

# Efficiently Updating the Discovered High Average-Utility Itemsets with Transaction Insertion

Jerry Chun-Wei Lin<sup>1,2</sup>, Shifeng Ren<sup>2</sup>, Philippe Fournier-Viger<sup>3</sup>  
Jeng-Shyan Pan<sup>1</sup>, and Tzung-Pei Hong<sup>4,5</sup>

<sup>1</sup>Fujian Provincial Key Laboratory of Big Data Mining and Applications  
Fujian University of Technology, Fujian, China

<sup>2</sup>School of Computer Science and Technology  
Harbin Institute of Technology Shenzhen Graduate School, Shenzhen, China

<sup>3</sup>School of Natural Sciences and Humanities  
Harbin Institute of Technology Shenzhen Graduate School, Shenzhen, China

<sup>4</sup>Department of Computer Science and Engineering  
National University of Kaohsiung, Kaohsiung, Taiwan

<sup>5</sup>Department of Computer Science and Engineering  
National Sun Yat-sen University, Kaohsiung, Taiwan

jerrylin@ieee.org; renshifeng@stmail.hitsz.edu.cn; philfv@hitsz.edu.cn  
jengshengpan@fjut.edu.cn; tphong@nuk.edu.tw

## Abstract

High-utility itemset mining (HUIM) is an extension of frequent-itemset mining (FIM) but considers the unit profit and quantity of items to discover the set of high-utility itemsets (HUIs). Traditionally, the utility of an itemset is the summation of the utilities of the itemset in all the transactions regardless of its length. This approach is, however, inappropriate in real-world applications since the utility of the itemset increases along with the number of items within it. High average-utility itemset mining (HAUIM) was designed to provide more reasonable utility measure by taking the size of the itemset into account. Several algorithms were respectively designed to efficiently and effectively mine the set of high average-utility itemsets (HAUIs). Those algorithms can only handle, however, the static database and unsuitable for the dynamic environment since the size of data is frequently changed in real-life situations. In this paper, an incremental high-average utility pattern mining (IHAUPM) algorithm is presented to handle the incremental database with transaction insertion. The well-known fast updated (FUP) concept in the FIM is modified to adopt the designed algorithm, thus efficiently updating the discovered HAUIs. Based on the designed model for HAUIM with transaction insertion, the proposed IHAUPM algorithm can easily only handle the inserted transactions. Experiments are carried on six datasets and the results showed that the designed algorithm has better performance than the state-of-the-art algorithms performing in the batch manner.

**Keywords:** high average-utility itemsets; FUP; dynamic; insertion; incremental.

# 1 Introduction

Data mining or knowledge discovery can be referred to discover useful, potential, and implicit information from a very large database [1, 2, 13, 38]. The discovered information can thus be used to aid managers or decision makers for making the significant or efficient strategies. Based on different requirements in various applications, varied types of knowledge are respectively presented in different ways. Apriori algorithm [1] was the first algorithm to mine the correlations between the purchase itemsets, which can be referred as association-rule mining (ARM). It uses the level-wise approach to generate-and-test candidates at each level for finding the relationships among the items. The set of frequent itemsets (FIs) is first discovered by the Apriori algorithm based on the pre-defined minimum support threshold. After that, the set of FIs is then combined together, forming the set of association rules (ARs) based on the pre-defined minimum confidence threshold. This approach requires the amounts of computation and easily causes the memory leakage if the size of transactions in the database is too long, the frequent pattern (FP)-growth [13] algorithm was then further designed to mine the FIs based on its designed FP-tree structure. Several algorithms were widely studied to enhance the mining performance of the Apriori-based [2, 6, 36, 38] and tree-based [3, 8, 13] algorithms.

Since the size of database is dynamically changed in real-life situations, such as some transactions could be inserted [7] into the original one; the discovered information is necessary to be updated since some new information may be arisen but some information may be lost. Most algorithms of FIM or ARM can only, however, handle the static database. When the data is changed in the database, the algorithms must be performed in the batch manner, which indicates that the already discovered information is useless and updated database is required to be scanned and examined to mine the required information. To handle the situation of transaction insertion, Cheung et al. proposed the Fast UPdated (FUP) [7] algorithm to adopt the pruning techniques used in in the DHP (Direct Hashing and Pruning) algorithm [37] for updating the discovered information. It divides the original databases and new inserted transactions into four parts and each part is performed by its own algorithm to update the required information. Since some unpromising itemsets in the FUP concept can be eliminated, the updating performance can thus be significantly improved.

In real-world situations, the frequent itemsets in FIM or ARM may only contribute a small portion to the overall profit, while non-frequent ones may contribute a large portion to the profit. High-utility itemset mining (HUIM) [31, 42, 43] was thus proposed to concern more factors such as unit profit and quantity of the items to reveal the high-utility itemsets (HUIs), which is an extension of FIM and widely applied in real-world situations. The utility of an itemset is to sum up the utilities of the itemset in all transactions regardless of its length [32, 42]. It is based on the measure of local transaction utility (quantity) and external utility (unit profit of items) to evaluate whether an itemset is a HUI. Liu et al. [31] then proposed the transaction-weighted utilization (TWU)-model to overestimate the upper value of the designed high-transaction-weighted utilization itemsets (HTWUIs), thus maintaining

1 the transaction-weighted downward closure (TWDC) property for mining the complete HUIs.  
2 To speed up mining performance, Lin et al. then presented a high utility pattern (HUP)-tree  
3 algorithm [21] for mining the set of HUIs. Several algorithms were also extensively developed  
4 to mine the HUIs [5, 17, 39, 43].

5 In HUIM, the itemset within more number of items may increase its utility regardless of its  
6 length. Thus, it is unfair to judge whether the itemset is a HUI with different lengths according  
7 to the same user-specified minimum threshold. To solve this limitation, Hong et al. proposed  
8 a high average-utility itemset mining (HAUIM) [15] to alleviate the effect of the length of  
9 the itemset and identify a better utility effect by considering the size of the itemset than the  
10 original utility standard. The average-utility of an itemset is defined as the total utility of an  
11 itemset divided by its number of items within it. An itemset is called a high average-utility  
12 itemset (HAUI) if its average utility is no less than pre-defined minimum high average-utility  
13 count. Thus, HAUIM is a totally different issue compared to the traditional HUIM since the  
14 new downward closure property, upper-bound model, pruning strategy, as well as the mining  
15 procedure are required to be developed. Lin et al. presented a high average-utility pattern  
16 (HAUP)-tree algorithm [21] to mine the set of HAUIs based on the efficient tree structure. The  
17 projection-based algorithm called PAI [19] was also designed to speed up mining performance.  
18 Lin et al. [30] then designed a average-utility (AU)-list structure to efficiently mine the HAUIs  
19 from a static database, which is the state-of-the-art algorithm for mining HAUIs.

20 Since the previous works of HAUIM can only handle the static database, but in real-life  
21 situations, database size is dynamically changed; for example, some transactions are frequently  
22 inserted into the original database. As we mentioned above, some information may be thus  
23 arisen, but some information may be lost in the updated progress. Although the batch pro-  
24 cessing of the updated database can solve this problem, it causes, however, multiple database  
25 scans and the already discovered information become invalid and useless. In this paper, we  
26 propose an incremental algorithm to efficiently update the discovered HAUIs when some trans-  
27 actions are frequently inserted into the original database. Major contributions are summarized  
28 as follows.

- 29 • We modified FUP concept [7] used in the FIM and ARM for updating the set of the  
30 discovered HAUIs. We have also maintained the correctness and completeness of the  
31 updated HAUIs to ensure that all the required information is completely discovered and  
32 updated.
- 33 • We adopted the efficient HAUP-tree structure to maintain the necessary information for  
34 later mining progress. An incremental IHAUPM algorithm is developed to divide the  
35 original database and new transactions into four cases for maintenance. The itemsets of  
36 each case can be correctly maintained and processed.
- 37 • The designed algorithm is sometimes unnecessary to multiply rescan the original database  
38 for updating the discovered HAUIs compared to the existing state-of-the-art algorithms

1        **running in the batch manner.**

2        The rest of the paper is organized as follows. Related work is reviewed in Section 2. The  
3 preliminaries and problem statement of the designed algorithm are stated in Section 3. The  
4 developed models and designed algorithm are described in Section 4. An illustrated example is  
5 provided in Section 5. Experiments are conducted and evaluated in Section 6. The conclusion  
6 and future work are finally given in Section 7.

## 7    **2   Related Work**

8    In this session, the related works of incremental concept, the high-utility itemset mining  
9 (HUIM), and the high average-utility itemset mining (HAUIM) are respectively reviewed and  
10 discussed.

### 11   **2.1   Incremental Concept**

12    Data mining techniques are used to extract meaningful and implicit patterns or rules from a  
13 very large dataset, and the most common approach is called association-rule mining (ARM).  
14 Agrawal and Srikant first proposed the Apriori algorithm [1] to mine association rules (ARs)  
15 from a set of transactions. It applies the level-wise approach to first mine the frequent itemsets  
16 (FIs) based on the minimum support threshold and performs the combinational mechanism  
17 from the mined FIs to retrieve the ARs based on the minimum confidence threshold. Since  
18 the Apriori algorithm applies the level-wise approach, thus the huge amounts of runtime and  
19 memory usage are necessary. To solve this problem, Han et al. proposed the Frequent-Pattern-  
20 tree structure (FP-tree) structure [13] and FP-growth algorithm for efficiently mining FIs  
21 without generation of candidate itemsets. **To further improve the efficiency of FP-growth,**  
22 **Grahne and Zhu [12] proposed a novel FP-array technique to greatly reduce traversing time of**  
23 **the FP-tree especially in the sparse datasets. Another NFP-tree structure [8] was also designed**  
24 **to reduce the number of tree nodes, as well as the size of Header\_Table by keeping two counts**  
25 **of each tree node. Thus, fewer tree nodes are traversed in the tree structure to mine the FIs.**

26    The above algorithms can only process the static database. In real-life situations, trans-  
27 actions usually grow over time and the mined information must be re-evaluated. For example  
28 in ARM, some new ARs may be generated and some old ones may become invalid when new  
29 transactions are inserted into the original database. **The conventional algorithms are processed**  
30 **in the batch manner to re-process the entirely updated database whether the transactions are**  
31 **inserted, deleted or modified from the original database.** This procedure requires the amounts  
32 of computational cost and waste the discovered knowledge. Cheung et al. thus proposed the  
33 Fast UPdated (FUP) concept [7] to effectively handle the incremental database with transac-  
34 tion insertion. The FUP concept can be used to divide the original database and newly inserted  
35 transactions into four cases and each case is performed by its own procedure to maintain and  
36 update the discovered information. After all items in four cases are respectively performed,

1 the up-to-date information is thus correctly maintained. The FUP concept can be concluded  
 2 with four cases for handling the transaction insertion, shown in Figure 1.

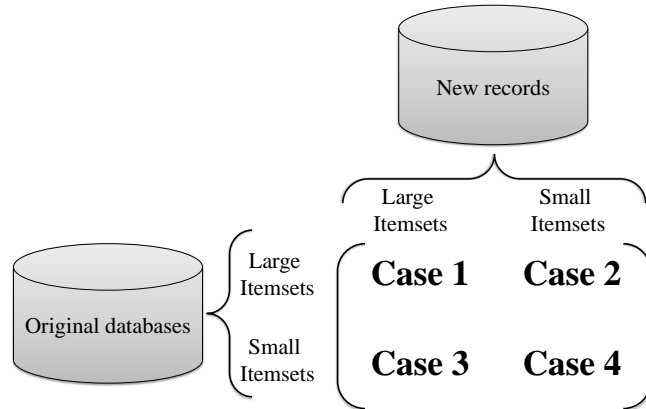


Figure 1: The FUP concept with transaction insertion.

