

An Emphasized Apriori Algorithm for Huge Sequence of Datasets

Ravi Kumar.V

Associate Professor,
Department of IT,
TKRCET.

Dr.Purna Chander Rao

Research Scholar- JNTUH
Professor,
SVIET.

Abstract:

In this paper, an Emphasized Apriori algorithm is proposed based on the most popular Apriori algorithm to overcome its drawbacks which are nothing but time consumption and the memory space. Since the Apriori algorithm scans the entire database based on minimum support and minimum confidence, it consumes more time and also more space. In this approach, new candidate set was prepared by considering the minimum support such that the total number of items will be reduced intern reduces the time taken and also memory space. This approach also includes the concept of Frequent Pattern (FP) growth to delete the items which are not frequent. Finally, a FP tree is established based on the frequent itemsets such that the items which are not frequent were delete intern to reduce the space consumption.

Keywords: Sequential Pattern Mining, Apriori, Frequent Pattern, Itemsets, Time Consumption.

1. Introduction

With the quick development of the Internet and information, mining valuable and important learning or data is a basic issue as of late. Contingent upon particular applications [1], [2], the found learning can be named affiliation manage mining [3], characterization [4], bunching [5], and successive example mining [6], [7] among others [8]. Consecutive example mining (SPM) concerns the requested arrangement information, for example, DNA successions, utilization of Web log, Web-click streams, or the logs of system stream. It can likewise be utilized to anticipate the obtained practices of the clients in crate examination. For instance, a client

purchases a few socks in one exchange at shopping center, he or she will get a few shoes in a later exchange. In prior, Agrawal et al. proposed the AprioriAll calculation [6] to produce and-test contender for mining the successive examples from a static database. Pei et al. planned the PrefixSpan calculation to effectively mine the successive examples in view of the projection component [9]. A grouping database is recursively anticipated into a few littler arrangements of anticipated database to accelerate the calculations for mining consecutive examples. Zaki et al. proposed a SPADE calculation to quick mine the successive examples [10]. The SPADE calculation uses the combinational properties in view of the effective grid seek methods with join operations.

In view of SPADE calculation, the consecutive examples can be determined with three database filters. Numerous calculations have been proposed to mine the consecutive examples, however the greater part of them are performed to handle the static database. At the point when the arrangements are changed whether succession inclusion [11] or erasure [12] in the first database, the found consecutive examples may get to be invalid or new successive examples may emerge. A natural approach to overhaul the successive examples is to re-handle the redesigned database in cluster mode, which is wasteful in certifiable applications. To handle the dynamic database with succession inclusion, Lin et al. proposed a FASTUP calculation [13] to incrementally keep up and overhaul the found successive examples with arrangement addition. The first database is still, nonetheless, required to be rescanned if the found successive example is expansive in the additional groupings yet little in the first database in view of the

FASTUP idea. Hong et al. at that point amplified the pre-expansive idea of affiliation lead mining to separately keep up and overhaul the found consecutive examples with succession inclusion [14] and arrangement erasure [15] in a level-wise manner, which requires more calculations of numerous database rescans. Lin initially planned a quick redesigned consecutive example (FUSP)- tree and built up the calculations for productively taking care of an incremental database with grouping inclusion [16]. The FUSP tree is implicit progress before the arrangements are embedded into the first database.

Two sections with four cases are then separated in light of the FASTUP idea to keep up and overhaul the FUSP tree for later mining process. Lin et al. additionally proposed the support calculation for grouping cancellation [17]. The first database is, be that as it may, required to be rescanned in the event that it is important to keep up a succession which is little in the first database yet substantial in the embedded groupings with arrangement inclusion or a succession is little both in the first database and in the erased successions with grouping erasure. In this paper, an Enhanced Apriority approach is proposed to diminish the time utilization and memory space in consecutive example mining framework. The proposed approach depends on the most well known Apriority algorithm which is having an issue of exercise in futility and parcel of memory utilization. The improved Apriority algorithm beats these two issues by erasing the itemsets which are not visit. Whatever is left of paper is sorted out as takes after: area II depicts the points of interest of Apriority algorithm. Area III depicts the points of interest of the proposed Enhanced Apriority algorithm. The test results are represented in area IV lastly the conclusions are given in segment V.

