



High Utility Item Sets Mining for Transactional Databases

G Kurubindu¹, P Lalitha kumari²

#1. M.Tech (CSE) in Department of Computer Science Engineering,

#2. Assoc.Prof, Department of Computer Science and Engineering, Sri Sivani College Of Engineering,
Chilakapalem, Srikakulam, A.P, India.

Abstract:

Mainstream issue in data mining, which is called "high-utility itemset mining" or all the more for the most part utility mining. High Utility Itemsets which are itemsets having an utility gathering a client determined least utility edge value i.e min_util. The principle target of utility mining is to discover thing sets with highest utilities, by thinking about benefit, amount, cost or some other client inclinations. Research has been done in region of mining HUI's. Different procedures have been connected. The fundamental issue with setting edge value which is for the most part client particular, is it should be proper. In Order to set most fitting or right Threshold value for mining HUI's, user needs to do trial and mistake which thus is tedious and repetitive process, in light of the fact that if min_util is set too low, framework will bring about getting substantial data of HUI, which thus makes framework incapable with the end goal of HUI. In the event that we set min_util too high, this will bring about getting little sum or no HUI's. Consequently setting least edge value is troublesome. The proposed framework is following Top-k framework for mining top-k HUI's, which is utilizing two algorithms TKU (mining top-k utility itemsets) and TKO (mining top-k in one phase), without setting min_util edge.

Keywords: Frequent Itemset Mining, High Utility Itemset, Closed High Utility Itemsets, Top-k Mining, and Transaction Utility.

I. Introduction

Data mining is the productive disclosure of important and distinctive data from a substantial accumulation of data.. There are different subdivisions for it. Frequent itemset mining (FIM) find just regular things. In any case, benefit of the things are not considered. This is on the grounds that buy amount not considered, all things saw as having same significance and furthermore find frequent designs

that are not fascinating. A thing can be available or missing in an exchange. The generally utilized algorithms are Apriori, Eclat, LCM, Pre Post, FIN and FPGrowth algorithm. Frequent itemset mining (FIM) doesn't fulfill the prerequisite of the client and dealer. Certain affiliation rules are related with this mining. For instance the individuals who purchase bread and spread are probably going to purchase drain as well. In FIM computational cost is low however miss more vital examples to clients. High utility mining is an expansion to the issue of continuous example mining. High utility mining implies distinguishing itemsets having high benefits when they are sold together. Here the utility of an itemset is estimated regarding weight, benefit, cost, amount or other data relying upon the client inclination. The benefit produced and exchange tally is considered. An exchange database implies a database with a rundown of exchanges like butter, bread, milk and so forth. In High Utility Itemset (HUI) mining the itemsets that produce a benefit higher than least limit is called as high utility itemsets. So there will be an exchange database and a unit benefit table for each of the things in the exchange database. The hui mining is intriguing a result of two primary reasons. One is that find itemsets that create a high benefit in client exchanges than those that are purchased often. Additionally no Apriori property or against monotone property in hui mining. The counter monotone property is that if an itemset is rare then all its superset likewise rare and can be pruned. To beat this the idea of exchange weighted utilization (twu) presented. Since there is no hostile to monotone property wasteful as far as time and memory necessity or even come up short on memory? Distinctive algorithms for hui mining are UP-Growth, UpGrowth+, IHUP, IIDS, d2HUP, HUI-Miner and so forth. The algorithms for hui mining can be two or one stages. Hui mining takes additional time and space for execution. So fused an idea called shut hui mining. An itemset is said to be shut if there

exists no other itemset with the end goal that the first is the subset of other and furthermore there help tally ought not be equivalent. The different algorithms utilized are CHUI-Miner, CLSMiner[15] and so forth. Setting the edge is a dull activity in hui mining. On the off chance that the limit picked isn't an appropriate value there is a shot of missing a few things or itemsets. Top-k high utility itemset mining beats this trouble. In top-k hui mining, k is the coveted number of HUIs to be mined. So just less values must be put away sparing time and memory. There are different algorithms for top-k mining like TKU,TKO and so on. A couple of top k frequent design algorithms are BTK(TB-tree)[13], TFP [14](Top-k Frequent Pattern) and so on. In top-k no need of setting the limit, rather it is at first set to zero and raise consequently. Utilizing top-k we can without much of a stretch locate the top k sets of items that contribute the highest benefits to the organization. The top k can be connected to CHUD algorithm after stage I to diminish the quantity of hopefuls put away with the goal that putting away just top k values from the IT-tree.IT-tree is utilized as a part of CHUD algorithm. Here the k is the base utility edge determined by us. Just top k values are put away in light of the base utility limit.

