

A Naïve Approach to High Utility Rare Itemset Mining Algorithm using Temporal Concept – THURI

Jyothi Pillai¹, O.P. Vyas²

Associate Professor, Bhilai Institute of Technology, Durg, Chhattisgarh, India¹

Professor, Indian Institute of Information Technology, Allahabad, Uttar Pradesh, India²

Abstract: Business strategies use information about any organizations past performance that can be used to predict its future performance. The right business strategy can be formulated by clearly understanding the dynamically changing business environment. Temporal data mining can be used for obtaining temporal measures of various operations of company. Temporal mining can be instrumental for tracking the temporal changes in the business activities which allows the company for gaining insights to process improvement and optimization. An appealing question in temporal mining that is concerned with business strategy is to study the contribution of past decisions on the success of the organization. Temporal Data Mining is an important step in the Knowledge Discovery process that discovers temporal patterns or temporal rules from Temporal Databases. Temporal Data Mining Algorithms are the algorithms which consider temporal patterns from temporal data or fit models to temporal databases. Traditional Frequent itemset mining discovers frequent itemsets from transactional databases using only items occurrence frequency and not considering items utility. But in many real world situations, utility of itemsets based upon user's perspective such as cost, profit or revenue is of significant importance.

One of the latest data mining research areas is Utility Mining which emphasis on all types of utility factors and incorporates utility concepts in data mining tasks. The utility-based descriptive data mining which aims at discovering itemsets having high total utility is termed as High utility itemset mining. High Utility itemsets may contain frequent as well as rare itemsets. Temporal data mining is a very fast expanding field with many new research results reported and many new temporal data mining analysis methods or prototypes developed recently. The temporal significant rare utility itemsets are those itemsets which appear infrequently in the current time window of large databases but are highly associated with specific data. In this paper, a novel method is proposed, namely THURI (High Utility Rare Itemset Mining Algorithm using Temporal Concept), for efficiently and effectively mine high utility rare itemsets from databases with temporal consideration of utility values. The novel contribution of THURI is that it can effectively extract high utility rare itemsets from temporal transaction databases.

Keywords: Frequent Itemset Mining, Rare Itemset Mining, Utility Mining, High utility Rare Itemset Mining, Temporal Mining

I. INTRODUCTION

The business strategies can be enhanced by keeping track of all business activities and updating them accordingly with time. Temporal data mining has recently received increasing attention, as many processes in business and science have interesting time changing aspects. Applications include the search for patterns in large time series databases and streaming time series.

Temporal data mining can be defined as the activity of looking for interesting correlations or patterns in large sets of temporal data accumulated for other purposes. In many temporal data mining scenarios, there is a need to incorporate timing information more explicitly into the patterns. This would give the patterns (and the rules generated from them) greater descriptive and inferential power. All techniques mentioned above treat events in the sequence as instantaneous. However, in many applications different events persist for different amounts of time and the durations of events carry important information.

The field of temporal data mining is relatively young and one expects to see many new developments in the near

future. In all data mining applications, the primary constraint is the large volume of data.

Temporal Data Mining is a rapidly evolving area of research that is at the intersection of several disciplines, including statistics, temporal pattern recognition, temporal databases, optimization, visualization, high-performance computing, and parallel computing. Basically temporal data mining is concerned with the analysis of temporal data and for finding temporal patterns and regularities in sets of temporal data. Also temporal data mining techniques allow for the possibility of computer driven, automatic exploration of the data.

Temporal association rules mining doesn't consider the utility of every item. Temporal utility mining is a research which is extended from temporal association rules mining and utility mining [10]. Utility of an itemset is considered as the value of this itemset, and utility mining aims at identifying the itemsets with high utilities [3]. The goal of utility mining is to identify high utility itemsets which drive a large portion of the total utility. The temporal high

utility itemsets are the itemsets whose support is larger than a pre-specified threshold in current time window of the data stream. An important question in retail marketing that can be addressed by temporal mining of business operations is which areas can be temporally optimized so as to increase business profitability and customer satisfaction. The most profitable products or services of the company can be found out by applying temporal data mining techniques to business data. Also the changing purchasing behavior of customers according to time can be analyzed by using temporal data mining. Accordingly, high utility values can be assigned to products which are more preferred by customers in a particular time period.

