

ProUM: Projection-Based Utility Mining on Sequence Data

Wensheng Gan^{1,5}, Jerry Chun-Wei Lin^{1,2*}, Jiexiong Zhang¹, Han-Chieh Chao³, Hamido Fujita⁴ and Philip S. Yu⁵

¹*School of Computer Sciences and Technology, Harbin Institute of Technology (Shenzhen), Shenzhen, Guangdong 518055, China*

²*Department of Computing, Mathematics and Physics, Western Norway University of Applied Sciences (HVL), Bergen 5050, Norway*

³*Department of Electrical Engineering, National Dong Hwa University, Hualien 97401, Taiwan*

⁴*Faculty of Software and Information Science, Iwate Prefectural University, Morioka 020-8550, Japan*

⁵*Department of Computer Sciences, University of Illinois at Chicago, Chicago, IL 60607, USA*

Email: wsgan001@gmail.com, jerrylin@ieee.org, jiexiong.zhang@foxmail.com, hcc@ndhu.edu.tw, HFujita-799@acm.org, psyu@uic.edu

Abstract

Utility is an important concept in economics. A variety of applications consider utility in real-life situations, which has led to the emergence of utility-oriented mining (also called utility mining) in the recent decade. Utility mining has attracted a great amount of attention, but most of the existing studies have been developed to deal with itemset-based data. Time-ordered sequence data is more commonly seen in real-world situations, which is different from itemset-based data. Since they are time-consuming and require large amount of memory usage, current utility mining algorithms still have limitations when dealing with sequence data. In addition, the mining efficiency of utility mining on sequence data still needs to be improved, especially for long sequences or when there is a low minimum utility threshold. In this paper, we propose an efficient **Projection-based Utility Mining** (ProUM) approach to discover high-utility sequential patterns from sequence data. The utility-array structure is designed to store the necessary information of the sequence-order and utility. ProUM can significantly improve the mining efficiency by utilizing the projection technique in generating utility-array, and it effectively reduces the memory consumption. Furthermore, a new upper bound named sequence extension utility is proposed and several pruning strategies are further applied to improve the efficiency of ProUM. By taking utility theory into account, the derived high-utility sequential patterns have more insightful and interesting information than other kinds of patterns. Experimental results showed that the proposed ProUM algorithm significantly outperformed the state-of-the-art algorithms in terms of execution time, memory usage, and scalability.

Keywords: economics, utility mining, sequence, projection, sequential pattern

1. Introduction

In the era of big data, data mining and data analytics [9] are some of the fundamental technologies for discovering knowledge in data, and they have become more prevalent in our life due to the rapid growth of massive data [15]. Up to now, a handful of methods have been proposed for discovering useful and interesting patterns [3, 20] from different types of data. For example, frequent pattern mining (FPM) [20] and association rule mining (ARM) [3] from the transaction data have been extensively studied. One of the well-known applications of FPM and ARM is in market basket analysis. In addition, mining sequence data, which is common for many real-life applications, has also attracted a lot of attention. Sequential pattern mining (SPM) is one of the well-studied research fields for mining sequence data [2, 35, 33]. Knowing the useful patterns and auxiliary knowledge from sequences/events can benefit a number of applications, such as web access analysis, event prediction, time-aware recommendation, and DNA detection [11]. Up to

now, research has been conducted on mining interesting patterns from transaction or sequential data [11, 15, 20, 33]. However, most of them are based on the co-occurrence frequency of patterns.

Motivation. In reality, massive data, for example, a sequence, contains valuable but hidden auxiliary information. However, the measures of support [20] and confidence [2, 3] cannot be effectively utilized to discover implicit or potential information. For instance, implicit factors such as the utility, interest, risk, and profit of objects in the data are not considered in traditional ARM or SPM. The main goal of tasks for data mining and analytics is generally to achieve utility maximization. However, most of the existing algorithms of FPM, ARM, and SPM are unable to discover the valuable targeted patterns that benefit utility maximization. For example, support/frequency-based data mining models might be insufficient for achieving a time-aware recommendation for users based on the users' click-stream or purchase behavior. Specifically, a variety of applications consider utility. Typical examples include the profit of products in supermarkets and retail stores, the satisfaction feedbacks of different restaurants, and the popularity of hot showing movies.

*Corresponding author. Email: jerrylin@ieee.org.

Thus, in these circumstances, the frequency framework loses its adaptiveness.

Utility [32] is an important concept in Economics. The emergence of a new mining and computing framework, called utility mining, can be realized by taking utility theory [32] from economics into account [18]. Utility mining is in the cross-domain of information technology and economics. In the past decade, utility mining has been extensively studied in areas like high-utility itemset mining (HUIM) [6, 30, 31, 37], high-utility sequential pattern mining (HUSPM) [7, 41, 44], and high-utility episode mining (HUEM) [42]. It has been successfully applied to discover utility-driven knowledge in the cross-domain of information technology and business. The identified patterns, which can bring valuable profits for retailers or managers, are more useful than those frequent-based patterns in business. HUIM addresses the itemset-based data to mine high-utility itemsets (HUIs), and HUSPM deals with the sequence data to discover high-utility sequential patterns (HUSPs). In general, the *utility* can be a user-specified subjective measure, such as satisfaction, profit, risk, interest, and so on. Utility mining [18] has become an important branch of data science, which is aimed at utilizing the auxiliary information from data (e.g., itemsets, events, and sequences). Time-ordered sequence data is more commonly seen in real-world situations, which is different from itemset-based data. To a certain extent, itemset data is a special case of sequence data.

Challenges. Several prior studies have focused on improving the mining efficiency of HUSPM, such as USpan [44], HuspExt [7], PHUS [21], and HUS-Span [41]. Among them, USpan [44], HuspExt [7], and HUS-Span [41] are all based on a lexicographic sequence tree with two concatenation mechanisms and several pruning strategies w.r.t. upper bounds on utility. In general, some challenges remain in addressing the problem of HUSPM, which are described below.

First, the computing mechanism of the utility of a pattern is different from that of frequency of a pattern. The former is more complicated and the utility of a sequence is not downward closed. This means that the support-based pattern mining techniques and models cannot be directly applied to discover utility-driven patterns, and the pruning search space in HUSPM is more difficult.

Second, the HUSPM problem is intrinsically more complex than HUIM and FPM. HUSPM may easily face a critical combinatorial explosion of search space without powerful pruning strategies w.r.t. upper bounds on utility because of the inherent time order embedding in sequence data.

Third, a common way to identify the interesting HUSPs is to recursively generate the projected sub-databases and then scan these whole sub-databases. However, this is very inefficient and costly in memory when the number of sequences in processed database is large-scale. Therefore, how to efficiently reduce the size of the databases that need to be projected and scanned is a crucial problem to be solved for efficiently discovering HUSPs.

In summary, speeding up the execution time and reducing memory consumption without losing HUSPs are critical in HUSPM. How to improve the mining efficiency of utility min-

ing on sequence data is still an open problem.

Contributions. In light of the aforementioned challenges, we propose a novel utility mining framework, called **Projection-based Utility Mining** on sequence data (ProUM). Based on the developed utility-array with the projection mechanism, the utility-driven mining model, ProUM, can not only extract the insightful high-utility patterns but also achieve better efficiency. The effectiveness and efficiency of the proposed ProUM is evaluated by comparing it with the well-known frequency-based SPM model and state-of-the-art utility mining algorithms. The major contributions of this paper can be summarized as follows:

- We adopt utility and time-order significance as the key criterion for evaluating utility mining on sequence data. By considering the utility factor and time-order relations among items/objects, we design an efficient method, called **Projection-based Utility Mining** on sequence data (ProUM).
- A compact data structure, namely, utility-array, is presented to store the compact information (e.g., utility, position, and time order) of sequences from the processed sequence database. ProUM can quickly discover a set of high-utility sequential patterns based on the developed utility-array with the projection mechanism.
- This projection-based approach utilizes several pruning techniques in a depth-first search manner, which consist of the utilization of utility property and the proposed upper bound named sequence extension utility (*SEU*). Therefore, ProUM is able to filter a large number of unpromising patterns at an early stage and return the significant patterns in the mining process.
- Experiments on both real and synthetic datasets show that the proposed utility-array representation achieves a lossless compression capability of quantitative sequence data. Moreover, ProUM significantly outperforms the state-of-the-art algorithms, such as USpan and HUS-Span.

The remainder of this paper is organized as follows: Some related works of support-based sequence mining and utility-based mining on itemset-based data and sequence data are briefly reviewed in Section 2. Some basic preliminaries and the problem statement of HUSPM are given in Section 3. Details of the proposed data structure, pruning strategies with upper bound, and the main procedure of the ProUM algorithm are described in Section 4. The experimental evaluation of the proposed ProUM algorithm is provided in Section 5. Finally, conclusions are drawn in Section 6.

2. Literature Review

Much research has been conducted on frequency-based sequential pattern mining and utility mining of itemset-based data, but fewer works have integrated utility theory for mining high-utility patterns from event/sequence data. In this section, we separately present the prior works on SPM, HUIM, and HUSPM.

2.1. Frequency-Based Mining on Sequences

Pattern (i.e., itemset, rule, and sequence) mining [20, 33] is a kind of well-studied data mining and analytics model. The applications of pattern mining models are very extensive, and details can be referred to in the survey literature [11, 15, 18, 19]. A great effort has been put forth by the data mining community to discover frequent patterns from itemset-based data, such as Apriori [3] and FP-growth [20] methods. Different from itemset-based data, timely ordered sequence data is more commonly seen in the real-world, in areas such as traffic data, web access, customer shopping data, travel routes, stock market trends, DNA chains, and so on [2, 11, 33, 35]. The problem of sequential pattern mining (SPM) from sequence data was first presented by Agrawal and Srikant [2]. Frequent pattern mining (FPM) from itemset-based data is closely related to SPM [2, 11, 33, 35]. A number of algorithms have been proposed to discover the complete frequent sequential patterns from sequential databases, including SPADE [46], SPAM [8], and PrefixSpan [33]. These algorithms use many strategies to make the mining of sequential patterns more efficient and practical. Unfortunately, before discovering the final result sets, SPM may produce a huge amount of candidates, partially due to the combinatorial nature of the mining task and the timely ordered information embedding in the sequences. Analysis is difficult when dealing with long sequences because most of the SPM algorithms may generate an exponential number of sequences, especially when using a lower minimum support threshold. Among them, the well-known PrefixSpan [33] algorithm follows a pattern growth mechanism that uses a series of projected databases for achieving better mining performance in terms of execution time and memory cost. It recursively extracts the prefix sub-sequences, and then projects the postfix sub-sequences into the sub-databases [33]. Comprehensive overviews of sequential pattern mining have been given by Fournier-Viger et al. [11] as well as Gan et al. [19].

2.2. Utility-Driven Mining on Transaction Data

The key measure for discovering patterns in the aforementioned FPM and SPM is the frequency (aka relevant co-occurrence) [20]. In general, the statistical frequency is an objective measure while some subjective measures and useful factors (e.g., utility, business profit, risk, and preference) are ignored. Therefore, the support/frequency-based data mining approaches cannot return the real useful knowledge, which decreases the effectiveness of mining task. Up to now, a variety of applications have considered utility. A new mining and computing framework named utility mining [18] has been proposed by taking the utility theory from economics into account. Utility mining has been developed to successfully applied to discover utility-driven knowledge in many real-world applications. The early works of utility mining were related to high-utility itemset mining (abbreviated as HUIM) [6, 30, 31, 37], which addresses itemset-based data. In general, the *utility* can be any user-specified subjective measure, such as satisfaction, profit, risk, interest, and so on.

Many previous studies of utility mining have focused on developing efficient algorithms that can achieve better mining

performance, such as Apriori-like approaches (e.g., Two-Phase [31]), tree-based approaches (e.g., IHUP [6], UP-growth [39], UP-growth+ [37]), list-based approaches (e.g., HUI-Miner [30], FHM [12]), and other hybrid algorithms (e.g., EFIM [47]). In addition to efficiency, the effectiveness of the data mining and analytic models is also very important. Therefore, a number of studies have been developed to improve the effectiveness for mining utility-oriented patterns, and the current state-of-the-art approaches have been provided in [18]. For instance, Lin et al. studied the problem of dynamic high-utility itemset mining on different types of dynamic data with record insertion [18, 25], record deletion [26], and record modification [25]. In some applications, the collected data is not precise and contains uncertainty. There have been some interesting works that deal with uncertain data for mining high-utility patterns [24]. At the same time, other interesting issues of utility mining also have been studied, such as utility mining with discount strategies [26] or negative values [22], discovering top- k high-utility patterns [38], correlated utility mining [13, 17], and HUIM in big data [29]. Recently, a new utility measure, called utility occupancy [16], has been proposed to solve the drawbacks of the existing utility mining models.