2. Apriori Algorithm

Agrawal and Srikant [6] firstly proposed Apriority algorithm . This calculation depends on Apriori property which expresses each sub (k-1)- Itemset of regular k-Itemset must be visit. Two primary process

are executed in Apriority algorithm : one is applicant era handle, in which the bolster tally of the comparing sensor things is computed by checking value-based database and second is substantial itemset era, which is produced by pruning those hopeful Itemsets which has a bolster number not as much as least limit. These procedures are iteratively rehashed until competitor Itemsets or huge Itemsets gets to be unfilled as in illustration appeared in Fig 1. Unique database is filtered first time for the competitor set, comprises of one sensor thing and there support has checked, then these 1-Itemset applicants are pruned by just evacuating those things that has a thing tally not as much as client determined edge (in above case threshold=30%). In second pass database is filtered again to create 2-Itemset hopefuls comprise of two things, on the other hand pruned to delivered expansive 2-Itemset utilizing apriori property. As per apriori property each sub 1-Itemset of 2 regular Itemsets must be visit. This procedure closes as in fourth sweep of database 4-Itemset hopeful will be pruned and huge itemset will be unfilled.

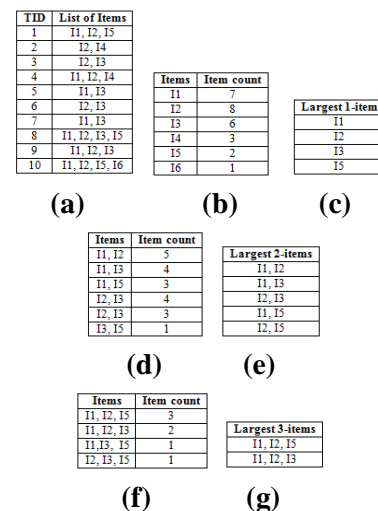


Figure.1 (a) original database, (b) candidate-1, (c) large-1 items, (d) candidate-2, (e) large-2 items, (f) candidate-3, (g) largest-3 items

There are two constraints of this calculation: one is unpredictable hopeful itemset era handle which expends substantial memory and tremendous execution time and second issue is over the top database checks

for competitor era. By and large there are two approaches to beat these constraints: one route is to investigate distinctive pruning and separating systems to make competitor Itemset littler. Second approach is either supplant unique database with subset of exchange in light of huge incessant Itemset or minimizes the quantity of sweeps over the database.

3. Enhanced Apriori Algorithm

Despite being straightforward and clear, the Apriority algorithm experiences some shortcoming. Squandering of expensive time for checking all out datasets and the lower least support or bigger itemsets is the principle confinement couch priori calculation. Since, it has exponential multifaceted nature; it devours parcel of memory furthermore reaction time. For instance, if the quantity of exchanges is 100, then the aggregate number of itemsets will be 2100 furthermore it does mining twice. This issue can be illuminated by decreasing the quantity of itemsets by continuous itemsets mining (FIM) assistant diminishes the time required for mining. Be that as it may, the fundamental impediment with FIM is the utilization of a great deal of space and gets to be wasteful for ongoing itemset applications.

In this paper, another example mining methodology is proposed to beat the issue of ordinary Apriori and FIM calculations. We should characterize a few definitions, for example, be the exchange set, be the arrangement of things in every single exchange and k-thing set is additionally an itemset with the end goal that it is a subset to thing set I, . Here another parameter s characterized as the support or the recurrence of event of thing, characterized as , applicant itemset of size k and be the successive itemset of size k. We first output every one of the exchanges to produce F1 which contains the things, their bolster check and Transaction IDs where the things are found. And afterward we utilize F1 as a partner to create F2, F3 ... , Fk. When we need to produce C2, we make a self-join F1 * F1 to develop 2-itemset C (x, y), where x and y are the things of C2. Before filtering all the exchange records

