

## Mining Fuzzy Sequential Patterns with Fuzzy Time-Intervals in Quantitative Sequence Databases

Truong Duc Phuong<sup>1</sup>, Do Van Thanh<sup>2</sup>, Nguyen Duc Dung<sup>3</sup>

<sup>1</sup>Hanoi Metropolitan University, Vietnam

<sup>2</sup>Faculty of Information Technology, Nguyen Tat Thanh University, Vietnam

<sup>3</sup>Institute of Information Technology, Vietnamese Academy of Science and Technology, Vietnam

E-mails: tdphuong@daihocthudo.edu.vn dvthanh@ntt.edu.vn nddung@ioit.ac.vn

**Abstract:** The main objective of this paper is to introduce fuzzy sequential patterns with fuzzy time-intervals in quantitative sequence databases. In the fuzzy sequential pattern with fuzzy time-intervals, both quantitative attributes and time distances are represented by linguistic terms. A new algorithm based on the Apriori algorithm is proposed to find the patterns. The mined patterns can be applied to market basket analysis, stock market analysis, and so on.

**Keywords:** Data mining, fuzzy sequential pattern, fuzzy time-interval, sequence database.

### 1. Introduction

Mining sequential patterns are one of the most important domains in data mining. Mining sequential patterns from transaction databases (events present or not) is the first introduced in 1995 [1] and there are many related works [8-10, 17, 22, 23]. The sequential pattern presented a relationship between events in a chronological sequence with the same object. For instance, “**If** a customer buys a Laptop and later a Modem, **then** he will buy a Printer”. Mining fuzzy sequential patterns in sequence databases, in which values of attributes are numeric or categorical, is also introduced in [2, 5, 12, 14, 15, 18]. In these works, the values are transformed into fuzzy sets. A relationship between events in the fuzzy sequential patterns is like “**If** a customer buys a Large number of Laptops and later an Average number of Modems, **then** he will buy a Small number of Printers”.

In many cases, the values of time distances among the events in a sequence are preferred. The time-interval sequential patterns in transaction databases were investigated by authors such as Chen and Huang [6], Chen, Chiang and Ko [7], Chang, Chueh and Lin [3], Chang, Chueh and Luo [4]. In [7], time-interval sequential patterns were presented such as  $\langle \text{Laptop}, I_1, \text{Modem}, I_2, \text{Printer} \rangle$ , which meant “**If** a customer buys a Laptop and later a Modem an interval of  $I_1$ , **then**

he will buy a Printer an interval of  $I_2$ ”, where  $I_1$  and  $I_2$  were predetermined time ranges. For instance,  $I_1$  was ranged from 3 to 5 days,  $I_2$  from 10 to 12 days. To solve sharp boundary problems when a time interval is near the boundary of two time ranges in [7], the fuzzy theory was applied to time intervals [6]. The fuzzy time-interval sequential patterns in [6] showed the relationship among events such as  $\langle \text{Laptop, Short, Modem, Long, Printer} \rangle$ , which meant “**If** a customer buys a Laptop and later a Modem an interval of Short, **then** he will buy a Printer an interval of Long”, where Short and Long were linguistic terms for time intervals. In [6], time intervals among events were transformed into fuzzy sets and the two algorithms were proposed. They were the FTI-Apriori algorithm based on the idea of an Apriori algorithm and the FTI-PrefixSpan algorithm based on the PrefixSpan algorithm. Papers [3, 4] also revealed fuzzy time-interval sequential patterns in sequence databases such as  $\langle \text{Laptop, } \mu_{\text{Laptop\_Modem}}, \text{Modem, } \mu_{\text{Modem\_Printer}}, \text{Printer} \rangle$  that meant “**If** a customer buys a Laptop and later a Modem an interval of  $\mu_{\text{Laptop\_Modem}}$ , **then** he will buy a Printer an interval of  $\mu_{\text{Modem\_Printer}}$ ”, where  $\mu_{\text{Laptop\_Modem}}$  and  $\mu_{\text{Modem\_Printer}}$  were trapezoidal fuzzy numbers which present time intervals between events of pair (Laptop, Modem) and (Modem, Printer). The fuzzy numbers were computed by frequency of time intervals of events in a sequence. For example,  $\mu_{\text{Laptop\_Modem}} = (6, 6, 12, 12)$  and  $\mu_{\text{Modem\_Printer}} = (2, 2, 7, 15)$ . The main idea of the SPFTI algorithm [3] and the ISPFTI algorithm [4] were taking the advantage of the idea of Apriori algorithm, but the trapezoidal fuzzy numbers were received from databases unlike predetermined fuzzy sets in [6].

In [21], we considered to mine fuzzy association rules with fuzzy time-intervals from quantitative databases. The results showed that the patterns such as  $\langle \text{Laptop\_Large, Short, Modem\_Average, Long, Printer\_Small} \rangle$  that meant “**If** a Large number of Laptops are sold and later an Average number of Modems an interval of Short, **then** a Small number of Printers will be sold an interval of Long”. In the pattern, Laptop\_Large, Modem\_Average, and Printer\_Small were the fuzzy sets of attributes; Short and Long were the pre-defined fuzzy sets for time intervals.

Works [3, 4, 6, 7] have indicated only the fuzzy time-interval sequential patterns in transaction sequence databases, not yet in quantitative sequence databases. While paper [21] has been applied to only quantitative databases, not yet quantitative sequence databases.

So, the main objective of the paper is to introduce an algorithm for mining fuzzy sequential patterns with fuzzy time-intervals in quantitative sequence databases. They are like  $\langle \langle \text{Laptop\_Large, Short, Modem\_Average, Long, Printer\_Small} \rangle \rangle$  that means “**If** a customer buys a Large number of Laptops and later an Average number of Modems an interval of Short, **then** he will buy a Small number of Printers an interval of Long”, where Short and Long are predetermined linguistic terms for time intervals. These patterns are different from the ones in [3, 4, 6, 7] by quantitative attributes and from the ones in [21] by interesting in the objects. The main idea of the algorithm is to convert quantitative attributes and time intervals to linguistic terms in the same way as in [21] and then improving the Apriori algorithm [1] to find fuzzy sequential patterns with fuzzy time-intervals. The comparison of datasets and patterns of the selected algorithms is described in Table 1.

Table 1. Comparison of datasets and patterns of the selected algorithms

Algorithm	Dataset	Pattern	Description
AprioriAll [1]	Transactional sequence databases	$\langle \text{Laptop, Modem, Printer} \rangle$	<b>If</b> a customer buys a Laptop and later a Modem, <b>then</b> he will buy a Printer
FGBSPMA [5]	Quantitative sequence databases	$\langle \text{Laptop}_{\text{Large}}, \text{Modem}_{\text{Average}}, \text{Printer}_{\text{Small}} \rangle$	<b>If</b> a customer buys a Large number of Laptops and later an Average number of Modems, <b>then</b> he will buy a Small number of Printers
I-Apriori algorithm, I-PrefixSpan [7]	Transactional sequence databases	$\langle \text{Laptop}, I_1, \text{Modem}, I_2, \text{Printer} \rangle$	<b>If</b> a customer buys a Laptop and later a Modem an interval of $I_1$ , <b>then</b> he will buy a Printer an interval of $I_2$
FTI-Apriori, FTI-PrefixSpan [6]	Transactional sequence databases	$\langle \text{Laptop}, \text{Short}, \text{Modem}, \text{Long}, \text{Printer} \rangle$	<b>If</b> a customer buys a Laptop and later a Modem an interval of Short, <b>then</b> he will buy a Printer an interval of Long
SPFTI [3], ISPFTI [4]	Transactional sequence databases	$\langle \text{Laptop}, \mu_{\text{Laptop\_Modem}}, \text{Modem}, \mu_{\text{Modem\_Printer}}, \text{Printer} \rangle$	<b>If</b> a customer buys a Laptop and later a Modem an interval of $\mu_{\text{Laptop\_Modem}}$ <b>then</b> he will buy a Printer an interval of $\mu_{\text{Modem\_Printer}}$
FTQ [21]	Quantitative (not sequence) databases	$\langle \text{Laptop}_{\text{Large}}, \text{Short}, \text{Modem}_{\text{Average}}, \text{Long}, \text{Printer}_{\text{Small}} \rangle$	<b>If</b> a Large number of Laptops are sold and later an Average number of Modems an interval of Short, <b>then</b> a Small number of Printers will be sold an interval of Long
<b>Proposed</b>	<b>Quantitative sequence databases</b>	$\langle \langle \text{Laptop}_{\text{Large}}, \text{Short}, \text{Modem}_{\text{Average}}, \text{Long}, \text{Printer}_{\text{Small}} \rangle \rangle$	<b>If</b> a customer buys a Large number of Laptops and later an Average number of Modems an interval of Short, <b>then</b> he will buy a Small number of Printers an interval of Long

The rest of the paper is organized as follows: Section 2 defines problems. Section 3 develops the FSPFTIM algorithm to find out the fuzzy sequential patterns with fuzzy time-intervals and gives an example. Section 4 presents experimental results. Conclusions are drawn in Section 5.