Different from the above-mentioned approaches, there are several genetic algorithms (e.g., HUIM-BPSO [27] and ACO-based HUIM-ACS [43]) methods that have been applied to deal with the utility mining problem. However, evolutionary computation techniques for HUIM do not provide any benefit for improving the mining efficiency. Besides, the interesting topic called privacy preserving utility mining [14] also has been extensively studied. A detailed survey of current development of HUIM was reported by Gan et al. [18].

2.3. Utility-Driven Mining on Sequences

In addition to itemset-based data, sequence data also has been addressed in utility mining, which is called high-utility sequential pattern mining (HUSPM) [7, 41, 44]. In FPM and SPM, the Apriori property [2, 3, 35] is widely adopted as the downward closure property to prune the search space. However, the Apriori property does not hold in HUSPM, and this makes the analysis of HUSPM difficult. Due to the absence of the downward closure property in sequence utility, the sequence-weighted utilization (*SWU*) [4, 44] is utilized in HUSPM to prune the search space. It has been proved that the *SWU* value is an upper bound of the utilities of a sequence and all its super-sequences. For extracting high-utility sequential patterns, Ahmed et al. [5] first proposed two algorithms, UtilityLevel and UtilitySpan. UtilityLevel is an Apriori-like algorithm, and UtilitySpan is based on PrefixSpan [33]. UL and US discover HUSPs in two phases by using *SWU* to prune the search space. They first find the candidate sequences with a high *SWU* value, then they compute the actual utilities of each sequence in candidates, and finally all the HUSPs can be identified. Next, the UMSP algorithm [34] was developed to discover high-utility mobile sequences, and the UWAS-tree and IUWAS-tree [4] were designed to find high-utility web log sequences.

Unfortunately, all these algorithms only consider single-item sequences (i.e., itemset data that was addressed in HUIM) but

not the element-based sequences. Yin et al. [44] presented a generic definition of the HUSP mining framework and proposed a new mining algorithm named USpan. In the USpan model, the quantitative sequences, along with their utility and time order information, are represented as the utility-matrix structure. Then, two upper bounds (*SWU* and sequence-projected utilization (*SPU*) [44]) are applied to prune the search space, which is represented as a lexicographic tree. To prune the lexicographic tree for a better efficiency, *SWU* is utilized in a depth pruning strategy and *SPU* is utilized in a breath pruning strategy. However, the high-utility sequential patterns that are mined by USpan are not complete. Alkan et al. proposed HuspExt [7] with a Cumulate Rest of Match (CRoM) based pruning technique to improve the mining performance. HuspExt also utilizes an upper bound on the utilities of the candidate sequences to prune the search space. Lan et al. [21] further introduced a projection-based PHUS algorithm to mine HUSPs using a sequence-utility upper-bound (*SUUB*) model. The maximum utility measure is introduced in the *SUUB* model to obtain a tighter upper bound on the utility of a sequence.

The above algorithms (e.g., USpan [44] and HuspExt [7]) adopt *SWU* to prune the search space, however, they usually suffer from the problem of an exponential number of candidate sequences, especially when the user-defined minimum utility threshold is small. The HUS-Span algorithm [41], which was proposed recently, utilizes a new upper bound, called the prefix extension utility (*PEU*). Although the results discovered by HUS-Span are complete, HUS-Span is not efficient enough. The generate-and-test approach creates an overflow of candidate sequences. Recently, some interesting issues of HUSPM have been extensively studied that can improve the effectiveness of mining high-utility sequential patterns. For example, the problems of mining top-*k* high-utility sequential patterns [45, 41], discovering periodic HUSPs [10], mining HUSPs with multiple minimum utility thresholds [28], and incrementally mining HUSPs on a dynamic database [40] have been addressed. It should be noted that several genetic algorithms have been developed for HUIM (e.g., HUIM-BPSO [27] and HUIM-ACS [43]), but they have not been proposed to deal with HUSPM yet. More current development of HUSPM can be referred to in literature reviews [14, 18, 36].

3. Preliminaries and Problem Formulation

This section introduces some basic concepts and principles of sequence mining and utility-oriented sequence mining. Some definitions from prior works are adopted to present the problem clearly. More details about the background of sequence data and sequence mining can be found in [11, 33, 46].

3.1. Sequence Data

Let $I = \{i_1, i_2, i_3, \dots, i_n\}$ be a set of distinct items. An *itemset* X is a subset of items, that is, $X \subseteq I$. A sequence s is an ordered list of *itemsets* (also called *elements* or *events*). Note that the items within each element can be unordered, without loss of generality, and it is assumed that they are sorted alphabetically. Additionally, “<” is used to represent that one item

occurs before another item. In SPM, an item occurs once at most in an element of a sequence. An item can occur multiple times (also called occurred quantity) in an element while in HUSPM. A group of the sequences stored with their identifiers (*sid*) is called a sequence database, denoted as D . Thus, a sequence database $D = \{s_1, s_2, \dots, s_m\}$ is a set of sequences/tuples (sid, e_1, \dots, e_j), where *sid* is a sequence id and e_j is an element that contains a set of items belonging to I .

The total number of items in a sequence is called its length. A sequence with length k is called a k -sequence. The size of a sequence is the number of itemsets/elements within this sequence. For example, a sequence $s = \langle [ac], [abc], [abd], [ef] \rangle$ consists of six distinct items and four elements. Thus, s is called 10-sequence since the length of s is ten, and its size is four. A k -itemset, also called k - q -itemset, is an itemset that contains exactly k items. A k -sequence (k - q -sequence) is a sequence having k items.

Definition 1. (*sub-sequence and sup-sequence*) Given two sequences, $\alpha = \langle a_1, \dots, a_m \rangle$ and $\beta = \langle b_1, \dots, b_n \rangle$, α is called a sub-sequence of β iff¹ each a_j ($1 \leq j \leq m$) can be mapped by b_{i_j} ($a_j \subseteq b_{i_j}$) and preserves the order as $1 \leq i_1 < i_2 < \dots < i_m \leq n$. In other words, if α or β contains α , then β can be called a super-sequence.

For example, in considering three sequences, $\alpha = \langle [ac], [c], [abd] \rangle$, $\beta = \langle [ac], [d] \rangle$, and $\gamma = \langle [cd] \rangle$, α is said to be a super-sequence of β while γ is not a sub-sequence of α . The reason for this is that the sequence $\gamma = \langle [cd] \rangle$ cannot be mapped to any sequence in α .

Definition 2. (*sequence containing*) Given two sequences, t_1 and t_2 , then t_2 uniquely contains t_1 iff there is only one $t_1 \subseteq t_2$ such that $t_1 = t_2$, denoted as $t_1 \sqsubseteq t_2$. Similarly, for a sequence t and a sequence s , s uniquely contains t , denoted as $t \sqsubseteq s$, iff there is only one s' and $s' \subseteq s$ such that $s' \sim t$.

A quantitative sequential database (shown in Table 1), is used as a running example in this paper. Table 1 has five sequences/transactions and six items. In addition, each item i_j in D is associated with a unit utility (also called *external utility*), which is denoted as $pr(i_j)$. The unit utility (e.g., price and profit) for each item is provided in Table 2, which can be called the *profit-table*. In general, the profit-table is based on the prior knowledge of similar users or contents. In the running example, the unit utility of an item (e) is \$6.

Definition 3. (*quantitative sequence*) For the addressed HUSPM problem, the processed database is the quantitative sequence database (q -database) that each item $i_j \in I$ ($1 \leq j \leq n$) in an element/itemset v is associated with a *quantity* (also called *internal utility*), denoted as $q(i_j, v)$. For convenience, “ q ” is used to refer to the object associated with quantity throughout this paper. Thus, the term “ q -sequence” means a sequence with quantities, and “sequence” means a sequence without quantities. Similarly, the “ q -itemset” means an itemset having quantities while an “itemset” does not have quantities.

¹In this paper, the term “iff” means “if and only if”.

Table 1: A quantitative sequence database

SID	Q-sequence
S_1	$\langle [(a:2) (c:1)], [(c:2)], [(b:10) (f:3)], [(a:2) (e:1)] \rangle$
S_2	$\langle [(f:2)], [(a:5) (d:2)], [(c:2)], [(b:4)], [(a:4) (d:1)] \rangle$
S_3	$\langle [(a:4)], [(b:4)], [(f:5)], [(a:1) (b:2) (e:1)] \rangle$
S_4	$\langle [(a:3) (b:4) (d:5)], [(c:2) (e:1)] \rangle$
S_5	$\langle [(b:1) (e:1)], [(c:1)], [(f:2)], [(d:2)], [(a:4) (e:2)] \rangle$

Table 2: A profit-table

Item	a	b	c	d	e	f
Profit	\$3	\$2	\$10	\$4	\$6	\$1

3.2. Utility Mining on Sequence Data

Utility mining incorporates the utility theory and mining techniques to deal with complex data, such as quantitative sequence data. Some definitions in the utility framework on sequence data are briefly introduced below.

Definition 4. (*utility of q -item*) Let $q(i_j, v)$ be the quantity of (i_j) in a q -itemset v , and $pr(i_j)$ be the unit profit of (i_j). The utility of a q -item (i_j) in a q -itemset v is denoted as $u(i_j, v)$ and defined as:

$$u(i_j, v) = q(i_j, v) \times pr(i_j). \quad (1)$$

Definition 5. (*utility of q -itemset*) The utility of a q -itemset v is denoted as $u(v)$ and defined as:

$$u(v) = \sum_{i_j \in v} u(i_j, v). \quad (2)$$

Definition 6. (*utility of q -sequence*) The utility of a q -sequence $s = \langle v_1, v_2, \dots, v_d \rangle$ is defined as:

$$u(s) = \sum_{v \in s} u(v). \quad (3)$$

Definition 7. (*utility of q -database*) The utility of a quantitative sequential database D is the sum of the utility of each of its q -sequences:

$$u(D) = \sum_{s \in D} u(s). \quad (4)$$

For instance, consider the running example in Table 1. The utility of the item (a) in the first q -itemset in S_1 is calculated as: $u(a, [(a:2) (c:1)]) = q(a, [(a:2) (c:1)]) \times pr(a) = 2 \times \$3 = \$6$. In addition, the utility of the first q -itemset $\langle [(a:2) (c:1)] \rangle$ is $u([(a:2) (c:1)]) = u(a, [(a:2) (c:1)]) + u(c, [(a:2) (c:1)]) = 2 \times \$3 + 1 \times \$10 = \16 . We have that $u(S_1) = u([(a:2) (c:1)]) + u([(c:2)]) + u([(b:10) (f:3)]) + u([(a:2) (e:1)]) = \$16 + \$20 + \$23 + \$12 = \71 . Therefore, the overall utility in Table 1 is $u(D) = u(S_1) + u(S_2) + u(S_3) + u(S_4) + u(S_5) = \$71 + \$69 + \$38 + \$63 + \$52 = \$293$.

Definition 8. (*sequence matching*) Given a q -sequence $s = \langle v_1, v_2, \dots, v_d \rangle$ and a sequence $t = \langle w_1, w_2, \dots, w_{d'} \rangle$, if $d = d'$ and the items in v_k are the same as the items in w_k for $1 \leq k \leq d$, then t matches s , which is denoted as $t \sim s$.

For instance, in Table 1, the sequences $\langle [ac] \rangle$, $\langle [ac], [b] \rangle$, $\langle [a], [b], [e] \rangle$ all match S_1 . A sequence in a q -sequence database may have more than one match in a q -sequence. For instance, $\langle [a], [b] \rangle$ has two matches in S_3 , such as $\langle [a:4], [b:4] \rangle$ and $\langle [a:4], [b:2] \rangle$. The measure of the utility of sequences for HUSPM is more challenging than that for SPM and HUIM due to the multiple matching cases.

Definition 9. (*q -itemset containment*) Given two itemsets w and w' , the itemset w is contained in w' (denoted as $w \subseteq w'$) if w is a subset of w' or w is the same as w' . Given two q -itemsets v and v' , v is said to be contained in v' if for any item in v there exists the same item having the same quantity in v' . This is denoted as $v \subseteq v'$.

