applications like market-basket problem, decision support, web usage mining, and bioinformatics.

*Keywords:*

Association Rule Mining , Hadoop, Frequent itemset, Algorithm, Mapreduce.

[*] **HOD MCA,  D. Y. P. I. M. R Pimpri, Pune-18**
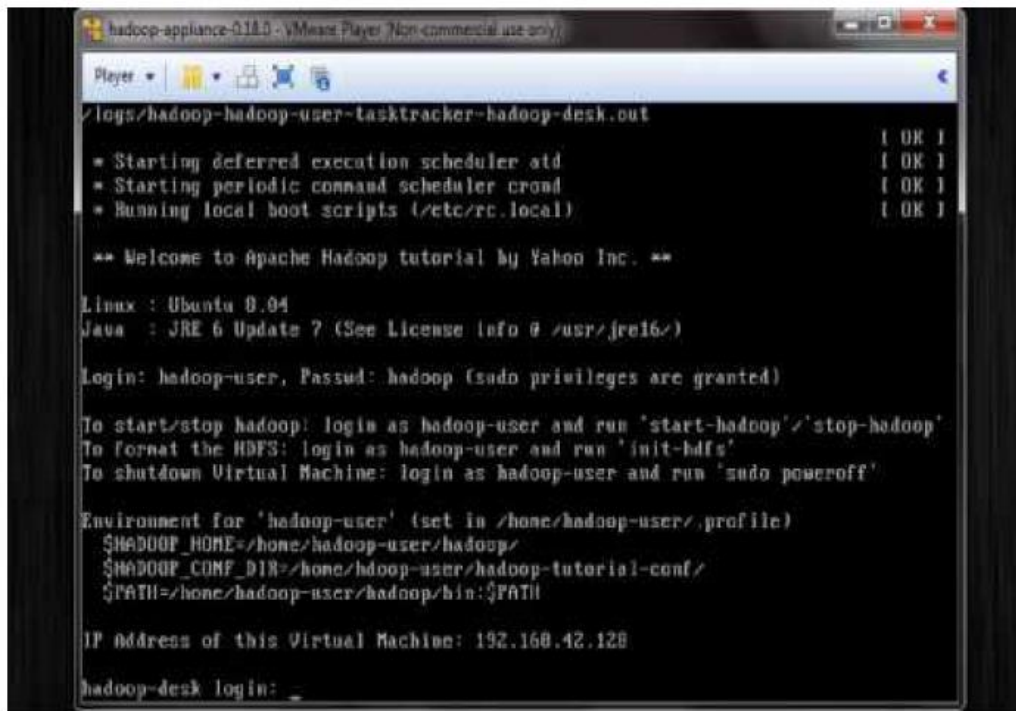
## 1. Introduction

To provide great user experience to users in their day to day activity the Big Data needs to be analyzed. But because of development of IT industries, services, technologies and data, the tremendous measure of complex data is produced from the different sources that can be in different frame. Such complex and massive data is hard to handle and process. So Size, Complexity, and variability of big data are the real difficulties to perceive association rules and frequent item sets.

Various existing data mining techniques are produced displayed to infer association rule and frequently happening item sets, yet with the quick entry of time of big data customary data mining algorithm have not been able to meet huge datasets analysis requirements. There is need to improve performance and accuracy of parallel processing with minimizing execution time complexity. Also assuring the output of a computation is insensitive to changes in any one personal record. So that it will restricting privacy leaks from results. Hence, there is need to provide better frequent item set mining approach using HDFS framework with privacy Preservation techniques.

Formally, frequent item set mining is to perform the following task: we are given a set $B = i1$, in of items, called the item base and a database $T = (t1.........tm)$ of transactions. An item may, for example, represent a product. In this case, the item base represents the set of all products offered by a supermarket. The term item set refers to any subset of the item base B. Each transaction is an item set and may represent, in the supermarket setting, a set of products that has been bought by a customer. As several customers may have bought the same set of products, the total of all transactions must be represented as a vector (as above) or as a multi set. Alternatively, each transaction may be enhanced by a transaction identifier (Tid). Note that the item base B is usually not given explicitly, but only implicitly as the union of all transactions, that is, $B = k$ (-{1.....m} tk.

