

A COMPREHENSIVE STUDY ON HIGH UTILITY ITEM SET MINING WITH ADVANCE MINING ALGORITHMS

^{1*} V CHANDRA PRAKASH, CH PAVANI²

¹² ASSISTANT PROFESSOR, DEPT OF CSE , SRI INDU COLLEGE OF ENGINEERING & TECHNOLOGY(AUTONOMOUS), SHERIGUDA IBRAHIMPATNAM, RANGAREDDY, TELANGANA ,INDIA

Abstract— mining high utility itemset from exchange database alludes to the revelation of itemset with a high utility like benefits. Itemset Utility Mining is an expansion of Frequent Itemset mining, which finds itemsets that happen frequently. In High Utility Itemset Mining, the objective is to perceive itemsets which have utility qualities over a given utility edge. The high utility itemset is the itemset with an utility at the very least a client indicated least help edge esteem; else that itemset is treated as a low utility itemset. In this paper, we present a writing study of the current situation with research and the different algorithms and its restrictions for high utility itemset mining.

Keywords—Association rules mining, frequent itemset mining, high utility itemset mining.

I. INTRODUCTION

Data mining is the way toward uncovering non - trifling, already obscure and conceivably valuable data from huge databases. Association Rule Mining (ARM) is an imperative data mining method that is utilized to find the examples/rules among items in a substantial database [1]. The objective of ARM is to recognize a gathering of items which happen together, for instance in a market bin examination. Mining association rules can be deteriorated into two stages: the first is creating frequent itemsets. The second is producing association rules. The fundamental test in association rule mining is to recognize frequent itemsets. Finding frequent itemset is one of the imperative strides in association rule mining. Since the arrangement of the second sub - the issue is clear, the greater part of the specialists had concentrated on the most proficient method to create frequent itemsets.

Frequent itemsets are the itemsets which happen frequently in the exchange database. The target of Frequent Itemset Mining is to distinguish all the frequent itemsets in an exchange database. Additionally, items having high and low moving frequencies may have low and high-benefit esteems, separately. For example, some frequently sold items, for example, bread, drain, and pen may have bring down benefit esteems when contrasted with that of infrequently sold higher benefit esteem items, for example, gold, platinum and precious stone. Consequently, finding just customary frequent examples in a database can't satisfy the necessity of finding the most important itemsets/clients that add to the real piece of the aggregate benefits in reality retail database. This gives the inspiration to build up a mining model to find the itemsets/clients adding to most of the benefit. A frequent itemset is an itemset having recurrence bolster more prominent a base client determined edge. [1]

II. HIGH UTILITY ITEMS ET MINING

1.1 Data Mining

Data mining indicates the movement of separating new, profitable and nontrivial data from substantial volumes of data. Most ordinarily, the point is to discover examples or assemble models utilizing explicit algorithms from different logical controls including man-made consciousness, machine learning, and database frameworks. The data mining undertakings can be characterized into two classes:

- Predictive data mining where the objective is to assemble an executable model from data which can be utilized for arrangement, forecast or estimation.
- Descriptive data mining where the objective is to find fascinating examples and associations with regards to data.

1.2 Frequent Pattern Mining

Frequent itemsets are the items that show up frequently in the exchanges. The primary objective of frequent itemset mining is to distinguish all the itemsets in the exchange data set, which are frequently acquired. Item sets are characterized as a non-void set of items. In the event that itemset is with k -distinctive items is named as a k -itemset. For ex {bread, margarine, milk} may signify as a 3-itemset in a grocery store transaction [1].

Let $I = \{i\}$ be a set of items and D be an exchange database $\{ \}$ where every exchange $T \in D$ is a subset of I . The help or recurrence of an example $X \{ \}$ is the quantity of exchange contained the example in the value-based database.

The Apriori [1], the calculation is the underlying answer for the frequent example mining issue. To beats the issues of Apriori, which produces more competitor sets and require more sweeps of database FP-Growth has been proposed [2], Uses FP-Tree data structure with no hopeful age and utilizing just two database examines. In the structure of frequent itemsets mining, the significance of an item are not considered.