2  
 3 In Figure 1, it can be seen that the itemsets in Case 1 are frequent both in the original  
 4 database and in the new transactions; those itemsets will still be frequent after the database is  
 5 updated. Similarly, the itemsets in Case 4 will still remain small after the database is updated.  
 6 Thus, the results in Cases 1 and 4 will not affect the final frequent itemsets. The itemsets in  
 7 Cases 2 and 3 may, however, be removed from the existing frequent itemsets or some frequent  
 8 itemsets may be arisen.

9 Based on the FUP concept and the FP-tree structure, Hong et al. presented an incremental  
 10 algorithm [14] for fast updating the frequent patterns. Although the FP-tree structure can be  
 11 easily maintained in the memory, this algorithm can only be used handle the binary database.  
 12 However, several realistic factors such as weight, interestingness, or unit profit of items can  
 13 seriously influence the discovered patterns. It is thus a critical issue to develop more realistic  
 14 algorithm to handle the incremental database.

## 15 2.2 High Utility Itemset Mining

16 High utility itemset mining (HAUIM) [5, 17, 39, 42, 43], an extension of high utility itemset  
 17 mining (HUIM), is based on the measurement of local transaction utility and external utility.  
 18 The local utility can be considered as the quantity of the items and the external utility can be  
 19 considered as the unit profit of the items in the database. The utility of an itemset is to sum  
 20 up the utilities of the items in the whole database. An itemset is concerned as the high-utility  
 21 itemset (HUI) if its utility is no less than the minimum high-utility threshold.

22 Since the measure of HUIM is totally different than the measure used in the ARM, the  
 23 downward closure (DC) property cannot be directly applied into the HUIM. In order to discover  
 24 the complete and correct high utility itemsets, Liu et al. [32] first proposed a transaction-  
 25 weighted downward closure property (TWDC) of itemsets to limit the search space and then  
 26 designed the Two-Phase algorithm to level-wisely mine the HUIs. This approach requires

1 to generate-and-test candidates level-by-level, which suffers the problem of multiple database  
2 scans. Besides, all the potential itemsets at each level are necessary kept for finding the actual  
3 HUIs. Lin et al. then presented the high-utility pattern (HUP)-tree structure [21] to keep  
4 the necessary information in the main memory for further mining process. Tseng et al. then  
5 presented a compact UP-tree structure, the UP-Growth and UP-Growth+ algorithms and five  
6 pruning strategies for efficiently mining the HUIs [39]. Although the UP-Growth+ outperforms  
7 the Two-Phase algorithm, the memory leakage is still a problem in a large scale dataset. Liu  
8 and Qu then designed a novel list structure and developed a HUI-Miner algorithm [33] to mine  
9 the HUIs. The HUI-Miner is unnecessary to generate the amounts of candidates but still can  
10 discover the complete set of HUIs. Liu et al. then designed an algorithm [35] to mine the  
11 next itemsets using the suffix itemset for later combination. The reverse enumeration tree was  
12 designed as the search space and a liner CAUL data structure was used to keep the necessary  
13 information to facilitate the mining performance based on the designed d2HUP algorithm.  
14 Krishnamoorthy [17] proposed several pruning strategies to reduce the cost for constructing  
15 the utility-list structure to improve the performance of HUI-Miner. Similar as the mining  
16 limitation of FIs or ARs, it necessitates to frequently update the discovered information in the  
17 dynamic database. Ahmed et al. [4] then designed an incremental and interactive algorithms  
18 for mining HUIs. Lin et al. also presented an incremental algorithm [20] to update the  
19 discovered HUIs by employing the Two-Phase model and the FUP concept. Yun and Ryang  
20 designed a novel HUPID-Tree tree [44] and the HUPID-Growth algorithm to mine the HUIs.  
21 The over-estimated utility of the items is reduced and each node in the tree structure keeps  
22 more information for later mining process. A Header\_Table keeps not only the utility value of  
23 each item but also its frequency.

24 Traditional algorithms of HUIM have to cope the exponential problem of the search space,  
25 Lin et al. then designed the PSO-based algorithm [28] to reduce the search space for mining  
26 the HUIs in a limit of time. To find the first  $k$  HUIs as the interesting patterns, Tseng et  
27 al. designed a top- $k$  [40] approach for mining the HUIs. Duong et al. [10] then presented  
28 three pruning strategies to speed up the mining performance of [40]. In real-life applications,  
29 the HUIM can be applied in different situations. For example, in the basket analysis, the  
30 profit of items may be changed during the sales promotion. Lin et al. [25] then presented an  
31 one-phase HUI-DMiner algorithm and two pruning strategies to handle this issue for mining  
32 the precisely HUIs. Considering the positive and negative unit profit of the products during  
33 the promotion season [9], Lin et al. presented the FHN algorithm [24] and developed several  
34 pruning strategies to efficiently mine the HUIs. Based on the uncertainty constraint, for  
35 example, in the environment of wireless sensor network, the received data may consist of the  
36 inaccurate or uncertain information. Lin et al. then designed the algorithm [22] to mine the  
37 HUIs over the uncertain database by considering the excepted support model. The above  
38 algorithms consider that all items use a single minimum utility threshold to identify HUIs,  
39 however, it is inadequate and unfair to find the information with a single threshold since

1 all items are different and importance of each item should be concerned. Lin et al. then  
2 presented an algorithm to mine the HUIs under varied minimum high-utility thresholds [27].  
3 This approach finds more reasonable HUIs than that of the traditional approaches.

4 The issue of privacy-preserving data mining (PPDM) has also emerging as a critical topic  
5 in recent decades, which can also be extended to the HUIM since some profitable products are  
6 considered as the confidential information. Lin et al. then surveyed the works of PPUM [23] and  
7 designed the algorithm [26] to hide the confidential HUIs and minimize three side effects in the  
8 sanitization process. Other issues related to HUIM are also developed in progress [29, 41, 18].

### 9 **2.3 High Average-Utility Itemset Mining**

10 In the traditional HUIM, the utility of an itemset is the summation of the utilities of the  
11 itemset in all the transactions regardless of its length. Thus, the utility of an itemset in a  
12 transaction increases along with the increasing of its length. This process will result that a  
13 longer itemset comes up with higher utility. Using the same minimum high utility threshold  
14 to identify the itemsets with different lengths is thus unfair. The high average-utility itemset  
15 mining (HAUIM) [15] was designed to alleviate the effect of the length of itemsets and identify  
16 really good utility itemsets. The average-utility measure was designed to reveal a better utility  
17 effect of combining several items than the original utility measure. The average-utility of an  
18 itemset is defined by calculating the total utility of an itemset divided by the number of items  
19 within the itemset. Several algorithms were extensively studied to mine the high average-utility  
20 itemsets (HAUIs). The Apriori-like TPAU [15] was developed to level-wisely mine the HAUIs  
21 under the generate-and-test approach. This approach requires the amounts of computation  
22 to find the HAUIs. Lan et al. then provided a projection-based approach [19] to tight the  
23 upper-bound of the promising candidates to find the HAUIs.

24 To speed up mining performance, Lin et al. proposed a high average-utility pattern  
25 (HAUP)-tree algorithm [21] for integrating the TPAU and FP-tree structure to firstly con-  
26 struct a condensed tree structure then efficiently mining HAUIs. To improve the mining  
27 performance, the HAUI-Miner [30] algorithm was presented to adopt the HUI-Miner structure  
28 for efficiently discovering the HAUIs without candidate generation. It uses the average-utility  
29 (AU)-list to keep the necessary information for the later join operation and mining progress.  
30 Although the above approaches can be used to efficiently mine the HAUIs, all of them fail to  
31 efficiently update the discovered HAUIs when the size of database is frequently changed. For  
32 example, when the transactions are frequently inserted into the original database, the above  
33 approaches have to process the entire database in the batch manner, which indicates that the  
34 updated database is scanned again to mine the updated information even though the size of  
35 inserted transactions is quite small. Besides, the computational cost is high since the set of the  
36 discovered HAUIs becomes old information; some HAUIs become invalid and some itemsets  
37 may arise as the new HAUIs. It is thus an important issue to efficiently updated the discovered  
1 information especially in the big data era under the dynamic environment.

## 3 Preliminaries and Problem Statement

### 3.1 Preliminaries

Let  $I = \{i_1, i_2, \dots, i_m\}$  be a finite set of  $m$  distinct items. A quantitative database is a set of transactions  $D = \{T_1, T_2, \dots, T_n\}$ , where each transaction  $T_q \in D$  ( $1 \leq q \leq m$ ) is a subset of  $I$  and has a unique identifier  $q$ , called its *TID*. Besides, each item  $i_j$  in a transaction  $T_q$  has a purchase quantity denoted as  $q(i_j, T_q)$ . A profit table *PTable* indicates the unit profit value of each item in the database as  $PTable = \{p(i_1), p(i_2), \dots, p(i_m)\}$ , where profit values are positive integers. A set of  $k$  distinct items  $X = \{i_1, i_2, \dots, i_k\}$  such that  $X \subseteq I$  is said to be a  $k$ -itemset, where  $k$  is the length of the itemset. An itemset  $X$  is said to be contained in a transaction  $T_q$  if  $X \subseteq T_q$ . A minimum high average-utility threshold  $\delta$  is set according to the user's preference (a positive integer). A quantitative database is shown in Table 1, which will be used as running example for the rest of this paper. This database contains five transactions and seven distinct items, denoted with letters from (*a*) to (*g*). The profit table indicates the unit profit of each item appearing in the database, and is shown in Table 2.

Table 1: A quantitative database.

TID	Items with their quantities
$T_1$	$a:5, b:2, c:3, d:2$
$T_2$	$b:5, d:4, e:2, g:12$
$T_3$	$a:1, c:5, f:4$
$T_4$	$c:2, d:3, e:3$
$T_5$	$a:3, d:2, f:6$

Table 2: A profit table.

Item	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>
Profit	4	6	5	9	12	8	3

**Definition 1.** The average-utility of an item  $i_j$  in a transaction  $T_q$  is denoted as  $au(i_j, T_q)$ , and defined as:

$$au(i_j, T_q) = \frac{q(i_j, T_q) \times p(i_j)}{1}, \quad (1)$$

where  $q(i_j, T_q)$  is the quantity of  $i_j$  in  $T_q$ , and  $p(i_j)$  is the unit profit value of  $i_j$ .

For example, the average-utility of items (*b*), (*c*), and (*d*) in  $T_1$  is respectively calculated as  $au(b, T_1) (= \frac{2 \times 6}{1}) (= 12)$ ,  $au(c, T_1) (= \frac{3 \times 5}{1}) (= 15)$ , and  $au(d, T_1) (= \frac{2 \times 9}{1}) (= 18)$ .

**Definition 2.** The average-utility of a  $k$ -itemset  $X$  in a transaction  $T_q$  is denoted as  $au(X, T_q)$ , and defined as:

$$au(X, T_q) = \frac{\sum_{i_j \subseteq X \wedge X \subseteq T_q} q(i_j, T_q) \times p(i_j)}{|X|} = \frac{\sum_{i_j \subseteq X \wedge X \subseteq T_q} q(i_j, T_q) \times p(i_j)}{k}, \quad (2)$$



2 where  $k$  is the number of items in  $X$ .

3 For example, the average-utility of itemsets  $(ac)$  and  $(acd)$  in  $T_1$  are respectively calculated  
 4 as  $au(ac) = \frac{5 \times 4 + 3 \times 5}{2} = (17.5)$  and  $au(acd) = \frac{5 \times 4 + 3 \times 5 + 2 \times 9}{3} (= 17.6)$ .

