



HAL
open science

Combining CSP and Constraint-based Mining for Pattern Discovery

Medhi Khiari, Patrice Boizumault, Bruno Crémilleux

► **To cite this version:**

Medhi Khiari, Patrice Boizumault, Bruno Crémilleux. Combining CSP and Constraint-based Mining for Pattern Discovery. International Conference on Computational Science and Its Applications (ICCSA'10), Mar 2010, Fukuoka, Japan, France. pp.432–447. hal-01016582

HAL Id: hal-01016582

<https://hal.science/hal-01016582>

Submitted on 30 Jun 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Combining CSP and Constraint-Based Mining for Pattern Discovery

Mehdi Khiari, Patrice Boizumault, and Bruno Crémilleux

GREYC, Université de Caen Basse-Normandie, Campus Côte de Nacre,
F-14032 Caen Cedex, France
{Forename.Surname}@info.unicaen.fr

Abstract. A well-known limitation of a lot of data mining methods is the huge number of patterns which are discovered: these large outputs hamper the individual and global analysis performed by the end-users of data. That is why discovering patterns of higher level is an active research field. In this paper, we investigate the relationship between local constraint-based mining and constraint satisfaction problems and we propose an approach to model and mine patterns combining several local patterns, i.e., patterns defined by n-ary constraints. The user specifies a set of n-ary constraints and a constraint solver generates the whole set of solutions. Our approach takes benefit from the recent progress on mining local patterns by pushing with a solver on local patterns all local constraints which can be inferred from the n-ary ones. This approach enables us to model in a flexible way *any* set of constraints combining several local patterns. Experiments show the feasibility of our approach.

1 Introduction

In current scientific, industrial or business areas, the critical need is not to generate data, but to derive knowledge from huge datasets produced at high throughput. Extracting or discovering knowledge from large amounts of data is at the core of the Knowledge Discovery in Databases task, often also named “data mining”. This involves different challenges, such as designing efficient tools to tackle data and the discovery of patterns of a potential user’s interest. There is a large range of methods to discover the patterns but it is well-known that the “pattern flooding which follows data flooding” is an unfortunate consequence in exploratory Knowledge Discovery in Databases processes and the most significant patterns are lost among too much trivial, noisy and redundant information.

Many works propose methods to reduce the collection of patterns, such as the constraint-based paradigm [23], the pattern set discovery approach [9, 17], the so-called condensed representations [5] as well as the compression of the dataset by exploiting the Minimum Description Length Principle [25]. The constraint-based pattern mining framework is a powerful paradigm to discover new highly valuable knowledge [23]. Constraints provide a focus on the most promising knowledge by reducing the number of extracted patterns to those of potential interest for user. There are now generic approaches to discover *local patterns*

(cf. Section 2.1) under constraints [826] and this issue is rather well-mastered, at least for data described by items (i.e., boolean attributes). We call *local constraints* the constraints addressing local patterns. Here, locality refers to the fact that checking whether a pattern satisfies or not a constraint can be performed independently of the other patterns holding in the data. Nevertheless, even if the number of produced local patterns is reduced thanks to the constraint, the output still remains too large for individual and global analysis by the end-user.

On the other hand, the interest of a pattern also depends on the other patterns which are mined. A lot of patterns which are expected by the user (cf. Section 2.2) or models such as classifiers or clustering require to consider simultaneously several patterns to combine the fragmented information conveyed by the local patterns. Local constraints, by considering only one pattern, are insufficient to define and discover such higher patterns. There are few attempts on particular cases by using devoted methods [2818] but there is no generic approach. That is why we claim that discovering patterns under constraints involving comparisons between local patterns is a major issue. In the following of this paper, we call *n-ary constraints* such constraints.

