

A Survey on PPCT-FIM: Prepost Computation Tree Based Frequent Itemset Mining

Mr. Phalke Sagar Balkrishna¹, Prof. Amol Subhash Rajpure²

PG Student, Department of Computer Engineering, Dattakala Institute,

Faculty of Engineering, Bhigwan, Swami Chincholi, Daund, Pune, India¹

Assistant Professor, Department of Computer Engineering, Dattakala Institute,

Faculty of Engineering, Bhigwan, Swami Chincholi, Daund, Pune, India²

Abstract: Mining incessant itemsets is a fundamental issue in information mining and plays a significant job in numerous information mining applications. As of late, some itemset portrayals in light of hub sets have been proposed, which have demonstrated to be proficient for mining visit itemsets. In this paper, we propose a PrePost Computation Tree based Frequent Itemset Mining (PPCT-FIM), calculation for mining continuous itemsets. To accomplish high productivity, PPCT-FIM finds visit itemsets utilizing a set-list tree with a half breed search procedure and legitimately identifies visit itemsets without competitor age under some case. For assessing the exhibition of PPCT-FIM, we have direct broad examinations to contrast it against and existing driving calculations on an assortment of genuine and engineered datasets. The exploratory outcomes show that PPCT-FIM is altogether quicker than PFIM calculations.

Keywords: Data mining, Frequent itemset, Mining Massive Data, Pruning Rule, Incremental Update

I. INTRODUCTION

Frequent itemset mining is a significant activity that has been generally examined in numerous common sense applications, for example, information mining [1]-[3], programming bug discovery [4], spatiotemporal information examination and organic investigation [5]. Given an exchange table, where every exchange contains a lot of things, visit itemset mining restores all arrangements of things whose frequencies (likewise alluded to as help of the arrangement of things) in the table are over guaranteed limit. Because of its down to earth significance, since right off the bat proposed in [6], visit itemset mining has gotten broad considerations and numerous calculations are proposed [7]-[9]. The existing incessant itemset mining calculations can be characterized into two gatherings: candidate generation- based calculations [10]-[14] and design development based calculations [15]-[17]. The applicant age based calculations initially create competitor itemsets and these up-and-comers are approved against the exchange table to distinguish visit itemsets. The hostile to monotone property is used in competitor age based calculations to prune search space. In any case, the competitor age based calculations require various pass table sweeps and this will cause a high I/O cost on huge information. The example development based calculations try not to create up-and-comers expressly. They build the exceptional tree-based information structures to keep the fundamental data about the continuous itemsets of the exchange table. By utilization of the built information structures, the continuous itemsets can be registered effectively. In any case, design development based calculations have the issue that the built information structures are unpredictable and as a rule surpass the accessible memory on enormous information.

To summarize, the current calculations can't process frequent itemsets on gigantic information effectively. In frequent itemset mining, the quantity of the incessant itemsets ordinarily is delicate to the estimation of the help edge. On the off chance that the help edge is little, there will be countless regular itemsets and it is hard for the clients to make productive choices. Despite what might be expected, if the help limit is huge, it is conceivable that no successive itemsets can be found or the fascinating itemsets might be missed. Along these lines, an appropriate bolster limit is pivotal for the useful regular itemset mining and the clients frequently need to perform frequent itemset digging for a few times before the good help limit is resolved. The procedure regularly is intuitive. On huge information, the current calculations regularly need a long execution time to process frequent itemsets and this will influence clients' working productivity genuinely [18]. The focal point of this work is to locate another productive calculation to register frequent itemsets on monstrous information rapidly.

II. RELATED WORK

In this work, we want to utilize this reuse idea to a much larger degree. In typical massive data applications, with the increasing data volume and the disk I/O bottleneck, data usually is stored in read/append-only mode [19]. Therefore, the overall data set can be divided into two parts: the much larger old data set storing the historical data, and the relative

small new data set storing the newly generated data. Based on the description above, this work devises a new PFIM algorithm (Precomputation-based Frequent Itemset Mining algorithm) on massive data, which utilizes the pre-constructed frequent itemsets on the old data set to return the frequent itemsets quickly. Since the too small value of support threshold will generate too many frequent itemsets, we assume in this work that there exists a lower bound ω of the support threshold specified by the users in practical applications. Because of the real/append-only mode, given the old table T_0 , PFIM first pre-constructs the frequent itemsets (refer to as quasi-frequent itemsets in this work) whose supports are no less than ω . The new transactions are accumulated in the new table T_1 . Taking advantage of the pre-constructed quasi-frequent itemsets, given the specified support threshold, PFIM can compute the frequent itemsets on $T_0 \cup T_1$ quickly.

