

Mining Time-Interval Sequential Patterns Using Clustering Analysis

Hao-En Chueh¹, Nancy P. Lin²

¹Department of Information Management, Yuanpei University

²Department of Computer Sciences and Information Engineering, Tamkang University

¹hechueh@mail.ypu.edu.tw, ²nancylin@mail.tku.edu.tw

Abstract—A time-interval sequential pattern is a sequential pattern with the information about the time intervals between itemsets. Most algorithms of mining time-interval sequential patterns find the time intervals between itemsets by predefining some non-overlap time partitions, however, a predefined set of non-overlap time partitions cannot be suitable for every pair of successive itemsets. Therefore, in this paper, we present a new algorithm to mine time-interval sequential patterns without defining any time partitions in advance. The algorithm first adopts the clustering analysis to automatically generate the suitable time partitions for frequent occurring pairs of successive itemsets, and then uses these generated time partitions to extend typical algorithms to discover the time-interval sequential patterns. Our result of experiment verifies that this algorithm outperforms than the algorithms which use some predefined and non-overlap time partitions.

Keywords: Data Mining, Sequential Pattern, Time Interval, Clustering Analysis.

1. Introduction

Data mining is the procedure of discovering hidden, useful, previously unknown information from databases. The existing techniques of data mining include association analysis, classification, clustering, etc, and one of the important techniques is mining sequential patterns. Mining sequential patterns first introduced by Agrawal and Srikant (1995) is the task of finding frequently occurring patterns related to time or other sequences from a sequence database [1]. An example of a sequential pattern is “A customer who bought a digital camera will return to buy an extra memory card”. Mining sequential patterns is widely used in the field of retail business to assist in making various marketing decisions since many transaction records are stored as sequence data [3,5,7].

Many algorithms of mining sequential patterns

have been proposed [1,4,6,9], and most of these algorithms only focus on finding the order of the itemsets, but ignore the time intervals between successive itemsets. An example of a sequential pattern with time intervals between successive itemsets is “A customer who bought a digital camera will buy an extra memory card within one month”. In business field, actually, a sequential pattern which includes the time intervals between successive itemsets is more valuable than a traditional sequential pattern without any time information, because the time intervals between successive itemsets can offer useful information for businesses to sell the appropriate products to their customers at the right time. To this end, some researches start to propose algorithms to discover the sequential patterns with time intervals between itemsets, called time-interval sequential pattern [2].

To find the time intervals between every pair of successive itemsets, most proposed algorithms of mining time-interval sequential patterns usually predefine some non-overlap time partitions, and assume that the time interval between each pair of successive itemsets can match one time partition of this predefined set, but a predefined set of non-overlap time partitions, in fact, cannot be suitable for every pair of successive itemsets. Generating the suitable time partitions for every pair of successive itemsets directly from the real data is more reasonable. Accordingly, in this paper, we present a new algorithm to mine time-interval sequential patterns without predefining any time partitions. This algorithm tries to use the clustering analysis first to automatically generate the suitable time partitions between frequent occurring pairs of successive itemsets, and then uses these generated time partitions to extend typical algorithms to discover the time-interval sequential patterns.

The rest of this paper is organized as follows. Some researches related to time-interval sequential patterns are reviewed in section 2. The proposed time-interval sequential patterns mining algorithm is presented in section 3. A simple experiment is displayed in section 4. The conclusions are given

in section 5.

2. Time-Interval Sequential Patterns

The problem of sequential patterns mining which can be described as the task of discovering frequently occurring ordered patterns from a given sequence database was first introduced by Agrawal and Srikant in the mid 1990s [1].

A sequence is an ordered list of itemsets. Let $I = \{i_1, i_2, \dots, i_m\}$ be a set of items, $S = \langle s_1, s_2, \dots, s_k \rangle$ is a sequence, where $s_i \subseteq I$ is called an itemset. Length of a sequence means the number of itemsets in the sequence, and a sequence contains k itemsets is called a k -sequence. The support of a sequence S is denoted by $\text{supp}(S)$ and means the percentage of total number of records containing sequence S in the sequence database. If $\text{supp}(S)$ is greater than or equal to a predefined threshold, called minimal support, then sequence S is considered as a frequent sequence and called a sequential pattern.

