

SIMULATED ANALYSIS OF EFFICIENT ALGORITHMS FOR MINING TOP-K HIGH UTILITY ITEMSETS

Surbhi Choudhary¹, Devendra Nagal², Swati Sharma³

¹PhD Research Scholar, Dept. of Computer Applications, JNU Jodhpur

^{2,3}Faculty, Dept. of Electrical Engineering, JNU Jodhpur

Abstract— High utility sequential pattern mining is an emerging topic in the data mining community. Compared to the classic frequent sequence mining, the utility framework provides more informative and actionable knowledge since the utility of a sequence indicates business value and impact. However, the introduction of “utility” makes the problem fundamentally different from the frequency-based pattern mining framework and brings about dramatic challenges. Although the existing high utility sequential pattern mining algorithms can discover all the patterns satisfying a given minimum utility, it is often difficult for users to set a proper minimum utility. A too small value may produce thousands of patterns, whereas a too big one may lead to no findings. In this paper, we propose a novel framework called top-k high utility sequential pattern mining to tackle this critical problem. Accordingly, an efficient algorithm, Top-k high Utility Sequence (TUS for short) mining, is designed to identify top-k high utility sequential patterns without minimum utility. In addition, three effective features are introduced to handle the efficiency problem, including two strategies for raising the threshold and one pruning for filtering unpromising items. Our experiments are conducted on both synthetic and real datasets..

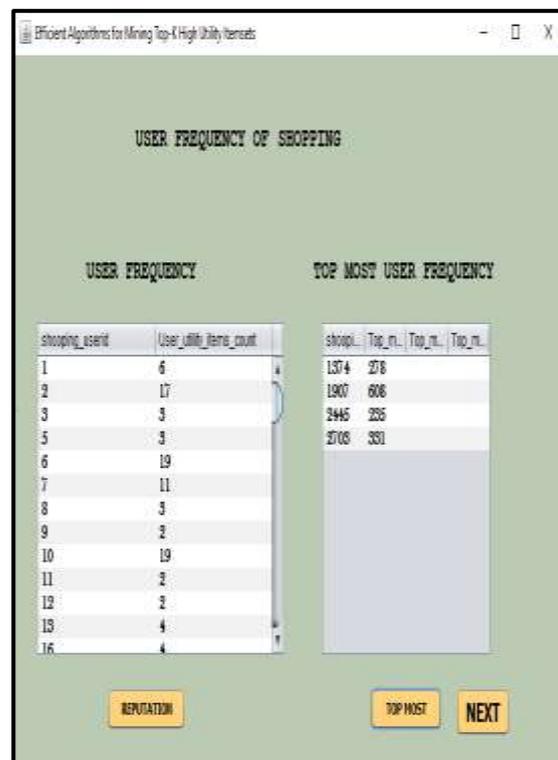
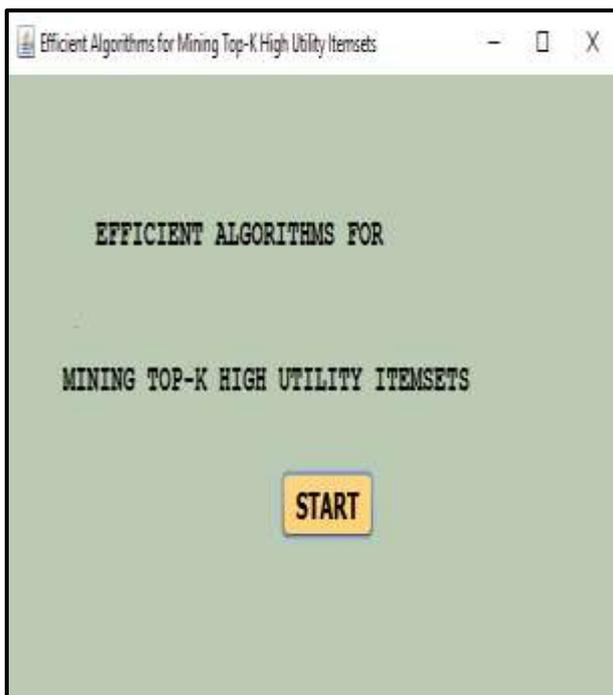
I. INTRODUCTION