In many real-life applications, high-utility itemsets consist of rare items. Rare itemsets provide useful information in different decision-making domains; customers purchase microwave ovens or LEDs rarely as compared to bread, butter, milk, etc. The former may yield more profit for the supermarket than the latter [7]. Jyothi et al proposed High Utility Rare Itemset Mining HURI algorithm [7], for generating high utility rare itemsets of users' interest. HURI is a two-phase algorithm.

In this paper, A Naïve Approach to HURI using temporal concept (THURI) is proposed. The novel contribution of THURI is that it can effectively identify the temporal high utility rare itemsets from temporal databases.

The rest of paper is organized as follows. In section 2, some related works are discussed: section 3 presents the THURI algorithm and section 4 presents conclusion and future work.

II. RELATED WORK

As defined by Weiqiang Lin et al, Temporal Data Mining is a single step in the process of Knowledge Discovery in Temporal Databases that enumerates structures (temporal patterns or models) over the temporal data, and any algorithm that enumerates [WMG2002] Temporal data mining analyzes temporal data for finding temporal patterns and regularities in temporal data sets. Temporal data mining tasks include characterization and comparison, clustering analysis, classification, association rules mining, pattern analysis, prediction and trend analysis of temporal data. Weiqiang Lin et al have proposed "The Additive Distributional Recursion Algorithm (ADRA)" in General Hidden Distribution-based Analysis Theory to build temporal data models for discovering temporal patterns. The authors have also proposed a framework of Temporal Clustering method and Distribution-based Temporal Clustering Algorithm to discover temporal patterns. A framework of Temporal Classification is also proposed by using Temporal Clustering method. A framework of Temporal Feature Selection is put forward for discovering temporal patterns. In Temporal data mining, useful information is harvested from temporal data [15]. Today the importance of temporal information has increased in health care and business organizations. Temporal Data Mining has wide applications in various fields such as classification, clustering, similarity computation, pattern discovery, and prediction containing temporal data.

Swati Soni et al proposed a portfolio management solution with business intelligence characteristics in [13]. The authors proposed a novel algorithm for temporal association mining with utility approach to find the temporal high utility itemset by generating less candidate itemsets.

M. Sulaiman Khan et al presented a novel algorithm Dual Support Apriori for Temporal data (DSAT) which is used to discover Jumping Emerging Patterns (JEPs) from time series data [11]. The technique discovers the itemset variations over time. The authors conclude that DSAT utilizes less memory, less computational cost and minimum dataset scans.

Temporal time-stamped data is collected in large amounts for Stock exchanges, regarding different business transactions like quotations, trading, payment, delivery etc. On-line surveillance systems are not capable of detecting or preventing mal-practices hidden deep in databases because of short-term data analysis within time constraints. Girish Keshav Palshikar et al proposed a pattern recognition method for scrutiny of the trading databases for detecting or preventing mal-practices by extracting temporal pattern and then analysis of patterns by investigation experts [5]. A fuzzy temporal logic notation is described for specifying such patterns. A tool called SNIFFER is used by the authors for detecting the impact of surveillance pattern on the given temporal databases.

Temporal data mining is concerned with analyzing large volumes of ordered data streams having temporal interdependencies for automatically discovering interesting patterns, trends or relationships. Srivatsan Laxman et al presented a survey report of different algorithms for discovering patterns in sequential data streams [12]. An overview of temporal data mining methods for analyzing large volumes of sequential data streams to uncover hidden patterns is put forward by the authors.

Tarek F. Gharib et al present temporal association rules concept for handling time series by considering time expressions in association rules [14]. Due to dynamic nature of temporal databases, the discovered rules have to be updated frequently. An incremental algorithm, ITARM, is presented by the authors for maintaining and updating the temporal association rules of a transaction database. The benefit of the algorithm is that the earlier mining results are used for deriving the final mining output which reduces the time taken to generate new candidates. Experimental analysis on both synthetic and the real dataset have shown a considerable improvement over the traditional method of mining temporal large database.