## 2. Problem definition

This section presents some definitions related to the proposed problem.

**Definition 1.** Let  $E = \{e_1, e_2, \dots, e_u\}$  be a set of attributes,  $s = \langle (a_1, q_1, t_1), (a_2, q_2, t_2), (a_3, q_3, t_3), \dots, (a_n, q_n, t_n) \rangle$  be a quantitative sequence, where  $a_k \in E$  is an attribute ( $1 \leq k \leq n$ ),  $t_k$  is time of  $a_k$ ,  $t_{k-1} \leq t_k$  with  $2 \leq k \leq n$  and  $a_k(t_k) = q_k$ ,  $q_k$  is numeric or categorical. A transaction sequence is  $\langle \text{Sid}, s \rangle$  where  $s$  is a quantitative sequence and Sid is the identifier of the sequence. A quantitative sequence database is a set of all transaction sequences.

**Example 1.** A quantitative sequence database is shown in Table 2.

In Table 2,  $E = \{a, b, c, d, e, f, g, h, i\}$  is set of attributes, time is started at 0. The first transaction sequence is  $\langle 1, (a, 2, 1), (b, 2, 4), (e, 5, 29) \rangle$ . It means that this sequence includes Sid with value 1 and three transactions: The  $a$  with value 2

(quantitative of  $a$  is 2) at time 1, the  $b$  with value 2 at time 4 and the  $e$  with value 5 at time 29.

Table 2. An example of a quantitative sequence database

Sid	Quantitative sequence
1	$\langle\langle a, 2, 1 \rangle, \langle b, 2, 4 \rangle, \langle e, 5, 29 \rangle\rangle$
2	$\langle\langle d, 2, 1 \rangle, \langle a, 5, 2 \rangle, \langle d, 4, 24 \rangle\rangle$
3	$\langle\langle b, 1, 1 \rangle, \langle d, 2, 11 \rangle, \langle e, 5, 28 \rangle\rangle$
4	$\langle\langle f, 6, 1 \rangle, \langle b, 6, 5 \rangle, \langle c, 1, 19 \rangle, \langle c, 2, 25 \rangle\rangle$
5	$\langle\langle a, 1, 4 \rangle, \langle b, 1, 5 \rangle, \langle d, 2, 10 \rangle, \langle e, 5, 28 \rangle\rangle$
6	$\langle\langle a, 2, 0 \rangle, \langle b, 1, 5 \rangle, \langle e, 1, 30 \rangle\rangle$
7	$\langle\langle i, 5, 2 \rangle, \langle a, 3, 17 \rangle, \langle h, 2, 17 \rangle\rangle$
8	$\langle\langle c, 6, 3 \rangle, \langle i, 5, 10 \rangle, \langle f, 3, 18 \rangle\rangle$
9	$\langle\langle h, 3, 4 \rangle, \langle a, 1, 10 \rangle, \langle b, 6, 21 \rangle\rangle$
10	$\langle\langle a, 2, 0 \rangle, \langle g, 5, 0 \rangle, \langle b, 2, 3 \rangle, \langle e, 1, 30 \rangle\rangle$

**Definition 2.** Let  $FE = \{F^{e_1}, F^{e_2}, \dots, F^{e_u}\}$  be a set of fuzzy sets of attributes of  $E$ ,  $F^{e_k} = \{f_{h_k,1}^{e_k}, f_{h_k,2}^{e_k}, \dots, f_{h_k,h_k}^{e_k}\}$  be a set of fuzzy sets of  $e_k$  attribute ( $k=1, 2, \dots, u$ ), where  $f_{h_k,j}^{e_k}$  is the  $j$ -th fuzzy set ( $1 \leq j \leq h_k$ ),  $h_k$  is the number of fuzzy sets of  $e_k$ ;  $f_{h_k,j}^{e_k}$  is called a fuzzy attribute. Each fuzzy set has its membership function  $\mu: X \rightarrow [0, 1]$ . Sequence  $fs = \langle (fa_1, fq_1, t_1), (fa_2, fq_2, t_2), \dots, (fa_n, fq_n, t_n) \rangle$  is called a fuzzy one, where  $fa_i \in F^{a_i}$  ( $1 \leq i \leq n$ ) is a fuzzy set,  $fq_i$  is the value of the membership function  $\mu_{fa_i}$  of  $fa_i$  at  $q_i$  ( $fq_i = \mu_{fa_i}(q_i)$ ),  $fa_i$  is called a fuzzy attribute.  $\langle Sid, fs \rangle$  is called a *fuzzy transaction sequence*. A fuzzy sequence database is a set of all fuzzy transaction sequences.

**Example 2.** Given fuzzy sets are in [13],  $A_{K,i_m}^{x_m}$  is  $i_m$ -th linguistic value ( $1 \leq i_m \leq K$ ) of the  $x_m$  attribute ( $x_m \in E$ ),  $K$  is the number of partitions of  $x_m$ ,  $\mu_{K,i_m}^{x_m}$  is the membership function of  $A_{K,i_m}^{x_m}$  that is determined as follows:

$$(1) \quad \mu_{K,i_m}^{x_m}(v) = \max\{1 - |v - a_{i_m}^K| / b^K, 0\},$$

where

$$a_{i_m}^K = mi + (ma - mi)(i_m - 1) / (K - 1),$$

$$b^K = (ma - mi) / (K - 1),$$

where  $mi$  is the minimum value of  $x_m$  attribute domain, and  $ma$  is the maximum value. With  $K=3$  for all attributes in the database of the Example 1, we get membership functions as in the Fig. 1 and the fuzzy transaction sequence database  $D'$  is described in Table 3.

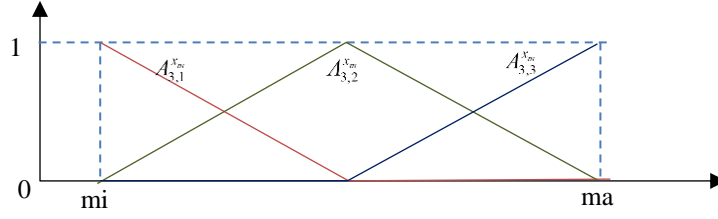


Fig. 1. The membership functions of  $x_m$  with  $K=3$

Table 3. The fuzzy transaction sequence database  $D'$

Sid	Fuzzy sequence
1	$\langle (f_{3,1}^a, 0.5, 1), (f_{3,2}^a, 0.5, 1), (f_{3,1}^b, 0.6, 4), (f_{3,2}^b, 0.4, 4), (f_{3,3}^e, 1, 29) \rangle$
2	$\langle (f_{3,1}^d, 1, 1), (f_{3,3}^a, 1, 2), (f_{3,1}^d, 1, 24) \rangle$
3	$\langle (f_{3,1}^b, 1, 1), (f_{3,1}^d, 1, 11), (f_{3,3}^e, 1, 28) \rangle$
4	$\langle (f_{3,3}^f, 1, 1), (f_{3,3}^b, 1, 5), (f_{3,1}^c, 1, 19), (f_{3,1}^c, 0.6, 25), (f_{3,2}^c, 0.4, 25) \rangle$
5	$\langle (f_{3,1}^a, 1, 4), (f_{3,1}^b, 1, 5), (f_{3,1}^d, 1, 10), (f_{3,3}^e, 1, 28) \rangle$
6	$\langle (f_{3,1}^a, 0.5, 0), (f_{3,2}^a, 0.5, 0), (f_{3,1}^b, 1, 5), (f_{3,1}^e, 1, 30) \rangle$
7	$\langle (f_{3,1}^i, 1, 2), (f_{3,2}^a, 1, 17), (f_{3,1}^h, 1, 17) \rangle$
8	$\langle (f_{3,3}^c, 1, 3), (f_{3,1}^i, 1, 10), (f_{3,1}^f, 1, 18) \rangle$
9	$\langle (f_{3,3}^h, 1, 4), (f_{3,1}^a, 1, 10), (f_{3,3}^b, 1, 21) \rangle$
10	$\langle (f_{3,1}^a, 0.5, 0), (f_{3,2}^a, 0.5, 0), (f_{3,1}^g, 1, 0), (f_{3,1}^b, 0.6, 3), (f_{3,2}^b, 0.4, 3), (f_{3,1}^e, 1, 30) \rangle$