Many algorithms have been proposed to mine sequential pattern [1,4,6,9], and most algorithms only focus on finding the frequently occurring order of the itemsets, but ignore the time intervals between itemsets. The time intervals between itemsets, in fact, can offer useful information for businesses to sell the appropriate products to their customers at the right time. Due to the value of the time intervals between itemsets, some researchers start to propose algorithms to mine various sequential patterns with time information between itemsets recently [2,8,10,11].

Srikant et al. [10] utilize three restrictions, the maximum interval (*max-interval*), the minimum interval (*min-interval*), and the time window size (*window-size*) to find sequential patterns related to time intervals, and the discovered pattern is like $((A, B), (C, D))$, where (A, B) and (C, D) are called subsequences of $((A, B), (C, D))$. The *max-interval* and the *min-interval* are respectively used to indicate the maximal and minimal interval within subsequence. The *window-size* is used to indicate the interval among subsequences. For example, let the *max-interval* is set to 10 hours, the *min-interval* is set to 3 hours, and *window-size* is set to 24 hours. That is, the time interval between A and B lies in $[2, 10]$, the time interval between C and D also lies in $[2, 10]$, and the time interval between (A, B) and (C, D) lies in $[0, 24]$.

Mannila et al. [8] use a window width (*win*) to

find frequent episodes in sequences of events, and the discovered episode is like (A, B, C) . Assume that the *win* is set to 3 days, then the episode (A, B, C) means that, in 3 days, A occurs first, B follows, and C happens finally.

Wu et al. [11] utilize a window (d) as well to find the sequential pattern likes (A, B, C) , such that, in a sequential pattern, the interval between adjacent events is within the window d . Assume d is set to 5 hours, then the discovered pattern (A, B, C) means that A occurs first, B follows, and C happens finally and the interval between A and B , and between B and C is within 5 hours.

Chen et al. [2] use a predefined set of non-overlap time partitions to discover potential time-interval sequential patterns, and the discovered pattern is like (A, I_0, B, I_2, C) , where I_0, I_2 belong to the non-overlap set of time partitions. Assume that, in the set of time partitions, I_0 denotes the time interval t satisfying $0 \leq t \leq 1$ day; I_2 denotes the time interval t satisfying $3 < t \leq 7$ days, and then the pattern (A, I_0, B, I_2, C) means that A, B and C happen in this order, and the interval between A and B is within 1 day, and the interval between B and C lies between 3 days and 7 days.

Although these preceding discussion researches can discover the sequential patterns with the time information between itemsets by using a or some predefined and fixed time partitions, the patterns whose intervals between itemsets lie outside these predefined time ranges cannot be found yet.

Therefore, this work presents a new algorithm to mine time-interval sequential patterns without defining any time partitions in advance, and the main concept of this algorithm is to generate the suitable time partitions directly from the real data. This algorithm first uses the clustering analysis to automatically generate the suitable time partitions for frequent occurring pairs of successive itemsets, and then adopts these generated time partitions to extend typical algorithms to mine the sequential patterns with time intervals between every pair of successive itemsets.

Details of the proposed time-interval sequential mining patterns algorithm are presented in the following section.

3. The Proposed Algorithm

The proposed time-interval sequential patterns mining algorithm is presented in this section. This

algorithm uses clustering analysis to automatically generate the suitable time partitions directly from the real data for the frequent occurring pairs of successive itemsets, and then uses these generated time partitions to extend the typical algorithms to discover the sequential patterns with time intervals between every pair of successive itemsets. Steps of the algorithm are described as follows.

Notation:

$I = \{ i_1, i_2, \dots, i_m \}$: The set of items.

$S_i = \langle s_1, s_2, \dots, s_n \rangle$: A sequence, where each

$$s_k \subseteq I.$$

$D = \{ S_1, S_2, \dots, S_k \}$: The set of sequences.

$\text{supp}(S_i)$: The support of S_i .

min-supp : The minimal support threshold.

CS_k : The set of candidate k-sequences.

FS_k : The set of frequent k-sequences.

$CTIS_k$: The set of candidate time-interval k-sequences.