Impediments of Frequent Pattern Mining :

- i) The buy amounts are not considered. Accordingly, an item may just show up once or zero time in an exchange.
- ii) All items are seen as having a similar significance, the utility of weight. For instance, if a client purchases an over the top expensive jug of wine or only a bit of bread, it is seen as being similarly critical.

In this way, frequent example mining may discover many frequent examples that are not fascinating. For instance, one may find that {bread, milk} is a frequent example. Be that as it may, from a business point of view, this example might be uninteresting on the grounds that it doesn't create much benefit. In addition, frequent example mining algorithms may miss the uncommon examples that create a high benefit, for example, maybe {caviar, wine}

1.3 Utility Mining

Utility mining is a standout amongst the most difficult data mining assignments is the mining of high utility itemsets productively. Distinguishing proof of the itemsets with high utilities is called Utility Mining. The utility can be estimated according to the client inclinations utility can be estimated as far as cost, benefit or different articulations. The confinements of frequent or uncommon itemset mining spurred specialists to imagine an utility based mining approach, which enables a client to helpfully express his or her viewpoints concerning the value of itemsets as utility qualities and after that find itemsets with high utility qualities higher than a limit. In utility based mining the term utility alludes to the quantitative portrayal of client inclination i.e. as indicated by an itemsets utility esteem is the estimation of the significance of that itemset in the client's point of view.

The high utility itemset mining issue is to discover all itemsets that have utility bigger than a client determined estimation of least utility. The esteem or benefit Associated with each item in a database is known as the utility of that itemset. For instance, the PC framework is progressively beneficial them phone as far as profit.[2]

Utility characterizes as Interestingness, gainfulness or significance of the item. Utility estimated regarding cost benefit or other client inclination.

The utility of items in the exchange database includes the accompanying two viewpoints:

- (1) The significance of unmistakable items, called outer utility (e), i.e. unit benefit and
- (2) The significance of items in exchanges, called inside utility (I), i.e. amount

Utility of Itemset (U) = outside utility (e) * inner utility (I).

In numerous applications like a cross-advertising in retail locations, online internet business the board, site clickstream investigation and finding the vital example in biomedical applications High utility mining is generally utilized. [2]

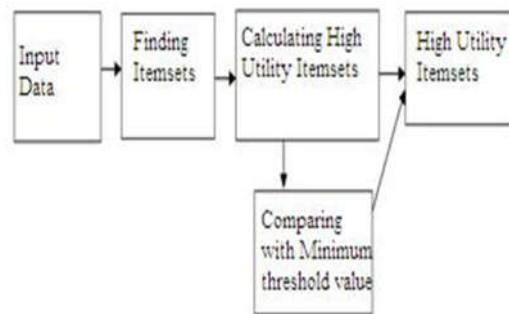


Fig. 1 Flow diagram of HUIM

II. III.LITERATURE SURVEY

Frequent itemsets assume a fundamental job in numerous data mining undertakings that attempt to discover fascinating examples from databases, for example, association rules, connections, groupings, scenes, classifiers, bunches and some more, of which the mining of association rules is a standout amongst the most famous issues. The first inspiration for seeking association rules emerged from the need to break down alleged market exchange data, that is, to inspect client conduct regarding the obtained items. Association rules depict how regularly items are obtained together. For instance, an association rule "bread \Rightarrow spread (80%)" states that 80% of the clients that purchased bread additionally purchased margarine. Such rules can be valuable for choices concerning item evaluating, advancements, store design, and numerous others.

Numerous specialists imagined thoughts to create frequent itemsets. The time required for creating frequent itemsets assumes a vital job.

Investigations of frequent itemset (or example) mining are recognized in the data mining field due to its expansive applications in mining association rules, relationships, and diagram design limitation dependent on frequent examples, consecutive examples, and numerous other data mining errands. Productive algorithms for mining frequent itemsets are significant for mining association rules and in addition for some, other data mining undertakings.