**2 Research Methodology:** The experiment is carried out in a machine with Intel core i5-M480 with 2.67GHz. And RAM of 4GB. VMware workstation is installed and two VMs are setup. One VM is for Hadoop 0.18.0 setup with configuration of 250GB hard disk and 512 MB main

memory. The other VM is for Eclipse Juno, the virtual configuration is 20GB hard disk and 2 GB main memory. Ubuntu 14.04 is installed in it. Virtual machines are created using VMware Workstation 7.1.3. Necessary packages are installed for eucalyptus and virtual cloud setup is done in VMWare workstation. Eucalyptus platform is installed with Ubuntu packages in a single system with VMWare. Ubuntu image is used in this research work. Bridged connection will be set as 0 because of the same machine. If multiple machines are used IP address of the respective machines are to be used. Hadoop is installed in one virtual machine and eucalyptus is installed in another virtual machine with eclipse environment. Soon after the configuration is completed eucalyptus will be started. Figure 1.1 shows the screenshot of Hadoop login. It has been given to understand the clear picture of virtual machine setup of Hadoop.



**Figure 1.1 Screenshot of Hadoop login**

**Figure 1.2 gives the screenshot of Hadoop environment that is appearing soon after the setup process of Hadoop login.**

**Figure 1.2 Screenshot of Hadoop screen**

Environment variables and IP address of the virtual machine are being set. The subsequent figure 1.3 presents with the Hadoop screen shot.



**Figure 1.3 Screenshot of Hadoop process running state**

It is showing the Name Node, Secondary Name Node, Job Tracker and Task Trackers running. The algorithms are implemented which gives most frequent item set within very less time.
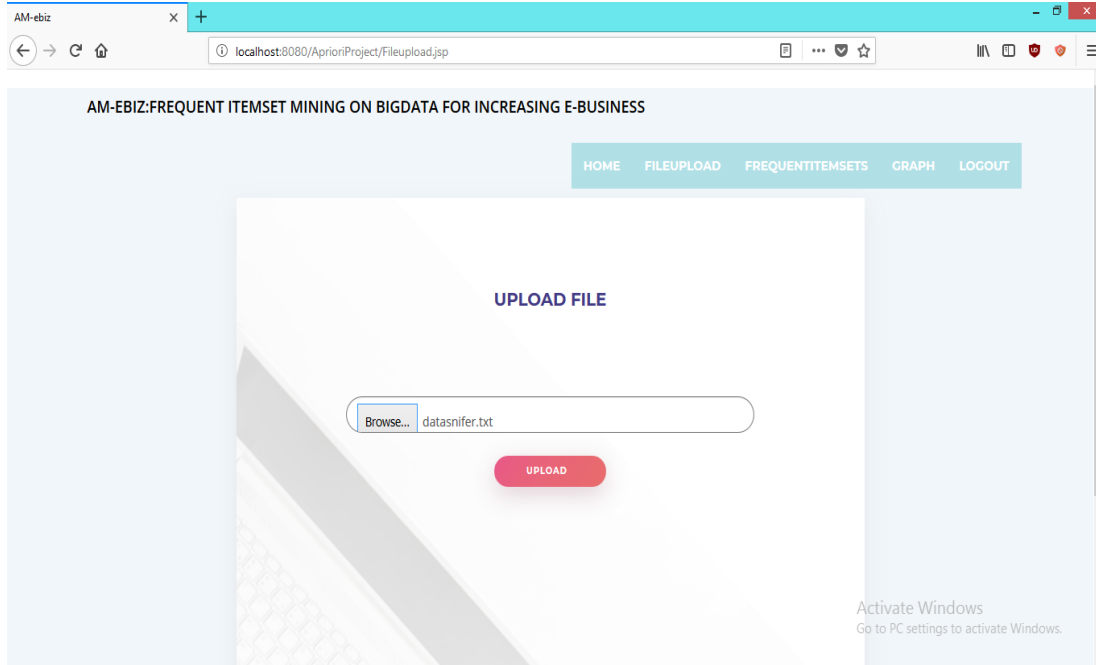
## 3. Results and Analysis

The following figure shows the main page of the system.



