



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 5 Issue: XII Month of publication: December 2017

DOI:

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Survey On Mining Top K High Utility Itemsets in One Phase

G Naga Kalyani¹, Prof. S. Viziananda Row²

¹ M. Tech Student, Dept. of Computer Science and Systems Engineering, AU College of Engineering (A), Visakhapatnam, India.

²Professor, Dept. of Computer Science and Systems Engineering, AU College of Engineering (A), Visakhapatnam, India.

Abstract: This paper presents a study on finding Top K itemsets with high utility. High Utility Item sets (HUI) mining has emerged as an interesting and challenging research topic in data mining. It finds applications in web-click analysis, cross-marketing in retail stores and bio medical analysis etc. The number of high-utility itemsets that can be extracted from a transactional database depends upon the value of minimum utility threshold. It is often difficult for a user to find a suitable threshold value which fits their purpose. The database can generate many high-utility itemsets at low threshold value and very few itemsets at higher threshold values. In order to relieve the user from this tedious task, we use an efficient algorithm named TKO for mining top-k high utility itemsets from a large transactional database. The parameter k can be set by the user according to his/her needs. We conduct extensive experiments on real and synthetic datasets and the experimental results demonstrate the effectiveness of our Threshold raising strategies in terms of total execution time and the memory usage.

Keywords: Utility mining, high utility mining, top k high utility itemset mining

I. INTRODUCTION

Frequent pattern mining finds patterns from a database, which have frequency no less than a given minimum support threshold. Frequent pattern mining finds applications in market-basket analysis, mining association rules [1], plagiarism detection and biomedical data analysis. The algorithms developed for mining frequent patterns have mostly employed the monotone/anti-monotone property to prune the exponential search space effectively. The monotone property states that the subsets of a frequent pattern are also frequent and the anti-monotone property states that the supersets of an infrequent pattern are also infrequent. However, the frequent patterns extracted can be of low profit value. The concept of high-utility pattern mining was introduced to capture the notion of utility, which has been observed in real life. High-utility pattern mining finds patterns from a database which have their utility value no less than a given minimum utility-threshold. The utility function measures the importance of a pattern and varies according to the application. For example, in a retail store domain, a manager may be interested in finding combinations of products with high profits, which relates to the unit profits and purchased quantities of products that are not considered in frequent pattern mining. High-utility pattern mining has wide range of applications in cross-marketing in retail stores, web-click stream analysis, medicine etc. Be that as it may, effectively mining HUIs in databases is not a simple assignment on the grounds that the downward closure property utilized as a part of FIM does not hold for the utility of itemsets. At the end of the day, pruning scan space for HUI mining is troublesome on the grounds that a superset of a low utility itemsets set can be high utility. To handle this issue, the idea of exchange weighted usage (TWU) model was acquainted with encourage the execution of the mining assignment. Although many algorithms have been introduced in HUI mining, it is very difficult for users to set an appropriate minimum support because it highly depends on data types. If it is set too high, no result itemsets are found while too small value makes an enormous number of result patterns which cause inefficiencies in terms of computation time and memory usage. to address this issue, and to control the itemsets with the highest utilities without setting the thresholds, a better solution is to change the task of mining HUIs as mining top-k high utility itemsets. Here the users specify k. Here k is the number of desired itemsets, instead of specifying the minimum utility threshold. Top k HUI mining is used to find, what are the top-k sets of products that contribute the highest profit to the company and how to efficiently found these itemsets without setting the min-utility threshold. A naive approach for extracting top-k high-utility itemsets can be to set the minimum utility threshold to zero and apply any high-utility itemset mining algorithm to and the complete set of high-utility itemsets. Top-k itemsets can be then chosen from the result set.

However, this approach is computationally very inefficient as the search space is exponential in the number of different items. In order to improve the efficiency, we propose effective strategies to raise the minimum utility threshold from zero as quickly as possible. top-k HUI mining is essential to many applications, developing efficient algorithms for mining such patterns is not an easy task. It poses four major challenges those are,

- A. The utility of itemsets is neither monotone nor anti monotone.
- B. How to incorporate the concept of top-k pattern mining with the TWU model.
- C. The min_util threshold is not given in advance in top-k HUI mining.
- D. How to effectively raise the min_util Border threshold without missing any top-k HUIs.