The real test found in frequent example mining is countless examples. As the base edge moves toward becoming lower, an exponentially substantial number of itemsets are created. In this manner, pruning irrelevant examples should be possible successfully in the mining procedure and that winds up one of the primary subjects in frequent example mining. Thusly, the principle point is to improve the way toward discovering designs which ought to be effective, versatile and can recognize the vital examples which can be utilized in different ways. A couple of frequent itemset mining algorithms [18] are thought about

Since its presentation in 1993 [20], the assignment of association rule mining has gotten a lot of consideration. Today the mining of such rules is as yet a standout amongst the most mainstream design revelation techniques in Knowledge disclosure strategies.

Association rule mining, a standout amongst the most critical and very much inquired about procedures of data mining, was first presented in [20]. It intends to remove fascinating relationships, frequent examples, associations or causal structures among sets of items in the exchange databases or other data stores [24].

Association rules are generally utilized in different regions, for example, media transmission systems, market and hazard the executives, stock control and so forth. The connection among information and data mining and learning disclosure in database process are depicted in [19].

Examination of the hypothesis of information revelation in the database and the fact of the mechanical strategy is given in [25].

Apriori [13], as an established calculation of association rules mining, embraces an iterative strategy to find frequent itemsets. It comprises of two stages: the joining step and the prune step. In the joining step, a hopeful k-itemsets is created by joining two frequent (k-1) itemsets. At that point in the prune step, all itemsets whose (k-1) subset is anything but a frequent k-itemsets, are expelled from the applicant k-itemsets. At that point the database is checked to register the help of the competitor k-itemsets. This procedure is rehashed until no new competitor k-itemsets is produced. An enormous number of outputs of the database and a

DIC Algorithm [5] is contrasted and the Apriori calculation. It demonstrates that the DIC calculation is superior to Apriori Algorithm. DIC utilizes parcel way to deal with find frequent itemsets.

Hine [12] can be scaled up to substantial databases by database dividing and when the dataset winds up thick, (restrictive) FP-trees can be developed powerfully as a major aspect of the mining procedure. H-mine, which exploits another data structure called H-struct and powerfully changes connects in the mining procedure. H-mine re-modifies the connections at mining diverse "anticipated" databases and has space overhead. H-mine ingests highlights of FP-development. It is basically a frequent-design development approach since it segments its pursuit space as per the two examples to be hunt down and the data set to be looked on, by a separation and-vanquish procedure, without creating and testing applicant designs. Time and number of outputs are high amid the calculation.

Segment system is utilized in dealing with extensive datasets [2] to discover frequent itemsets. One of the key highlights of this calculation is that it requires various ignores the database. For circle occupant databases, this requires perusing the database totally for each pass bringing about an extensive number of plate I/O.

An enhanced Apriori calculation [23] receives a strategy to lessen the repetitive age of sub-itemsets amid pruning the competitor itemsets. This can frame specifically the set of frequent itemsets and dispense with competitors having a subset that isn't frequent meanwhile. This calculation can't decrease the quantity of database examining and the excess of the rules.

Frequent item diagram calculation [22] is a change of FP-tree. The commitment of this methodology is to tally the frequent 2-itemsets and to shape a graphical structure which separates all conceivable frequent itemsets in the database. This graphical structure requires finish filter all frequent-2 itemsets, frequent-3 itemsets, ... frequent n-itemsets. It needs various database sweeps to discover frequent itemsets.

FP-tree based Approximate Frequent Itemsets mining calculation [11], FP-AFI is produced to find surmised frequent itemsets from a FP-tree like structure. It utilizes a recursive capacity for getting the set of exchanges. The continued checking of the database is required with the goal that execution time is high.

Another Mining Algorithm dependent on frequent itemsets [26] was presented. This new calculation can get the new frequent itemsets by examining the undirected itemset diagram. It needs re-output of the undirected data items to get the frequent itemsets.

HVCFPMINETREE [1] (Horizontal and vertical Compact Frequent Itemset Pattern Mining Tree) calculation consolidates all the most extreme event of frequent itemsets before changing over into the tree structure. It prompts a FPTree structure.

Frequent itemset Mining calculation [28] utilizing a Compressed Prefix Tree with example development (CT-ITL) was presented. It utilizes the Item-Trans Link (ITL) data structure that consolidates the advantages of both level and vertical data designs for association rule mining.