Each fuzzy tuple  $(f_{K,i_m}^x, fq, t)$  in the  $D'$  above means that the fuzzy attribute  $f_{K,i_m}^x$  is the  $i_m$ -th fuzzy set of  $x$  attribute,  $fq \in [0, 1]$  is the value of the membership function of  $f_{K,i_m}^x$  at  $q$  and  $t$  is the time of the event. For example, with  $(f_{3,2}^a, 0.5, 1)$ , the value of membership function of the fuzzy set  $f_{3,2}^a$  at  $a=2$  is 0.5, the transaction time is 1.

**Definition 3.** Let  $LT = \{lt_j | j=1, 2, \dots, p\}$  be fuzzy sets of time-intervals,  $\mu_{lt_j} : X \rightarrow [0, 1]$  be the membership function of the fuzzy set  $lt_j$  [13]. Then,  $\alpha = \langle \langle b_1, ltb_1, b_2, ltb_2, \dots, b_{r-1}, ltb_{r-1}, b_r \rangle \rangle$  is called a *fuzzy sequence with fuzzy time-intervals* if  $b_i, 1 \leq i \leq r$  is a fuzzy set and  $ltb_i \in LT, 1 \leq i \leq r-1$ . A fuzzy sequence with fuzzy time-intervals whose length is  $r$  referred to  $r$ -fuzzy sequential pattern with fuzzy time-intervals.

**Example 3.** Let  $LT = \{\text{Short, Medium, Long}\}$  be fuzzy sets of time intervals. Its membership functions as in the Fig. 2 [6] and the fuzzy sequence database  $D'$  in Table 3, then  $\alpha = \langle \langle f_{3,1}^a, \text{Short}, f_{3,1}^b, \text{Medium}, f_{3,1}^e \rangle \rangle$  is a fuzzy sequence with fuzzy time-intervals whose length is 3 or 3-fuzzy sequence with fuzzy time-intervals.

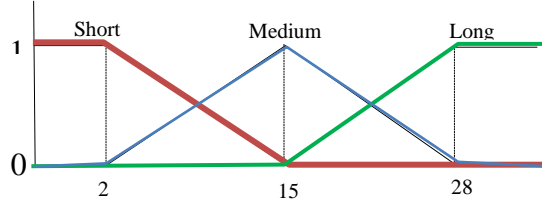


Fig. 2. The membership functions of the fuzzy sets in LT

**Definition 4.** Given a fuzzy sequence  $B = \langle (b_1, bq_1, bt_1), (b_2, bq_2, bt_2), \dots, (b_r, bq_r, bt_r) \rangle$  and a fuzzy sequence with fuzzy time-intervals  $\alpha = \langle \langle b_1, ltb_1, b_2, ltb_2, \dots, b_{r-1}, ltb_{r-1}, b_r \rangle \rangle$ , we define:

- The support of  $B$  for  $\alpha$ , denoted by  $\gamma_B(\alpha)$ , is

$$(2) \quad \gamma_B(\alpha) = \begin{cases} bq_1 & \text{if } r = 1, \\ \left( \prod_{i=1}^r bq_i \right) \times \min_{1 \leq j \leq r-1} \{ \mu_{lbt_j}(bt_{j+1} - bt_j) \} & \text{if } r > 1. \end{cases}$$

- A fuzzy sequence  $B$  belongs to a fuzzy transaction sequence  $S = \langle \text{Sid}, fs \rangle$  where  $fs = \langle (fa_1, fq_1, t_1), (fa_2, fq_2, t_2), \dots, (fa_n, fq_n, t_n) \rangle$  if there exists a integer  $r$  such that  $b_k = fa_{i_k} \wedge bq_k = fq_{i_k} \wedge bt_k = t_{i_k}$ ,  $1 \leq k \leq r$ , and  $i_1 < i_2 < \dots < i_r$ . The support of the fuzzy transaction sequence  $S$  for  $\alpha$ , denoted by  $\text{Supp}_S(\alpha)$ , is the maximum of the supports of  $B$  which belong to  $S$ ,

$$(3) \quad \text{Supp}_S(\alpha) = \max_{B \in S} (\gamma_B(\alpha)).$$

- The support of a fuzzy sequence with fuzzy time-intervals  $\alpha$ , denoted  $\text{Supp}(\alpha)$ , is the average of supports of all fuzzy transaction sequences in  $D'$  for  $\alpha$ .

$$(4) \quad \text{Supp}(\alpha) = \frac{\sum_{i=1}^{NS} \text{Supp}_{S_i}(\alpha)}{NS},$$

where  $S_i$  is the  $i$ -th fuzzy transaction sequence in  $D'$ ,  $NS$  is the number of fuzzy transaction sequences in  $D'$ .

- A fuzzy sequential pattern with fuzzy time-intervals is the fuzzy sequence with fuzzy time-intervals whose support is not less than a user-defined threshold.

**Example 4.** Given the fuzzy transaction sequence database  $D'$ , the membership functions of fuzzy sets of LT mentioned in the Example 3, we compute the support of two fuzzy sequences with fuzzy time-intervals  $\langle \langle f_{3,1}^a, \text{Short}, f_{3,1}^b \rangle \rangle$  and  $\langle \langle f_{3,3}^f, \text{Short}, f_{3,2}^b, \text{Medium}, f_{3,1}^c \rangle \rangle$ .

**Case 1.**  $\alpha = \langle \langle f_{3,1}^a, \text{Short}, f_{3,1}^b \rangle \rangle$ .

$\text{Supp}_{S_1}(\alpha) = 0.5 \times 0.667 \times \mu_{\text{Short}}(4 - 1) = 0.5 \times 0.667 \times 0.923 = 0.308$ , because the first fuzzy transaction sequence has only the fuzzy sequence  $\langle \langle f_{3,1}^a, 0.5, 1 \rangle, \langle f_{3,1}^b, 0.667, 4 \rangle \rangle$  fitting  $\alpha$ . Similarly, the supports of the 5th fuzzy transaction sequence and the 6th fuzzy transaction sequence are  $\text{Supp}_{S_5}(\alpha) = 1 \times 1 \times \mu_{\text{Short}}(5 - 4) = 1 \times 1 \times 1 = 1$  and  $\text{Supp}_{S_6}(\alpha) = 0.5 \times 1 \times \mu_{\text{Short}}(5 - 0) = 0.5 \times 0.769 = 0.385$ , respectively. Not all the rest fuzzy

transaction sequences support for  $\alpha$ . So that, the support of the fuzzy sequence with fuzzy time-intervals  $\alpha$  is computed as follows:

$$\text{Supp}(\alpha) = (0.308 + 0 + 0 + 0 + 1 + 0.385 + 0 + 0 + 0 + 0) / 10 = 0.1693.$$

**Case 2.**  $\beta = \langle \langle f_{3,3}^f, \text{Short}, f_{3,2}^b, \text{Medium}, f_{3,1}^c \rangle \rangle$ .

In  $D'$ , there is only the 4th fuzzy transaction sequence,  $S_4$ , containing fuzzy sequences fitting  $\beta$ ,

$S_4 = \langle 4, (f_{3,3}^f, 1, 1), (f_{3,2}^b, 0.333, 5), (f_{3,3}^b, 0.667, 5), (f_{3,1}^c, 1, 19), (f_{3,1}^c, 0.6, 25), (f_{3,2}^c, 0.4, 25) \rangle$  has two fuzzy sequences, fitting  $\beta$  that are  $B_1, B_2$ :