$FTIS_k$: The set of frequent time-interval k-sequences.

Algorithm:

Step 1: Producing FS_1 . Each itemset can be regarded as a candidate 1-sequence. A candidate 1-sequence whose support is greater than or equal to min-supp is a frequent 1-sequence, and the set of all frequent 1-sequences is FS_1 .

Step 2: Producing CS_2 . For any two frequent 1-sequences s_1 and s_2 , where $s_1, s_2 \in FS_1$ and $s_1 \neq s_2$, then we can generate 2 candidate 2-sequences $\langle s_1, s_2 \rangle$ and $\langle s_2, s_1 \rangle$ that belong to CS_2 .

Step 3: Producing FS_2 . A candidate 2-sequence whose support is greater than or equal to min-supp is a frequent 2-sequence, and the set of all frequent 2-sequences is FS_2 .

Step 4: Finding the set of all the suitable time partitions for each frequent 2-sequences of FS_2 . For any frequent 2-sequence of $FS_2 \langle s_p, s_q \rangle$, all the time intervals between s_p and s_q appears in D are listed in increasing order, then the following clustering steps are used to find the suitable time partitions of all these intervals.

Step 4(a): Assume that $T(1, z) = [t_1, t_2, \dots, t_z]$ is the increasingly ordered list of the time intervals of $\langle s_p, s_q \rangle$. Let $T \langle s_p, s_q \rangle = \{ T(1, z) \}$ be the set of all the suitable time partitions of $\langle s_p, s_q \rangle$. The

first step of the clustering analysis is to find the maximal difference between two adjacent intervals from all partitions of $T \langle s_p, s_q \rangle$, and then divide the partition with the maximal difference into 2 partitions. At the beginning, only one partition is in $T \langle s_p, s_q \rangle$, therefore, the chosen partition is $T(1, z)$. Assume that the different between t_i and t_{i+1} is maximal, then $T(1, z)$ is divided into $T(1, i)$ and $T(i+1, z)$, where $T(1, i) = [t_1, \dots, t_i]$, $T(i+1, z) = [t_{i+1}, \dots, t_z]$.

Step 4(b): The second step is to calculate the support of $\langle s_p, s_q \rangle$ that respectively includes time intervals within these two partitions. If the support of $\langle s_p, s_q \rangle$ that includes time intervals within $T(1, i)$ is greater than or equal to min-supp , $T(1, i)$ is a suitable time partition of $\langle s_p, s_q \rangle$, and then $T(1, i)$ is reserved, otherwise $T(1, i)$ is deleted. Similarly, if the support of $\langle s_p, s_q \rangle$ that includes time intervals within $T(i+1, z)$ is greater than or equal to min-supp , $T(i+1, z)$ is a suitable time partition, and then $T(i+1, z)$ is reserved, otherwise $T(i+1, z)$ is deleted. Thus, the chosen partition is placed by the reserved partitions. If no partition is reserved, this chosen partition is considered as non-dividable. If all of the differences between two adjacent intervals in a partition are equal, this partition is considered as non-dividable as well.

Step 4(c): Repeating step 4(a) and step 4(b), until all time partitions in $T \langle s_p, s_q \rangle$ are non-dividable.

Step 5: Producing $FTIS_2$. Each 2-sequence of FS_2 is extended by all its suitable time partitions to form $FTIS_2$. Let $T \langle s_p, s_q \rangle = \{ T^1, T^2, \dots, T^R \}$ is the set of all the suitable time partitions of the sequence $\langle s_p, s_q \rangle$, then $T \langle s_p, T^i, s_q \rangle, i = 1 \dots R$, is a frequent time-interval 2-sequences.