The value of an itemset is called the utility of an itemset and the process of identifying high utilities itemsets is termed as utility mining [4]. Temporal high utility itemsets mining is an important process for Discovery of temporal high utility itemsets which have support greater than a pre-specified threshold in existing time window of the data stream. A novel method, named Temporal High Utility

Itemsets (THUI-Mine), is proposed by Chun-Jung Chu et al to efficiently mine temporal high utility itemsets present in data streams. The authors conclude that the proposed THUI-Mine effectively discover all temporal high utility itemsets within different time windows of data streams with less memory space and execution time.

Temporal sequence is a series of nominal symbols of particular alphabets whereas a time series is a sequence of continuous, real-valued elements. Temporal sequences are present in different domains such as engineering, medicine, finance, etc., hence one of the most crucial data mining tasks is to model and extract information from these temporal sequences. Cláudia M. Antunes et al presented an overview of methods for mining temporal sequences. Three steps are involved in finding associations between sequences of events: First, represent and model the data sequence in a suitable form; second, similarity measures between sequences are defined. Finally, the models and representations are applied to the actual mining problems.

[2]Chang-Hung Lee et al explored general temporal association rules mining problem in publication databases. A set of transactions consisting of a set of items in each transaction where each item consists of individual exhibition period is termed as a publication database. Chang-Hung Lee et al propose a novel algorithm Progressive-Partition-Miner (abbreviated as PPM) for discovering temporal association rules present in a publication database. In PPM, the publication database is first partitioned with respect to exhibition periods of items. Then the occurrence count for each candidate 2-itemset is progressively accumulated on the basis of partitioning characteristics. The authors conclude that by using progressive counting concepts and scan reduction techniques in Algorithm PPM, I/O and CPU cost are reduced. Also memory utilization is effectively controlled due to proper partitioning. Algorithm PPM was found to effectively mine large publication databases like bookstore database, video rental store database, library-book rental databases and electronic commerce transactions.

Junheng-Huang et al surveyed a new method for temporal association rules mining from large databases where items have different exhibition periods[9]. An efficient algorithm named Standing for Segmented Progressive Filter Algorithm (SPFA) was also presented for effectively discovering general temporal association rules. In SPFA, at first the database is segmented into sub-databases where item of each sub-database will have either common starting time or common ending time. Then, SPFA progressively filters candidate 2-itemsets for each sub-database using thresholds with respect to time.

As stock market is dynamic and volatile, Temporal Data mining is widely used in financial markets and stock-price forecasting. Gerasimos Marketos et al propose an Intelligent Stock Market Assistant which serves as a portfolio management solution having business intelligence characteristics for finding all possible relations between stocks [7]. The proposed tool uses a sequence mining algorithm consisting of pre-processing

and pattern evaluation steps. The technical analysis focuses on the stock chart and discovers common or correlated behavior between different stocks.

Time series data analysis has to be done on multiple instances of the same record. In traditional clustering methods, instances are exclusively classified by attaching an event to a specific cluster. Richi et al proposed a method for investigating the predictive power of the clustering technique of stock market data using temporal pattern recognition. In the proposed method, according to the supplied prediction confidence of each transaction, the prediction accuracy was considered for developing trading strategies.

III. THURI-TEMPORAL CONCEPT IN MININT HIGH UTILITY RARE ITEMSETS

An important issue extended from Association Rule Mining is the discovery of temporal association rules from temporal databases. Temporal data mining can be defined as the activity of discovering interesting correlations or patterns in large sets of temporal data. The information generated from temporal mining helps organizations to stay up to date by providing right knowledge about the current environmental changes and about the right products at right time. Also temporal high utility rare itemsets mining is an important process for mining interesting non-frequent patterns from temporal databases.

A. HURI Algorithm

Rare itemset mining is very important as rare itemsets may bring adequate profits to the business. In [7], Jyothi et al proposed High Utility Rare Itemset Mining [HURI] to find high utility rare-itemsets based on minimum threshold values and user preferences. The utility of items is decided by considering factors such as profit, sale, temporal aspects, etc. of items

B. Extraction of HURI using THURI Algorithm

THURI algorithm (Figure 1) is an extension of HURI algorithm which incorporates temporal concept.