FREQUENT sequential pattern mining [1], as one of the fundamental research topics in data mining, discovers frequent subsequences in sequence databases. It is very useful for handling order-based business problems, and has been successfully adopted to various domains and applications such as complex behaviour analysis [1] and gene sequence analysis [2], [3], [4]. In the frequency-based framework for typical sequence analysis, the downward closure property (also known as Apriori property) [1] plays a fundamental role in identifying frequent sequential patterns. However, taking the frequency to measure pattern interestingness may be insufficient for selecting actionable sequences associated with expected quality and business impact, because the patterns identified under the frequency (support) framework do not disclose the business value and impact. To solve the above problems, the concept utility is introduced into sequential pattern mining to select sequences of high utility by considering the quality and value (such as profit) of itemsets. This leads to an emerging area, high utility pattern mining [1], [2], [3], [8], [9] and high utility sequential pattern mining, which selects interesting patterns / sequential patterns based on minimum utility rather than minimum support. The utility-

based patterns are proven to be more informative and actionable for decision-making than the frequency-based ones [2]. For instance, in [4], [5], the authors discuss the extraction of profitable behaviours from the mobile commerce environments. [4] Proposes methods to mine high utility sequences from web logs by assigning each page an impact/significance. In [6], a US pan algorithm is built for utility-based sequential pattern mining satisfying a predefined minimum utility. Although algorithms such as USpan can obtain high utility sequences based on a given minimum utility, it is very difficult for users to specify an appropriate minimum utility threshold and to directly obtain the most valuable patterns. This is because the complexity of utility-based sequencedatabases (which may be different from the classic itemsets), determining multiple factors including the distribution of the items and utilities, density of the database, lengths of these sequences, and so on. Consequently, it is not surprising that, with a same minimum utility threshold, some datasets may produce millions of patterns while others may contribute nothing. The challenge here is that it may not be doable to tune the threshold to capture the expected number of patterns. This is because the sensitivity of the threshold makes it hard to tune for a variety of databases. It may be very costly and time consuming to achieve the proper threshold for the desired patterns. In fact, the classic frequency/support based pattern mining also faces the same challenge. Accordingly, the concept of extracting top-k patterns has been proposed in [2], [4], [7], [8] to select the patterns with the highest frequency. In the top-k frequent pattern mining, instead of letting a user specify a threshold, the top-k pattern selection algorithms allow a user to set the number of top-k high frequency patterns to be discovered. This makes it much easier and more intuitive and practical than determining a minimum support; also the determination of k by a user is more straightforward than considering data characteristics, which are often invisible to users, for choosing a proper threshold. The easiness for users to determine k does not indicate the simplicity of developing an efficient algorithm for selecting top-k high utility sequential patterns. In the utility framework, TKU is the only method for mining top-k high utility itemsets, to the best of our knowledge. Nowork is reported on mining top-k high utility sequences. There is significant difference between top-k utility itemset mining and top-k utility sequence mining in which the order between itemsets is considered. In fact, the problem of top-k high utility sequence mining is much more challenging than mining top-k high utility itemsets. First, as

with high utility Itemset mining, the downward closure property does not hold in the utility-based sequence mining. This means that the existing top-k frequent sequential pattern mining algorithms [7] cannot be directly applied. Second, compared to top-k high utility itemset mining [8], utility-based sequence analysis faces the critical combinational explosion and computational complexity caused by sequencing between itemsets. This means that the techniques in [9] cannot be directly transferred to top-k high utility sequential pattern mining either. Third, since the minimum utility is not given in advance, the algorithm essentially starts the searching from 0 minimum supports. This not only incurs very high computational costs, but also the challenge of how to raise the minimum threshold without missing any top-k high utility sequences.