5 **Definition 3.** The average-utility of an itemset  $X$  in  $D$  is denoted as  $au(X)$ , and is defined  
 6 as:

$$au(X) = \sum_{X \subseteq T_q \wedge T_q \in D} au(X, T_q). \quad (3)$$

7 For example, the average-utility of an itemset  $(ac)$  in the database depicted in Table 1 is  
 8 calculated as  $au(ac) = au(ac, T_1) + au(ac, T_3) = \frac{35}{2} + \frac{29}{2} (= 32)$ .

9 **Definition 4.** The transaction utility of a transaction  $T_q$  is denoted as  $tu(T_q)$ , and defined as:

$$tu(T_q) = \sum_{i_j \subseteq T_q} u(i_j, T_q). \quad (4)$$

10 For example, the utilities of transactions in Table 1 are respectively calculated as  $tu(T_1) =$   
 11  $20 + 12 + 15 + 18 (= 65)$ ,  $tu(T_2) (= 126)$ ,  $tu(T_3) (= 61)$ ,  $tu(T_4) (= 73)$ , and  $tu(T_5) (= 78)$ .

12 **Definition 5.** The total utility of a database  $D$  is denoted as  $TU^D$ , and defined as the sum  
 13 of all transaction utilities, that is:

$$TU^D = \sum_{T_q \in D} tu(T_q). \quad (5)$$

14 For example, the total utility in the running example of Table 1 is calculated as  $TU^D =$   
 15  $65 + 126 + 61 + 73 + 78 (= 403)$ .

16 To maintain the downward closure (DC) property, the average-utility upper bound ( $auub$ )  
 17 property [15] was designed to over-estimate the utility of the promising itemsets, thus holding  
 18 the DC property of HAUIM.

19 **Definition 6** (transaction-maximum utility). The transaction-maximum utility of a transac-  
 20 tion  $T_q$  denoted as  $tmu(T_q)$  and defined as:

$$tmu(T_q) = \max\{u(i_j) | i_j \subseteq X \wedge X \subseteq T_q\}. \quad (6)$$

21 For example in Table 1, the transaction-maximum utilities of transactions  $T_1$  to  $T_5$  are  
 22 respectively calculated as  $tmu(T_1) (= 20)$ ,  $tmu(T_2) (= 36)$ ,  $tmu(T_3) (= 32)$ ,  $tmu(T_4) (= 36)$  and  
 23  $tmu(T_5) (= 48)$ .

24 **Definition 7** (Average-utility upper-bound,  $auub$  property). The average utility upper bound  
 1 of an itemset  $X$  is denoted as  $auub(X)$  and defined as:

$$auub(X) = \sum_{X \subseteq T_q \wedge T_q \in D} tmu(T_q), \quad (7)$$

2 where  $tmu(T_q)$  is the maximum utility of transaction  $T_q$  such that  $i_j \subseteq X \wedge X \subseteq T_q$ .

3 For example in Table 1, the *auub* values of (*a*) and (*ac*) are respectively calculated as  
 4  $auub(a) = (20 + 32 + 48)(= 100)$  and  $auub(ac) (= 20 + 32)(= 52)$ .

5 Based on the *auub* property, the DC property can be held, thus reducing the size of un-  
 6 promising candidates. Thus, the set of high average-utility upper bound itemsets (HAUUBIs)  
 7 can be found to ensure the **completeness** and **correctness** of the *auub* property.

8 **Definition 8** (High average-utility upper bound itemset, *HAUUBI*). An itemset  $X$  is a HAU-  
 9 UBI if it satisfies the condition as:

$$HAUUBI \leftarrow \{X | auub(X) \geq TU^D \times \delta\} \quad (8)$$

10 Assume that the minimum average-utility threshold is set as 14.4%. For example in Table  
 11 1, (*a*) is a HAUUBI since its  $auub(a)(= 100 > 403 \times 14.4\% = 58.023)$  but  $auub(ac)$  is not a  
 12 HAUUBI since  $auub(ac)(= 52 < 58.023)$ .

### 13 3.2 Problem Statement

14 Assume that the newly inserted database named (*d*) is inserted into the original database (*D*),  
 15 and its given example is shown in Table 3. The unit profit of the items of Table 3 is the same  
 16 as the original database shown in Table 2.

Table 3: Newly inserted transactions.

TID	Items
$T_6$	$a:4, b:6, c:4, d:2$
$T_7$	$a:5, c:3, e:1, f:2, g:6$

17 The problem of incremental HAUIM with transaction insertion is to efficiently update the  
 18 set of the discovered HAUIs. Thus, an itemset  $X$  is a HAUI in the updated database ( $D + d$ )  
 19 if its utility is no less than the updated minimum average-utility count. The set of HAUIs is  
 20 thus formally defined as:

$$HAUIs \leftarrow \{X | au(X) \geq (TU^D + TU^d) \times \delta\}, \quad (9)$$

21 where  $TU^d$  is the total utility in the inserted database (*d*), and  $\delta$  is the minimum high average-  
 22 utility threshold, which can be specified by user's preference.

## 23 4 Proposed Incremental Algorithm

24 In the past, Lin et al. proposed a HAUP-growth mining algorithm to mine th HAUIs from  
 1 the constructed HAUP-tree [21]. Each node in the HAUP-tree keeps the quantities of the

2 prefix items of the processed node in the branch, which can speed up mining performance  
3 without building the enormous conditional trees as the FP-growth does. The HAUP-growth  
4 algorithm is designed, however, for mining the HAUIs in the batch manner. Thus, when some  
5 transactions are inserted into the original database, the HAUP-growth algorithm needs to  
6 process the updated database again to mine the updated information since some rules may  
7 thus be arisen but some rules may be lost in the updated database. Several existing algorithms  
8 can only handle the static databases. In this paper, we present an incremental high-average-  
9 utility pattern mining (IHAUPM) algorithm to handle the updated databases for mining the  
10 HAUIs with transaction insertion. The designed IHAUPM algorithm has two phases. In the  
11 first phase, the built HAUP-tree from the original database is maintained and updated by  
12 the inserted transactions. This process divides the itemsets into four cases according to the  
13 modified FUP concept used in the FIM but the designed algorithm considers the *auub* property  
14 of high average-utility upper bound itemsets (HAUUBIs) for the incremental process. Each  
15 case is performed by the designed procedure to maintain the HAUP-tree. Based on the FUP  
16 concept and *auub* property of HAUUBIs, the four cases is presented in Figure 2.

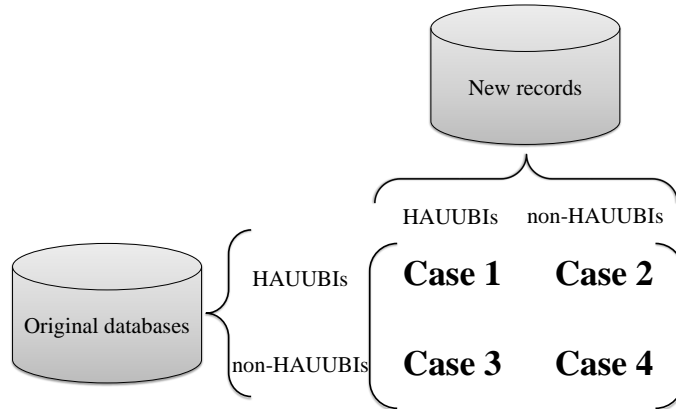


Figure 2: Proposed four cases model of incremental HAUM.

17 The four cases model of the incremental HUIM can be described below.

- 18 • **Case 1:** An itemset is a HAUUBI both in the original database and in the inserted  
19 transactions. For an itemset  $X$  in case 1, we can obtain that:

$$auub(X)^U \geq (TU^D + TU^d) \times \delta. \quad (10)$$

- 20 • **Case 2:** An itemset is a HAUUBI in the original database, but is a non-HAUUBI in the  
21 inserted transactions. For an itemset  $X$  in case 2, we can obtain two situations without  
1 rescanning the original database  $D$  as:

$$\begin{cases} auub(X)^U \geq (TU^D + TU^d) \times \delta & \text{if } auub(X)^d \geq TU^d \times \delta; \\ auub(X)^U < (TU^D + TU^d) \times \delta & \text{if } auu(X)^d < TU^d \times \delta. \end{cases} \quad (11)$$

- 2 • **Case 3:** An itemset is a non-HAUUBI in the original database, but is a HAUUBI in  
 3 the inserted transactions. For an itemset  $X$  in case 3, we can obtain two situations by  
 4 rescanning original database  $D$  to get the  $awub(X)^D$  as:

$$\begin{cases} awub(X)^D + awub(X)^d = awub(X)^U \geq (TU^D + TU^d) \times \delta; \\ awub(X)^D + awub(X)^d = awub(X)^U < (TU^D + TU^d) \times \delta. \end{cases} \quad (12)$$

- 5 • **Case 4:** An itemset is a non-HAUUBI both in the original database and in the inserted  
 6 transactions. For an itemset  $X$  in case 4, we can obtain that:

$$awub(X)^U < (TU^D + TU^d) \times \delta. \quad (13)$$

7 From the above four cases, a summary of the final results in four cases is given in Table 4.

Table 4: Four cases of the incremental HAUIM.

Case: Original - New	Results
Case 1: HAUUBI - HAUUBI	Always frequent
Case 2: HAUUBI - non-HAUUBI	Decided from the existing information
Case 3: non-HAUUBI - HAUUBI	Decided by rescanning the original database
Case 4: non-HAUUBI - non-HAUUBI	Always non-HAUUBI

8 To ensure the **correctness** and **completeness** of the discovered HAUIs by the designed  
 9 IHAUPM, the lemmas and proofs of four cases are given below.

10 **Lemma 1** (Itemsets in Case 1). An itemset  $X$  remains as the HAUUBI in the updated  
 11 database ( $U$ ) if its  $awub(X)$  is no less than the updated minimum average-utility count both  
 12 in  $D$ (= original database) and  $d$ (= inserted transactions).

13 *Proof.* Since  $awub(X)^D \geq TU^D \times \delta$  and  $awub(X)^d \geq TU^d \times \delta$ , thus  $awub(X)^U = awub(X)^D +$   
 14  $awub(X)^d \geq TU^D \times \delta + TU^d \times \delta = (TU^D + TU^d) \times \delta \Rightarrow awub(X)^U \geq (TU^D + TU^d) \times \delta.$

15 □

16 **Corollary 1.** *An itemset  $X$  still remains a HAUUBI in the updated database if it belongs to*  
 17 **Case 1**; *the itemsets in Case 1 do not affect the number of the final updated itemsets, but*  
 18 *affect the actual average-utilities of them.*

19 **Lemma 2** (Itemsets in Case 2). If an itemset  $X$  is a HAUUBI (appearing in HAUI-tree) in  $D$   
 20 but small in the inserted transactions  $d$ , two possible situations exist as: 1) it still remains a  
 21 HAUUBI in  $U$ (= updated database) or 2) it becomes a non-HAUUBI in  $U$  since its  $awub(X)$   
 22 is less than the updated minimum average-utility count.

23 *Proof.* Since  $awub(X)^D \geq TU^D \times \delta$ , the itemset  $X$  in  $d$  has two conditions as  $awub(X)^d \geq$   
 24  $TU^d \times \delta$  or  $awub(X)^d < TU^d \times \delta$ , thus:

1 (1) If  $auub(X)^d \geq TU^d \times \delta$ , we can obtain that  $auub(X)^U = auub(X)^D + auub(X)^d \geq (TU^D +$   
2  $TU^d) \times \delta$ .