Step 6: Producing $CTIS_k, k \geq 3$. For any two frequent time-interval (k-1)-sequences S_1 and S_2 , where $S_1 = \langle s_{1,1}, T_{1,1}, s_{1,2}, \dots, s_{1,k-2}, T_{1,k-2}, s_{1,k-1} \rangle$, $S_2 = \langle s_{2,1}, T_{2,1}, s_{2,2}, \dots, s_{2,k-2}, T_{2,k-2}, s_{2,k-1} \rangle \in FTIS_{k-1}$; $s_{1,2} = s_{2,1}, s_{1,3} = s_{2,2}, \dots, s_{1,k-1} = s_{2,k-2}$; $T_{1,2} = T_{2,1}, T_{1,3} = T_{2,2}, \dots, T_{1,k-2} = T_{2,k-3}$, then we can generate a candidate time-interval k-sequences $S_1 = \langle s_{1,1}, T_{1,1}, s_{1,2}, \dots, s_{1,k-2}, T_{1,k-2}, s_{1,k-1}, T_{2,k-2}, s_{2,k-1} \rangle$.

Step 7: Producing $FTIS_k, k \geq 3$. A candidate

time-interval k-sequence whose support is greater than or equal to $min-supp$ is a frequent time-interval k-sequence, and the set of all frequent time-interval k-sequences is $FTIS_k$.

Step 8: Repeating step 6 and step 7, until no next $CTIS_k$ can be generated.

According to these steps, an experiment using a simple sequence database will be displayed in the next section.

4. Experiment

In this section, we use the sequence database shown as in Table 1 to discover the time-interval sequential patterns. In Table 1, Id denotes the record number of a sequence, and each sequence is represented as $\langle (s_1, t_1), (s_2, t_2), \dots, (s_n, t_n) \rangle$, where s_i denotes an itemset, and t_i denotes the time stamp that s_i occurs; $min-supp$ is set as 0.3.

| Id | Sequence |
|------|--|
| 01 | $(s_5, 8), (s_4, 15), (s_6, 20)$ |
| 02 | $(s_1, 2), (s_3, 7), (s_2, 11), (s_6, 18)$ |
| 03 | $(s_2, 3), (s_1, 4), (s_3, 7), (s_6, 16), (s_7, 19)$ |
| 04 | $(s_1, 2), (s_2, 8), (s_6, 10), (s_7, 15)$ |
| 05 | $(s_5, 4), (s_6, 16), (s_1, 20), (s_3, 24)$ |
| 06 | $(s_7, 7), (s_1, 13), (s_5, 18), (s_2, 25), (s_6, 28)$ |
| 07 | $(s_5, 4), (s_1, 8), (s_3, 12), (s_6, 16), (s_7, 20)$ |
| 08 | $(s_1, 3), (s_5, 6), (s_2, 9), (s_4, 18), (s_6, 21)$ |
| 09 | $(s_2, 5), (s_1, 10), (s_3, 15), (s_6, 20), (s_7, 25)$ |
| 10 | $(s_6, 2), (s_7, 8), (s_5, 12), (s_2, 17)$ |

First, we need to calculate the supports of all itemsets to produce FS_1 . Supports of all itemsets are shown in Table 2. Here, we can obtain $FS_1 = \{s_1, s_2, s_3, s_5, s_6, s_7\}$.

| itemset | supp |
|---------|------|
| s_1 | 0.8 |
| s_2 | 0.7 |
| s_3 | 0.5 |
| s_4 | 0.2 |
| s_5 | 0.6 |
| s_6 | 1 |
| s_7 | 0.6 |

Next, CS_2 is generated by jointing $FS_1 \times FS_1$; Supports of the sequences in CS_2 are calculated and shown in Table 3. Therefore, we obtain $FS_3 = \{ \langle s_1, s_2 \rangle, \langle s_1, s_3 \rangle, \langle s_1, s_6 \rangle, \langle s_1, s_7 \rangle, \langle s_2, s_6 \rangle, \langle s_2, s_7 \rangle, \langle s_3, s_6 \rangle, \langle s_3, s_7 \rangle, \langle s_5, s_2 \rangle, \langle s_5, s_6 \rangle, \langle s_6, s_7 \rangle \}$.