II. SIMULATION RESULTS



Efficient Algorithms for Mining Top-K High Utility Itemsets

USER TRANSACTION (BOUGHT) ITEMS

SHOPPING

| user_id | product_id | user_util | shopping_util | shopping_id | shopping_util | id | Name | Total |
|---------|------------|-----------|---------------|-------------|---------------|----|---------------|--------|
| 0 | 1 | 5 | 0 | 1 | 5 | 1 | BEAUTYES.. | 3194 |
| 0 | 7465 | 3 | 1 | 2 | 4 | 2 | HOME_SER.. | 5179 |
| 0 | 7466 | 5 | 1 | 3 | 5 | 3 | HOTELRY.. | 2165 |
| 0 | 7467 | 5 | 1 | 4 | 5 | 4 | NIGHT_LL.. | 9867 |
| 0 | 7468 | 4 | 1 | 5 | 4 | 5 | RESTAURA.. | 9194 |
| 0 | 33686 | 3 | 1 | 6 | 5 | 6 | SHOPPING | 33351 |
| 0 | 33689 | 2 | 1 | 7 | 5 | 7 | PKTS | 1671 |
| 0 | 33690 | 3 | 2 | 8 | 5 | 8 | Total_user... | 263776 |
| 0 | 33691 | 5 | 2 | 9 | 5 | | | |
| 0 | 33692 | 3 | 2 | 10 | 1 | | | |
| 0 | 33693 | 3 | 2 | 11 | 5 | | | |
| 0 | 33694 | 4 | 2 | 12 | 3 | | | |
| 0 | 9985 | 4 | 2 | 13 | 5 | | | |

VIEW RATINGS VIEW RATINGS TOTAL

TOTAL 263776 TOTAL 33351 NEXT

Efficient Algorithms for Mining Top-K High Utility Itemsets

UTILITY ITEMS OF SHOPPING

UTILITY ITEMS TOP MOST UTILITY ITEMS

| shopping_itemid | shopping_item_na.. | item_utility_count | shopping_item.. | shopping_item.. | Topmost_utili.. |
|-----------------|--------------------|--------------------|-----------------|-----------------|-----------------|
| 1 | Department Sto... | 3 | 24 | Musicals A... | 83 |
| 2 | America (Track... | 7 | 29 | Kitchen & B... | 71 |
| 3 | Art Supplies To... | 12 | 36 | Department... | 71 |
| 4 | Coffee & Tea A... | 7 | 48 | Art Gallere... | 63 |
| 5 | Hardware Stores... | 1 | 81 | Home Deco... | 57 |
| 6 | Arts & Crafts B... | 1 | 93 | Electronics ... | 72 |
| 7 | Bookstore Vary... | 1 | 196 | Computers ... | 56 |
| 8 | Accessories Lin... | 5 | 236 | Women's CL... | 66 |
| 9 | Adult Lingerie | 9 | 333 | Bookstore... | 50 |
| 10 | Drugsstore Com... | 1 | 536 | Shopping Ce... | 84 |
| 11 | Toy Stores Card... | 18 | 544 | Men's Cloth... | 136 |
| 12 | Adult Entertain... | 12 | 751 | Electronics ... | 69 |
| 13 | Used, Vintage &... | 3 | 932 | Farmers Ma... | 68 |

Utility Items Top-K Utility Itemsets in One Phase NEXT

Efficient Algorithms for Mining Top-K High Utility Itemsets

UTILITY ITEMS OF SHOPPING

UTILITY ITEMS TOP MOST UTILITY ITEMS

| shopping_itemid | UTILITY_ITEMS... | shopping_item.. | TOPMOST UTIL.. |
|-----------------|------------------|-----------------|----------------|
| 3 | 12 | 24 | 83 |
| 11 | 18 | 29 | 71 |
| 12 | 12 | 36 | 71 |
| 22 | 41 | 48 | 63 |
| 24 | 83 | 93 | 72 |
| 25 | 11 | 236 | 66 |
| 27 | 45 | 536 | 84 |
| 29 | 71 | 544 | 136 |
| 32 | 39 | 751 | 69 |
| 34 | 16 | 932 | 68 |
| 36 | 71 | | |
| 37 | 13 | | |
| 38 | 95 | | |

FREQUENCY ITEMS Top-K Utility Itemsets NEXT

Efficient Algorithms for Mining Top-K High Utility Itemsets

SHOPPING

NEGATIVE PROFIT UTILITY ITEMS

| shopping_itemid | shopping_item_na.. | Topmost_utili.. |
|-----------------|--------------------|-----------------|
| 81 | Home Decor F... | 18 |
| 147 | Women's Cloth... | 17 |
| 236 | Women's Cloth... | 18 |
| 544 | Men's Clothing ... | 13 |
| 765 | Home Decor F... | 12 |
| 1163 | Men's Clothing ... | 15 |
| 1888 | Men's Clothing ... | 15 |

