# A NOVEL APPROACH FOR FREQUENT DATA PARTITIONING USING PARALLEL MINING ITEMSETS AND HADOOP

C.Yosepu[1], K Ganapathi Babu[2], D Harith Reddy[3]

[1]Associate Professor, [2,3]Assistant Professor, CSE Dept, St.Martins Engineering College

## Abstract

**Parallel traditional algorithms are used for mining frequent itemsets. Traditional parallel mining algorithms partition the data equal among a group of computing nodes. Existing parallel Frequent Itemset Mining algorithms have serious performance problems. Data partitioning strategy is used to resolve the problem Given a huge dataset, data partitioning strategies in the existing solutions endure high communication and mining overhead provoked by redundant transactions transmitted among computing nodes. We address this problem by Hadoop The core of Apache Hadoop contains a storage part, known as Hadoop Distributed File System (HDFS), and a processing part called as Map Reduce. Hadoop divides files into large chunks. It dispense them across computing nodes in a cluster. By using this approach the performance of existing parallel frequent-pattern increases. This paper shows the different parallel mining algorithms for frequent itemsets mining. We summarize the different algorithms that were developed for the frequent itemsets mining, like candidate key generation algorithm, such as Apriori algorithm and without candidate key generation algorithm, such as FP-growth algorithm. These algorithms does't have mechanisms like load balancing, data distribution I/O overhead, and fault tolerance.The efficient method is the FiDoop using ultrametric tree (FIUT) and Mapreduce programming model. FIUT scans the database two times. FIUT has four advantages. First: It reduces the I/O overhead as it scans the database two times. Second: only frequent item sets in each transaction are inserted as computing nodes for compressed storage. Third: FIU is enhanced method to partition database, which extensively reduces the search space. Fourth: frequent itemsets are created by examining only the leaves of tree rather than traversing complete tree, which decrease the computing time.**

**Index Terms: Data Mining, Recommender Systems, Social Network**

## I. INTRODUCTION

Parallel Frequent Itemset mining is looking for series of actions and load balancing of dataset. Constructing Hadoop cluster is particularly for storage and analyzing data. Through frequent Itemset mining we can extract knowledge from data. Instance of this technique is Market Basket Algorithm. It also have an effect on load balancing. It helps to enhance the speed of performance. This parallel Frequent Itemset mining is done by using map reduce programming model. Partitioning of data in large dataset through algorithm making data more efficient. This data partitioning is conceded out on Hadoop clusters. Data partitioning essential for scalability and high efficiency in cluster.

Frequent Itemsets Mining data partition have an effect on computing nodes and the traffic in network. Data partition may be spread over various nodes, and users at the computing node can execute local transactions on the partition.This enhances performance for sites that have regular transactions relating certain views of data, even as maintaining availability and security.Pperformance of parallel Frequent Itemset Mining on Hadoop clusters increases by using Fidoop-DP concept.Fidoop-DP is voronoi diagram. It is conceptualized on data partitioning

technique. Data mining is a process of extracting the required pattern from the large amount of data. There are many data mining techniques like clustering, classification and association rule.

The most popular one is the association rule that is divided into two parts i) generating the frequent itemset ii) generating association rule from all itemsets. Frequent itemset mining (FIM) is the core problem in the association rue mining. Sequential FIM algorithm suffers from performance deterioration when it operated on huge amount of data on a single machine.to address this problem parallel FIM algorithms were proposed. There are two types of algorithms that can be used for mining the frequent itemsets first method is the candidate itemset generation approach and without candidate itemset generation algorithm. The example for candidate itemset generation approach is the Apriori algorithm and for, without candidate itemsets generation is the FPgrowth algorithm. The important data-mining problem is discovering the association rule between the frequent itemset.in order to find best method for mining in parallel, we explore a spectrum for trade-off between computation, synchronization, communication, memory usage. Count distribution, data distribution, candidate distribution are three algorithms for discovering the associate rule between frequent itemsets. Minimizing communication is the focus of the count distribution algorithm.it will thus even at the expense of winding up redundant duplication computation in parallel. The data distribution effectively utilizes the main memory of the system.it is communication-happy algorithm. Here nodes to all other nodes broadcast the local data. The candidate distribution algorithm for both, to segment the database upon the different transaction support and the patterns, exploits linguistics of a particular problem. Load balancing is also incorporated by this algorithm.[1]

The most popular technique is the association rule that isdivided into two parts i) Finding the frequent itemset ii) generating association rule from all itemsets. Frequent itemset mining (FIM) is the main problem in the association rule mining. Sequential FIM algorithm have

performance deterioration problem when it is operated on large amount of data on a single machine.Parallel FIM algorithms were proposed to solve this problem. There are two kinds of algorithms that can be used for mining the frequent itemsets first approach is the candidate itemset generation approach and without candidate itemset generation algorithm. Apriori algorithm is an example for candidate itemset generation approach and FPgrowth algorithm.is an example for, without candidate itemsets generation.The important data-mining problem is finding the association rule between the frequent itemset. In order to find the best method for mining in parallel, we explore a spectrum for trade-off among computation, synchronization, communication, memory usage. Data distribution, Count distribution,, candidate distribution are three different algorithms for finding the associate rule between frequent itemsets. Minimizing communication is the main focus of the count distribution algorithm.It will thus even at the expense of winding up redundant replication computation in parallel. The data distribution effectively uses the main memory of the system.It is communication-happy algorithm. Herecomputing nodes to all other nodes broadcast the local data. The candidate distribution algorithm for both, to fragment the database upon the various transaction support and the patterns, exploits linguistics of a particular problem. Load balancing is also integrated by this algorithm.[1]