III. CANDIDATE GENERATION BASED ALGORITHMS

The candidate-generation-based algorithms firstly generate the candidates of the frequent itemsets, then the candidates are validated against the transaction table, and the frequent itemsets are discovered. Apriori algorithm [11], [20] adopts a level-wise execution mode. It uses the downward closure property, i.e. any superset of an infrequent itemset must also be infrequent, to prune the search space. By a pass of scan on the transaction table, it first counts the item occurrences to find the frequent 1-itemsets F_1 . Subsequently, the frequent k -itemsets in F_k are used to generate the candidates C_{k+1} of the frequent $(k+1)$ -itemsets. Another pass of scan is needed to compute the supports of candidates in C_{k+1} to find the frequent $(k+1)$ -itemsets F_{k+1} . This process iterates similarly until the F_{k+1} is empty. Apriori algorithm often needs multiple passes over table; it will incur a high I/O cost on massive data. Savasere et al. [12] propose Partition algorithm to generate frequent itemsets by reading the transaction table at most two times. The execution of Partition consists of two stages. In the first stage, Partition algorithm divides the table into a number of non-overlapping partitions in terms of the allocated memory, and the local frequent itemsets for each partition are computed. All the local frequent itemsets are merged at the end of first stage to generate the candidates of frequent itemsets. In the second phase, another pass over table is performed to acquire the support of the candidates and the global frequent itemsets can be discovered.

The useful property adopted in Partition is that, every global frequent itemsets must be appeared in local frequent itemsets of at least one partition. Partition algorithm utilizes vertical table representation of transaction table and the support counting is performed by recursive TID (Transaction Identifier) list intersection. In the first phase, Partition may generate many false positives, i.e. the itemsets are frequent locally but not frequent globally. Therefore, it needs another table scan to remove the false positives. Zaki [13] proposes another vertical mining algorithm Eclat. Eclat decomposes the original search space by a lattice theoretic approach into smaller sub lattices, each of which is a group of itemsets with a common prefix (referred to as prefix-based equivalence class). Depending on the allocated memory size, Eclat can recursively partition large classes into smaller ones until each class can be maintained entirely in the memory. Then, each class is processed independently in the breath-first fashion to compute the frequent itemsets. Eclat processes the sub lattices sequentially one by one and does not need post-processing overhead as Partition algorithms. The main problem of Eclat is that when the intermediate results of vertical TID lists can become too large for memory, especially in dense database, the performance of Eclat starts to suffer.

IV. PATTERN GROWTH BASED ALGORITHMS

Pattern-growth-based algorithms do not generate candidate itemsets explicitly but compress the required information for frequent itemsets in specific data structure. The frequent itemsets can be acquired quickly with the notion of projected databases, a subset of the original transaction database relevant to the enumeration node. Agarwal et al. [26] present DepthProject algorithm to mine long itemsets in databases. DepthProject analyzes the hubs of the lexicographic tree inside and out first request. The assessment procedure of a hub alludes to the help tallying of the competitor expansion of the node. During the inquiry, the anticipated exchange sets are kept up for a portion of the hubs on the way from the root to the hub P at present being expanded. During the search, the projected transaction sets are maintained for some of the nodes on the path from the root to the node P currently being extended. Normally, the projected transaction sets only contain the relevant part of the transaction database for counting the support at the node P . In the process of depth-first search, the projected database can be reduced further at the children of P and DepthProject can reuse the counting work of its previous exploration. At the lower levels of the lexicographic tree, a specialized counting technique called bucketing is used to substantially improve the counting time.

Han et al. [16] propose a FP-tree-based FP-development calculation to mine the total set of regular examples by design part development. FP-tree (visit design tree) is a minimized prefix-based trie structure to store the basic data about continuous examples. In every exchange, just continuous length-1 things, which are arranged with the plunging request of help, are utilized to build the FP-tree. At that point the FP-development calculation works on FP-tree as opposed to on the first database to mine regular examples. FP-development calculation begins with a regular length-1 example (introductory addition design), and the arrangement of continuous things co-happening with the addition design is separated as contingent example base, which is then developed as contingent FP-tree. With the present addition design

and the restrictive FP-tree, if the contingent FP-tree isn't unfilled, FP-development performs mining recursively. The incessant examples are procured by connecting the new ones produced from the contingent FP-tree and the addition design. FP-development changes the issue of finding long successive examples to searching for shorter ones and afterward linking the postfix. An extra streamlining is proposed for FP-development, for example on the off chance that every one of the hubs of the FP-tree lie on a solitary way, the successive examples can be produced by list of every one of the mixes of the sub-ways with the help being the base help of the itemsets contained in the sub-way.

V. CONCLUSION AND FUTURE WORK

In this work, we present a novel structure called PPCT-FIM to facilitate the process of mining frequent itemsets. Based on PFIM, an algorithm named PPCT-FIM is proposed to fast find all frequent itemsets in databases. Compared with Noderset, the key advantage of PPCT-FIM lies in that its size much smaller. This makes PPCT-FIM more suitable for mining frequent itemsets. The extensive experiments show that PPCT-FIM is favorable. PPCT-FIM proves to be state-of-the-art since it always runs fastest on all datasets with different minimum supports and occupied less memory when compared with previous leading algorithms.

