

## Soft set based association rule mining



Feng Feng<sup>a,b,\*</sup>, Junghoo Cho<sup>b</sup>, Witold Pedrycz<sup>c,d</sup>, Hamido Fujita<sup>e</sup>, Tutut Herawan<sup>f</sup>

<sup>a</sup> Department of Applied Mathematics, School of Science, Xi'an University of Posts and Telecommunications, Xi'an, 710121, China

<sup>b</sup> Department of Computer Science, University of California, Los Angeles, CA, 90095, USA

<sup>c</sup> Department of Electrical and Computer Engineering, University of Alberta, AB T6R 2G7, Canada

<sup>d</sup> Systems Research Institute, Polish Academy of Sciences, Warsaw, Poland

<sup>e</sup> Faculty of Software and Information Science, Iwate Prefectural University, Iwate, 020-0193, Japan

<sup>f</sup> Faculty of Computer Science and Information Technology, University of Malaya, 50603 Pantai Valley, Kuala Lumpur, Malaysia

### ARTICLE INFO

#### Article history:

Received 13 March 2016

Revised 12 July 2016

Accepted 20 August 2016

Available online 26 August 2016

#### Keywords:

Soft set

Association rule

Maximal association rule

Data mining

Information system

Transactional dataset

### ABSTRACT

Association rules, one of the most useful constructs in data mining, can be exerted to capture interesting dependencies between variables in large datasets. Herawan and Deris initiated the investigation of mining association rules from transactional datasets using soft set theory. Unfortunately, some existing concepts in the literature were unable to realize properly Herawan and Deris's initial idea. This paper aims to offer further detailed insights into soft set based association rule mining. With regard to regular association rule mining using soft sets, we refine several existing concepts to improve the generality and clarity of former definitions. Regarding maximal association rule mining based on soft sets, we point out the drawbacks of some existing definitions and offer some way to rectify the problem. A number of new notions, such as transactional data soft sets, parameter-taxonomic soft sets, parameter cosets, realizations and  $M$ -realizations of parameter sets are proposed to facilitate soft set based association rule mining. Several algorithms are designed to find  $M$ -realizations of parameter sets or extract  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules in parameter-taxonomic soft sets. We also present an example to illustrate potential applications of our method in clinical diagnosis. Moreover, two case studies are conducted to highlight the essentials of soft set based association rule mining approach.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

With the development of computer science and information technology, very huge amounts of data have been collected and stored in the memory of computers. By analyzing the collected data, it is possible to arrive at a sophisticated understanding of the corresponding facts, behaviors, and natural phenomena in the real world. In response to this need, data mining becomes an increasingly important field which aims to extract useful information from the collected data and transform it into an understandable form of knowledge for further use. Data mining tasks involve a number of different methods at the intersection of artificial intelligence, database systems, mathematics, machine learning and statistics.

One of the most important issues in data mining is to discover association rules, which was initially introduced by Agrawal et al.

[1]. In Agrawal et al.'s seminal work [2], three novel algorithms namely *Apriori*, *Apriori<sub>Tid</sub>* and *AprioriHybrid* were presented to extract association rules from large databases. Since then, association rules have received considerable attention from researchers and practitioners around the world. As a powerful data mining tool, association rule mining has been successfully applied to various domains such as bioinformatics, e-business, epidemiology, finance, health science, marketing and so forth. In particular, association rules are extremely useful for analyzing transactional data. A typical scenario of direct application is to analyze consumers' purchase records, commonly known as the "market basket data". Using association rule mining, one can discover some unexpected patterns of purchase behavior which may help to design effective marketing strategies accordingly.

Molodtsov's soft set theory [22] was proposed in 1999 as a general mathematical tool to deal with uncertainty. The rationale behind soft sets is founded on the idea of parameterization, which suggests that complicated objects should be perceived from various points of view. Each aspect provides an approximate description of the whole entity with high complexity. Without the limitation caused by inadequacy of parameterization tools, this theory comes

\* Corresponding author.

E-mail addresses: [fengnix@hotmail.com](mailto:fengnix@hotmail.com) (F. Feng), [cho@cs.ucla.edu](mailto:cho@cs.ucla.edu) (J. Cho), [wpedrycz@ualberta.ca](mailto:wpedrycz@ualberta.ca) (W. Pedrycz), [HFujita-799@acm.org](mailto:HFujita-799@acm.org) (H. Fujita), [tutut@um.edu.my](mailto:tutut@um.edu.my) (T. Herawan).

with an ability to represent and manipulate data in a convenient and meaningful way [8]. Maji et al. introduced several algebraic operations in soft set theory and examined their basic properties [19]. Ali et al. [7] proposed several new operations in soft set theory to further consolidate the algebraic basis of soft sets. Based on these new operations, Qin and Hong [24] introduced some congruence relations on soft sets and discussed certain lattice structures. Xiao et al. [26] proposed exclusive disjunctive soft sets and some related operations. Gong et al. [15] initiated the concept of bijective soft sets and explored decision rules in bijective soft decision systems [14]. Xu et al. [27] incorporated statistical logistic regression into soft set decision theory and proposed a novel parameter reduction method to select financial ratios for business failure prediction. Qin et al. [25] presented a soft set model of equivalence classes of information systems and applied it to the selection of clustering attributes for categorical datasets. Mamat et al. [21] designed a soft set based algorithm for clustering attribute selection. In addition, soft sets and their extensions have been successfully applied to various algebraic structures [3–5,17,18] and decision-making problems [6,10,20].

Although association rules are very useful for discovering associations hidden in the collected data, sometimes they might fail to identify other, no less interesting data connections. For this reason, Amir et al. [9] proposed maximal association rules to extract associations that are frequently lost when using regular association rules. In maximal association rule mining, the item domain is partitioned into pairwise disjoint categories. Intuitively, a maximal association rule  $X \xrightarrow{M} Y$  says that whenever the itemset  $X$  is *alone* in a transaction  $t$  (i.e.,  $X$  is contained in  $t$  and meanwhile there is no other item in  $t$  from the same category of  $X$ ), then the itemset  $Y$  also appears in the same transaction, with some confidence. As pointed out by Amir et al. [9], maximal association rules are not designed to replace regular association rules, but rather to complement them.

Inspired by interesting relationships among transactional datasets, Boolean-valued information systems and soft sets, Herawan and Deris [16] came up with the innovative idea about mining association rules from transactional datasets using soft set theory. Their pioneering work established a new research direction of applying soft sets to data mining. Unfortunately, some basic notions (particularly those developed for mining maximal association rules) in [16] failed to fulfill Herawan and Deris's initial idea, which hindered the further research in this new direction. Moreover, Herawan and Deris's soft set based association rule mining method relies on the prerequisite that the soft set must be transformed from a transactional dataset, which makes it difficult to be understood in terms of soft set theory and also restricts its potential applications to more diverse domains.

In the present study, we revisit Herawan and Deris's initial idea and manage to address the issues mentioned above. We first introduce some new notions for mining association rules in transactional datasets, such as realizations of itemsets, realizations of association rules,  $M$ -realizations of itemsets and  $M$ -realizations of maximal association rules. To improve the generality and clarity of the theory for mining regular association rules using soft sets, we refine some fundamental concepts and compare them with the former ones. In addition, we also define notions like transactional data soft sets,  $\sigma$ -frequent sets of parameters,  $\sigma$ -strong association rules and  $\gamma$ -reliable association rules, which are useful for mining regular association rules using soft sets. With regard to maximal association rule mining based on soft sets, we point out some difficulties suffered by existing definitions. To rectify some improper definitions in [16], we introduce a new concept called parameter-taxonomic soft sets, based on which corrective notions are presented to fulfill Herawan and Deris's initial idea

amply. We also propose some new notions, such as  $\sigma$ - $M$ -frequent sets of parameters,  $\sigma$ - $M$ -strong maximal association rules and  $\gamma$ - $M$ -reliable maximal association rules, which are designed for mining maximal association rules using soft sets. Algorithms for finding  $M$ -realizations of parameter sets or mining  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules from a given parameter-taxonomic soft set are presented. We propose an illuminating example to show potential applications of soft set based association rule mining in clinical diagnosis. In addition, we conduct two case studies to highlight some essential points of the newly proposed concepts and algorithms.

This paper is organized as follows. To facilitate our discussion, Section 2 recalls some basic concepts regarding information systems and soft sets. Section 3 gives a brief introduction to some essential points concerning transactional datasets and association rules. Section 4 is devoted to improving some existing definitions for mining association rules based on soft sets. In Section 5, we rectify some improperly defined existing notions and design algorithms for mining maximal association rules with parameter-taxonomic soft sets. Section 6 covers an example to illustrate potential applications of soft set based association rule mining in clinical diagnosis. In Section 7, two case studies are carried out to highlight the essentials of association rule mining based on soft sets. Finally, conclusions and some possible directions for future work are given in Section 8.

## 2. Preliminaries

In this section, we recall some basic notions regarding information systems and soft sets.

**Definition 2.1** [23]. An information system is a pair  $\mathcal{I} = (U, A)$  of non-empty finite sets  $U$  and  $A$ , where  $U$  is a set of objects and  $A$  is a set of attributes; each attribute  $a \in A$  is a function  $a: U \rightarrow V_a$  and  $V_a$  is called the *domain* of the attribute  $a$ .

If  $V_a = \{0, 1\}$  for all  $a \in A$ , then  $\mathcal{I} = (U, A)$  is called a *Boolean-valued information system*. If the set  $A$  of attributes is partitioned into two disjoint subset of attributes, namely *condition* and *decision attributes*, then the information system  $\mathcal{I} = (U, C, D)$  is called a *decision system*, where  $C$  and  $D$  are sets of condition and decision attributes, respectively.

Molodtsov [22] initiated the theory of soft sets, which provides a general framework for uncertainty modelling from a parameterization point of view. Let  $U$  be an initial universe of objects and  $E_U$  (or simply  $E$ ) be the set of all parameters associated with objects in  $U$ , called a *parameter space*. In most cases parameters are considered to be attributes, characteristics or properties of objects in  $U$ . The pair  $(U, E)$  is also known as a *soft universe*. We denote the power sets of  $U$  by  $\mathcal{P}(U)$ .

**Definition 2.2** [22]. A pair  $\mathfrak{S} = (F, A)$  is called a *soft set* over  $U$ , where  $A \subseteq E$  and  $F: A \rightarrow \mathcal{P}(U)$  is a set-valued mapping, called the *approximate function* of the soft set  $\mathfrak{S}$ .

It is clear that a soft set  $\mathfrak{S} = (F, A)$  over  $U$  can be seen as a parameterized family of subsets of  $U$ . For any parameter  $e \in A$ , the subset  $F(e) \subseteq U$  could be interpreted as the set of *e-approximate elements*. Note that  $F(e)$  may be arbitrary: some of them may be empty, and some may have nonempty intersections [22]. The absence of any restrictions on the approximate description in soft set theory also facilitates its applications to problems arising from various domains. As suggested by Molodtsov [22], one can use any suitable parametrization—with the help of words and sentences, real numbers, functions, mappings, etc. In what follows, the collection of all soft sets over  $U$  with parameter sets contained in  $E$  is denoted by  $\mathcal{S}^E(U)$ .

**Definition 2.3** [13]. A soft set  $S = (F, A)$  over  $U$  is said to be full if  $\bigcup_{a \in A} F(a) = U$ . A full soft set  $S = (F, A)$  over  $U$  is called covering soft set if  $F(a) \neq \emptyset$  for all  $a \in A$ .

**Definition 2.4** [13]. A soft set  $S = (F, A)$  over  $U$  is called a partition soft set if  $\{F(a) : a \in A\}$  forms a partition of  $U$ .

There exist some fundamental connections between information systems and soft sets. Note first that a soft set  $S = (F, A)$  over  $U$  gives rise to an information system  $\mathcal{S}_S = (U, A)$  in a natural way. In fact, for all  $a \in A$ , one can define a corresponding function  $a : U \rightarrow V_a = \{0, 1\}$  by

$$a(x) = \begin{cases} 1, & \text{if } x \in F(a), \\ 0, & \text{otherwise.} \end{cases}$$

This justifies the tabular (or matrix) representation of soft sets widely used in the literature.

On the other hand, soft sets also provide an efficient representation of information systems as follows.

**Proposition 2.5** [12]. Let  $\mathcal{S} = (U, A)$  be an information system and  $a \in A$ . Define a soft set  $S_a = (F_a, B_a)$  such that  $B_a = \{a\} \times V_a$  and

$$F_a(a, t) = \{x \in U : a(x) = t\},$$

for all  $(a, t) \in B_a$ . Then  $S_a = (F_a, B_a)$  is a partition soft set over  $U$ .

**Proposition 2.6** [12]. Let  $\mathcal{S} = (U, A)$  be an information system. Define a soft set  $S_{\mathcal{S}} = (F, B)$  such that  $B = \bigcup_{a \in A} \{a\} \times V_a$  and

$$F(a, v) = \{x \in U : a(x) = v\},$$

for all  $(a, v) \in B$ . Then  $S_{\mathcal{S}} = (F, B)$  is a covering soft set over  $U$ .

### 3. Transactional datasets and association rules

In this section, we briefly introduce some essential points concerning transactional datasets, association rules and maximal association rules. Most of these notions come from the pioneering work of Agrawal et al. [1] and Amir et al. [9], with some minor modifications. Nevertheless, it is worth noting that some new notions such as realizations of itemsets, realizations of association rules,  $M$ -realizations of itemsets and  $M$ -realizations of maximal association rules are also initiated to facilitate our discussion hereinafter.

To recognize consumers' behaviors, the association rules method has been developed particularly for the analysis of transactional datasets. Let  $I = \{i_1, \dots, i_{|I|}\}$  (where  $|\cdot|$  denotes the cardinality of a set) be the set of items, called the item domain. A transaction  $t$  is a nonempty set of items chosen from  $I$ . Each transaction can be identified by a transaction identifier (TID). The collection  $D = \{t_1, \dots, t_{|D|}\}$  consisting of all transactions under consideration is called a transactional dataset. Any given nonempty subset  $X$  of the item domain  $I$  is called an itemset. If the relation  $X \subseteq t$  holds, we say that the itemset  $X$  appears in the transaction  $t$ , or simply  $t$  supports  $X$ . An itemset with  $k$  items is called a  $k$ -itemset. For simplicity, we identify a 1-itemset  $\{i_s\}$  with the single item  $i_s$ .