Mining patterns under local constraints requires the exploration of a large search space, even in the case of the simplest patterns, i.e., data described by items. Obviously, mining patterns under n-ary constraints is even harder because we have to take into account and compare the solutions satisfying each pattern involved in a n-ary constraint. In this paper, we investigate the relationship between constraint-based mining and constraint programming and we propose an approach to model and mine patterns under n-ary constraints. As Constraint Satisfaction Problem (CSP) has the ability to define constraints on several variables [1], it is a natural way to model n-ary constraints. We show that each pattern of a n-ary constraint can be assimilated to a variable in the CSP framework. The great advantage of this modeling is its flexibility, it enables us to define a large broad of n-ary constraints. Basically, with our approach, the user specifies the model, that is, the set of n-ary constraints which has to be satisfied, and a constraint solver generates the correct and complete set of solutions. The CSP community has developed several efficient constraint solvers that we can reuse and the resolution can be performed at the level of this global modeling. But we think that it would be a pity not to take benefit from the recent progress on mining local patterns. That is why a key point of our approach is to divide a n-ary constraint in two parts, i.e., a set of local constraints \mathcal{C}_{loc} which is solved by a solver on local patterns and a set of n-ary constraints \mathcal{C}_{n-ary} which is solved by a CSP solver (cf. Section 4 for more details). We claim that is this combination between the local and n-ary levels which enables us the discovery of patterns under n-ary constraints. In other words, the contribution of this paper is to propose an approach joining local constraint mining and set constraint programming in order to model n-ary constraints and discover patterns under such constraints. More generally, the paper investigates the relationship between constraint-based mining and set constraint programming.

This paper is organized as follows. Section 2 sketches definitions and presents the problem statement. The background on pattern discovery and set constraint programming is given in Section 3. We propose our approach to model and mine patterns under n-ary constraints in Section 4. Section 5 details experiments and deals with a discussion and research issues related to our approach.

2 Definitions and Motivations

Below we give definitions used in the paper and the context and motivations.

2.1 Definitions

Let \mathcal{I} be a set of distinct literals called *items*, an itemset (or pattern) is a non-null subset of \mathcal{I} . The language of itemsets corresponds to $\mathcal{L}_{\mathcal{I}} = 2^{\mathcal{I}} \setminus \emptyset$. A transactional dataset is a multi-set of itemsets of $\mathcal{L}_{\mathcal{I}}$. Each itemset, usually called transaction or object, is a database entry. For instance, Table 1 gives a transactional dataset \mathbf{r} where 9 objects o_1, \dots, o_9 are described by 6 items A, \dots, c_2 .

Table 1. Example of a transactional context \mathbf{r}

Trans.	Items
o_1	$A \ B \ c_1$
o_2	$A \ B \ c_1$
o_3	$C \ c_1$
o_4	$C \ c_1$
o_5	$C \ c_1$
o_6	$A \ B \ C \ D \ c_2$
o_7	$C \ D \ c_2$
o_8	$C \ c_2$
o_9	$D \ c_2$

Let X be a local pattern. Pattern mining aims at discovering information from all the patterns or a subset of $\mathcal{L}_{\mathcal{I}}$. More precise, constraint-based mining task selects all the itemsets of $\mathcal{L}_{\mathcal{I}}$ present in \mathbf{r} and satisfying a predicate which is named *constraint*. Local patterns are regularities that hold for a particular part of the data. A local pattern is of special interest if it exhibits a deviating behavior w.r.t. the underlying global model of the data [14] because we are seeking for surprising knowledge which deviates from the already known background model. There are a lot of constraints to evaluate the relevance of local patterns. A well-known example is the frequency constraint which focuses on patterns having a frequency in the database exceeding a given minimal threshold $\gamma > 0$: $freq(X) \geq \gamma$. Many works [23] replace the frequency by other interestingness measures to evaluate the relevance of patterns such as the area of a pattern ($area(X)$ is the product of the frequency of the pattern times its length, i.e., $area(X) = freq(X) \times count(X)$ where $count(X)$ denotes the cardinality of X).

In practice, the user is often interested in discovering more complex patterns such as the simplest rules in the classification task based on associations [30], pairs of exception rules [28] which may reveal global characteristics from the database. The definition of such patterns relies on properties involving several local patterns [6]. These patterns are formalized by the notion of *n-ary constraint*:

Definition 1 (n-ary constraint). *A constraint q is said n-ary if several local patterns have to be compared to check if q is satisfied or not.*

The next section provides more precise examples of n-ary constraints.