NEGATIVE NEXT

Efficient Algorithms for Mining Top-K High Utility Items

SHOPPING

NEGATIVE UTILITY ITEMS

| shopping_item... | shopping_item... | Topmost_utili... |
|------------------|------------------|------------------|
| 3 | Art Supplie... | 1 |
| 10 | Drugstores ... | 1 |
| 11 | Toy Stores ... | 1 |
| 12 | Adult Enter... | 1 |
| 17 | Museums A... | 1 |
| 19 | Women's CL... | 1 |
| 21 | Lingerie Co... | 1 |
| 22 | Music & DV... | 1 |
| 25 | Skin Care C... | 2 |
| 27 | Department... | 3 |
| 29 | Kitchen & B... | 9 |
| 32 | Desserts To... | 5 |
| 36 | Department | 8 |

POSITIVE UTILITY ITEMS

| shopping_it... | shopping_it... | Items_utili... |
|----------------|----------------|----------------|
| 1 | Departme... | 3 |
| 2 | America... | 7 |
| 3 | Art Suppl... | 11 |
| 4 | Coffee & ... | 7 |
| 5 | Hardware ... | 1 |
| 6 | Arts & Cr... | 1 |
| 7 | Bookstore... | 1 |
| 8 | Accessori... | 5 |
| 9 | Adult Lin... | 9 |
| 11 | Toy Store... | 17 |
| 12 | Adult Ent... | 11 |
| 13 | Used, Vin... | 3 |
| 14 | Adult Lea... | 1 |

NEGATIVE
POSITIVE
NEXT

Efficient Algorithms for Mining Top-K High Utility Items

COMPARE WITH POSITIVE AND NEGATIVE

SHOPPING

NEXT

| shopping_item... | shopping_item... | Topmost_utili... |
|------------------|------------------|------------------|
| 81 | Home Decor... | 18 |
| 147 | Women's CL... | 17 |
| 236 | Women's CL... | 18 |
| 544 | Men's Cloth... | 13 |
| 765 | Home Decor... | 12 |
| 1163 | Men's Cloth... | 15 |
| 1888 | Men's Cloth... | 15 |

| shopping_it... | shopping_it... | Items_utili... |
|----------------|----------------|----------------|
| 24 | Museums ... | 83 |
| 29 | Kitchen & ... | 62 |
| 36 | Departmen... | 63 |
| 93 | Electronics... | 69 |
| 536 | Shopping ... | 84 |
| 544 | Men's Clo... | 123 |
| 751 | Electronics... | 61 |
| 932 | Farmers M... | 67 |

NEGATIVE
POSITIVE

Efficient Algorithms for Mining Top-K High Utility Items

SHOPPING

POSITIVE PROFIT UTILITY ITEMS

| shopping_itemid | shopping_item_names | Items_utility count |
|-----------------|-----------------------|---------------------|
| 24 | Museums Art Galle... | 83 |
| 29 | Kitchen & Bath Fu... | 62 |
| 36 | Department Stores... | 63 |
| 93 | Electronics Photog... | 69 |
| 536 | Shopping Centers S... | 84 |
| 544 | Men's Clothing Wo... | 123 |
| 751 | Electronics Compu... | 61 |
| 932 | Farmers Market Sb... | 67 |

NEXT

POSITIVE

Efficient Algorithms for Mining Top-K High Utility Items

OVER ALL RESULT ANALYSIS

SHOPPING

| id | NAME | ITEMS |
|----|------------------|--------|
| 51 | TOTALITEMS | 263776 |
| 52 | SHOPPING | 33351 |
| 53 | USERUTILITY | 2366 |
| 54 | TOPMOSTUTIL... | 4 |
| 55 | UTILITY_ITEMS | 16153 |
| 56 | TOPMOST_UTL... | 14 |
| 57 | POSITVIE | 3071 |
| 58 | NEGATIVEW | 14632 |
| 59 | NEGATIVE_UTL... | 7 |
| 60 | POSITIVE_UTIL... | 8 |

VIEW_ALL_PRODUCT
GRAPH

Efficient Algorithms for Mining Top-K High Ut...