**Figure 1.4: Main page AM_ebiz**

After this the system requests for uploading the data which can be browsed from the available files which is shown in Figure 1.5.Figure 1.6 shows the choice of dates to be selected by the user. Later the figures shows how 1 item frequent item set ,2 item frequent item set, three item frequent item set is retrieved.

**Figure 1.5: Dataset Uploading Phase**



**Figure 1.6: Find Frequent item set from specific Range**

**Figure 1.7: First Phase extracted item set group**



**Figure 1.8: Two extracted pairs of item set**
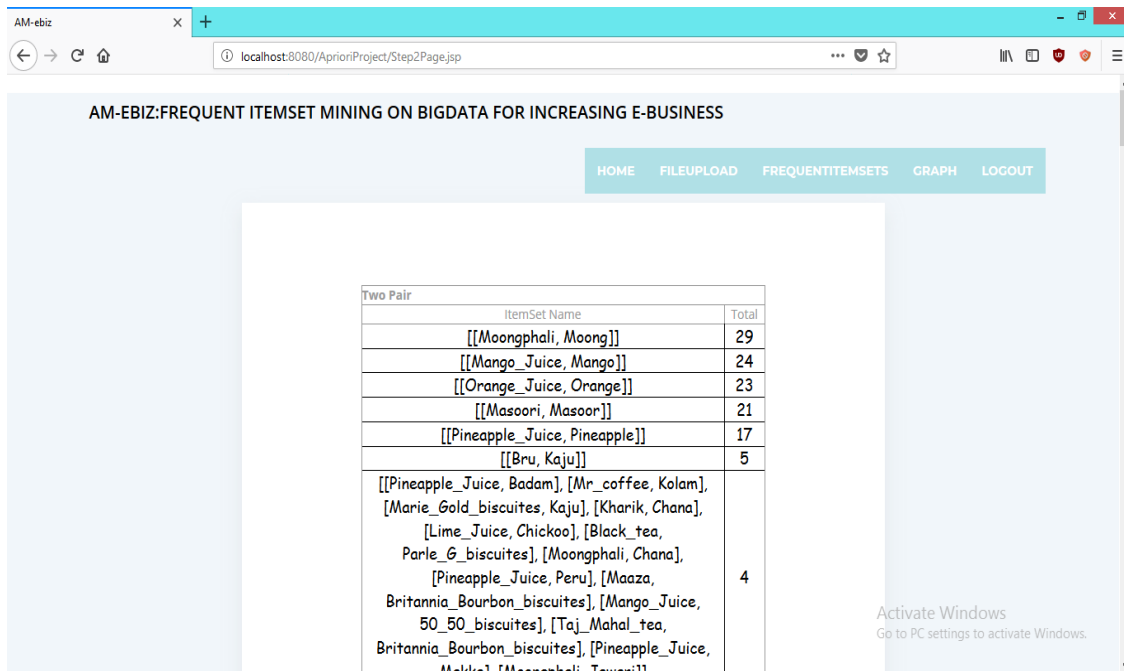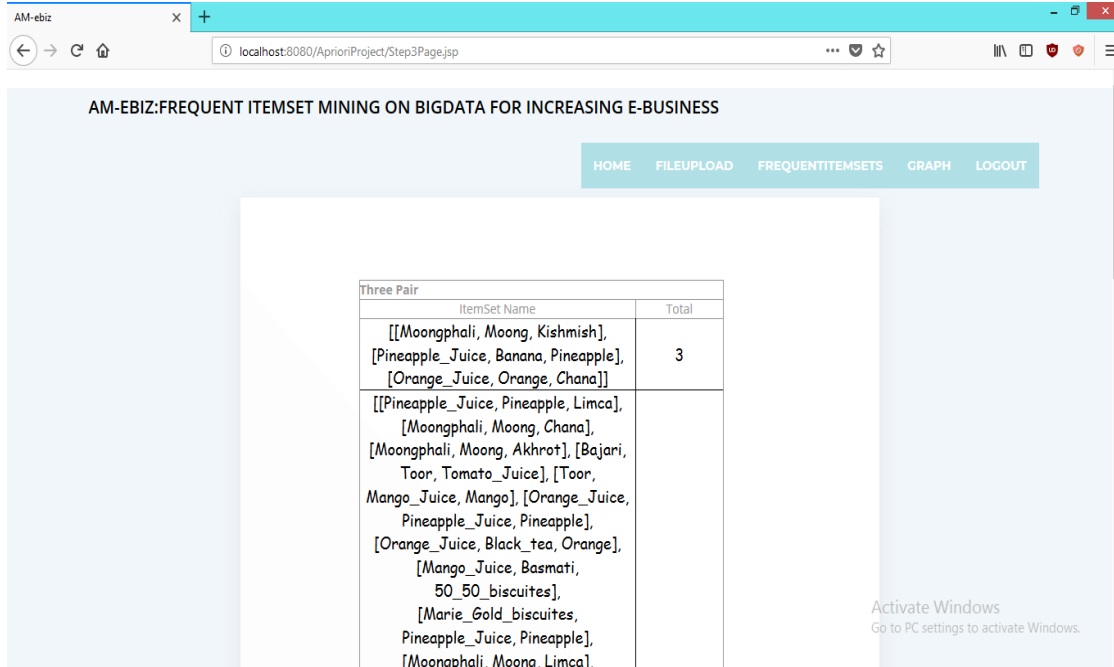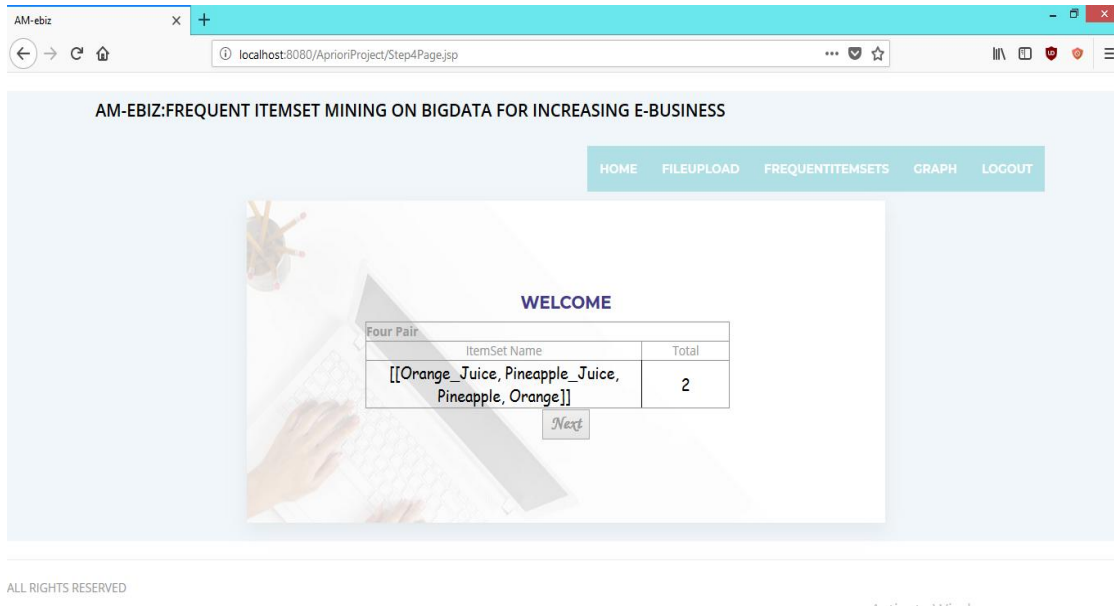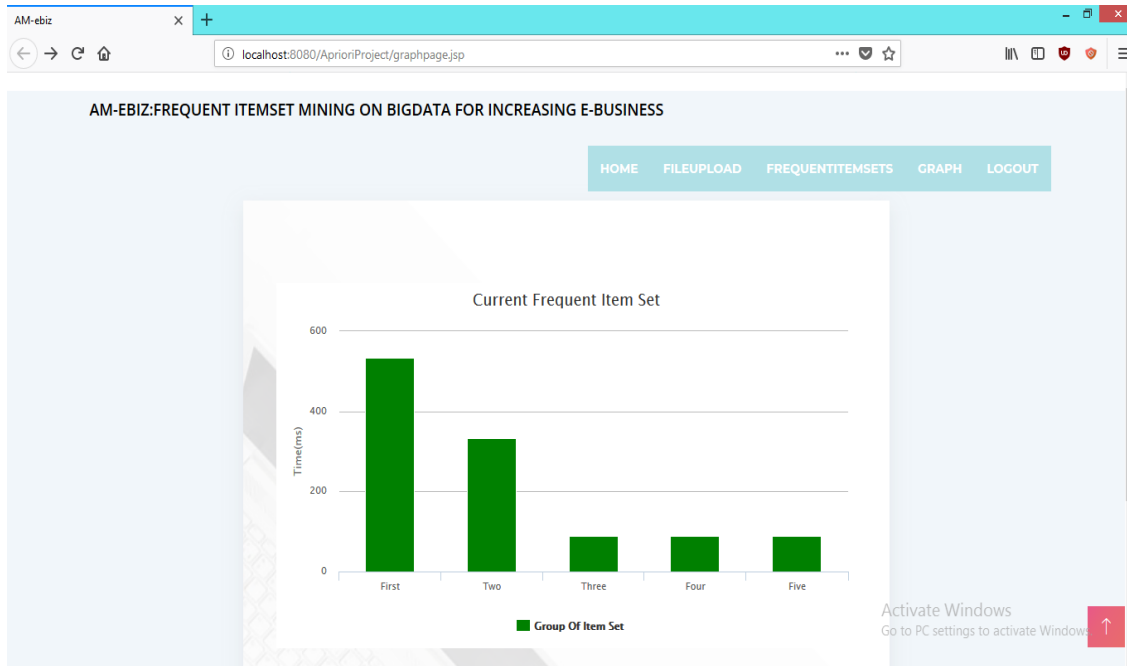
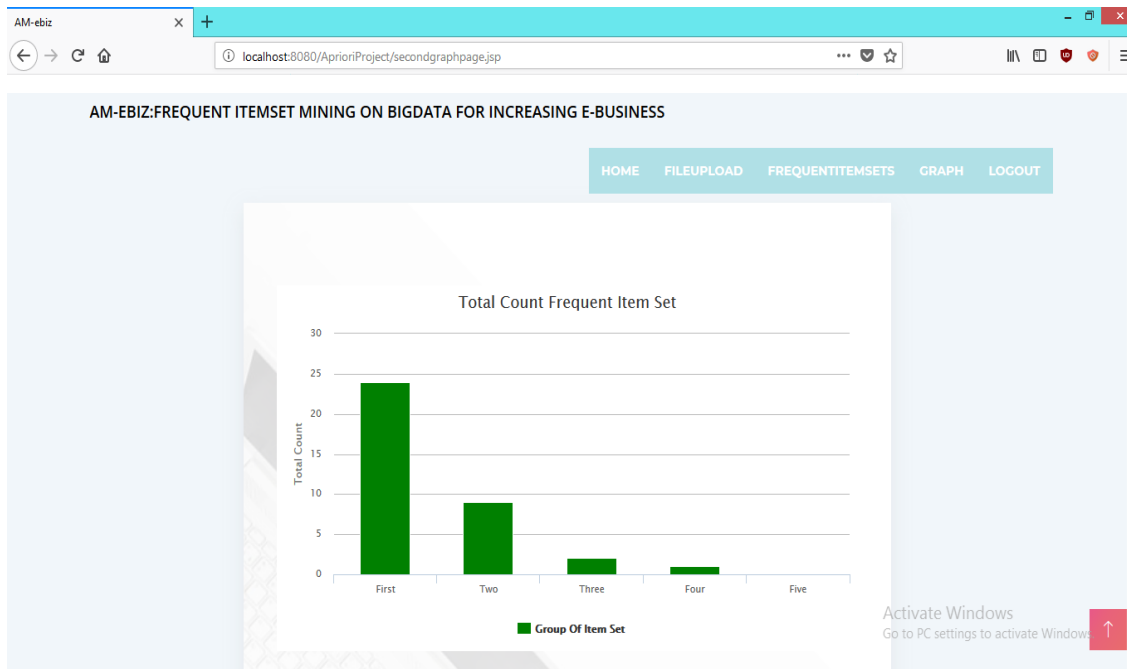**Figure 1.9: Three extracted pairs of item set**



