# Constraint-based mining of fault-tolerant patterns from Boolean data

Jérémy Besson[1,2], Ruggero G. Pensa[1], Céline Robardet[3], and
Jean-François Boulicaut[1]

[1] INSA Lyon, LIRIS CNRS UMR 5205, F-69621 Villeurbanne cedex, France
[2] UMR INRA/INSERM 1235, F-69372 Lyon cedex 08, France
[3] INSA Lyon, PRISMA, F-69621 Villeurbanne cedex, France
{Firstname.Name}@insa-lyon.fr

**Abstract.** Thanks to an important research effort the last few years, inductive queries on local patterns (e.g., set patterns) and complete solvers which can evaluate them on large data sets have been proved extremely useful. The more we use such queries on real-life data, e.g., biological data (and thus intrinsically dirty and noisy), the more we are convinced that inductive queries should return fault-tolerant patterns. In this work, we consider user-defined constraints for a declarative specification of fault-tolerance. We discuss the design of such constraints on bi-sets extracted from Boolean data sets. Our starting point is the fundamental limitation of formal concept discovery (i.e., closed set mining) from noisy data and we propose a constraint-based mining approach for relevant fault-tolerant bi-set mining. Formalizing three recent proposals, our framework enables a better understanding of the needed trade-off between extraction feasibility, completeness, relevancy, and ease of interpretation of these fault-tolerant patterns. An original empirical evaluation on both synthetic and real-life medical data is given. It enables a comparison of the various proposals and it motivates further directions of research.

## 1 Introduction

According to the inductive database approach, mining queries can be expressed declaratively in terms of constraints on the desired patterns or models [1–3]. Thanks to an important research effort the last few years, inductive queries on local patterns (e.g., set or sequential patterns) and complete solvers which can evaluate them on large data sets (Boolean or sequence databases) have been proved extremely useful. Properties of user-defined constraints have been studied in depth (e.g., monotonicity, succinctness, convertibility) and sophisticated pruning strategies enable to compute complete answer sets for many constraints (i.e., Boolean combination of primitive constraints) of practical interest. However, the more we use these techniques on real-life data, e.g., biological or medical data (and thus intrinsically dirty and noisy), the more we are convinced that inductive queries should return fault-tolerant patterns. One interesting direction of research is to introduce softness w.r.t. constraint satisfaction [4, 5]. We consider

in this paper another direction which leads to crispy user-defined constraints in which fault-tolerance is declaratively specified.

Our starting point is the fundamental limitation of formal concept discovery (i.e., connected closed sets) from noisy data. Formal concept analysis has been developed for more than two decades [6] as a way to extract knowledge in Boolean data sets. Informally, formal concepts are maximal bi-sets/rectangles of true values[4]. For instance, Table 1 is a toy example data set $\mathbf{r}_1$ and the bi-set $(\{t_6,\ t_7\}, \{g_1,\ g_2,\ g_3,\ g_4,\ g_5\})$ is a formal concept in $\mathbf{r}_1$.

|       | $g_1$ | $g_2$ | $g_3$ | $g_4$ | $g_5$ | $g_6$ | $g_7$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| $t_1$ | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| $t_2$ | 1 | 1 | 1 | 1 | 0 | 1 | 0 |
| $t_3$ | 0 | 1 | 1 | 1 | 1 | 1 | 1 |
| $t_4$ | 0 | 0 | 0 | 1 | 1 | 1 | 0 |
| $t_5$ | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| $t_6$ | 1 | 1 | 1 | 1 | 1 | 0 | 0 |
| $t_7$ | 1 | 1 | 1 | 1 | 1 | 0 | 0 |

**Table 1.** A Boolean context $\mathbf{r}_1$

Some algorithms are dedicated to the computation of complete collections of formal concepts [7]. Since, by construction, formal concepts are built on closed sets, the extensive research on (frequent) closed set computation (see [8] for a survey) has obviously open new application domains for formal concept discovery. For instance, the formal concept $(\{t_6,\ t_7\}, \{g_1,\ g_2,\ g_3,\ g_4,\ g_5\})$ in $\mathbf{r}_1$ is built on the closed set $\{g_1,\ g_2,\ g_3,\ g_4,\ g_5\}$ whose frequency is 2 (i.e., $|\{t_6,\ t_7\}|$). When considering very large and/or dense Boolean matrices, constraint-based mining of formal concepts has been studied [9, 10]: every formal concept which furthermore satisfies some other user-defined constraints (e.g., a minimal size for set components) is computed. A formal concept associates a maximal set of objects to a maximal set of properties which are all in relation. The strength of such an association is often too strong in real-world data. Even though the extraction might remain tractable, the needed post-processing and interpretation phases turn to be tedious or even impossible. Indeed, in noisy data, not only the number of formal concepts explodes but also many of them are not relevant enough. It has motivated new directions of research where interesting bi-sets are considered as dense rectangles of true values [11–15].

In this paper, we consider a constraint-based mining approach for relevant fault-tolerant formal concept mining. We decided to look for an adequate formalization of three of recent proposals (i.e., CBS [11], FBS [15], and DRBS [14]) which have been motivated by a fault-tolerance extension to formal concepts. We do not provide the algorithms which have been recently published for solving

---

[4] We might say combinatorial bi-sets/rectangles since it is up to arbitrary permutations of lines and columns in the Boolean matrix.

inductive queries on such patterns [11, 15, 14]. The contribution of this paper is to propose a simple framework to support a better understanding of the needed trade-off between extraction feasibility, completeness, relevancy, and ease of interpretation of these various pattern types. This formalization enables to predict part of the behavior of the associated solvers and some formal properties can be established. An original empirical evaluation on both synthetic and real-life medical data is given. It enables to compare the pros and cons of each proposal. An outcome of these experiments is that fault-tolerant bi-set mining is possible. Used in conjunction with other user-defined constraints, this should support the dissemination of relevant local set pattern discovery techniques for intrinsically noisy data.