$B_1 = \langle (f_{3,3}^f, 1, 1), (f_{3,2}^b, 0.333, 5), (f_{3,3}^c, 1, 19) \rangle$  has the support for  $\beta$ ,

$$\gamma_{B_1}(\beta) = 1 \times 0.333 \times 1 \times \min\{\mu_{\text{Short}}(4), \mu_{\text{Medium}}(14)\} = 0.333 \times \min\{0.846, 0.923\} = 0.333 \times 0.846 = 0.285,$$

$B_2 = \langle (f_{3,3}^f, 1, 1), (f_{3,2}^b, 0.333, 5), (f_{3,1}^c, 0.6, 25) \rangle$  has the support for  $\beta$ ,

$$\gamma_{B_2}(\beta) = 1 \times 0.333 \times 0.6 \times \min\{\mu_{\text{Short}}(4), \mu_{\text{Medium}}(20)\} = 0.2 \times \min\{0.846, 0.615\} = 0.2 \times 0.615 = 0.123.$$

Hence, the support of  $S_4$  for  $\beta$  is  $\text{Supp}_{S_4}(\beta) = \max\{\gamma_{B_1}(\beta), \gamma_{B_2}(\beta)\} = 0.285$ .

Consequently,  $\text{Supp}(\beta) = (0 + 0 + 0 + 0.285 + 0 + 0 + 0 + 0 + 0 + 0) / 10 = 0.0285$ .

If the user-defined threshold is of value 0.1 then the fuzzy sequence with fuzzy time-intervals  $\langle \langle f_{3,1}^a, \text{Short}, f_{3,1}^b \rangle \rangle$  is a fuzzy sequential pattern with fuzzy time-intervals because its support 0.1693 is greater than  $\text{min\_sup}$ , whereas the fuzzy sequence with fuzzy time-intervals  $\langle \langle f_{3,3}^f, \text{Short}, f_{3,2}^b, \text{Medium}, f_{3,1}^c \rangle \rangle$ , having the support of  $0.0285 < 0.1$ , is not a sequential pattern.

### 3. Algorithm of Mining Fuzzy Sequential Patterns with Fuzzy Time-Intervals – FSPFTIM Algorithm

#### 3.1. The problem

**Input:**  $D$  is a quantitative sequence database,  $\text{min\_sup}$  is a user-defined threshold,  $FE$  is a set of fuzzy sets with its membership functions of attributes in  $D$ ,  $LT$  is a set of fuzzy sets with its membership functions of time intervals.

**Output:**  $k$ -fuzzy sequential patterns with fuzzy time intervals,  $k > 1$

#### 3.2. Proposed approach

First, all the quantitative attributes in  $D$  are partitioned into fuzzy sets, then the transaction sequences in  $D$  are transformed into fuzzy transaction sequences based on these fuzzy sets and their corresponding membership functions. And we get a fuzzy sequence database, called  $D'$ . Next, the FSPFTIM Algorithm is applied to find out fuzzy sequential patterns with fuzzy time-intervals. This algorithm is improved from the Apriori algorithm.

The algorithm is an iterative process consisting of several steps. At the  $k$ -th step, the algorithm generates a set denoted  $C_k$  of candidate  $k$ -fuzzy sequences with fuzzy time intervals and then calculates the support of these candidate sequences. The sequences having the support greater than  $\min\_sup$  will be added into  $L_k$  as a set of  $k$ -fuzzy sequential patterns with fuzzy time-intervals. This process will be ended when it is unable to generate any new set of fuzzy sequential patterns with fuzzy time-intervals.

More specifically, the process of generating the sets  $C_k$  is as follows:

**For  $k=1$ .** All fuzzy attributes of  $D'$  are added into the set  $C_1$  of candidate 1-fuzzy sequences with fuzzy time-intervals. By calculating the support of sequences in  $C_1$ , the set  $L_1$  of 1-fuzzy sequential patterns with fuzzy time-intervals is created.

**For  $k=2$ .** The set  $C_2$  of 2-fuzzy sequences with fuzzy time-intervals are generated by joining two sequences in  $L_1$  together and with LT in the form of  $L_1 \times LT \times L_1$ . For example, if  $L_1 = \{f_a, f_b\}$ ,  $LT = \{lt_1, lt_2, lt_3\}$ , then there are nine candidate 2-fuzzy sequences with fuzzy time-intervals such as  $\langle\langle f_a, lt_1, f_a \rangle\rangle$ ,  $\langle\langle f_a, lt_2, f_a \rangle\rangle$ ,  $\langle\langle f_a, lt_3, f_a \rangle\rangle$ ,  $\langle\langle f_a, lt_1, f_b \rangle\rangle$ ,  $\langle\langle f_a, lt_2, f_b \rangle\rangle$ ,  $\langle\langle f_a, lt_3, f_b \rangle\rangle$ ,  $\langle\langle f_b, lt_1, f_b \rangle\rangle$ ,  $\langle\langle f_b, lt_2, f_b \rangle\rangle$ ,  $\langle\langle f_b, lt_3, f_b \rangle\rangle$ .

**For  $k>2$ .** Let  $\langle\langle b_1, lt_1, b_2, lt_2, \dots, lt_{k-2}, b_{k-1} \rangle\rangle$  and  $\langle\langle b_2, lt_2, b_3, lt_3, \dots, lt_{k-1}, b_k \rangle\rangle$  be 2  $(k-1)$ -fuzzy sequential patterns with fuzzy time-intervals in  $L_{k-1}$ , then  $\alpha = \langle\langle b_1, lt_1, b_2, lt_2, b_3, lt_3, \dots, b_{k-1}, lt_{k-1}, b_k \rangle\rangle$  is a candidate  $k$ -fuzzy sequence with fuzzy time-intervals [6]. In the same way, all candidate  $k$ -fuzzy sequences with fuzzy time-intervals are generated and set  $C_k$  is created. The support of each candidate  $k$ -fuzzy sequence with fuzzy time-intervals in  $C_k$  is calculated by the formula (4) and the set  $L_k$  is created.

The result of the algorithm is a set of all  $k$ -fuzzy sequential patterns with fuzzy time-intervals  $L_k$  for  $k>1$ .

### 3.3. FSPFTIM Algorithm

Pseudo code of the algorithm is shown in Fig. 3.

In the algorithm, the operator “\*” joins the values into a fuzzy sequence with fuzzy time-intervals. For example,  $fe_1 * \text{Short} * fe_2$  is a presentation of the fuzzy sequence with fuzzy time-intervals  $\langle\langle fe_1, \text{Short}, fe_2 \rangle\rangle$  where  $\text{Short} \in LT$  and  $fe_1, fe_2$  are fuzzy sets of quantitative attributes.

