

AN IMPROVED APRIORI ALGORITHM RUCHA KALE¹, SHARAYU FUKEY²

¹ V.J.T.I.MUMBAI-19

² V.J.T.I.MUMBAI-19

Abstract

Since the development in the networks, information technology, Internet facilities is attracting more and more data generation from the users, data is being generated massively and increasing like a fire. Due to spreading of this information with high speed, more emphasis is being given on handling of such data and using the data mining rules for deriving associations. The provided study describes the Apriori algorithm in association rule mining algorithm in detail and also the algorithm implementation is discussed. Also it proposes a new optimization algorithm called APRIORI-IMPROVED which is based on the insufficient of traditional Apriori. It illustrates the optimization in terms of generation, transactions, compression, and performance and so on. APRIORI-IMPROVED provides a hash structure an efficient horizontal data representation and optimized strategy of storage for saving time and space. The performance study shows that APRIORI-IMPROVED is much faster than simple Apriori.

Keywords: Author Guide, Article, Camera-Ready Format, Paper Specifications, Paper Submission.

1. Introduction

In spite of rapid development in modern computer technology and Database technology could support the store and quickly retrieve the huge scale databases or data warehouses, but these tactics were only to gather these massive data and not to organize them effectively and use the knowledge hidden with them, which eventually lead to today's world's phenomenon of "rich data but poor knowledge"[1]. The emergence of data mining technology met people needs by providing them only the data specific to them and facilitating them suggestions to make faster choices in selecting some data item sets. Association rule mining is

to find relationships and correlations amongst these items in large databases consisting transactions[1].

An influential algorithm Apriori, was first discussed by R. Agrawal and R. Shrikant. Apriori provides an iterative approach known as level-wise and breadth first approach, in which k- data item sets are used to generate (k+1)-data item sets. In terms of the feature of Apriori property, called anti monotone, one can efficiently generate candidate item sets, by discarding unnecessary remaining ones. Apriori algorithm uses a two-step process Join and Prune[2]. However there are two major drawbacks of the algorithms based on generated and tested candidate item sets (1) Database has to be scanned multiple number of times in order to generate candidate item sets. Multiple scans indeed increases the I/O load and it's respective time consumption. (2) A lot of CPU time will be consumed by the generation of huge candidate sets and calculation of their support [4].

1.1 2. Apriori Algorithm Summary

Apriori algorithm is one of the most efficient algorithms used for mining the frequent data sets of Boolean association rules. The algorithm is an approach based on the two-stage frequency, where the design of Association rule mining algorithm can be distinguished into two sub-issues:

- (1) Find all the item sets those having support greater than the minimum support, which is known as the frequent items.
- (2) Based on these obtained frequent sets, all the association rules will be derived, and for each frequent item set A, all non-

empty subset of A will be further derived if it satisfies the condition as stated as below

Support(A) / Support(a) >= minimum_confidence, to generate the association rule A-a.

That is from the frequent sets found in the first step, to exploit the rules with confidence not less than minimum confidence minimum_confidence has to be specified by the user.[1]

Apriori algorithm consists two sub processes named as Apriori-generate () and subset ().

Apriori-generate() actually generates a candidate, further uses an Apriori property that is all non-empty frequent sets must also be frequent item sets; in order to delete all those candidates of the non-frequent subsets obtained. After the completion of generation of all the candidates, the database D will be scanned, as well as for each of the transactions, use the subset () sub procedure for identifying all the candidate subsets and make cumulative count for each of the subsets. Finally all candidates meet the minimum support of frequent item set L[2].

3. Apriori Illustration

The following example represents general representation of Apriori Algorithm.

(A) Generating Frequent item sets.

Following table (1) shows the transaction table consisting 4 datasets |D|=4
 Let us consider minimum support count be 2.

100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

The specific process is described as below:

- (1) Scan the complete database in order to initialize the source data and to form the candidate-1 with all the itemsets C1 in an aggregation.
- (2) To find the frequent 1-itemsets L1, each of the candidate itemsets are checked for the occurrences and the condition for minimum_support count is tested. This condition enables to eliminate all other non-required itemsets that cannot participate into generation of association rules.

(B) Generate association rules from frequent item sets

After obtaining frequent item sets, the next step is to find the association rules from them respectively.