Section 2 provides the needed definitions. Thanks to the chosen constraint-based mining approach, Section 3 is a discussion on important principles for fault-tolerant formal concept mining. Section 4 provides not only experimental results on synthetic data when various levels of noise are added but also experiments on a real-life medical data sets. Section 5 is a short conclusion.

## 2    Pattern domains

We now define the different classes of patterns to be studied in this paper. Assume a set of objects $\mathcal{O} = \{t_1, \ldots, t_m\}$ and a set of Boolean properties $\mathcal{P} = \{g_1, \ldots, g_n\}$. The Boolean context to be mined is $\mathbf{r} \subseteq \mathcal{O} \times \mathcal{P}$, where $r_{ij} = 1$ if property $g_j$ is satisfied by object $t_i$, 0 otherwise. Formally, a bi-set is an element $(X, Y)$ where $X \subseteq \mathcal{O}$ and $Y \subseteq \mathcal{P}$. $\mathcal{L} = 2^{\mathcal{O}} \times 2^{\mathcal{P}}$ denotes the search space for bi-sets. We say that a bi-set $(X, Y)$ is included in a bi-set $(X', Y')$ (denoted $(X, Y) \subseteq (X', Y')$) iff $(X \subseteq X' \wedge Y \subseteq Y')$.

**Definition 1.** *Let us denote by $\mathcal{Z}_l(x, Y)$ the number of false values of a row $x$ on the columns in $Y$: $\mathcal{Z}_l(x, Y) = \sharp\{y \in Y | (x, y) \notin \mathbf{r}\}$ where $\sharp$ denotes the cardinality of a set. Similarly, $\mathcal{Z}_c(y, X) = \sharp\{x \in X | (x, y) \notin \mathbf{r}\}$ denotes the number of false values of a column $y$ on the rows in $X$.*

Let us now give an original definition of formal concepts (see [6] for a classical one). Sub-constraint 2.1 expresses that a formal concept contains only true values. Sub-constraint 2.2 denotes that formal concept relevancy is enhanced by a maximality property.

**Definition 2 (FC).** *A bi-set $(X, Y) \in \mathcal{L}$ is a formal concept in $\mathbf{r}$ iff*
*(2.1) $\forall x \in X, \ \mathcal{Z}_l(x, Y) = 0 \ \wedge \ \forall y \in Y, \ \mathcal{Z}_c(y, X) = 0$*
*(2.2) $\forall x \in \mathcal{O} \setminus X, \ \mathcal{Z}_l(x, Y) \geq 1 \ \wedge \ \forall y \in \mathcal{P} \setminus Y, \ \mathcal{Z}_c(y, X) \geq 1$*

**Example 1** *Given $\mathbf{r}_1$, we have $\mathcal{Z}_l(t_6, \{g_4, g_5, g_6\}) = 1$ and $\mathcal{Z}_c(g_5, \mathcal{O}) = 2$. $(\{t_3, t_4, t_6, t_7\}, \{g_4, g_5\})$ and $(\{t_3, t_4\}, \{g_4, g_5, g_6\})$ are FC patterns.*

**Definition 3 (DRBS [14]).** *Given integer parameters $\delta \geq 0$ and $\epsilon > 0$, a bi-set $(X, Y) \in \mathcal{L}$ is called a DRBS pattern (Dense and Relevant Bi-Set) in $\mathbf{r}$ iff*
*(3.1) $\forall x \in X, \ \mathcal{Z}_l(x, Y) \leq \delta \ \wedge \ \forall y \in Y, \ \mathcal{Z}_c(y, X) \leq \delta$*

*(3.2)* $\forall e \in \mathcal{O} \setminus X, \ \forall x \in X, \ \mathcal{Z}_l(e, Y) \geq \mathcal{Z}_l(x, Y) + \epsilon$
$\quad\quad\quad \wedge \ \ \forall e' \in \mathcal{P} \setminus Y, \ \forall y \in Y, \ \mathcal{Z}_c(e', X) \geq \mathcal{Z}_c(y, X) + \epsilon$
*(3.3) It is maximal, i.e., $\not\exists (X', Y') \in \mathcal{L}$ s.t. $(X', Y')$ is a DRBS pattern and $(X, Y) \subseteq (X', Y')$.*

DRBS patterns have at most $\delta$ false values per row and per column (Sub-constraint 3.1) and are such that each outside row (resp. column) has at least $\epsilon$ false values plus the maximal number of false values on the inside rows (resp. columns) according to Sub-constraint 3.2. The size of a DRBS pattern increases with $\delta$ such that when $\delta > 0$, it happens that several bi-sets are included in each other. Only maximal bi-sets are kept (Sub-constraint 3.3). Notice that $\delta$ and $\epsilon$ can be chosen differently on rows and on columns.