**Definition 3.1.** Let  $D$  be a transactional dataset and  $X$  be an itemset. Then  $\Delta_D(X) = \{t \in D : X \subseteq t\}$  is called the realization of  $X$  in  $D$ .

The realization  $\Delta_D(X)$  is the set consisting of all the transactions which support the itemset  $X$  in the transactional dataset  $D$ . The cardinality of this set (i.e., the number of transactions supporting  $X$  in  $D$ ) is called the support of  $X$  in  $D$  and denoted by  $S_D(X)$ .

Given two disjoint nonempty itemsets  $X, Y \subseteq I$ , an association rule is a formal expression of the form  $X \Rightarrow Y$ . The itemsets  $X, Y$  are referred to as antecedent and consequent of the rule, respectively.

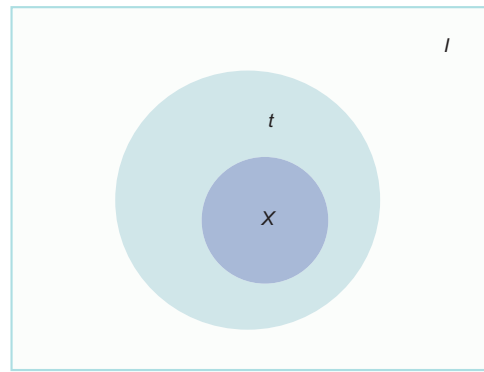


Fig. 1. Illustration of an itemset  $X$  supported by a transaction  $t$ .

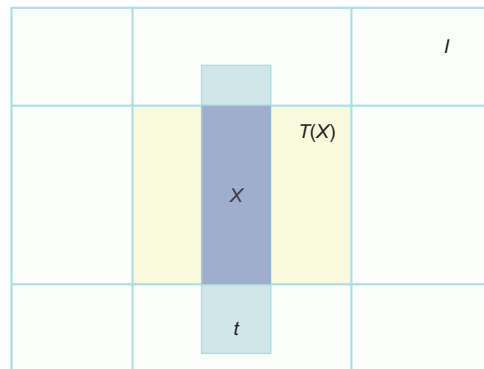


Fig. 2. Illustration of an itemset  $X$   $M$ -supported by a transaction  $t$ .

**Definition 3.2.** Let  $D$  be a transactional dataset and  $X \Rightarrow Y$  be an association rule. The realization of  $X \Rightarrow Y$  in  $D$  is defined as

$$\Delta_D(X \Rightarrow Y) = \{t \in D : (X \cup Y) \subseteq t\}.$$

Clearly,  $\Delta_D(X \Rightarrow Y) = \Delta_D(X \cup Y)$ . The cardinality of this set is called the support of  $X \Rightarrow Y$  and denoted by  $S_D(X \Rightarrow Y)$ . An association rule  $X \Rightarrow Y$  holds with some confidence  $C_D(X \Rightarrow Y)$ , which is defined as follows:

$$C_D(X \Rightarrow Y) = \frac{S_D(X \Rightarrow Y)}{S_D(X)} = \frac{S_D(X \cup Y)}{S_D(X)}. \tag{1}$$

In particular, we define  $C_D(X \Rightarrow Y) = 0$  if  $S_D(X) = 0$ .

Next, we discuss the recognition of maximal association rules in transactional datasets as initiated in Amir et al. [9]. To this end, we should first consider the partition of the item domain  $I$ . A partition  $T$  of  $I$  is called a taxonomy of  $I$ . Each block in  $T$  is called a category. For any given item  $i_k \in I$ , there exists a unique category containing  $i_k$ , which is denoted by  $T(i_k)$ . If an itemset  $X$  is from some category  $T_i \in T$  (i.e.,  $X \subseteq T_i$ ), then this unique category is denoted by  $T(X)$  in what follows.

**Definition 3.3** [9]. Let  $t$  be a transaction in a transactional dataset  $D$  and  $X$  be an itemset from a category. Then  $X$  is said to be alone in  $t$  if  $t \cap T(X) = X$ . In this case, we say that the transaction  $t$   $M$ -supports the itemset  $X$  in  $D$ .

It is worth noting that  $X$  is the largest subset of  $t$  contained in the category  $T(X)$  if it is  $M$ -supported by the transaction  $t$ . In addition, the difference between the notions of supporting and  $M$ -supporting is illustrated by Figs. 1 and 2.

**Definition 3.4.** Let  $D$  be a transactional dataset and  $X$  be an itemset from a category. Then  $\Delta_D^M(X) = \{t \in D : t \cap T(X) = X\}$  is called the  $M$ -realization of  $X$  in  $D$ .

The  $M$ -realization  $\Delta_D^M(X)$  comprises all the transactions that  $M$ -supports  $X$  in  $D$ . The cardinality of  $\Delta_D^M(X)$  is called the  $M$ -support of  $X$  in  $D$  and denoted by  $S_D^M(X)$ .

**Definition 3.5.** Let  $D$  be a transactional dataset and  $X, Y$  be nonempty itemsets from two distinct categories. Then the formal expression  $X \xrightarrow{M} Y$  is called a *maximal association rule*. The  $M$ -realization of  $X \xrightarrow{M} Y$  is defined as

$$\Delta_D^M(X \xrightarrow{M} Y) = \{t \in D : (t \cap T(X) = X) \wedge (Y \subseteq t)\}.$$

The cardinality of  $\Delta_D^M(X \xrightarrow{M} Y)$  is called the  $M$ -support of  $X \xrightarrow{M} Y$  in  $D$  and denoted by  $S_D^M(X \xrightarrow{M} Y)$ .

The intuitive meaning of the rule  $X \xrightarrow{M} Y$  is that whenever a transaction  $M$ -supports  $X$ , then  $Y$  also appears in the transaction, with some probability. To measure this probability, one only needs to take those transactions with at least one item from the category of  $T(Y)$  into consideration. Accordingly, we have the following definition.

Let  $D(X, T(Y)) = \{t \in D : (t \cap T(X) = X) \wedge (t \cap T(Y) \neq \emptyset)\}$ . The  $M$ -confidence of the maximal association rule  $X \xrightarrow{M} Y$  is defined as

$$C_D^M(X \xrightarrow{M} Y) = \frac{S_D^M(X \xrightarrow{M} Y)}{\text{card}(D(X, T(Y)))}, \quad (2)$$

where  $\text{card}(X)$  denotes the cardinality of a set  $X$ . In addition, we define  $C_D^M(X \xrightarrow{M} Y) = 0$  if  $D(X, T(Y)) = \emptyset$ .

#### 4. Mining regular association rules using soft sets

In this section, we focus on the investigation of mining association rules using soft sets. To facilitate our discussion, we first recall several definitions directly quoted from [16].

**Definition 4.1 [16].** Let  $(F, E)$  be a soft set over the universe  $U$  and  $u \in U$ . An items co-occurrence set in a transaction  $u$  can be defined as

$$\text{Coo}(u) = \{e \in E : f(u, e) = 1\}.$$

Obviously,  $\text{Coo}(u) = \{e \in E : F(e) = 1\}$ .

**Definition 4.2 [16].** Let  $(F, E)$  be a soft set over the universe  $U$  and  $X \subseteq E$ . A set of attributes  $X$  is said to be supported by a transaction  $u \in U$  if  $X \subseteq \text{Coo}(u)$ .

**Definition 4.3 [16].** Let  $(F, E)$  be a soft set over the universe  $U$  and  $X, Y \subseteq E$ , where  $X \cap Y = \emptyset$ . An association rule between  $X$  and  $Y$  is an implication of the form  $X \Rightarrow Y$ . The itemsets  $X$  and  $Y$  are called antecedent and consequent, respectively.

**Definition 4.4 [16].** Let  $(F, E)$  be a soft set over the universe  $U$  and  $X, Y \subseteq E$ , where  $X \cap Y = \emptyset$ . The support of a association rule  $X \Rightarrow Y$ , denoted by  $\text{sup}(X \Rightarrow Y)$  is defined by

$$\text{sup}(X \Rightarrow Y) = \text{sup}(X \cup Y) = |\{u : X \cup Y \subseteq \text{Coo}(u)\}|.$$

**Definition 4.5 [16].** Let  $(F, E)$  be a soft set over the universe  $U$  and  $X, Y \subseteq E$ , where  $X \cap Y = \emptyset$ . The confidence of a association rule  $X \Rightarrow Y$ , denoted by  $\text{conf}(X \Rightarrow Y)$  is defined by

$$\text{conf}(X \Rightarrow Y) = \frac{\text{sup}(X \cup Y)}{\text{sup}(X)} = \frac{|\{u : X \cup Y \subseteq \text{Coo}(u)\}|}{|\{u : X \subseteq \text{Coo}(u)\}|}.$$

In order to make Herawan and Deris's initial idea regarding soft set based association rule mining clearly understood, we slightly revise some concepts and notations in what follows.

As pointed out by Herawan and Deris [16], a transactional dataset can be transformed into a Boolean-valued information system; hence it can be represented by a soft set in a natural way.

**Table 1**  
A transactional dataset  $U$  from [9].

TID	Transaction
$t_0$	A, B, x, y, 1
$t_1$	A, B, D, u, z, 1, 2, 3
$t_2$	A, B, C, z, 1
$t_3$	A, B, x, y, z, 2, 3, 4
$t_4$	C, z, 2, 3
$t_5$	A, B, u, 1, 3
$t_6$	C, D, z, 1, 2
$t_7$	A, B, u, x, y, 4
$t_8$	A, D, z, 2, 4
$t_9$	A, B, x, y, z, 1

More specifically, let  $I = \{i_1, \dots, i_{|I|}\}$  be the item domain and  $D = \{t_1, \dots, t_{|D|}\}$  be a transactional dataset. Then  $D$  can be viewed as a Boolean-valued information system  $\mathcal{S}_D = (D, I)$  such that

$$i_k(t_i) = \begin{cases} 1, & \text{if } i_k \text{ appears in } t_i, \\ 0, & \text{otherwise,} \end{cases}$$

for all  $i_k \in I$  and  $t_i \in D$ .

Furthermore, the information system  $\mathcal{S}_D = (D, I)$  can be represented by a soft set  $\mathfrak{S}_D = (F, A)$  over  $U$  with  $U = D$ ,  $A = I$  and  $F(i_k) = \Delta_D(i_k)$  for all  $i_k \in A$ . That is, the set of  $i_k$ -approximate elements coincides with the realization of the item  $i_k$  in  $D$ . In what follows, the soft set  $\mathfrak{S}_D = (F, A)$  is called the *transactional data soft set* induced by the transactional dataset  $D$ .

**Definition 4.6** (See Definition 4.1 for comparison). Let  $\mathfrak{S} = (F, A) \in \mathcal{S}^E(U)$  and  $u \in U$ . Then  $\text{Co}_{\mathfrak{S}}(u) = \{a \in A : u \in F(a)\}$  is called the *parameter coset* of the object  $u$  in  $\mathfrak{S}$ .

**Definition 4.7** (See Definition 4.2 for comparison). Let  $\mathfrak{S} = (F, A) \in \mathcal{S}^E(U)$  and  $X$  be a nonempty subset of  $A$ . Then  $X$  is *supported* by an object  $u \in U$  in  $\mathfrak{S}$  if  $X \subseteq \text{Co}_{\mathfrak{S}}(u)$ . The set of objects

$$\Delta_{\mathfrak{S}}(X) = \{u \in U : X \subseteq \text{Co}_{\mathfrak{S}}(u)\}$$

is called the *realization* of  $X$  in  $\mathfrak{S}$ .

Intuitively, the parameter coset  $\text{Co}_{\mathfrak{S}}(u)$  consists of all the parameters satisfied by the object  $u$ . On the contrary, the realization  $\Delta_{\mathfrak{S}}(X)$  points out those objects in  $U$  which satisfy all the parameters given by  $X$ . The cardinality of  $\Delta_{\mathfrak{S}}(X)$  is called the *support* of  $X$  in  $\mathfrak{S}$ , which is denoted by  $\text{supp}_{\mathfrak{S}}(X)$ .

As an illustration of the above notions, let us consider an example initially proposed by Amir et al. [9].

**Example 4.8.** Consider the following transactional dataset  $U = \{t_0, t_1, \dots, t_9\}$  consisting of ten transactions given in Table 1. Here the item domain is

$$I = \{A, B, C, D, u, x, y, z, 1, 2, 3, 4\}.$$

Let  $U$  be the universe and  $I$  be the set of parameters. Then the transactional data soft set induced by the dataset  $U$  is a soft set  $\mathfrak{S}_U = (F, I)$  over  $U$  with its tabular representation given in Table 2. For  $t_4 \in U$ , the parameter coset of the object  $t_4$  in the transactional data soft set  $\mathfrak{S}_U$  is a set of parameters as follows:

$$\text{Co}_{\mathfrak{S}_U}(t_4) = \{a \in A : t_4 \in F(a)\} = \{C, z, 2, 3\}.$$

Taking a nonempty set of parameters  $X_0 = \{C, 1\}$ , it is easy to verify that  $X_0$  is supported by  $t_2$  and  $t_6$ . Thus the realization of  $X_0$  in  $\mathfrak{S}_U$  is the following set of objects:

$$\Delta_{\mathfrak{S}_U}(X_0) = \{u \in U : X_0 \subseteq \text{Co}_{\mathfrak{S}_U}(u)\} = \{t_2, t_6\}.$$

In addition, the support of  $X_0$  in  $\mathfrak{S}_U$  is  $\text{supp}_{\mathfrak{S}_U}(X_0) = 2$ .