Here, the conditional property can be calculated using the support for item sets

$$Confidence (X \rightarrow Y) = \frac{\partial (X \rightarrow Y)}{\partial x} * 100\% = \frac{support(X \cup Y)}{support(X)} * 100\%$$

TID	ITEMS
-----	-------

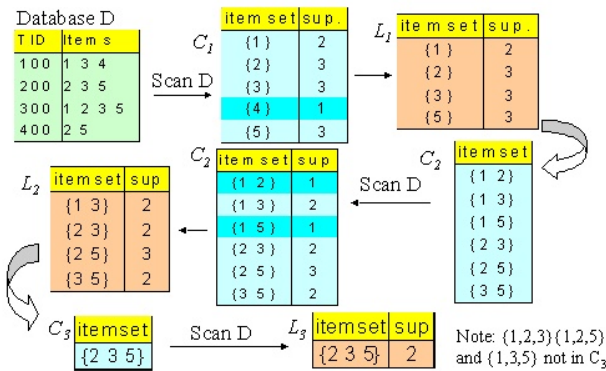


Fig (1) Selection of frequent item sets in Apriori transitions.

Depending on the above selected frequent items and their respective eliminations depending on the minimum support criterion, association rules are derived based on the conditions of condition.

4. Basic Apriori algorithm evaluation

The basic Apriori algorithm works efficiently and suitably but it consists of certain shortcomings that are difficult to overcome. Those are enlisted below:

- (1) Repeated scanning of the database which indeed increases the CPU overhead.
 - (2) This can generate comparatively large candidate set, from $L(k-1)$ to generate k item sets, and C_k proves to be the exponential growth.
 - (3) Reducing the affairs need to be scanned in later circle, the one affair does not contain all $(k+1)$ -item sets
 - (4) The fitness landscape of the algorithm is very narrow. It only considers a single Boolean association rule mining.
- However in practical applications, there may be more complex scenarios such as multi-dimensional database, multivolume and multi-layer association rules.

5. APRIORI-IMPROVED ALGORITHM

The following section will address the improved Apriori ideas.

In the process of Apriori, the following definitions are needed:

Definition 1: Suppose $T = \{T_1, T_2, \dots, T_m\}$, ($m \geq 1$) is a set of transactions, $T_i = \{I_1, I_2, \dots, I_n\}$, ($n \geq 1$) is the set of items, and k -item set $= \{i_1, i_2, \dots, i_k\}$, ($k \geq 1$) is also the set of k items, and k -item set $\subseteq I$.

Definition 2: Suppose σ (item set), is the support count of item set or the frequency of occurrence of an item set in transactions. **Definition 3:** Suppose C_k is the candidate item set of size k , and L_k is the frequent item set of size k [4].

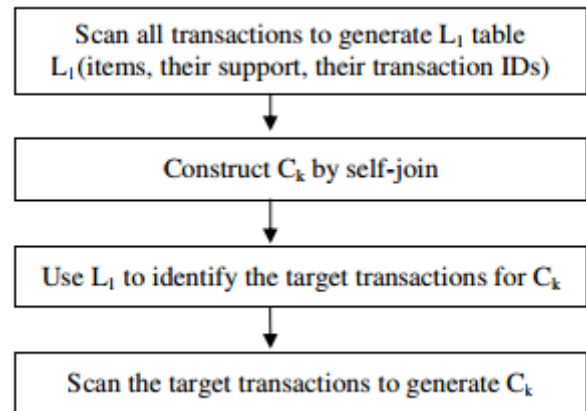


Fig (2) Steps representing C_k generation. The improvement of algorithm can be elaborated as follows[4]:

- //Generate items, items support, their transaction ID
- (1) $L_1 = \text{find_frequent_1_itemsets}(T)$;
 - (2) For ($k = 2; L_{k-1} \neq \Phi; k++$) {
 //Generating the C_k from the L_{k-1}
 - (3) $C_k = \text{candidates that are generated from } L_{k-1}$;
 //getting the item with minimum support in C_k using L_1 , ($1 \leq w \leq k$).
 - (4) $x = \text{Get_item_minimum_support}(C_k, L_1)$;
 // getting the target transaction IDs that contain item x .

- (5) Tgt = get_Transaction_ID(x);
- (6) For each transaction t in Tgt Do
- (7) Increment the count of all the items in Ck that are found in Tgt;
- (8) Lk= items in Ck \geq minimum_support;
- (9) End;
- (10) }

6. IMPROVED-APRIORI IMPLEMENTATION

Assume we have transaction set D has 9 transactions, and the minimum support = 3. The transaction set is shown in Table.2

T_ID	Items
T ₁	I ₁ , I ₂ , I ₅
T ₂	I ₂ , I ₄
T ₃	I ₂ , I ₄
T ₄	I ₁ , I ₂ , I ₄
T ₅	I ₁ , I ₃
T ₆	I ₂ , I ₃
T ₇	I ₁ , I ₃
T ₈	I ₁ , I ₂ , I ₃ , I ₅
T ₉	I ₁ , I ₂ , I ₃

Table(2) Transaction table considered for Improved-Apriori implementation

Items	support	
I ₁	6	
I ₂	7	
I ₃	5	
I ₄	3	
I ₅	2	deleted

Table(3) The candidate-1 item set

Initially, all the transaction are scanned to get frequent 1-itemset L1 which consists of the items as well as their support count and the transactions ids that contain these items, and then eliminate the candidates that are not frequent or their respective support are less than the minimum_support. The frequent 1-itemset is shown in table 4.