## II. LITERATURE SURVEY

YalingXun, Jifu Zhang, Xiao Qin, "FiDoop-DP: Data Partitioning in Frequent Itemset Mining on Hadoop Clusters", 2016.It defines, A data partitioning approach called FiDoop-DP using the Map Reduce programming model. The primary goal of FiDoop-DP is to increase the performance. A similarity metric to facilitate data-aware partitioning. As a forthcoming research direction, we will apply this metric to examine advanced load balancing strategies on a heterogeneous Hadoop cluster. I.Pramudiono&amp;M.Kitsuregwa," Fp-tax: Tree structure based generalized association rule mining",2004.This paper describes, exploration

of data partioning issues in parallel FIM.Main emphasis is on map-reduce. Forthcoming work is development of Fidoop which deeds correlation among traction to partition large datasets in Hadoop. X.Lin," Mr-apriori:Association rules algorithm based on mapreduce",2014.It explain, Main Emphasis on classical Algorithm connecting and pruning step using prefix Itemset based storage using has table. It points some of the limitations of Apriori algorithm. S. Hong, Z.Huaxuan, C. Shiping, and H.Chunyan," The study of improved fp-growth algorithm in mapreduce",2013.This describes, Build cloud platform to implement the parallel FPgrowth algorithm based on linked list and PLFPG.PLFPG algorithm compared higher efficiency and scalability. M. Liroz-Gistau, R. Akbarinia, D. Agrawal, E. Pacitti, and P. Valduriez," Data partitioning for minimizing transferred data in mapreduce",2013.It state that, Map Reduce jobs are executed over distributed system composed of a master and set of workers. Input is dividing into several splits and assigned to map tasks. Future work is evasion to perform the repartitioning in parallel

Sandy moen's at al,[2] proposed two methods for mining frequent itemset in parallel on the Mapreduce framework. Dist-Eclat is the first method which distributes the search space evenly as possible among mapper. This technique is use to mine large dataset but not massive datasets. This algorithm operates in three steps: We use vertical database relatively than transaction database. In the first step the vertical database is partitioned into equal sized chunks called shards and distributed to available mappers. Each mapper mine the frequent singletons from each block and give to the reducer. The reducer accumulates all the frequent tested. In the second step the set of frequent itemsets of size K are generated (Pk). Frequent singleton itemsets are disseminated to the mappers. Each mapper runs Éclat [3] to determine frequent K-sized superset of items. The reducer will collects all the frequent K-sized supersets of items and dispense it to the next batch of mappers. Round Robin is a method which is used for the distribution of the frequent itemset. The third step is the mining the prefix tree.

The mutual information between the mappers are independent, so mapper complete each step independently Every mapper takes the database and gives itemsets for which, we want to know the support .the reducer takes all itemsets and returns only the global frequent itemsets. These itemsets are considered as candidates and distributed to the mappers for breath-first search.
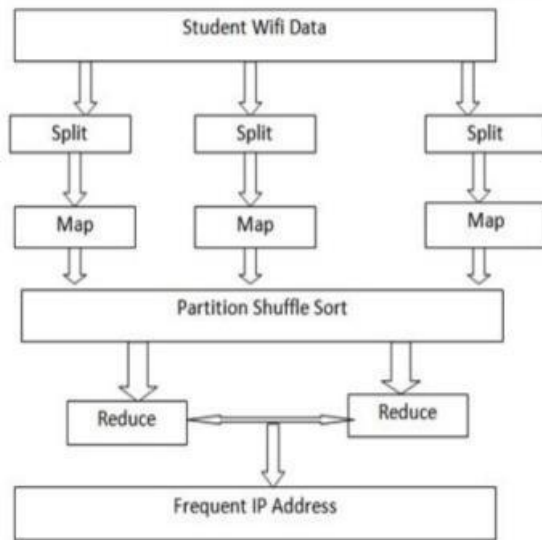
This process continues K-times to generate K-FI's.next step is computing the possible extension. The mapper gives local Tid-list to the reducer the reducer combines the local Tid-lists, to one Tid-list, and assigns prefix to mappers. The mapper in the final step works on individual prefix group. A prefix group fits in the memory as a conditional database. The diffsets are used to mine the frequent itemsets in the conditional database.