3 (2) If  $auub(X)^d < TU^d \times \delta$ , we can obtain that  $auub(X)^U = auub(X)^D + auub(X)^d < (TU^D +$   
4  $TU^d) \times \delta$ .

5 □

6 Thus, for the itemsets in case 2, some itemsets still remain as the HAUUBI but some  
7 itemsets may become the non-HAUUBI in the updated database. For the itemsets in case 2, it  
8 is easier to maintain the updated information of the itemsets since those itemsets exist in the  
9 *Header\_Table* of the HAUP-tree. Besides, after one database scan of the inserted transactions,  
10 the *auub* values of all itemsets in new transactions are discovered.

11 **Corollary 2.** *An itemset  $X$  may remain as a HAUUBI or become a non-HAUUBI in the*  
12 *updated database if it belongs to **Case 2**; the items in **Case 2** may be removed, and thus affect*  
13 *the final rules.*

14 **Lemma 3** (Itemsets in Case 3). If an itemset  $X$  is a small itemset (not existing in the  
15 *Header\_Table* of HAUP-tree) but its  $auub(X)^d$  in  $d$  is no less than the minimum high average-  
16 utility count of  $d$ , two possible conditions may exist when the database is updated as: 1) it  
17 becomes a HAUUBI or 2) it remains the non-HAUUBI in  $U$ .

18 *Proof.* If an itemset  $X$  does not exist in the *index\_table* of HAUP-tree, the  $auub(X)^D <$   
19  $TU^D \times \delta$ .

20 Since  $auub(X)^d \geq TU^d \times \delta$ ,  $D$  is then scanned to find the  $auub(X)^D$ .

21 For the updated database  $U$ , we can obtain that  $auub(X)^U = auub(X)^D + auub(X)^d \geq$   
22  $auub(X)^D + TU^d \times \delta$ . Thus,  $auub(X)^U \geq (TU^D + TU^d) \times \delta$  or  $auub(X)^U < (TU^D + TU^d) \times \delta$ . □

23 **Corollary 3.** *An itemset  $X$  may become a HAUUBI or remain as a non-HAUUBI in the*  
24 *updated database if it belongs to **Case 3**; the itemsets in **Case 3** affect the final rules since*  
25 *some itemsets may become the HAUIs.*

26 The itemsets in Case 4 are unnecessary processed since they will never be the HAUUBIs,  
27 neither as the HAUIs. The lemma is given below.

28 **Lemma 4** (Itemsets in Case 4). An itemset  $X$  remains a non-HAUUBI in  $U$  if it is a non-  
29 HAUUBI both in  $D$  and  $d$ .

30 *Proof.* An itemset  $X$  is a non-HAUUBI in both  $D$  and  $d$  if its  $auub(X)^D < TU^D \times \delta$  and  
31  $auub(X)^d < TU^d \times \delta$ . Thus,

32  $auub(X)^U = auub(X)^D + auub(X)^d < TU^D \times \delta + TU^d \times \delta (= TU^D + TU^d) \times \delta$ . □

33 **Corollary 4.** *An itemset  $X$  remains a non-HAUUBI in the updated database if it belongs to*  
34 ***Case 4**. The itemsets in **Case 4** do not affect the final results.*

1 According to the above lemmas and corollaries, the **correctness** and **completeness** of  
 2 the proposed algorithm is obtained. In the designed model, the necessary information (the  
 3 set of HAUUBIs) from the original database is kept in the main memory as a compressed  
 4 tree structure for later mining progress. In real-world situation, the size of newly inserted  
 5 transactions is quite small compared with the original database. Thus, it is easily to scan  
 6 the newly inserted transactions to obtain the required information for later updating progress.  
 7 When the transactions are inserted into the original database, we only need to scan the inserted  
 8 parts to maintain and update the discovered information in the main memory. Although the  
 9 original database is sometimes required to be scanned again while some itemsets belong to  
 10 Case 3, this situation is rarely happened since the size difference is quite huge between the  
 11 original database and newly inserted transactions. Thus, the multiple database scan can be  
 12 successfully avoided. With the proposed algorithm for handling three cases (Case 4 can be  
 13 directly ignored), the HAUP-tree can be efficiently and correctly updated and the required  
 14 HAUIs can be successfully obtained by utilizing HAUP-growth method. The reason is that  
 15 the *auub* property is used as the upper bound to maintain the downward closure property of  
 16 the HAUUBIs. Thus, if an item/set is not a HAUUBI, its superset will not be the HAUUBI,  
 17 neither as the HAUI. Therefore, the results of the proposed IHAUPM algorithm are correct  
 18 and complete. Details of the designed incremental approach are described in Algorithm 1.

19 First, the new inserted database is scanned to find the transaction-maximum utility (*tmu*)  
 20 of each transaction (Lines 1 to 2), and the *auub* of each item in  $d$  is then calculated (Lines 3 to  
 21 4). After that, the total utility of  $d$  is then calculated, as well as the summed up total utility  
 22 of the updated database (Lines 5 to 6). The items in four cases are respectively processed.  
 23 Notice that the itemsets in case 4 can be ignored since they cannot be the HAUIs according  
 24 to our proofs. If the updated *auub* of each itemset is greater than or equal to the updated  
 25 minimum high average-utility count, the itemset belong to case 1; the *auub* of each itemset in  
 26 the *Header\_Table* of the HAUP-tree is then updated, and the satisfied itemset is then put into  
 27 the set of *insert\_set* for later processing (Lines 9 to 12). If the updated *auub* of the itemset  
 28 is less than the updated minimum high average-utility count, the itemset is directly removed  
 29 from the HAUP-tree, and the child nodes of the processed node in the HAUP-tree are directly  
 30 linked to the parent node of it (Lines 14 to 16). Otherwise, if the *auub* of the itemset in  $d$  is  
 31 no less than the minimum high average-utility count and the processed itemset does not exist  
 32 in the HAUP-tree, the original database is scanned again to find the actual *auub* of the itemset  
 33 and put it in the set of *rescan\_set* for updating the corresponding branch in the HAUP-tree  
 34 (Lines 18 to 21). After that, the itemsets in the *rescan\_set* are sorted by their *auub-descending*  
 35 order (Line 22) and they are respectively inserted into the end of *Header\_Table* (Line 23). The  
 36 **UpdateBranch** function is then performed on two different sets (Lines 24 to 25) to update  
 37 the HAUP-tree. Details are shown in Algorithm 2.

38 For the processed itemsets in the set (*rescan\_set* or *insert\_set*), if the itemset exists in the  
 39 branch of the HAUP-tree (Lines 3 to 5), the *auub* value of the node is then updated by addition

---

**Algorithm 1:** Incremental high-average utility pattern mining (IHAUPM)

---

**Input:**  $D$ , a quantitative database;  $PTable$ , a profit table; the total utility  $TU^D$  in  $D$ ;  $\delta$ , the minimum average-utility threshold;  $htree$ , the built HAUP-tree from  $D$ ;  $d$ , a set of inserted transactions.

**Output:** The updated HAUP-tree.

```
1 for each  $T_q \in d$  do
2   └ calculate  $tmu(T_q)$ ;
3 for each  $i_j$  in  $d$  do
4   └ calculate  $auub(i_j)$ ;
5 calculate total utility in  $d$  as  $TU^d$ ;
6 calculate updated total utility  $TU^U = TU^D + TU^d$ ;
7 for each  $i_j \in index\_table$  of  $htree$  do
8   └ for each  $i_j$  in  $d$  do
9     └ if  $auub(i_j)^D + auub(i_j)^d \geq TU^U \times \delta$  then
10       └  $auub(i_j)^U = auub(i_j)^D + auub(i_j)^d$ ;
11       └ update  $auub(i_j) \in index\_table$  as  $auub(i_j)^U$ ;
12       └  $insert\_set \leftarrow insert\_set \cup i_j$ ;
13     └ else if  $auub(i_j)^D + auub(i_j)^d < TU^U \times \delta$  then
14       └  $parent\_node(i_j).child \leftarrow child\_node(i_j).parent$ ;
15       └ remove  $i_j$  from  $htree$ ;
16       └ remove  $i_j$  from  $index\_table$ ;
17     └ else
18       └ if  $auub(i_j)^d \geq TU^d \times \delta \wedge i_j \notin index\_table$  then
19         └ rescan  $D$  to obtain  $auub(i_j)^D$ ;
20         └ if  $auub(i_j)^D + auub(i_j)^d \geq TU^U \times \delta$  then
21           └  $rescan\_set \leftarrow rescan\_set \cup i_j$ ;
22 sort  $rescan\_set$  in current  $auub$ -descending order;
23 insert items in  $rescan\_set$  to the end of  $index\_table$ ;
24 UPdateBranch( $rescan\_set$ );
25 UPdateBranch( $insert\_set$ );
```

---

---

**Algorithm 2:** UPdateBranch(set)

---

**Input:** the itemsets in the *set*.

**Output:** The updated *htree*.

```
1 for each  $i_j$  in the set do
2   for each  $i_j \subseteq T_q \wedge T_q \in D \parallel T_q \in d$  do
3     if  $i_j \in htree.branch \wedge i_j \in insert\_set$  then
4        $node(i_j).auub+ = tmu(T_q)$ ;
5        $node(i_j).quan\_Ary += T_q.quantity[1..j]$ ;
6     else
7       if  $i_j \in insert\_set \wedge i_j \in rescan\_set$  then
8         find the  $i_j.position$  in  $T_q \in D$ ;
9          $node(i_j).auub = tmu(T_q)$ ;
10         $node(i_j).quan\_Ary = T_q.quantity[1..j]$ ;
```

---

1 (Line 4), as well as the attached quantity array of the node (Line 5); otherwise, the position of  
2 the itemset in the transaction is then found and inserted into its corresponding position of the  
3 branch in the HAUP-tree (Line 8). The *auub* value of the node is then added into the node  
4 (Line 9), and its quantity array is then corresponding updated (Line 10).

5 In the second phase, the updated HAUP-tree is then processed by the constructed *Header\_Table*  
6 to respectively mine the HAUIs using the HAUP-growth approach [21]. This is a trivial task  
7 since the complete information is kept in the HAUP-tree structure. After that, the set of  
8 HAUIs is then maintained and updated.

## 9 5 An Illustrated Example

10 In this section, an example is given to demonstrate the proposed IHAUPM algorithm. Suppose  
11 that the original quantitative database and its profit table were respectively shown in Tables  
12 1 and 2. The minimum high average-utility threshold is initially set as 14.4%, which can be  
13 defined by users' preference. The HAUP-tree was first built from Table 1 to keep the necessary  
14 information for later mining process. From the given example, the constructed HAUP-tree is  
15 then shown in Figure 3, and the total utility of Table 1 was calculated as 403.

16 We also assume that some transactions are then inserted into the original database, and the  
17 inserted transactions were shown in Table 3. In this example, the total utility of the inserted  
18 database was calculated as 171. Thus, in the updated database, the updated total utility is  
19  $(403 + 171)(= 574)$ , and the minimum high average-utility count is then calculated as  $(574 \times$   
20  $14.4\%)(= 82.656)$ .

21 First, the *auub* values of all items in Table 3 are respectively calculated and the results are  
22 shown in Table 5.