At line 2, each fuzzy tuple, considered as a 1-fuzzy sequence with fuzzy time-intervals, is added into  $C_1$ . If the support of a sequence is greater than or equal to  $\min\_sup$  then the sequence will be added into  $L_1$  at line 3. The set of candidates 2-fuzzy sequences with fuzzy time-intervals,  $C_2 = L_1 \times LT \times L_1$ , is generated by the lines 4-12. The support of candidate 2-fuzzy sequences with fuzzy time-intervals in  $C_2$  are calculated (lines 13-15) and 2-fuzzy sequences with fuzzy time-intervals having the support greater than or equal to  $\min\_sup$ , are added into  $L_2$  (line 16). In the lines 17-24, candidate  $k$ -fuzzy sequences with fuzzy time-intervals are added into  $C_k$  (with  $k > 2$ ) and by calculating their support,  $k$ -fuzzy patterns with fuzzy time-intervals are included in  $L_k$ . The loop for generating candidate sequences is ended when  $L_k$  is empty. The result of the algorithm is the set of all  $k$ -fuzzy sequential patterns with fuzzy time-intervals (line 24) with  $k>1$ .



```

Procedure FSPMFTI (D, min_sup)
1.   Creating a fuzzy sequence database  $D'$  from the sequence database  $D$ 
2.    $C_1 \leftarrow \{fe \mid fe \text{ is an attribute of } D'\}$ 
3.    $L_1 \leftarrow \{\alpha \in C_1 \mid \text{Supp}(\alpha) \geq \text{min\_sup}\}$ 
4.    $C_2 \leftarrow \emptyset$ ;
5.   for each  $fe_1 \in L_1$ 
6.     for each  $fe_2 \in L_1$ 
7.       for each  $ltd \in LT$ 
8.          $\alpha \leftarrow fe_1 * ltd * fe_2$ ;
9.         add  $\alpha$  to  $C_2$ ;
10.      end for
11.    end for
12.  end for
13.  for each  $\alpha \in C_2$ 
14.    Computing the support of  $\alpha$  ( $\text{Supp}(\alpha)$ );
15.  end for
16.   $L_2 \leftarrow \{\alpha \in C_2 \mid \text{Supp}(\alpha) \geq \text{min\_sup}\}$ 
17.  for ( $k > 2; L_{k-1} \neq \emptyset; k++$ )
18.     $C_k \leftarrow \text{fuzzy\_apriori\_gen}(L_{k-1})$ ;
19.    for each  $c \in C_k$ 
20.      Computing the  $\text{Supp}(c)$ 
21.    end for
22.     $L_k \leftarrow \{\alpha \in C_k \mid \text{Supp}(\alpha) \geq \text{min\_sup}\}$ 
23.  end for
24.  return  $\cup L_k \ // k > 1$ 
Procedure fuzzy_apriori_gen( $L_{k-1}$ ) //generating the candidates of  $C_k$ 
25.   $C_k \leftarrow \emptyset$ ;
26.  for each  $a \in L_{k-1}$ 
27.    for each  $b \in L_{k-1}$ 
28.       $\alpha \leftarrow \emptyset$ ;
29.      for ( $i=2; i \leq k-2; i++$ )
30.        if ( $a_i \neq b_{i-1}$  or  $alt_i \neq blt_{i-1}$ )
31.          break;
32.        end if
33.         $\alpha \leftarrow c * a_i * alt_i$ ;
34.      end for
35.      if ( $i=k-1$  and  $a_{k-1} = b_{k-2}$ )
36.         $\alpha \leftarrow a_1 * alt_1 * c * a_{k-1} * blt_{k-2} * b_{k-1}$ ;
37.        add  $\alpha$  to  $C_k$ ;
38.      end if
39.    end for
40.  end for
41.  return  $C_k$ 

```

Fig. 3. The FSPFTIM Algorithm

The function generating candidate  $k$ -fuzzy sequences with fuzzy time-intervals in  $C_k$  from  $L_{k-1}$  is presented in the lines 25-41. The check of the conditions for combining two  $(k-1)$ -fuzzy sequential patterns with fuzzy time-intervals in  $L_{k-1}$  to creates a candidate  $k$ -fuzzy sequence with fuzzy time-intervals in lines 29-35. A candidate  $k$ -fuzzy sequence with fuzzy time-intervals  $\alpha$  is created and added into  $C_k$  in lines 36-37. At line 41, the set of candidate  $k$ -fuzzy sequences with fuzzy time-intervals,  $C_k$ , is created.

### The complexity of the FSPFTIM Algorithm

The parameters used to evaluate the computational complexity of the proposed algorithm are as follows:

- $N$  is the number of transaction sequences in  $D$ ;
- $M$  is the total number of attributes in  $D$ ;
- $l$  is the average length of transaction sequences in  $D$ ;
- $h$  is the number average of of fuzzy sets associated with each attribute in  $D$ ;
- $|LT|$  is the number of fuzzy sets of time-intervals  $LT$ ;
- $\text{min\_supp}$  is a user-defined threshold.

Based on the similar way as in [20], the computational complexity of the FSPFTIM algorithm is

$$O(N.l.h + M.N.l.h^2 + |L_1|^2 \cdot \binom{2}{l,h} \cdot N \cdot |LT|) + \sum_{k=3}^{l,h} |L_{k-1}|^2 \cdot \binom{k}{l,h} \cdot N \cdot |LT| \blacksquare$$

Details are given in the Appendix.

#### 3.4. An example of executing the FSPFTIM Algorithm

##### Input:

- $\text{min\_sup} = 0.11$ ;
- Quantitative sequence database  $D$  in Table 2;
- Fuzzy sets of attributes and their membership functions in the Example 2;
- Set of fuzzy sets of time-intervals  $LT = \{\text{Short}, \text{Medium}, \text{Long}\}$  and their membership functions in the Example 3.

**Output:**  $k$ -fuzzy sequential patterns with fuzzy time-intervals,  $k > 1$ .

##### Execution:

- $D'$  is created as the Table 3;
  - $C_1 = \left\{ \begin{array}{l} f_{3,1}^a, f_{3,2}^a, f_{3,3}^a, f_{3,1}^b, f_{3,2}^b, f_{3,3}^b, f_{3,1}^c, f_{3,2}^c, f_{3,3}^c, \\ f_{3,1}^d, f_{3,1}^e, f_{3,3}^e, f_{3,1}^f, f_{3,3}^f, f_{3,1}^g, f_{3,1}^h, f_{3,3}^h, f_{3,1}^i \end{array} \right\}$ ;
  - $L_1 = \{f_{3,1}^a, f_{3,2}^a, f_{3,1}^b, f_{3,3}^b, f_{3,1}^d, f_{3,1}^e, f_{3,3}^e, f_{3,1}^i\}$ ;
  - $C_2$  (included 147 elements) =
- $$= \left\{ \begin{array}{l} \langle \langle f_{3,1}^a, \text{Short}, f_{3,1}^a \rangle \rangle, \langle \langle f_{3,1}^a, \text{Medium}, f_{3,1}^a \rangle \rangle, \langle \langle f_{3,1}^a, \text{Long}, f_{3,1}^a \rangle \rangle, \\ \langle \langle f_{3,1}^a, \text{Short}, f_{3,2}^a \rangle \rangle, \langle \langle f_{3,1}^a, \text{Medium}, f_{3,2}^a \rangle \rangle, \langle \langle f_{3,1}^a, \text{Long}, f_{3,2}^a \rangle \rangle, \\ \dots, \\ \langle \langle f_{3,1}^i, \text{Short}, f_{3,1}^i \rangle \rangle, \langle \langle f_{3,1}^i, \text{Medium}, f_{3,1}^i \rangle \rangle, \langle \langle f_{3,1}^i, \text{Long}, f_{3,1}^i \rangle \rangle \end{array} \right\};$$

$$\begin{aligned}
& \bullet L_2 = \left\{ \left\langle \left\langle f_{3,1}^a, \text{Short}, f_{3,1}^b \right\rangle \right\rangle, \left\langle \left\langle f_{3,1}^a, \text{Long}, f_{3,3}^e \right\rangle \right\rangle, \left\langle \left\langle f_{3,1}^b, \text{Long}, f_{3,3}^e \right\rangle \right\rangle, \right. \\
& \quad \left. \left\langle \left\langle f_{3,1}^b, \text{Short}, f_{3,1}^d \right\rangle \right\rangle, \left\langle \left\langle f_{3,1}^b, \text{Long}, f_{3,1}^e \right\rangle \right\rangle, \left\langle \left\langle f_{3,1}^d, \text{Medium}, f_{3,3}^e \right\rangle \right\rangle \right\}. \\
& \bullet C_3 = \left\{ \left\langle \left\langle f_{3,1}^a, \text{Short}, f_{3,1}^b, \text{Long}, f_{3,1}^e \right\rangle \right\rangle, \left\langle \left\langle f_{3,1}^a, \text{Short}, f_{3,1}^b, \text{Short}, f_{3,1}^d \right\rangle \right\rangle, \right. \\
& \quad \left. \left\langle \left\langle f_{3,1}^a, \text{Short}, f_{3,1}^b, \text{Long}, f_{3,3}^e \right\rangle \right\rangle, \left\langle \left\langle f_{3,1}^a, \text{Short}, f_{3,1}^d, \text{Medium}, f_{3,3}^e \right\rangle \right\rangle \right\}; \\
& \bullet L_3 = \left\{ \left\langle \left\langle f_{3,1}^a, \text{Short}, f_{3,1}^d, \text{Medium}, f_{3,3}^e \right\rangle \right\rangle \right\}; \\
& \bullet C_4 = \emptyset; \\
& \bullet L_4 = \emptyset. \\
& \bullet \text{Output} = \left\{ \left\langle \left\langle f_{3,1}^a, \text{Short}, f_{3,1}^b \right\rangle \right\rangle, \left\langle \left\langle f_{3,1}^a, \text{Long}, f_{3,3}^e \right\rangle \right\rangle, \right. \\
& \quad \left\langle \left\langle f_{3,1}^b, \text{Long}, f_{3,3}^e \right\rangle \right\rangle, \left\langle \left\langle f_{3,1}^b, \text{Short}, f_{3,1}^d \right\rangle \right\rangle, \\
& \quad \left\langle \left\langle f_{3,1}^b, \text{Long}, f_{3,1}^e \right\rangle \right\rangle, \left\langle \left\langle f_{3,1}^d, \text{Medium}, f_{3,3}^e \right\rangle \right\rangle, \\
& \quad \left. \left\langle \left\langle f_{3,1}^a, \text{Short}, f_{3,1}^d, \text{Medium}, f_{3,3}^e \right\rangle \right\rangle \right\}.
\end{aligned}$$

## 4. Experimental results

The algorithm is implemented in the C# programming language and run on Chip Intel Core i5 2.5 GHz, RAM 4 GB, Windows 7 OS.