**Definition 4.9** (See Definition 4.3 for comparison). Let  $\mathfrak{S} = (F, A) \in \mathcal{S}^E(U)$  and  $X, Y$  be nonempty subsets of  $A$  with  $X \cap Y = \emptyset$ . Then

**Table 2**  
Tabular representation of the transactional data soft set  $\mathfrak{S}_U$ .

$U$	A	B	C	D	u	x	y	z	1	2	3	4
$t_0$	1	1	0	0	0	1	1	0	1	0	0	0
$t_1$	1	1	0	1	0	0	1	1	1	1	1	0
$t_2$	1	1	1	0	0	0	0	1	1	0	0	0
$t_3$	1	1	0	0	0	1	1	1	0	1	1	1
$t_4$	0	0	1	0	0	0	0	1	0	1	1	0
$t_5$	1	1	0	0	1	0	0	0	1	0	1	0
$t_6$	0	0	1	1	0	0	0	1	1	1	0	0
$t_7$	1	1	0	0	1	1	1	0	0	0	0	1
$t_8$	1	0	0	1	0	0	0	1	0	1	0	1
$t_9$	1	1	0	0	0	1	1	1	1	0	0	0

a formal expression  $X \Rightarrow Y$  is called an *association rule* between  $X$  and  $Y$  in the soft set  $\mathfrak{S}$ . The sets of parameters  $X, Y$  are referred to as the *antecedent* and *consequent* of the association rule, respectively. The *realization* of  $X \Rightarrow Y$  in  $\mathfrak{S}$  is defined by  $\Delta_{\mathfrak{S}}(X \Rightarrow Y) = \Delta_{\mathfrak{S}}(X \cup Y)$ .

The *support* of an association rule  $X \Rightarrow Y$  is defined as the number of objects in  $U$  supporting both  $X$  and  $Y$ , which is denoted by  $\text{supp}_{\mathfrak{S}}(X \Rightarrow Y)$ . In other words, we have

$$\text{supp}_{\mathfrak{S}}(X \Rightarrow Y) = \text{supp}_{\mathfrak{S}}(X \cup Y) = \text{card}(\Delta_{\mathfrak{S}}(X \cup Y)).$$

**Definition 4.10** (See Definition 4.5 for comparison). Let  $\mathfrak{S} = (F, A) \in \mathcal{S}^E(U)$  and  $X \Rightarrow Y$  be an association rule in  $\mathfrak{S}$ . The *confidence* of  $X \Rightarrow Y$  in  $\mathfrak{S}$  is defined by

$$\text{conf}_{\mathfrak{S}}(X \Rightarrow Y) = \frac{\text{supp}_{\mathfrak{S}}(X \cup Y)}{\text{supp}_{\mathfrak{S}}(X)}.$$

In addition, we define  $\text{conf}_{\mathfrak{S}}(X \Rightarrow Y) = 0$  if  $\text{supp}_{\mathfrak{S}}(X) = 0$ .

To select potentially interesting association rules from a given soft set  $\mathfrak{S} = (F, A)$ , the users are requested to specify a minimum support rate (denoted by  $\sigma$ ) and a minimum confidence level (denoted by  $\gamma$ ) in advance. Accordingly,  $X \subseteq A$  is said to be a  $\sigma$ -frequent set of parameters if  $\text{supp}_{\mathfrak{S}}(X) \geq \sigma \cdot |U|$ . We say that an association rule  $X \Rightarrow Y$  is  $\sigma$ -strong in the soft set  $\mathfrak{S}$  if  $\text{supp}_{\mathfrak{S}}(X \Rightarrow Y) \geq \sigma \cdot |U|$ . Moreover,  $X \Rightarrow Y$  is said to be  $\gamma$ -reliable in  $\mathfrak{S}$  if  $\text{conf}_{\mathfrak{S}}(X \Rightarrow Y) \geq \gamma$ .

**Remark 4.11.** Here we compare the above notions with those defined in [16]. Note first that the above notions are purely established in the framework of soft sets, although they are strongly inspired by association rule mining in transactional datasets. Moreover, in our new definitions, the soft set  $\mathfrak{S} = (F, A)$  is arbitrary, needless to be a soft set transformed from a transactional dataset as required in [16]. Hence, our new definitions not only guarantee the utmost generality of these concepts, but make them more clearly understood from a viewpoint of soft set theory. Last but not least, the corresponding notions regarding association rule mining in transactional datasets can easily be concretized from the above definitions if the soft set is taken to be a transactional data soft set  $\mathfrak{S}_D = (F, A)$  induced by a transactional dataset  $D$ .

To further illustrate several notions mentioned above, let us consider an example as follows.

**Example 4.12** (Continuation of Example 4.8). Consider the transactional dataset  $U = \{t_0, t_1, \dots, t_9\}$  and the soft set  $\mathfrak{S}_U = (F, I)$  induced by the dataset  $U$  in Example 4.8. Let  $X_1 = \{A, B\}$ ,  $Y_1 = \{x\}$  and  $Z_1 = \{1\}$ . Clearly,  $X_1, Y_1$  and  $Z_1$  are pairwise disjoint sets of parameters. Thus by definition,  $X_1 \Rightarrow Y_1, X_1 \Rightarrow Z_1$  and  $Z_1 \Rightarrow Y_1$  are association rules in the soft set  $\mathfrak{S}_U$ . The realizations of these association rules are as follows:

$$\Delta_{\mathfrak{S}_U}(\{A, B\} \Rightarrow \{x\}) = \{t_0, t_3, t_7, t_9\},$$

$$\Delta_{\mathfrak{S}_U}(\{A, B\} \Rightarrow \{1\}) = \{t_0, t_1, t_2, t_5, t_9\}$$

and

$$\Delta_{\mathfrak{S}_U}(\{1\} \Rightarrow \{x\}) = \{t_0, t_9\}.$$

It follows that

$$\text{supp}_{\mathfrak{S}_U}(\{A, B\} \Rightarrow \{x\}) = 4,$$

$$\text{supp}_{\mathfrak{S}_U}(\{A, B\} \Rightarrow \{1\}) = 5$$

and

$$\text{supp}_{\mathfrak{S}_U}(\{1\} \Rightarrow \{x\}) = 2.$$

By calculation, the realizations of  $X_1, Y_1$  and  $Z_1$  in  $\mathfrak{S}_U$  are as follows:

$$\Delta_{\mathfrak{S}_U}(\{A, B\}) = \{t_0, t_1, t_2, t_3, t_5, t_7, t_9\},$$

$$\Delta_{\mathfrak{S}_U}(\{x\}) = \{t_0, t_3, t_7, t_9\}$$

and

$$\Delta_{\mathfrak{S}_U}(\{1\}) = \{t_0, t_1, t_2, t_5, t_6, t_9\}.$$

Hence the supports of  $X_1, Y_1$  and  $Z_1$  in  $\mathfrak{S}_U$  are

$$\text{supp}_{\mathfrak{S}_U}(\{A, B\}) = 7,$$

$$\text{supp}_{\mathfrak{S}_U}(\{x\}) = 4$$

and

$$\text{supp}_{\mathfrak{S}_U}(\{1\}) = 6.$$

Thus we have

$$\text{conf}_{\mathfrak{S}_U}(\{A, B\} \Rightarrow \{x\}) = 4/7,$$

$$\text{conf}_{\mathfrak{S}_U}(\{A, B\} \Rightarrow \{1\}) = 5/7$$

and

$$\text{conf}_{\mathfrak{S}_U}(\{1\} \Rightarrow \{x\}) = 1/3.$$

Now, suppose that the user has specified the minimum support rate  $\sigma_1 = 0.45$  and the minimum confidence level  $\gamma_1 = 0.55$ . Then  $X_1$  and  $Z_1$  are  $\sigma_1$ -frequent sets of parameters since

$$\text{supp}_{\mathfrak{S}_U}(\{A, B\}) = 7 > \sigma_1 \cdot |U| = 4.5$$

and

$$\text{supp}_{\mathfrak{S}_U}(\{1\}) = 6 > \sigma_1 \cdot |U| = 4.5.$$

However,  $Y_1$  is not  $\sigma_1$ -frequent since

$$\text{supp}_{\mathfrak{S}_U}(\{x\}) = 4 < \sigma_1 \cdot |U| = 4.5.$$

It can be seen that  $X_1 \Rightarrow Y_1$  is not  $\sigma_1$ -strong but  $\gamma_1$ -reliable in  $\mathfrak{S}_U$  since

$$\text{supp}_{\mathfrak{S}_U}(\{A, B\} \Rightarrow \{x\}) = 4 < \sigma_1 \cdot |U| = 4.5$$

and

$$\text{conf}_{\mathfrak{S}_U}(\{A, B\} \Rightarrow \{x\}) = 4/7 > \gamma_1 = 0.55.$$

Note also that  $X_1 \Rightarrow Z_1$  is both  $\sigma_1$ -strong and  $\gamma_1$ -reliable in  $\mathfrak{S}_U$  since

$$\text{supp}_{\mathfrak{S}_U}(\{A, B\} \Rightarrow \{1\}) = 5 > \sigma_1 \cdot |U| = 4.5$$

and

$$\text{conf}_{\mathfrak{S}_U}(\{A, B\} \Rightarrow \{1\}) = 5/7 > \gamma_1 = 0.55.$$

In addition, the association rule  $Z_1 \Rightarrow Y_1$  is neither  $\sigma_1$ -strong nor  $\gamma_1$ -reliable in  $\mathfrak{S}_U$  since

$$\text{supp}_{\mathfrak{S}_U}(\{1\} \Rightarrow \{x\}) = 2 < \sigma_1 \cdot |U| = 4.5$$

and

$$\text{conf}_{\mathfrak{S}_U}(\{1\} \Rightarrow \{x\}) = 1/3 < \gamma_1 = 0.55.$$

### 5. Mining maximal association rules using soft sets

Herawan and Deris [16] also initiated the study on mining maximal association rules by virtue of soft sets. To make the present study self-contained, we first recall the following definitions which are directly quoted from [16].

**Definition 5.1** [16]. Let  $(F, E)$  be a soft set over the universe  $U$  and  $X \subseteq E_i$ . A set of attributes  $X$  is said to be maximal supported by a transaction  $u$  if  $X = \text{Coo}(u) \cap E_i$ .

**Definition 5.2** [16]. Let  $(F, E)$  be a soft set over the universe  $U$  and  $X \subseteq E_i$ . The maximal support of a set of parameters  $X$ , denoted by  $\text{Msup}(X)$  is defined by the number of transactions  $U$  maximal supporting  $X$ , i.e.

$$\text{Msup}(X) = |\{u : \text{Coo}(u) \cap E_i\}|,$$

where  $|X|$  is the cardinality of  $X$ .

**Definition 5.3** [16]. Let  $(F, E)$  be a soft set over the universe  $U$  and two maximal itemsets  $X, Y \subseteq E_i$ , where  $X \cap Y = \emptyset$ . A maximal association rule between  $X$  and  $Y$  is an implication of the form  $X \xrightarrow{M} Y$ . The itemsets  $X$  and  $Y$  are called maximal antecedent and maximal consequent, respectively.

**Definition 5.4** [16]. Let  $(F, E)$  be a soft set over the universe  $U$  and two maximal itemsets  $X, Y \subseteq E_i$ , where  $X \cap Y = \emptyset$ . The maximal support of a maximal association rule  $X \xrightarrow{M} Y$ , denoted by  $\text{Msup}(X \xrightarrow{M} Y)$  is defined by

$$\text{Msup}(X \xrightarrow{M} Y) = \text{Msup}(X \cup Y) = |\{u : X \cup Y = \text{Coo}(u) \cap E_i\}|. \quad (3)$$

**Definition 5.5** [16]. Let  $(F, E)$  be a soft set over the universe  $U$  and two maximal itemsets  $X, Y \subseteq E_i$ , where  $X \cap Y = \emptyset$ . The confidence of a maximal association rule  $X \xrightarrow{M} Y$ , denoted by  $\text{Mconf}(X \xrightarrow{M} Y)$  and is defined by

$$\text{Mconf}(X \xrightarrow{M} Y) = \frac{\text{Msup}(X \cup Y)}{\text{Msup}(X)} = \frac{|\{u : X \cup Y = \text{Coo}(u) \cap E_i\}|}{|\{u : X = \text{Coo}(u) \cap E_i\}|}. \quad (4)$$

It should be noted that the above definitions suffer from some drawbacks, which might hinder further research following this line of exploration. In what follows, we point out the shortcomings stemming from the above definitions and offer some modifications to the original definitions.

**Definition 5.6.** Let  $\mathfrak{S} = (F, A) \in \mathcal{S}^E(U)$  and  $T = \{C_1, \dots, C_{|T|}\}$  be a partition of the parameter set  $A$ . Then the triple  $\mathfrak{S} = (F, A, T)$  is called a *parameter-taxonomic soft set* over  $U$ .

The partition  $T$  is also referred to as the *taxonomy* of the parameter set  $A$ . Each member of  $T$  is called a *category*. For every parameter  $a \in A$ , the unique category containing  $a$  is denoted by  $C_a$ . If  $X$  is a nonempty set of parameters from a single category  $A_i \in T$ , then  $A_i$  is also denoted by  $C_X$  in what follows.

**Definition 5.7** (See Definitions 5.1 and 5.2 for comparison). Let  $\mathfrak{S} = (F, A, T)$  be a parameter-taxonomic soft set over  $U$ ,  $u \in U$  and  $X$  be a nonempty set of parameters in the category  $C_X$ . Then  $X$  is said to be *M-supported* by  $u \in U$  in  $\mathfrak{S}$  if  $\text{Co}_{\mathfrak{S}}(u) \cap C_X = X$ . The set  $\Delta_{\mathfrak{S}}^M(X) = \{u \in U : \text{Co}_{\mathfrak{S}}(u) \cap C_X = X\}$  is called the *M-realization* of  $X$  in the soft set  $\mathfrak{S}$ . The cardinality of  $\Delta_{\mathfrak{S}}^M(X)$  is called the *M-support* of  $X$  in  $\mathfrak{S}$  and denoted by  $\text{supp}_{\mathfrak{S}}^M(X)$ .

**Proposition 5.8.** Suppose that  $\mathfrak{S} = (F, A, T)$  is a parameter-taxonomic soft set over  $U$ . Let  $X_1, X_2$  be two sets of parameters in some category  $C \in T$ . If  $\Delta_{\mathfrak{S}}^M(X_1) \cap \Delta_{\mathfrak{S}}^M(X_2) \neq \emptyset$ , then  $X_1 = X_2$ .

**Proof.** Note first that  $C_{X_1} = C_{X_2} = C$  by the hypothesis. Assume that  $u \in \Delta_{\mathfrak{S}}^M(X_1) \cap \Delta_{\mathfrak{S}}^M(X_2) \neq \emptyset$ . Then we have

$$\text{Co}_{\mathfrak{S}}(u) \cap C_{X_1} = \text{Co}_{\mathfrak{S}}(u) \cap C = X_1$$

and

$$\text{Co}_{\mathfrak{S}}(u) \cap C_{X_2} = \text{Co}_{\mathfrak{S}}(u) \cap C = X_2.$$

It follows that  $X_1 = X_2$  as required.  $\square$

The above result indicates that in a parameter-taxonomic soft set, each object  $u$  could  $M$ -supports at most one set of parameters from any chosen category. In addition, for any object  $u$  and any category  $C$ , we can find the unique set of parameters  $X$  which is  $M$ -supported by  $u$  whenever it does exist.