II. Related work

As of late, high utility itemset mining has gotten bunches of consideration and numerous effective algorithms have been proposed. Some earlier methods presented for mining High Utility Itemsets are quickly recorded. Philippe Fournier-Viger, Cheng-Wei Wu, Souleymane Zida, Vincent S.Tseng [3] proposed an algorithm named FHM which diminishes the quantity of join tasks by presenting a novel structure named EUCS(Estimated Utility Cooccurrence Structure). However FHM underpins just static databases and the time has come devouring. Vincent S. Tseng, Bai-En Shie, Cheng-Wei Wu, and Philip S. Yu, proposed a productive algorithm for mining high utility itemsets from value-based databases. In this work, two algorithms, in particular utility example development (UP-Growth) and UP-Growth+, is proposed for mining high utility itemsets with an arrangement of successful systems for pruning hopeful itemsets. The data of high utility itemsets is kept up in a tree-based data structure named utility example tree (UP-Tree) to such an extent that competitor itemsets can be produced proficiently with just two outputs of database. Since

this framework requires different stages to examine the database it is highly perplexing. Sen Su, Shengzhi Xu, Xiang Cheng, Zhengyi Li, and Fangchun Yang proposed a private FP-development algorithm, which is alluded to as PFP-development. The PFP-development algorithm comprises of a pre-handling stage and a mining stage. In the pre-preparing stage, to enhance the utility and protection exchange off, a novel savvy part strategy is proposed to change the database. For a given database, the pre-handling stage should be performed just once. In the mining stage, to counterbalance the data misfortune caused by exchange part, we devise a run-time estimation technique to evaluate the genuine help of itemsets in the first database. The constraint is that if itemsets of a similar length are produced at the same time, the quantity of FP-trees put away in memory will develop at an exponential rate. The primary normal for one-stage algorithms is that they find high utility itemsets utilizing just a single stage and create no competitors. Junqiang Liu, Ke Wang, Benjamin C.M. Fung proposed an algorithm, d2 HUP, in particular Direct Discovery of High Utility Patterns, for mining high utility itemsets in one stage without competitor age. It is a reconciliation of the profundity first pursuit of the turn around set specification tree, the pruning methods that radically lessens the quantity of examples to be identified, and a novel data structure that empowers productive calculation of utilities and upper limits. HUI-Miner considers a database of vertical configuration and changes it into utility-records. The utility-list structure utilized as a part of HUI-Miner permits specifically figuring the utility of created itemsets in fundamental memory without examining the first database. In spite of the fact that the above investigations may perform well in a few applications, they are not created for top-k high utility itemset mining regardless they experience the ill effects of the unobtrusive issue of setting suitable limits.

III. Methodology

About CHUD and IT-tree

It consists of two phases. Phase I - potential closed high utility itemsets found(PCHUI).PCHUIs are defined as a set of itemsets having estimated utility no less than absolute minimum utility. Phase II-By scanning the database once, CHUIs are identified from the set of PCHUIs found in Phase I and their utility unit array are computed. Utility unit array for lossless representation. CHUD [Closed+ High Utility

Itemset Discovery] is an augmentation of Eclat[6] and DCI-Closed[7] calculation. It considers vertical database and mines CHUIs in a profundity first way. An itemset is closed+ high utility itemset(CHUI) in the event that it is a shut high utility itemset and that itemset ought to be incorporated into utility unit cluster.