SHOPPING

| id | NAME | ITEMS |
|----|--------------|--------|
| 1 | TOTALIT... | 263776 |
| 2 | SHOPPING | 33351 |
| 3 | USERUTIL... | 2366 |
| 4 | TOPMOST... | 4 |
| 5 | UTILITY_I... | 16153 |
| 6 | TOPMOST... | 14 |
| 7 | POSITVIE | 3071 |
| 8 | NEGATIV... | 14632 |
| 9 | NEGATIV... | 7 |
| 10 | POSITIVE_... | 8 |
| 11 | TOTALIT... | 263776 |
| 12 | TOTALIT... | 263776 |
| 13 | BEAUTY& | 21344 |

UTILITY CALCULATION NEXT

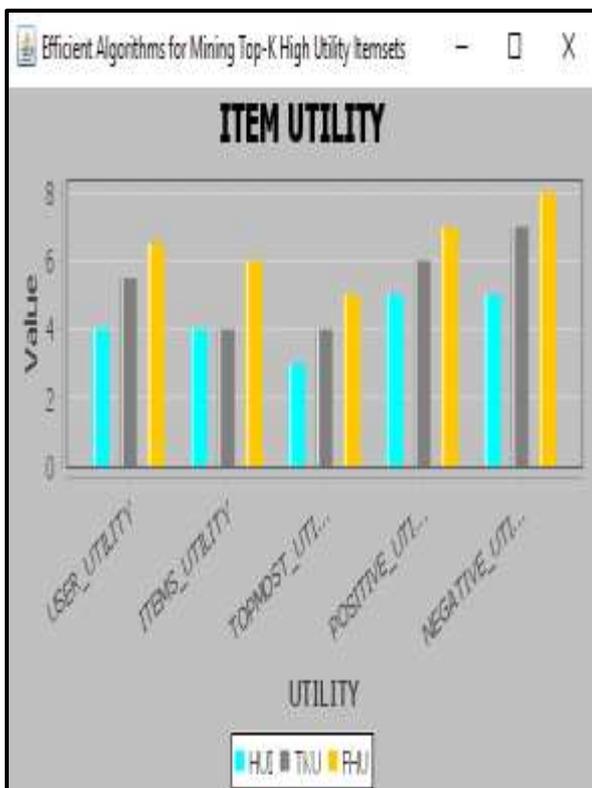
Input Window for UP Growth and UP Growth+ Algorithm

Transaction:

All Transactions:

Items:

Item Utility:



Message: 27

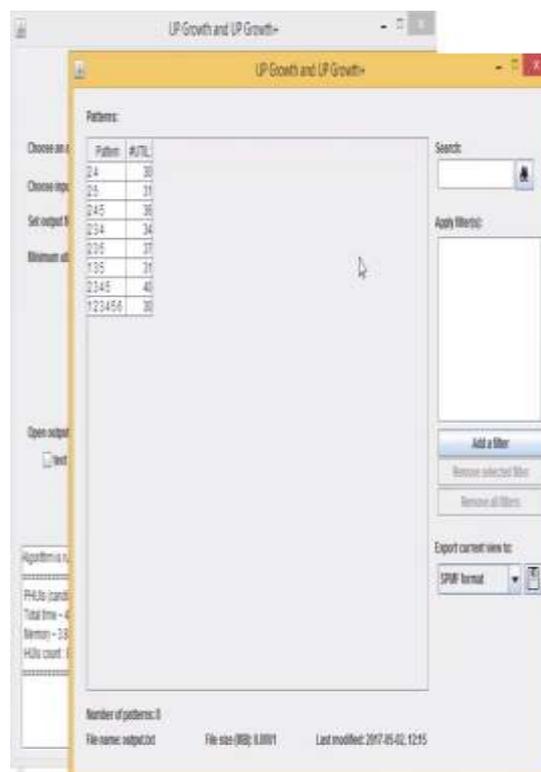
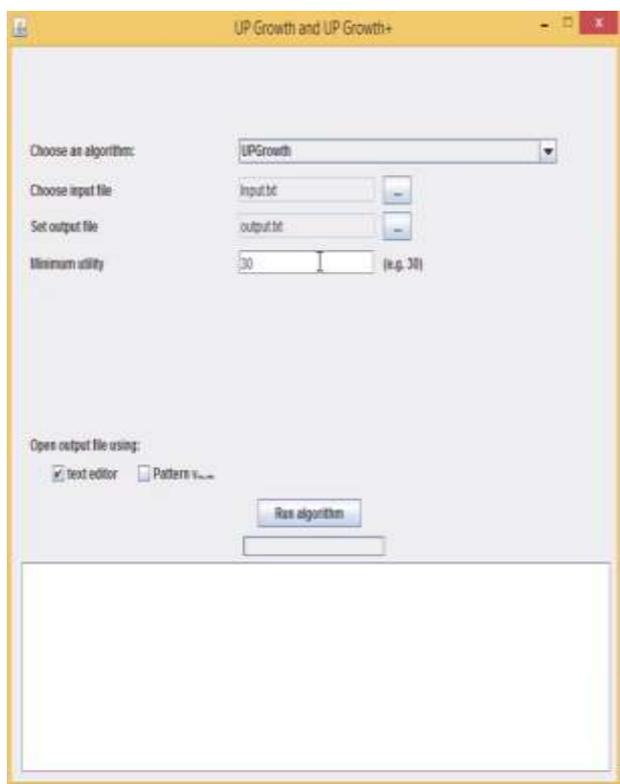
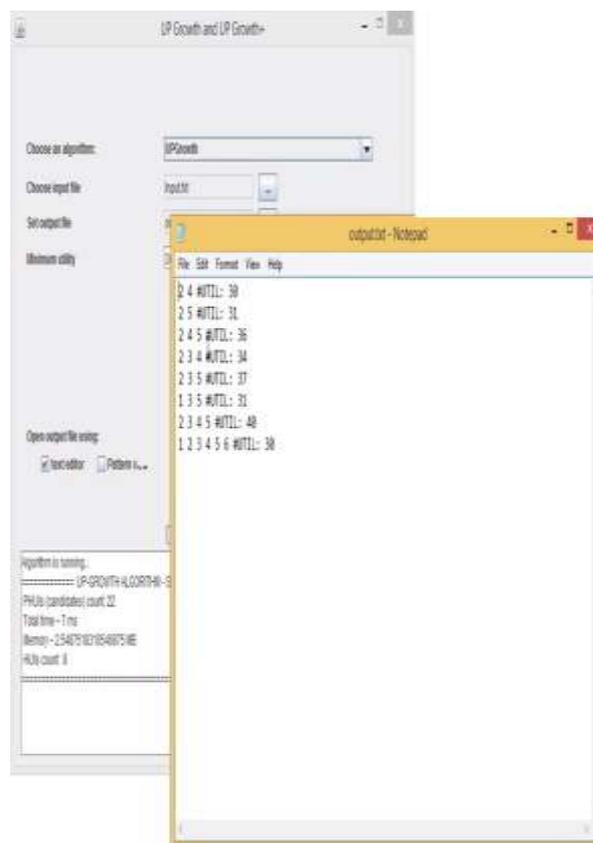
and UP Growth+ Algorithm