### 4.1. Datasets

Two datasets used for experiment of the proposed algorithm are S100I1000T3D341K and Online Retail\_France. These datasets are shown in Table 4. The dataset S100I1000T3D341K is generated according to the theory of *R. Agrawal* and *R. Srikant* [19]. Input parameters for generating data include the number of transactions  $D$ , the average of transactions  $T$ , the number of attributes  $I$ , and the number of transaction sequences  $S$ . The length of the transactions is generated by the Poisson distribution in which the parameter is the average of transactions ( $T$ ). The quantitative values and the time are random integers from 1 up to 1000. The Online Retail\_France dataset is gathered from Online Retail [24] with the Country="France" from 01/12/2010 to 09/12/2011. The dataset includes information as follows:

- Customer ID;
- Invoice Date;
- Stock Code;
- Quantity.

Table 4. Datasets

Dataset	Number of attributes ( $I$ )	Number of transactions ( $D$ )	Number of transaction sequences ( $S$ )	Average of transactions ( $T$ )	Average of transaction sequences
S100I1000T3D341K	1000	341	100	29.3	3.41
Online Retail_France	1523	365	87	22.8	4.20

The test dataset is transformed into a quantitative sequence database fitted for the input of the algorithm as follows:

- CusId: Customer ID;
- Time: is an integer number starting at 1, which is the number of days of Invoice Date differ from the first day (01/12/2010);
- StockCode: Descriptions of stock codes bought;
- Quantity: Quantitative value of each StockCode.

Each StockCode is partitioned into fuzzy sets and their membership functions are defined according to the Formulae (1) in Example 2 with  $K=3$  and  $A_{3,1}^{x_m} = x_{m\_Small}$ ,  $A_{3,2}^{x_m} = x_{m\_Average}$ ,  $A_{3,3}^{x_m} = x_{m\_Large}$ . The  $ma$ ,  $mi$  are the maximum, minimum of the quantitative values of  $x_m$  StockCode. Quantitative values of the StockCode are used to calculate fuzzy values.

Fuzzy sets of the time-intervals and their membership functions are defined as in formula (1) in the Example 2, too. The  $ma$ ,  $mi$  are the first time and the last time in order and  $A_{3,1}^t = \text{Short}$ ,  $A_{3,2}^t = \text{Medium}$ ,  $A_{3,3}^t = \text{Long}$  where  $t$  is time-interval. The time distances are used to calculate fuzzy time-intervals.

## 4.2. Results

4.2.1. Relationships between the number of fuzzy sequential patterns with fuzzy time-intervals and  $\min\_sup$ , and between the execution time and  $\min\_sup$  according to the different numbers of partitions of quantitative attributes.

In this case, the number of fuzzy sets of time intervals is fixed. Namely, the time intervals are partitioned into 3 fuzzy sets ( $K_i=3$ ). The relationships between the number of fuzzy sequential patterns with fuzzy time-intervals and  $\min\_sup$ , and between the execution time and  $\min\_sup$  according to different numbers of partitions of quantitative attributes for the S100I1000T3D341K dataset are shown in the Figs 4a and b, respectively, and for the Online Retail\_France datasets are also shown in the Figs 5a and b.

From these Figs, we can see that when the number of partitions of time intervals is fixed, the number of fuzzy sequential patterns with fuzzy time-intervals as well as the execution time will be changed in inverse direction with changes of  $\min\_sup$ . This is true for both the S100I1000T3D341K and Online Retail\_France datasets, and does not depend on the number of partitions of quantitative attributes ( $K$ ). Also, as  $K$  increases, the number of fuzzy sequential patterns with fuzzy time-intervals decreases, but with  $\min\_sup$  greater than a certain threshold, when  $K$  increases then the execution time increases. In addition, the difference between the numbers of

sequence patterns, and between the times of execution (depending on the number of partitions of the quantitative attributes) will be narrowed very rapidly according to the increase of `min_sup`.

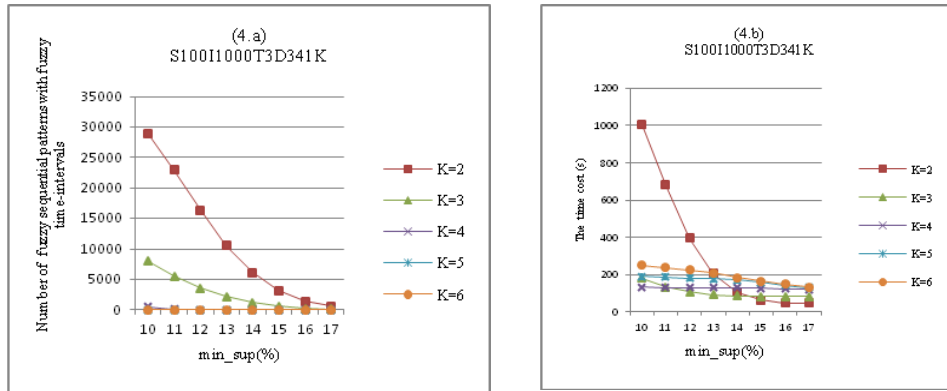


Fig. 4. Relationship between the number of fuzzy sequential patterns with fuzzy time-intervals and `min_sup` (a), Relationship between the execution time and `min_sup` (b) according to the different numbers of partitions of quantitative attributes ( $K$ ) for the S100I1000T3D341K dataset

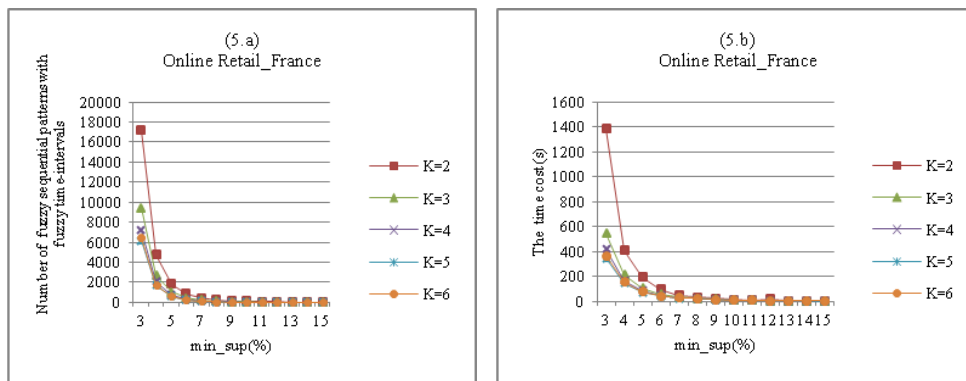


Fig. 5. Relationship between the number of fuzzy sequential patterns with fuzzy time-intervals and `min_sup` (a), Relationship between the execution time and `min_sup` (b) according to different numbers of partitions of quantitative attributes ( $K$ ) for the Online Retail\_France dataset

#### 4.2.2. Relationships between the number of fuzzy sequential patterns with fuzzy time-intervals and `min_sup`, and between the execution time and `min_sup` according to the different numbers of partitions of time-intervals

In this case, the number of partitions of quantitative attributes is 3 ( $K=3$ ). Figs 6a and 7a, show the relationship between the number of fuzzy sequential patterns with fuzzy time-intervals and `min_sup` according to the different numbers of partitions of the time-intervals ( $K_t$ ) for the S100I1000T3D341K and Online Retail\_France datasets, respectively. Similarly, Figs. 6b and 7b show the relationship between the time of execution and `min_sup` for the two test datasets mentioned above.