**Property 1** *When $\delta = 0$ and $\epsilon = 1$, DRBS $\equiv$ FC.*

**Example 2** *If $\delta = \epsilon = 1$, $(X, Y) = (\{t_1, t_2, t_3, t_4, t_6, t_7\}, \{g_3, g_4, g_5\})$ is a DRBS pattern in $\mathbf{r}_1$. Columns $g_1$, $g_2$, $g_6$ and $g_7$ contain at least two false values on $X$, and $t_5$ contains three false values on $Y$.*

The whole collection of DRBS can be computed (in rather small data sets) by using the correct and complete algorithm DR-MINER described in [14]. It is a generic algorithm for bi-set constraint-based mining which is an adaptation of DUAL-MINER [16]. It is based on an enumeration strategy of bi-sets which enables efficient anti-monotonic or monotonic pruning (Sub-constraint 3.1 in conjunction with other user-defined constraints which have monotonicity properties), and partial pruning for Sub-constraint 3.2. Sub-constraint 3.3 is checked in a post-processing phase.

We now consider a preliminary approach for specifying symmetrical fault-tolerant formal concepts. Indeed, the DRBS type has been designed afterwards.

**Definition 4 (CBS [11]).** *Given an integer parameter $\delta$, a bi-set $(X, Y) \in \mathcal{L}$ is called a CBS pattern (Consistent Bi-Set) iff*
*(4.1) $\forall x \in X, \ \mathcal{Z}_l(x, Y) \leq \delta \ \ \wedge \ \ \forall y \in Y, \ \mathcal{Z}_c(y, X) \leq \delta$*
*(4.2) No row (resp. column) outside $(X, Y)$ is identical to a row (resp. column) inside $(X, Y)$*
*(4.3) It is maximal, i.e., $\not\exists (X', Y') \in \mathcal{L}$ s.t. $(X', Y')$ is a CBS pattern and $(X, Y) \subseteq (X', Y')$.*

Notice that again, parameter $\delta$ can be chosen with different values on rows and on columns.

**Example 3** *If $\delta = 1$, $(X, Y) = (\{t_1, t_2, t_3, t_6, t_7\}, \{g_1, g_3, g_5\})$ is a CBS pattern in $\mathbf{r}_1$. Columns $g_6$ and $g_7$ contain more than one false value on $X$, $t_4$ and $t_5$ contain more than one false value on $Y$. $g_2$ and $g_4$ contain only one false value, but as they are identical on $X$, either we add both or they are both excluded. As there are two false values on $t_1$, we do not add them.*

**Property 2** *When $\delta = 0$, CBS $\equiv$ FC. Furthermore, when $\epsilon = 1$, each DRBS pattern is included in one of the CBS patterns.*

[11] proposes an algorithm for computing CBS patterns by merging formal concepts which have been extracted beforehand. The obtained bi-sets are then processed to keep only the maximal ones having less than $\delta$ false values per row and per column. This principle is however incomplete: every bi-set which satisfies the above constraints can not be extracted by this principle. In other terms, some CBS patterns can not be obtained as a merge between two formal concepts. CBS patterns might be extracted by a straightforward adaptation of the DR-MINER generic algorithm but the price to pay for completeness might be expensive.

Let us finally consider another extension of formal concepts which is not symmetrical. It has been designed thanks to some previous work on one of the few approximate condensed representations of frequent sets, the so-called $\delta$-free sets [17, 18]. $\delta$-free sets are well-specified sets whose counted frequencies enable to infer the frequency of many sets (sets included in their so-called $\delta$-closures) without further counting but with a bounded error. When $\delta = 0$, the 0-closure on a 0-free set $X$ is the classical closure and it provides a closed set. The context here is different but the idea is now to consider bi-sets built on $\delta$-free sets with the intuition that it will provide strong associations between sets of rows and sets of columns. It has been introduced for the first time in [15] as a potentially interesting local pattern type for bi-cluster characterization.

Due to space limitation, we do not provide details on $\delta$-freeness and $\delta$-closures [17, 18]. A set $Y \subseteq \mathcal{P}$ is $\delta$-free for a positive integer $\delta$ if its absolute frequency in $\mathbf{r}$ differs from the frequency of all its strict subsets by at most $\delta$. For instance, in $\mathbf{r}_1$, $\{g_2\}$ is a 1-free set. The $\delta$-closure of a set $Y \subseteq \mathcal{P}$ is the superset $Z$ of $Y$ such that every added property ($\in Z \setminus Y$) is almost always true for the objects which satisfy the properties from $Y$: at most $\delta$ false values are enabled. For instance, the 1-closure of $\{g_2\}$ is $\{g_1, g_2, g_3, g_4, g_5\}$. It is possible to consider bi-sets which can be built on $\delta$-free sets and their $\delta$-closures on one hand, on the sets of objects which support the $\delta$-free set on the properties on another hand.

**Definition 5 (FBS).** *A bi-set $(X, Y) \in \mathcal{L}$ is a FBS pattern (Free-set based Bi-Set) iff $Y$ can be decomposed into $Y = K \cup C$ such that $K$ is a $\delta$-free set in $\mathbf{r}$, $C$ is its associated $\delta$-closure and $X = \{t \in \mathcal{O} \mid \forall k \in K, (t, k) \in \mathbf{r}\}$. By construction, $\forall y \in Y, \mathcal{Z}_c(y, X) \leq \delta$ and $\forall y \in K, \mathcal{Z}_c(y, X) = 0$.*

**Property 3** *When $\delta = 0$, FBS $\equiv$ FC.*

**Example 4** *If $\delta = 1$, $\{g_2\}$ is a $\delta$-free set and $(\{t_2, t_3, t_6, t_7\}, \{g_1, g_2, g_3, g_4, g_5\})$ is a FBS pattern in $\mathbf{r}_1$. Another one is $(\{t_3, t_4\}, \{g_2, g_3, g_4, g_5, g_6, g_7\})$. We get at most one false value per column but we have three false values on $t_4$.*

The extraction of FBS can be extremely efficient thanks to $\delta$-freeness anti-monotonicity. The implementation described in [18] can be straightforwardly extended to output FBS patterns. Notice that FBS patterns are bi-sets with

a bounded number of exception per column but every bi-set with a bounded number of exception per column is not necessarily a FBS pattern. An example of a bi-set with at most 1 false value per column which is not a FBS pattern in $\mathbf{r}_1$ is $(\{t_1, t_2, t_3, t_4, t_6, t_7\}, \{g_3, g_4, g_5\})$.