IT-Tree[9] implies Itemset-Tidset match tree. It is for discovering CHUIs.Each hub has an itemset, Tid, two requested arrangements of things named PREVIOUS-SET and POST-SET. Figure1 demonstrates an IT-tree for the exchange database appeared in Table2. Additionally here every hub is appended with an expected utility esteem. In any case, it isn't appeared in the figure1,but clarified in past segments. An information structure called exchange utility table (TU-table) is embraced for putting away the exchange utilities of exchanges. It is a couple with Tid and the second esteem is the exchange utility. So if Tid given, we can productively recover the TU esteem [8].The working of IT-tree like Eclat (Equivalence Class Transformation) calculation. Given an itemset X ,t(X) is the arrangement of all tids that contains X .From the above value-based database; Example:- (ABE)=12.Difference amongst Eclat and Apriori is that how they crosses prefix tree and how they decide the help of an itemset. Apriori is crossed in expansiveness first request and bolster computed by checking the entire database. Éclat is profundity first requesting and the help of another itemset by figuring the crossing point between tidsets.Given a tidset Y ,i(Y) is the arrangement of every single regular thing to all the tids in Y.

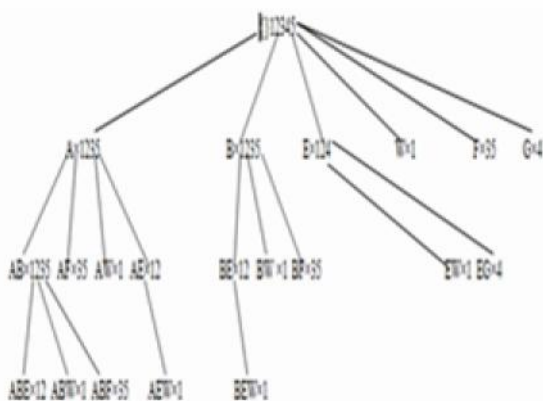


Figure 1: A case IT-tree

Conclusion of an itemset X [c(X)] is the littlest shut set that contains X[20,21].To discover the conclusion of an itemset X:-

- 1) Compute the picture of X in the exchange space to get t(X).
- 2) Next guide t(X) to its picture in the itemset space utilizing the mapping to get i(t(X)). At that point that subsequent itemset must be shut. That is, $c(X) = I \circ t(X) = i(t(X))$.

From Table2,c(ABE) = i(t(ABE)) = i(12) = ABE. In this manner ABE shut. At that point think about AE; c(AE) = i(t(AE)) =i(12).As as of now observed i(12) can be A, B, AB, AE, ABE. So it isn't shut. The help of an itemset X is equivalent to help of closure[1,2].In an IT-tree system, for a given hub or prefix class, one can perform crossing points of the tidsets of all sets of components in a class, and check if least help is met; bolster tallying is concurrent with age. Each subsequent continuous itemset is a class unto itself, with its own particular components, that will be recursively extended. In other words, for a given class of itemsets with prefix P, [P] = {I1, I2, ..., In}, one can play out the crossing point of t(Pl_i) with all t(Pl_j) with j > I, to get another class of continuous extensions,{ [Pl_i] = {I_j | j > I and σ(Pl_i * I_j) ≥ least support}. Additionally subsumption checking performed. Give X_i and X_j a chance to be two itemsets,then X_i subsumes another itemset X_j,if and just if X_j subset of X_i and bolster count(X_i) = bolster count(X_j).The requirement for subsumption checking is that in the wake of including a shut set when we investigate consequent branches, we may produce another set, which can't be expanded further and the new set is the subset or equivalent to the previous and their help tallies are equivalent. For this situation the new set is non-shut set subsumed by the previous and it ought not be added to shut set. For instance From Table2,{F} is subsumed by {ABF} in light of the fact that {F} ⊂ {ABF} and SC({F}) = SC({ABF}). A thing is a promising thing if exchange weighted usage of the thing is more noteworthy than or equivalent to supreme least utility. Otherwise,it is an unpromising item[1].