The relationship between the number of fuzzy sequential patterns with fuzzy time-intervals and  $\min\_sup$ , and between the time of execution and  $\min\_sup$  according to changes of the number of partitions of time-intervals in the context the number of partitions of the quantitative attributes is fixed they are very similar to the relationship mentioned in Section 4.2.1. For example, Figs 6b and 7b also show that when the number of partitions of quantitative attributes is fixed, if the number of partitions of time-intervals increases then the execution time also increases.

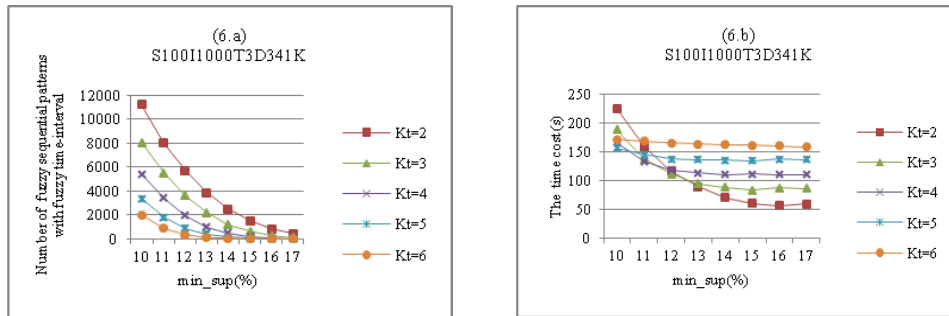


Fig. 6. Relationship between the number of fuzzy sequential patterns with fuzzy time-intervals and  $\min\_sup$  (a), between the execution time and  $\min\_sup$  (b) according to the different numbers of partitions of time-intervals ( $K_t$ ) for the S100I1000T3D341K dataset

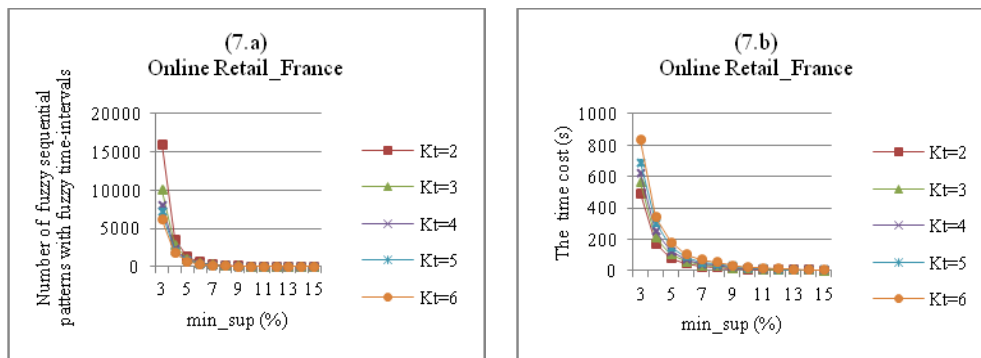


Fig. 7. Relationship between the number of fuzzy sequential patterns with fuzzy time-intervals and  $\min\_sup$  (a), between the execution time and  $\min\_sup$  (b) according to the different numbers of partitions of time-intervals ( $K_t$ ) for the Online Retail\_France dataset

#### 4.3. Meaning of fuzzy sequential patterns with fuzzy time-intervals

Table 5 shows all fuzzy sequence patterns with fuzzy time-intervals that were found out using the FSPFTIM algorithm in the Online Retail\_France dataset when the quantitative attributes are partitioned into three linguistic terms {Small, Average, Large}, the time interval into three linguistic terms {Short, Medium, Long}, the user – defined threshold  $\min\_sup$  is 0.14.

Some of these patterns can be explained as follows:

$\langle\langle\text{POSTAGE}_{\text{Small}}, \text{Short}, \text{POSTAGE}_{\text{Small}}\rangle\rangle$  means “**If** a customer buys a Small number of POSTAGE products, **then** he will buy a Small number of these products an interval of **Short**”;

$\langle\langle\text{POSTAGE}_{\text{Small}}, \text{Medium}, \text{RABBIT-NIGHT-LIGHT}_{\text{Small}}\rangle\rangle$  means “**If** a customer buys a Small number of POSTAGE products, **then** he will buy a Small number of RABBIT-NIGHT-LIGHT an interval of **Medium**”;

$\langle\langle\text{POSTAGE}_{\text{Small}}, \text{Medium}, \text{POSTAGE}_{\text{Small}}, \text{Short}, \text{POSTAGE}_{\text{Small}}\rangle\rangle$  means “**If** a customer buys a Small number of POSTAGE products and later a Small number of POSTAGE products an interval of **Medium**, **then** he will buy a Small number of POSTAGE products an interval of **Short**”.

Table 5. The fuzzy sequential patterns with fuzzy time-intervals in Online Retail\_France dataset

Fuzzy sequential patterns with fuzzy time-intervals	Support
$\langle\langle\text{POSTAGE}_{\text{Small}}, \text{Short}, \text{POSTAGE}_{\text{Small}}\rangle\rangle$	0.236
$\langle\langle\text{POSTAGE}_{\text{Small}}, \text{Medium}, \text{POSTAGE}_{\text{Small}}\rangle\rangle$	0.257
$\langle\langle\text{POSTAGE}_{\text{Small}}, \text{Short}, \text{RABBIT-NIGHT-LIGHT}_{\text{Small}}\rangle\rangle$	0.151
$\langle\langle\text{POSTAGE}_{\text{Small}}, \text{Medium}, \text{RABBIT-NIGHT-LIGHT}_{\text{Small}}\rangle\rangle$	0.161
$\langle\langle\text{POSTAGE}_{\text{Average}}, \text{Short}, \text{POSTAGE}_{\text{Small}}\rangle\rangle$	0.141
$\langle\langle\text{POSTAGE}_{\text{Small}}, \text{Short}, \text{POSTAGE}_{\text{Small}}, \text{Short}, \text{POSTAGE}_{\text{Small}}\rangle\rangle$	0.146
$\langle\langle\text{POSTAGE}_{\text{Small}}, \text{Medium}, \text{POSTAGE}_{\text{Small}}, \text{Short}, \text{POSTAGE}_{\text{Small}}\rangle\rangle$	0.144

## 5. Conclusions and future work

The paper proposes the FSPFTIM algorithm to mine fuzzy sequential patterns with fuzzy time-intervals in quantitative sequence databases. The proposed algorithm is presented in both idea and pseudo code. An example of application of the algorithm is also illustrated in this paper. The experimental results show the influence of some input parameters to the number of fuzzy sequential patterns with fuzzy time-intervals and the execution time. In this paper, the FSPFTIM algorithm is developed based on the Apriori algorithm with the approach of search in breadth first. In the future, we will continue to work on developing more efficient algorithms to mine fuzzy sequence patterns with fuzzy time – intervals. These algorithms will be developed based on the approach of search in depth first, such as based on the projected database organization [10] or the FPTree algorithm [11].

## References

1. Agrawal, R., R. Srikant. Mining Sequential Patterns. – In: Proc. of 11th International Conference on Data Engineering, 1995.
2. Cao, H., et al. A Fuzzy Sequential Pattern Mining Algorithm Based on Independent Pruning Strategy for Parameters Optimization of Ball Mill Pulverizing System. – Information Technology and Control, Vol. **43**, 2014, No 3, pp. 303-314.
3. Chang, C.-I., H.-E. Chueh, N. P. Lin. Sequential Patterns Mining with Fuzzy Time-Intervals. – In: 6th International Conference on FSKD’09, Vol. **3**, 2009.
4. Chang, C.-I., H.-E. Chueh, d Y.-C. Luo. An Integrated Sequential Patterns Mining with Fuzzy Time-Intervals – In: International Conference on Systems and Informatics (ICSAI’12), 2012.

