



INTERNATIONAL JOURNAL OF ENGINEERING SCIENCES & RESEARCH TECHNOLOGY

Infrequent ITEMSET Mining Using Temporal Frequent Scheme

G.Divya*, N.Kamalraj

M.Phil Scholar, Computer Science, Dr.SNS Rajalakshmi College of Arts & Science, India
Head, Department of Computer Technology, Dr.SNS Rajalakshmi College of Arts & Science, India

Abstracts

Discovering rare data correlations is more interesting than mining frequent ones. The Proposed system focuses on the issue of discovering infrequent weighted itemsets by using weights for differentiating between relevant items and not within each transaction. To reduce the complexity of the mining process in high dimensional databases, the temporal infrequent weighted itemset TIWI Miner has been proposed. The system also performs the temporal partition of the database for timely analysis. Using the temporal scheme the user can identify the frequent and infrequent item weights for the specific time interval. Using the above the user can take decision according to the item transaction. The proposed system implements the aggregation techniques in the transactional database using FP-Growth algorithm.

Keywords:-correlation,temporal partition,mining.

Introduction

Frequent weighted itemsets indicates relationships frequently storing in data in which items may weight differently. This paper handles the problems of discovering infrequent weighted itemset mining problem. The Proposed system Mainly focus on the issues of discovering rare itemsets by using weights for differentiating between relevant items and not within each transaction. To reduce the complexity of the mining process in high dimensional databases, the temporal infrequent weighted itemset TIWI Miner adopts an FP-tree node pruning strategy to early discard items (nodes) that could never belong to any itemset satisfying the TIWI-support threshold. With the help of weighted value calculation the Proposed system system performs the infrequent itemset mining. Using the temporal scheme the user can identify the frequent and infrequent item weights for the specific time interval. Using the above the user can take decision according to the item transaction. The proposed system implements the aggregation techniques in the transactional database using FP-Growth algorithm.

- It is more efficient than the existing algorithm.
- Mining of negative association rules from rare itemsets

Previous work

In existing system the authors focus on discovering more informative association rules, i.e., the weighted association rules (WAR), which include weights denoting item significance. The main drawback of the WAR is weights are introduced only during the rule generation step after performing the traditional frequent itemset mining process. To overcome the above issue the Weighted Support and Significance Framework proposed. This framework attempted to push item weights into the item set mining process. It proposes to exploit the anti-monotonicity of the proposed weighted support constraint to drive the Apriori-based itemset mining phase. But the drawback of the technique is weights have to be pre-assigned, while, in many real-life cases, this might not be the case. To address this issue, some existing technique that analyzed transactional data set is represented and evaluated by means of a well-known indexing strategy which is named as HITS. The HITS helps to automate item weight assignment.

Problem definition

This paper addresses the discovery of infrequent and weighted itemsets from transactional weighted data sets. The problem of mining itemsets by considering weights associated with each item is known as the weighted itemset mining problem. The key issues in mining infrequent patterns are:

Advantages of the Proposed System

- Addresses the infrequent itemset mining task.
- This evaluates the frequent and infrequent association.
- Performs aggregation functions such as sum(), count(), MAX(), Min() etc.,
- Overcomes the decision making problem by applying temporal partition of database.
- Fast and accurate

- (1) Identifying interesting infrequent patterns, and
- (2) Efficiently discover infrequent patterns in large data sets.

To address this issue, probabilistic models have been constructed and integrated in Apriori-based or projection-based algorithms. The authors first addressed the issue of discovering minimal infrequent itemsets the itemsets that satisfy a maximum support threshold and do not contain any infrequent subset, from transactional data sets. An FP-Growth-like algorithm for mining minimal infrequent itemsets has also been proposed. To reduce the computational time the authors introduce the concept of residual tree.

The algorithms

FP-Growth like TIWI Algorithm, TIWI support mining algorithm and TIWI Mining Algorithms are used in proposed system.

FP-Growth like TIWI Algorithm

Input: filtered transaction from temporal scheme

Output: In Frequent item set

Description: FP-Growth: Allows frequent itemset discovery without candidate itemset generation. On like FP growth algorithm.