## 3  Discussion

This section discusses the desired properties for formal concept extensions towards fault-tolerant patterns. It enables to consider the pros and the cons of the available proposals and to better understand related open problems.

- **Fault tolerance** Can we control the number of false values inside the bi-sets?

- **Relevancy** Are they consistent w.r.t. the outside rows and columns? At least two views on consistency exist. We might say that a bi-set $B$ is weakly consistent if it is maximal and if we have no row (resp. column) outside $B$ identical to one row (resp. column) inside $B$. $B$ is called strongly consistent if we have no row (resp. column) outside $B$ with at most the same number of false values than one row (resp. column) of $B$.

- **Ease of interpretation** For each bi-set $(X, Y)$, do we have a function which associates $X$ and $Y$ or even better a Galois connection? If a function exists which associates to each set $X$ (resp. $Y$) at most a unique set $Y$ (resp. $XT$), the interpretation of each bi-set is much easier. Furthermore, if the two functions are monotonically decreasing, when the size of $X$ (resp. $Y$) increases, the size of its associated set $Y$ (resp. $X$) decreases. This property is meaningful since the more we have rows inside a bi-set, the less there are columns that can be associated to describe them (or vice versa). One of the appreciated properties of formal concepts is clearly the existence of such functions. If $f_1(X, \mathbf{r}) = \{g \in \mathcal{P} \mid \forall t \in X, (t, g) \in \mathbf{r}\}$ and $f_2(Y, \mathbf{r}) = \{t \in \mathcal{O} \mid \forall g \in Y, (t, g) \in \mathbf{r}\}$, $(f_1, f_2)$ is a Galois connection between $\mathcal{O}$ and $\mathcal{P}$: $f_1$ and $f_2$ are decreasing functions w.r.t. set inclusion.

- **Completeness and efficiency** Can we compute the whole collection of specified bi-sets, i.e., can we ensure a completeness w.r.t. the specified constraints? Is it tractable in practice?

The formal concepts satisfy these properties except the first one. Indeed, we have an explicit Galois connection which enables to compute the complete collection in many data sets of interest. These bi-sets are maximal and consistent but they are not fault-tolerant.

In a FBS pattern, the number of false values are only bounded on columns. They are not strongly consistent because we can have rows outside the bi-set with the same number of false values than a row inside (one of this false value must be on the $\delta$-free set supporting set). On the columns, the property is satisfied. These bi-sets are however weakly consistent. There is no function from column to row sets (e.g., using $\delta = 1$ in $\mathbf{r}_1$, $(\{t_2, t_6, t_7\}, \{g_1, g_2, g_3, g_4, g_5\})$ and $(\{t_1, t_6, t_7\}, \{g_1, g_2, g_3, g_4, g_5\})$ are two FBS with the same set of columns). However, we have a function between $2^{\mathcal{O}}$ to $2^{\mathcal{P}}$. The definition of this pattern is not

symmetrical. In many data sets, including huge and dense ones, complete collections of FBS can be extracted efficiently. Further research is needed for a better characterization of more relevant FBS patterns which might remain easy to extract from huge databases, e.g., what is the impact of different $\delta$-thresholds for the $\delta$-free-set part and the $\delta$-closure computation? how can we avoid an unfortunate distribution of the false values among the same rows?

CBS are symmetrical on rows and columns. Indeed, the number of exceptions is bounded on rows and on columns. CBS are weakly consistent but not strongly consistent (see Example 3). There are neither a function from $2^{\mathcal{O}}$ to $2^{\mathcal{P}}$ nor from $2^{\mathcal{P}}$ to $2^{\mathcal{O}}$ (e.g., $(\{t_1, t_2, t_3, t_4\}, \{g_1, g_3, g_4\})$ and $(\{t_1, t_2, t_3, t_4\}, \{g_2, g_3, g_4\})$ are two CBS with $\delta = 2$ having the same set of rows in the reduction of $\mathbf{r}_1$ to the black rectangle in Table 1). According to the implementation proposal in [11], extracting these patterns can be untractable even in rather small data sets and its extraction strategy is not complete w.r.t. the specified constraints.

By definition, a DRBS has a bounded number of exceptions per row and per column and they are strongly consistent. Two new properties can be considered.

**Property 4 (Existence of functions $\phi$ and $\psi$ on DRBS ($\epsilon > 0$))** *For $\epsilon > 0$, DRBS patterns are embedded by two functions $\phi$ (resp. $\psi$) which associate to $X$ (resp. $Y$) a unique set $Y$ (resp. $X$).*

**Property 5 (Monotonicity of $\phi$ and $\psi$ on DRBS ($\delta$ fixed))** *Let $\mathcal{L}_{\delta,\epsilon}$ the collection of DRBS patterns and $\mathcal{L}'_{\tau\tau'}$ the subset of $\mathcal{L}_{\delta,\epsilon}$ s.t. $(X, Y) \in \mathcal{L}'_{\tau\tau'}$ iff $(X, Y)$ contains at least a row (resp. column) with $\tau$ (resp. $\tau'$) false values in $Y$ (resp. $X$), and such that no row (resp. column) contains more. Then, $\phi$ and $\psi$ are decreasing functions on $\mathcal{L}'_{\tau\tau'}$.*

Unfortunately, the functions loose this property on the whole DRBS collection. Furthermore, we did not identified yet an intentional definition of these functions. As a result, it leads to a quite expensive computation of the complete collection. Looking for such functions is clearly one of the main challenges for further work.