| CS_2 | supp | CS_2 | supp |
|----------------------------|------|----------------------------|------|
| $\langle s_1, s_2 \rangle$ | 0.4 | $\langle s_5, s_1 \rangle$ | 0.2 |
| $\langle s_1, s_3 \rangle$ | 0.5 | $\langle s_5, s_2 \rangle$ | 0.3 |
| $\langle s_1, s_5 \rangle$ | 0.2 | $\langle s_5, s_3 \rangle$ | 0.2 |
| $\langle s_1, s_6 \rangle$ | 0.7 | $\langle s_5, s_6 \rangle$ | 0.5 |
| $\langle s_1, s_7 \rangle$ | 0.4 | $\langle s_5, s_7 \rangle$ | 0.1 |
| $\langle s_2, s_1 \rangle$ | 0.2 | $\langle s_6, s_1 \rangle$ | 0.1 |
| $\langle s_2, s_3 \rangle$ | 0.2 | $\langle s_6, s_2 \rangle$ | 0.1 |
| $\langle s_2, s_5 \rangle$ | 0.0 | $\langle s_6, s_3 \rangle$ | 0.1 |
| $\langle s_2, s_6 \rangle$ | 0.6 | $\langle s_6, s_5 \rangle$ | 0.1 |
| $\langle s_2, s_7 \rangle$ | 0.3 | $\langle s_6, s_7 \rangle$ | 0.5 |
| $\langle s_3, s_1 \rangle$ | 0.0 | $\langle s_7, s_1 \rangle$ | 0.1 |
| $\langle s_3, s_2 \rangle$ | 0.1 | $\langle s_7, s_2 \rangle$ | 0.2 |
| $\langle s_3, s_5 \rangle$ | 0.0 | $\langle s_7, s_3 \rangle$ | 0.0 |
| $\langle s_3, s_6 \rangle$ | 0.4 | $\langle s_7, s_5 \rangle$ | 0.2 |
| $\langle s_3, s_7 \rangle$ | 0.3 | $\langle s_7, s_6 \rangle$ | 0.1 |

For each frequent 2-sequence of FS_2 , all its time intervals are recorded and listed in increasing order (Table 4).

| FS_2 | time intervals |
|----------------------------|---|
| $\langle s_1, s_2 \rangle$ | $T\langle s_1, s_2 \rangle = \{6, 9, 12\}$ |
| $\langle s_1, s_3 \rangle$ | $T\langle s_1, s_3 \rangle = \{3, 4, 5\}$ |
| $\langle s_1, s_6 \rangle$ | $T\langle s_1, s_6 \rangle = \{8, 10, 12, 16, 18\}$ |
| $\langle s_1, s_7 \rangle$ | $T\langle s_1, s_7 \rangle = \{12, 13, 15\}$ |
| $\langle s_2, s_6 \rangle$ | $T\langle s_2, s_6 \rangle = \{2, 3, 7, 12, 13, 15\}$ |
| $\langle s_2, s_7 \rangle$ | $T\langle s_2, s_7 \rangle = \{7, 16, 20\}$ |
| $\langle s_3, s_6 \rangle$ | $T\langle s_3, s_6 \rangle = \{4, 5, 9, 11\}$ |
| $\langle s_3, s_7 \rangle$ | $T\langle s_3, s_7 \rangle = \{8, 10, 12\}$ |
| $\langle s_5, s_2 \rangle$ | $T\langle s_5, s_2 \rangle = \{3, 5, 7\}$ |
| $\langle s_5, s_6 \rangle$ | $T\langle s_5, s_6 \rangle = \{10, 12, 15\}$ |
| $\langle s_6, s_7 \rangle$ | $T\langle s_6, s_7 \rangle = \{3, 4, 5, 6\}$ |

According to the step 4 described in the above section, the set of all suitable time partitions for

each sequences of FS_2 are obtained as in Table 5. Next, each 2-sequence of FS_2 is extended by all its suitable time partitions to form $FTIS_2$ (Table 6).