**Proposition 5.9.** Let  $\mathfrak{S} = (F, A, T)$  be a parameter-taxonomic soft set over  $U$ ,  $C \in T$  and  $u \in U$ . If  $\text{Co}_{\mathfrak{S}}(u) \cap C = X \neq \emptyset$ , then  $u \in \Delta_{\mathfrak{S}}^M(X)$ .

**Proof.** Assume that  $\text{Co}_{\mathfrak{S}}(u) \cap C = X \neq \emptyset$ . Then  $X$  is a nonempty set of parameters from a single category  $C_X = C$ . Moreover, we have  $\text{Co}_{\mathfrak{S}}(u) \cap C_X = \text{Co}_{\mathfrak{S}}(u) \cap C = X$ , which implies  $u \in \Delta_{\mathfrak{S}}^M(X)$ .  $\square$

To illustrate the above new notions, let us revisit the example proposed in previous section.

**Example 5.10.** Consider the transactional dataset  $U = \{t_0, t_1, \dots, t_9\}$  and the soft set  $\mathfrak{S}_U = (F, I)$  induced by the dataset  $U$  in Example 4.8. Let  $T = \{C_1, C_2, C_3\}$  be a partition of the set  $I$  of parameters, where  $C_1 = \{A, B, C, D\}$ ,  $C_2 = \{u, x, y, z\}$  and  $C_3 = \{1, 2, 3, 4\}$ . Intuitively,  $C_1$ ,  $C_2$  and  $C_3$  represent for three distinct categories, which are ‘‘Capitals’’, ‘‘Lowercase’’ and ‘‘Digits’’, respectively. It can be seen that  $\mathfrak{T} = (F, I, T)$  is a parameter-taxonomic soft set over  $U$ . Let  $t_0 \in U$  and  $X_1 = \{A, B\} \subseteq C_1$ . Then we have

$$\text{Co}_{\mathfrak{T}}(t_0) = \{A, B, x, y, 1\}.$$

Thus  $X_1$  is  $M$ -supported by  $t_0$  in  $\mathfrak{T}$  since  $\text{Co}_{\mathfrak{T}}(t_0) \cap C_1 = X_1$ . In addition, the  $M$ -realization of  $X_1$  in  $\mathfrak{T}$  is the set

$$\Delta_{\mathfrak{T}}^M(\{A, B\}) = \{t_0, t_3, t_5, t_7, t_9\},$$

and so the  $M$ -support of  $X_1$  in  $\mathfrak{T}$  is  $\text{supp}_{\mathfrak{T}}^M(X_1) = 5$ .

Based on the above results, we present the following algorithms to find all  $M$ -realizations of parameter sets in a given parameter-taxonomic soft set.

**Remark 5.11.** It is clear that the execution time of the above algorithms depends on the size of the input parameter-taxonomic soft set  $\mathfrak{S} = (F, A, T)$  over  $U$ . Theoretically, the time complexity of Algorithm 1 is  $O(|U| \cdot |T|)$  since it needs to search all the objects

---

**Algorithm 1** Construct a set-valued matrix from a parameter-taxonomic soft set.

---

**Procedure:** Construct-Matrix( $\mathfrak{S}$ )

- 1 **Input:** a parameter-taxonomic soft set  $\mathfrak{S} = (F, A, T)$  over  $U$ .
  - 2 **foreach**  $u_i \in U$  **do**
  - 3   **foreach**  $C_j \in T$  **do**
  - 4      $X_{ij} \leftarrow \text{Co}_{\mathfrak{S}}(u_i) \cap C_j$
  - 5   **end foreach**
  - 6 **end foreach**
  - 7 **Output:** a set-valued matrix of parameters  $\mathcal{M} = (X_{ij})_{|U| \times |T|}$ .
- 

in the universe of discourse  $U$  as well as all the categories in the taxonomy  $T$ . In a similar fashion, it can be seen that Algorithm 2 has a complexity of  $O(|U|^2 \cdot |T|)$ .

The following example can help to illustrate some basic ideas about the above algorithms.

**Algorithm 2** Calculate all  $M$ -realizations in a parameter-taxonomic soft set.

**Procedure:** Calculate- $M$ -Realization( $\mathfrak{S}$ )  
 1 **Input:** a parameter-taxonomic soft set  $\mathfrak{S} = (F, A, T)$  over  $U$ .  
 2 Construct the matrix  $\mathcal{M}$  from  $\mathfrak{S}$  using the procedure Construct-Matrix( $\mathfrak{S}$ ).  
 3 Initialize  $i \leftarrow 1, j \leftarrow 1, k \leftarrow 1$  and  $\Delta_{\mathfrak{S}}^M(X_{ik}) \leftarrow \emptyset$ .  
 4 **foreach**  $1 \leq k \leq |T|$  **do**  
 5   **foreach**  $1 \leq i \leq |U|$  **do**  
 6     **foreach**  $1 \leq j \leq |U|$  **do**  
 7       **if**  $X_{ik} = X_{jk}$  and  $X_{ik} \neq \emptyset$  **then**  
 8          $\Delta_{\mathfrak{S}}^M(X_{ik}) \leftarrow \Delta_{\mathfrak{S}}^M(X_{ik}) \cup \{u_i, u_j\}$   
 9       **end foreach**  
 10   **end foreach**  
 11 **end foreach**  
 12 **Output:** all  $M$ -realizations  $\Delta_{\mathfrak{S}}^M(X_{ik})$  of parameter sets in  $\mathfrak{S}$ .

**Table 3**  
A set-valued matrix  $\mathcal{M}$ .

$\mathcal{M}$	$C_1$	$C_2$	$C_3$
$t_0$	{A, B}	{x, y}	{1}
$t_1$	{A, B, D}	{u, z}	{1, 2, 3}
$t_2$	{A, B, C}	{z}	{1}
$t_3$	{A, B}	{x, y, z}	{2, 3, 4}
$t_4$	{C}	{z}	{2, 3}
$t_5$	{A, B}	{u}	{1, 3}
$t_6$	{C, D}	{z}	{1, 2}
$t_7$	{A, B}	{u, x, y}	{4}
$t_8$	{A, D}	{z}	{2, 4}
$t_9$	{A, B}	{x, y, z}	{1}

**Example 5.12** (Continuation of Example 5.10). Let us consider the parameter-taxonomic soft set  $\mathfrak{T} = (F, I, T)$  over  $U$  in Example 5.10. Using Algorithm 1, we can construct a set-valued matrix  $\mathcal{M} = (X_{ij})_{10 \times 3}$  of parameters from  $\mathfrak{T}$ , as shown in Table 3. From the first column of the matrix  $\mathcal{M}$ , we have

$$\Delta_{\mathfrak{T}}^M(\{A, B\}) = \{t_0, t_3, t_5, t_7, t_9\},$$

since  $X_{11} = X_{41} = X_{61} = X_{81} = X_{10,1} = \{A, B\}$ . Similarly, we can obtain other  $M$ -realizations such as

$$\Delta_{\mathfrak{T}}^M(\{A, B, C\}) = \{t_2\},$$

$$\Delta_{\mathfrak{T}}^M(\{A, B, D\}) = \{t_1\},$$

$$\Delta_{\mathfrak{T}}^M(\{A, D\}) = \{t_8\},$$

$$\Delta_{\mathfrak{T}}^M(\{C\}) = \{t_4\}$$

and

$$\Delta_{\mathfrak{T}}^M(\{C, D\}) = \{t_6\}.$$

Using Algorithm 2, one can get all nonempty  $M$ -realizations of parameter sets in  $\mathfrak{T}$  as listed in Table 4.

**Definition 5.13** (See Definition 5.3 for comparison). Let  $\mathfrak{S} = (F, A, T)$  be a parameter-taxonomic soft set over  $U$  and  $X, Y$  be nonempty parameter sets from two distinct categories (i.e.,  $C_X \neq C_Y$ ). Then the formal expression  $X \xrightarrow{M} Y$  is called a *maximal association rule* between  $X$  and  $Y$  in the parameter-taxonomic soft set  $\mathfrak{S}$ . The parameter sets  $X, Y$  are referred to as the *antecedent* and *consequent* of  $X \xrightarrow{M} Y$ , respectively.

**Definition 5.14** (See Definition 5.4 for comparison). Let  $\mathfrak{S} = (F, A, T)$  be a parameter-taxonomic soft set over  $U$  and  $X \xrightarrow{M} Y$  be a maximal association rule in  $\mathfrak{S}$ . The  $M$ -realization of  $X \xrightarrow{M} Y$  is defined as

$$\Delta_{\mathfrak{S}}^M(X \xrightarrow{M} Y) = \{u \in U : (\text{Co}_{\mathfrak{S}}(u) \cap C_X = X) \wedge (Y \subseteq \text{Co}_{\mathfrak{S}}(u))\}.$$

**Table 4**  
Nonempty  $M$ -realizations of parameter sets in  $\mathfrak{T}$ .

Parameter set	$M$ -realization
{A, B}	{ $t_0, t_3, t_5, t_7, t_9$ }
{A, B, C}	{ $t_2$ }
{A, B, D}	{ $t_1$ }
{A, D}	{ $t_8$ }
{C}	{ $t_4$ }
{C, D}	{ $t_6$ }
{u}	{ $t_5$ }
{u, x, y}	{ $t_7$ }
{u, z}	{ $t_1$ }
{x, y}	{ $t_0$ }
{x, y, z}	{ $t_3, t_9$ }
{z}	{ $t_2, t_4, t_6, t_8$ }
{1}	{ $t_0, t_2, t_9$ }
{1, 2}	{ $t_6$ }
{1, 3}	{ $t_5$ }
{1, 2, 3}	{ $t_1$ }
{2, 3}	{ $t_4$ }
{2, 4}	{ $t_8$ }
{2, 3, 4}	{ $t_3$ }
{4}	{ $t_7$ }

The cardinality of  $\Delta_{\mathfrak{S}}^M(X \xrightarrow{M} Y)$  is called the  $M$ -support of  $X \xrightarrow{M} Y$  and denoted by  $\text{supp}_{\mathfrak{S}}^M(X \xrightarrow{M} Y)$ .

**Definition 5.15** ((See Definition 5.5 for comparison). Let  $\mathfrak{S} = (F, A, T)$  be a parameter-taxonomic soft set over  $U$  and  $X \xrightarrow{M} Y$  be a maximal association rule in  $\mathfrak{S}$ . The  $M$ -confidence of  $X \xrightarrow{M} Y$  is defined as

$$\text{conf}_{\mathfrak{S}}^M(X \xrightarrow{M} Y) = \frac{\text{supp}_{\mathfrak{S}}^M(X \xrightarrow{M} Y)}{\text{card}(U(X, Y))},$$

where  $U(X, Y) = \{u \in U : (\text{Co}_{\mathfrak{S}}(u) \cap C_X = X) \wedge (\text{Co}_{\mathfrak{S}}(u) \cap C_Y \neq \emptyset)\}$  is called the *relevant domain* of  $X \xrightarrow{M} Y$ . Besides, we define  $\text{conf}_{\mathfrak{S}}^M(X \xrightarrow{M} Y) = 0$  if  $U(X, Y) = \emptyset$ .

**Remark 5.16.** It is worth noting that the above Definitions 5.13–5.15 revise Herawan and Deris's Definitions 5.3–5.5 in [16], respectively. These amendments are of vital importance for both theory and applications. Note first that the antecedent and consequent of a maximal association rule  $X \xrightarrow{M} Y$  are required to be chosen from the same category (i.e.,  $C_X = C_Y$ , denoted by  $E_i$  in [16]). But according to Amir et al.'s original ideas regarding maximal associations,  $X$  and  $Y$  should be subsets of two distinct categories (see Definition 3.5). Hence it is required that  $C_X \neq C_Y$  in Definition 5.13.

The definitions of  $M$ -support and  $M$ -confidence, as given in [16], can cause even more difficulties. According to Herawan and Deris's Definition 5.4, the  $M$ -support is given by Eq. (3), where  $E_i = C_X = C_Y$  is a category. In other words, the  $M$ -support of a maximal association rule  $X \xrightarrow{M} Y$  is defined to be the number of objects in  $U$  that  $M$ -supports  $X \cup Y$  in  $\mathfrak{S}$ . However, Amir et al. [9] pointed out that a maximal association rule  $X \xrightarrow{M} Y$  intuitively means that whenever  $X$  is the only item of its type in a transaction, than  $Y$  also appears in the transaction with some confidence. In other words, Eq. (3) fails to express Amir et al.'s original ideas exactly.

In addition, according to Herawan and Deris's Definition 5.5, the  $M$ -confidence is given by Eq. (4), where  $E_i = C_X = C_Y$  is a category. However, it can be true that an object  $u \in U$   $M$ -supports  $X \cup Y$  but it does not  $M$ -support  $X$ . In this case, the ratio in Eq. (4) has a zero denominator but a nonzero numerator, which turns out to be a serious problem. These important issues will be further illustrated by a cased study in Section 7.

Like mining regular association rules, minimum  $M$ -support rate and minimum  $M$ -confidence (which are denoted by  $\sigma$  and  $\gamma$ , respectively) should be predetermined for finding useful maximal association rules from a given parameter-taxonomic soft set  $\mathfrak{S} = (F, A, T)$ . We say that a set  $X$  of parameters is  $\sigma$ - $M$ -frequent in  $\mathfrak{S}$  if  $\text{supp}_{\mathfrak{S}}^M(X) \geq \sigma \cdot |U|$ . Noted that a subset of a  $\sigma$ - $M$ -frequent set of parameters might not be  $\sigma$ - $M$ -frequent. A maximal association rule  $X \xrightarrow{M} Y$  is  $\sigma$ - $M$ -strong in the parameter-taxonomic soft set  $\mathfrak{S}$  if  $\text{supp}_{\mathfrak{S}}^M(X \xrightarrow{M} Y) \geq \sigma \cdot |U|$ . In addition,  $X \xrightarrow{M} Y$  is said to be  $\gamma$ - $M$ -reliable in  $\mathfrak{S}$  if  $\text{conf}_{\mathfrak{S}}^M(X \xrightarrow{M} Y) \geq \gamma$ . Then our main purpose is to find all maximal association rules which are  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable from the given parameter-taxonomic soft set.