In this paper we are presents the literature survey study over the concept of Top K high utility itemset mining using the concepts of data mining. we study an efficient algorithm named TKO (mining Top K utility itemsets in One phase) is proposed for mining the complete set of top k HUIs in databases without the need of specify the min_util threshold. We conduct extensive experiments on real as well as synthetic datasets and the experiment results demonstrate the effectiveness of our approach. The paper is organized as follows. Section 2 reviews therelated work and background knowledge and definitions related to utility mining is explained in Section 3. TKO algorithm and strategies forimproving the performance in Section 4. The experimentalresults are presented in Section 5 and Section 6 concludesthe paper.

II. REVIEW OF LITERATURE

High Utility Itemset Mining is a popular concept and many algorithms have been proposed for HUI mining such as two phase[6], IHUP[10], IIDS, UP-Growth[11], D2HUP and HUIMiner[8]. These algorithms can be generallyclassified in two types: Two-phase and one-phase algorithms. The characteristics of two-phase algorithm is that itconsist of two phases. In the first phase, they create a set of candidates that are potential high utility itemsets. In thesecond phase, they calculate the precise utility of each candidate found in the first phase to identify high utilityitemsets. Two-phase, IHUP, IIDS, and UP-Growth are two-phase based algorithms. D2HUP and HUIminer are one phase based algorithms. et al in [6] proposes two-phase algorithm for finding high utility itemsets is proposed. Utility mining is to find all the itemsets whose utility values are beyond a user specified threshold. This paper explain transaction weighted utilization in Phase I, only the combinations of high transaction weighted utilization itemsetsare added into the candidate set at each level during the level-wise search. In phase II, another database scan is performed to removes the overestimated itemsets. Two-phase requires fewer database scans, less computational cost and less memory space. It performs efficiently in terms of speed and memory cost both on synthetic and real databases, even on large databases.

Shuning Xing et alIn [4] proposed a process of UP-Tree by introducing a Fast Utility Tree (FU-Tree) is proposed. In this method, they introduce the Link Queue to reduce the number of scanning the original database and adopt prefix utility to minimize the overestimated utility. The theoretical analyses and experimental results show that FU-Tree outperforms UP-Tree in the time consumption of construction trees, and enhances the efficiency of mining high utility itemsets.

Ahmed et al. [10] proposed a tree-based algorithm, called IHUP, for mining high utility itemsets. They use an IHUP-Tree to maintain the information of high utility itemsets and transactions. Every node in IHUP-Tree consists of an item name, a support count, and a TWU value. The algorithm finds out High Utility Itemsets in three steps. In the first step the algorithm constructs IHUP-Tree in the second step it generates HTWUIs. Finally from these HTWUIs it identifies the high utility itemsets.

Tseng et al. [11] proposed an algorithm called UP-Growth (Utility Pattern Growth) which uses data structure, named UP-Tree to maintain the information of high utility itemsets and transactions. The framework of the algorithm consists of three parts: (1) construction of UP-Tree, (2) generation of potential high utility itemsets from the UP-Tree by UP-Growth, and (3) identification of high utility itemsets from the set of potential high utility itemsets. UP-Growth uses two strategies, namely *DLU* (Discarding local unpromising items) and *DLN* (Decreasing local node utilities) for efficiently generating PHUIs from the global UP-Tree with Although strategies DGUand DGN can effectively reduce thenumber of candidates in phase I, they are applied during theconstruction of the global UP-Tree. The UP-Growth outperforms the state-of-the-art algorithms substantially, especially when the database contains lots of long transactions.

M. Liu et In [8] proposed a novel data structure called utility-list, and developed an efficient algorithm HUI Miner, for high utility itemset mining. Utility-lists provide not only utility information about itemsets but also important pruning information for HUI-Miner. The framework of the algorithm HUI-Miner does not generate candidate high utility itemsets. After constructing the initial utility-lists from a original database, HUI-Miner can mine high utility itemsets from these utility-lists without scan the original database. The initial utility lists are considered by scans the database twice. HUI-Miner gains significant performance Improvement over state art algorithms in terms of both running time and memory consumption.

Vincent S. Tseng et al In [2] proposes a novel framework for top-k high utility itemset mining, where k is the desired number of HUIs to be mined is proposed. Two types of efficient algorithms named TKU (mining Top-K Utility itemsets) and TKO (mining Top-K utility itemsets in one phase) are proposed for mining such itemsets without the need to set min_util.