A. Two-stage calculation

This strategy keeps up a Transaction - weighted Downward Closure Property [3]. In this way, just the mixes of high exchange weighted usage itemsets are included into the hopeful set at each dimension amid the dimension shrewd inquiry. Stage I may overestimate some low utility itemsets, however it never disparages any itemsets. In stage II, just a single additional database check is performed to channel the overestimated itemsets. Two-Phase requests various databases check and create an immense number of hopeful itemsets on account of a dimension astute strategy.

B. Compacted Transaction Utility (CTU-Mine)

Erwin et al in [4] saw that the customary competitor - create and-test approach for distinguishing high utility itemsets isn't appropriate for thick data sets. Their work proposes a novel calculation CTU - Mine which mines high utility itemsets utilizing the example development approach. A comparable contention is displayed by Yu et al. Existing algorithms for high utility mining are section list based receiving an apriori like competitor set age and test approach and in this way are lacking in datasets with high measurements

C. Transient High utility itemset (THUI)

A tale technique, in particular THUI (Temporal High Utility Itemsets) – Mine was proposed by V.S. Tseng et al in [3] for mining transient high utility itemsets from data streams proficiently and successfully. The tale commitment of THUI-Mine is that it can viably distinguish the worldly high utility itemsets by producing less fleeting high exchange weighted usage 2 - itemsets with the end goal that the execution time can be decreased considerably in mining all high utility itemsets in data streams. Along these lines, the way toward finding all worldly high utility itemsets under record-breaking windows of data streams can be accomplished successfully with restricted memory space, less competitor itemsets and CPU I/O time. Th is meets the basic necessities on existence productivity for mining data streams. The exploratory outcomes demonstrate that THUI-Mine can find the worldly high utility itemsets with high execution and less competitor itemsets when contrasted with different algorithms under different trial conditions. Also, it performs adaptable as far as execution time under expansive databases. Subsequently, THUI - Mine is promising for mining worldly high utility itemsets in data streams.

D Utility pattern growth (UP -growth)

To address issue of generating a large number of candidates, UP-Growth [5] (V.S Tseng et al., 2010) has recently been proposed and it uses PHU (Potential High Ut ility) model. For reducing the number of candidate itemsets, the UP -Growth applies four strategies, DGU (Discarding Global Unpro mising items), DGN (Decreasing Global Node utilities), DLU (Discarding Lo cal Unpromising items), and DLN (Decreasing Local Node utilities). Besides, it constructs a tree structure, named UP -Tree, with two database scans and conducts mining high utility itemsets. In other words, it demands three database scans for discovering high utility itemsets. In the first database scan, TWU values of each item are accu mulated. In the second database scan, items hav ing less TWU values than the user-specified minimu m utility threshold are removed fro m each transaction. In addition, items in transactions are arranged according to TWU descending order and the transactions are inserted into the UP -Tree. In this stage, DGU and DGN are applied for reducing overestimated utilities. After that, high utility itemsets are generated fro m the UP -Tree with DLU and DLN.

E. High utility itemset miner (HUI-Miner)

Liu & Qu (2012) proposed HUI -Miner algorithm in [6]. It is a high utility itemsets with a list data structure, called utility list. It first creates an initial utility list for itemsets of the length 1 for promising items. Then, HUI-M iner constructs recursively a utility list for each itemset of the length k using a pair of utility lists for itemsets of the length k-1. For mining high utility itemsets, each utility list for an itemset keeps the informat ion of TIDs for all of transactions containing the itemset, utility values of the item set

in the transactions, and the sum of utilities of the remain ing items that can be included to super itemsets of the itemset in the transactions. The distinct advantage of HUI-M iner is that it avoids the costly candidate generation and utility co mputation.