**Proposition 5.17.** *Let  $\mathfrak{S} = (F, A, T)$  be a parameter-taxonomic soft set over  $U$ . Let  $X \xrightarrow{M} Y$  be a maximal association rule in  $\mathfrak{S}$  and  $\sigma \in [0, 1]$ . If  $X \xrightarrow{M} Y$  is  $\sigma$ - $M$ -strong, then  $X$  is  $\sigma$ - $M$ -frequent and  $Y$  is  $\sigma$ -frequent in  $\mathfrak{S}$ .*

**Proof.** Assume that  $X \xrightarrow{M} Y$  is a  $\sigma$ - $M$ -strong maximal association rule in  $\mathfrak{S}$ . By definitions of realizations and  $M$ -realizations, we have

$$\Delta_{\mathfrak{S}}^M(X \xrightarrow{M} Y) = \Delta_{\mathfrak{S}}^M(X) \cap \Delta_{\mathfrak{S}}(Y).$$

It follows that

$$\sigma \cdot |U| \leq \text{supp}_{\mathfrak{S}}^M(X \xrightarrow{M} Y) \leq \min\{\text{supp}_{\mathfrak{S}}^M(X), \text{supp}_{\mathfrak{S}}(Y)\}.$$

This implies that  $X$  is  $\sigma$ - $M$ -frequent and  $Y$  is  $\sigma$ -frequent in  $\mathfrak{S}$ .  $\square$

However, it should be noted that the converse of the above implication does not hold in general as illustrated by the following example.

**Example 5.18.** Consider the parameter-taxonomic soft set  $\mathfrak{T} = (F, I, T)$  over  $U$  in Example 5.10. Let  $X_1 = \{A, B\}$  and  $Y_2 = \{z\}$ . Clearly,  $X_1$  and  $Y_2$  are nonempty parameter sets from two distinct categories. Thus  $X_1 \xrightarrow{M} Y_2$  is a maximal association rule. Now, assume that the user has specified the minimum  $M$ -support rate  $\sigma_2 = 0.4$ . Note first that  $X_1$  is  $\sigma_2$ - $M$ -frequent since

$$\text{supp}_{\mathfrak{T}}^M(X_1) = |\{t_0, t_3, t_5, t_7, t_9\}| = 5 \geq \sigma_2 \cdot |U| = 4.$$

Also we deduce that  $Y_2$  is  $\sigma_2$ -frequent since  $\text{supp}_{\mathfrak{T}}(Y_2) = 7 \geq \sigma_2 \cdot |U| = 4$ . Nevertheless, we have

$$\text{supp}_{\mathfrak{T}}^M(X_1 \xrightarrow{M} Y_2) = |\{t_3, t_9\}| = 2 < \sigma_2 \cdot |U| = 4,$$

which indicates that  $X_1 \xrightarrow{M} Y_2$  is not a  $\sigma_2$ - $M$ -strong maximal association rule.

Based on the notions and results given above, we propose the following algorithm for mining all  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules from a parameter-taxonomic soft set.

**Remark 5.19.** The above algorithm has three parts. The first part establishes the matrix of parameter sets. The second part generates the classes  $\mathcal{C}_{\sigma,k}^M$  and  $\mathcal{D}_{\sigma,j}$ . The last part produces  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules. The execution time of Algorithm 3 depends on the size of the input parameter-taxonomic soft set  $\mathfrak{S} = (F, A, T)$  over  $U$ . In fact, it can be seen that

$$|\mathcal{C}_{\sigma,k}^M| = |\{X_{ik} \in \mathcal{M} \mid \text{supp}_{\mathfrak{S}}^M(X_{ik}) \geq \sigma \cdot |U|\}| \leq |U|$$

for all  $1 \leq k \leq |T|$ . Note also that

$$|\mathcal{D}_{\sigma,j}| = |\{Y_{sj} \in \mathcal{P}(C_j) \mid \text{supp}_{\mathfrak{S}}(Y_{sj}) \geq \sigma \cdot |U|\}| \leq 2^{|C_*|},$$

where  $1 \leq j \leq |T|$  and  $C_*$  is the largest category in the taxonomy  $T$ . Hence the maximal association rule extracting part of Algorithm 3 has the complexity of  $O(2^{|C_*|} \cdot |U| \cdot |T|^2)$ , which is still

**Algorithm 3** Find all  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules in a parameter-taxonomic soft set

---

**Procedure:** Find- $\sigma$ - $M$ -Strong- $\gamma$ - $M$ -Reliable-Maximal-Rule( $\mathfrak{S}, \sigma, \gamma$ )

- 1 **Input:** a parameter-taxonomic soft set  $\mathfrak{S} = (F, A, T)$  over  $U$  and the thresholds  $\sigma, \gamma \in [0, 1]$ .
- 2 Construct the matrix  $\mathcal{M}$  from  $\mathfrak{S}$  using the procedure Construct-Matrix( $\mathfrak{S}$ ).
- 3 Calculate all  $M$ -realizations  $\Delta_{\mathfrak{S}}^M(X_{ik})$  of sets of parameters in  $\mathfrak{S}$  using the procedure Calculate- $M$ -realization( $\mathfrak{S}$ ) and put  $X_{ik}$  into the class  $\mathcal{C}_{\sigma,k}^M$  if  $\text{supp}_{\mathfrak{S}}^M(X_{ik}) \geq \sigma \cdot |U|$ , where  $1 \leq k \leq |T|$  and  $T = \{C_1, \dots, C_{|T|}\}$ .
- 4 Calculate all realizations  $\Delta_{\mathfrak{S}}(Y_{sj})$  for all  $Y_{sj} \in \mathcal{P}(C_j)$  and put  $Y_{sj}$  into the class  $\mathcal{D}_{\sigma,j}$  if  $\text{supp}_{\mathfrak{S}}(Y_{sj}) \geq \sigma \cdot |U|$ , where  $1 \leq j \leq |T|$  and  $\mathcal{P}(C_j)$  denotes the power set of the category  $C_j$ .
- 5 **foreach**  $1 \leq k \leq |T|$  **do**
- 6   **foreach**  $1 \leq j \leq |T|$  **do**
- 7     **foreach**  $X_{ik} \in \mathcal{C}_{\sigma,k}^M$  **do**
- 8       **foreach**  $Y_{sj} \in \mathcal{D}_{\sigma,j}$  **do**
- 9          **if**  $k \neq j$ , **calculate**  $\text{supp}_{\mathfrak{S}}^M(X_{ik} \xrightarrow{M} Y_{sj}) = |\Delta_{\mathfrak{S}}^M(X_{ik}) \cap \Delta_{\mathfrak{S}}(Y_{sj})|$ .
- 10          **if**  $\text{supp}_{\mathfrak{S}}^M(X_{ik} \xrightarrow{M} Y_{sj}) \geq \sigma \cdot |U|$ , **calculate**  $|U(X_{ik}, Y_{sj})|$ .
- 11          **if**  $\text{supp}_{\mathfrak{S}}^M(X_{ik} \xrightarrow{M} Y_{sj}) \geq \gamma \cdot |U(X_{ik}, Y_{sj})|$ , **put**  $X_{ik} \xrightarrow{M} Y_{sj}$  into  $\mathcal{SRR}_{\sigma,\gamma}^M$ .
- 12       **end foreach**
- 13     **end foreach**
- 14   **end foreach**
- 15 **end foreach**
- 16 **Output:** the class  $\mathcal{SRR}_{\sigma,\gamma}^M$  of all  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules.

---

under  $O(2^{|A|} \cdot |U| \cdot |T|^2)$  even in the worst case since  $|C_1| + |C_2| + \dots + |C_{|T|} = |A|$ . It is worth noting that in some real-world applications one can restrict the classes  $\mathcal{C}_{\sigma,k}^M$  and  $\mathcal{D}_{\sigma,j}$  to some proper subclasses so as to further reduce the number of maximal association rules extracted by Algorithm 3. This can help to significantly reduce the running time of Algorithm 3 as illustrated by Example 6.1 in the next section.

## 6. An application to clinical diagnosis

In this section, an illustrating example is presented to show potential applications of the proposed method in clinical diagnosis.

**Example 6.1.** Suppose that there are 1000 patients  $p_i$  ( $i = 1, 2, \dots, 1000$ ) suffering from several symptoms  $s_j$  ( $j = 1, 2, \dots, 10$ ), which stand for “high fever”, “cough”, “rash”, “diarrhea”, “sore throat”, “runny nose”, “fatigue”, “headache”, “vomiting” and “nasal congestion”, respectively. These symptoms are possibly associated with several diseases  $d_j$  ( $j = 1, 2, 3$ ), which stand for “influenza”, “dengue” and “common cold”, respectively. Let  $U = \{p_1, p_2, \dots, p_{1000}\}$ ,  $S = \{s_1, s_2, \dots, s_{10}\}$ ,  $D = \{d_1, d_2, d_3\}$  and  $A = S \cup D$ . Clearly,  $S \cap D = \emptyset$  and  $T = \{C_1, C_2\} = \{S, D\}$  forms a taxonomy of the parameter set  $A$ . This taxonomy contains two categories  $S$  and  $D$ , namely “Symptoms” and “Diseases”. Clinical diagnosis information regarding these patients is stored in a parameter-taxonomic soft set  $\mathfrak{S} = (F, A, T)$  over  $U$  with its tabular representation given by Table 5, where  $p_{ik}$  represent a certain type of patients in exactly the same condition and “ $|p_{ik}|$ ” gives the number of patients corresponding to the type  $p_{ik}$ .

Now we endeavor to find all  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules in  $\mathfrak{S}$  using Algorithm 3. To do this, let us assume first that the user has chosen the thresholds  $\sigma = 0.15$



**Table 5**  
Tabular representation of the parameter-taxonomic soft set  $\mathfrak{S}$ .

$U$	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$	$s_7$	$s_8$	$s_9$	$s_{10}$	$d_1$	$d_2$	$d_3$	$ p_{i_k} $
$p_{i_0}$	0	1	0	0	1	1	0	0	0	1	0	0	1	70
$p_{i_1}$	1	0	1	1	0	1	1	1	0	0	0	1	0	35
$p_{i_2}$	1	1	0	0	1	0	0	1	0	0	1	0	0	100
$p_{i_3}$	1	1	0	0	0	0	1	1	0	1	0	0	1	50
$p_{i_4}$	1	1	0	1	1	0	0	0	1	0	1	0	0	45
$p_{i_5}$	1	0	1	0	0	0	0	1	1	0	0	1	0	200
$p_{i_6}$	0	1	0	0	1	0	1	1	0	1	0	0	1	30
$p_{i_7}$	1	1	0	0	0	0	1	1	0	1	1	0	0	210
$p_{i_8}$	0	1	0	0	1	0	0	0	0	1	0	0	1	220
$p_{i_9}$	0	1	0	0	1	1	0	0	0	1	1	0	0	40

**Table 6**  
A set-valued matrix  $\mathcal{M}_1$ .

$\mathcal{M}_1$	Symptoms	Diseases
$p_{i_0}$	$\{s_2, s_5, s_6, s_{10}\}$	$\{d_3\}$
$p_{i_1}$	$\{s_1, s_3, s_4, s_6, s_7, s_8\}$	$\{d_2\}$
$p_{i_2}$	$\{s_1, s_2, s_5, s_8\}$	$\{d_1\}$
$p_{i_3}$	$\{s_1, s_2, s_7, s_8, s_{10}\}$	$\{d_3\}$
$p_{i_4}$	$\{s_1, s_2, s_4, s_5, s_9\}$	$\{d_1\}$
$p_{i_5}$	$\{s_1, s_3, s_8, s_9\}$	$\{d_2\}$
$p_{i_6}$	$\{s_2, s_5, s_7, s_8, s_{10}\}$	$\{d_3\}$
$p_{i_7}$	$\{s_1, s_2, s_7, s_8, s_{10}\}$	$\{d_1\}$
$p_{i_8}$	$\{s_2, s_5, s_{10}\}$	$\{d_3\}$
$p_{i_9}$	$\{s_2, s_5, s_6, s_{10}\}$	$\{d_1\}$

**Table 7**  
Nonempty  $M$ -realizations of parameter sets in  $\mathfrak{S}$ .

Parameter set	$M$ -realization	$M$ -support
$\{s_2, s_5, s_6, s_{10}\}$	$\{p_{i_0}, p_{i_9}\}$	110
$\{s_1, s_3, s_4, s_6, s_7, s_8\}$	$\{p_{i_1}\}$	35
$\{s_1, s_2, s_5, s_8\}$	$\{p_{i_2}\}$	100
$\{s_1, s_2, s_7, s_8, s_{10}\}$	$\{p_{i_3}, p_{i_7}\}$	260
$\{s_1, s_2, s_4, s_5, s_9\}$	$\{p_{i_4}\}$	45
$\{s_1, s_3, s_8, s_9\}$	$\{p_{i_5}\}$	200
$\{s_2, s_5, s_7, s_8, s_{10}\}$	$\{p_{i_6}\}$	30
$\{s_2, s_5, s_{10}\}$	$\{p_{i_8}\}$	220
$\{d_1\}$	$\{p_{i_2}, p_{i_4}, p_{i_7}, p_{i_9}\}$	395
$\{d_2\}$	$\{p_{i_1}, p_{i_5}\}$	235
$\{d_3\}$	$\{p_{i_0}, p_{i_3}, p_{i_6}, p_{i_8}\}$	370

and  $\gamma = 0.85$ . Then we can use the procedure ConstructMatrix( $\mathfrak{S}$ ) to construct a set-valued matrix  $\mathcal{M}_1 = (X_{ij})_{10 \times 2}$  from  $\mathfrak{S}$ , as shown in Table 6. Using this matrix, we can calculate  $M$ -realizations of parameter sets in  $\mathfrak{S} = (F, A, T)$  and the obtained results are listed in Table 7.