| FS_2 | suitable time partitions |
|----------------------------|---|
| $\langle s_1, s_2 \rangle$ | $T_{1,2}^1 = [6,12]$ |
| $\langle s_1, s_3 \rangle$ | $T_{1,3}^1 = [3,5]$ |
| $\langle s_1, s_6 \rangle$ | $T_{1,6}^1 = [8,12], T_{1,6}^2 = [16,18]$ |
| $\langle s_1, s_7 \rangle$ | $T_{1,7}^1 = [12,15]$ |
| $\langle s_2, s_6 \rangle$ | $T_{2,6}^1 = [2,7], T_{2,6}^2 = [12,15]$ |
| $\langle s_2, s_7 \rangle$ | $T_{2,7}^1 = [16,20]$ |
| $\langle s_3, s_6 \rangle$ | $T_{3,6}^1 = [4,11]$ |
| $\langle s_3, s_7 \rangle$ | $T_{3,7}^1 = [8,12]$ |
| $\langle s_5, s_2 \rangle$ | $T_{5,2}^1 = [3,7]$ |
| $\langle s_5, s_6 \rangle$ | $T_{5,6}^1 = [10,12]$ |
| $\langle s_6, s_7 \rangle$ | $T_{6,7}^1 = [3,6]$ |

| |
|--|
| $\langle s_1, T_{1,2}^1, s_2 \rangle$ |
| $\langle s_1, T_{1,3}^1, s_3 \rangle$ |
| $\langle s_1, T_{1,6}^1, s_6 \rangle, \langle s_1, T_{1,6}^2, s_6 \rangle$ |
| $\langle s_1, T_{1,7}^1, s_7 \rangle$ |
| $\langle s_2, T_{2,6}^1, s_6 \rangle, \langle s_2, T_{2,6}^2, s_6 \rangle$ |
| $\langle s_2, T_{2,7}^1, s_7 \rangle$ |
| $\langle s_3, T_{3,6}^1, s_6 \rangle$ |
| $\langle s_3, T_{3,7}^1, s_7 \rangle$ |
| $\langle s_5, T_{5,2}^1, s_2 \rangle$ |
| $\langle s_5, T_{5,6}^1, s_6 \rangle$ |
| $\langle s_6, T_{6,7}^1, s_7 \rangle$ |

$CTIS_3$, the set of candidate time-interval 3-sequences is generated by jointing $FTIS_2 \times FTIS_2$. Supports of the sequences of $CTIS_3$ are calculated and shown in Table 7. A candidate time-interval 3-sequence whose support is greater than or equal to $min-supp$ is called as a frequent time-interval 3- sequence. Therefore, we can obtain the set of all the frequent time-interval 3-sequences,

$FTIS_4 = \{ \langle s_1, T_{1,2}^1, s_2, T_{2,6}^1, s_6 \rangle, \langle s_1, T_{1,3}^1, s_3, T_{3,6}^1, s_6 \rangle, \langle s_1, T_{1,3}^1, s_3, T_{3,7}^1, s_7 \rangle, \langle s_1, T_{1,6}^1, s_6, T_{6,7}^1, s_7 \rangle, \langle s_3, T_{3,6}^1, s_6, T_{6,7}^1, s_7 \rangle \}$.

| $CTIS_3$ | supp |
|---|------|
| $\langle s_1, T_{1,2}^1, s_2, T_{2,6}^1, s_6 \rangle$ | 0.3 |
| $\langle s_1, T_{1,2}^1, s_2, T_{2,6}^2, s_6 \rangle$ | 0.1 |
| $\langle s_1, T_{1,2}^1, s_2, T_{2,7}^1, s_7 \rangle$ | 0.1 |
| $\langle s_1, T_{1,3}^1, s_3, T_{3,6}^1, s_6 \rangle$ | 0.4 |
| $\langle s_1, T_{1,3}^1, s_3, T_{3,7}^1, s_7 \rangle$ | 0.3 |
| $\langle s_1, T_{1,6}^1, s_6, T_{6,7}^1, s_7 \rangle$ | 0.4 |
| $\langle s_2, T_{2,6}^1, s_6, T_{6,7}^1, s_7 \rangle$ | 0.1 |
| $\langle s_2, T_{2,6}^2, s_6, T_{6,7}^1, s_7 \rangle$ | 0.2 |
| $\langle s_3, T_{3,6}^1, s_6, T_{6,7}^1, s_7 \rangle$ | 0.3 |
| $\langle s_5, T_{5,2}^1, s_2, T_{2,7}^1, s_7 \rangle$ | 0 |
| $\langle s_5, T_{5,6}^1, s_6, T_{6,7}^1, s_7 \rangle$ | 0.1 |