23 Based on the modified FUP model, the *auub* of each item in the inserted transactions is



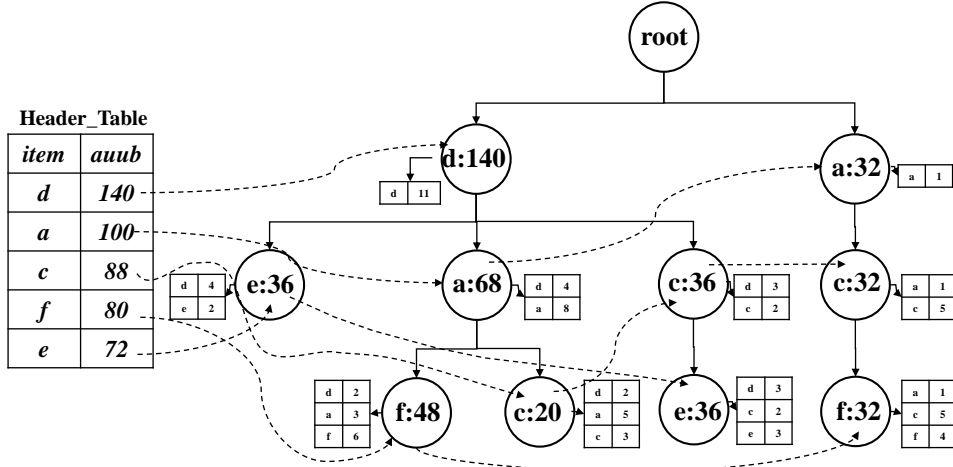


Figure 3: Constructed HAUP-tree.

Table 5: The  $auub$  values of all items in the inserted database.

Item	a	b	c	d	e	f	g
<b>auub</b>	56	36	56	36	20	20	20

1 checked against the minimum high average-utility count in the inserted transactions. In this  
 2 example, the items ( $a$ ), ( $c$ ) and ( $d$ ) belong to case 1; the items ( $e$ ) and ( $f$ ) belong to case  
 3 2, and the item ( $b$ ) belongs to case 3. The remaining item ( $g$ ) belongs to case 4 and can be  
 4 directly ignored since it certainly cannot be a HAUI.

5 First, the items in case 1 are respectively processed to update their  $auub$  values in the  
 6 *Header\_Table* of HAUP-tree. For example, the  $auub$  value of ( $a$ ) in the HAUP-tree was  
 7  $auub(a)^D (= 100)$ , and its  $auub$  value in the inserted database was calculated as  $auub(a)^d (= 56)$ ;  
 8 the updated  $auub$  value of ( $a$ ) is summed up as  $auub(a)^U (= 100 + 56) (= 156 > 82.656)$ .  
 9 The item ( $a$ ) is then put into the set of *insert\_set* and its  $auub$  value in the *Header\_Table* of  
 10 HAUP-tree is then updated. This procedure is then processed for the items ( $c$ ) and ( $d$ ). After  
 11 that, the *insert\_set* = { $a, c, d$ }. The items in case 2 are then processed. In this example, the  
 12  $auub$  values of ( $e$ ) and ( $f$ ) are respectively calculated and the results are  $auub(e)^U (= 72 +$   
 13  $20) (= 92)$  and  $auub(f)^U (= 80 + 20) (= 100)$ . In this example, they remains as the potential  
 14 HAUIs and the  $auub$  values of them in the *Header\_Table* are respectively updated; ( $e$ ) and ( $f$ )  
 15 are both put into the set of *insert\_set*. After that, the items in case 3 are then processed. In  
 16 this example, only an item ( $b$ ) exists in case 3; the original database is required to be scanned  
 17 to find the  $auub$  value of ( $b$ ) in the original database, and in this example, the  $auub(b)^U (= 92)$ ,  
 18 which is larger than the minimum high average-utility count; ( $b$ ) is then put into the two  
 19 sets of *insert\_set* and *rescan\_set*. Thus, *insert\_set* = { $a, b, c, d, e, f$ }, and *rescan\_set* = { $b$ }.

20 After that, the initially built HAUP-tree is updated based on the items in the *rescan\_set*.  
 21 In this example, it can be found that an item ( $b$ ) appears in transaction  $T_1$  and  $T_2$  of Table  
 22 1; the corresponding branches of ( $b$ ) in  $T_1$  and  $T_2$  are extracted. The corresponding branches

1 indicate that the items appear in the set of *insert\_set* and corresponding to the items in the  
 2 set of *rescan\_set* in the original database, which can be used to update the HAUP-tree of the  
 3 items in the *rescan\_set* since those items did not exist in the HAUP-tree. The results are then  
 4 shown in Table 6.

Table 6: Corresponding branches in the original database.

<b>Tid</b>	<b>Items with their quantities</b>	<b>Corresponding branch</b>
$T_1$	$a:5, b:2, c:3, d:2$	$a:5, \mathbf{b:2}, c:3, d:2$
$T_2$	$b:5, d:4, e:2, g:12$	$\mathbf{b:5}, d:4, e:2$

5 From the discovered corresponding branches of (*b*) shown in Table 6, the HAUP-tree is  
 6 then updated and the item (*b*) is then inserted to the HAUP-tree according to its correct  
 7 position in the branches. Besides, the attached array of the items (nodes) in the HAUP-tree is  
 8 then updated. Since the items (*a*), (*c*), (*d*), and (*e*) were already inserted into the HAUP-tree  
 9 except the item (*b*), it is unnecessary to inserted those items in the HAUP-tree.

10 After that, the items in *insert\_set* are processed for handling the newly inserted transac-  
 11 tions. For the inserted transactions, the corresponding branches of the items in *insert\_set* are  
 12 then found and the results are shown in Table 7.

Table 7: Corresponding branches in the inserted transactions.

<b>Tid</b>	<b>Items with their quantities</b>	<b>Corresponding branch</b>
$T_6$	$a:4, b:6, c:4, d:2$	$a:4, b:6, c:4, d:2$
$T_7$	$a:5, c:3, e:1, f:2, g:6$	$a:5, c:3, e:1, f:2$

13 Since the corresponding branches of Table 7 are from the newly inserted transactions;  
 14 those branches are respectively inserted into the HAUP-tree following the construction phase  
 15 of HAUP-tree. The final updated HAUP-tree is then shown in Figure 4.

16 After that, the HAUP-growth [21] is then processed to find the set of HAUIs. Notice that  
 17 the HAUIs can be thus discovered while all the insertion operations are completely processed.  
 18 From Figure 4, the discovered HAUIs are  $\{a:156, b:92, c:144, d:176, e:92, f:100\}$ .

## 19 6 Experimental Results

20 In this section, the performance of the proposed IHAUPM algorithm is compared with the  
 21 state-of-the-art PAI [19], TPAU [15], HAUI-tree [35], HAUI-Miner [33] and the HAUP-growth [21]  
 22 algorithms on several datasets. All algorithms were implemented in Java and experiments were  
 23 carried on a computer having an Intel(R) Core(TM) i7-4790 3.60GHz processor with 8 GB of  
 24 main memory, running the 64 bit Microsoft Windows 7 operating system. Experiments were  
 25 conducted on six datasets, including four real-world datasets [11] and two synthetic dataset  
 26 generated using the IBM Quest Synthetic Data Generator [16]. A simulation model [31] was  
 27 developed to generate quantities (internal utilities) and unit profit values (external utilities)

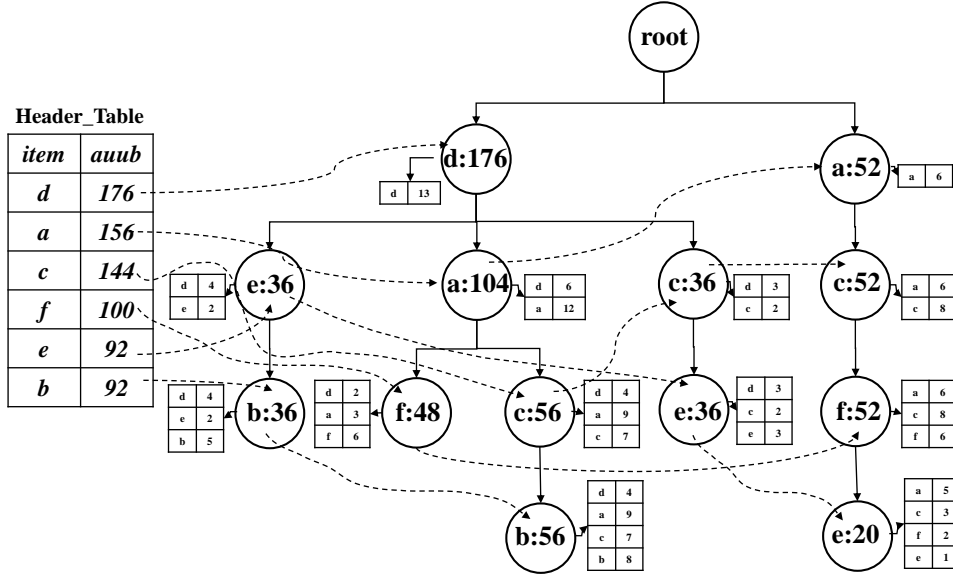


Figure 4: The final updated HAUP-tree.

1 of items in transactions for all datasets. External utilities have been generated in the [0.01,  
2 10] interval using a log-normal distribution, and internal utilities have been randomly chosen  
3 in the [1,5] interval. Parameters used to describe the six datasets are shown in Table 8, and  
4 the characteristics of these datasets are shown in Table 9. In the experiments, the minimum  
5 average-utility threshold is denoted as **TH**, and the insertion ratio of the new transactions is  
6 denoted as **IR**. The **IR** is used to set the size of the newly inserted transactions compared to  
7 the original database. For example, if the **IR** is set as 5% and the size of the original dataset  
8 has 100 transactions, the number of newly inserted transactions is thus set as 5. To simulate  
9 the dynamic situation in real-world applications, the newly inserted transactions are randomly  
10 generated five times and then respectively inserted into the original dataset.

Table 8: Parameters of used datasets.

<b># D </b>	Total number of transactions
<b># I </b>	Number of distinct items
<b>AveLen</b>	Average length of transactions
<b>MaxLen</b>	Maximal length of transactions
<b>Type</b>	Database type (sparse or dense)

11 To assess the performance of the compared algorithms, the runtime, memory usage, number  
12 of determined candidates, and scalability were respectively analyzed. Results are reported  
13 below.

Table 9: Characteristics of used datasets.

Dataset	# D	# I	AvgLen	MaxLen	Type
foodmart	21,556	1,559	4	11	sparse
kosarak	990,002	41,270	8	2,498	sparse
retail	88,162	16,407	10	76	sparse
mushroom	8,125	120	23	23	dense
T10I4D100K	100,000	1,000	10	29	sparse
T40I10D100K	100,000	1,000	39.6	77	dense

## 1 6.1 Runtime

2 In this section, the runtime for mining HAUIs are compared under various minimum high  
 3 average-utility thresholds on six datasets. Notice that the runtime is the sum up of I/O time,  
 4 interval time and mining time. Results are shown in Figure 5.