Since  $\sigma = 0.15$  and  $|U| = 1000$ , one can easily observe from Table 9 that the parameter sets  $X_{41} = \{s_1, s_2, s_7, s_8, s_{10}\}$ ,  $X_{61} = \{s_1, s_3, s_8, s_9\}$  and  $X_{91} = \{s_2, s_5, s_{10}\}$  are  $\sigma - M$ -frequent in  $\mathfrak{S}$ . In other words, we get the following class:

$$C_{\sigma,1}^M = \{\{s_1, s_2, s_7, s_8, s_{10}\}, \{s_1, s_3, s_8, s_9\}, \{s_2, s_5, s_{10}\}\}.$$

In view of the background of our application, it is meaningful to restrict our discussion to those maximal associations whose antecedents are only from the category  $S$ . Hence the other class  $C_{\sigma,2}^M$  will not be calculated and used in what follows. Due to the same reason, we only need to consider the class  $\mathcal{D}_{\sigma,2}$  which consists of the sets of parameters from the category  $D$ . Specifically, we have

$$\mathcal{D}_{\sigma,2} = \{Y_1, Y_2, Y_3\} = \{\{d_1\}, \{d_2\}, \{d_3\}\},$$

since  $Y_1 = \{d_1\}$ ,  $Y_2 = \{d_2\}$  and  $Y_3 = \{d_3\}$  are all the  $\sigma$ -frequent sets of parameters in  $\mathfrak{S}$ .

Next, we calculate the  $M$ -support of those maximal association rules consisting of antecedents from  $C_{\sigma,1}^M$  and consequents from

**Table 8**  
A matrix of  $M$ -supports.

$M$ -support	$Y_1$	$Y_2$	$Y_3$
$X_{41}$	210	0	50
$X_{61}$	0	200	0
$X_{91}$	0	0	220

$\mathcal{D}_{\sigma,2}^M$ . The results are summarized in Table 8, from which we deduce that  $X_{41} \xrightarrow{M} Y_1$ ,  $X_{61} \xrightarrow{M} Y_2$  and  $X_{91} \xrightarrow{M} Y_3$  are  $\sigma - M$ -strong maximal association rules. By further calculation, we have  $|U(X_{41}, Y_1)| = 260$ ,  $|U(X_{61}, Y_2)| = 200$  and  $|U(X_{91}, Y_3)| = 220$ . In addition, we have

$$\begin{aligned} \text{supp}_{\mathfrak{S}}^M(X_{41} \xrightarrow{M} Y_1) &= 210 < \gamma \cdot |U(X_{41}, Y_1)| = 221, \\ \text{supp}_{\mathfrak{S}}^M(X_{61} \xrightarrow{M} Y_2) &= 200 > \gamma \cdot |U(X_{61}, Y_2)| = 170 \end{aligned}$$

and

$$\text{supp}_{\mathfrak{S}}^M(X_{91} \xrightarrow{M} Y_3) = 220 > \gamma \cdot |U(X_{91}, Y_3)| = 187.$$

Therefore,  $X_{61} \xrightarrow{M} Y_2$  and  $X_{91} \xrightarrow{M} Y_3$  are  $\sigma - M$ -strong and  $\gamma - M$ -reliable maximal association rules in the parameter-taxonomic soft set  $\mathfrak{S}$ .

### 7. Case studies

In this section, we conduct two case studies to highlight the essentials of the newly proposed concepts and algorithms.

#### 7.1. Case study 1: a dataset derived from Reuters-21578

In the first case study, we consider a real dataset derived from Reuters-21578 which is widely used as a benchmark for text mining [11,16]. It is a corpus consisting of 30 articles, in which 10 articles discuss the product corn associated with the countries Canada and USA, while other 20 articles are mainly about the product fish and countries Canada, France and USA.

Specifically, we have a parameter set

$$I = \{\text{Corn, Fish, Canada, France, USA}\}$$

and a dataset  $D = \{t_1, t_2, \dots, t_{30}\}$  as shown in Table 9. We divide the parameter set into two categories

$$C_1 = \text{“Products”} = \{\text{Corn, Fish}\}$$

and

$$C_2 = \text{“Countries”} = \{\text{Canada, France, USA}\},$$

which form a taxonomy  $T = \{C_1, C_2\}$ . Using this dataset and taxonomy, we can construct a parameter-taxonomic soft set  $\mathfrak{S}_D = (F, I, T)$  over  $D$ . The approximate function of  $\mathfrak{S}_D$  is defined as

$$F(\text{Corn}) = \{t_1, t_2, \dots, t_{10}\},$$

$$F(\text{Fish}) = F(\text{France}) = \{t_{11}, t_{12}, \dots, t_{30}\},$$

**Table 9**  
A dataset  $D$  from Reuters-21578.

TID	Transaction
$t_1$	Canada, USA, Corn
$t_2$	Canada, USA, Corn
$\vdots$	$\vdots$
$t_{10}$	Canada, USA, Corn
$t_{11}$	Canada, France, USA, Fish
$t_{12}$	Canada, France, USA, Fish
$\vdots$	$\vdots$
$t_{30}$	Canada, France, USA, Fish

**Table 10**  
Tabular representation of the parameter-taxonomic soft set  $\mathfrak{S}_D$ .

$D$	Corn	Fish	Canada	France	USA
$t_1$	1	0	1	0	1
$t_2$	1	0	1	0	1
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$t_{10}$	1	0	1	0	1
$t_{11}$	0	1	1	1	1
$t_{12}$	0	1	1	1	1
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$t_{30}$	0	1	1	1	1

**Table 11**  
Nonempty  $M$ -realizations of parameter sets in  $\mathfrak{S}_D$ .

Parameter set	$M$ -realization	$M$ -support
{Corn}	$\{t_1, t_2, \dots, t_{10}\}$	10
{Fish}	$\{t_{11}, t_{12}, \dots, t_{30}\}$	20
{Canada, USA}	$\{t_1, t_2, \dots, t_{10}\}$	10
{Canada, France, USA}	$\{t_{11}, t_{12}, \dots, t_{30}\}$	20

and  $F(\text{Canada}) = F(\text{USA}) = D$ . Tabular representation of the parameter-taxonomic soft set  $\mathfrak{S}_D = (F, I, T)$  is shown in Table 10.

By using Algorithm 3, it is convenient to find all  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules from the parameter-taxonomic soft set  $\mathfrak{S}_D = (F, I, T)$ . The following steps are performed to achieve this goal.

**Step 1:** We first need to specify certain thresholds. For instance, let  $\sigma = 0.3$  and  $\gamma = 0.9$ .

**Step 2:** Then a set-valued matrix

$$\mathcal{M}_2 = (X_{ij})_{30 \times 2}$$

$$= \begin{pmatrix} X_{11} = \{\text{Corn}\} & X_{12} = \{\text{Canada, USA}\} \\ X_{21} = \{\text{Corn}\} & X_{22} = \{\text{Canada, USA}\} \\ \vdots & \vdots \\ X_{10,1} = \{\text{Corn}\} & X_{10,2} = \{\text{Canada, USA}\} \\ X_{11,1} = \{\text{Fish}\} & X_{11,2} = \{\text{Canada, France, USA}\} \\ X_{12,1} = \{\text{Fish}\} & X_{12,2} = \{\text{Canada, France, USA}\} \\ \vdots & \vdots \\ X_{30,1} = \{\text{Fish}\} & X_{30,2} = \{\text{Canada, France, USA}\} \end{pmatrix}$$

is constructed from the parameter-taxonomic soft set  $\mathfrak{S}_D$  by executing the procedure Construct-Matrix( $\mathfrak{S}_D$ ).

**Step 3:** Based on the matrix  $\mathcal{M}_2$ , we can further calculate  $M$ -realizations of parameter sets in  $\mathfrak{S}_D$  using the procedure Calculate- $M$ -realization( $\mathfrak{S}_D$ ) in Algorithm 2. The obtained results are shown in Table 11. Moreover, we get the following two classes

$$\mathcal{C}_{\sigma,1}^M = \{\{\text{Corn}\}, \{\text{Fish}\}\}$$

and

$$\mathcal{C}_{\sigma,2}^M = \{\{\text{Canada, USA}\}, \{\text{Canada, France, USA}\}\},$$

which consist of all  $\sigma$ - $M$ -frequent parameter sets in  $\mathfrak{S}_D$ .

**Table 12**  
The  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules in  $\mathfrak{S}_D$ .

Maximal association rule	$M$ -support	$M$ -confidence
$\{\text{Canada, USA}\} \xrightarrow{M} \{\text{Corn}\}$	10	100%
$\{\text{Canada, France, USA}\} \xrightarrow{M} \{\text{Fish}\}$	20	100%
$\{\text{Corn}\} \xrightarrow{M} \{\text{Canada, USA}\}$	10	100%
$\{\text{Fish}\} \xrightarrow{M} \{\text{Canada, France, USA}\}$	20	100%
$\{\text{Corn}\} \xrightarrow{M} \{\text{Canada}\}$	10	100%
$\{\text{Corn}\} \xrightarrow{M} \{\text{USA}\}$	10	100%
$\{\text{Fish}\} \xrightarrow{M} \{\text{Canada, France}\}$	20	100%
$\{\text{Fish}\} \xrightarrow{M} \{\text{Canada, USA}\}$	20	100%
$\{\text{Fish}\} \xrightarrow{M} \{\text{France, USA}\}$	20	100%
$\{\text{Fish}\} \xrightarrow{M} \{\text{Canada}\}$	20	100%
$\{\text{Fish}\} \xrightarrow{M} \{\text{France}\}$	20	100%
$\{\text{Fish}\} \xrightarrow{M} \{\text{USA}\}$	20	100%

**Step 4:** Since  $\sigma = 0.3$  and  $|D| = 30$ , by calculation we obtain the following two classes

$$\mathcal{D}_{\sigma,1}^M = \{\{\text{Corn}\}, \{\text{Fish}\}\}$$

and

$$\mathcal{D}_{\sigma,2}^M = \{\{\text{Canada, France, USA}\}, \{\text{Canada, France}\}, \{\text{Canada, USA}\}, \{\text{France, USA}\}, \{\text{Canada}\}, \{\text{France}\}, \{\text{USA}\}\},$$

which consist of all  $\sigma$ -frequent parameter sets in  $\mathfrak{S}_D$ .

**Step 5:** We construct potentially interesting maximal association rules either by taking the parameter sets from  $\mathcal{C}_{\sigma,1}^M$  as the antecedents and the parameter sets from  $\mathcal{D}_{\sigma,2}^M$  as the consequents, or by taking the parameter sets from  $\mathcal{C}_{\sigma,2}^M$  as the antecedents and the parameter sets from  $\mathcal{D}_{\sigma,1}^M$  as the consequents. Since  $\sigma = 0.3$ ,  $\gamma = 0.9$  and  $|D| = 30$ , by calculation we can identify all  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules in  $\mathfrak{S}_D$ , as shown in Table 12. It is interesting to see that all these maximal association rules have the same  $M$ -confidence of 100%.

**Remark 7.1.** The  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules captured by Algorithm 3 include all those rules discovered in [16], namely the first four rules in Table 12 (see Fig. 16 in [16] for comparison). In addition, we also find eight rules with  $M$ -confidence of 100%, which are not reported in [16]. However, it should be noted that all these obtained rules are not consistent with Herawan and Deris's initial definition of maximal association rules (i.e., Definition 5.3). For instance, we consider the first rule

$$\{\text{Canada, USA}\} \xrightarrow{M} \{\text{Corn}\}$$

in Table 12. It is easy to see that

$$\{\text{Canada, USA}\} \subseteq C_{\{\text{Canada, USA}\}} = C_2 = \{\text{Canada, France, USA}\}$$

and

$$\{\text{Corn}\} \subseteq C_{\{\text{Corn}\}} = C_1 = \{\text{Corn}, \text{Fish}\}.$$

Thus  $\{\text{Canada, USA}\} \xrightarrow{M} \{\text{Corn}\}$  is not a maximal association rule according to Definition 5.3, since it is clear that  $C_1 \neq C_2$ . Nevertheless, this rule is indeed a maximal association rule according to Definition 5.13. On the other hand, let us consider the formal expression  $\{\text{USA}\} \xrightarrow{M} \{\text{Canada, France}\}$ , which is not a maximal association rule according to Definition 5.13. But according to Herawan and Deris's initial definition,  $\{\text{USA}\} \xrightarrow{M} \{\text{Canada, France}\}$  is a maximal association rule due to the fact that

$$C_{\{\text{USA}\}} = C_{\{\text{Canada, France}\}} = C_2$$

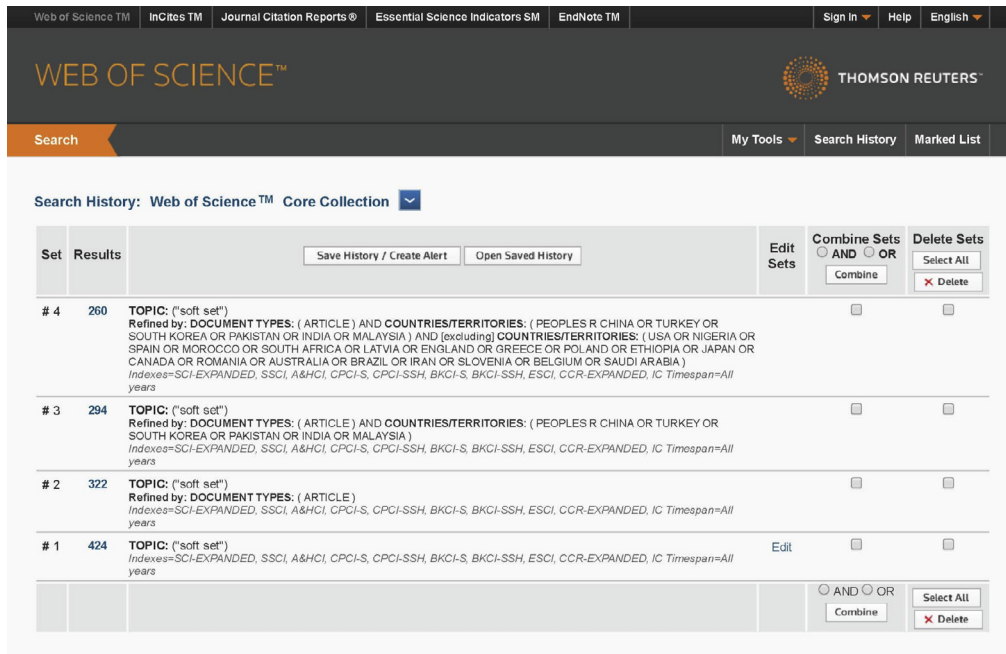


Fig. 3. Process of collecting data in the WoS Core Collection.

and  $\{USA\} \cap \{Canada, France\} = \emptyset$ . Next, let us calculate the maximal support and confidence of this rule. According to Definition 5.4,

$$Msup(\{USA\} \overset{M}{\Rightarrow} \{Canada, France\}) = 20,$$

while by Definition 5.2, we have

$$Msup(\{USA\}) = 0.$$

In this case, Herawan and Deris’s *M*-support (i.e., Definition 5.4) and *M*-confidence (i.e., Definition 5.5) can suffer a major difficulty since the fraction in Eq. (4) has a zero denominator but a nonzero numerator.