III. TERMS AND DEFINITIONS

In this section, define the related mathematical definitions for discovering interesting and useful information from a transaction database.

TABLE 1: Transactional Database

Tid	Transaction	Count
T1	{a,c,d}	{1,1,1}
T2	{a,c,e,g}	{2,6,2,5}
T3	{a,b,c,d,e,f}	{1,2,1,6,1,5}
T4	{b,c,d,e}	{4,3,3,1}
T5	{b,c,e,g}	{2,2,1,2}

TABLE 2: Utility Table

Item	a	b	c	d	e	f	g
Profit	5	2	1	2	3	1	1

TABLE 3: Transaction Utility Table

TID	TU
T1	8
T2	27
T3	30
T4	20
T5	11

TABLE 4: Transaction Weighted Utility Table

ITE	TW
M	U
A	65
B	61
C	96
D	58
E	88
F	30
G	38

A. Transaction Database

Let I be a set of items. A transaction database is a set of transactions $D = \{T_1, T_2, \dots, T_n\}$ such that for each transaction $T \in I$ and T_c has a unique identifier c called its Tid. Each item $i \in I$ is associated with a positive number $eu(i)$, called its external utility (e.g. unit profit). For each transaction T_c , T such that $i \in T_c$, a positive number $iu(i, T_c)$ is called the internal utility of i (e.g. purchase quantity).

B. Absolute utility of an item

Utility of an item i_b in a transaction t_v is denoted as $U(i_b, t_v)$ and defined as

$$U(i_b, t_v) = Q(i_j, t_v) * P(i_j, I)$$

C. Absolute utility of an itemset in a transaction

Utility of itemset c in transaction t_v is denoted as $U(c, t_v)$ and defined as follows $U(X, T) = \sum_{i \in X \wedge X \subseteq T} U(i, T)$.

For example, in TABLE 1, $U(\{a, e\}, T_3) = U(a, T_3) + U(e, T_3) = 4 \times 1 + 1 \times 4 = 8$, and $U(\{a, e\}) = U(\{a, e\}, T_3) + U(\{a, e\}, T_2) = 8 + 13 = 21$.

D. Transaction utility and total utility:

The transaction utility (TU) of a transaction T_r is defined as $TU(T_r) = U(T_r, T_r)$. the total utility of DB is the sum of the utilities of all the transactions in DB.

The total utility of a database D is denoted as $TU(D)$ and defined as $\sum_{T_r \in D} TU(T_r)$.

E. Transaction-weighted utilization:

The transaction-weighted utility of itemset X in DB, denoted as $TWU(X)$, is the sum of the utilities of all the transactions containing X in DB, where

$$TWU(X) = \sum_{T \in DB \wedge X \subseteq T} TU(T).$$

F. Remaining utility of an itemset in a transaction:

The remaining utility of itemset X in transaction T, denoted as $RU(X, T)$, is the sum of the utilities Σ of all the items in T/X in T, where $RU(X, T) = \sum_{i \in (T-X)} u(i, T)$.

G. Z-element

An element is called Z-element if and only if its remaining utility value (RU) is equal to zero. Otherwise, the element is called NZ-element. The set of all Z-elements in the utility list of X is denoted as $ZE(X)$.

H. High Utility Itemset

An item set X is called high utility item set (HUI) if $U(X)$ is no less than a user-specified minimum utility threshold min-util otherwise X is a low utility item set.

I. Top-k high utility item set

An item set X is called top-k high utility itemset in D iff there are less than k itemsets whose utilities are larger than $EU(X)$ in $fHUI(D, 0)$.

J. Optimal minimum utility threshold

An absolute minimum utility threshold d is called optimal minimum utility threshold iff there does not exist a threshold.

IV. SYSTEM OVERVIEW AND METHODOLOGY

The problems with previous HUI miner algorithms are: It generates a huge set of HUIs. These huge number of HUIs forms a challenging problem to the mining performance and higher processing time it consumes. It also generates huge number of candidate itemsets, then higher processing time it consumes. and setting an appropriate threshold value to mine HUIs is difficult task to user. The burden of threshold setting arises a new area called top-k mining. In top-k no threshold is set externally. It is set internally by the algorithm itself. Initially the minimum utility threshold is set to zero. The threshold is set automatically as the work progresses. There are different methods for raising the threshold. Using k instead of minimum utility threshold helps in finding only top-k sets of products that contribute highest profits to the company. Also an efficient way to find itemsets without setting minimum utility threshold.