5. Chen, R.-S. et al. Discovery of Fuzzy Sequential Patterns for Fuzzy Partitions in Quantitative Attributes. – In: ACS/IEEE International Conference on Computer Systems and Applications, 2001, pp. 144-150.
6. Chen, Y.-L., T. C.-K. Huang. Discovering Fuzzy Time-Interval Sequential Patterns in Sequence Databases. – IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics), Vol. **35**, 2005, No 5, pp. 959-972.
7. Chen, Y.-L., M.-C. Chiang, M.-T. Ko. Discovering Time-Interval Sequential Patterns in Sequence Databases. – Expert Systems with Applications, Vol. **25**, 2003, No 3, pp. 343-354.
8. Fowkes, J., C. Sutton. A Subsequence Interleaving Model for Sequential Pattern Mining. – arXiv preprint arXiv:1602.05012, 2016.
9. Garofalakis, M. N., R. Rastogi, K. Shim. SPIRIT: Sequential Pattern Mining with Regular Expression Constraints. – VLDB, Vol. **99**, 1999.
10. Han, J., et al. Prefixspan: Mining Sequential Patterns Efficiently by Prefix-Projected Pattern Growth. – In: Proc. of 17th International Conference on Data Engineering, 2001.
11. Han, J., J. Pei, Y. Yin. Mining Frequent Patterns without Candidate Generation. – In: ACM Sigmod Record, Vol. **29**, 2000.
12. Hong, T.-P., C.-S. Kuo, S.-C. Chi. Mining Fuzzy Sequential Patterns from Quantitative Data. – In: Proc. of Conference IEEE SMC'99, Vol. **3**, 1999.
13. Hu, Y.-C., et al. A Fuzzy Data Mining Algorithm for Finding Sequential Patterns. – International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, Vol. **11**, 2003, No 2, pp. 173-193.
14. Hu, Y.-C., G.-H. Tzeng, C.-M. Chen. Deriving Two-Stage Learning Sequences from Knowledge in Fuzzy Sequential Pattern Mining. – Information Sciences, Vol. **159**, 2004, No 1, pp. 69-86.
15. Huang, T., et al. Extracting Various Types of Informative Web Content via Fuzzy Sequential Pattern Mining. – Asia-Pacific Web (APWeb) and Web-Age Information Management (WAIM) Joint Conference on Web and Big Data. Springer, Cham, 2017, pp. 230-238.
16. Huang, T. C.-K. Discovery of Fuzzy Quantitative Sequential Patterns with Multiple Minimum Supports and Adjustable Membership Functions. – Information Sciences, Vol. **222**, 2013, pp. 126-146.
17. Kieu, T., et al. Mining Top-K Co-Occurrence Items With Sequential Pattern. – Expert Systems with Applications, Vol. **85**, 2017, pp. 123-133.
18. Kuo, R. J., C. M. Chao, C. Y. Liu. Integration of k-Means Algorithm and AprioriSome Algorithm for Fuzzy Sequential Pattern Mining. – Applied Soft Computing, Vol. **9**, 2009, No 1, pp. 85-93.
19. Agrawal, R., R. Srikant. Fast Algorithms for Mining Association Rules. – In: J. Bocca, M. Jarke, C. Zaniolo, Eds. Proc. of 20th International Conference on Very Large Data Bases (VLDB'94). Santiago, Morgan Kaufmann, 1994, pp. 487-499.
20. Tan, P.-N., M. Steinbach, V. Kumar. Association Analysis: Basic Concepts and Algorithms. – Introduction to Data Mining, 2005, pp. 327-414.
21. Phuong, T. D., D. V. Thanh, N. D. Dung. An Effective Algorithm for Association Rules Mining from Temporal Quantitative Databases. – Indian Journal of Science and Technology, Vol. **9**, 2016, No 17.
22. Wright, A. P., et al. The Use of Sequential Pattern Mining to Predict Next Prescribed Medications. – Journal of Biomedical Informatics, Vol. **53**, 2015, pp. 73-80.
23. Yu, C.-C., Y.-L. Chen. Mining Sequential Patterns from Multidimensional Sequence Data. – IEEE Transactions on Knowledge and Data Engineering, Vol. **17**, 2005, No 1, pp. 136-140.
24. UCI-Machine Learning Repository. June 2017.  
<http://archive.ics.uci.edu/ml/datasets.html>



## Appendix

The computational complexity of the FSPFTIM algorithm is estimated as follows:

- *Creating the fuzzy sequence database  $D'$* : Each fuzzy transaction sequence in  $D'$  has  $l.h$  fuzzy attributes, so the computational complexity of converting the quantitative sequence database  $D$  to the fuzzy sequence database  $D'$  is  $O(N.l.h)$ .

- *Generating  $C_1$  and calculating the support of sequences in  $C_1$* :  $C_1$  includes the  $M.h$  candidate 1-fuzzy sequences with fuzzy time-intervals (in this case, they are the fuzzy attributes in  $D'$ ). To calculate the support for each candidate sequence in  $C_1$ , it is necessary to look at all fuzzy transaction sequences in  $D'$  as well as all fuzzy attributes in these transaction sequences. So, the computational complexity is  $O(M.h)$ .  $O(N.l.h) = O(M.N.l.h^2)$ .

- *Generating  $C_2$  and calculating the support of sequences in  $C_2$  (or creating  $L_2$  for short)*:

- $C_2$  – the set of candidate 2-fuzzy sequences with fuzzy time-intervals is created by joining two sequences in  $L_1$  together and with  $LT$ , namely,  $C_2=L_1 \times LT \times L_1$ . The number of candidate sequences in  $C_2$  can be  $|L_1|^2 \cdot |LT|$ . To calculate the support of each candidate sequence in  $C_2$ , we have to look at all fuzzy transaction sequences in  $D'$  as well as all 2-fuzzy sequences in these transaction sequences, i.e., the computational complexity of calculating the support of each candidate sequence in  $C_2$  is  $O(N \cdot \binom{2}{l,h})$ . Hence the computational complexity of creating  $L_2$  is  $O(|L_1|^2 \cdot |LT|) \cdot O(N \cdot \binom{2}{l,h}) = O(|L_1|^2 \cdot \binom{2}{l,h} \cdot N \cdot |LT|)$ .

- *Generating  $C_k$  and calculating the support of sequences in  $C_k$  with  $k > 2$  (or creating  $L_k$  for short)*:  $C_k$  is generated by merging two  $(k-1)$ -fuzzy sequential patterns with fuzzy time-intervals in  $L_{k-1}$ , so the number of sequences in  $C_k$  can be  $|L_{k-1}|^2 \cdot |LT|$ . Check of merging conditions requires a maximum of  $2k-5$  comparisons, including  $k-2$  comparisons for fuzzy attributes and  $k-3$  comparisons for fuzzy time-intervals. Hence, the cost of generating  $C_k$  is  $O((2k-5) \cdot |L_{k-1}|^2 \cdot |LT|)$ . And similar to the above, the computational complexity of each sequence in  $C_k$  is  $O(N \cdot \binom{k}{l,h})$ . Therefore the computational complexity of creating  $L_k$  ( $k > 2$ ) is

$$O((2k-5) \cdot |L_{k-1}|^2 \cdot |LT|) \cdot O(N \cdot \binom{k}{l,h}) = O((2k-5) \cdot |L_{k-1}|^2 \cdot \binom{k}{l,h} \cdot N \cdot |LT|).$$

Finally, the computational complexity of the proposed algorithm is  $O(N.l.h + M.N.l.h^2 + |L_1|^2 \cdot \binom{2}{l,h} \cdot N \cdot |LT|) + \sum_{k=3}^{l,h} (2k-5) \cdot |L_{k-1}|^2 \cdot \binom{k}{l,h} \cdot N \cdot |LT|) = O(N.l.h + M.N.l.h^2 + |L_1|^2 \cdot \binom{2}{l,h} \cdot N \cdot |LT|) + \sum_{k=3}^{l,h} |L_{k-1}|^2 \cdot \binom{k}{l,h} \cdot N \cdot |LT|) \blacksquare$

**Note:**  $L_k$  ( $k \geq 1$ ) depends so much on  $\min\_sup$ .

*Received 20.08.2017; Second Version 20.05.2018; Accepted 28.05.2018*