TIWI allows to identify the infrequent itemset discovery without candidate itemset generation

Two step approach:

Steps:

- (i) Build a compact data structure called the TIWI tree built using 2 passes over the data-set.
- (ii) Extracts infrequent itemsets directly from the TIWI-tree traversal through the Tree.

TIWI support mining algorithm

Input: Weighted transactional dataset (td)

Maximum TIWI support threshold (St)

Maximum Temporal threshold (T)

Step 1: read all transaction data (td) from the shopping cart

Step 2: get the temporal threshold (T) and segment the td.

$T(td) = \sum(I \text{ to } n) \{Trans_i(\text{split}(T))\}$

Step 3: scan transaction data Ttd and count TIWI support of each item

Step 4: count item TIWI support (td)

Step 5: create the initial FP tree from the td.

Step 6: for each transaction t_i in td
Insert t_i in Tree

Step 7: store tree, St, null in a tree (which satisfies St)

Step 8: return the output from step 7.

Read the Transaction database to get the support S of each 1-itemset, compare S with $\max_sup(St)$, and get a set of infrequent 1-itemsets, L1 (Line 1-4). If St satisfied

by every trans t_i from td, store data in the tree. Finally return the F tree (Line 7-8).

Algorithm: TIWIMining

Description: The TIWIMining algorithm takes three parameters one is the tree from the TIWI support algorithm. Another one is maximum TIWI support threshold. And finally perform prefix process.

Input: a FP tree (Htree)

Maximum TIWI support threshold (St)

The set of items with patterns (prefix)

Output: F values which is a set of TIWI extending prefix

Steps:

1. Initially assign 0 for F.
 - A. $F=0$
2. For each item I in the header treetable table HTree
 - I=prefix U{i}-generate a new itemset I by joining prefix and I with TIWI support set to the TIWI support item i
3. If I is infrequent
 - A. Store I.
4. End if
5. If $TIWI\text{-support}(I) \leq St$ then
 - A. $F=FU\{I\}$
6. End if
7. $Conditional_pattern(P)=generate(Htree, ,I)$
8. $HTree_I=createFPtree(Conditional_pattern)$
9. Perform pruning
 - A. Prune=identify(Htree I, St)
 - B. Htree=remove(HtreeI, prune)
10. If $HTree_I \neq 0$ then
 - A. $F=F \cup TIWIMining(HTree, St, I)$
11. End
12. Return the output F.

Generate a new itemset by joining prefix and I with TIWI support set to the TIWI support item i (Line 2). If I is infrequent store in the tree. If the support is less than or equal to threshold perform union (Line 5). Perform pruning. If $HTree_I \neq 0$ update Tree F.

Implementation and experiments

System Development

Infrequent Itemset Mining Using Association Rules

Many researchers are taking place in calculating In-FREQUENT item set in web mining using association rules. For example, non In-FREQUENT rules falsely calculated and spurious rules falsely generated may be produced in the In-FREQUENT pattern mining process. The main challenge of In-FREQUENT pattern mining is how to select the items and transactions to find rare frequency. The TIWI algorithm is implemented to process the In-FREQUENT pattern mining mechanism initially. Then FP growth like algorithm will be applied.

trans_id	trans_items	date_of_trans
2	A	8/25/2014 10:59:47 PM
3	A,B	8/25/2014 10:59:52 PM
4	A,B,C	8/25/2014 10:59:55 PM
5	B,C	8/25/2014 10:59:58 PM
6	A,B	8/25/2014 11:00:03 PM
7	B,C	8/25/2014 11:00:07 PM
8	C,D	8/25/2014 11:00:11 PM
9	C,D,E	8/25/2014 11:00:14 PM
10	D	8/25/2014 11:00:20 PM

FigurTransactional details

Properties for Infrequent Pattern Mining

Property 1

Let \sum_{XUY} be the set of all transactions containing XUY. To perform $X \rightarrow Y$ by removing items in XUY from the transactions in \sum_{XUY} , the maximal number of transactions that should be modified, called the minus support count, is computed as

$$MSC_{X \rightarrow Y} = C_{XUY} - [|D| \times MST] + 1 \quad (1)$$

Proof - Removing items in XUY from the transactions in \sum_{XUY} will decrease Sup_{XUY} . Let θ be the number of modified transactions when $X \rightarrow Y$ is In-FREQUENT. $(C_{XUY} - \theta) / |D| < MST \Rightarrow C_{XUY} - |D| \times MST < \theta$