### 7.2. Case study 2: a dataset derived from the Web of Science Core Collection

The first case study serves as a running example and shows that our newly proposed notions and algorithms do well in mining maximal association rules from a widely used benchmark dataset. In the following case study, we will focus on collecting and analyzing a real dataset to figure out whether international collaboration can help researchers enhance their research impact in the field of soft sets. This case study also highlights the importance of joint exertion of regular association rules, maximal association rules and many other concepts proposed in Sections 4 and 5.

The dataset used in this case study is generated from the Web of Science (WoS) Core Collection database. We choose WoS Core Collection since it is thought to be the most influential citation index for scientific research worldwide.

As shown in Fig. 3, this real dataset can be collected by taking the following steps in WoS Core Collection:

1. Choose the search field “Topic” and search the term “soft set” in WoS Core Collection. It returns 424 results.
2. Use the search filter “Document Types” and specify it as “Article” to refine the search results. The number of results reduces to 322.
3. Use the search filter “Countries/Territories” to refine the search results by selecting the top 6 countries/territories with respect

to the descending order of record count. The number of results reduces to 294.

4. Use the search filter “Countries/Territories” to refine the search results by excluding all countries/territories except for the top 6 countries/territories selected above. The number of results reduces to 260.

To facilitate the mathematical modelling and analysis of this practical problem using soft set based association rule mining, we need to first establish a parameter-taxonomic soft set using the collected data from the above process. For this purpose, let us consider the universe

$$P = \{p_1, p_2, \dots, p_{260}\}$$

consisting of the papers obtained after going through the above procedure in WoS Core Collection. All information retrieving items associated with these papers constitute the parameter space *E*. Since this case study focuses on investigating the dependency between international collaboration and research impact, it is natural to consider two categories of parameters. The first one is

$$C_1 = \text{“Countries”} = \{c_1, c_2, c_3, c_4, c_5, c_6\}$$

where  $c_j$  ( $j = 1, 2, \dots, 6$ ) stand for China, India, Korea, Malaysia, Pakistan and Turkey, respectively. As mentioned above, we choose these countries since they are top 6 countries according to the number of published articles indexed in WoS Core Collection. The second category under consideration is

$$C_2 = \text{“ESI”} = \{T, F\}.$$

This category “ESI” contains two parameters T and F, indicating whether a paper is an ESI (Essential Science Indicators) highly cited paper. It is well known that ESI highly cited papers are regarded as the most influential papers in a particular research field since they are the top 1% most highly cited articles published in each discipline during a decade. Intuitively, we use ESI highly cited papers to represent significant research impact in this case study.

Now we construct the parameter-taxonomic soft set  $\mathfrak{W} = (G, B, R)$  over *P* with the parameter set  $B = C_1 \cup C_2 \subseteq E$  and the taxonomy  $R = \{C_1, C_2\}$ . The approximate function  $G: B \rightarrow \mathcal{P}(U)$  of  $\mathfrak{W}$

**Table 13**  
Simplified tabular representation of the parameter-taxonomic soft set  $\mathfrak{W}$ .

$P$	China	India	Korea	Malaysia	Pakistan	Turkey	T	F	Count
$p_{[1-3]}$	1	0	0	0	0	0	1	0	3
$p_{[4-109]}$	1	0	0	0	0	0	0	1	106
$p_{[110-115]}$	0	0	0	0	0	1	1	0	6
$p_{[116-165]}$	0	0	0	0	0	1	0	1	50
$p_{[166-168]}$	0	0	1	0	0	0	1	0	3
$p_{[169-185]}$	0	0	1	0	0	0	0	1	17
$p_{[186-187]}$	0	1	0	0	0	0	1	0	2
$p_{[188-215]}$	0	1	0	0	0	0	0	1	28
$p_{[216-229]}$	0	0	0	0	1	0	0	1	14
$p_{[230-239]}$	0	0	0	1	0	0	0	1	10
$p_{[240-243]}$	1	0	1	0	0	0	1	0	4
$p_{[244-250]}$	1	0	1	0	0	0	0	1	7
$p_{251}$	1	0	1	0	1	0	1	0	1
$p_{[252-253]}$	1	0	0	0	1	0	0	1	2
$p_{[254-255]}$	1	0	0	1	0	0	0	1	2
$p_{256}$	1	0	0	0	0	1	0	1	1
$p_{257}$	1	1	0	0	0	0	0	1	1
$p_{258}$	0	0	0	0	1	1	0	1	1
$p_{259}$	0	0	1	0	1	0	1	0	1
$p_{260}$	0	0	1	0	1	0	0	1	1

**Table 14**  
Supports and  $M$ -supports of nontrivial parameter sets in  $\mathfrak{W}$ .

Parameter set	Support	$M$ -support
{China}	127	109
{India}	31	30
{Korea}	34	20
{Malaysia}	12	10
{Pakistan}	20	14
{Turkey}	58	56
{China, India}	1	1
{China, Korea}	12	11
{China, Malaysia}	2	2
{China, Pakistan}	3	2
{China, Turkey}	1	1
{Korea, Pakistan}	3	2
{Pakistan, Turkey}	1	1
{China, Korea, Pakistan}	1	1

is given by

- $p \in G(T) \Leftrightarrow p$  is an ESI highly cited paper,
- $p \in G(F) \Leftrightarrow p$  is not an ESI highly cited paper,
- $p \in G(c_j) \Leftrightarrow$  at least one of the coauthors of  $p$  is from the country  $c_j$ ,

where  $c_j \in C_1 = \{\text{China, India, Korea, Malaysia, Pakistan, Turkey}\}$ .

Table 13 gives a simplified tabular form of the parameter-taxonomic soft set  $\mathfrak{W} = (G, B, R)$ , in which papers of the same type will be incorporated into a single row and the last column counts the number of each type of papers. For instance, the first row indicates that there are three ESI highly cited papers  $p_1, p_2, p_3$  authored by researchers solely from China.

Using parameter-taxonomic soft set  $\mathfrak{W}$  and its simplified tabular form, one can efficiently calculate both supports and  $M$ -supports of parameter sets in the parameter-taxonomic soft set  $\mathfrak{W}$ . For instance, by adding the numbers of the seventh and eighth row in last column, we obtain the  $M$ -support  $\text{supp}_{\mathfrak{W}}^M(\{\text{India}\}) = 30$ . Similarly, if we sum the numbers of the tenth and fifteenth row in last column, we get the support  $\text{supp}_{\mathfrak{W}}(\{\text{Malaysia}\}) = 12$ . Note also that  $M$ -realizations of parameter sets in  $\mathfrak{W}$  can be calculated by running the procedure in Algorithm 2. All the supports and  $M$ -supports of nontrivial parameter sets in  $\mathfrak{W}$  are listed in Table 14.

In light of the context related to this case study, the intuitive meaning of the supports and  $M$ -supports of parameter sets

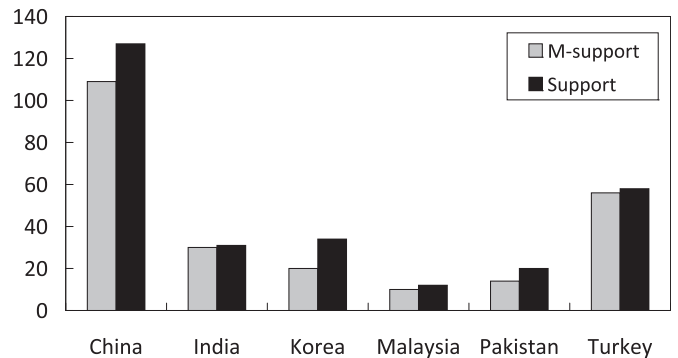


Fig. 4. Comparison between supports and  $M$ -supports of parameter sets in  $C_1$ .

in Table 14 can be interpreted in a natural way. For instance,  $\text{supp}_{\mathfrak{W}}^M(\{\text{Korea}\}) = 20$  means that there are 20 papers authored solely by Korean researchers, while  $\text{supp}_{\mathfrak{W}}(\{\text{Korea}\}) = 34$  shows that there are 34 papers authored by at least one researcher from Korea. The difference

$$\text{supp}_{\mathfrak{W}}(\{\text{Korea}\}) - \text{supp}_{\mathfrak{W}}^M(\{\text{Korea}\}) = 14$$

gives that the number of papers involving international collaboration between researchers from Korea and other countries. Fig. 4 illustrate the comparison between supports and  $M$ -supports of parameter sets corresponding to countries in the category  $C_1$ .

Moreover, using supports and  $M$ -supports of parameter sets in  $\mathfrak{W}$ , a fuzzy set  $F: C_1 \rightarrow [0, 1]$  can be defined as follows:

$$F(c_j) = \frac{\text{supp}_{\mathfrak{W}}(\{c_j\}) - \text{supp}_{\mathfrak{W}}^M(\{c_j\})}{\sum_{i=1}^6 (\text{supp}_{\mathfrak{W}}(\{c_i\}) - \text{supp}_{\mathfrak{W}}^M(\{c_i\}))},$$

where  $c_j \in C_1 = \{\text{China, India, Korea, Malaysia, Pakistan, Turkey}\}$ . It is easy to see that the fuzzy set  $F$  describes the willingness of researchers from each country to participate in international research collaboration. By calculation, we have

$$F = \begin{pmatrix} \text{China} & \text{India} & \text{Korea} & \text{Malaysia} & \text{Pakistan} & \text{Turkey} \\ 0.857 & 0.048 & 0.667 & 0.095 & 0.286 & 0.095 \end{pmatrix},$$

with its plot shown by Fig. 5. It is clear that researchers from China, Korea or Pakistan are more likely to collaborate with foreign researchers in the field of soft sets, while researchers from India, Malaysia or Turkey tend to work with colleagues from their own country.

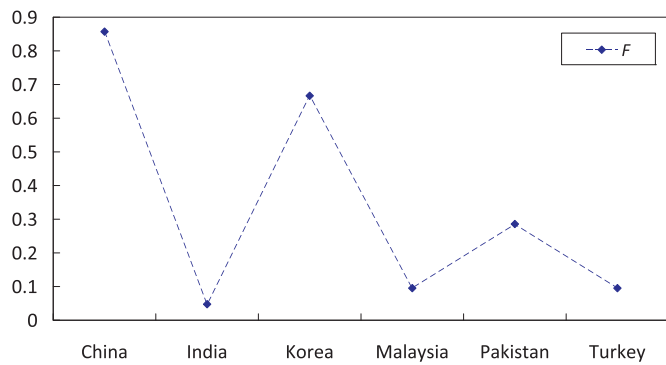


Fig. 5. Illustration of the fuzzy set  $F$ .

Table 15

Maximal association rules in the parameter-taxonomic soft set  $\mathfrak{W}$ .

Rule No.	Maximal association rule	$M$ -support	$M$ -confidence
$M_1$	{China} $\overset{M}{\Rightarrow}$ {T}	3	2.75%
$M_2$	{India} $\overset{M}{\Rightarrow}$ {T}	2	6.67%
$M_3$	{Korea} $\overset{M}{\Rightarrow}$ {T}	3	15%
$M_4$	{Malaysia} $\overset{M}{\Rightarrow}$ {T}	0	0%
$M_5$	{Pakistan} $\overset{M}{\Rightarrow}$ {T}	0	0%
$M_6$	{Turkey} $\overset{M}{\Rightarrow}$ {T}	6	10.71%
$M_7$	{China, Korea} $\overset{M}{\Rightarrow}$ {T}	4	36.36%
$M_8$	{Korea, Pakistan} $\overset{M}{\Rightarrow}$ {T}	1	50%
$M_9$	{China, Korea, Pakistan} $\overset{M}{\Rightarrow}$ {T}	1	100%

Table 16

Regular association rules in the parameter-taxonomic soft set  $\mathfrak{W}$ .

Rule No.	Regular association rule	Support	Confidence
$R_1$	{China} $\Rightarrow$ {T}	8	6.3%
$R_2$	{India} $\Rightarrow$ {T}	2	6.45%
$R_3$	{Korea} $\Rightarrow$ {T}	9	26.47%
$R_4$	{Malaysia} $\Rightarrow$ {T}	0	0%
$R_5$	{Pakistan} $\Rightarrow$ {T}	2	10%
$R_6$	{Turkey} $\Rightarrow$ {T}	6	10.34%
$R_7$	{China, Korea} $\Rightarrow$ {T}	5	41.67%
$R_8$	{Korea, Pakistan} $\Rightarrow$ {T}	2	66.67%
$R_9$	{China, Korea, Pakistan} $\Rightarrow$ {T}	1	100%

Inspired by the discovery mentioned above, an interesting question comes into our mind: would different attitudes toward international collaboration lead to different results regarding research impact? To answer this question, instead of mining all the  $\sigma$ - $M$ -strong and  $\gamma$ - $M$ -reliable maximal association rules in  $\mathfrak{W}$  using Algorithm 3, we only need to comparatively analyse some maximal association rules and regular association rules of the type “{Country}  $\overset{M}{\Rightarrow}$  {T}” and “{Country}  $\Rightarrow$  {T}”. Some important maximal associations and regular associations of this type are listed in Tables 15 and 16, respectively. Each of them essentially discloses some real connections between international collaboration and research impact in the field of soft sets.