Utility unit exhibit is to make the portrayal lossless.Each shut HUI is clarified with an uncommon structure called utility unit cluster to such an extent that the subsequent itemset is known as a closed+ high utility itemset[1,4]. The possibility of utility unit cluster makes the set of CHUIs lossless

because HUIs and their utilities can be derived from this set without accessing the original database. Let $X=[v_1, v_2, \dots, v_k]$ and contains k utility values. The i th utility value v_i in $v(X)$ is denoted as $v(X, a_i)$. The utility value of X can be expressed as $u(X) = \text{sum of all}(X, a_i)$. From Table 2, the first utility value in $V\{ABF\}$ is $V\{ABF, A\} = \text{absolute utility}(A, T3) + \text{absolute utility}(A, T5) = 2$. $V\{ABF, B\} = \text{absolute utility}(B, T3) + \text{absolute utility}(B, T5) = 2$. $V\{ABF, F\} = \text{absolute utility}(F, T3) + \text{absolute utility}(F, T5) = 15$. Therefore the utility unit array of $\{ABF\}$ is $V\{ABF\} = [2, 2, 15]$.

During the first phase subsume check. Compute closure of itemsets. Then exploration of the IT-tree is done.

ALGORITHM Subsume Check

Input: Node X , Prev-set(X).

Output:- True: If X is non-closed and subsumed by other itemsets.

False: If X is not subsumed by other itemsets.

Step1: For each item i an element of PREV-SET(X) do
Step2: If $g(X)$ Subset or equal to $g(a)$ then true else false.

ALGORITHM Closure-itemsets

Input: The node of $X; N(X), \text{POST-SET}(X)$.

Output: The closure of $X; X_c$.

Step1: Initialize $X_c = X$.

Step2: For each item $i \in \text{POST-SET}(X)$ do

Step3: If $g(X)$ subset or equal to $g(i)$ then

Step4: $\text{POST-SET}(X) = \text{POST-SET}_X \cup \{i\}$.

Step5: $X_c = X_c \cup \{a\}$.

Step6: Finally X_c is returned.

ALGORITHM: Top-k CHUD(TCHUD)

Input: The transactional database, Minimum utility value. Output: The complete set of CHUIs.

Step1: Scan the database to convert into vertical database. Step2: After scan promising items are collected into an ordered list, in the increasing order of support. Unpromising Step8: Absolute utility(Y) \geq Absolute minimum utility do step

9 to 13.

Step9: Again Y an element of HC_{k-1} and Support(X) > Support(Y) then

Step10: Support(Y) = Support(X).

Step11: Otherwise, if Y not an element of HC_{k-1} , then

Step12: $HC_{k-1} = HC_{k-1} \cup Y$.

Step13: Support(Y) = Support(X).

IV. Proposed System

In the proposed system Top-K High Utility Itemsets is mined by using a novel algorithm named TKO in one phase. In TKO min util is not set. Instead a border minimum utility threshold is used. TKO can raise the min_util threshold as quickly as possible, and further reduce as much as possible the number of candidates and intermediate low utility itemsets produced in the mining process. A new framework for top-k high utility itemset mining, where k is the desired number of HUIs. TKO (mining Top-K utility itemsets in One phase) are proposed for mining the complete set of top-k HUIs in databases without any need to specify the min_util threshold. the TKO algorithm uses a list-based structure named utility-list to store the utility information of itemsets in the database.

TKO uses vertical data representation techniques to discover top-k HUIs in only one phase. It utilizes the search procedure of HUI-Miner and its utility-list structure. When an itemset is generated by TKO, its utility is calculated by its utility-list without scanning the original database. We first describe a basic version of TKO named TKO Base and then the advanced version, which includes several strategies to increase its efficiency. Many different types of data structure and algorithm have been proposed to extract frequent pattern from a large given database. CR algorithm is one of the fastest frequent pattern mining algorithm, which can efficiently represent whole data structure over single scan of the database. We have proposed an efficient tree based structure CR Tree in terms of execution time and memory usage. Corelation Tree or CR Tree algorithms. Even a huge database can be processed by CR Tree if out-of-date transactions are removed concurrently. CR Tree is better than FP tree.