Definition 10. (*q -sequence containment*) Given two sequences $t = \langle w_1, w_2, \dots, w_d \rangle$ and $t' = \langle w'_1, w'_2, \dots, w'_{d'} \rangle$, the sequence t is contained in t' (denoted as $t \subseteq t'$) if there exists an integer sequence $1 \leq k_1 \leq k_2 \leq \dots \leq d'$ such that $w_j \subseteq w'_{k_j}$ for $1 \leq j \leq d$. In addition, consider two q -sequences $s = \langle v_1, v_2, \dots, v_d \rangle$ and $s' = \langle v'_1, v'_2, \dots, v'_{d'} \rangle$. We say s is contained in s' (denoted as $s \subseteq s'$) if there exists an integer sequence $1 \leq k_1 \leq k_2 \leq \dots \leq d'$ such that $v_j \subseteq v'_{k_j}$ for $1 \leq j \leq d$. In the following, $t \subseteq s$ is used to indicate that $t \sim s_k \wedge s_k \subseteq s$ for convenience.

For example, the itemset $[ab]$ is contained in the itemset $[abe]$, while $[abe]$ does not contain the item $[f]$. The q -itemset $[(a:1) (b:2)]$ is contained in $[(a:1) (b:2) (e:1)]$, but it is not contained in $[(a:3) (b:4) (d:5)]$. In Table 1, S_3 contains $[(a:1) (b:2)]$, but S_4 does not contain it. Consequently, the sequences $\langle [(a:4)], [(b:4)] \rangle$ and $\langle [(a:4)], [(b:2)] \rangle$ are contained in S_3 , but $\langle [(a:3)], [b:2] \rangle$ is not contained in S_3 .

It should be noted that the definition of utility of sequence originally proposed in [4, 5] is too specific. Therefore, the later studies [41, 44] have adopted “the maximum utility of all occurrences of a sequence t in a q -sequence s ” as the real utility of t in s . The proposed model in this paper follows this definition of utility.

Definition 11. (*maximal utility of t in s*) Consider a sequence t and a q -sequence s . The utility of t in s , denoted as $u(t, s)$, may have different utility values. The maximum utility is chosen among these utility values as the utility of t in s , as defined below:

$$u(t, s) = \max\{u(s_k) | t \sim s_k \wedge s_k \subseteq s\}. \quad (5)$$

Definition 12. (*utility of a sequence in D*) Let $u(t)$ denote the overall utility of a sequence t in a quantitative sequential database D . It is defined as:

$$u(t) = \sum_{t \subseteq s \wedge s \in D} u(t, s). \quad (6)$$

For instance, consider two sequences $\langle [a], [b] \rangle$ and $\langle [f], [ad] \rangle$ in Table 1. $\langle [a], [b] \rangle$ has two utility values in S_3 , and thus $u(\langle [a], [b] \rangle, S_3) = \max\{u(\langle [a:4], [b:4] \rangle), u(\langle [a:4], [b:2] \rangle)\} = \max\{\$20, \$16\} = \20 . $\langle [f], [ad] \rangle$ also has two utility values in S_2 , such that $u(\langle [f], [ad] \rangle, S_2) = \max\{u(\langle [f:2], [(a:5) (d:2)] \rangle), u(\langle [f:2], [(a:5) (d:2)] \rangle)\} = \max\{\$23, \$12\} = \23 .

$(d:2)\rangle$, $u(\langle [f:2], [(a:4) (d:1)] \rangle) = \max\{\$25, \$18\} = \25 . Intuitively, the calculation of the overall utility of a sequence in the sequence database is quite a bit more complicated than that of HUIM and SPM.

Therefore, the overall utility of $\langle [a], [b] \rangle$ in Table 1 can be obtained as $u(\langle [a], [b] \rangle) = u(\langle [a], [b] \rangle, S_1) + u(\langle [a], [b] \rangle, S_2) + u(\langle [a], [b] \rangle, S_3) = \$26 + \$23 + \$20 = \$69$. Notice that $\langle [a], [b] \rangle$ is different for $\langle [b], [a] \rangle$ because HUSPM considers the time orders embedding in sequences. Thus, $\langle [b], [a] \rangle$ is contained in S_5 while $\langle [a], [b] \rangle$ is not contained in S_5 .

3.3. Problem Definition

Definition 13. (*high-utility sequential pattern, HUSP*) In a quantitative sequential database D , a sequence t is said to be a high-utility sequential pattern (denoted as *HUSP*) if its overall utility in D satisfies:

$$HUSP \leftarrow \{t | u(t) \geq \delta \times u(D)\}. \quad (7)$$

where δ is the minimum utility threshold δ (usually given as a percentage).

In the running example, it is assumed that δ is set as 25%, and then $\delta \times u(D) = 25\% \times \$293 = \$73.25$. Thus, since $u(\langle [a], [b] \rangle) = \$69 < \$73.25$, $\langle [a], [b] \rangle$ is not a HUSP. Based on the above-stated concepts, the formal definition of the utility mining on sequence data (also called high-utility sequential pattern mining) problem can be defined below.

Problem Statement: Given a quantitative sequential database D (with a profit-table) and a user-defined minimum utility threshold δ , the utility-driven mining problem of high-utility sequential pattern mining (HUSPM) consists of enumerating all HUSPs whose overall utility values in this database are no less than the prespecified minimum utility account, such as $\delta \times u(D)$.

Therefore, the goal of HUSPM is to search for the set of sequences that achieves the highest utility score and their utility values are not less than the minimum utility value.

4. Proposed Utility Mining Algorithm: ProUM

This section describes the proposed projection-based ProUM algorithm for discovering high-utility sequence-based patterns by recursively projecting the utility-array based on the prefix sequences. ProUM utilizes the utility-array data structure to avoid multiple scans of the original database and projecting of the sub-databases. Only the compact utility-array is needed to be projected and scanned in each mining process. The framework of the proposed ProUM algorithm is presented in Figure 1. First, details of the search space, the utility-array structure, and the projection mechanism are presented below.

4.1. Lexicographic Sequence Tree

According previous studies (e.g., SPAM [8], PrefixSpan [33], USpan [44]), and the complete search space of SPM and HUSPM can be represented abstractly as the lexicographic sequence tree [8]. For the addressed problem for mining high-utility sequential patterns, a lexicographic q -sequence tree (LQS-tree) that

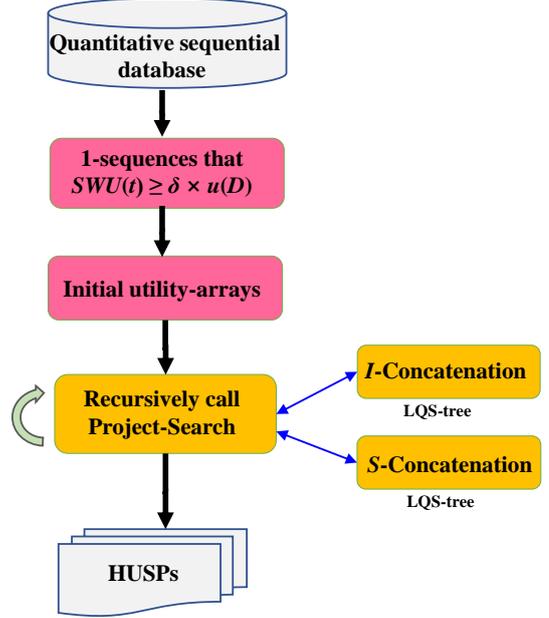


Figure 1: Framework of the ProUM algorithm

was used in USpan [44] is adopted to present the search space of ProUM.

Definition 14. (*lexicographic q -sequence tree*) A lexicographic q -sequence tree (LQS-tree) is a variant of the lexicographic tree structure [8] satisfying the following conditions:

- Each node in LQS-tree is associated with a q -sequence along with the utility of the sequence while the root is empty and labeled with “ $\langle root \rangle$ ”.
- All the nodes except the root in LQS-tree are the prefix-based extension of its parent node; in other words, any node’s child is an extension sequence/node of the node itself.
- All the children of any node in LQS-tree are listed in an increasing order, for example, lexicographic order.

It should be noted that the LQS-tree is just a conceptual structure, and the real visited/searched space may be different depending on different cases. Additionally, if $\delta = 0$, then the complete set of the found high-utility sequential patterns from the sequence data is equal to a complete LQS-tree, which covers the complete search space. It is important to notice that the utility property in LQS-tree is different from the frequent/support property in the previous lexicographic tree [8, 33].

Definition 15. (*I-Concatenation and S-Concatenation*) There are two operations, called *I-Concatenation* and *S-Concatenation*, to generate new sequences based on prefix node in LQS-tree.

- Given a sequence t and an item i_j , the *I-Concatenation* of t with i_j consists of appending i_j to the last itemset of t , denoted as $\langle t \oplus i_j \rangle_{I-Concatenation}$.

- An *S-Concatenation* of t with an item i_j consists of adding i_j to a new itemset appended after the last itemset of t , denoted as $\langle t \oplus i_j \rangle_{S\text{-Concatenation}}$.

For example, given a sequence $t = \langle [a], [c] \rangle$ and a new item (d) , $\langle t \oplus c \rangle_{I\text{-Concatenation}} = \langle [a], [cd] \rangle$ and $\langle t \oplus c \rangle_{S\text{-Concatenation}} = \langle [a], [c], [d] \rangle$, then based on the definition of sequence length and sequence size, the *I-Concatenation* does not increase the length of sequence t while the *S-Concatenation* increases the length of sequence t since the number of itemsets in t increases by one. All the candidate sequences in the search space w.r.t. LQS-tree can be enumerated for the purpose of mining HUSPs based on the two prefix-based operations.

Without loss of generality, it is assumed that the concatenated nodes in LQS-tree are listed in alphabetical order. In the running example in Table 1, the sequence in S_1 is listed as $\langle [(a) (c)], [(c) (f)], [(a) (e)] \rangle$ instead of $\langle [(c) (a)], [(c)], [(f) (b)], [(e) (a)] \rangle$. The expression of a sequence is unique using such a convention.

It should be noted that each LQS-tree node represents a candidate of the search space of HUSPs. A part of the LQS-tree can be referred to in [28, 44]. All the child nodes of a parent node are assumed to be ordered lexicographically with *I-extension*² sequences before the *S-extension*³ sequences. The most straightforward way for spanning the tree structure is to traverse and determine the lexicographic tree node one-by-one using either the Depth-First-Search (DFS) strategy or the Breadth-First-Search (BFS) strategy. Three questions remain to be answered: (1) how to compact the necessary information (e.g., utility, order of sequences, and position in each sequence) from the sequence database into the LQS-tree; (2) how to effectively traverse the nodes in LQS-tree and then quickly calculate their utility values; and (3) how to effectively prune the search space without spanning the complete LQS-tree. To address these problems, we propose a new data structure, utility-array, and several pruning strategies, which are respectively described as the following subsections.

4.2. Utility-Array and Projection Mechanism

Real sequences may often be very long in some application scenarios; for example, web search sequences, DNA, network intrusion access/log, and numerous other sequence data are all complex and long. This may easily lead to a huge search space for HUSPM. By adopting the concept of remaining utility [44, 41], we propose a new compact data structure, called utility-array, for storing the necessary information from the sequence data, including the utility information of item/sequence as well as their time/sequence order.

Definition 16. (*remaining utility* [44, 41]) Given a sequence t and a sequence database D , the remaining utility of t in a q -sequence s is the overall utilities of all items whose positions

²Notice that the terms *I-Concatenation* and *I-extension* are used interchangeably in this paper

³Notice that the terms *S-Concatenation* and *S-extension* are used interchangeably in this paper

are after t in s , and defined as: $u_{rest}(t, s) = \max\{u_{rest}(t, p_k, s)\}$, where p_k is the k -position, and $u_{rest}(t, p_k, s) = \sum_{i' \in s \wedge t \prec i'} u(i')$. Thus, the overall remaining utility of t in D is defined as: $u_{rest}(t) = \sum_{s \in D} u_{rest}(t, s)$.