F. Faster high utility itemset (FHM)

Philippe Fournier-Viger (2014) proposed FHM algorithm in [6]. It extends the Hui-Miner Algorithm. It is a Depth-first search Algorithm. It relies on utility-lists to calculate the exact utility of itemsets. This algorithm integrates a novel strategy named EUCP (Estimated Utility Co-occurrence Pruning) to reduce the number of joins operations when min ing high -utility itemsets using the utility list data structure. Estimated Utility Co -Occurrence Structure (EUCS) stores the transaction weighted utility (TWU) of all 2-itemsets. It built during the initial database scans. EUCS represented as a triangular matrix or hash map. The memo ry footprint of the EUCS structure is small. FHM is up to 6 times faster than HUI -Miner.

G. Efficient high utility itemset (EFIM)

EFIM (Efficient high-utility Item set Mining), wh ich introduces several new ideas to more efficiently discovers high -utility it em sets both in terms of execution time and memory [7]. EFIM relies on two upper-bounds named sub-tree utility and local utility to more effectively prune the search space. It also introduces a novel array -based utility counting technique called Fast Utility Counting to calculate these upper-bounds in linear t ime and space. Transaction merg ing is obviously desirable. However, a key problem is to imp lement it efficiently. To find identical transactions in $O(n)$ time,, sort the original database according to a new total order T on transactions. Sorting is achieved in time, and is performed only once. Projected databases generated by EFIM are often very small due to transaction merging.

Sr. no.	Studies	Year	Dataset	Method used	algorithm	Limitation
1	Ying Liu, Wei-keng Liao, A lok Choudhary	2005	Transaction dataset	Level wise approach	Two phase	Multiple scan of database and generate many candidate itemset
2	Alva Erwin, Raj P. Gopalan, N.R. Achuthan	2007	Transaction dataset	Pattern growth approach	Co mpressed Transaction Utility(CTU-M ine)	Co mplex for Evaluation due to the Tree structure
3	Vincent S. Tseng, Chun-Jung Chu, Tyne Liang	2008	Transaction dataset	Pattern growth	Temporal h igh utility itemset mining(THUI)	Huge memory requirement and a lot of false candidate itemset
4	Vincent S. Tseng, Cheng-Wei Wu, Bai-En Shie, and Philip S. Yu	2010	Transaction dataset	Pattern growth	Utility pattern growth(UP-growth)	Co mplex for Evaluation due to the Tree structure
5	Mengchi Liu, Junfeng Qu	2012	Transaction dataset	Level wise approach	High utility itemset miner(HUI-Miner)	Calculating the utility of an itemset joining utility list is very costly.
6	Philippe Fournier-Viger, Cheng-Wei Wu, Souleymane Zida, Vincent S. Tseng	2014	Transaction dataset	Level wise approach	Faster high utility itemset mining(FHM)	Static database, large memo ry overhead

Table 1.Performance summary of a survey

CONCLUSION

In data mining Association Rule Mining is a standout amongst the most critical assignments. An extensive number of proficient algorithms are accessible for association rule mining, which considers mining of frequent itemsets. Be that as it may, a developing point in Data Mining is Utility Mining, which fuses utility contemplations amid itemset mining. Utility Mining covers all parts of financial utility in data mining and aides in location of itemset having high utility, similar to benefit. High Utility itemset mining is extremely valuable in a few genuine applications. In this paper study paper, we give the different strategy for high utility itemset mining and correlation of all procedure with exchange dataset, which technique is utilized and constraint of every calculation.