In this system using TKO algorithm to mine High Utility Itemsets without using minimum threshold value in one phase from large transactional databases. In the previous algorithm analysis we study the high utility algorithms are divided into one phase and two phase algorithms i.e. tree based and list based data structure algorithms.

A. Data Structure

The TKO algorithm uses utility-list data structure to store the utility information about a database. The utility lists of each itemset is called as initial utility-lists, which are constructed by scanning the database twice.

The utility-list of an item(set) X consists of one or more tuples. Each tuple represents the information of X in a transaction T_r and has three fields: $\langle Tid, iutil, rutil \rangle$. Fields Tid and iutil respectively contains the identifier of T_r and the utility of X in T_r . Field rutil indicates the remaining utility of X in T_r .

inserted in the list and the (k+1)th episode is removed from the buffer. After the algorithm terminates, top-K list contains the desired output of top-k high utility itemset mining in the large transactional database.

C. Algorithm 1: TKO (D, K) Algorithm

- 1) *Input:* D: a transaction database, K: number of HUIS
- 2) *Output:* The set of Top k high-utility itemsets
- // step 1: construct initial utility list for candidate 1 itemsets.
 - a) Scan database D to calculate the TWU of single items.
 - b) Initially δ : Min_util set to 0.
 - c) $I^* \leftarrow$ each item i such that $TWU(i) \geq \delta$.
 - d) Let be the total order of TWU ascending values on I^* .
 - e) Scan database D to build the utility-list of each item $i \in I^*$.
- //step 2: explores search space by calling search procedure.
 - f) Search ($\emptyset, I^*, \text{null}, \text{min_util}, k, \text{TopK-CL-List}$);
 - i) For each $X = \{x_1, x_2, \dots, x_L\} \in \text{Class}[P]$ do
 - ii) If ($\text{SUM}(X.\text{iutils}) \geq \delta$) then
 - iii) Raise min_util Border by the strategy RUC
 - iv) $\delta \leftarrow \text{RUC}(X, k, \text{TopK-CI-List})$
 - v) End if
 - vi) If ($\text{NZEU}(X) + \text{RU}(X) \geq \delta$) then
 - vii) Initialize $\text{Class}[X]$ and $\text{ULS}[X]$ to NULL.
 - viii) For each $Y = \{y_1, y_2, \dots, y_L\} \in \text{Class}[P] \ y_L > x_L$ do
 - ix) Concatenate X and Y then store in Z $Z \leftarrow XUY$
 - //step 3: construct candidate L+1 itemsets by construct procedure
 - x) Construct utility list of Z by calling construct method $\text{ul}(Z) \leftarrow \text{Construct}(\text{ul}(P), X, Y, \text{ULS}[P])$
 - xi) $\text{Class}[X] \leftarrow \text{Class}[X] \cup Z$
 - xii) $\text{ULS}[X] \leftarrow \text{ULS}[X] \cup \text{ul}(Z)$ revised utility list after construct of z itemset for further search process.
 - xiii) End for
 - xiv) Recursive call Top K-HUI-Search($X, \text{ULS}[X], \text{Class}[X], \delta, \text{Top K-CI-List}$)
 - xv) End if
 - xvi) End for
- g) Output each item $i \in \text{TopK-CL-List}$;

D. Algorithm 2: RUC Algorithm

- 1) *Input:* X an itemset, TopK-CL-List: list that contains Top K HUIS, K: number of HUIS
- 2) *Output:* δ , optimal threshold value.
 - i) if $\text{size}(\text{TopK-CL-List}) < k$ then
 - ii) Add X item (which utility $>$ min_util Border) to Topk-CL-List.
 - iii) Else
 - iv) if $\text{utility}(X) >$ min utility then
 - v) Remove kth high utility itemset i.e. itemset(s) having least utility.
 - vi) Add itemset to the List.
 - vii) Sort the list in decreasing order of utility values of itemsets.
 - viii) min utility = least utility in TopK-CL-List.
 - ix) End if
 - x) End if

The Top K itemsets for the above transactional database in Table 1 where k=3 are, {<dbec>: 40, <bec>: 37, <dbe>:36}. The search space was shown in Fig. 2.