Basically, the remaining utility of a sequence means the sum of the utilities after this sequence. Intuitively, the remaining utility is based on matching position. For example, a sequence $\langle s \rangle$ has two remaining utility values in S_1 : $u_{rest}(\langle a \rangle, p_1, s_1) = \$10 + \$20 + (\$20 + \$3) + (\$6 + \$6) = \65 and $u_{rest}(\langle a \rangle, p_2, s_1) = \6 . In addition, the remaining utility of $\langle [a], [b] \rangle$ in s_1 is $u_{rest}(\langle [a], [b] \rangle, s_1) = \$3 + (\$6 + \$6) = \$15$ and $u_{rest}(\langle [a], [b] \rangle) = u_{rest}(\langle [a], [b] \rangle, s_1) + u_{rest}(\langle [a], [b] \rangle, s_2) + u_{rest}(\langle [a], [b] \rangle, s_3) = \$15 + \$16 + \$36 = \$67$. Based on the concept of remaining utility, we introduce a data structure to represent the necessary information (both the utility and sequence order w.r.t. position) of each q -sequence.

Definition 17. (*utility-array*) Suppose that all the items in a q -sequence s in a q -sequence database have different unique occurred positions are $\{p_1, p_2, \dots, p_k\}$, where $\{p_1 < p_2 < \dots < p_k\}$, and the total number of positions is equal to the length of s . The utility-array of a q -sequence $s = \langle e_1, e_2, \dots, e_n \rangle$ (e_n is the n -element in s) consists of a set of arrays from left to right in s . Each array is related to an item i_j in each position p_k and contains the following fields: $array_{p_k} = [eid, item, u, ru, next_pos, next_eid]$. Details are given below:

- Field *eid* is the element ID of an element containing i_j ;
- Field *item* is the name of item i_j ;
- Field *u* is the actual utility of i_j in position p_k ;
- Field *ru* is the remaining utility of i_j in position p_k ;
- Field *next_pos* is the next position of i_j in s ;
- Field *next_eid* is the position of the first item in next element ($eid+1$) after current element (eid).

In addition, the utility-array records the first occurred position of each distinct item in s . In summary, a utility-array of a q -sequence s is a set of arrays related to each item in s and contains the position, utility, and sequence order information.

In the running example in Table 1, the constructed utility-array can be described in Table 3, and it can be seen that the first occurring position of each distinct item in s is also recorded in the utility-array, for example, the first occurring positions of a and f are 1 and 5, respectively. Intuitively, the utility-array contains all the necessary information of each sequence t , including not only the utility values of each item in each position/element but also the sequence order and position information⁴.

Definition 18. (*position in utility-array*) Each array has a unique position as an index in the designed utility-array structure. Moreover, the size of the arrays in the utility-array of a q -sequence s

⁴Note that “-” means the position is empty.

Table 3: The utility-array structure of S_1

	<i>eid</i>	<i>item</i>	<i>u</i>	<i>ru</i>	<i>next_pos</i>	<i>next_eid</i>
<i>array</i> ₁	1	<i>a</i>	\$6	\$65	6	3
<i>array</i> ₂	1	<i>c</i>	\$10	\$55	3	3
<i>array</i> ₃	2	<i>c</i>	\$20	\$35	-	4
<i>array</i> ₄	3	<i>b</i>	\$20	\$15	-	6
<i>array</i> ₅	3	<i>f</i>	\$3	\$12	-	6
<i>array</i> ₆	4	<i>a</i>	\$6	\$6	-	-
<i>array</i> ₇	4	<i>e</i>	\$6	\$0	-	-

is equal to the length of s . This position represents a match of an item i_j in s and can be used as an index for quickly retrieving the utility-array and calculating the detailed information of i_j .

In the implementation details, an array is used to store the set of information of the compact utility-array. Thus, position pos_k indexes $array_{pos_k}$. For example, in Table 3, the $array_1$ is indexed by position 1, and the $array_2$ is indexed by position 2. First, the 1-sequences with the low SWU values in each sequence are removed when calculating the compact utility-array of original sequence database D . Then, the new (revised) transactions are used to construct the initial utility-arrays. Inspired by the database projection idea of PrefixSpan [33], we present the following prefix-projected and span mechanism in utility-array.

Definition 19. (*prefix, suffix, and projection [33]*) Assume all the items in an element of a sequence database D are listed alphabetically. Given two sequences, $\alpha = \langle e_1, e_2, \dots, e_n \rangle$, and $\beta = \langle e'_1, e'_2, \dots, e'_m \rangle$ ($m \leq n$), β is called a *prefix* of α iff it meets the following conditions: (1) $e'_i = e_i$ for $i \leq m - 1$; (2) $e'_m \subseteq e_m$; and (3) all the items in $(e_m - e'_m)$ are alphabetically after those in e'_m . Additionally, the remaining part/elements after the prefix β in a sequence are called *suffix* with regards to prefix β . Let α be a sequence in D , then the α -projected sub-database, denoted as $D|_{\alpha}$, is the collection of suffixes of the sequences (which contains α) in D with regards to prefix α .

The ProUM algorithm recursively partitions the processed utility-array based on the projection mechanism [33]. Specifically, each subset of the extracted sequential patterns are further divided when necessary. Thus, the projection mechanism [33] forms a divide-and-conquer framework. Correspondingly, ProUM recursively constructs the corresponding projected utility-arrays but not the projected sub-databases.

Definition 20. (*projected utility-array*) Let t be a sequence in D , and the utility-array of t is denoted as $t.ua$. The t -projected utility-array, denoted as $(D.ua)|_t$, is the collection of *suffix* of arrays in $D.ua$ w.r.t. prefix t .

For example, the projected utility-array of $\langle [a] [c] \rangle$ in S_1 is shown in Table 4. Similarly, the other utility-arrays of $\langle [a] [c] \rangle$ can be projected in other sequences, for example, S_2 and S_4 . As an accurate representation, utility-array provides provably equivalent decomposition as projected database from the original sequence data, but it requires much less memory space. It

Table 4: The projected utility-array of $\langle [a] [c] \rangle$ in S_1 .

	<i>eid</i>	<i>item</i>	<i>u</i>	<i>ru</i>	<i>next_pos</i>	<i>next_eid</i>
<i>array</i> ₄	3	<i>b</i>	\$20	\$15	-	6
<i>array</i> ₅	3	<i>f</i>	\$3	\$12	-	6
<i>array</i> ₆	4	<i>a</i>	\$6	\$6	-	-
<i>array</i> ₇	4	<i>e</i>	\$6	\$0	-	-

proceeds by dividing the initial utility-arrays into smaller ones projected on the subsequences that were obtained so far, and only their corresponding suffixes are kept. The number of transactions/sequences in the projected utility-array is less than original database. This can substantially reduce the cost of the projection operation when the projected utility-arrays can be held in the main memory. By combining sequences from a series of projected utility-arrays, all the HUSPs and the candidates can be acquired.

It is important to notice that, instead of constructing the projected sub-database that only contain the updated sequences, a set of the projected compact utility-array of each sequence is only constructed $s \in D$. In other words, only the projected utility-arrays in each projection process are constructed and then scanned for constructing the next updated ones. Different from previous studies [41, 44] that require scanning the projected sub-database to construct the data structure (e.g., utility-matrix [44] and utility-chain [41]), the proposed ProUM algorithm does not need to construct and scan the projected sub-database.

Based on the designed utility-array and its construction process, the following desirable properties of the utility-array can be obtained: (1) The obtained information from utility-array is exact. Since the utility-array of $(l+1)$ -sequence is constructed based on the built utility-array of l -sequences, it is parameter-free and contains the complete information. (2) It is space efficient because it requires an inconsequential space overhead to construct and project a series of utility-arrays, which allows massive sequences to be processed in main memory (for most data mining, disk is death). (3) It has simplicity and intuitiveness because it can be constructed in deterministic time and regarded as the representation of quantitative sequences.

4.3. Proposed Upper Bound and Pruning Strategies

It is known that the complete search space of SPM or HUSPM is much more difficult than FIM. For HUSPM, it has $2^{m \times n}$ possible candidates in total, where m is the total number of all the possible distinct items in the database, and n is the number of elements in the longest sequence. However, in HUSP mining, the downward closure property (e.g., the Apriori property [2, 35]) does not hold for the utility of sequence patterns. Because SWU has the downward closure property, the current algorithms (e.g., USpan [44] and HuspExt [7]) adopt SWU to prune the search space. However, they usually suffer from the problem of an exponential number of candidate sequences since SWU is a loose upper bound to over-estimate the true utility of a sequence.

An optimization with a new upper bound is proposed below to further improve ProUM's efficiency. The search space of ProUM can be systematically explored by utilizing the presented pruning strategies.

Definition 21. The sequence-weighted utilization (*SWU*) [44] of a sequence t in a quantitative sequential database D is denoted as $SWU(t)$ and defined as follows:

$$SWU(t) = \sum_{t \subseteq s \wedge s \in D} u(s). \quad (8)$$

For example, in Table 1, $SWU(\langle a \rangle) = u(S_1) + u(S_2) + u(S_3) + u(S_4) + u(S_5) = \$71 + \$69 + \$38 + \$63 + \$52 = \$293$, and $SWU(\langle [a] [c] \rangle) = u(S_1) + u(S_2) + u(S_4) = \$71 + \$69 + \$63 = \$203$.

Theorem 1. (*global downward closure property* [44]) Given a quantitative sequential database D and two sequences t and t' , if $t \subseteq t'$, then:

$$SWU(t') \leq SWU(t). \quad (9)$$

Theorem 2. Given a quantitative sequential database D and a sequence t , it can be obtained that:

$$u(t) \leq SWU(t). \quad (10)$$

The proof for Theorem 1 and Theorem 2 can further referred to in [7, 44]. To improve the performance of utility mining, the USpan algorithm [44] introduces an upper bound based on the remaining utility concept. Details are introduced below.

Definition 22. (*first match*) Given two q -sequences s and s' , if $s \subseteq s'$, then the extension of s in s' is said to be the rest of s' after s , and is denoted as $\langle s' - s \rangle_{rest}$. Given a sequence t and a q -sequence s , if $t \sim s_k \wedge s_k \subseteq s$ ($t \subseteq s$), the rest of t in s is the rest part of s after s_k , which is denoted as $\langle s - t \rangle_{rest}$, where s_k is the first match of t in s .

As an example, consider the sequence $t = \langle [a], [b] \rangle$, q -sequences $s = \langle [a:4], [b:2] \rangle$, and S_3 in Table 1. The remaining part of s in S_3 is $\langle S_3 - s \rangle_{rest} = \langle [(e:1)] \rangle$, and thus it is unique. However, two remaining parts of t exist in S_3 since it has two matches of t in S_3 , and the first one is $\langle [a:4], [b:4] \rangle$. Based on our definition, $\langle S_1 - t \rangle_{rest} = \langle [(f:5)], [(a:1) (b:2) (e:1)] \rangle$.

Definition 23. (*sequence-projected utilization, SPU* [44]) The sequence-projected utilization (*SPU*) of a sequence t in a sequence database D is denoted as $SPU(t)$ and defined as follows:

$$SPU(t) = \sum_{i \in s' \wedge s' \subseteq s \wedge s \in D} (u_{rest}(i, s) + u_p(t, s)), \quad (11)$$

where i is the pivot of t in s , and $u_{rest}(i, s)$ is referred to the remaining utility at q -item i (exclusive) in q -sequence s , such as $u_{rest}(i, s) = \sum_{i' \in s \wedge i < i'} u(i')$. $u_p(t, s)$ is the utility of t in position pivot (p) in s . Note that the pivot is the first place where the q -subsequences match t .

Thus, the *SPU* value of sequence t is the sum of the remaining utilities and utilities of the far left subsequences that match t . However, this is not a true upper bound on utility. When using *SPU* to prune the search space in the LQS-tree with DFS strategy, it may miss some of the real HUSPs. Other reports of this serious problem can be referred to in [36]. Therefore, we propose a real upper-bound on utility for mining HUSPs and the details are given below.

Definition 24. (*sequence extension utility of t in s*) The sequence extension utility (*SEU*) is used to present the maximum utility of the possible extensions that based on the prefix t . Let $SEU(t, s)$ denote the *SEU* of a sequence t in s , and it indicates how much of the sequence's overall utility remains to be extended/concatenated. It is defined as follows:

$$SEU(t, s) = u_{rest}(i, s) + u(t, s), \quad (12)$$

where $u_{rest}(i, s)$ is the remaining utility of i in s , i is the pivot of t in s , w.r.t. is the first occurring position of $s' \sim t$, and $u(t, s)$ is the maximum utility of t in s .