REFER ENCES

- 1.U Kan imozhi, J K Kavitha, D Manjula, Mining High Utility Itemsets – A Recent Survey, International Journal of Scientific Engineering and Technology, Vo lu me No.3 Issue No.11, pp: 1339-1344, 2014.
- 2.Maya Joshi, Mansi Patel, A Survey on High Utility Itemset Mining Using Transaction Databases, International Journal of Computer Science and Information Technologies, Vol. 5 (6), 7407 -7410,2014.
- 3.Jyothi Pillai, O.P. Vyas, Overview Of Itemset Mining And its Application, International Journal of Computer Applications (0975 – 8887), Volume 5– No.11, August 2010.
- 4.Sudip Bhattacharya, Deepty Dubey, High Utility Itemset Mining, International Journal of Emerging Technology and Advanced Engineering, Vo lu me 2, Issue 8, August 2012.
- 5.Hua-Fu Li, Hasin-Yug Huang, Yi-Cheng Chen, Yu -Jiun Liu, Suh-Yin Lee, Fast and Memory Efficient Mining of High Utility Itemsets in Data Streams, IEEE International Conference on Data Mining,2008.
- 6.Shekhar Patel B Madhushree, A Survey on Discovering High Utility Itemset Mining from Transactional Database, Information and Knowledge Management, Vo l.5, No.12, 2015
- 7.R.Shyamala Devi, D.Shanthi, A Survey Mining High Utility Item Sets And Frequent Item Sets, International Journal of Innovative Research in Science, Engineering and Technology, Vol. 4, Issue 12, December 2015.
- 8.Smita R. Londhe,, Rupali A. Mahajan,, Bhagyashree J. Bhoyar ,”Overview on Methods for Mining High Utility Itemset from Transactional Database”, International Journal of Scientific Engineering and Research (IJSER), Vo lu me 1 Issue4, Deccember2013.
- 9.R. Agrawal and R. Srikant, 1994, —Fast Algorithms for Mining Association Rulesl, in Proceedings of the 20th International Conference Very Large Databases, pp. 487-499.
- 10.Attila Gyenesei, —Mining Weighted Association Rules for Fuzzy Quantitative ItemsI, Lecture Notes in ComputerScience, Springer, Vol.1910/2000, pages 187-219, TUCS Technical Report No.346, ISBN 952-12-659-4, ISSN 1239-1891, May 2000.
- 11.R. Chan, Q. Yang, Y. D. Shen, —Mining High utility ItemsetsI, In Proc. of the 3rd IEEE Intel.Conf. on Data Mining(ICDM), 2003.
- 12.H. Yun, D. Ha, B. Hwang, and K. Ryu. —Mining association rules on significant rare data using relative supportI. Journal of Systems and Software, 67(3):181–191, 2003.
- 13.H.Yao, H. J. Hamilton, and C. J. Butz, —A Foundational Approach to Mining Itemset Utilities from
- 14.DatabasesI,Proceedings of the Third SIAM International Conference on Data Mining, Orlando, Florida, pp. 482 -486,2004.
- 16.G. Weiss. —Mining with rarity: a unifying frameworkI, .SIGKDD Explor. Newsl., 6(1):7–19, 2004.
- 17.Liu, Y., Liao, W., and A. Choudhary, A., —A Fast High Utility Itemsets Mining AlgorithmI, In Proceedings of the Utility-Based Data Mining Workshop, August 2005.

19. Lu, S., Hu, H. and Li, F. 2005. —Mining weighted association rules. *Intelligent Data Analysis*, 5(3):211–225.
20. V. S. Tseng, C.J. Chu, T. Liang, —Efficient Mining of Temporal High Utility Itemsets from Data streams, *Proceedings of Second International Workshop on Utility-Based Data Mining*, August 20, 2006
22. H. Yao, H. Hamilton and L. Geng, —A Unified Framework for Utility-Based Measures for Mining Itemsets, In *Proc. Of the ACM Intel. Conf. on Utility-Based Data Mining Workshop (UBDM)*, pp. 28-37, 2006.
23. Erwin, R.P. Gopalan and N. R. Achuthan, 2007, —A Bottom-up Projection based Algorithm for mining high utility itemsets, in *Proceedings of 2nd Workshop on Integrating AI and Data Mining (AIDM 2007)*, Australia, *Conferences in Research and Practice in Information Technology (CRPIT)*, Vol. 84.
24. J. Hu, A. Mojsilovic, —High-utility pattern mining: A method for discovery of high-utility item sets, *Pattern Recognition* 40 (2007) 3317–3324.
25. L. Szathmary, A. Napoli, P. Valtchev, —Towards Rare Itemset Mining, *Proceedings of the 19th IEEE International Conference on Tools with Artificial Intelligence*, 2007, Volume 1, Pages: 305-312, ISBN ~ ISSN: 1082-3409, 0-7695-3015-X
28. Kriegel, H-P et al. 2007. —Future Trends in Data Mining, *Data Mining and Knowledge Discovery*, 15:87–97.