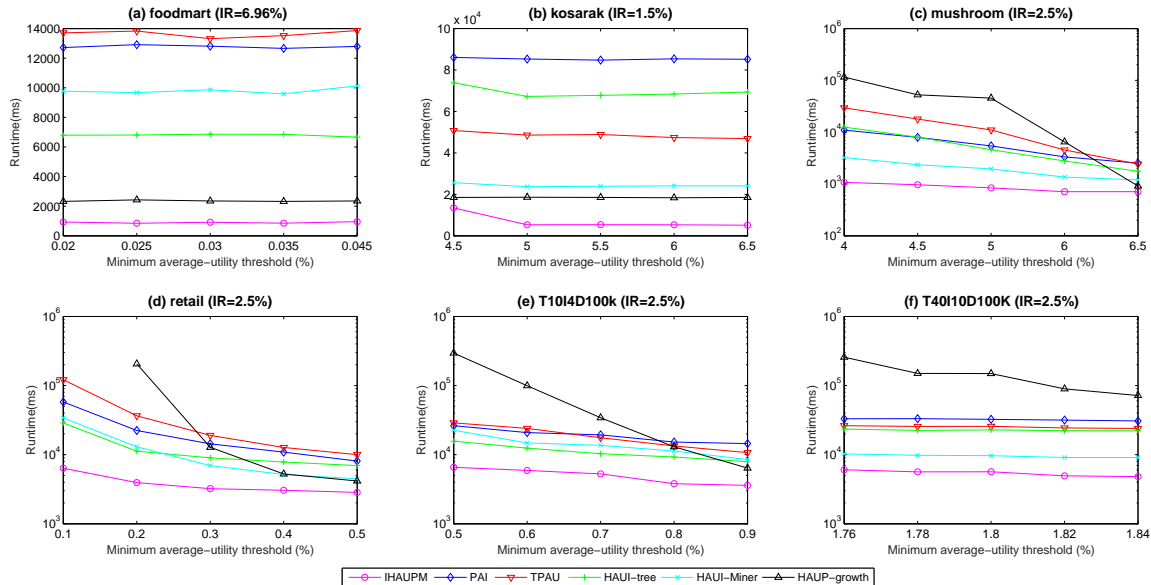


Figure 5: Runtime under various minimum high average-utility thresholds.

5 From Figure 5, it can be seen that the runtime decreases along with the increasing of mini-  
 6 mum high average-utility thresholds. This is reasonable since when the minimum high average-  
 7 utility threshold is set higher, fewer HAUIs are then discovered; thus runtime of the compared  
 8 algorithms decreased. Moreover, it is obvious to see that the proposed IHAUPM algorithm  
 9 outperforms the other algorithms under various minimum high average-utility thresholds. For  
 10 example in Figure 5(a), the proposed IHAUPM has one order of magnitude faster than the  
 11 TPAU and PAI algorithms and 10 times faster than the HAU-Miner algorithm, 6 times faster  
 12 than the HAU-tree algorithm. This is reasonable since the TPAU and PAI algorithms apply  
 13 the Two-Phase model to level-wisely discover the potential candidates for deriving the HAUIs.  
 14 **Although the HAU-Miner is the state-of-the-art algorithm for mining HAUIs, it still requires**

1 to handle the updated dataset in the batch manner, which means that the discovered infor-  
 2 mation is thus ignored and the partial parts of the original dataset are scanned with multiple  
 3 times. In addition, the runtime of the compared algorithms is then evaluated under various  
 4 insertion ratios and the results are then shown in Figure 6.

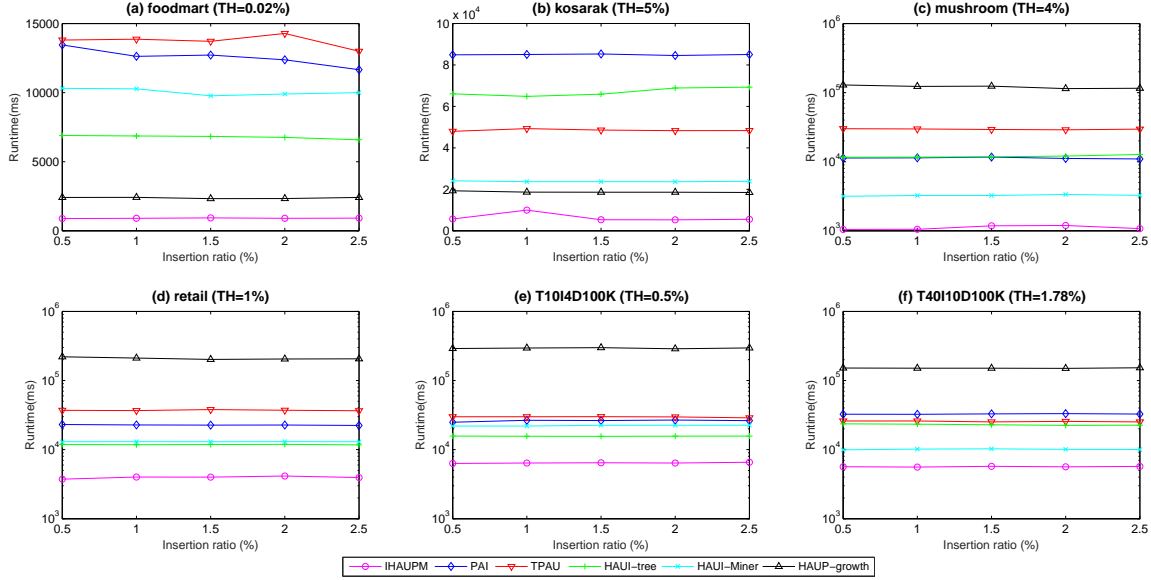


Figure 6: Runtime under various insertion ratios.

5 As shown in Figure 6, it can be concluded that the proposed IHAUPM performs well under  
 6 various insertion ratios and has one or two orders of magnitude faster than the other algorithms.  
 7 For example in Figure 6(f), the designed IHAUPM performs better than the state-of-the-art  
 8 HAUl-Miner algorithm. The reason is that IHAUPM is designed for particularly handling  
 9 the incremental database whereas the HAUl-Miner is performed on the batch manner and  
 10 the database is required to be scanned with multiple times. Overall, the designed IHAUPM  
 11 performs well than the other approaches.

## 12 6.2 Memory Usage

13 In this section, the memory usage for deriving the HAUlS from the updated dataset is com-  
 14 pared. Results are shown in Figure 7.

15 It can be observed in Figure 7 that the designed IHAUPM algorithm has the least memory  
 16 usage than the others under various minimum high average-utility thresholds on six datasets.  
 17 For example in Figure 7(a), the memory usage of the designed algorithm is about 7 to 8 times  
 18 less than PAI, TPAU, HAUl-tree and HAUl-growth and 1 times less than the state-of-the-art  
 19 HAUl-Miner. The reason is that the PAI and TPAU are performed by the Two-Phase model,  
 20 thus a large number of candidates is kept in the memory. The HAUl-Miner has no valid pruning  
 21 strategies and it sometimes generates a large number of unpromising candidates, which leads  
 22 to massive memory usage. The designed IHAUPM algorithm has, however, to maintain the

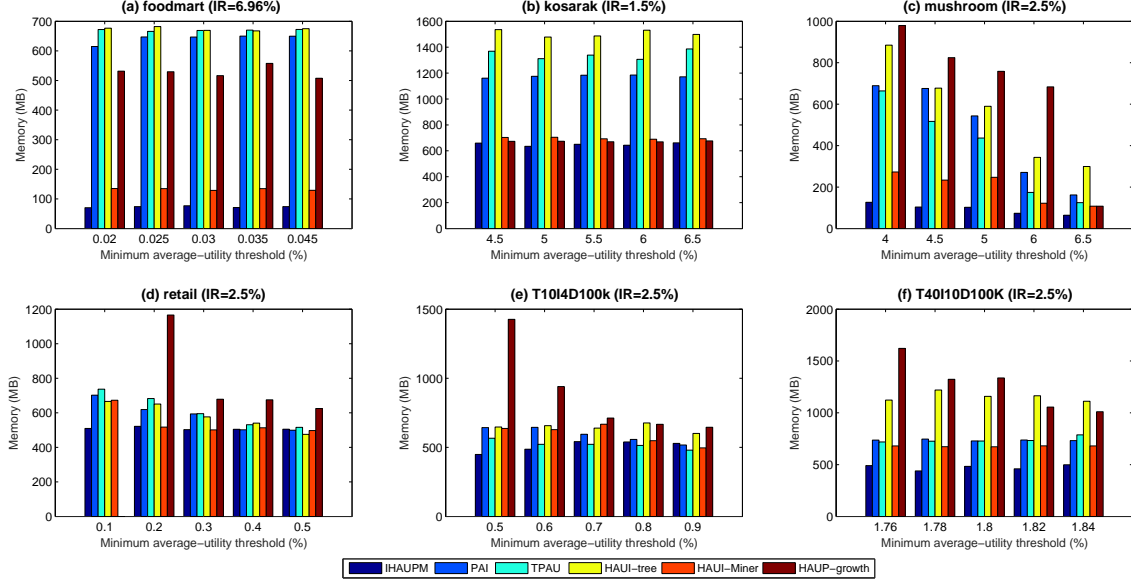


Figure 7: Memory usage under various minimum average-utility thresholds.

1 incremental parts of the inserted transactions, and the several unpromising candidates can be  
 2 early pruned by the designed algorithm. From Figures 7(c), 7(d), 7(e), and 7(f), it can be  
 3 found that the HAUP-growth needs a large number of memory usage. Although the HAUP-  
 4 growth is efficient in mining process but it necessitates to keep the quantity information of  
 5 the items in a branch, thus requiring more memory usage. In particular, when the minimum  
 6 high average-utility threshold is set lower such as 0.1%, the HAUP-growth sometimes has  
 7 the problem of memory leakage, which can be found in Figure 7(d). The memory usage of  
 8 compared algorithms under various insertion ratios are then shown in Figure 8.

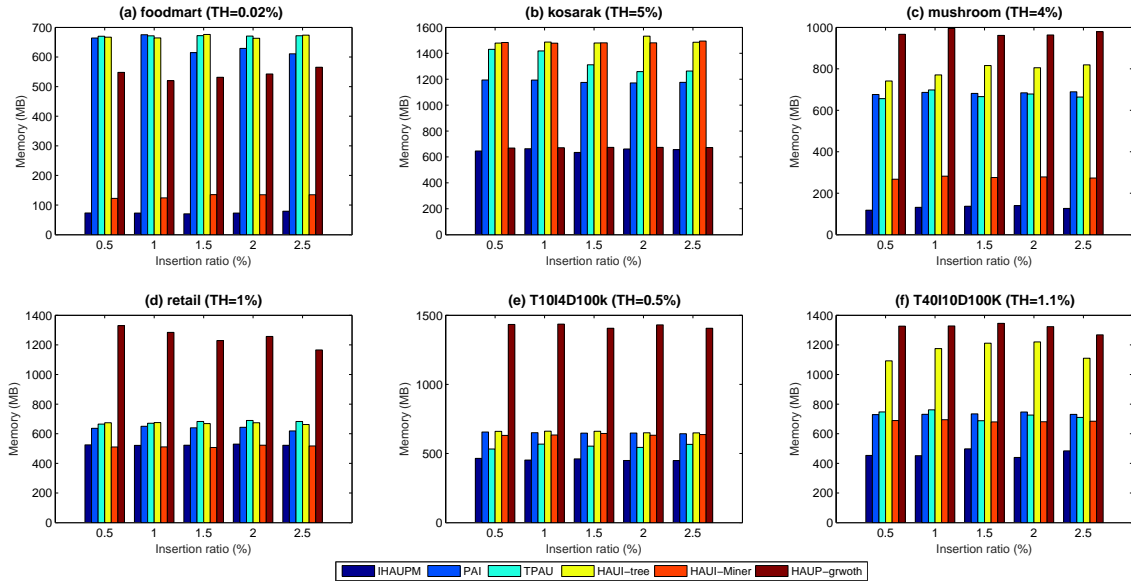


Figure 8: Memory usage under various insertion ratios.

1 It can be observed in Figure 8 that the designed IHAUPM algorithm outperforms the other  
 2 algorithms under various insertion ratios. The HUAP-growth, PAI and TPAU almost have the  
 3 most memory usage on six datasets. Although the HAU-Miner has relatively low memory  
 4 usage, the designed IHAUPM algorithm still performs well than HAU-Miner. It also can be  
 5 found that the HAU-Miner is sometimes not efficiency since it is not designed for incremental  
 6 database. Overall, the designed IHAUPM algorithm for handling the incremental database  
 7 with transaction insertion performs well among the compared algorithms.