Algorithm THURI

Description: Finding High Utility Rare Itemsets of users' interest from each time partition

Ck: Candidate itemset of size *k*

Lk: Rare itemset of size *k*

n: Number of time partitions

begin for *k* = 1 to *n*

For each transaction *t* in database

begin

increment support for each item *i* present in *t*

End

LI = {Rare 1-itemset with support less than user provided max_sup}

for(*k*= 1; *Lk*!=∅; *k*++)

begin

C k+I = candidates generated from *Lk*;

//loop to calculate total utility of each item

For each transaction *t* in database

begin

Calculate total quantity of each item *i* in *t*

Find total utility for item i using following formula:-

$u(i,t) = \text{quantity}[i] * \text{user_provided_utility for } i$
End

//loop to find rare itemsets and their utility

For each transaction t in database

begin

increment the count of all candidates in C_{k+1} that are contained in t

$L_{k+1} = \text{candidates in } C_{k+1} \text{ less than } \text{min_support}$

Add L_{k+1} to the Itemset_Utility Table in database by calculating rare itemset utility Using following formula:

Utility(R,t) = \sum for each individual item i in R ($u(i,t)$);
End

//loop to find high utility rare itemset

For each itemset $iset$ in rare itemset Table R

begin

If (Utility($iset$) >

user_provided_threshold_for_high_utility_rare_itemset)

then $iset$ is a rare_itemset that is of user interest
i.e.high_utility_rare_itemset

else $iset$ is a rare itemset but is not of user interest

End

Return high_utility_rare_itemsets

End

END

Fig.1. Pseudo Code for THURI

C. Performance Evaluation of THURI

The temporal aspect is incorporated by dividing the data set into four quarters and then mining has been performed accordingly. In item utility table (Table 2), each item is assigned an external utility and internal utility is calculated from database D .

In THURI Algorithm, high utility rare itemsets are generated in three phases:-

- In first phase, different values of Utilities are assigned to rare itemsets in different time periods (monthly, bimonthly or quarterly).

In Table 2, EUQ1, EUQ2, EUQ3, EUQ4 are external utilities of different items in 4 different quarters. Similarly IUQ1, IUQ2, IUQ3, IUQ4 are internal utilities of items in 4 different quarters.

- In second phase, rare itemsets are generated from temporal database by considering those itemsets which have support value less than the maximum support threshold

For example, on application of THURI (Figure 1) on Temporal Transactional dataset described in Table 1 and by setting the value of maximum support threshold to 40%, the rare itemsets generated in Quarter 3 are listed in Table 3.

- In third phase, by inputting the utility threshold value according to users' interest, rare itemsets having utility value greater than the minimum utility threshold are generated from different time periods.

For example, by setting high utility threshold as 60, the high utility rare itemsets generated in Quarter 3 are listed in Table 4.

IV. CONCLUSION

For achieving the preset mission and vision of the business new business strategies have to be developed frequently. To formulate the right business strategy, the key part is to understand the dynamic nature of the environment in which the business operates. Temporal data mining can be instrumental in tracking the changes in the business environment over time and in enhancing the quality of business strategies. In this paper, a novel method is proposed, namely THURI (High Utility Rare Itemset Mining Algorithm using Temporal Concept), for efficiently and effectively mine high utility rare itemsets from databases with temporal consideration of utility values. The novel contribution of THURI is that it can effectively extract high utility rare itemsets from different quarters. The temporal aspect is incorporated by dividing the data set into quarters, months or seasonal time windows and then mining has been performed accordingly.

In this paper, only calendric (quarterly, seasonal, bimonthly, monthly time windows) temporal type have been considered. More time windows such as cyclic or trend analysis, will be considered for the temporal concept in the future.