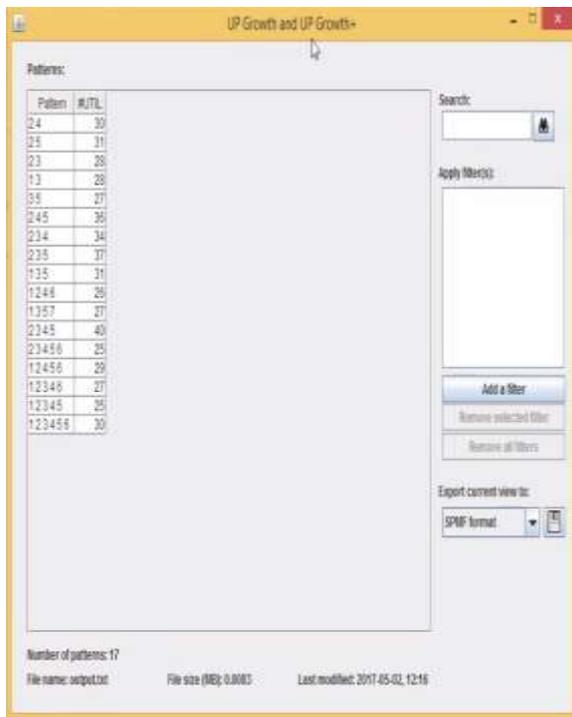
Transaction:

All Transactions:

Items:

Item Utility:





III. COMPARISON WITH UPGROWTH AND UPGROWTH+ ALGORITHM

In this project, two algorithms utility pattern growth (UP-Growth) and UP-Growth+, for mining high utility itemsets are used to compare. Performance of UP-Growth and UP-Growth+ become more efficient since database contain long

transactions and generate fewer number of candidates than FP-Growth.

Two important problems are always in consideration first is how minimize number no of candidates and another is how to remove space and time complexity. Also, choosing an appropriate minimum utility threshold is a difficult task for application users: if the threshold is high, there might be no HUI; if the threshold is low, there might result too many HUIs, and the mining performance might be severely affected, even leading to memory overflow. It would also be a time-consuming task if one tries to determine the threshold value through various testing calculations.