Based on the definition of $\langle s - t \rangle_{rest}$ (cf. Definition 24), it can be seen that $u_{rest}(i, s)$ is equal to $\langle s - t \rangle_{rest}$. For consistency, $\langle s - t \rangle_{rest}$ is used in the following contents.

Definition 25. (*sequence extension utility t in D*) The overall sequence extension utility of a sequence t in a quantitative sequential database D is denoted as $SEU(t)$ and defined as follows:

$$SEU(t) = \sum_{t \subseteq s \wedge s \in D} (u(t, s) + u(\langle s - t \rangle_{rest})), \quad (13)$$

where $u(t, s)$ is the maximum utility value of t in s .

It should be noted that $u(t, s)$ is the maximum utility of t in s (also can be referred to Definition 13) while the $u_p(t, s)$ cannot guarantee the maximum utility of t in s . *SEU* is different from *SPU*. Intuitively, in each q -sequence s that $t \subseteq s$, *SEU* contains less utility values than *SWU*.

Note that $u(\langle s - t \rangle_{rest})$ can be obtained from the constructed utility-array of t in s , which contains the remaining utility of t in each position. For example, in the sequence $t = \langle [a], [b] \rangle$ in Table 1, $SEU(t) = u(t, S_1) + u(\langle S_1 - t \rangle_{rest}) + u(t, S_2) + u(\langle S_2 - t \rangle_{rest}) + u(t, S_3) + u(\langle S_3 - t \rangle_{rest}) = (\$18 + \$15) + (\$23 + \$16) + (\$20 + \$18) = \110 .

Theorem 3. (*local downward closure property*) Given a quantitative sequential database D and two sequences t and t' , if $t \subseteq t'$, it can be obtained that:

$$SEU(t') \leq SEU(t). \quad (14)$$

Proof. Suppose that $s_{q'}$ is a q -sequence that satisfies $u(s_{q'}) = u(t', s)$, where $t' \sim s_{q'} \wedge s_{q'} \subseteq s \wedge s \in D$. The sequence t' can be divided into two parts as the prefix t and the extension e such that $t + e = t'$. Correspondingly, the sequence $s_{q'}$ can also be divided into two parts as the prefix $s_{q'_t}$ matching t and the extension $s_{q'_e}$ matching e such that $s_{q'_t} + s_{q'_e} = s_{q'}$. Then, we have:

$$\begin{aligned} SEU(t', s) &= u(t', s) + u(\langle s - t' \rangle_{rest}) \\ &= u(s_{q'_t}) + u(s_{q'_e}) + u(\langle s - t' \rangle_{rest}) \\ &\leq u(t, s) + u(s_{q'_e}) + u(\langle s - t' \rangle_{rest}) \\ &\leq u(t, s) + u(\langle s - t \rangle_{rest}) \\ &= SEU(t, s). \end{aligned}$$

Thus, $SEU(t', s) \leq SEU(t, s)$. According to $t \subseteq t'$, it is obtained that the set of sequences where $t' \subseteq s$ is a subset of that of $t \subseteq s$. Therefore, $SEU(t') = \sum_{t' \subseteq s \wedge s \in D} \{u(t', s) + u(\langle s - t' \rangle_{rest})\} \leq \sum_{t' \subseteq s \wedge s \in D} \{u(t, s) + u(\langle s - t \rangle_{rest})\} \leq \sum_{t \subseteq s \wedge s \in D} \{u(t, s) + u(\langle s - t \rangle_{rest})\} = SEU(t)$. So far, this theorem holds.

Theorem 4. The SEU value of a sequence t is an upper bound on the utility of this sequence in a quantitative sequential database D . It always has the following relationship:

$$u(t) \leq SEU(t). \quad (15)$$

Proof. We have that $u(t) = \sum_{t \subseteq s \wedge s \in D} \{u(t, s)\} \leq \sum_{t \subseteq s \wedge s \in D} \{u(t, s) + u(\langle s - t \rangle_{rest})\} = SEU(t)$.

Definition 26. (promising HUSP) A sequence t in D is called a *promising* high-utility sequential pattern iff: 1) if the node for t is an *I-Concatenation* node and satisfies $SWU(t) \geq \delta \times u(D)$ or $SEU(t) \geq \delta \times u(D)$; and 2) if the node for t is an *S-Concatenation* node and satisfies $SWU(t) \geq \delta \times u(D)$ or $SEU(t) \geq \delta \times u(D)$; otherwise, this sequence/node is called an invalid or *unpromising* pattern.

Based on Theorem 4, if the upper bound SEU is less than $\delta \times u(D)$, then ProUM can be directly stopped from going deeper and the search procedure can be backtracked. The SEU value of a sequence is an upper bound on the utilities of its extension (the part of its super-sequences). It should be noted that the SEU has the local downward closure property but not the global downward closure property. The reason for this is that some super-sequences of a node/sequence t are not in the subtree of t , and they may be the promising HUSPs even though t 's SEU value is less than $\delta \times u(D)$. SEU can be used to effectively prune the search space in finding HUSPs. In general, SEU is tighter than the previous upper bound SWU . Note that for any non-root node N in the LQS-tree, the SEU can be quickly obtained as an upper bound of all the nodes in the subtree rooted at node N .

Strategy 1. (Pruning of the unpromising one- q -sequences by SWU , called the PUO strategy): Let t be the sequence represented by a node N in the LQS-tree, t' be represented as a child node of N , and δ be the minimum utility threshold. If $SWU(t) \geq \delta \times u(D)$, then ProUM can be stopped from exploring node N . The reason for this is that the sequence t' is always a super-sequence of t . Hence, $u(t') \leq SWU(t') \leq SWU(t) < \delta \times u(D)$. The upper bound SWU has the global downward closure property, and thus any super-sequence t' and its extensions cannot be a desired HUSP. Note that the PUO is a global pruning strategy.

Strategy 2. (Pruning of the unpromising k - q -sequences by SEU , called the PUK strategy): The upper bound SEU of a sequence t can be utilized to prune the unpromising k - q -sequence in its subtree at an early stage when traversing the LQS-tree with the DFS strategy. Thus, if $SEU(t) < \delta \times u(D)$, then the generation of the utility-arrays of its *I-Concatenation* and *S-Concatenation* can be stopped, and traversing all the subtrees from t can be stopped. This is because the utility of t and any of

t 's offspring would not more than $SEU(t)$. Note that the PUK is a local pruning strategy that can be used in the depth pruning in LQS-tree, for example, *I*-extension pruning and *S*-extension pruning.

For example, to avoid constructing the utility-array of the unpromising $(k+1)$ -sequences, the PUK strategy can be applied when scanning the projected sub-databases for k -sequences. This filter operation can reduce both the execution time and memory cost. In summary, the PUO strategy can be used for both width pruning and depth pruning, but it is only used for width pruning in the proposed ProUM algorithm; ProUM utilizes the more powerful PUK strategy for depth pruning. The former also affects the performance of the later. With the PUK strategy, ProUM can easily be stopped from going deeper and the search procedure can be backtracked.

4.4. Proposed ProUM Algorithm

The details of LQS-tree, utility-array, the projection mechanism, and the pruning strategies with SWU and SEU have been introduced as far. To summarize, the pseudocode of main procedure of ProUM is shown in Algorithm 1. A quantitative sequence database D , a profit-table *ptable*, and a minimum utility threshold δ are contained in the input for ProUM; the output includes all the high-utility sequential patterns (HUSPs). Without loss of generality, it is assumed that the proposed ProUM traverses the LQS-tree using the Depth-First-Search (DFS) strategy. By deleting the 1-sequences that $SWU(t) < \delta \times u(D)$ (Line 3), it first scans the original database once to obtain the SWU value of each 1-sequence $t \in D$ (Line 2), and then the revised database D' is obtained. Then, the revised database D' is scanned once to construct the initial utility-arrays for all the sequences in D' (Line 4). After that, ProUM recursively projects a series of sub-utility-arrays based on the prefix sequences (Line 5) by traversing the LQS-tree with DFS strategy.

Algorithm 1 The ProUM algorithm

Input: D ; *ptable*; δ .

Output: *HUSPs*: the complete set of high-utility sequential patterns.

- 1: initialize $D.ua = \emptyset$;
 - 2: scan the original database once to get the SWU value of each 1-sequence;
 - 3: get the revised database D' , by deleting the 1-sequences that $SWU(t) < \delta \times u(D)$ (the PUO strategy);
 - 4: scan the revised database D' once to construct the initial utility-arrays $D.ua$ for all sequences in D' ;
 - 5: **call Project-Search**($\emptyset, I^*, D.ua, \delta$).
 - 6: **return** *HUSPs*
-

The details of the projection and searching procedure are presented in Algorithm 2. When visiting a node/sequence t , ProUM first initializes two sets, $iItem = \emptyset$ and $sItem = \emptyset$ (Line 1). Then, to obtain the promising items for *I-Concatenation* and *S-Concatenation* (Line 2), it scans the projected utility-array $(D.ua)_t$ once. Note that the SEU value of each item is calculated simultaneously during the utility-array scanning (Line

Algorithm 2 The Project-Search procedure

Input: t : a sequence as prefix; $(D.ua)|_t$: the projected utility-array of t ; δ : the minimum utility threshold.

Output: *HUSPs*: the set of high-utility sequential patterns with prefix t .

```
1: initialize  $iItem = \emptyset$  and  $sItem = \emptyset$ ;
2: scan the projected utility-array  $(D.ua)|_t$  once to:
   1) put I-Concatenation items of  $t$  into  $iItem$ ;
   2) put S-Concatenation items of  $t$  into  $sItem$ ;
   3) calculate the SEU values of these items form  $(D.ua)|_t$ ;
3: remove unpromising items  $i_j \in iItem$  that have  $SEU(i_j) < \delta \times u(D)$  (the PUK strategy);
4: remove unpromising items  $i_j \in sItem$  that have  $SEU(i_j) < \delta \times u(D)$  (the PUK strategy);
5: for each item  $i \in iItem$  do
6:    $t' \leftarrow I\text{-Concatenation}(t, i)$ ;
7:   construct the projected utility-array  $(D.ua)|_{t'}$ ;
8:   if  $SEU(t') \geq \delta \times u(D)$  then
9:     if  $u(t') \geq \delta \times u(D)$  then
10:      output  $t'$  into HUSPs;
11:     end if
12:     call Project-Search( $t', (D.ua)|_{t'}, \delta$ ).
13:   end if
14: end for
15: for each item  $i \in sItem$  do
16:    $t' \leftarrow S\text{-Concatenation}(t, i)$ ;
17:   construct the projected utility-array  $(D.ua)|_{t'}$ ;
18:   if  $SEU(t') \geq \delta \times u(D)$  then
19:     if  $u(t') \geq \delta \times u(D)$  then
20:      output  $t'$  into HUSPs;
21:     end if
22:     call Project-Search( $t', (D.ua)|_{t'}, \delta$ ).
23:   end if
24: end for
25: return HUSPs
```

2). After obtaining the updated two sets of $iItem$ and $sItem$, ProUM removes unpromising items that have $SEU(i_j) < \delta \times u(D)$ in $iItem$ and $sItem$, respectively (Lines 3 to 4, the PUK strategy). Next, all these items in $iItem$ and $sItem$ may be used to generate the promising extensions t' as descendant of t . The process of items in $iItem$ is shown in Lines 5 to 14. For a new extension t' whose prefix is t (Line 6), it first constructs the projected sub-utility-array $(D.ua)|_{t'}$ based on the previous utility-array $(D.ua)|_t$ (Line 7). At the same time, the *SEU* value of this *I-Concatenation* can be calculated. Then, ProUM checks this t' whether is able to be the extension as descendant of t . The PUK strategy is used (Line 8, using *SEU* upper bound). If its $SEU < \delta \times u(D)$, then ProUM backtracks to the parent of t ; otherwise, ProUM continues to check the overall utility of this extension t' and outputs t' as a final HUSP if t' satisfies $\delta \times u(D)$ (Lines 9 to 11). Additionally, ProUM calls the **Project-Search** procedure to begin the next projection and search with respect to prefix t' (Line 12). Finally, ProUM recursively explores the other extension nodes in $iItem$ (Line 5) in a similar manner. Similarly, ProUM performs the above procedure to handle each *S-extension* item in $sItem$ (Lines 15 to 24).