Items	support	T_IDs	
I ₁	6	T ₁ , T ₄ , T ₅ , T ₇ , T ₈ , T ₉	
I ₂	7	T ₁ , T ₂ , T ₃ , T ₄ , T ₆ , T ₈ , T ₉	
I ₃	5	T ₅ , T ₆ , T ₇ , T ₈ , T ₉	
I ₄	3	T ₂ , T ₃ , T ₄	
I ₅	2	T ₁ , T ₈	deleted

Table(4) Frequent_1 item-set

The next step is generating candidate 2-itemset from L1. To get support count for every item set, each item set is splitted in 2-itemset into two elements then L1 table is used to determine the transactions where you can find the item set in, rather than searching for them in all transactions. for example, let's take the first item in table.4 (I₁, I₂), in the original Apriori we scan all 9 transactions to find the item (I₁, I₂); but in our proposed improved algorithm we will split the item (I₁, I₂) into I₁ and I₂ and get the minimum support between them using L1, here i1 has the smallest minimum support. After that we search for item set (I₁, I₂) only in the transactions T₁, T₄, T₅, T₇, T₈ and T₉.

Items	support	Min	Found in	
I ₁ , I ₂	4	I ₁	T ₁ , T ₄ , T ₅ , T ₇ , T ₈ , T ₉	
I ₁ , I ₃	4	I ₃	T ₅ , T ₆ , T ₇ , T ₈ , T ₉	
I ₁ , I ₄	1	I ₄	T ₂ , T ₃ , T ₄	deleted
I ₂ , I ₃	3	I ₃	T ₅ , T ₆ , T ₇ , T ₈ , T ₉	
I ₂ , I ₄	3	I ₄	T ₂ , T ₃ , T ₄	
I ₃ , I ₄	0	I ₄	T ₂ , T ₃ , T ₄	deleted

Table (5) Frequent_2 item sets

The same process is carried out to generate 3-itemset depending on L1 table, as it is shown in table 6.

Items	support	Min	Found in	
I ₁ , I ₂ , I ₃	2	I ₃	T ₅ , T ₆ , T ₇ , T ₈ , T ₉	deleted
I ₁ , I ₂ , I ₄	1	I ₄	T ₂ , T ₃ , T ₄	deleted
I ₁ , I ₃ , I ₄	0	I ₄	T ₂ , T ₃ , T ₄	deleted
I ₂ , I ₃ , I ₄	0	I ₄	T ₂ , T ₃ , T ₄	deleted

Table (6) Frequent_3 item sets

For a given frequent item set LK, all non-empty subsets that satisfy the minimum confidence are

found, and then all candidate association rules are generated. Using the original Apriori and our improved Apriori, the obvious difference between number of scanned transactions with our improved Apriori and the original Apriori will be considerably observed. From the table 7, number of transactions in 1-itemset is the same in both of sides, and whenever the k of k-item set increases, the gap between our improved Apriori and the original Apriori increase from view of time consumed, and hence the time consumed to generate candidate support count will be considerably reduced.

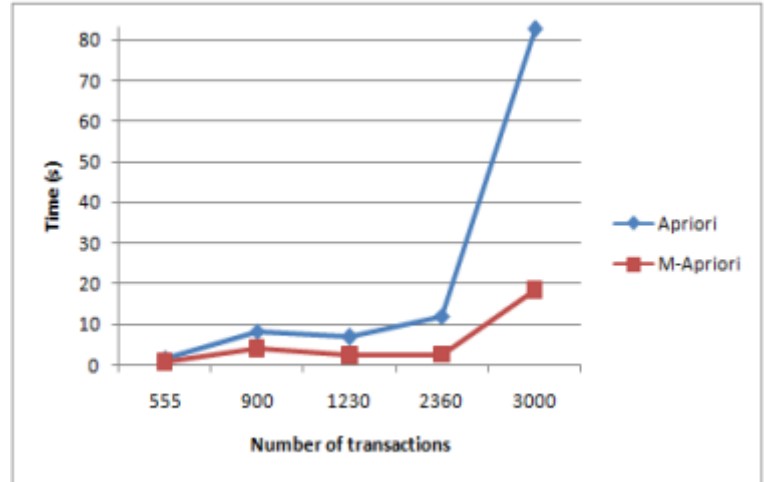
	Original Apriori	Our improved Apriori
1-itemset	45	45
2-itemset	54	25
3-itemset	36	14
sum	135	84

Table (7) Number of transaction scanned elements

We developed an implementation for original Apriori and our improved Apriori, and we collect 5 different groups of transactions as the follow:

- T1: 556 transactions.
- T2: 800 transactions.
- T3: 1250 transactions.
- T4: 2500 transactions.
- T5: 3100 transactions.

