# Frequent Subgraph Mining Algorithm with Load Balancing for Enhanced Iterative Map reducing Techniques

## M.Gokilavani*, B.Anitha, R.Rohini

**Department of Computer Science, Vivekanandha College of Engineering for Women, Namakkal, Tamil Nadu, India**

**Abstract :** Mining frequent sub graphs has attracted a great deal of attention in many areas for data analysis. Above the years, a lot of algorithms have been projected to respond this problem. In a lot of algorithms suppose that the information construction of the extracting job is little enough to well in the chief recollection of a computer. Statistics grows under together size and quantity also; such a statement does not grip any longer. For the purpose of they should be use a Parallel iterative *Map-Reduce* based frequent sub graph mining algorithm (*PFSM-H*),it is also have the some disadvantages. In this work, we propose a frequent sub graph mining algorithm called *LB-FSM*. Existing problem will be solved by LB-FSM method, which uses an iterative Map Reduce based framework. The Load Balancing Frequent Sub graph Mining algorithm will be very useful for mining process, as it returns all the frequent sub graphs for a known consumer-distinct hold up and it is well-organized as it apply every one optimizations that the newest *FSM* algorithms. The experiments with genuine existence and great artificial datasets legalize the success of LB-FSM for mining *frequent sub-graphs* as of huge graph datasets.

**Keywords :** FSM; Mapreduce; PFSM-H; Frequent sub-graph; LB-FSM.

## Introduction

Bigdata includes data mining, computational biology, environmental sciences, e-commerce, web mining and social network analysis. In a few domains, analyze and extract the huge amount of data for extracting novel has a turn into a schedule job. There are lots of efficient algorithms for find a recurrent item sets in extremely great contract database. They can make use of some relational kind with more item set for discover the association rules, for extracting prevalent patterns that survive in datasets, or for the classification also. Though that it can't be applying to these techniques more than datasets which are not an item sets. In the fresh time, there has been an augmented notice in rising data mining algorithms that function on graph. Such type of graphs occurs logically in numeral different application domains, as well as computer networks, chemical compound, semantic web, bioinformatics and social networks. Everyone domain mentioned wants extracting of recurrent sub graphs more than bulky information sets. However, the algorithms urbanized as a result distant are not scalable.

The majority jobs complete on FSM contain listening carefully on algorithms that take for granted graph information is stored in chief recollection. Reminiscence contingent algorithms might be extremely nescient if the information set is huge. Characteristic graph extracting algorithms noiselessly take for granted that the graph fit in the recollection of a distinctive workplace, otherwise at smallest amount on a solitary floppy; the on top of graphs infringe these assumption, on both sides of multiple Giga-bytes, and caption to Tera and Peta-bytes of information. Process of a map reduce algorithm base diagram will be the bigger one. The

difficulty is computationally concentrated, as results we shall think on principles that can be executed by a solitary around of map reduce. We examine how to reduce two significant events of difficulty. HADOOP is the unlock basis completion of map reduce. HADOOP provide the Distributed File System (HDFS) and PIG, an elevated height words for information examination. Solve the duty of frequent sub graph taking out on a dispersed display place similar to Map Reduce is demanding for a variety of reasons.

## Mapreduce

Map Reduce is a dispensation technique and an agenda replica for dispersed compute by the java basis. The mappers in the training chapter also create the extracting duty by emit the recurrent solitary border pattern as < key, value> couple. The mapper job is the primary duty which takes the contribution data and converts into a set of information where entity rudiments are out of order downward into tuples. The lessen chore is the mixture of drag your feet phase and the reducer task is to procedure the statistics that come beginning mapper. Subsequent to dispensation it produce a novel set of production <key,value> couple which will be stored in HDFS.

## Frequent Subgraph Mining

Frequent Sub graph Mining (FSM) is the spirit of grid withdrawal. The aim of FSM is to remove each and every one the recurrent sub graphs, in a known information set, whose incidence count are on top of a particular doorsill. Additional than the investigate movement linked by means of FSM the significance of FSM is as well reflect in its a lot of areas of its request. The simple thought at the back FSM is to "cultivate" applicant sub graphs, in moreover a width primary or deepness initial way (applicant age group), and then decide if the recognized contender sub graphs happen regularly sufficient in the grid information set for them to be careful motivating (hold up including).