Implementation details. During the t -projected utility-array scan with respect to s , to calculate the $SEU(t, s)$ value, an intuitive method is to scan all arrays in the utility-array of s . ProUM has an efficient implementation. The *SEU* value of the sequence is obtained simultaneously when constructing the utility-array of a sequence/transaction. Thus, *SEU* of s is added into the utility-array of s . By also storing this *SEU* value, we can avoid scanning all the elements in the utility-array for calculating the upper bound *SEU* of t in D . Thus, for a given sequence t , its *SEU* value can be quickly obtained since it can be accumulated from the stored *SEU* values by a set of sequences/transactions w.r.t. *sid*.

5. Experiments

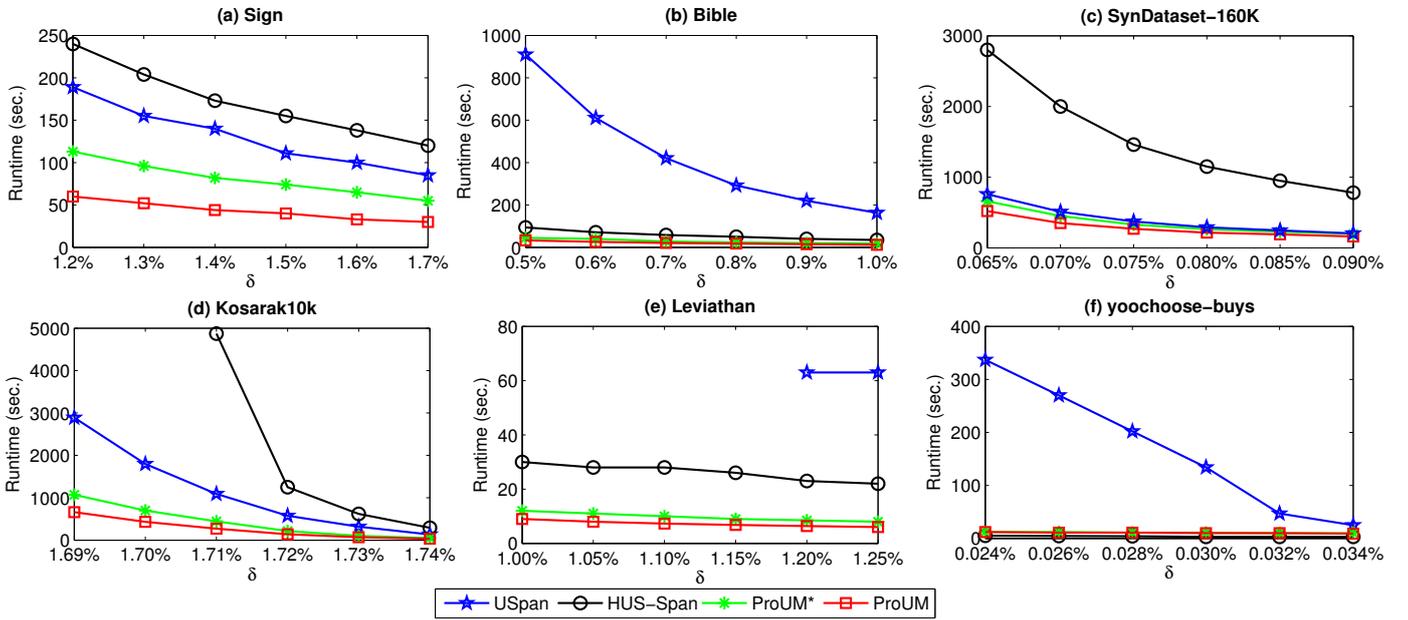
Several experiments were conducted to demonstrate the effectiveness and efficiency of the proposed projection-based utility mining ProUM algorithm.

Evaluation metric. In general, the evaluation metric of comparison for utility-oriented sequential mining algorithms consists of effectiveness analysis with derived patterns, efficiency analysis with execution time and memory consumption, and scalability evaluation. In the following subsections, we use these metrics to evaluate the performance of ProUM.

Compared baselines. For efficiency analysis, USpan [44] (replacing the original *SPU* by *SEU*) and the state-of-the-art HUS-Span [41] algorithm were selected as the baselines. It should be noted that the CRoM that was used in HuspExt [7] is not a true upper bound, and the discovered results by HuspExt are not complete. Therefore, HuspExt is not compared in the following experiments. Two variants of ProUM (respectively denoted as ProUM*, and ProUM) were compared to evaluate the effect of the proposed pruning strategies. ProUM is a hybrid optimization algorithm, as shown in Algorithm 1 and Algorithm 2. In addition, the difference between ProUM and

Table 5: Dataset features

	Dataset	# D	# I	avg(#S)	max(#S)	avg(#Seq)	ave(#Ele)	description
1.	Sign	730	267	52	94	51.99	1.0	language utterance
	Bible	36,369	13,905	21.64	100	17.85	1.0	text
	SynDataset-160k	159,501	7,609	6.19	20	26.64	4.32	synthetic sequences
	Kosarak10k	10,000	10,094	8.14	608	8.14	1.0	web click stream
	Leviathan	5,834	9,025	33.81	100	26.34	1.0	text
	yoochoose-buys	234,300	16,004	1.13	21	2.11	1.97	purchase data
2.	C8S6T4I3D X K (10k)	10,000	7,312	6.22	18	26.99	4.35	synthetic dataset
	C8S6T4I3D X K (80k)	79,718	7,584	6.19	18	26.69	4.32	synthetic dataset
	C8S6T4I3D X K (160k)	159,501	7,609	6.19	20	26.64	4.32	synthetic dataset
	C8S6T4I3D X K (240k)	239,211	7,617	6.19	20	26.66	4.32	synthetic dataset
	C8S6T4I3D X K (320k)	318,889	7,620	6.19	20	26.64	4.32	synthetic dataset
	C8S6T4I3D X K (400k)	398,716	7,621	6.18	20	26.64	4.32	synthetic dataset

Figure 2: Runtime by varying δ

ProUM* is that in Algorithm 2 Lines 3 to 4, ProUM* adopts the PUO strategy to filter the unpromising items in $item$ and $sItem$.

5.1. Data Description and Experimental Configuration

Datasets. For the performance tests, a total of seven datasets were chosen for the different characteristics they displayed. The goal was to show the efficiency of the developed algorithm in a wide range of situations. The chosen datasets and their characteristics are displayed in Table 5. Note that #|D| is the number of sequences, #|I| is the number of different symbols/items in the dataset, #S is the length of a sequence s , #Seq is the number of elements per sequence, and #Ele is the average number of items per element/itemset. SynDataset-160K is a synthetic sequential dataset generated by IBM Quest Dataset Generator [1]. The original yoochoose-buys⁵ dataset contains the quantity and unit

profit of each object/item while other datasets⁶ do not contain the quantity and unit profit. Therefore, we adopted a simulation method that is widely used in previous studies [23, 30, 37] to generate the quantitative and profit information for each object/item in datasets except for yoochoose-buys.

Experimental configuration. All the compared algorithms in the experiments were implemented in Java language. Note that the original USpan algorithm with SPU upper bound may cause incomplete mining results. Thus, the USpan code used here is a revised and optimized version. Furthermore, SPU was replaced with SEU in USpan so that it could discover the complete HUSPs. All the experiments were performed on a personal ThinkPad T470p computer with an Intel(R) Core(TM) i7-7700HQ CPU @ 2.80 GHz 2.81 GHz processor, 32 GB of RAM, and with 64-bit Microsoft Windows 10 operating system.

⁵<https://recsys.acm.org/recsys15/challenge/>

⁶<http://www.philippe-fournier-viger.com/spmf/index.php>

Table 6: Number of patterns (candidates and final results) under various δ values

		# of patterns under various δ values					
		δ_1	δ_2	δ_3	δ_4	δ_5	δ_6
(a) Sign	#P1	6,598,217	5,265,824	4,250,363	3,490,867	2,886,277	2,418,799
	#P2	6,598,217	5,265,825	4,250,363	3,490,869	2,886,278	2,418,799
	#P3	24,586,054	19,345,657	15,424,452	12,536,127	10,269,881	8,534,821
	#P4	6,598,215	5,265,822	4,250,359	3,490,865	2,886,274	2,418,798
	#HUSPs	78,336	56,395	41,151	30,440	22,702	17,274
(b) Bible	#P1	92,563	59,041	40,766	29,183	21,227	16,488
	#P2	95,012	60,588	41,786	29,940	21,746	16,887
	#P3	262,465	163,564	109,449	76,979	56,443	42,205
	#P4	100,706	64,091	43,983	31,630	23,486	17,831
	#HUSPs	2,760	1,714	1,124	764	553	411
(c) SynDataset-160K	#P1	8,751,355	5,357,157	3,313,183	2,122,646	1,394,602	948,156
	#P2	8,752,654	5,357,847	3,313,681	2,123,007	1,394,913	948,397
	#P3	23,917,337	13,187,663	8,101,205	5,412,965	3,735,188	2,645,849
	#P4	8,753,634	5,358,865	3,314,487	2,123,613	1,395,359	948,520
	#HUSPs	58,710	17,903	4,794	1,344	394	172
(d) Kosarak10k	#P1	124,833,772	82,478,688	51,535,423	24,542,377	12,100,103	5,524,463
	#P2	-	-	51,535,979	24,542,940	12,100,669	5,525,033
	#P3	478,850,673	321,429,317	204,661,580	100,469,791	42,141,209	20,427,300
	#P4	124,833,676	82,478,593	51,535,330	24,542,295	12,100,024	5,524,390
	#HUSPs	23	22	22	22	22	21
(e) Leviathan	#P1	-	-	-	-	46,193	41,610
	#P2	76,549	68,058	60,557	54,084	48,162	43,361
	#P3	205,381	181,067	159,208	140,946	125,177	111,605
	#P4	82,625	73,315	65,076	58,140	52,181	47,031
	#HUSPs	1,802	1,520	1,322	1,152	996	869
(f) yoochoose-buys	#P1	325,473	304,534	278,131	250,138	201,440	158,219
	#P2	312,780	304,737	278,318	250,307	201,620	158,409
	#P3	314,346	305,226	280,022	251,896	203,867	159,858
	#P4	313,724	304,642	279,481	251,434	203,422	159,463
	#HUSPs	317,682	296,890	273,926	238,273	191,203	146,761

5.2. Efficiency Analytics

As previously mentioned, a good high-utility sequence mining method should be efficient and able to scale well to handle long sequence data. Thus, the running time of the compared methods were compared under different parameter settings. We increased the minimum utility threshold from δ_1 to δ_6 on each dataset while keeping the tested data size fixed. To obtain accurate experimental results under each setting, each compared approach was ran three times, and the average running times are plotted in Figure 2. As shown, the runtime of USpan exceeded 10,000 seconds in Leviathan when the minimum utility threshold was lower than 1.20%, and thus USpan only has two points in Figure 2(e).

As shown in each sub-figure of Figure 2, ProUM is intuitively the most efficient except in yoochoose-buys. Both the proposed algorithm with or without using the PUK strategy (ProUM and ProUM*) to prune the unpromising candidates before constructing the utility-arrays consistently outperformed the state-of-the-art HUS-Span approach, even by up to 3 orders of magnitude. For the Kosarak10k data in Figure 2(d), the performance of ProUM decreased when δ increased, but it decreased slowly afterwards when $\delta = 1.72\%$. USpan always required a longer execution time than ProUM and ProUM*, from 2,890 seconds to 130 seconds. In particular, HUS-Span had the longest execution time in this dataset, and it consumed 5,000 seconds when δ was smaller than 1.71%. In general, ProUM outperformed ProUM* in all the test datasets under different parameter settings. For example, in Figure 2(a), the difference of the runtime between ProUM* and ProUM can be observed.

When $\delta = 1.2\%$ on the Sign dataset, the runtime of ProUM* closed to 115 seconds while the runtime of ProUM was approximately 60 seconds. These observations can also be intuitively seen on other datasets, such as Figure 2(c), Figure 2(d), and Figure 2(e). These observations indicate that the *local downward closure* property of the SEU upper bound plays an active role in pruning the search space of the projection-based ProUM algorithm.