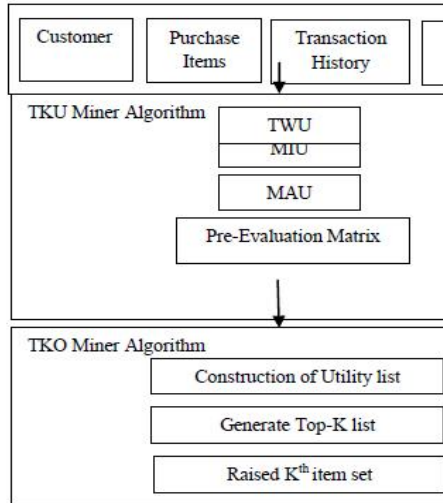


Figure 2. proposed Architecture diagram

FHM (Fast High utility Mining)

Our proposition depends on the perception that despite the fact that HUI-Miner plays out a solitary stage and in this manner don't create competitors according to the meaning of the two-stage show, HUI-Miner investigates the inquiry space of itemsets by producing itemsets and an expensive join activity must be performed to assess the utility of each itemset. To decrease the quantity of joins that are performed, we propose a novel pruning system named EUCP (Estimated Utility Cooccurrence Pruning) that can prune itemsets without performing joins. This procedure is anything but difficult to execute and exceptionally powerful. We name the proposed calculation joining this technique FHM (Fast High-utility Miner). We look at the execution of FHM and HUI-Miner on four genuine datasets. Results demonstrate that FHM performs up to 95 % less join tasks than HUI-Miner and is up to six times speedier than HUI-Miner.

The FHM computation

In this segment, we display our proposition, the FHM calculation. The principle methodology (Algorithm 1) takes as information an exchange database with utility esteems and the $minutil$ edge. The calculation first outputs the database to compute the TWU of every thing. At that point, the calculation recognizes the set I^* of all things having a TWU no not exactly $minutil$ (different things are disregarded since they can't be a piece of a high-utility itemsets by Property

3). The TWU estimations of things are then used to set up an aggregate request \neg on things, which is the request of climbing TWU esteems (as proposed in [7]). A moment database check is then performed. Amid this database filter, things in exchanges are reordered by the aggregate request \neg , the utility-rundown of every thing $I \in I^*$ is assembled and our novel structure named EUCS (Estimated Utility Co-Occurrence Structure) is manufactured. This last structure is characterized as an arrangement of triples of the frame $(a, b, c) \in I^* \times I^* \times R$. A triple (a,b,c) shows that $TWU(\{a, b\}) = c$. The EUCS can be actualized as a triangular lattice or as a hashmap of hashmaps where just tuples of the frame (a, b, c) with the end goal that $c \neq 0$ are kept. In our usage, we have utilized this last portrayal to be more memory proficient in light of the fact that we have watched that couple of things co-happens with different things. Building the EUCS (Estimated Utility Co-Occurrence Structure) is quick (it is performed with a solitary database examine) and involves a little measure of memory, limited by $|I^*| \times |I^*|$, in spite of the fact that practically speaking the size is considerably littler on the grounds that a set number of sets of things co-happens in exchanges (cf. segment 5). After the development of the EUCS, the profundity first pursuit investigation of itemsets begins by calling the recursive methodology Search with the purge itemset \emptyset , the arrangement of single things I^* , $minutil$ and the EUCS structure. Algorithm Steps

Algorithm 1: The FHM algorithm

Input : D: a transaction database, $minutil$: a user-specified threshold

Output: the set of high-utility itemsets

Step:1 Scan D to calculate the TWU of single items;

Step:2 $I^* = \{i \mid TWU(i) < minutil\}$;

Step:3 Let σ be the total order of TWU ascending values on I^* ;

Step:4 Scan D to build the utility-list of each item $i \in I^*$ and build the EUCS structure;

Step:5 Search $(\emptyset, I^*, minutil, EUCS)$;

GRAPHS ANALYSIS

Comparison between level and number of candidate
On the x axis we have taken level number and on y axis numbers of candidates are displayed.