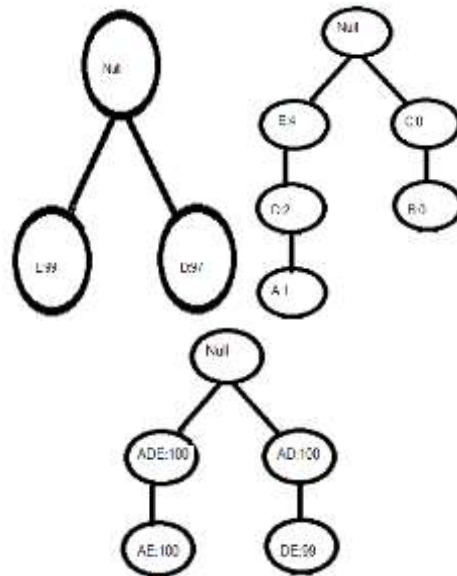


Figure: Tree generation for infrequent items

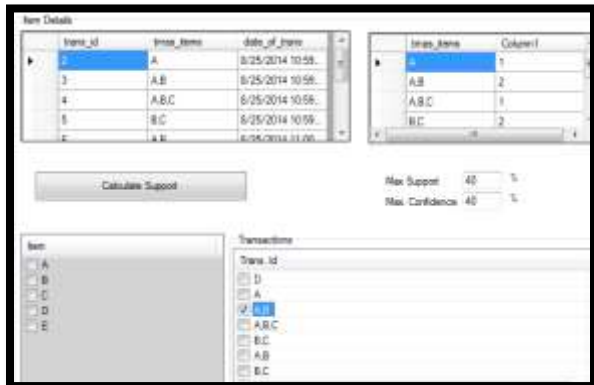


Figure:Support Calculation

Item	Support
A	96
D	97
E	99
DE	99
AE	100
AD	100
ADF	100

Figure:Support of items

Performance evaluation

The system performs TIWI mining process for identifying the rare item set using FP growth like items. The proposed system scans the database only twice. This helps to reduce the communication cost than existing system.



Figure: Execution Time Analysis

The Proposed System takes the few execution time. The scalability of this approach is first evaluated in terms of the database size, the number of infrequent items, and the number of strong rules, respectively. After that, the infrequent items are selected in such a way that all of them have at least one item in common to evaluate the effectiveness of this approach on the four measures. Finally, the overlapping degree of a rule is defined and experiments were made to observe how the correlation among rules influences the performance.



Figure: Comparison between TIWI, MIWI Miner and MINIT in terms of execution time.

Scalability approach

The processing time reported includes the CPU time efficiency consumed in the preprocessing steps (after infrequent items have been selected), the template generation, and the complete process for calculating infrequent items. The I/O time spent on the index construction and the database modification is excluded in order to highlight the impact of the database scale on this indexing mechanisms and the proposed method for infrequent pattern mining.

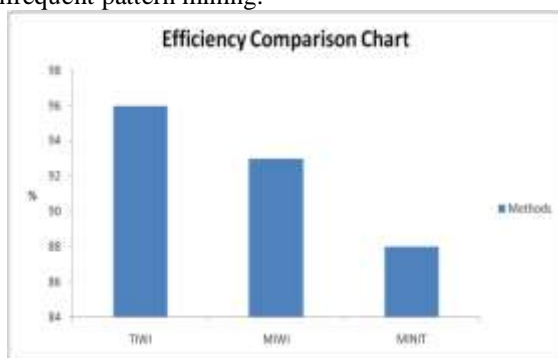


Figure: Efficiency Comparison chart

The results of the CPU time under varied database sizes are plotted in the above Figure. Each of the variations is scalable in terms of the database size. In this implementation, the table keeping the template information can be fully loaded into the main memory. The indexing mechanisms and the prime-

number representation are the major reasons for the good scalability of our approach. The former supports fast data access with hashing techniques. Moreover, the infrequent items are mapped to the prime numbers in a reverse order of their frequencies. In this way, the product of prime numbers for representing a infrequent itemset will not be too large.