In addition, it is interesting to observe that USpan sometimes ran even faster than HUS-Span. In many cases, however, it is not clear whether the recently proposed HUS-Span was faster than the USpan method that was optimized for this experiment. For example, when the experiment was conducted on Sign, SynDataset-160K, and Kosarak10k, it seems that HUS-Span had a longer running time than USpan. For the SynDataset-160K shown in Figure 2(c), HUS-Span was the most time consuming among the four algorithms. It required 2,824 seconds when δ was set to 0.065%, which was quite a bit longer than the others. The performance of USpan, ProUM*, and ProUM declined before $\delta = 0.075\%$, and it nearly remained stable afterwards. In other datasets, such as Bible, Leviathan, and yoochoose-buys, USpan performed worse than HUS-Span as well as the two variants of the proposed ProUM model. A possible reason for this is that HUS-Span needs additional time to scan the projected sub-databases for calculating the utility information from the built utility-chains. In addition, in some datasets, the upper bound PEU sometimes had a similar effect to that of the proposed SEU upper bound.

The projection mechanism of utility-array makes a contri-

bution to the improvement, which can be observed in SynDataset-160K and Kosarak10k. This is because the small size of the utility-array creates a favorable *SEU* value that enhances the computation of the later processes. Specifically, based on an observation of the runtime between ProUM and ProUM*, the PUK strategy plays an active role in filtering the unpromising patterns before constructing the set of utility-arrays. In summary, the enhanced ProUM algorithm that utilizes powerful pruning strategies always had the best performance compared to the baseline ProUM* as well as USpan and the state-of-the-art HUS-Span algorithm. The designed ProUM algorithm is acceptable and efficient in discovering high-utility sequential patterns on different types of datasets.

Summary of efficiency study. The above-stated results demonstrate the efficiency of ProUM. Under different parameter settings (when δ is large), ProUM always required less time than the existing HUSPM algorithms. In addition, the divide-and-conquer strategy was applied to project the utility-arrays during the recursive mining processes. All the experimental results demonstrate the suitability of the proposed ProUM models for dealing with both real or synthetic datasets.

5.3. Candidate Analysis

The generated patterns of the four compared algorithms were investigated to evaluate the effect of pruning strategies by conducting under the same parameter settings, as shown in Figure 2. The results of the different kinds of generated patterns, both generated candidates and final HUSPs, are plotted in Table 6. Note that the #HUSPs is the number of final discovered HUSPs, and #P1, #P2, #P3, and #P4 are the numbers of the candidates generated by USpan, HUS-Span, ProUM*, and ProUM, respectively.

As shown in Table 6, it can be clearly observed that, on all the tested datasets, the number of HUSPs was always quite a bit less than that of the candidate patterns (e.g., #P1, #P2, #P3, and #P4) under various minimum utility thresholds. For example, in Kosarak10k, the discovered HUSPs were changed from 23 to 21 while the related candidates were increased from 8,753,634 up to 23,917,337. These results reflect the fact that there are a huge number of candidate patterns that are generated in a HUSPM algorithm but very few of them are the final interesting desired patterns. As mentioned previously, there are several challenges in utility mining when dealing with sequence data. How to effectively prune the search space in HUSPM is more difficult due to the absence of the downward closure property in the sequence utility.

Intuitively, on all tested datasets with different parameter settings, #P1 and #P4 was nearly equal to #P2 while #P2 had the most number of candidates among all the compared candidate patterns. This is because the previously mentioned upper bound error in the USpan algorithm was replaced by the proposed *SEU* upper bound in our experiments. Thus, both USpan and ProUM used the same upper bounds, *SWU* and *SEU*, to prune the search space. This results in nearly the same results of the number of the candidate patterns. It is interesting to observe that #P4 was nearly equal to #P2, which indicates that the *SEU* used in ProUM had a similar powerful pruning ability

to *PEU* that was used in HUS-Span for the addressed HUSPM problem.

Specifically, the difference between #P3 and #P4 proves the effectiveness of the PUK strategy for ProUM. That is, the proposed *SEU* upper bound has a better ability to prune the search space than the loose *SWU*. Although both *SWU* and *SEU* affect the candidate patterns for mining HUSPs, in general, the numbers of #P4 were always quite smaller than those of #P3. For example, as shown in Bible, #P4 changed from 262,465 to 42,205 while #P3 had its number decreased from 100,706 to 17,831 when varying δ from 0.5% to 1.0%. Therefore, ProUM* adopts the PUO strategy to filter the unpromising items in *item*, and *sItem* is not more powerful than ProUM, which utilizes PUK strategy (w.r.t. the *SEU* upper bound) to remove these unpromising items.

Discussion. These results of the patterns indicate that the proposed upper bound *SEU* is more suitable than *SWU* to prune the subtrees of LQS-tree for mining HUSPs. The results demonstrate the positive effect of pruning strategies in ProUM for discovering utility-driven sequential patterns.

5.4. Memory Evaluation

In this subsection, the mining efficiency is evaluated in terms of memory consumption. All parameters are set to the default values shown in Figure 2 unless otherwise stated. Figures 3(a) to (f) respectively show plots of the results of the peak memory usage of all the compared algorithms. Note that Java API was used to calculate the peak memory consumption of each compared algorithm during the whole mining process.

As shown, the projection utility-array-based models, both ProUM* and ProUM, performed significantly better than the baselines. Although the HUS-Span algorithm also utilizes the projection technique, it needed to project the sub-databases before the construction of utility-chains. This consumes more execution time and memory cost than ProUM, which only projects and scans the sub-utility-arrays. For example, as shown in Figures 3 (a) and (f), the peak memory consumption for ProUM was significantly less than that of HUS-Span because ProUM consumes the reasonable memory to store the compact utility-arrays and generate promising patterns. In addition, the improved variant ProUM consumes less memory than the baseline ProUM* that adopts the PUO strategy to remove the unpromising items.

Figure 3 shows the effects of the parameter - minimum utility threshold δ on the memory performance of ProUM. As shown on all datasets, the memory usage of ProUM* and ProUM did not change much when δ increased while the memory usage of USpan and HUS-Span may change more substantially in most cases. For example, when $\delta = 1.20\%$, ProUM consumed only 1,000 MB on Leviathan while HUS-Span consumed more than 2,000 MB. The performance gap is more obvious on yoochoose-buys, mainly because USpan and HUS-Span were ineffective on Leviathan and yoochoose-buys. The decrease with δ was quite rapid because a small δ makes all the compared HUSPM algorithms execute searching in LQS-tree more times, which makes it harder to return the mined results, especially for the existing USpan and HUS-Span algorithms. In addition, a very

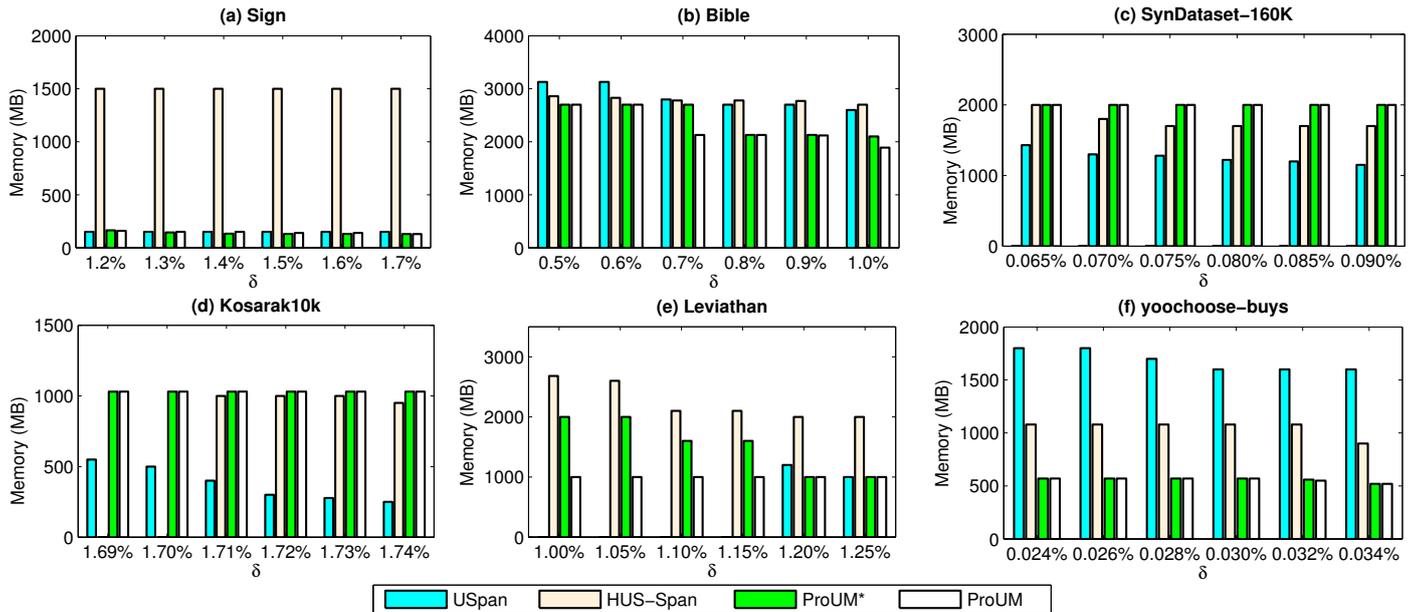


Figure 3: Memory usage by varying δ

large δ brings no extra benefit to the discovered results of HUSPM. Hence, in practice, δ does not need to be too large.

Summary. The proposed ProUM model with several pruning strategies consumed less memory than HUS-Span for all the parameter settings and also less than that of the optimized USpan in most cases. USpan had the least memory consumption in Figures 3(c) and 3(d). Nonetheless, the best performing USpan on memory consumption for these cases had much worse execution times in Figure 2(d). As previously mentioned, one of the advantages of ProUM is that it is able to filter a large amount of unpromising patterns at an early stage by building the projected utility-arrays.

5.5. Scalability Test

Scalability is important mainly because many real-world data is massive, especially large-scale sequence data. Therefore, scalability is an important acceptance criteria for a designed data mining model. In this subsection, ProUM is analytically compared with the existing methods on a large-scale dataset, and the experimental results are shown in Figures 4 (a) to (c), respectively.

Figure 4 shows the results of scalability performance in the synthetic dataset C8S6T4I3D|X|K in terms of different data sizes from 10K to 400K sequences. The running time of each algorithm increased linearly as the number of sequences grew. ProUM showed superior scalability with respect to the dataset size among the compared methods. For the test dataset that contained element-based sequences, USpan performed better than the HUS-Span algorithm that utilizes the *PEU* upper bound. For example, the running time of HUS-Span exceeded 2,000 seconds when dealing with 400K sequences. According to the memory usage, ProUM required more memory than other baselines, as shown in the center sub-figure. In addition, the amount

of memory used was bounded in ProUM through the pruning strategies that were employed. The candidate patterns still show that *#P3* was similar to *#P1* and *#P2* (Figure 4(c)). Note that here *#P3* is number of generated candidates in ProUM.

Discussion. In summary, the scalability results confirm the intuition that the projected ProUM method using utility-array with a series of indexing positions is scalable for large-scale datasets, and it is superior to the existing algorithms.

6. Conclusions

In general, utility-based sequence analytics are more useful than other support-based data mining techniques. However, utility mining on sequence data can easily suffer from several problems, not just with critical combinational explosion but also from computational complexity caused by sequencing between itemsets/elements. How to improve the mining efficiency of utility mining on sequence data is still an open problem. To this end, we developed a projection-based utility mining algorithm, called ProUM, for the fast mining of high-utility sequential patterns. A new data structure, called utility-array, was also proposed, which can be directly used to calculate the utility and remaining utility of a sequence without scanning the database. Based on the projection mechanism in applying in utility-array, the presented solutions, including two utility bounds, the corresponding pruning strategies, and the ProUM algorithm, are proposed. ProUM was compared with USpan and HUS-Span, which are the state-of-the-art algorithms for mining HUSPs in sequence data. Extensive experimental results on both synthetic and real-life datasets demonstrated that ProUM had a better efficiency compared to the state-of-the-art baselines.