E. Threshold raising Strategies

The TKO algorithm generates many candidates since the minimum threshold start from zero. Algorithm raises the minimum threshold before mining high-utility itemsets from a large transactional database by bellow specified strategies.

- 1) RUC (Raising threshold by Utility of Candidates)
- 2) RUZ (Reducing estimated utility values by using Z elements)
- 3) EPB (Exploring the most Promising Branches first)

V. EXPERIMENTAL RESULTS

A. Experiments

To evaluate the performance of TKO Algorithm, we have done extensive experiments on various databases, in which TKO algorithm is compared with the state-of-the-art mining algorithms UP-growth and HUI-Miner with optimal thresholds.

B. Experimental setup

For experiments we use both real and synthetic datasets. Foodmart is a real-life sparse dataset from a retail store, with real utility values. which is acquired from Microsoft Food-mart 2000 database and mushroom dataset is type of dense dataset which is obtained from the FIMI repository. The databases do not provide item utility (external utility) and item count for each transaction (internal utility). Here internal utilities for items are generated randomly ranging from 1 to 10 by transaction utility values Generation code. Table 5 shows the statistical information about these databases, the number of transactions, the number of distinct items, the average number of items in a transaction, and the maximal number of items in the longest transaction(s).

Table5. characteristics of datasets

Dataset	Transactions	Avg. Length Of Transactions	No. Of Items	Type
Foodmart	4141	4.4	1559	Sparse
Mushroom	8124	23	119	Dense
DB_Utility	2880	4	10	Synthetic

C. Time Comparison

The running time of algorithm on foodmart and mushroom datasets with varied K respectively. We compare these execution times with UP Growth, HUI Miner algorithms with optimal thresholds in Fig.3 From the above experiments we observe that TKO algorithm was executes well and take less than second when the K value is low. If the k value increases the execution time also increases for sparse dataset such as foodmart. For dense dataset the algorithm takes almost equal time to execute for various K values and also higher than the hui miner because dense datasets contains large size transactions. Running time was recorded by the “Timestamp” command, and it contains starting time, and ending time of the algorithm. For almost all databases, proposed algorithm performs the best and faster than the existing algorithm.

D. Memory Comparison

Fig. 4 shows memory usage of algorithms. TKO uses list data structure to maintain the utility information and high utility itemsets and also it is one phase algorithm. For that reason TKO takes less memory compare to up-growth and little bit higher or equivalent to the HUI Miner because we maintain and update extra Topk-cl-list for top k high utility itemsets. TKO algorithm generally uses less memory compare to TKU, two phase and state art algorithms.