### 8 6.3 Number of Candidates

9 In this section, the number of candidates for finding the HAUIs under the various minimum  
 10 high average-utility thresholds is compared. Results are shown in Figure 9.

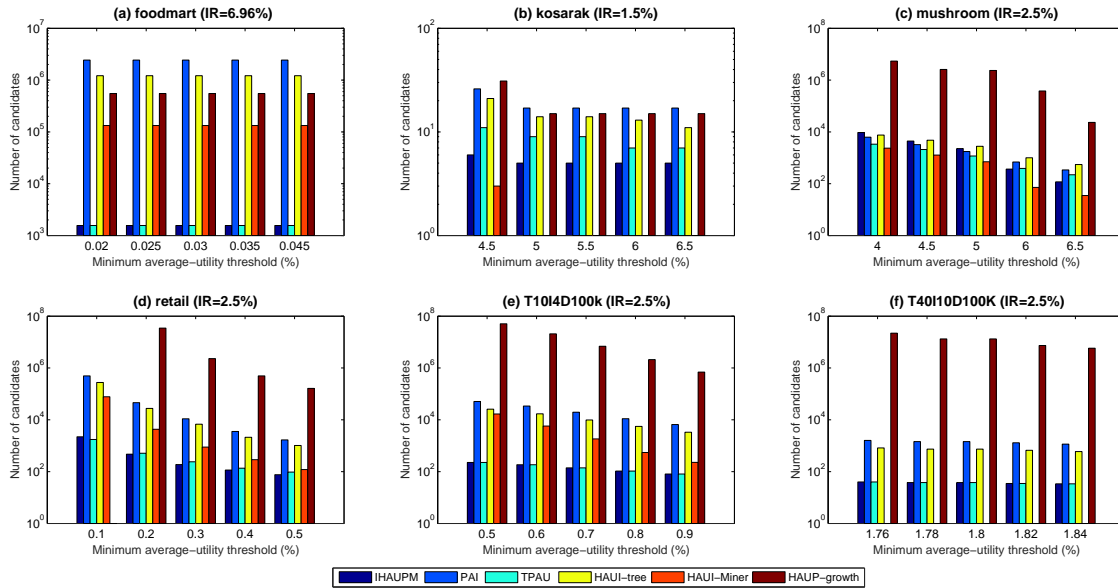


Figure 9: Number of candidates under various minimum high average-utility thresholds.

11 It can be observed in Figure 9 that the number of candidates generated by the designed  
 12 algorithm is much less than the PAI, HAU-tree and HAUP-growth algorithms. For example  
 13 in Figure 9(a), the number of candidates generated by the designed algorithm is about 2 to 3  
 14 orders of magnitude less than PAI, HAU-tree and HAUP-growth and the same as the TPAU.  
 15 In Figures 9(b) and 9(c), the number of candidates of the designed algorithm is less than  
 16 TPAU. The reason is that PAI, HAU-tree and HAUP-growth has no valid pruning strategies  
 17 to reduce the unpromising candidates for the whole updated dataset but the designed IHAUPM  
 18 algorithm handles only the incremental parts; the size of the incremental dataset is much less  
 19 than that of the updated one. Although the number of candidates generated by TPAU and  
 20 HAU-Miner is sometimes the same as or slightly more than designed algorithm, the TPAU  
 21 suffers the inefficiency mining performance since it mines the HAUIs in the level-wise manner  
 22 and the HAU-Miner is designed for the static dataset but not for the incremental one. Besides,

1 the algorithms are also compared in terms of number of candidates under different insertion  
 2 ratios and the results are shown in Figure 10.

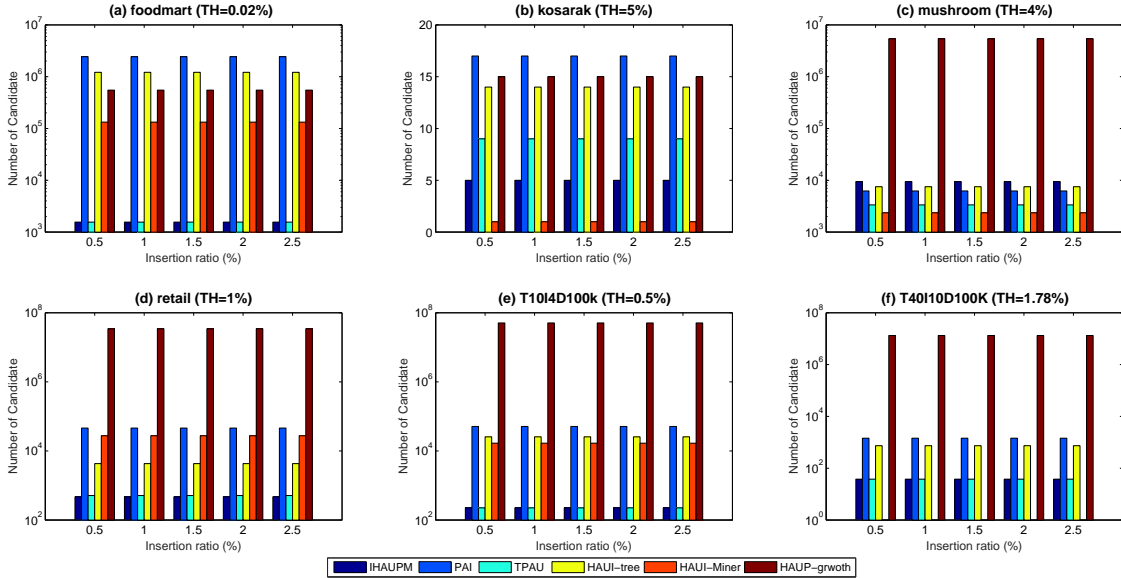


Figure 10: Number of candidates under various insertion ratios.

3 It can be observed in Figure 10 that the designed IHAUPM algorithm almost has the least  
 4 number of candidates on six datasets under various insertion ratios, while the HAUP-growth  
 5 and PAI almost have the most number of candidates. For example in Figure 10(a), the PAI  
 6 has the two or three orders of magnitude more than the designed algorithm, and in Figure  
 7 10(c), the number of candidate of HAUP-growth is about 8 to 9 times more than others. Also,  
 8 the number of candidates remains steady for all compared algorithms. It is reasonable since  
 9 the size of the inserted transactions is much less than the original database; the designed  
 10 IHAUPM algorithm only handles the small parts but the others need to maintain the whole  
 11 updated database even though only few transactions are inserted into the database.

## 12 6.4 Scalability

13 In this section, the scalability of the algorithms were compared on a real-world kosarak  
 14 dataset [11] and a series of synthetic dataset T10I4N4KD|X|K generated by IBM Genera-  
 15 tor [16] using the same parameter other than dataset size. |X| represents the dataset size,  
 16 varying from 100K to 500K transactions with 100K increment in our experiments. To test the  
 17 scalability of the proposed algorithm on real-world datasets, the original dataset is divided into  
 18 5 datasets kosarak|X|K, where the the dataset size is vary from 198K to 990K, with increments  
 19 of 198K. The **TH** is set as 3.5% for the synthetic dataset and 4.5% for the kosarak dataset,  
 20 and the **IR** is set as 5% for the synthetic datasets and 2.5% for the kosarak dataset. The  
 21 results are respectively shown in Figures 11 and 12.

22 It can be observed in Figures 11 and 12 that the runtime, memory usage, and number of



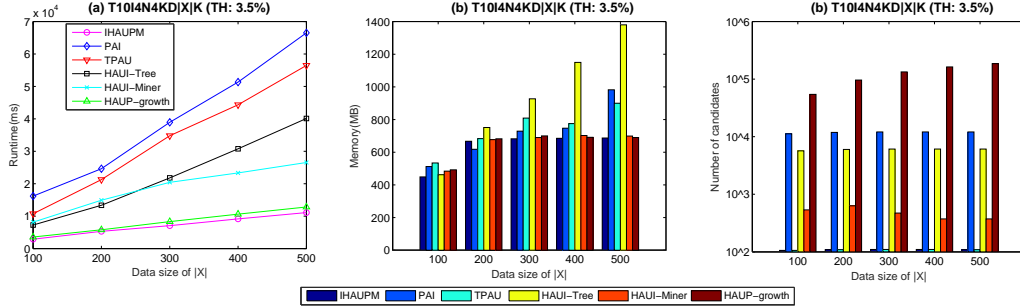


Figure 11: Scalability under various synthetic dataset sizes.

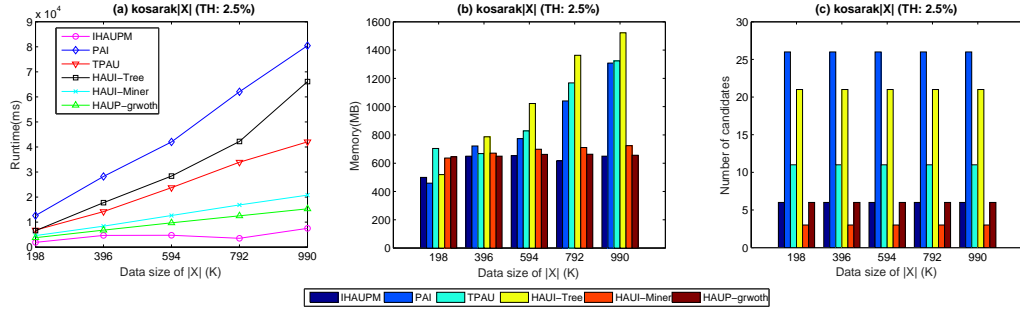


Figure 12: Scalability under various kosarak dataset sizes.

1 candidates increase along with the increasing of dataset size. This is reasonable since when  
 2 the dataset size increases, the number of HAUIs increases, as well as the runtime, memory  
 3 usage and number of candidates. The designed IHAUPM algorithm performs well than the  
 4 other algorithms on both synthetic datasets and real-world datasets, and the PAI and TPAU  
 5 algorithm have worse results of runtime. This is reasonable since they are performed in the  
 6 level-wise manner to mine the HAUIs. The HAUI-tree requires the most memory usage since  
 7 it keeps a lot of information in the tree structure. The HAUP-growth needs the most number  
 8 of candidates since the potential and complete HAUIs are maintained by the HAUP-tree struc-  
 9 ture. Overall, the designed IHAUPM algorithm outperforms the others in terms of runtime,  
 10 memory usage and the number of candidates. The effectiveness and efficiency of the designed  
 11 algorithm is thus quite acceptable for real-world applications.

12 From the experiments, we can also see that the proposed algorithm performs than other  
 13 algorithms in terms of runtime, memory usage and scalability in the incremental environment.  
 14 This is reasonable since the designed algorithm is unnecessary to rescan the entire database with  
 15 multiple times but only focusing on the incremental parts. For the batch-manner algorithms,  
 16 the database is multiply rescanned whenever updated records come. Besides, the presented  
 17 HAUP-tree keeps the sufficient information for later updating with only few additional memory  
 18 usages to update the HAUIs. Thus, the proposed algorithm scales well whether in the synthetic  
 19 or real-world datasets.

## 7 Conclusion and Future Work

High average utility mining (HAUIM) has attracted a lot of attention since it provides a fair measure than the HUIM to evaluate the discovered patterns. However, most algorithms of HAUIM are designed for handling the static database. Thus, when transactions are inserted into the original database, the updated database is required to be re-mined and the previous discovered information becomes useless.