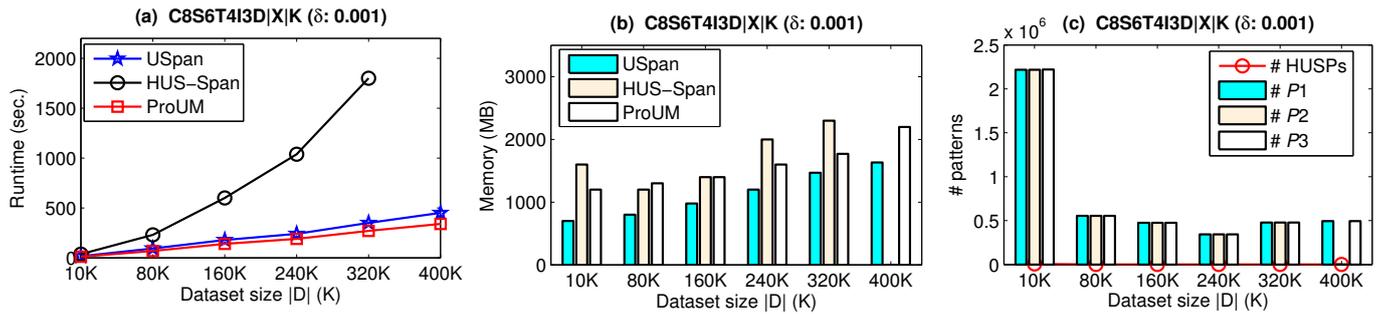


Figure 4: Scalability test.

Acknowledgment

This work was partially supported by the Shenzhen Technical Project under project No. KQJSCX 20170726103424709 and No. JCYJ 20170307151733005. Specifically, Wensheng Gan was supported by the CSC (China Scholarship Council) Program during the study at University of Illinois at Chicago, IL, USA.

References

References

- [1] Rakesh Agrawal and Ramakrishnan Srikant. Quest synthetic data generator. <http://www.Almaden.ibm.com/cs/quest/syndata.html>, 1994.
- [2] Rakesh Agrawal and Ramakrishnan Srikant. Mining sequential patterns. In *The International Conference on Data Engineering*, pages 3–14. IEEE, 1995.
- [3] Rakesh Agrawal, Ramakrishnan Srikant, et al. Fast algorithms for mining association rules. In *Proceedings of the 20th International Conference on Very Large Data Bases*, volume 1215, pages 487–499, 1994.
- [4] Chowdhury-Farhan Ahmed, Syed-Khairuzzaman Tanbeer, and Byeong-Soo Jeong. Mining high utility web access sequences in dynamic web log data. In *11th ACIS International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing*, pages 76–81. IEEE, 2010.
- [5] Chowdhury-Farhan Ahmed, Syed-Khairuzzaman Tanbeer, and Byeong-Soo Jeong. A novel approach for mining high-utility sequential patterns in sequence databases. *ETRI Journal*, 32(5):676–686, 2010.
- [6] Chowdhury-Farhan Ahmed, Syed-Khairuzzaman Tanbeer, Byeong-Soo Jeong, and Young-Koo Lee. Efficient tree structures for high utility pattern mining in incremental databases. *IEEE Transactions on Knowledge and Data Engineering*, 21(12):1708–1721, 2009.
- [7] Ozgur Kirmemis Alkan and Pinar Karagoz. CROm and HuspExt: Improving efficiency of high utility sequential pattern extraction. *IEEE Transactions on Knowledge and Data Engineering*, 27(10):2645–2657, 2015.
- [8] Jay Ayres, Jason Flannick, Johannes Gehrke, and Tomi Yiu. Sequential pattern mining using a bitmap representation. In *Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 429–435. ACM, 2002.
- [9] Ming-Syan Chen, Jiawei Han, and Philip S. Yu. Data mining: an overview from a database perspective. *IEEE Transactions on Knowledge and Data Engineering*, 8(6):866–883, 1996.
- [10] Duy-Tai Dinh, Bac Le, Philippe Fournier-Viger, and Van-Nam Huynh. An efficient algorithm for mining periodic high-utility sequential patterns. *Applied Intelligence*, 48(12):4694–4714, 2018.
- [11] Philippe Fournier-Viger, Jerry Chun-Wei Lin, Rage-Uday Kiran, and Yun-Sing Koh. A survey of sequential pattern mining. *Data Science and Pattern Recognition*, 1(1):54–77, 2017.
- [12] Philippe Fournier-Viger, Cheng-Wei Wu, Souleymane Zida, and Vincent S Tseng. FHM: Faster high-utility itemset mining using estimated utility co-occurrence pruning. In *International Symposium on Methodologies for Intelligent Systems*, pages 83–92. Springer, 2014.

- [13] Wensheng Gan, Jerry Chun-Wei Lin, Han-Chieh Chao, Tzung-Pei Hong, and Philip S Yu. CoUPM: Correlated utility-based pattern mining. In *Proceedings of the IEEE International Conference on Big Data*, pages 2607–2616. IEEE, 2018.
- [14] Wensheng Gan, Jerry Chun-Wei Lin, Han-Chieh Chao, Shyue-Liang Wang, and Philip S Yu. Privacy preserving utility mining: a survey. In *Proceedings of the IEEE International Conference on Big Data*, pages 2617–2626. IEEE, 2018.
- [15] Wensheng Gan, Jerry Chun-Wei Lin, Han-Chieh Chao, and Justin Zhan. Data mining in distributed environment: a survey. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 7(6):e1216, 2017.
- [16] Wensheng Gan, Jerry Chun-Wei Lin, Philippe Fournier-Viger, and Han-Chieh Chao. Exploiting high utility occupancy patterns. In *Asia-Pacific Web and Web-Age Information Management Joint Conference on Web and Big Data*, pages 239–247. Springer, 2017.
- [17] Wensheng Gan, Jerry Chun-Wei Lin, Philippe Fournier-Viger, Han-Chieh Chao, and Hamido Fujita. Extracting non-redundant correlated purchase behaviors by utility measure. *Knowledge-Based Systems*, 143:30–41, 2018.
- [18] Wensheng Gan, Jerry Chun-Wei Lin, Philippe Fournier-Viger, Han-Chieh Chao, Vincent S Tseng, and Philip S Yu. A survey of utility-oriented pattern mining. *arXiv preprint arXiv:1805.10511*, 2018.
- [19] Wensheng Gan, Jerry Chun-Wei Lin, Philippe Fournier-Viger, Han-Chieh Chao, and Philip S Yu. A survey of parallel sequential pattern mining. *ACM Transactions on Knowledge Discovery from Data*, 13(3):25, 2019.
- [20] Jiawei Han, Jian Pei, Yiwen Yin, and Runying Mao. Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *Data Mining and Knowledge Discovery*, 8(1):53–87, 2004.
- [21] Guo-Cheng Lan, Tzung-Pei Hong, Vincent S Tseng, and Shyue-Liang Wang. Applying the maximum utility measure in high utility sequential pattern mining. *Expert Systems with Applications*, 41(11):5071–5081, 2014.
- [22] Jerry Chun-Wei Lin, Philippe Fournier-Viger, and Wensheng Gan. FHN: An efficient algorithm for mining high-utility itemsets with negative unit profits. *Knowledge-Based Systems*, 111:283–298, 2016.
- [23] Jerry Chun-Wei Lin, Wensheng Gan, Philippe Fournier-Viger, Tzung-Pei Hong, and Han-Chieh Chao. FDHUP: Fast algorithm for mining discriminative high utility patterns. *Knowledge and Information Systems*, 51(3):873–909, 2017.
- [24] Jerry Chun-Wei Lin, Wensheng Gan, Philippe Fournier-Viger, Tzung-Pei Hong, and Vincent S Tseng. Efficient algorithms for mining high-utility itemsets in uncertain databases. *Knowledge-Based Systems*, 96:171–187, 2016.
- [25] Jerry Chun-Wei Lin, Wensheng Gan, and Tzung-Pei Hong. A fast updated algorithm to maintain the discovered high-utility itemsets for transaction modification. *Advanced Engineering Informatics*, 29(3):562–574, 2015.
- [26] Jerry Chun-Wei Lin, Wensheng Gan, and Tzung-Pei Hong. A fast maintenance algorithm of the discovered high-utility itemsets with transaction deletion. *Intelligent Data Analysis*, 20(4):891–913, 2016.
- [27] Jerry Chun-Wei Lin, Lu Yang, Philippe Fournier-Viger, Tzung-Pei Hong, and Miroslav Voznak. A binary PSO approach to mine high-utility itemsets. *Soft Computing*, 21(17):5103–5121, 2017.
- [28] Jerry Chun-Wei Lin, Jiexiong Zhang, and Philippe Fournier-Viger. High-utility sequential pattern mining with multiple minimum utility thresholds. In *Asia-Pacific Web and Web-Age Information Management Joint*

- Conference on Web and Big Data*, pages 215–229. Springer, 2017.
- [29] Ying-Chun Lin, Cheng-Wei Wu, and Vincent S Tseng. Mining high utility itemsets in big data. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 649–661, 2015.
- [30] Mengchi Liu and Junfeng Qu. Mining high utility itemsets without candidate generation. In *Proceedings of the 21st ACM International Conference on Information and Knowledge Management*, pages 55–64. ACM, 2012.
- [31] Ying Liu, Wei-Keng Liao, and Alok Choudhary. A two-phase algorithm for fast discovery of high utility itemsets. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 689–695. Springer, 2005.
- [32] Alfred Marshall. From principles of economics. In *Readings in the Economics of the Division of Labor: the Classical Tradition*, pages 195–215. World Scientific, 2005.
- [33] Jian Pei, Jiawei Han, Behzad Mortazavi-Asl, Helen Pinto, Qiming Chen, Umeshwar Dayal, and Mei Chun Hsu. PrefixSpan: Mining sequential patterns efficiently by prefix-projected pattern growth. In *The International Conference on Data Engineering*, pages 215–224. IEEE, 2001.
- [34] Bai-En Shie, Hui-Fang Hsiao, Vincent S Tseng, and Philip S Yu. Mining high utility mobile sequential patterns in mobile commerce environments. In *Proceedings of International Conference on Database Systems for Advanced Applications*, pages 224–238. Springer, 2011.
- [35] Ramakrishnan Srikant and Rakesh Agrawal. Mining sequential patterns: generalizations and performance improvements. In *Proceedings of International Conference on Extending Database Technology*, pages 1–17. Springer, 1996.
- [36] Tin Truong-Chi and Philippe Fournier-Viger. A survey of high utility sequential pattern mining. In *High-Utility Pattern Mining*, pages 97–129. Springer, 2019.
- [37] Vincent S Tseng, Bai-En Shie, Cheng-Wei Wu, and Philip S Yu. Efficient algorithms for mining high utility itemsets from transactional databases. *IEEE Transactions on Knowledge and Data Engineering*, 25(8):1772–1786, 2013.
- [38] Vincent S Tseng, Cheng-Wei Wu, Philippe Fournier-Viger, and Philip S Yu. Efficient algorithms for mining top- k high utility itemsets. *IEEE Transactions on Knowledge and Data Engineering*, 28(1):54–67, 2016.
- [39] Vincent S Tseng, Cheng-Wei Wu, Bai-En Shie, and Philip S Yu. UP-Growth: an efficient algorithm for high utility itemset mining. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 253–262. ACM, 2010.
- [40] Jun-Zhe Wang and Jiun-Long Huang. On incremental high utility sequential pattern mining. *ACM Transactions on Intelligent Systems and Technology*, 9(5):55, 2018.
- [41] Jun-Zhe Wang, Jiun-Long Huang, and Yi Cheng Chen. On efficiently mining high utility sequential patterns. *Knowledge and Information Systems*, 49(2):597–627, 2016.
- [42] Cheng-Wei Wu, Yu-Feng Lin, Philip S Yu, and Vincent S Tseng. Mining high utility episodes in complex event sequences. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 536–544. ACM, 2013.
- [43] Jimmy Ming-Tai Wu, Justin Zhan, and Jerry Chun-Wei Lin. An ACO-based approach to mine high-utility itemsets. *Knowledge-Based Systems*, 116:102–113, 2017.
- [44] Junfu Yin, Zhigang Zheng, and Longbing Cao. USpan: an efficient algorithm for mining high utility sequential patterns. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 660–668. ACM, 2012.
- [45] Junfu Yin, Zhigang Zheng, Longbing Cao, Yin Song, and Wei Wei. Efficiently mining top- k high utility sequential patterns. In *Proceedings of the IEEE 13th International Conference on Data Mining*, pages 1259–1264. IEEE, 2013.
- [46] Mohammed J Zaki. SPADE: an efficient algorithm for mining frequent sequences. *Machine Learning*, 42(1-2):31–60, 2001.
- [47] Souleymane Zida, Philippe Fournier-Viger, Jerry Chun-Wei Lin, Cheng-Wei Wu, and Vincent S Tseng. EFIM: a fast and memory efficient algorithm for high-utility itemset mining. *Knowledge and Information Systems*, 51(2):595–625, 2017.