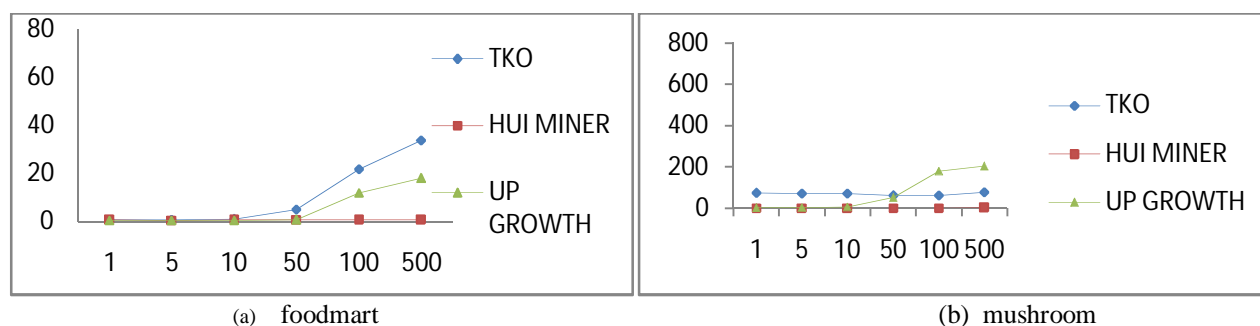


Fig. 3. Runtime of TKO, HUI Miner, UP Growth

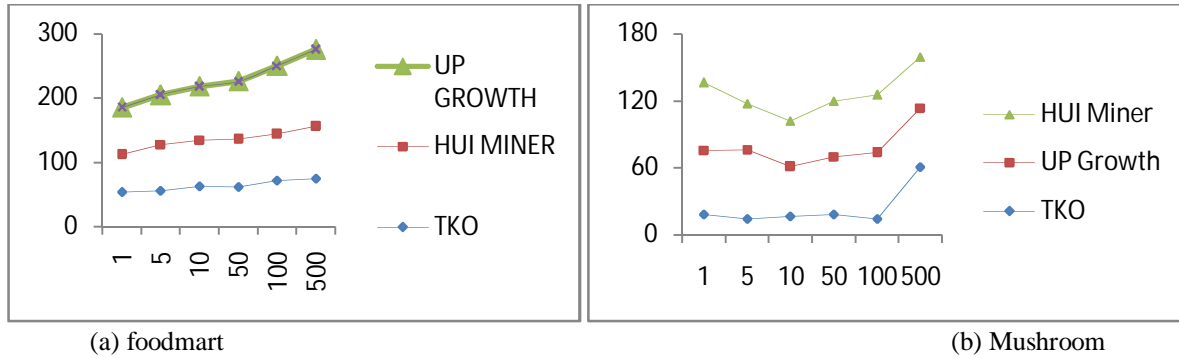


Fig4 Runtime of TKO, HUI Miner, UP Growth

VI. CONCLUSION

In this we are presenting a literature survey on various algorithms used for mining high utility itemsets. we studied the difficulty of setting threshold value in high utility mining. To address this problem the authors has proposed a novel representation of high-utility itemsets named Top k High-Utility Itemsets (Top k HUI). In our framework the top k high utility itemsets are find out by using TKO algorithm which was faster than all top k HUI algorithms. With our framework we find out several k value high utility itemsets on both real and synthetic datasets efficiently. The algorithm was implemented for different types of datasets. Performance factors like time, memory space of TKO algorithm are compared with HUI Miner and UP growth algorithms.

REFERENCES

- [1] R. Agrawal and R. Srikant, T. Imielinski, A. Swami, "Mining association rules between sets of items in large databases", in proceedings of the ACM SIGMOD International Conference on Management of data, pp.207-216, 1993.
- [2] Vincent S. Tseng, Cheng-wei Wu, Philippe Fournier-Viger and Philip S. Yu, "Efficient algorithms for mining Top-K high utility itemsets", in IEEE Transactions on Knowledge and data engineering, vol.28 no.1 January 2016.
- [3] Serin Lee, Jong Soo Park, "Top-K high utility itemset mining based on Utility List Structures", In Proceedings of IEEE International Conference on Data Mining(ICDM), Maebashi, pp. 101 - 108, 2016.
- [4] Shuning Xing, Fangai Liu, Jiwei Wang, Lin Pang, ZhenguoXu, "Utility Pattern Mining Algorithm bases on Improved Utility Pattern Tree", in 8thinternational symposium on computational intelligence and design, pp.258 – 261,2015.
- [5] Junqiang Liu, Benjamin C.M. Fung "Mining High Utility Patterns in One Phase without Generating Candidates" in , IEEE Transaction Knowledge in Data Engineering, vol. 10, no. 12, pp. 1-14, Dec.2015.
- [6] Y. Liu, W. Liao, and A. Choudhary, "A fast high utility itemsets mining algorithm," in Proc. Utility-Based Data Mining Workshop, 2005, pp. 90–99.
- [7] R. Chan, Q. Yang, and Y. Shen, "Mining high-utility itemsets," in Proc. IEEE Int. Conf. Data Mining, 2003, pp. 19–26.
- [8] M. Liu and J. Qu, "Mining high utility itemsets without candidate generation," in Proc. ACM Int. Conf. Inf. Knowl. Manag., 2012, pp. 55–64.
- [9] H. Ryang and U. Yun, "Top-k high utility pattern mining with effective threshold raising Strategies," Knowl.-Based Syst., vol. 76, pp. 109–126, 2015.
- [10] 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.
- [11] V. S. Tseng, C. Wu, B. Shie, and P. S. Yu, "UP-Growth: An efficient algorithm for high utility itemset mining," in Proc. ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2010, pp. 253–262.
- [12] Frequent Itemset Mining Implementations Repository [Online]. Available: <http://fimi.cs.helsinki.fi/>
- [13] FoodMart2000, Microsoft Developer Network (MSDN) [Online]. Available: [https://technet.microsoft.com/en-us/library/aa217032 \(v=sql.80\).aspx](https://technet.microsoft.com/en-us/library/aa217032 (v=sql.80).aspx)



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)