to check the bolster tally of every competitor, utilize F1 to get the exchange IDs of the base bolster tally amongst x and y, and in this way examine for C2 just in these particular exchanges. A similar thing applies for C3, build 3-itemset C(x, y, z), where x, y and z are the things of C3 and utilize F1 to get the exchange IDs of the base bolster tally between x, y and z and after that output for C3 just in these particular exchanges and rehash these means until no new continuous itemsets are distinguished. Presently to decrease the memory space when extensive exchanges are there a straightforward administer can be taken after: Let n be the quantity of hubs in the FP-tree and k be the shade of the groups of the exchanges in the database. Presently, absolutely $n > k$. If so then k is at most $n - 1$.

Assume we have 1000 exchanges then k will be at generally 999. There are such a large number of conceivable outcomes of hues and every one of the hues will be picked by us. All things considered, unmistakably that prompts to a terrible decision. Presently, let $n \geq k$ as this can likewise be conceivable then k will be at most n yet at the same time the run applies as n can't be not as much as k since then at every level hubs will have a similar shading. It must be same if the tree is completely needy. Since it consumes exponential memory room, the potential outcomes of hues getting produced ought to be minimized. This should be possible by utilizing another scientific recipe for looking at the quantity of hubs and hues i.e. $n > 2k$.

For this situation hues will be minimized definitely eg; if $n = 1000$ now then k will be around $\log_2(1000) = 10$. The base 2 connotes that the bunch is getting apportioned into 2 sections and selecting implies out of the two just 1 is getting chosen. This can be any number of segments relying upon client's decision. Client will have the decision of choosing the base. The estimation of the base is equivalent to the quantity of allotments of the group. Utilizing this approach less memory space is devoured at once and things can be mines in a lesser measure of time. Consequently, it fills the need.

Give us a chance to accept an extensive grocery store tracks deals information by stock-putting away unit (SSU) for every thing, for example, "Sugar", "Dal", "Drain", "Wheat", "Oil", "Rice" is distinguished by a numerical SSU. The general store has a database of exchanges where every exchange is an arrangement of SSUs that were purchased together. Give the database of exchanges a chance to comprise of taking after itemsets: The exchange set as appeared in Table 1. At first, filter all exchanges to get visit 1-itemset I1 which contains the things and their bolster tally and the exchanges ids that contain these things, and after that dispense with the competitors that are not visit or the thing having bolster not exactly the base support as appeared in table 2.

Table 1. The Transactions

Transactions	Item sets
T1	Rice, oil
T2	Rice, wheat, sugar
T3	Milk, dal
T4	Dal, sugar, oil
T5	Wheat, rice
T6	Rice, dal, sugar, milk
T7	Rice, dal, sugar, milk, oil

Table 2. The candidate 1- itemset

Items	support
Rice	5
Oil	3
Wheat	2
Dal	4
Sugar	4
milk	3

The regular 1-itemset is appeared in table 3. The sets which are in intense will be erased in regular 2-itemset as appeared in table 4. The sets which are in strong will be erased in successive 3-itemsets as appeared in table 5.

Table 3. The frequent 1- itemset

Items	support	Transaction IDs (T_ID)
Rice	5	T1, T2, T5, T6, T7
Oil	3	T1, T4, T7
Wheat	2	T2, T5
Dal	4	T3, T4, T6, T7
Sugar	4	T2, T4, T6, T7
milk	3	T3, T6, T7

Table 4. The frequent 2- itemset

Items	support	Min	Found in
Rice, oil	2	Oil	T1, T7
Rice, Dal	2	Dal	T6, T7
Rice, Sugar	3	Sugar	T2, T6, T7
Rice, Milk	2	Milk	T6, T7
Oil, Dal	2	Dal	T4, T7
Oil, Sugar	2	Oil	T4, T7
Oil, Milk	1	Oil	T7
Dal, Sugar	3	Dal	T4, T6, T7
Dal, Milk	3	Dal	T3, T6, T7
Sugar, Milk	2	Milk	T6, T7

Table 5. The frequent 3- itemset

Items	support	Min	Found in
Rice, Sugar, Dal	2	Dal	T6, T7
Rice, Sugar, Milk	2	Milk	T6, T7
Dal, Sugar, Milk	2	Dal	T6, T7