The set of candidate time-interval 4-sequences, $CTIS_4$, is generated by jointing $FTIS_3 \times FTIS_3$. Here, only one sequence, $\langle s_1, T_{1,3}^1, s_3, T_{3,6}^1, s_6, T_{6,7}^1, s_7 \rangle$, is generated. The support of the sequence $\langle s_1, T_{1,3}^1, s_3, T_{3,6}^1, s_6, T_{6,7}^1, s_7 \rangle$ is 0.3, thus $\langle s_1, T_{1,3}^1, s_3, T_{3,6}^1, s_6, T_{6,7}^1, s_7 \rangle$ is a frequent time-interval 4-sequences, and we obtain $FTIS_4 = \{ \langle s_1, T_{1,3}^1, s_3, T_{3,6}^1, s_6, T_{6,7}^1, s_7 \rangle \}$. Because no next $CTIS_5$ can be generated, the algorithm stops here.

According to the experimental results, we can clearly see that the suitable time partitions for every pair of successive itemsets are different and overlap, therefore, for every pair of successive itemsets, it is more reasonable to generate the suitable time partitions directly from the real data when mining time-interval sequential patterns.

5. Conclusions

A sequential pattern with the time intervals between successive itemsets is more valuable than a traditional sequential pattern without any time information. Most existing time-interval sequential pattern mining algorithms reveal the time-interval between itemsets by predefining some non-overlap time partitions, but a predefined set of non-overlap time partitions, in fact, cannot be suitable for every

pair of successive itemsets.

In this paper, we present a new algorithm to mine time-interval sequential patterns without pre-defining any time partitions. This algorithm use the clustering analysis to automatically generate the suitable time partitions between frequent occurring pairs of successive itemsets, and then uses these generated time partitions to extend typical algorithms to discover the time-interval sequential patterns. The result of our experiment verifies that the concept of our algorithm is more reasonable than these algorithms which use some predefined and non-overlap time partitions.

References

- [1] R. Agrawal, and R. Srikant, "Mining sequential patterns," *Proceedings of the International Conference on Data Engineering*, pp. 3–14, 1995.
- [2] Y. L. Chen, M. C. Chiang, and M. T. Ko, "Discovering time-interval sequential patterns in sequence databases," *Expert Systems with Applications*, vol.25, no. 3, pp. 343–354, 2003.
- [3] M. S. Chen, J. Han, and P. S. Yu, "Data mining: An overview from a database perspective," *IEEE Transactions on Knowledge and Data Engineering*, vol.8, no.6, pp. 866–883, 1996.
- [4] M. S. Chen, J. S. Park, and P. S. Yu, "Efficient data mining for path traversal patterns," *IEEE Transactions on Knowledge and Data Engineering*, vol. 10, no. 2, pp. 209–221, 1998.
- [5] M. H. Dunham, *Data mining, Introductory and Advanced Topics*, Pearson Education Inc., 2003.
- [6] J. Han, G. Dong, and Y. Yin, "Efficient mining of partial periodic patterns in time series database," *Proceedings of 1999 International Conference on Data Engineering*, pp. 106–115, 1999.
- [7] J. Han, and M. Kamber, *Data mining: Concepts and Techniques*, Academic Press, 2001.
- [8] H. Mannila, H. Toivonen, and A. Inkeri Verkamo, "Discovery of frequent episodes in event sequences," *Data Mining and Knowledge Discovery*, vol. 1, no.3, pp. 259–289, 1997.
- [9] J. Pei, J. Han, H. Pinto, Q. Chen, U. Dayal, and M.-C. Hsu, "PrefixSpan: Mining sequential patterns efficiently by prefix-projected pattern growth," *Proceedings of 2001 International Conference on Data Engineering*, pp. 215–224, 2001.
- [10] R. Srikant, and R. Agrawal, "Mining sequential patterns: Generalizations and performance improvements," *Proceedings of the 5th International Conference on Extending Database Technology*, pp. 3–17, 1996.
- [11] P. H. Wu, W. C. Peng, and M. S. Chen, "Mining sequential alarm patterns in a telecommunication database," *Proceedings of Workshop on Databases in Telecommunications (VLDB 2001)*, pp. 37-51, 2001.