The two major investigate issues in FSM are thus how to professionally and efficiently :(i) produce the candidate recurrent sub graphs and (ii) decide the incidence of amount of the generate sub graphs. Successful applicant sub graph age group requires that the age band of photocopy or surplus candidate is avoided. Happening as well as requires repetitive judgment of contender sub graphs with sub graphs in the participation information, a procedure known as isomorphism examination. To additional make easy thoughtful of the pasture of FSM we differentiate flanked by recurrent sub tree taking out and the a great deal all-purpose field of FSM.

## Related Work

There stay alive a lot of algorithms for solve the in recollection account of FSM chore, the majority distinguished in the middle of them are AGM [28], FSG [1], gSpan[2], Gaston [3], and DMTL [4]. These methods take for granted so as to the information set is minute and the removal mission finish in a rational quantity of occasion by means of an in recollection technique. To believe the bulky data situation, a small number of customary folder based grid mining algorithms, such as, DB-Subdue [5], and DB-FSG [6] and OOFSG [7] are also projected. For important grid extracting farm duties, researchers careful shared recollection similar algorithms for FSM.obtainable a similar description of their recurrent sub graph taking out algorithm Subdue [8].Urbanized an equivalent toolkit[9] for their Motif- Miner [10] algorithm. Meinl shaped a software name Parmol[11] which include similar completion of Mofa[12], gSpan[2], FFSG [13] and Gaston [3]. ParSeMis[14] is another such tool that provides equivalent execution of span algorithm. To contract with the scalability trouble caused through the dimension of contribution graphs, there are couples of famous works, Part Miner [15] and Part Graph Mining[16], which are base on the thought of partition the grid information. Here in adding together exists a work [17] on adaptive corresponding grid extracting for CMP Architectures. Map Reduce framework has been used to mine frequent patterns where the transactions in the input database are simpler combinatorial objects such as, a set [18, 19, 20, 21], or a sequence [22]. The authors suppose FSM on Map Reduce, but their move toward is incompetent owing to a range of shortcoming[23]. And also the majority illustrious is so as to in their technique they do not take on some instrument to keep away from generate copy pattern. This reason an exponential augment in the bulk of the contender sub graph breathing space; in addition, the production place contain photocopy of the identical grid pattern that are solid to amalgamate as the customer has to supply a sub graph isomorphism schedule for this adjustment. A new trouble by means of the on top of technique is so as to it require the client to identify the numeral of Map Reduce repeated works. Author did not talk about how to decide the entirety repeated calculate

consequently that the principles are clever to calculate every one of recurrent patterns for an already existed bear. One practicable method power is to position the repeated counts to be alive the border count up of the chief business other than that will be an excess of the possessions. FSM-H does not undergo starting with some of the on top of confines. Throughout the amendment stage of this magazine, we become conscious of an additional labor[24] of FSM on Map Reduce. It is a non-repeated technique which runs Gaston [18] on every divider of the grid folder. There survive a quantity of mechanism [26,] that excavation sub graphs so as to be recurrent bearing in mind their induce occurrence in a solitary great grid. On the other hand, it should be a dissimilar as the aim of frequent sub graph mining algorithm using hadoop is to extract sub graphs that are recurrent in excess of a compilation of graph in a grid database.

## Proposed System

### Creating a new node

Every graph, node, and edge can grasp key/value characteristic pairs in an linked attribute lexicon (the keys must be hash able). By non-payment these are unfilled(Fig.1), but attributes can be additional or misrepresented by means of add_edge, add_node or straight treatment of the quality dictionary named G.graph, G.node and G.edge for a graph G.