To solve this problem, in this paper, we present an incremental high average-utility pattern mining (IHAUPM) algorithm for handling the transaction insertion in the dynamic databases. The FUP concept is thus modified and extended to handle original database and new transactions, and divide them into 4 cases. The compact HAUP-tree structure is used to maintain the 1-HAUUBIs, which can guarantee the correctness and completeness of the discovered HAUIs. The proof of each case in the designed algorithm is also provided. Based on the designed model, it rarely scans the original database except handling the itemsets in case 3. Theoretically, the designed algorithm can greatly reduce the computational cost while the transactions are frequently inserted into the original database. Thus, the designed algorithm can solve the realistic problem in the engineering, as well as provide the real-time decision-making system. Experimental results show that the proposed IHAUPM algorithm performs well than the state-of-the-art algorithms for maintaining the HAUIs.

Since transaction deletion and transaction modification also happen in real-life situations, it is necessary to design the efficient algorithms for updating the discovered information in the dynamic databases. Besides, a more compact data structure should be further designed to keep the necessary information especially in the limited memory.

## Acknowledgment

This research was partially supported by National Natural Science Foundation of China (NSFC) under grant No.61503092 and by the Tencent Project under grant CCF-Tencent IAGR20160115.

## References

- [1] R. Agrawal and R. Srikant, "Fast algorithms for mining association rules in large databases," *The International Conference on Very Large Databases*, pp. 487-499, 1994.
- [2] R. Agrawal, T. Imielinski, and A. Swami, "Mining association rules between sets of items in large databases," *ACM SIGMOD International Conference on Management of Data*, pp. 207-216, 1993.
- [3] R. Agrawal, C. Aggarwal, and V. Prasad, "A tree projection algorithm for generation of frequent itemsets," *Journal of Parallel and Distributed Computing*, vol. 61, pp. 350-371, 2001.

- 1 [4] C. F. Ahmed, S. K. Tanbeer, B. S. Jeong, and Y. K. Lee, "Efficient tree structures for  
2 high utility pattern mining in incremental databases," *IEEE Transactions on Knowledge  
3 and Data Engineering*, vol. 21, no. 12, pp. 17081721, 2009.
- 4 [5] R. Chan, Q. Yang, and Y. D. Shen, "Mining high utility itemsets," *IEEE International  
5 Conference on Data Mining*, pp. 19-26, 2003.
- 6 [6] C. Chang and C. Lin, "Perfect hashing schemes for mining association rules," *The  
7 Computer Journal*, vol. 48, pp. 168-179, 2005.
- 8 [7] D. Cheung, J. Han, V. Ng, and C. Wong, "Maintenance of discovered association rules in  
9 large databases: an incremental updating approach," *IEEE International Conference on  
10 Data Engineering*, pp. 106-114, 1996.
- 11 [8] C. Chi and K. Lain, "A new fp-tree algorithm for mining frequent itemsets," *The Content  
12 Computing, Advanced Workshop on Content Computing*, pp. 266-277, 2004.
- 13 [9] C. J. Chu, V. S. Tseng, and T. Liang, "An efficient algorithm for mining high utility item-  
14 sets with negative item values in large databases," *Applied Mathematics and Computation*,  
15 vol. 215(2), pp. 767778, 2009.
- 16 [10] Q. Duong, B. Liao, P. Fournier-Viger, T. L. Dama, "An efficient algorithm for mining  
17 the top-k high utility itemsets, using novel threshold raising and pruning strategies,"  
18 *Knowledge-Based Systems*, vol. 104, pp. 106-122, 2015.
- 19 [11] P. Fournier-Viger, J. C. W. Lin, A. Gomariz, T. Gueniche, A. Soltani, Z. Deng, and H.  
20 T. Lam, "The SPMF Open-Source Data Mining Library Version 2 and beyond," *The  
21 European Conference on Machine Learning and Principles and Practice of Knowledge  
22 Discovery*, 2016.
- 23 [12] G. Grahne and J. Zhu, "Fast algorithms for frequent itemset mining using fp-trees," *IEEE  
24 Transactions on Knowledge and Data Engineering*, vol. 17, pp. 1347-1362, 2005.
- 25 [13] J. Han, J. Pei, Y. Yin, and R. Mao, "Mining frequent patterns without candidate gener-  
26 ation: a frequent-pattern tree approach," *Data Mining and Knowledge Discovery*, vol. 8,  
27 pp. 53-87, 2004.
- 28 [14] T. P. Hong, C. W. Lin, and Y. L. Wu, "Incrementally fast updated frequent pattern trees,  
29 *Expert Systems with Applications*, vol. 34, pp. 2424-2435, 2008.
- 30 [15] T. P. Hong, C. H. Lee, and S. L. Wang, "Effective utility mining with the measure of  
31 average utility," *Expert Systems with Application*, vol. 38(7), pp.8259-8265, 2011.
- 32 [16] IBM Quest Data Mining Project, Quest Synthetic Data Generation Code, [http://www.  
33 almaden.ibm.com/cs/quest/syndata.html](http://www.almaden.ibm.com/cs/quest/syndata.html), 1996.

- 1 [17] S. Krishnamoorthy, "Pruning strategies for mining high utility itemsets," *Expert Systems*  
2 *with Applications*, vol. 42(5), pp. 2371-2381, 2015.
- 3 [18] D. Kim and U. Yun, "Mining high utility itemsets based on the time decaying model,"  
4 *Intelligent Data Analysis*, vol. 20(5), pp. 1157-1180, 2016.
- 5 [19] G. C. Lan, T. P. Hong, and V. S. Tseng, "Efficiently mining high average-utility itemsets  
6 with an improved upper-bound strategy," *International Journal of Information Technol-*  
7 *ogy & Decision Making*, vol. 11(5), 1009-103, 2012.
- 8 [20] C. W. Lin, G. C. Lan, and T. P. Hong, "An incremental mining algorithm for high utility  
9 itemsets," *Expert Systems with Applications*, vol. 39(8), pp. 7173-7180, 2012.
- 10 [21] C. W. Lin, T. P. Hong and W. H. Lu, "An effective tree structure for mining high utility  
11 itemsets," *Expert Systems with Applications*, vol. 38(6), pp. 7419-7424, 2011.
- 12 [22] J. C. W. Lin, W. Gan , P. Fournier-Viger, T. P. Hong, V. S. Tseng, "Efficient algorithms  
13 for mining high-utility itemsets in uncertain databases," *Knowledge-Based Systems*, vol.  
14 96, pp. 171-187, 2015.
- 15 [23] J. C. W. Lin, W. Gan, P. Fournier-Viger, L. Yang, Q. Liu, S. Lukas, M. Voznak, "High  
16 utility-itemset mining and privacy-preserving utility mining," *Perspectives in Science*,  
17 vol. 7, pp. 74-80, 2015.
- 18 [24] J. C. W. Lin, P. Fournier-Viger, W. Gan, "FHN: An efficient algorithm for mining high-  
19 utility itemsets with negative unit profits," *Knowledge-Based Systems*, vol. 111, pp.  
20 283-298, 2016.
- 21 [25] J. C. W. Lin, W. Gan, P. Fournier-Viger, T. P. Hong, and V. S. Tseng, "Fast algorithms  
22 for mining high-utility itemsets with various discount strategies," *Advanced Engineering*  
23 *Informatics*, vol. 30(2), pp. 109-126, 2016.
- 24 [26] J. C. W. Lin, W. T. Wu, P. Fournier-Viger, G. Lin, J. Zhan, and M. Voznak, "Fast  
25 algorithms for hiding sensitive high-utility itemsets in privacy-preserving utility mining,"  
26 *Engineering Applications of Artificial Intelligence*, vol. 55, pp. 269-284, 2016.
- 27 [27] J. C. W. Lin, W. Gan , P. Fournier-Viger, T. P. Hong, and J. Zhan, "Efficient mining  
28 of high-utility itemsets using multiple minimum utility thresholds," *Knowledge-Based*  
29 *Systems*, vol. 113, pp. 100-115, 2016.
- 30 [28] J. C. W. Lin, L. Yang, P. Fournier-Viger, J. M. T. Wu, T. P. Hong, Leon S. L. Wang, and J.  
31 Zhan, "Mining high-utility itemsets based on particle swarm optimization," *Engineering*  
32 *Applications of Artificial Intelligence*, vol. 5, pp. 320-330, 2016.
- 33 [29] Y. C. Lin, C. W. Wu, and V. S. Tseng, "Mining high utility itemsets in big gata,"  
34 *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pp. 649-661, 2016.

- 1 [30] J. C. W. Lin, T. Li, P. Fournier-Viger, T.P. Hong, J. Zhan, M. Voznak, “An efficient  
2 algorithm to mine high average-utility itemsets,” *Advanced Engineering Informatics*, vol.  
3 30 pp. 233-243, 2016.
- 4 [31] Y. Liu, W. K. Liao, and A. Choudhary, “A two-phase algorithm for fast discovery of high  
5 utility itemsets,” *Lecture Notes in Computer Science*, vol. 3518, pp. 689-695, 2005.
- 6 [32] Y. Liu, W. Liao, and A. Choudhary, “A fast high utility itemsets mining algorithm,” *The  
7 International Workshop on Utility-Based Data Mining*, pp. 90-99, 2005.
- 8 [33] M. Liu, J. Qu, “Mining high utility itemsets without candidate generation categories,”  
9 *International conference on Information and knowledge management*, pp. 55-64, 2012.
- 10 [34] T. Lu, B. Vo, H. T. Nguyen, and T. P. Hong, “A new method for mining high average  
11 utility itemsets,” *The IFIP International Conference on Computer Information Systems  
12 and Industrial Management*, pp. 33-42, 2014.
- 13 [35] J. Liu, K. Wang, B. Fung, “Mining high utility patterns in one phase without generating  
14 candidates,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, pp. 1245-  
15 1257, 2016.
- 16 [36] J. S. Park, M. Chen, and P. Yu, “An effective hash-based algorithm for mining association  
17 rules,” *ACM SIGMOD International Conference on Management of Data*, pp. 175-186,  
18 1995.
- 19 [37] J. S. Park, M. S. Chen, P. S. Yu, “Using a hash-based method with transaction trimming  
20 for mining association rules,” *IEEE Transactions on Knowledge and Data Engineering*,  
21 vol. 9(5), pp. 812-825, 1997.
- 22 [38] R. Srikant and R. Agrawal, “Mining generalized association rules,” *The International  
23 Conference on Very Large Data Bases*, pp. 407-419, 1995.
- 24 [39] V. S. Tseng, B. E. Shie, C. W. Wu, and P. S. Yu, “Efficient algorithms for mining high  
25 utility itemsets from transactional databases,” *IEEE Transactions on Knowledge and  
26 Data Engineering*, vol. 25(8), pp. 1772-1786, 2013.
- 27 [40] V. S. Tseng, C. W. Wu, P. Fournier-Viger, and P. S. Yu, “Efficient algorithms for mining  
28 Top- $K$  high utility itemsets,” *IEEE Transactions on Knowledge and Data Engineering*,  
29 vol. 28, pp. 54-67, 2016.
- 30 [41] J. M. T. Wu, J. Zhang, and J. C. W. Lin, “Mining of high-utility itemsets by ACO  
31 algorithm,” *The 3rd Multidisciplinary International Social Networks Conference on So-  
32 cialInformatics*, Article No. 44, 2016.
- 33 [42] H. Yao, H. Hamilton, and C. Butz, “A foundational approach to mining itemset utilities  
34 from databases,” *SIAM International Conference on Data Mining*, pp. 211-225, 2004.

- 1 [43] H. Yao, H. J. Hamilton, L. Geng, “A unified framework for utility based measures for  
2 mining itemsets,” *The International Workshop on Utility-Based Data Mining*, pp. 28-37,  
3 2006.
- 4 [44] U. Yun and H. Ryang, “Incremental high utility pattern mining with static and dynamic  
5 databases,” *Applied Intelligence*, vol. 42(2), pp. 323-352, 2015.