The following stride is to produce competitor 2-itemset from L1 split each itemset in 2-itemset into two components then utilize I1 table to decide the exchanges where you can discover the itemset in, instead of looking for them in all exchanges. For instance, we should take the primary thing in table.4 (Rice, Oil), in the first Apriori we check every one of the 7 exchanges to discover the thing (Rice, Oil); however in our proposed enhanced calculation we will part the thing (Rice, Oil), into Rice and Oil and get the base support between them utilizing L1 has the littlest least support. After that we hunt down itemset (Rice, Oil) just in the exchanges T1 the base certainty, and

afterward create all applicant affiliation rules. In the past case, on the off chance that we tally the quantity of examined exchanges to get (1, 2, 3)- itemset utilizing the first Apriori and our enhanced Apriori, we will watch the undeniable distinction between number of checked exchanges with our enhanced Apriori and the first Apriori. From the table 6, number of exchanges in1-itemset is the same in both of sides, and at whatever point the k of k-itemset increment, the crevice between our enhanced Apriori and the first Apriori increment from perspective of time expended, and thus this will decrease the time devoured to create applicant bolster tally. To get bolster mean each itemset, here Oil, and T7. For a given regular itemset LK, T4, discover all non-discharge subsets that fulfill the base certainty, and afterward produce all applicant affiliation rules. In the past case, on the off chance that we tally the quantity of examined exchanges to get (1, 2, 3)- itemset utilizing the first Apriori and our enhanced Apriori, we will watch the conspicuous contrast between number of checked exchanges with our enhanced Apriori and the first Apriori. From the table 6, number of exchanges in1-itemset is the same in both of sides, and at whatever point the k of k-itemset increment, the hole between our enhanced Apriori and the first Apriori increment from perspective of time devoured, and consequently this will lessen the time expended to produce competitor bolster check.

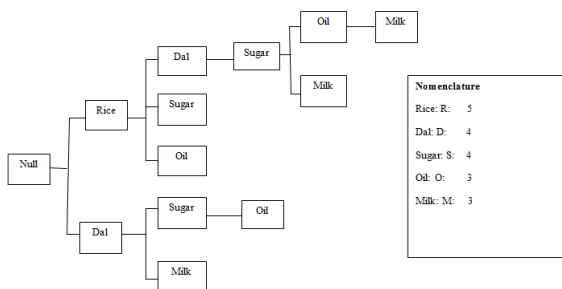


Figure.2. FP-tree of the above illustrated Example

The last yield of the FP-Tree is as appeared in Figure.2. Also, the base bolster tally is 3. Presently locate the incessant examples from the FP-Tree. It's minor. The things of the database and their recurrence of events are appeared in Table:2 for everything.

Above all else, we have to organize all the itemsets as indicated by their recurrence of events and afterward we will see everything one by one from base to beat. The things can be recorded as: Then we see Milk. To begin with we have to locate the contingent example base for Milk:3. This is because of the recurrence of event of Milk. Presently go to Graph 1 and check the Milks. There are 3 Milks and one event for each. Presently cross base to best and get the branches which have Milks with the event of Milk. We got 3 branches and they are RDSO: 1, D: 1, RDS: 1. To guarantee that you accurately got every one of the events of Milk in FP-Tree include events of every branch and contrast and the events recorded previously. For Milk we get 1+1+1 = 3 so it is right. At that point, consider Oil. What's more, along these lines the accuracy will be guaranteed for residual things. By doing this for all things we can erase all different branches aside from that and just that branch will stay in the FP-Tree which we can draw again for Milk and similarly for every single other thing.

