# Global Journals LATEX JournalKaleidoscope<sup>TM</sup>

Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.

# An Improved Apriori Algorithm based on Matrix Data Structure

Shalini Dutt<sup>1</sup>, Naveen Choudhary<sup>2</sup> and Naveen Choudhary<sup>3</sup> 2 <sup>1</sup> College of Technology and Engineering, Maharana Pratap University of Agriculture 3 and Technology, Udaipur, Rajasthan, India 4 Received: 11 December 2013 Accepted: 3 January 2014 Published: 15 January 2014 5

#### Abstract 7

Mining regular/frequent itemsets is very important concept in association rule mining which 8 shows association among the variables in huge database. the classical algorithm used for 9 extracting regular itemsets faces two fatal deficiencies .firstly it scans the database multiple 10 times and secondly it generates large number of irregular itemsets hence increases spatial and 11 temporal complexities and overall decreases the efficiency of classical apriori algorithm.to 12 overcome the limitations of classical algorithm we proposed an improved algorithm in this 13 paper with a aim of minimizing the temporal and spatial complexities by cutting off the 14 database scans to one by generating compressed data structure bit matrix(b matrix)-and by 15 reducing redundant computations for extracting regular itemsets using top down method. 16 theoritical analysis and experimental results shows that improved algorithm is better than 17 classical apriori algorithm. 18

19

Index terms— Apriori algorithm, frequent itemsets, association rule. 20

#### Introduction 1 21

n last few decades data has become so vast that extracting information from this huge data becomes very 22 important issue in data mining . hence data mining brought into scene from last few years. Mining association 23 rules is important process in data mining which shows relationship among the variables or affairs stored in data 24 warehouse, database and other information repositories. Association rule mining is two step process. First 25 it generates regular/frequent itemset set of items having count equal or greater than user specified parameter 26 i.e., minimum support and second it discovers association rules from these frequent itemsets. In this regard first 27 association rule mining algorithm apriori algorithm was proposed in 1994 to discover regular itemset. Limitations 28 of apriori results in lot of research in the field of data mining to build more efficient algorithms in respect of space 29 and time. This paper puts forward an improved algorithm using matrix data structure with simply counting 30 rows and columns and transaction reduction strategies using top down approach for finding out largest regular 31 itemset to smallest regular itemset. 32

In this way, it can greatly reduces complexity and increases the efficiency of improved algorithm. 33

The remaining section of this paper is organized as follows: Section 2 contains Apriori Algorithm. In section 3, 34 elaborate the proposed improved algorithm for extracting regular itemset and. Section 4 contains experimental 35 results and conclusion. 36

#### $\mathbf{2}$ II. 37

#### 3 Apriori Algorithm 38

The apriori algorithm is standard and classical algorithm for mining regular/frequent itemsets (if an itemset 39 satisfies minimum threshold i.e, min\_support, it is called regular itemset. The set of regular k-itemsets is 40 commonly denoted by LK.) brought by R. Agarwal and R. Shrikant in 1994, that leads to generate association 41

42 rules called association rule mining. It uses bottom-up and iterative approach known as breadth first search

45 space.

Apriori property? all non-empty subsets of a regular itemset must also regular. For example if  $\{1,2,3\}$  is a 46 regular itemset, then  $\{(1,2), (1,3), (2,3), (1), (??) \text{ and } (3)\}$  are must be regular itemset. it uses two key operations 47 first, Join operation? to discover regular k-itemset, a set of candidate k-itemset is generated by joining LK-1 with 48 itself. Second, Prune operation?discard the itemset if support count of itemset is less than minimum support 49 required and also discard the itemset if it is not following apriori property, remove itemset that have subset 50 that is not regular. apriori algorithm generates regular itemset as follows First, regular 1-itemsets are discovered 51 by scanning the transactional database to calculate the count of each item and select the items those satisfying 52 required minimum support , denoted by L1. Next L1, regular 1-itemset is used to discover L2 (set of regular 53 2itemset) and L2 is used to discover L3 (set of regular 3itemset) and so on until no more regular itemsets can be 54 discovered and generation of each Lk needs one complete scan of the database. Advantage of this algorithm is, 55 it is simple and easy to mine regular itemsets if database is small. also faces two fatal deficiencies. First, it scans 56 57 the database multiple times, so greatly high the I/O cost and second, generate large no. of candidate itemsets if 58 database is huge and overall decreases the efficiency of algorithm.