**Fig.1. Creating a new node**



**Fig.2.Appending edges**

### Appending edges

At the present that nodes are shaped, we can produce the edges, add a new edge in the grid and move toward reverse it. The edge is also additional in every one the super-graph of the grid to uphold the sub-graph relative flanked by grids(Fig.2). The starting stricture is the "basis node", and, of route, the second is the "aim node" (in tulip, all edge are sloping but you can decide not to believe the compass reading), the edges details arrange is the single in which they are additionally added one.

**Determining frequent sub graphs**

Recurrent sub graphs are strong-minded under the numeral era it appears in the entire information set. Our reason is to discover the sub graphs which happen additional than a precise number of period (min−sup will be provide by the consumer) in the information set(Fig.3). Let's believe the least amount support as 2, which mean a sub graphs must be appear in at slightest two grids. Subgraph−1 class contain the size-1 sub graphs, Observe that the Subgraph−1 class does not have the statistics of vertices. Only important particulars, the labels are stored.
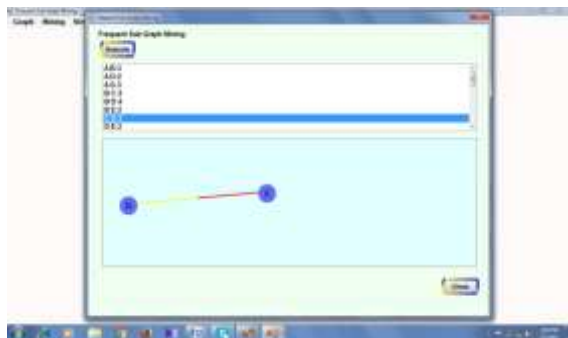


**Fig.3. Determining frequent sub graphs**

**Result and Discussion**

In this paper we present a work of fiction repeated Map Reduce based FSM algorithm, called LB-FSM. The presentation of LB-FSM more than real life and huge synthetic datasets for a variety of scheme and input configurations will be obtainable in an enhanced method in this paper, when compare to extra existing method. That period compare the execution time of LB-FSM also. It show the LB-FSM is significantly improved than the obtainable technique.

**References**

1. M. Kuramochi and G. Karypis, Frequent subgraph discovery, in Proc. Int. Conf. Data Mining, 2001, pp. 313–320.
2. X. Yan and J. Han, gSpan: Graph-based substructure pattern mining, in Proc. Int. Conf. Data Min., 2002, pp. 721–724.
3. S. Nijssen, and J. Kok, A quickstart in frequent structure mining can make a difference, in Proc. 10th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining, 2004, pp. 647–652.
4. V. Chaoji, M.Hasan, S. Salem, and M. Zaki, An integrated, generic approach to pattern mining: Data mining template library, Data Min. Knowl. Discov. J., vol. 17, no. 3, pp. 457–495,2008.
5. S. Chakravarthy, R. Beera, and R. Balachandran, Db-subdue: Database approach to graph mining, in Proc. Adv. Knowl. Discov. Data Mining, 2004, pp. 341–350.
6. S. Chakravarthy and S. Pradhan, Db-FSG: An SQL-based approach for frequent subgraph mining, in Proc. 19th Int. Conf. Database Expert Syst. Appl., 2008, pp. 684–692.
7. B. Srichandan and R. Sunderraman, Oo-FSG: An object-oriented approach to mine frequent subgraphs, in Proc. Australasian Data Mining Conf., 2011, pp. 221–228.
8. D. J. Cook, L. B. Holder, G. Galal, and R. Maglothin, Approaches to parallel graph-based knowledge discovery, J. Parallel Distrib. Comput., vol. 61, pp. 427–446, 2001.
9. C. Wang and S. Parthasarathy, Parallel algorithms for mining frequent structural motifs in scientific data, in Proc. 18th Annu. Int. Conf. Supercomput., 2004, pp. 31–40.
10. S.Parthasarathy and M. Coatney, Efficient discovery of common substructures in macromolecules, in Proc. IEEE Int. Conf. Data Mining, 2002, pp. 362–369.
11. T. Meinl, M. Worlein, O. Urzova, I. Fischer, and M. Philippsen, The parmol package for frequent subgraph mining. Electron. Commun. EASST, vol. 1, pp. 1–12, 2006.
12. C. Borgelt and M. Berthold, Mining molecular fragments: finding relevant substructures of molecules, in Proc. IEEE Int. Conf. Data Mining, 2002, pp. 51–58.