**Figure 1.10: Final extracted pairs of item set**

**Figure 1.11: Time required for pair wise item set extraction**



**Figure 1.12: Item set count of pair wise extraction**

These figures show the time requirement for retrieving the frequent item set by AM_ebiz Algorithms.

**Performance of AM_ebiz in terms of Time required in seconds with different support denominator with different dataset**

Following table indicates the time taken in seconds by Proposed Algorithm AM_ebiz for the three different datasets of Grocery, Electronic & Sports for different support counts.

**Table 1.1: Performance of Time required in seconds with different support denominator with different dataset**

| Support Count | Time in Seconds | Time in Seconds | Time in Seconds |
|---|---|---|---|
| Support 5 | 31 | 29 | 30 |
| Support 8 | 26 | 24 | 25 |
| Support 10 | 22 | 21 | 20 |
| support 15 | 15 | 16 | 14 |

## 4. Conclusion

The evolution of information technology generates the back volume of data at an explosive rate in many disciplines that increases the challenge of storing and using those massive data in an efficient way. There is a constant improvement in the development of new techniques and algorithms to convert the data into knowledge. In this paper we have shown the implementation of a Novel association rule mining algorithm AM e_biz .The results shows the Performance of AM_ebiz in terms of Time required in seconds with different support denominator with different dataset for three different item set of Grocery, Electronic & sports for various support values of 5,8,10 and 15 respectively. It shows the time required to extract the frequent item set in seconds with different dataset. Utilization is very good in applying the algorithm AM_ebiz for frequent item set mining with retail item set for various support values.

# References

[1] Wei Lu[1], YanyanShen[2], Su Chen[3], Beng Chin Ooi,"Efficient Processing of k Nearest Neighbor Joins using MapReduce", Proceedings of the VLDB Endowment, Vol. 5, No. 10, issued August 27th 31st 2012.

[2] Yaling Xun[1], Jifu Zhang[2], Xiao Qin[3], Senior Member, "FiDoop-DP: Data Partitioning in Frequent Itemset Mining on Hadoop Clusters", IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS, VOL. 28, NO. 1, JANUARY 2017

[3] RakeshAgrawal[1] Tomasz Imielinski[2] _ Arun Swami[3] : Mining Association Rules between Sets of Items in Large Databases IBM Almaden Research Center650 Harry Road, San Jose, CA 95120 2012 .

[4] Bay Vo [1],Bac Le [2:] fast algorithm for mining generalized association: International Journal of Database Theory and Application through IEEE proceedings. Vol. 3, No. 5, October 2012

[5]M.Jayakameswaraiah[1]: Design and Development of data mining system to estimate cars promotion using improved id3 algorithm,International Journal of Advanced Research in Computer and Communication Engineering, 2010

[6]Pappa, Gisele L.Automating The Design Of Data Mining Algorithms, 2010, XIII.

[7] Ramakrishna Shrikant [1] Rakesh Agarwal [2]: Comparison Of Various Association Rule Mining Algorithm On Frequent Itemsets in International Journal of Advanced Research in Computer Science and Software Engineering.