Let us come back to other related work. Co-clustering (bi-clustering) can be applied to Boolean data [19]. It provides linked partitions on both dimensions and it tends to compute rectangles with mainly true (resp. false) values. Heuristic techniques (i.e., local optimization) enable to compute one bi-partition, i.e., a quite restrictive collection of dense bi-sets. In fact, bi-clustering provides a global structure over the data while fault-tolerant formal concepts are typical local patterns which can lead to the discovery of unexpected but yet relevant local associations. Another approach for dense rectangle mining (geometric tiles) has been proposed in [12]. It is however limited to the special case where a built-in order exists on both dimensions. We could also consider previous approaches to fault-tolerant mono-dimensional set pattern mining [20, 21]. The extension of such dense sets to bi-sets is difficult: the connection which associates objects to properties and vice-versa is neither increasing nor decreasing.

## 4 Empirical evaluation

**Experiments on artificially noised data** Let us first discuss the evaluation method. We call $\mathbf{r}_2$ a reference data set, i.e., a data set which is supposed to be noise free and with built-in patterns. Then, we derive various data sets from it by adding some quantity of uniform random noise (i.e., for a X% noise level, each value is randomly changed with a probability of X%). Our goal is to compare the collection of formal concepts extracted from the reference data set with several collections of fault-tolerant formal concepts extracted from the noised matrices. To measure the relevancy of each extracted collection w.r.t the reference one, we test the presence of a subset of the reference collection in each of them. Since both sets of objects and properties of each formal concept can be changed when noise is introduced, we identify those having the largest area in common with the reference. Our measure, called $\sigma$, takes into account the common area and is computed as follows:

$$\sigma(\mathcal{C}_r, \mathcal{C}_a) = \frac{\sigma_1(\mathcal{C}_r, \mathcal{C}_a) + \sigma_2(\mathcal{C}_r, \mathcal{C}_a)}{2}$$

$\sigma_1$ and $\sigma_2$ are defined as follows:

$$\sigma_1(\mathcal{C}_r, \mathcal{C}_a) = \frac{1}{N_r} \sum_{i=1}^{N_r} max_j \left( \frac{(T_i, G_i)_r \cap (T_j, G_j)_a}{(T_i, G_i)_r \cup (T_j, G_j)_a} \right)$$

$$\sigma_2(\mathcal{C}_r, \mathcal{C}_a) = \frac{1}{N_a} \sum_{j=1}^{N_a} max_i \left( \frac{(T_i, G_i)_r \cap (T_j, G_j)_a}{(T_i, G_i)_r \cup (T_j, G_j)_a} \right)$$

where $\mathcal{C}_r$ is the collection of formal concepts computed on the reference data, $\mathcal{C}_a$ is a collection of patterns in a noised data set, $(T_i, G_i)_r$ and $(T_j, G_j)_a$ are bi-sets belonging to $\mathcal{C}_r$ and $\mathcal{C}_a$ respectively, and $N_r$ and $N_a$ are respectively the size of the reference formal concept collection and the size of the noised collection. When $\sigma_1(\mathcal{C}_r, \mathcal{C}_a) = 1$, all the bi-sets belonging to $\mathcal{C}_r$ have identical instances in the collection $\mathcal{C}_a$. Analogously, when $\sigma_2(\mathcal{C}_r, \mathcal{C}_a) = 1$, all the bi-sets belonging to $\mathcal{C}_a$ have identical instances in the collection $\mathcal{C}_r$. Indeed, when $\sigma = 1$, the two collections are identical. High values of $\sigma$, mean not only that we can find all the formal concepts of the reference collection in the noised matrix, but also that the noised collection does not contain many bi-sets that are too different from the reference ones. In this experiment, $\mathbf{r}_2$ concerns 30 objects (rows) and 15 properties (columns) and it contains 3 formal concepts of the same size which are pair-wise disjoints. In other terms, the formal concepts in $\mathbf{r}_2$ are $(\{t_1, \ldots, t_{10}\}, \{g_1, \ldots, g_5\})$, $(\{t_{11}, \ldots, t_{20}\}, \{g_6, \ldots, g_{10}\})$, and $(\{t_{21}, \ldots, t_{30}\}, \{g_{11}, \ldots, g_{15}\})$. Then, we generated 40 different data sets by adding to $\mathbf{r}_2$ increasing quantities of noise (from 1% to 40% of the matrix). A robust technique should be able to capture the three formal concepts concepts even in presence of noise. Therefore, for each data set, we have extracted a collection of formal concepts and different collections of fault-tolerant patterns with
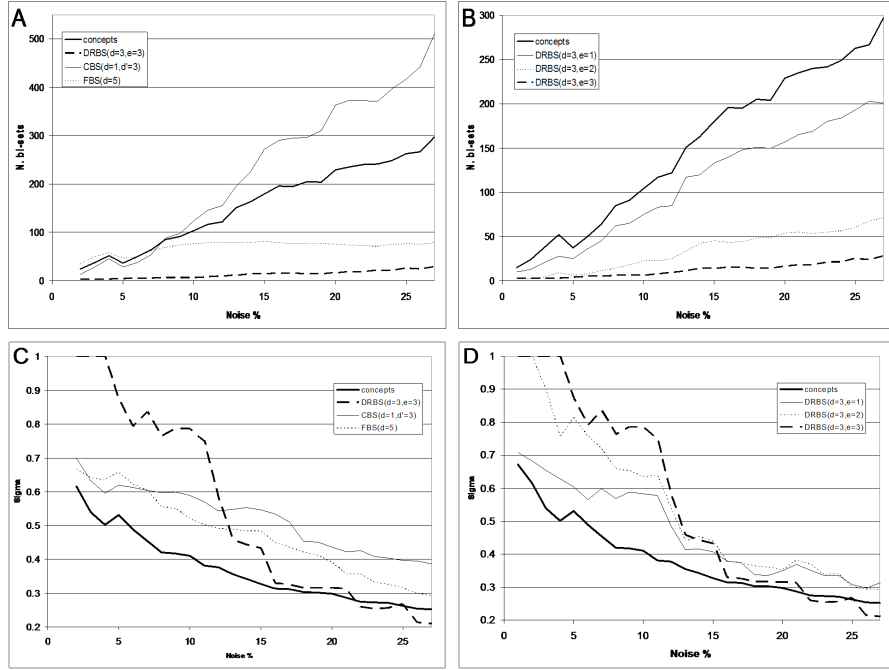
**Fig. 1.** Size of different collections of bi-sets and related values of $\sigma$ w.r.t. noise level for all types of bi-sets (a, b) and for different instances of DRBS collections (c, d)