ACKNOWLEDGEMENT

I profoundly grateful to **Prof. Rajpure A.S.** for his expert guidance and continuous encouragement throughout to see that this project rights its target since its commencement to its completion. I would like to express my deepest appreciation towards Principal **Dr. V. G. Arajpure**, HOD, **Prof. Rajpure A.S.** department of computer engineering and PG coordinator, **Prof. Bere S.S.** I must express my sincere heartfelt gratitude to all staff members of computer engineering department who helped me directly or indirectly during this course of work. Finally, I would like to thank my family and friends, for their precious support.

REFERENCES

- [1]. A. Ceglar and J. F. Roddick, "Association mining," *ACM Comput. Surv.*, vol. 38, no. 2, p. 5, 2006.
- [2]. H. Cheng, X. Yan, J. Han, and P. S. Yu, "Direct discriminative pattern mining for effective classification," in *Proc. 24th Int. Conf. Data Eng.*, Apr. 2008, pp. 169-178.
- [3]. H. Wang, W. Wang, J. Yang, and P. S. Yu, "Clustering by pattern similarity in large data sets," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, Jun. 2002, pp. 394-405.
- [4]. Z. Li and Y. Zhou, "PR-miner: Automatically extracting implicit programming rules and detecting violations in large software code," in *Proc. 10th Eur. Softw. Eng. Conf. Held Jointly 13th ACM SIGSOFT Int. Symp. Found. Softw. Eng.*, Sep. 2005, pp. 306-315.
- [5]. J. T. L. Wang, M. J. Zaki, H. Toivonen, and D. Shasha, Eds., *Data Mining in Bioinformatics*. London, U.K.: Springer, 2005.
- [6]. R. Agrawal, T. Imielinski, and A. Swami, "Database mining: A performance perspective," *IEEE Trans. Knowl. Data Eng.*, vol. 5, no. 6, pp. 914-925, Dec. 1993.
- [7]. C. C. Aggarwal, *Data Mining: The Textbook*. Cham, Switzerland: Springer, 2015.
- [8]. C. C. Aggarwal and J. Han, Eds., *Frequent Pattern Mining*. Cham, Switzerland: Springer, 2014.
- [9]. J. Han, H. Cheng, D. Xin, and X. Yan, "Frequent pattern mining: Current status and future directions," *Data Mining Knowl. Discovery*, vol. 15, no. 1, pp. 55-86, Aug. 2007.
- [10]. R. Agrawal, T. Imielinski, and A. N. Swami, "Mining association rules between sets of items in large databases," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 1993, pp. 207-216.
- [11]. R. Agrawal & R. Srikant, "Fast algorithms for mining association rules," in *Proc. 20th Int. Conf. Very Large Data Bases (VLDB)*, 1994, pp. 487-499.
- [12]. A. Savasere, E. Omiecinski, and S. B. Navathe, "An efficient algorithm for mining association rules in large databases," in *Proc. 21th Int. Conf. Very Large Data Bases (VLDB)*, 1995, pp. 432-444.
- [13]. M. J. Zaki, "Scalable algorithms for association mining," *IEEE Trans. Knowl. Data Eng.*, vol. 12, no. 3, pp. 372-390, May 2000.
- [14]. M. J. Zaki and K. Gouda, "Fast vertical mining using diffsets," in *Proc. 9th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2003, pp. 326-335.
- [15]. G. Grahne and J. Zhu, "Fast algorithms for frequent itemset mining using FP-trees," *IEEE Trans. Knowl. Data Eng.*, vol. 17, no. 10, pp. 1347-1362, Oct. 2005.
- [16]. J. Han, J. Pei, Y. Yin, and R. Mao, "Mining frequent patterns without candidate generation: A frequent-pattern tree approach," *Data Mining Knowl. Discovery*, vol. 8, no. 1, pp. 53-87, 2004.
- [17]. J. Pei, J. Han, H. Lu, S. Nishio, S. Tang, and D. Yang, "H-mine: Hyperstructure mining of frequent patterns in large databases," in *Proc. IEEE Int. Conf. Data Mining*, Nov./Dec. 2001, pp. 441-448.
- [18]. I. Triguero, J. A. Saez, J. Luengo, S. Garcia, and F. Herrera, "On the characterization of noise filters for self-training semi-supervised in nearest neighbor classification," *Neurocomputing*, vol. 132, pp. 30-41, 2014.
- [19]. R.C.Fernandez et al., "Liquid: Unifying nearline & offline big data integration," in *Proc. 7th Biennial Conf. Innov. Data Syst. Res. (CIDR)*, 2015