Moreover, the infrequent items are mapped to the prime numbers in a reverse order of their frequencies. In this way, the product of prime numbers for representing a infrequent itemset will not be too large.

To observe the impact of the number of noninfrequent items, the results on the CPU time of this approach under varied MIWI are plotted in above. Since more non-infrequent items lead to more costs on checking the constraint IWI, the two variations MINIT and MIWI and Both perform worse than the others.

Conclusion & future scope

This work proposes a new method named as TIWI, which is mainly focused to classify all the valid modifications such that every class of modifications is related with the infrequent items, infrequent weighted items, and spurious rules that can be affected after the modifications. This method proposes some innovation methods. With the methods proposed in this work, the transactions can be modified in an order so that both the numbers of infrequent items and modified entries are considered. The experimental results show that this exposition approach is scalable in terms of database size. Moreover, this approach and the efforts taken to the avoidance of undesired side effects in infrequent pattern mining is effective in two well-designed experiments. In most cases, all the infrequent items are hidden without false rules generated. In addition, it is observed and realistic that the common items and the overlapping degrees among infrequent items have a great impact on the performance of infrequent pattern mining.

This study preserves interest to discover the full set of rules that will be falsely hidden or generated as the side effects after infrequent pattern mining. This research study emphasis efficient mechanisms are required to speed up the infrequent pattern mining process for large databases. Another issue is the fast recognition of infrequent items that cannot be hidden according to the user-specified constraint. An ideal solution or goal is to build a system that can aid the database administrator to find the infrequent items for calculating. The other issue is to remove the threshold assumption. An infrequent pattern mining approach should be robust no matter how the adversary looks into the modified database, using a

higher threshold to reveal the hidden infrequent items. The challenge is to take into account both the above attacks and the undesired side effects.

Reference

1. R. Agrawal, T. Imielinski, and Swami, "Mining Association Rules between Sets of Items in Large Databases," Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '93), pp. 207-216, 1993.
2. M.L. Antonie, O.R. Zaiane, and A. Coman, "Application of Data Mining Techniques for Medical Image Classification," Proc. Second Intl. Workshop Multimedia Data Mining in Conjunction with seventh ACM SIGKDD (MDM/KDD '01), 2001.
3. G. Cong, A.K.H. Tung, X. Xu, F. Pan, and J. Yang, "Farmer: Finding Interesting Rule Groups in Microarray Datasets," Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '04), 2004.
4. W. Wang, J. Yang, and P.S. Yu, "Efficient Mining of Weighted Association Rules (WAR)," Proc. Sixth ACM SIGKDD Int'l Conf. Knowledge Discovery and data Mining (KDD '00), pp. 270-274, 2000.
5. F. Tao, F. Murtagh, and M. Farid, "Weighted Association Rule Mining Using Weighted Support and Significance Framework," Proc. ninth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (KDD '03), pp. 661-666, 2003.
6. K. Sun and F. Bai, "Mining Weighted Association Rules Without Preassigned Weights," IEEE Trans. Knowledge and Data Eng., vol. 20, no. 4, pp. 489-495, Apr. 2008.
7. A. Manning and D. Haglin, "A New Algorithm for Finding Minimal Sample Uniques for Use in Statistical Disclosure Assessment," Proc. IEEE Fifth Int'l Conf. Data Mining (ICDM '05), pp. 290-297, 2005.
8. A.M. Manning, D.J. Haglin, and J.A. Keane, "A Recursive Search Algorithm for Statistical Disclosure Assessment," Data Mining and Knowledge Discovery, vol. 16, no. 2, pp. 165-196, <http://mavdisk.mnsu.edu/haglin>, 2008.
9. D.J. Haglin and A.M. Manning, "On Minimal Infrequent Itemset Mining," Proc. Int'l Conf. Data Mining (DMIN '07), pp. 141-147, 2007.
10. J. Han, J. Pei, and Y. Yin, "Mining Frequent Patterns without Candidate Generation," Proc. ACM SIGMOD Int'l Conf. Management of Data, pp. 1-12, 2000.