13.  J. Huan, W. Wang, and J. Prins, Efficient mining of frequent subgraphs in the presence of isomorphism, in Proc. 3rd IEEE Int. Conf. Data Mining, 2003, pp. 549–552.
14.  M. Philippsen,M.Worlein,A. Dreweke, and T. Werth.Parsemis-The parallel and sequential mining suite. [Online]. Available:https://www2.cs.fau.de/ EN/research/ParSeMiS/ index.html,(2011).
15.  J. Wang, W. Hsu, M. L. Lee, and C. Sheng, A partition-based approach to graph mining, in Proc. 22nd Int. Conf. Data Eng., 2006, p. 74.
16.  S. N. Nguyen, M. E. Orlowska, and X. Li, Graph mining based on a data partitioning approach, in Proc. 19th Australasian Database Conf., 2008, pp. 31–37.
17.  G. Buehrer, S. Parthasarathy, and Y.-K. Chen, Adaptive parallel graph mining for CMP architectures, in Proc. 6th IEEE Int. Conf. Data Mining, 2006, pp. 97–106.
18.  G.-P. Chen, Y.-B. Yang, and Y. Zhang, Mapreduce-based balanced mining for closed frequent itemset, in Proc. IEEE 19th Int. Conf. Web Serv., 2012, pp. 652–653.
19.  S.-Q. Wang, Y.-B. Yang, Y. Gao, G.-P. Chen, and Y. Zhang, Mapreduce-based closed frequent itemset mining with efficient redundancy filtering, in Proc. IEEE 12th Int. Conf. Data Mining Workshops, 2012, pp. 49–453.
20.  L. Zhou, Z. Zhong, J. Chang, J. Li, J. Huang, and S. Feng, Balanced parallel FP-growth with Mapreduce, in Proc. IEEE Youth Conf. Inf. Comput. Telecommun, 2010, pp. 243–246.
21.  H. Li, Y. Wang, D. Zhang, M. Zhang, and E. Y. Chang, Pfp: parallel FP-growth for query recommendation, in Proc. ACM Conf. Recommender Syst., 2008, pp. 107–114.
22.  B.-S. Jeong, H.-J. Choi, M. A. Hossain, M. M. Rashid, and M. R. Karim, A MapReduce framework for mining maximal contiguous frequent patterns in large DNA sequence datasets, IETE Tech. Rev., vol. 29, pp. 162–168, 2012.
23.  S. Hill, B. Srichandan, and R. Sunderraman, An iterative Mapreduce approach to frequent subgraph mining in biological datasets, in Proc. ACM Conf. Bioinformat., Comput. Biol. Biomed., 2012, pp. 661–666.
24.  X. Xiao, W. Lin, and G. Ghinita, Large-scale frequent subgraph mining in Mapreduce, in Proc. Int. Conf. Data Eng., 2014, pp. 844–855.
25.  M. Kuramochi and G. Karypis, Finding frequent patterns in a large sparse graph*, Data Mining Knowl. Discov., vol. 11, pp. 243– 271, 2005.
26.  S. Skiadopoulos, M. Elseidy, E. Abdelhamid, and P. Kalnis, Grami: Frequent subgraph and pattern mining in a single large graph, Proc. Very Large Database Endow., vol. 7, pp. 517–528, 2014.
27.  M. Gokilavani, B. Anitha, R. Jayanthi ,A Survey On Mapreduce Using Frequent Subgraph Mining, ISSN: 0976-3104,IIOABJ/vol.7/6/1-6.
28.  A. Inokuchi, T. Washio, and H. Motoda, An apriori-based algorithm for mining frequent substructures from graph data, in Proc. 4th Eur. Conf. Principles Data Mining Knowl. Discov., 2000, pp. 13–23.

**\*\*\*\*\***