different parameters. For FBS collections, we considered $\delta$ values between 1 and 6. Then we extracted two groups of CBS collections given parameter $\delta$ (resp. $\delta'$) for the maximum number of false values per row (resp. per column): one with $\delta = 1$ and $\delta' = 1 \ldots 3$ and the second with $\delta' = 1$ and $\delta = 1 \ldots 3$. Finally we extracted DRBS collections for each combination of $\delta = 1 \ldots 3$ and $\epsilon = 1 \ldots 3$.

In Fig. 1(a,c), we only report the best results w.r.t. $\sigma$ for each class of patterns. Fig. 1a presents the number of extracted patterns in each collection. Fault-tolerant bi-set collections contain almost always less patterns than the collection of formal concepts. The only exception is the CBS class when $\delta = 1$. The DRBS class performs better than the other ones. The size of its collections is almost constant, even for rather high levels of noise. The discriminant parameter is $\epsilon$. In Fig. 1c, the values of the $\sigma$ measure for DRBS collections obviously decrease when the noise ratio increases. In general, every class of fault-tolerant bi-set performs better than the formal concept one. In terms of relevancy, the DRBS pattern class gives the best results as well. Notice that the results for FBS and CBS classes are not significantly different when their parameters change. The parameter that mostly influences the value of $\sigma$ is the $\epsilon$ parameter for the DRBS class. For reasonable levels of noise ($< 15\%$), it makes sense to use DRBS. For higher levels, CBS and FBS perform slightly better. In Fig. 1(b,d) we report the

experiments on the extraction of DRBS collections with $\delta = 3$ and $\epsilon = 1 \ldots 3$. Fig. 1b shows the number of extracted patterns. The size of the collections is considerably reduced when $\epsilon$ grows. Fig. 1d presents the $\sigma$ measure for these collections. Using a higher $\epsilon$ value improves the quality of the results because less patterns are produced. When the noise level is smaller than 5%, the collection of DRBS, with $\epsilon = 2..3$, is the same as the three formal concepts in the reference data $\mathbf{r}_2$. This experiment confirms that fault-tolerant bi-sets are more robust to noise than formal concepts, and that the provided collection for the crucially needed expert-driven interpretation is considerably reduced.

**Experiments on a medical data set** It is important to get a qualitative feedback about fault-tolerant pattern relevancy in a real case. For this purpose, we have considered the real world medical data set meningitis [22]. These data have been gathered from children hospitalized for acute meningitis over a period of 4 years. The pre-processed Boolean data set is composed of 329 patients described by 60 Boolean properties encoding clinical signs (e.g., consciousness troubles), cytochemical analysis of the cerebrospinal fluid (e.g., C.S.F proteins) and blood analysis (e.g., sedimentation rate). The majority of the cases are viral infections, whereas about one quarter of the cases are caused by bacteria. It is interesting to look at the bacterial cases since they need treatment with suitable antibiotics, while for the viral cases a simple medical supervision is sufficient. A certain number of attribute-variable pairs have been identified as being characteristic of the bacterial form of meningitis [22, 23]. In other terms, the quality of the fault-tolerant patterns can be evaluated w.r.t. available medical knowledge. Our idea is that by looking for rather large fault-tolerant bi-sets, the algorithms will provide some new associations between attribute-value pairs (Boolean properties) and objects. If the whole sets of objects and properties within bi-sets are compatible (e.g., all the objects are of bacterial type, and all the properties are compatible with bacterial meningitis), then we can argue that we got new relevant information.