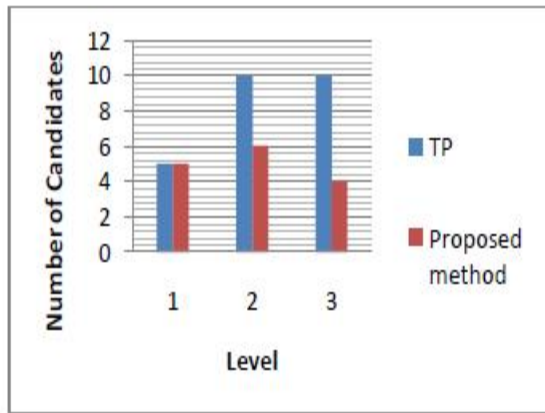


Figure 3 Comparison on the basis of level and number of candidate

V. Conclusion

In this paper, we have considered the issue of top k high utility itemsets mining, where k is the coveted number of high utility itemsets to be mined. Two productive calculations TKU (mining Top-K Utility itemsets) and TKO (mining Top-K utility itemsets in One stage) are proposed for mining such itemsets without setting least utility edges. TKU is the initial two-stage calculation for mining top-k high utility itemsets, which consolidates five procedures PE, NU, MD, MC and SE to successfully raise the fringe least utility limits and further prune the pursuit space. Then again, TKO is the first stage calculation created for top-k HUI mining, which coordinates the novel methodologies RUC, RUZ and EPB to enormously enhance its execution. Observational assessments on various kinds of genuine and manufactured datasets demonstrate that the proposed calculations have great adaptability on vast datasets and the execution of the proposed calculations is near the ideal instance of the condition of-the-art two-stage and one-stage utility mining calculations.

Despite the fact that we have proposed another system for top-k HUI mining, it has not yet been fused with other utility mining undertakings to find diverse sorts of top k high utility examples, for example, top-k high utility scenes, top-k shut high utility itemsets, top-k high utility web get to

examples and top k versatile high utility successive examples. These leave wide spaces for investigation as future work.

References

- [1] R. Agrawal and R. Srikant, "Fast algorithms for mining association rules," in Proc. Int. Conf. Very Large Data Bases, 1994, pp. 487–499.
- [2] C. Ahmed, S. Tanbeer, B. Jeong, and Y. Lee, "Efficient tree structures for high-utility pattern mining in incremental databases," IEEE Trans. Knowl. Data Eng., vol. 21, no. 12, pp. 1708–1721, Dec. 2009.
- [3] K. Chuang, J. Huang, and M. Chen, "Mining top-k frequent patterns in the presence of the memory constraint," VLDB J., vol. 17, pp. 1321–1344, 2008.
- [4] R. Chan, Q. Yang, and Y. Shen, "Mining high-utility itemsets," in Proc. IEEE Int. Conf. Data Mining, 2003, pp. 19–26.
- [5] P. Fournier-Viger and V. S. Tseng, "Mining top-k sequential rules," in Proc. Int. Conf. Adv. Data Mining Appl., 2011, pp. 180–194.
- [6] P. Fournier-Viger, C. Wu, and V. S. Tseng, "Mining top-k association rules," in Proc. Int. Conf. Can. Conf. Adv. Artif. Intell., 2012, pp. 61–73.
- [7] P. Fournier-Viger, C. Wu, and V. S. Tseng, "Novel concise representations of high utility itemsets using generator patterns," in Proc. Int. Conf. Adv. Data Mining Appl. Lecture Notes Comput. Sci., 2014, vol. 8933, pp. 30–43.
- [8] J. Han, J. Pei, and Y. Yin, "Mining frequent patterns without candidate generation," in Proc. ACM SIGMOD Int. Conf. Manag. Data, 2000, pp. 1–12.
- [9] J. Han, J. Wang, Y. Lu, and P. Tzvetkov, "Mining top-k frequent closed patterns without minimum support," in Proc. IEEE Int. Conf. Data Mining, 2002, pp. 211–218.
- [10] S. Krishnamoorthy, "Pruning strategies for mining high utility itemsets," Expert Syst. Appl., vol. 42, no. 5, pp. 2371–2381, 2015.

Authors



G Kurubindu completed her B.Tech at Sri Sivani College of engineering chilakapalem. She is pursuing m.tech in Sri Sivani College of engineering chilakapalem, Srikakulam, AP,

India.



P Lalitha kumari M.Tech (Ph.D) is working as Associate Professor in the Department of CSE at Sri Sivani College Of Engineering, Chilakapalem, Srikakulam, A.P, India.