4. Experimental Analysis

The performance evaluation of proposed approach is examined by applying on various datasets. For every itemset, the performance was measured by measuring the time taken for every transaction and the reduced amount of time at every transaction, the execution time will vary with number of transactions. As the number of transaction increases, the time taken for scanning also increases. Initially, the time taken for execution is evaluated for all transaction and the one transaction is considered for evaluation with varying minimum support. The obtained results are summarized as follows:

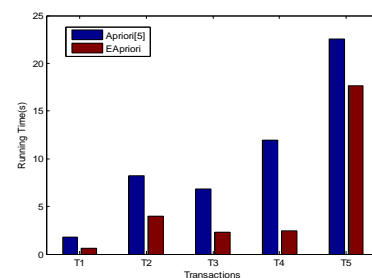


Figure.3. Running time for all transactions

Figure.3. shows the comparative analysis of the proposed approach with earlier Apriori algorithm with respect to the time taken for execution. The running time varies from transaction to transaction. As the number of transactions increases, the time required for execution also increases. From the above figure, it can be observed that the running time for proposed approach is less compared to Apriori algorithm.

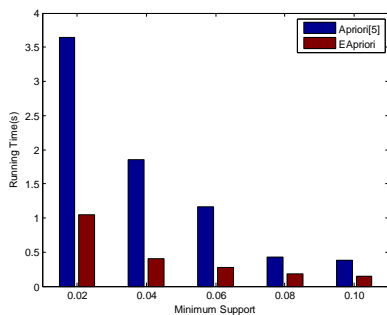


Figure.4. Running time variations w.r.t minimum support

Figure.4 shows the time taken for execution with respect to minimum support. As the minimum support increase, the time taken for execution will be decreased. The above figure compares the Apriori with the proposed approach with respect to execution time for varying minimum support. From the above figure, it can be observed that for the proposed approach, the time taken for execution at each and every minimum support is less compared to Apriori.

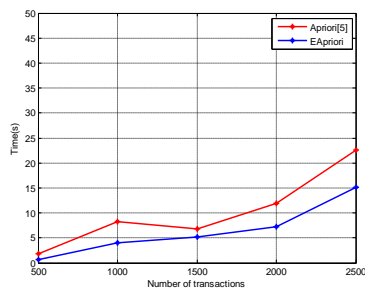


Figure.5. Time consuming comparison for different groups of transactions.

Figure.5 describes the details of the time taken for execution with varying number of transaction. The number of transactions is varied from 500 to 2500 and the respective time consumption results for both the

Apriori and the proposed EApriori are shown in figure.5. Since the proposed approach considers the frequent itemset as one more parameter to perform mining, the time taken by EApriori must be less when compare to Apriori. In the above figure, the time taken is increasing with number of transaction, but the increment is less for EApriori compared to Apriori.

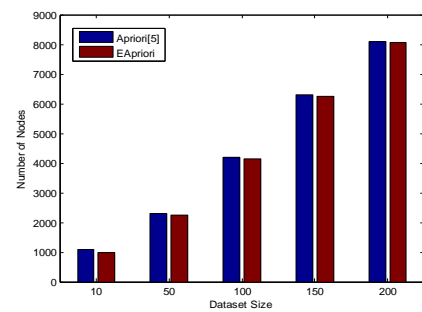


Figure.6. Number of nodes of tree for varying data size

Figure.6 describes the number of node obtained for varying data side such that the number of itemsets. Generally, the dataset size increases, the number of nodes also increases. From the above figure, the proposed approach almost approaches the earlier approach.

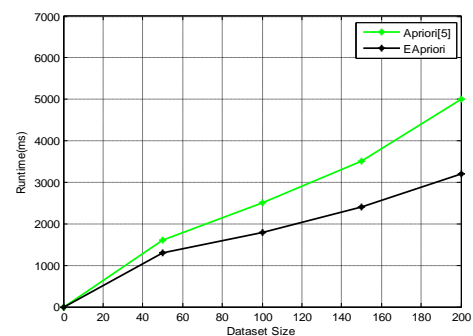


Figure.7. Runtime (ms) variations with varying dataset size

The runtime generally increases with an increase of dataset size. The above figure illustrates the variations of runtime in milliseconds for a varying dataset size. Form the above figure, it can be observed that the proposed approach has less increment in the runtime compared to earlier approach.