### 59 4 Improved Apriori Algorithm

The improved algorithm proposed in this paper works in two phases. In first phase required compressed data 60 structure i.e, b\_matrix is constructed and then this compressed data structure is used in second phase to 61 generate regular itemsets. This algorithm employs topdown approach to discover regular item sets from largest 62 regular item set to smallest regular itemset. Algorithm steps 1. In first phase b\_matrix is constructed for the 63 given transactional database. rows in b\_matrix represents each transaction and column represents items in 64 transactional database. In b\_matrix, each cell will contain values either 1 or 0 for showing the existence of items 65 in transactional database. Entry value will be 1, if the item is present in the respective row else 0, if the item is 66 67 absent in the row. With two more columns count and redundant transaction counter (TC). Here count column represents the size of row(the sum of total no of 1's in that particular row) and remove those columns whose sum 68 is not equal or greater than predefined min\_support value and then update count column. If row is duplicated in 69 database then it is represented by the value in the redundant transaction counter column and delete unnecessary 70 71 duplicate transaction/row and if row is not duplicate then redundant transaction counter column is set to 1. Then rearrange the b\_matrix in descending order based on count column. This is our required compressed data 72 structure and here the phase 1 of our improved algorithm completes. 2. Now generate regular itemsets directly 73 from b\_matrix. Select first row from b\_matrix and match its count value with next row count respectively. If 74 the next row count is more or equal to the processing row count then do AND operation among the rows, if 75 result is same to the processing row item set structure then increase the count value of support of processing 76 row item set by one and continue this procedure of matching and AND operation through rest of the rows in 77 b\_matrix and then check the value of total support. If it is greater or equal to predefined min\_support count 78 then extract the item set and its subsets and move them to frequent array list. The same procedure will be 79 repeated for rest of the transactions in b\_matrix until all rows are not checked. The gain of improved algorithm 80 is that it lessen the no. of comparisons to mine largest regular item set for duplicate transactions and transactions 81 having smaller item sets in size that is count value (since they do not have all the items of row under process) 82 and another major advantage is once largest regular item set is discovered then its subsets are searched and 83 moved into frequent array list. While searching for next largest regular item set it checks first, transaction under 84 processing is previously present in frequent array list because of prior largest itemset and its subsets, if itemset is 85 already in frequent array list, it avoids number of comparisons needed to calculate the support count of itemset. 86 Hence decreases number of scans and time needed to extract the regular itemset. 87

### 88 **5** IV.

### 89 6 Illustration

Consider the implementation of this improved algorithm through a sample below. TABLE-I shows a transactional
 database consist of 9 transactions. Set the minimum support counts as 3(min\_support=3).

### <sup>92</sup> 7 Table 1 a) Phase-1

93 Step1-Scan the transactional database and convert it into desired compressed data structure that is b\_matrix 94 M9\*6. Where each row represents one transaction and column represents distinct items in whole transactional 95 database and last column i.e, count represents the size of row. In b\_matrix entry value will be 1, if item is 96 present in the corresponding row else 0, if item is not present in the corresponding row.

Step2-After this rearrange the b\_matrix in desending order based on count column after removing , those columns whose sum is not equal or greater than required minimum support value. Here min\_support is 3, hence remove columns for items 4 and 5 and update count column and also merge the duplicate rows in TC column of b\_matrix to reduce the computations for redundant rows for finding regular item set.

# <sup>101</sup> 8 b) Phase-2

Step3-Now select first row TID-8 and extract its itemset {1, 2, 3} and calculate its support count in b\_matrix 102 using AND operation with rows having count value equal or greater than its own count value. If AND operation 103 results in same item set structure as processing row's item set structure, then increase its support count value. 104 after complete AND operation, check value of support count of item set, if it is equal or greater than required 105 min\_support than it is frequent/regular, then move itemset with its subset into frequent array list and move to 106 next row. Here item set support is less than required min support; hence it is not regular move to next row. TID 107 ITEMS T1 11,12,15 T2 12,14 T3 12,13 T4 11,12,14 T5 11,13 T6 12,13 T7 11,13 T8 11,12,13,15 T9 11,12,13 ( D D D D 108 D D D D ) Year 2014 c Table 2 : (b\_matrix) TID I1 I2 I3 COUNT TC T8 1 1 1 3 2 T1 1 1 0 2 2 T3 0 1 1 2 2 109 T5 1 0 1 2 2 T2 0 1 0 1 1 110