Enhilvathani et al [4] have used the Apriori algorithm for frequent item set generation on mapreduce programing model. For implementation of algorithm is given in five steps In the first step the transaction dataset is partitioned that is Divided into n subsets done that are of map phase.in the second step the data subsets are formatted as <key1, value1>pair, key is Tid(Transaction id). The mapreduce task is executed in third phase. The record of input item subsets are scanned by the Map function and candidate item sets input are generated by the map function.

## III.  PROPOSED WORK

This section is about objective of project, which tells detail of system. It reduces the complexity of data access and retrieval. When we have to dealing with big data i.e. huge amount of data traditional exisisting system seems inefficient. The alternative to this is apache Hadoop, which deals with big data with efficiency. Hadoop itself consists of Map Reduce and HDFS.We use Hadoop with concept called frequent Itemsets which makes it FiDoop(Frequent item Hadoop).It runs on Hadoop cluster Job of map reduce is partitioning of data, it splits the local input data to generate local 1-itemset and it further reduce to specific reduced data. It sorts the data in decreasing order of frequency. In

second step job of map reduce is FP-Growth based partitions and last job is last task to aggregate the result from previous stages to generate output. LSH based partitioning boost the performance of system by avoiding large number of comparisons. It uses bucket to keep similar transaction together.



System Architecture

## IV. METHODOLOGY

To deals with migration of high communication, and decrease computing cost in map reduce. We applied frequent item data partitioning technique which establishes correlation among transaction for data partitioning. It is based on: Similarity in data and transaction

•Group this highly correlated data.

•Group this highly correlated datacontent comes here Conclusion content comes here . Conclusion content comes here The existing references tell that frequent item mining improves the output up to 31% with18% average. We are working to develop system that investigates the detail of students. [Uses Wi-Fi].It allows generating result based on various parameters.

## V. CONCLUSION

Mapreduce programming model is useful for existing parallel mining algorithm for mining frequent itemsets from database and solves the load balancing and scalability. This paper gives the summary of algorithms planned for parallel mining of frequent itemsets .The Apriori and FP tree algorithms used for mining frequent itemsets. Main disadvantage of Apriori algorithm is that the database has to be scanned number of times and massive candidate keys wants to be exchanged between the processor. I/O and synchronization are the additional problems in the Apriori algorithm. The drawback of FP-growth, however, lies within the impracticableness to create in-memory FP trees to accommodate large-scale databases. This disadvantage becomes a lot of pronounced once it comes to massive and two-dimensional databases. To overcome these troubles, FiDoop, an parallel frequent itemset mining algorithm is designed. FiDoop include the ultra metric tree (FIU) rather than Apriori or FP-growth algorithm. The FIU tree achieves compressed storage. FiDoop runs three MapReduce jobs. The third MapReduce job is important. In third job the mapper separately decomposes itemsets and reducer built the ultra metric trees.

## REFERENCES

[1] "Parallel Mining of Association rule." Rakesh Agarwal ,John C Safer

[2] "Frequent Itemset Mining for Big Data Sandy Moens, Emin Aksehirli and Bart Goethals Universiteit Antwerpen, Belgium

[3] "ECLAT Algorithm for Frequent Itemsets Generation "Manjit kaur , Urvashi Grag Computer Science and Technology, Lovely Professional University Phagwara, Punjab, India International Journal of Computer Systems (ISSN: 2394-1065), Volume 01– Issue 03, December, 2014 Available at http://www.ijcsonline.com/

[4] "Implementation Of Parallel Apriori Algorithm On Hadoop Cluster" A. Ezhilvathani1, Dr. K. Raja. International Journal of Computer Science and Mobile Computing

[5] "Frequent Itemsets Parallel Mining Algorithms " Suraj Ghadge, Pravin Durge, Vishal Bhosale,Sumit Mishra. Department of Computer Engineering, JSPM's ICOER. International Engineering

Research Journal (IERJ) Volume 1 Issue 8 Page 599-604, 2015, ISSN 2395-1621

[6] "FiDoop: Parallel Mining of Frequent Itemsets Using MapReduce" Yaling Xun, Jifu Zhang, and Xiao Qin, Senior Member, IEEE

[7] Yaling Xun, Jifu Zhang, Xiao Qin,FiDoop-Dp Data Partitioning in Frequent Itemset Mining on Hadoop clusters,2016.

[8] S. Sakr, A. Liu, and A. G. Fayoumi, âœThe family of mapreduce and large-scale data processing systems, ACM Computing Surveys (CSUR), vol. 46, no. 1, p. 11, 2013.

[9] X. Lin, Mr-apriori: Association rules algorithm based on mapreduce,â in Software Engineering and Service Science (ICSESS), 2014 5th IEEE International Conference on. IEEE, 2014, pp. 141"144.

[10]S. Hong, Z. Huaxuan, C. Shiping, and H. Chunyan, âœThe study of improved fp-growth algorithm in mapreduce, in 1st International Workshop on Cloud Computing and Information Security. Atlantis Press, 2013. 11P. Uthayopas and N. Benjamas, Impact of i/o and execution scheduling strategies on large scale parallel data mining, Journal of Next Generation Information Technology (JNIT), vol. 5, no. 1, p. 78, 2014.