5. Conclusion

In this paper, a new approach was proposed for sequential pattern mining based on the earlier Apriori algorithm. This approach successfully reduces the memory space and also reduces the time required for execution even for large datasets. The enhancement if the proposed approach can be observed when there is an increment in the k itemsets. The time consumed to generate candidate support count in our enhanced Apriori is less than the time consumed in the original Apriori. The results also reveal the efficiency of proposed approach in the view of time consumption.

References

- [1] 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
- [2] M.-S. Chen, J. Han, and P. S. Yu, "Data mining: An overview from a database perspective," *IEEE Trans. Knowl. Data Eng.*, vol. 8, no. 6, pp. 866_883, Dec. 1996.
- [3] R. Agrawal and R. Srikant, "Fast algorithms for mining association rules in large databases," in *Proc. Int. Conf. Very Large Data Bases*, 1994, pp. 487_499.
- [4] S. B. Kotsiantis, "Supervised machine learning: A review of classification techniques," in *Proc. Conf. Emerg. Artif. Intell. Appl. Comput. Eng., Real Word AI Syst. Appl. eHealth, HCI, Inf. Retr. Pervasive Technol.*, 2007, pp. 3_24.
- [5] P. Berkhin, "A survey of clustering data mining techniques," in *Grouping Multidimensional Data*. Berlin, Germany: Springer-Verlag, 2006, pp. 25_71.
- [6] R. Agrawal and R. Srikant, "Mining sequential patterns," in *Proc. Int. Conf. Data Eng.*, 1995, pp. 3_14.
- [7] C. H. Mooney and J. F. Roddick, "Sequential pattern mining Approaches and algorithms," *ACM Comput. Surveys*, vol. 45, no. 2, pp. 1_39, Feb. 2013.
- [8] C.-W. Lin, T.-P. Hong, and W.-H. Lu, "An effective tree structure for mining high utility itemsets," *Expert Syst. Appl.*, vol. 38, no. 6, pp. 7419_7424, Jun. 2011.
- [9] J. Pei et al., "Mining sequential patterns by pattern-growth: The PrefixSpan approach," *IEEE Trans. Knowl. Data Eng.*, vol. 16, no. 11, pp. 1424_1440, Nov. 2004.
- [10] M. J. Zaki, "SPADE: An efficient algorithm for mining frequent sequences," *Mach. Learn.*, vol. 42, nos. 1_2, pp. 31_60, Jan. 2001.
- [11] D. W. Cheung, J. Han, V. T. Ng, and C. Y. Wong, "Maintenance of discovered association rules in large databases: An incremental updating technique," in *Proc. 25th Int. Conf. Data Eng.*, Mar. 1996, pp. 106_114.
- [12] D.W.-L. Cheung, S. D. Lee, and B. Kao, "A general incremental technique for maintaining discovered association rules," in *Proc. Int. Conf. Database Syst. Adv. Appl.*, Apr. 1997, pp. 185_194.
- [13] M.-Y. Lin and S.-Y. Lee, "Incremental update on sequential patterns in large databases," in *Proc. IEEE Int. Conf. Tools Artif. Intell.*, Nov. 1998, pp. 24_31.
- [14] T.-P. Hong, C.-Y. Wang, and S.-S. Tseng, "An incremental mining algorithm for maintaining sequential patterns using pre-large sequences," *Expert Syst. Appl.*, vol. 38, no. 6, pp. 7051_7058, Jun. 2011.
- [15] C.-Y. Wang, T.-P. Hong, and S.-S. Tseng, "Maintenance of sequential patterns for record deletion," in *Proc. IEEE Int. Conf. Data Mining*, Nov. 2001, pp. 536_541.



[16] C.-W. Lin, T.-P. Hong, W.-H. Lu, and W.-Y. Lin, "an incremental FUSP-tree maintenance algorithm," in Proc. 8th Int. Conf. Intell. Syst. Design Appl., Nov. 2008, pp. 445_449.

[17] C.-W. Lin, T.-P. Hong, and W.-H. Lu, "An efficient FUSP-tree update algorithm for deleted data in customer sequences," in Proc. Int. Conf. Innovative Comput., Inf. Control, Dec. 2009, pp. 1491-1494.