Step4select next row TID-1 and extract its item set  $\{1, 2\}$ . After AND operation its support count is 4, hence it is regular, move item set with its subset into frequent array list. So here regular /frequent array list is  $\{(1, 2), (2,$ 

## <sup>114</sup> 9 Experimental Results

All the experiments are carried out on core i7 Intel based PC machine with 2 GB main memory, running on 115 window 7 operating environment and the program code is written in java. Our experimental benchmark dataset 116 is taken from artificial data set of IBM that is, T1014D100K datasets. there are 100000 data records / affairs and 117 870 items in T1014D100K dataset. The improved algorithm is compared with classical apriori algorithm using 118 same hardware, dataset and minimum support requirement. The output of both the algorithms is same which 119 demonstrates that our algorithm is competent. the time computation is started from the spot when the file is 120 read into the memory until all the regular itemsets are generated. redundant rate in transactional database is 121 high hence using compressed data structure we eliminate it and reduces space and time complexities in proposed 122 algorithm. from the results shown in figure 1. we can conclude that time and space expenses are lesser as 123 compared to classical algorithm in proposed algorithm. 124

## 125 10 Conclusion

126 This proposed algorithm for mining regular itemset using bit matrix needs only single scan of whole transactional

database to construct compressed data structure . hence greatly reduces I/O cost and it also doesn't generate irregular itemset. So improved algorithm decreases temporal complexity and spatial complexity and have higher

 $^{129}$   $\,$  efficiency as compared to our classical apriori algorithm.  $^1$ 

 $<sup>^1 \</sup>odot$  2014 Global Journals Inc. (US) An Improved Apriori<br/> Algorithm based on Matrix Data Structure



Figure 1: Figure 1 :

- [Chen et al. ()] 'A Improved Aproiri Algorithm based on Pruning Optimization and Transaction Reduction'. Z
  Chen , S Cai , Zhu Song , C . AMISEC , 2nd IEEE International Conference, 2011. p. .
- [Agrawal and Srikant] R Agrawal, R Srikant. Fast Algorithms for Mining Association Rules in Large Databases
  Clin Proceedings of the 20th International Conference on Very Large Databases, (I) 994 p.
- [Chang and Liu] 'An Improved Apriori Algorithm'. Rui Chang , Zhiyi Liu . ICEOE 2011, IEEE International
  Conference, 1 p. . (Frequent array list)
- I36 [Jia et al.] 'An Improved Apriori Algorithm Based on Association Analysis'. Yubo Jia , Guanghu Xia , Hongdan
  Fan , Qian Zhang , Xu Li . ICNDC 2012, 3rd IEEE International Conference, p. .
- [Han ()] Data Mining Concepts and Techniques "SecondEdition, J Han . 2006. Morgan Kaufmann Publisher. p.
  .
- [Agrawal and Srikant (1994)] 'Fast algorithms for mining association rules in large databases'. R Agrawal, R
  Srikant . Proceedings of the 20th International Conference on Very Large Data Bases, VLDB, (the 20th
  International Conference on Very Large Data Bases, VLDB, the 20th
  International Conference on Very Large Data Bases, VLDB, the 20th
- <sup>143</sup> [Ma ()] Improved Algorithm based on Apriori Algorithm, Qiang Ma . 2010. 23 p. . (Development and application <sup>144</sup> of computer)
- [Hu] 'The Analysis on Apriori Algorithm Based on Interest Measure'. Jingyao Hu. ICCECT 2012, IEEE
  International Conference, p. .
- [Hu Ji and Xue-Feng ()] 'The Research and Improvement of Apriori for association rules mining'. -Ming Hu Ji ,
  Xian Xue-Feng . Computer Technology and Development 2006. 16 (4) p. .
- 149 [Wang and Liu] 'The Research of Improved Association Rules Mining Apriori Algorithm'. Huiying Wang ,
- 150 Xiangwei Liu . FSKD 2011, 8Th IEEE International Conference, 2 p. .