| Formal Concepts | | | | | |
|---|---|---|---|---|---|
| size | 354 366 | | | | |
| time | 5s | | | | |
| **FBS** | | | | | |
| $\delta$ | 1 | 2 | 3 | 4 | 5 | 6 |
| size | 141 983 | 67 898 | 39 536 | 25 851 | 18 035 | 13 382 |
| time | 19s | 10s | 6s | 4s | 3s | 2s |
| **DRBS ($\delta=1$)** | | | | | |
| $\epsilon$ | 1 | 2 | 3 | 4 | 5 | 6 |
| size | - | 75 378 | 22 882 | 8 810 | 4 164 | 2 021 |
| time | - | 1507s | 857s | 424s | 233s | 140s |

**Table 2.** Size and extraction time for FBS and DRBS in meningitis.

A straightforward approach to avoid some irrelevant patterns and to reduce the pattern collection size is to use size constraints on bi-set components. For this experiment we set a minimal size of 10 for sets of objects and 5 for sets of properties. Using D-Miner [10], we computed the collection of such large enough formal concepts and we got more than 300 000 formal concepts in a relatively short time (see Table 2). It is obviously hard to exploit such a quantity of patterns. For instance, we were not able to post-process this collection to produce CBS according to [11]. Then, we tried to extract different collections of FBS and DRBS. For FBS, with $\delta = 1$ (at most one exception per column), we got a 60% reduction on the size of the computed bi-sets. Using values of $\delta$ between 2 and 6, this size is reduced at each step by a coefficient between 0.5 and 0.3. We finally used DR-Miner to extract different collections of DRBS. The $\delta$ parameter was set to 1 (at most one exception per row and per column). The $\epsilon$ parameter enables to further reduce the size of the computed collection. Setting $\epsilon = 1$ leads to an untractable extraction but, with $\epsilon = 2$, the resulting collection is 80% smaller than the related formal concept collection. Moreover, with $\delta = 1$ and $\epsilon = 2$ the size of the DRBS collection is considerably smaller than the computed FBS collection for the same constraint (i.e., $\delta = 1$). On the other hand, computational times are sensibly higher.

We now consider relevancy issues. We have been looking for bi-sets containing the property "presence of bacteria detected in C.S.F. bacteriological analysis" with at least one exception. This property is typically true in the bacterial type of meningitis [22, 23]. By looking for bi-sets satisfying such a constraint, we expect to obtain associations between bacterial meningitis objects and properties characterizing this class of meningitis. First we analyzed the collection of FBS when $\delta = 1$. 763 FBS satisfy the chosen constraint. Among these, 124 FBS contain only one viral meningitis object. We got no FBS containing more than one viral object. Properties belonging to these FBS are either characteristic features of the bacterial cases or non discriminant (but compatible) features such as the age and sex of the patient. When $\delta = 2$, the number of FBS satisfying the constraint is 925. Among them, 260 contain at least one viral case of meningitis, and about 25 FBS contain more than one viral case. For $\delta = 5$ the obtained bi-sets are no longer relevant, i.e., the exceptions include contradictory Boolean properties (e.g., presence and absence of bacteria). We performed the same analysis on DRBS for $\epsilon = 2$. We found 24 rather large DRBS. Among them, 2 contain also one viral object. Only one DRBS seems irrelevant: it contains 3 viral and 8 bacterial cases. Looking at its Boolean properties, we noticed that they were not known as discriminant w.r.t. the bacterial meningitis. If we analyze the collection obtained for $\epsilon = 3$, there is only one DRBS satisfying the constraint. It is a rather large bi-set involving 11 Boolean properties and 14 objects. All the 14 objects belong to the bacterial class and the 11 properties are compatible with the bacterial condition of meningitis. It appears that using DRBS instead of FBS leads to a smaller number of relevant bi-sets for our analysis task (24 against 763). Notice however that DRBS are larger than FBS (for an identical number of exceptions): it means that the information provided by several FBS patterns

might be incorporated in only one DRBS pattern. Moreover we got no DRBS pattern whose set of properties is included in the set of properties of another one. This is not the case for FBS. To summarize this experiment, let us first note that using size constraints to reduce the size of the collection is not always sufficient. meningitis is a rather small data set which leads to the extraction of several hundreds of thousands of formal concepts (about 700 000 if no constraint is given). By extracting fault-tolerant bi-sets, we reduce the size of the collection to be interpreted and this is crucial for the targeted exploratory knowledge discovery processes. In particular, for DRBS, the $\epsilon$ parameter is more stringent than the $\delta$ parameter. Then, the relevancy of the extracted patterns can be improved if a reasonable number of exceptions is allowed. For instance, extracting FBS with a low $\delta$ (1 or 2) leads to relevant associations while a high $\delta$ (e.g., 5) introduces too many irrelevant bi-sets. From this point of view, the DRBS class leads to the most interesting results and their quality can be improved by tuning the $\epsilon$ parameter. On the other hand, FBS are easier to compute, even in rather hard contexts, while computing DRBS is in many cases untractable.

## 5   Conclusion

Looking for strong associations between sets of objects and sets of properties in possibly large and noisy Boolean data sets, we have discussed a fundamental limitation of formal concept mining. We lack from consensual extensions of formal concepts towards fault-tolerant patterns and it has given rise to several ad-hoc proposals. Also, relevancy issues are crucial to avoid too many irrelevant patterns during the targeted data mining processes. It is challenging to alleviate the expensive interpretation phases while we still want to promote unexpectedness of the discovered (local) patterns. Considering three recent proposals, we have formalized fault-tolerant bi-dimensional pattern mining within a constraint-based approach. It has been useful for a better understanding of the needed trade-off between extraction feasibility, completeness, relevancy, and ease of interpretation. An empirical evaluation on both synthetic and real-life medical data has been given. It shows that fault-tolerant formal concept mining is possible and this should have an impact on the dissemination of local set pattern discovery techniques in intrinsically noisy Boolean data. DRBS pattern class appears as a well-designed class but the price to pay is computational complexity. The good news are that (a) the submitted inductive queries on DRBS patterns might involve further user-defined constraints which can be used for efficient pruning, and (b) one can look for more efficient data structures and thus a more efficient implementation of the DR-MINER generic algorithm. A pragmatic usage of available algorithms is indeed to extract some bi-sets, e.g., formal concepts, and then select some of them (say $B = (X, Y)$) for further extensions towards fault-tolerant patterns: it becomes, e.g., the computation of a DRBS pattern (say $B' = (X', Y')$ such that the constraint $B \subseteq B'$ is enforced. Also, a better characterization of FBS pattern class might be useful for huge database processing.