The first experiment compares the time consumed of original Apriori, and the proposed improved algorithm by applying the above five groups of transactions in the implementation. The result is shown in following Figure 3.



Figure(3) : Time consuming comparison for different groups of transactions

The second experiment is used to compare the time consumed of original Apriori, and our proposed algorithm by applying the one group of transactions through various values for minimum support in the implementation. The corresponding result is shown in following Figure 4.

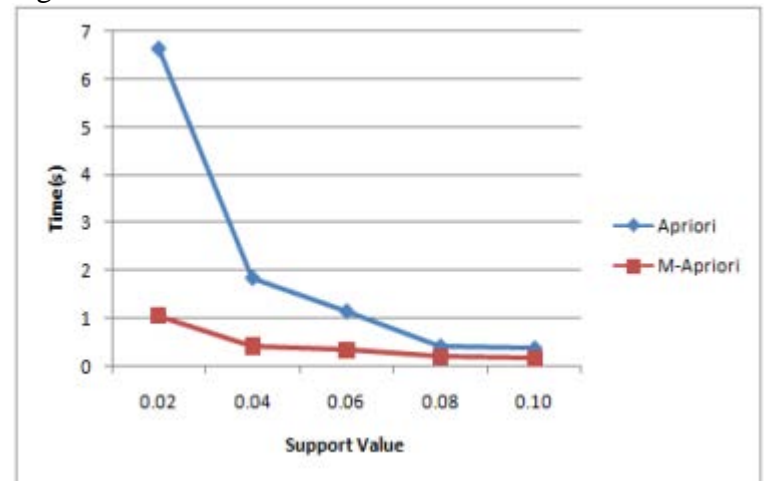


Figure (4): Time consuming comparison for different values based on minimum support

Evaluation of IMPROVED-APRIORI:

The time consuming in improved Apriori in each group of transactions is less than it in the original Apriori, and the difference gradually increases corresponding to the number of transactions[5]

T	Original Apriori (S)	Improved Apriori (S)	Time reducing rate (%)
T ₁	1.776	0.677	61.88%
T ₂	8.221	4.002	51.32%
T ₃	6.871	2.304	66.41%
T ₄	11.940	2.458	79.16%
T ₅	82.558	18.331	77.80%

Table (8) The time reducing rate of improved Apriori on the original Apriori as per the number of transactions

As the value of minimum support increase the rate is decreased also. The average of reducing time rate in the improved Apriori is 68.39%.

Min_Sup	Original Apriori (S)	Improved Apriori (S)	Time reducing rate (%)
0.02	6.638	1.056	84.09%
0.04	1.855	0.422	77.25%
0.06	1.158	0.330	71.50%
0.08	0.424	0.199	53.07%
0.10	0.382	0.168	56.02%

Table (9) The time reducing rate of improved Apriori on the original Apriori according to the values of minimum support obtained.

4. Conclusions

In this paper, an improved Apriori is proposed towards decreasing the time consumed in transactions scanning for candidate item sets by decreasing the total number of transactions to be scanned. The gap between our improved Apriori and the original Apriori increases from view of time consumed whenever the k of k-item set increases; and the gap between our improved Apriori and the original Apriori decreases whenever the value of minimum support increases from view of time consumed. The time consumed to generate candidate support count in our improved Apriori is less than the time consumed in the original Apriori; our improved Apriori reduces the time consuming by 67.38%. As this is proved and validated by the experiments and obvious in figure 3, figure 4, table 8 and table 9.

References

- (1) Yanxi Liu, "Study on Application of Apriori Algorithm in Data Mining", Second International Conference on Computer Modeling and Simulation, 2010
- (2) Yuntao Liu, Qing Liu, Danweng Zheng, Qingping Deng, Zhanpeng Tang, Jian Ying, "Application and Improvement Discussion about Apriori Algorithm of Association Rules Mining in Cases Mining of Influenza treated by contemporary famous old Chinese medicine", IEEE International Conference on Bioinformatics and Biomedicine Workshops, 2012
- (3) Yong Li, "The Java Implementation of Apriori Algorithm based on Agile Design Principles",.
- (4) Riu Chang, Zhiyi Liu, "An Approved Apriori Algorithm", International Conference on Electronics and Optoelectronics (ICOEO), 2011
- (5) Gang Yang, Hong Zhao, Lei Wang, Ying Liu, "An Implementation of Improved Apriori Algorithm", Proceedings of the Eighth International Conference on Machine Learning and Cybernetics, Baoding, July 2009
- (6) Mohammad Al-Maolegi , Bassam Arkok, "An Improved Apriori Algorithm" , International Journal on Natural Language Computing (IJNLC) Vol.03, February 2014