For instance, the maximal association rule {China}  $\overset{M}{\Rightarrow}$  {T} with  $M$ -confidence

$$\text{conf}_{\mathfrak{W}}^M(\{China\} \overset{M}{\Rightarrow} \{T\}) = 2.75\%$$

tells us that “If a paper is authored by researchers solely from China, then in 2.75% of the cases this paper is an ESI highly cited paper”. On the other hand, the regular association rule {China}  $\Rightarrow$  {T} with confidence

$$\text{conf}_{\mathfrak{W}}(\{China\} \Rightarrow \{T\}) = 6.3\%$$

reveals that “If a paper is authored by at least one Chinese researcher, then in 6.3% of the cases this paper is an ESI highly cited paper”. It is worth noting that the confidence of the regular association rule  $R_1$  is much higher than the corresponding  $M$ -confidence of the maximal association rule  $M_1$ . This indicates that the greater willingness to collaborate with foreign researchers definitely helps Chinese researchers to enhance their research impact. In fact, by comparing the support

$$\text{supp}_{\mathfrak{W}}(\{China\} \Rightarrow \{T\}) = 8$$

of the rule  $R_1$  with the  $M$ -support

$$\text{supp}_{\mathfrak{W}}^M(\{China\} \overset{M}{\Rightarrow} \{T\}) = 3$$

of the rule  $M_1$ , it can be seen that Chinese researchers get five extra ESI highly cited papers published through international collaboration.

As illustrated by Fig. 5, researchers from Korea (or Pakistan) also adopt a positive attitude toward international collaboration. Then it is natural to ask whether Korean (or Pakistani) researchers can also benefit from their greater willingness to collaborate with foreign researchers in the field of soft sets. One can get positive answers to these questions since we have the following results by comparison:

$$\text{conf}_{\mathfrak{W}}(\{Korea\} \Rightarrow \{T\}) = 26.47\% > \text{conf}_{\mathfrak{W}}^M(\{Korea\} \overset{M}{\Rightarrow} \{T\}) = 15\%$$

and

$$\text{conf}_{\mathfrak{W}}(\{Pakistan\} \Rightarrow \{T\}) = 10\% > \text{conf}_{\mathfrak{W}}^M(\{Pakistan\} \overset{M}{\Rightarrow} \{T\}) = 0\%.$$

In particular, by comparing the support

$$\text{supp}_{\mathfrak{W}}(\{Pakistan\} \Rightarrow \{T\}) = 2$$

of the rule  $R_5$  with the  $M$ -support

$$\text{supp}_{\mathfrak{W}}^M(\{Pakistan\} \overset{M}{\Rightarrow} \{T\}) = 0$$

of the rule  $M_5$ , one can see that positive attitude toward international collaboration helps Pakistan researchers to become authors of two ESI highly cited papers.

A deeper analysis by means of soft sets can reveal some more valuable information hidden in the collected data. For instance, since in Table 14

$$\text{supp}_{\mathfrak{W}}^M(\{China, Korea\}) = 11$$

is the largest one among all the  $M$ -supports of parameter sets consisting at least two countries, we find that the international collaboration between Chinese and Korean researchers is most active in the field of soft sets. Furthermore, by comparison the  $M$ -confidence of the rules  $M_1$ ,  $M_3$  and  $M_7$  in Table 15, it is clear to see that close international collaboration between researchers from China and Korea enormously help to increase the research impact of both sides. This fact is illustrated by Fig. 6.

Considering the contrary situation, one might wonder to know whether conservative attitude toward international collaboration will have a different effect on enhancing research impact in the field of soft sets. Let us take India as an example since Indian researchers are most reluctant to collaborate with foreign researchers as illustrated by Fig. 5. On one hand, from the comparison

$$\text{supp}_{\mathfrak{W}}^M(\{India\}) = 30 > \text{supp}_{\mathfrak{W}}^M(\{Pakistan\}) = 14,$$

it can be seen that there are much more papers authored solely by Indian researchers than by Pakistani researchers. Note also that without taking into account international collaboration, Indian researchers have greater research impact than Pakistani researchers since

$$\text{conf}_{\mathfrak{W}}^M(\{India\} \overset{M}{\Rightarrow} \{T\}) = 6.67\% > \text{conf}_{\mathfrak{W}}^M(\{Pakistan\} \overset{M}{\Rightarrow} \{T\}) = 0\%.$$

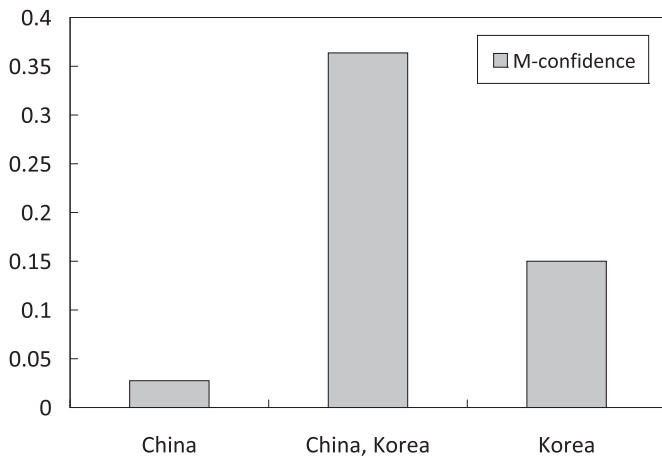


Fig. 6. Histogram illustration of the *M*-confidence of *M*<sub>1</sub>, *M*<sub>3</sub> and *M*<sub>7</sub>.

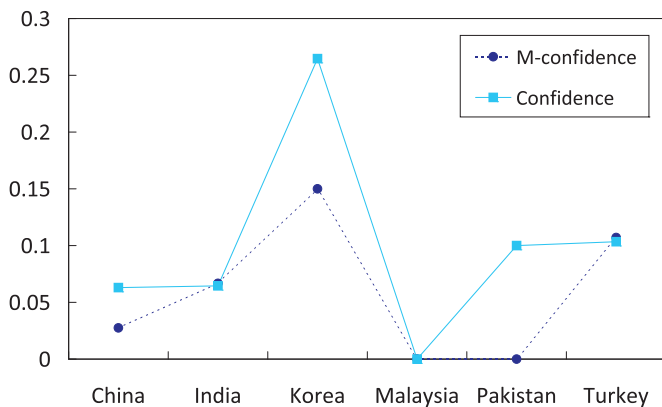


Fig. 7. Comparison between the confidence and *M*-confidence of rules.

On the other hand, it is interesting to see that compared with Pakistani researchers, less international collaboration seems to hinder Indian researchers from increasing their research impact since

$$\text{supp}_{2\mathbb{Y}}^M(\{\text{India}\} \xrightarrow{M} \{\text{T}\}) = \text{supp}_{2\mathbb{Y}}(\{\text{India}\} \Rightarrow \{\text{T}\}) = 2,$$

and

$$\text{conf}_{2\mathbb{Y}}(\{\text{India}\} \Rightarrow \{\text{T}\}) = 6.45\% < \text{conf}_{2\mathbb{Y}}(\{\text{Pakistan}\} \Rightarrow \{\text{T}\}) = 10\%.$$

Some other interesting results can be obtained by analysing both maximal and regular association rules related to Malaysia and Turkey in Tables 15 and 16.

The comparison between the *M*-confidence of maximal association rules *M*<sub>*i*</sub> and the confidence of regular association rules *R*<sub>*i*</sub> (*i* = 1, 2, ..., 5) is illustrated by Fig. 7. Taking all these into account, the above discussion reveals that positive attitude toward international collaboration definitely help Chinese, Korean and Pakistani researchers to enhance their research impact, while the opposite attitude seems to hinder Indian, Malay and Turkish researchers from further increasing their research impact in the field of soft sets.

### 8. Conclusions

The application range of soft set theory was further expanded by Herawan and Deris’s pioneer work on soft set approach to association rule mining. However, as shown above, some fundamental concepts for mining association rules using soft sets were defined improperly in the literature.

This study has systematically investigated soft set based association rule mining and proposed a number of new notions in order to achieve Herawan and Deris’s initial aim more amply. We designed several algorithms for calculating *M*-realizations of parameter sets or identifying  $\sigma$ -*M*-strong and  $\gamma$ -*M*-reliable maximal association rules in parameter-taxonomic soft sets. An illustrative example was presented to show applicability of the newly proposed concepts and algorithms in clinical diagnosis. In addition, we have conducted two case studies to highlight some essential points regarding soft set based association rule mining. The first case study successfully applied our new method to a benchmark dataset derived from Reuters-21578. In the second case study, a real dataset collected from the WoS Core Collection database has been analyzed by means of soft set based association rule mining. It put emphasis on the importance of joint exertion of regular association rules in soft sets, maximal association rules in parameter-taxonomic soft sets and other related notions. With the help of soft set based association rule mining, some interesting facts in the real world have been discovered. For instance, we have found that researchers from China, Korea or Pakistan are more likely to collaborate with foreign researchers in the field of soft sets, while researchers from India, Malaysia or Turkey tend to work with colleagues from their own country. It has been shown that the international collaboration between Chinese and Korean researchers is most active in the field of soft sets. Moreover, it has been revealed that positive attitude toward international collaboration can help researchers to enhance their research impact, while the opposite attitude has a somewhat negative effect.

As future work, it will be interesting to apply the proposed method to other practical cases or formalize soft set based association rule mining in terms of logical formulas over soft sets.

### Acknowledgments

The authors are highly grateful to the anonymous reviewers for their insightful comments and constructive suggestions which greatly improve the quality of this paper. This work was partially supported by National Natural Science Foundation of China (Program Nos. 11301415, 11271237, 11401469), Shaanxi Provincial Research Plan for Young Scientific and Technological New Stars (Program No. 2014KJXX-73), Natural Science Basic Research Plan in Shaanxi Province of China (Program No. 2013JQ1020), Scientific Research Program Funded by Shaanxi Provincial Education Department of China (Program Nos. 16JK1696, 2013JK1098, 2013JK1130, 2013JK1182) and New Star Team of Xi’an University of Posts and Telecommunications.

### References

- [1] R. Agrawal, T. Imielinski, A. Swami, Mining associations between sets of items in large databases, in: Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data, 1993, pp. 207–216.
- [2] R. Agrawal, R. Srikant, Fast algorithms for mining association rules, in: Proceedings of the 20th International Conference on Very Large Data Bases (VLDB), 1994, pp. 487–499.
- [3] M. Akram, N.O. Al-Shehrie, K.P. Shum, A. Farooq, Application of bipolar fuzzy soft sets in *k*-algebras, Ital. J. Pure Appl. Math. 32 (2014) 1–14.
- [4] M. Akram, F. Feng, Soft intersection lie algebras, Quasigroups Rel. Syst. 21 (2013) 1–10.
- [5] H. Aktaş, N. Çağman, Soft sets and soft groups, Inf. Sci. 177 (2007) 2726–2735.
- [6] J.C.R. Alcantud, A novel algorithm for fuzzy soft set based decision making from multiobserver input parameter data set, Inf. Fusion 29 (2016) 142–148.
- [7] M.I. Ali, F. Feng, X.Y. Liu, W.K. Min, M. Shabir, On some new operations in soft set theory, Comput. Math. Appl. 57 (2009) 1547–1553.
- [8] M.I. Ali, M. Shabir, M. Naz, Algebraic structures of soft sets associated with new operations, Comput. Math. Appl. 61 (2011) 2647–2654.
- [9] A. Amir, Y. Aumann, R. Feldman, M. Fresco, Maximal association rules: a tool for mining associations in text, J. Intell. Inf. Syst. 25 (2005) 333–345.
- [10] N. Çağman, S. Enginoğlu, Soft set theory and *uni – int* decision making, Eur. J. Oper. Res. 207 (2010) 848–855.

- [11] R. Feldman, Y. Aumann, A. Amir, A. Zilberstein, W. Klösgen, Maximal association rules: a new tool for mining for keywords cooccurrences in document collections, in: *Proceedings of the Third International Conference on Knowledge Discovery (KDD 1997)*, 1997, pp. 167–170.
- [12] F. Feng, M. Akram, B. Davvaz, V. Leoreanu-Fotea, Attribute analysis of information systems based on elementary soft implications, *Knowl.-Based Syst.* 70 (2014) 281–292.
- [13] F. Feng, C.X. Li, B. Davvaz, M.I. Ali, Soft sets combined with fuzzy sets and rough sets: a tentative approach, *Soft Comput.* 14 (2010) 899–911.
- [14] K. Gong, P. Wang, Z. Xiao, Bijective soft set decision system based parameters reduction under fuzzy environments, *Appl. Math. Modell.* 37 (2013) 4474–4485.
- [15] K. Gong, Z. Xiao, X. Zhang, The bijective soft set with its operations, *Computers and Mathematics with Applications* 60 (2010) 2270–2278.
- [16] T. Herawan, M.M. Deris, A soft set approach for association rules mining, *Knowl.-based Syst.* 24 (2011) 186–195.
- [17] Y.B. Jun, K.J. Lee, A. Khan, Soft ordered semigroups, *Math. Logic Q.* 56 (2010) 42–50.
- [18] Y.B. Jun, C.H. Park, Applications of soft sets in ideal theory of BCK/BCI-algebras, *Inf. Sci.* 178 (2008) 2466–2475.
- [19] P.K. Maji, R. Biswas, A.R. Roy, Soft set theory, *Comput. Math. Appl.* 45 (2003) 555–562.
- [20] P.K. Maji, A.R. Roy, R. Biswas, An application of soft sets in a decision making problem, *Comput. Math. Appl.* 44 (2002) 1077–1083.
- [21] R. Mamat, T. Herawan, M.M. Deris, MAR: Maximum attribute relative of soft set for clustering attribute selection, *Knowl.-Based Syst.* 52 (2013) 11–20.
- [22] D.A. Molodtsov, Soft set theory—first results, *Comput. Math. Appl.* 37 (1999) 19–31.
- [23] Z. Pawlak, A. Skowron, Rudiments of rough sets, *Inf. Sci.* 177 (2007) 3–27.
- [24] K.Y. Qin, Z.Y. Hong, On soft equality, *J. Comput. Appl. Math.* 234 (2010) 1347–1355.
- [25] H. Qin, X. Ma, J.M. Zain, T. Herawan, A novel soft set approach in selecting clustering attribute, *Knowl.-Based Syst.* 36 (2012) 139–145.
- [26] Z. Xiao, K. Gong, S.S. Xia, Y. Zou, Exclusive disjunctive soft sets, *Comput. Math. Appl.* 59 (2010) 2128–2137.
- [27] W. Xu, Z. Xiao, X. Dang, D. Yang, X. Yang, Financial ratio selection for business failure prediction using soft set theory, *Knowl.-Based Syst.* 63 (2014) 59–67.