# References

1. Imielinski, T., Mannila, H.: A database perspective on knowledge discovery. CACM **39** (1996) 58–64
2. De Raedt, L.: A perspective on inductive databases. SIGKDD Explorations **4** (2003) 69–77
3. Boulicaut, J.F.: Inductive databases and multiple uses of frequent itemsets: the cInQ approach. In: Database Technologies for Data Mining - Discovering Knowledge with Inductive Queries, Springer-Verlag (2004) 1–23
4. Antunes, C., Oliveira, A.: Constraint relaxations for discovering unknown sequential patterns. In: Revised selected and invited papers KDID'04. Volume 3377 of LNCS., Springer-Verlag (2005) 11–32
5. Bistarelli, M., Bonchi, F.: Interestingness is not a dichotomy: introducing softness in constrained pattern mining. In: Proceedings PKDD'05, Porto (PT) (2005) To appear as a Springer-Verlag LNAI volume.
6. Wille, R.: Restructuring lattice theory: an approach based on hierarchies of concepts. In Rival, I., ed.: Ordered sets. Reidel (1982) 445–470
7. Kuznetsov, S.O., Obiedkov, S.A.: Comparing performance of algorithms for generating concept lattices. JETAI **14 (2-3)** (2002) 189–216
8. Goethals, B., Zaki, M.: Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations FIMI 2003, Melbourne, USA (2003)
9. Stumme, G., Taouil, R., Bastide, Y., Pasqier, N., Lakhal, L.: Computing iceberg concept lattices with TITANIC. DKE **42 (2)** (2002) 189–222
10. Besson, J., Robardet, C., Boulicaut, J.F., Rome, S.: Constraint-based concept mining and its application to microarray data analysis. IDA **9(1)** (2005) 59–82
11. Besson, J., Robardet, C., Boulicaut, J.F.: Mining formal concepts with a bounded number of exceptions from transactional data. In: Revised selected and invited papers KDID'04. Volume 3377 of LNCS., Springer-Verlag (2004) 33–45
12. Gionis, A., Mannila, H., Seppänen, J.K.: Geometric and combinatorial tiles in 0-1 data. In: Proceedings PKDD'04. Volume 3202 of LNAI., Pisa, Italy, Springer-Verlag (2004) 173–184
13. Geerts, F., Goethals, B., Mielikäinen, T.: Tiling databases. In: Proceedings DS'04. Volume 3245 of LNAI., Padova, Italy, Springer-Verlag (2004) 278–289
14. Besson, J., Robardet, C., Boulicaut, J.F.: Approximation de collections de concepts formels par des bi-ensembles denses et pertinents. In: Proceedings CAp 2005, Nice, PUG (2005) 313–328 A major revision in English is currently submitted to IEEE ICDM 2005.
15. Pensa, R., Boulicaut, J.F.: From local pattern mining to relevant bi-cluster characterization. In: Proceedings IDA'05. Volume 3646 of LNCS., Madrid, Spain, Springer-Verlag (2005) 293–304
16. Bucila, C., Gehrke, J.E., Kifer, D., White, W.: Dualminer: A dual-pruning algorithm for itemsets with constraints. DMKD **7 (4)** (2003) 241–272
17. Boulicaut, J.F., Bykowski, A., Rigotti, C.: Approximation of frequency queries by mean of free-sets. In: Proceedings PKDD'00. Volume 1910 of LNAI., Lyon, F, Springer-Verlag (2000) 75–85

18. Boulicaut, J.F., Bykowski, A., Rigotti, C.: Free-sets: a condensed representation of boolean data for the approximation of frequency queries. Data Mining and Knowledge Discovery journal **7 (1)** (2003) 5–22

19. Dhillon, I.S., Mallela, S., Modha, D.S.: Information-theoretic co-clustering. In: Proceedings ACM SIGKDD 2003, Washington, USA, ACM Press (2003) 89–98

20. Yang, C., Fayyad, U., Bradley, P.S.: Efficient discovery of error-tolerant frequent itemsets in high dimensions. In: Proceedings ACM SIGKDD'01, San Francisco, USA, ACM Press (2001) 194–203

21. Seppänen, J.K., Mannila, H.: Dense itemsets. In: Proceedings ACM SIGKDD'04, Seattle, USA, ACM Press (2004) 683–688

22. François, P., Robert, C., Cremilleux, B., Bucharles, C., Demongeot, J.: Variables processing in expert system building: application to the aetiological diagnosis of infantile meningitis. Med Inform **15** (1990) 115–124

23. Robardet, C., Crémilleux, B., Boulicaut, J.F.: Characterization of unsupervized clusters by means of the simplest association rules: an application for child's meningitis. In: Proceedings IDAMAP'02 co-located with ECAI'02, Lyon, F (2002) 61–66