TABLE I. TEMPORAL TRANSACTION DATASET D

Date	Quarter	Tran- ID	A001	B002	C003	D004	E005	F006	G007	H008	I009	J010	K011	L012	M013	N014
02-Jan-13	1	A	1	2	2	1	1	5	1	0	2	6	1	0	0	0
23-Jan-13	1	A1	0	1	0	2	1	1	2	1	1	1	1	1	0	0
14-Feb-13	1	B	2	0	1	1	1	0	3	1	5	3	3	1	6	0
24-Oct-13	4	B2	6	3	1	1	2	6	6	4	2	4	3	1	2	2
04-Apr-13	2	C	1	0	0	1	1	0	4	2	0	3	1	0	3	4
22-May-13	2	C3	0	0	2	3	6	0	1	2	1	0	2	0	2	3
11-Jun-13	2	D	1	0	2	0	0	0	0	2	5	0	3	0	3	2
22-Feb-13	1	D4	1	7	1	0	4	0	0	0	2	0	0	2	0	8
16-Jul-13	3	E	2	0	3	2	1	0	1	6	1	4	1	1	3	1
17-Mar-13	1	E5	3	7	2	2	9	0	0	7	5	4	0	3	3	7
23-Aug-13	3	F	3	1	0	5	1	0	1	0	1	0	5	0	1	1
18-Jun-13	2	F6	0	3	0	3	0	0	2	0	0	5	1	0	2	1
17-Sep-13	3	G	5	2	6	6	0	0	0	1	0	0	5	3	2	1
21-Apr-13	2	G7	0	0	0	6	1	1	3	0	5	6	2	0	1	2
21-Sep-13	3	H	0	7	0	0	0	0	0	6	1	0	0	0	2	2
25-Jul-13	3	H8	1	8	5	8	0	4	0	1	1	0	6	1	1	2
03-Oct-13	4	I	1	0	1	1	1	0	4	0	2	0	1	0	4	4
06-May-13	2	I9	2	0	3	0	0	3	1	2	2	0	0	2	6	1
14-Oct-13	4	J	0	0	1	0	2	4	0	2	0	0	0	1	4	1
19-Sep-13	3	J10	4	0	6	8	0	8	0	3	3	0	0	2	6	2

TABLE II. ITEM UTILITY TABLE

Items	EUQ1	EUQ2	EUQ3	EUQ4	IUQ1	IUQ2	IUQ3	IUQ4	TUQ1	TUQ2	TUQ3	TUQ4
A001	3	2	1	1	7	4	15	7	21	8	15	7
B002	1	7	2	9	17	3	18	3	17	21	36	27
C003	2	1	4	2	6	7	20	3	12	7	80	6
D004	2	4	3	3	6	13	29	2	12	52	87	6
E005	7	3	5	5	16	8	2	5	112	24	2	25
F006	5	5	5	1	6	4	12	10	30	20	60	10
G007	6	6	3	3	6	11	2	10	36	66	6	30
H008	1	8	2	4	9	8	17	6	9	64	34	24
I009	8	3	1	2	15	13	7	4	120	39	7	8
J010	4	4	4	8	14	14	4	4	56	56	16	32
K011	3	5	1	5	5	9	17	4	15	45	17	20
L012	11	1	9	5	7	2	7	2	77	2	7	10
M013	4	3	2	5	9	17	15	10	36	51	30	50
N014	8	8	2	2	15	13	9	7	120	104	18	14

TABLE III. LIST OF RARE ITEMSETS GENERATED OF QUARTER 3

List_of_rareItemSet	Total_Utility
[E005]	2
[F006]	60
[E005, F006]	62
[G007]	6
[E005, G007]	8
[F006, G007]	66
[E005, F006, G007]	68
[J010]	16
[E005, J010]	18
[F006, J010]	76
[E005, F006, J010]	78
[G007, J010]	22
[E005, G007, J010]	24
[F006, G007, J010]	82
[E005, F006, G007, J010]	84

TABLE 4. LIST OF HIGH UTILITY RARE ITEMSETS OF QUARTER

List_of_high utility_rareItemSet	Total_Utility
[F006]	60
[E005, F006]	62
[F006, G007]	66
[E005, F006, G007]	68
[F006, J010]	76
[E005, F006, J010]	78
[F006, G007, J010]	82
[E005, F006, G007, J010]	84

REFERENCES

[1] Cláudia M. Antunes, Arlindo L. Oliveira, Temporal Data Mining: an overview, Lecture Notes in Computer Science, pp 1-15.

[2] Lee, C., Chen, M., Lin, C., Progressive Partition Miner: An Efficient Algorithm for Mining General Temporal Association Rules, IEEE Transactions On Knowledge And Data Engineering, 15(4), 1004-1017, 2003.

[3] Chun-Jung Chu , Vincent S. Tseng , Tyne Liang, An efficient algorithm for mining temporal high utility itemsets from data streams, The Journal of Systems and Software, Elsevier, 1105–1117, 2008.

[4] Chu, C., Tseng, V. S. & Liang, T., An efficient algorithm for mining temporal high utility itemsets from data streams, The Journal of Systems and Software 81, Elsevier, 1105–1117, 2008.