To address this issue, Wu [10] proposes top-k algorithm, mining the top k itemsets with the highest utility values without presetting the minimum threshold.

IV. CONCLUSION

In this paper, we have studied the problem of top-k high utility item sets mining, where k is the desired number of high utility item sets to be mined. Two efficient algorithms TKU (mining Top-K Utility item sets) and TKO (mining Top-K utility item sets in One phase) are proposed for mining such item sets without setting minimum utility thresholds. TKU is the first two-phase algorithm for mining top-k high utility item sets, which incorporates five strategies PE, NU, MD, MC and SE to effectively raise the border minimum utility thresholds and further prune the search space. On the other hand, TKO is the first one-phase algorithm developed for top-k HUI mining, which integrates the novel strategies RUC, RUZ and EPB to greatly improve its performance. Empirical evaluations on different types of real and synthetic datasets show that the proposed algorithms have good scalability on large datasets and the performance of the proposed algorithms is close to the optimal case of the state-of-the-art two-phase and one-phase utility mining algorithms. Although we have proposed a new framework for top-k HUI mining, it has not yet been incorporated with other utility mining tasks to discover different types of top-k high utility patterns such as top-k high utility episodes, top-k closed high utility item sets, top-k high utility web access patterns and top-k mobile high utility sequential patterns. These leave wide rooms for exploration as future work.

REFERENCES

- [1] R. Agrawal and R. Srikant, "Fast algorithms for mining association rules," in Proc. Int. Conf. Very Large Data Bases, 1994, pp. 487–499.
- [2] C. Ahmed, S. Tanbeer, B. Jeong, and Y. Lee, "Efficient tree structures for high-utility pattern mining in incremental databases," IEEE Trans. Knowl. Data Eng., vol. 21, no. 12, pp. 1708–1721, Dec. 2009.
- [3] K. Chuang, J. Huang, and M. Chen, "Mining top-k frequent patterns in the presence of the memory constraint," VLDB J., vol. 17, pp. 1321–1344, 2008.
- [4] R. Chan, Q. Yang, and Y. Shen, "Mining high-utility itemsets," in Proc. IEEE Int. Conf. Data Mining, 2003, pp. 19–26.
- [5] P. Fournier-Viger and V. S. Tseng, "Mining top-k sequential rules," in Proc. Int. Conf. Adv. Data Mining Appl., 2011, pp. 180–194.

- [6] P. Fournier-Viger, C. Wu, and V. S. Tseng, "Mining top-k association rules," in Proc. Int. Conf. Can. Conf. Adv. Artif. Intell., 2012, pp. 61–73.
- [7] P. Fournier-Viger, C. Wu, and V. S. Tseng, "Novel concise representations of high utility item sets using generator patterns," in Proc. Int. Conf. Adv. Data Mining Appl. Lecture Notes Comput. Sci., 2014, vol. 8933, pp. 30–43.
- [8] J. Han, J. Pei, and Y. Yin, "Mining frequent patterns without candidate generation," in Proc. ACM SIGMOD Int. Conf. Manag. Data, 2000, pp. 1–12.
- [9] J. Han, J. Wang, Y. Lu, and P. Tzvetkov, "Mining top-k frequent closed patterns without minimum support," in Proc. IEEE Int. Conf. Data Mining, 2002, pp. 211–218.
- [10] S. Krishnamoorthy, "Pruning strategies for mining high utility itemsets," Expert Syst. Appl., vol. 42, no. 5, pp. 2371–2381, 2015.
- [11] C. Lin, T. Hong, G. Lan, J. Wong, and W. Lin, "Efficient updating of discovered high-utility itemsets for transaction deletion in dynamic databases," Adv. Eng. Informat., vol. 29, no. 1, pp. 16–27, 2015.
- [12] G. Lan, T. Hong, V. S. Tseng, and S. Wang, "Applying the maximum utility measure in high utility sequential pattern mining," Expert Syst. Appl., vol. 41, no. 11, pp. 5071–5081, 2014.
- [13] 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.
- [14] M. Liu and J. Qu, "Mining high utility itemsets without candidate generation," in Proc. ACM Int. Conf. Inf. Knowl. Manag., 2012, pp. 55–64.