[5] Girish Keshav Palshikar, Arun Bahulkar, Fuzzy Temporal Patterns for Analyzing Stock Market Databases, Advances In Data Management 2000, Tata McGraw-Hill Publishing Company Ltd., ©CSI 2000

[6] Gerasimos Marketos, Konstantinos Padiaditakis, Yannis Theodoridis, BabisTheodoulidis, Intelligent Stock Market Assistant using Temporal Data Mining, pp 1-11, 2004.

[7] Jyothi Pillai, O.P. Vyas, High Utility Rare Item Set Mining (HURI): An Approach for Extracting High Utility Rare Item Sets, Journal on Future Engineering and Technology, Volume 7 (1), i-manager Publications – Oct 1, 2011.

[8] Jyothi Pillai, O. P. Vyas, Maybin Mueyba, HURI – A Novel Algorithm for Mining High Utility Rare Itemsets, Advances in Computing and Information Technology, Volume: 177, 2012, pp 531-540, @ Springer-Verlag Berlin Heidelberg, springer-link.com.

[9] Huang, J., Wei, W., Efficient Algorithm for Mining Temporal Association Rule, IJCSNS International Journal of Computer Science and Network Security, 7(4), 268-271, 2007.

[10] Rudresh Shah, Ravindra Gupta, Surendra Singh, Efficient Algorithm for High Utility Pattern Mining in Incremental Databases, International Journal of Electronics Communication and Computer Engineering, IJECCE, Volume 3, Issue 2, ISSN 2249 – 071X, 25-28, 2012.

[11] M. Sulaiman Khan, Frans Coenen, D. Reid, R. Patel, L. Archer, A sliding windows based dual support framework for discovering emerging trends from temporal data, Knowledge Based Systems, Volume 23, Number 4, May 2010, pp 316-322.

[12] Srivatsan Laxman, P S Sastry, A survey of temporal data mining, SADHANA, Academy Proceedings in Engineering Sciences, Vol. 31, Part 2, April 2006, pp 173–198.

[13] Swati Soni, Sini Shibu, Advance Mining of Temporal High Utility Itemset, International Journal of Information Technology and Computer Science(IJITCS),ISSN: 2074-9007 (Print), ISSN: 2074-9015 (Online),Publisher: MECS,IJITCS Vol. 4, No. 4, April 2012, pp 26-32.

[14] Tarek F. Gharib, Hamed Nassar, Mohamed Taha, Ajith Abraham, An efficient algorithm for incremental mining of temporal

association rules, Data & Knowledge Engineering, 69, 2010
Elsevier, pp 800–815.

- [15] Theophano Mitsa, Temporal Data Mining, CRC Press, A Chapman and Hall Book, 2009.
- [16] Lin, W., Orgun, M. A. & Williams, G. J., An Overview of Temporal Data Mining, Proceedings of the Australasian Data Mining Workshop, ADM02, Sydney, Australia, 83-90, 2002.

BIOGRAPHIES



Mrs. Jyothi Pillai is Associate Professor in Department of Computer Applications at Bhilai Institute of Technology, Durg (C.G.), India. She is a post-graduate from Barkatullah University, Bhopal, India. She is a Life member of Indian Society for Technical Education. She has a total teaching experience of 19 years. She has a total of 21 Research papers published in National and International Journals / Conferences into her credit. Presently, she is pursuing Ph.D. from Pt. Ravi Shankar Shukla University, Raipur under the guidance of Dr. O.P.Vyas, IIIT, Allahabad. Her current research interests are Business Intelligence and Soft Computing.



Dr.O.P.Vyas is currently working as Professor and Incharge Officer (Doctoral Research Section) in Indian Institute of Information Technology-Allahabad (Govt. of India's Center of Excellence in I.T.).Dr.Vyas has done B.Tech.(Computer Science) from IIT Kharagpur and has done Ph.D. work in joint collaboration with Technical University of Kaiserslautern (Germany) and I.I.T. Kharagpur. With more than 25years of academic experience Dr.Vyas has guided Four Scholars for the successful award of Ph.D. degree and has more than 80 research publications with two books to his credit. His current research interests are Linked Data Mining and Service Oriented Architectures.