2.2 Context and Motivations

N-ary constraints are very useful to design a lot of patterns requested by the users. For instance, the discovery of exception rules from a data set without domain-specific information is of a great interest [28]. An exception rule is defined as a deviational pattern to a strong rule and the interest of an exception rule is evaluated according to another rule. The comparison between rules means that these exception rules are not local patterns. More formally, an exception rule is defined within the context of a pair of rules as follows (I is an item, for instance a class value, X and Y are local patterns):

$$exception(X \rightarrow \neg I) \equiv \begin{cases} true & \text{if } \exists Y \in \mathcal{L}_{\mathcal{I}} \text{ such that } Y \subset X, \text{ one have} \\ & (X \setminus Y \rightarrow I) \wedge (X \rightarrow \neg I) \\ false & \text{otherwise} \end{cases}$$

Such a pair of rules is composed of a common sense rule $X \setminus Y \rightarrow I$ (the term “common sense rule” represents a user-given belief) and an exception rule $X \rightarrow \neg I$ since usually if $X \setminus Y$ then I . The exception rule isolates unexpected information. This definition assumes that the common sense rule has a high frequency and a rather high confidence and the exception rule has a low frequency and a very high confidence (the confidence of a rule $X \rightarrow Y$ is $freq(X \cup Y) / freq(X)$). Assuming that a rule $X \rightarrow Y$ holds iff at least 2/3 of the transactions containing X also contains Y , the rule $AC \rightarrow \neg c_1$ is an exception rule in our running example (cf. Table 1) because we jointly have $A \rightarrow c_1$ and $AC \rightarrow \neg c_1$. Note that Suzuki proposes a method based on sound pruning and probabilistic estimation [28] to extract the exception rules. Nevertheless, this method is devoted to this kind of patterns.

In the context of genomics, local patterns defined by groups of genes and satisfying the *area* constraint previously introduced above are at the core of the discovery of synexpression groups [15]. Nevertheless, in noisy data such as transcriptomic data, the search of fault-tolerant patterns is very useful to cope with the intrinsic uncertainty embedded in the data [3]. N-ary constraints are a way to design such fault-tolerant patterns: larger sets of genes with few exceptions are expressed by the union of several local patterns satisfying an area constraint

and having a large overlapping between them. From two local patterns, it corresponds to the following n-ary constraint: $area(X) > min_{area} \wedge area(Y) > min_{area} \wedge (area(X \cap Y) > \alpha \times min_{area})$ where min_{area} denotes the minimal area and α is a threshold given by the user to fix the minimal overlapping between the local patterns. The set of n-ary constraints can also be extended by the use of the universal quantifier (see Section 6).

Section 4 presents our approach to model patterns satisfying such n-ary constraints and how we combine local constraint mining and set constraint programming to extract these patterns.

3 Background: Related Works and Set CSP

3.1 Local Patterns and Pattern Sets Discovery

As said in the introduction, there are a lot of works to discover local patterns under constraints. A key issue of these works is the use of the property of monotonicity because pruning conditions are straightforwardly deduced [21]. A constraint q is anti-monotone w.r.t. the item specialization iff for all $X \in \mathcal{L}_{\mathcal{I}}$ satisfying q , any subset of X also satisfies q . In this paper, we use the MUSIC-DFS [1] prototype because it offers a set of syntactic and aggregate primitives to specify a broad spectrum of constraints in a flexible way [27]. MUSIC-DFS mines soundly and completely all the patterns satisfying a given set of input local constraints. The efficiency of MUSIC-DFS lies in its depth-first search strategy and a safe pruning of the pattern space exploiting the anti-monotonicity property to push the local constraints as early as possible. The pruning conditions are based on intervals representing several local patterns. The local patterns satisfying all the local constraints are provided in a condensed representation made of intervals (each interval represents a set of patterns satisfying the constraint and each pattern appears in only one interval). The lower bound of an interval is a prefix-free pattern and its upper bound is the prefix-closure of the lower bound [27].

There are also other approaches to combine local patterns. Recent approaches - pattern teams [17], constraint-based pattern set mining [9] and selecting patterns according to the added value of a new pattern given the currently selected patterns [4] - aim at reducing the redundancy by selecting patterns from the initial large set of local patterns on the basis of their usefulness in the context of the other selected patterns. Even if these approaches explicitly compare patterns between them, they are mainly based on the reduction of the redundancy or specific aims such as classification processes. We think that n-ary constraints are a flexible way to take into account a bias given by the user to direct the final set of patterns toward a specific aim such as the search of exceptions. General data mining frameworks based on the notion of local patterns to design global models are presented in [16][13]. These frameworks help to analyze and improve current methods in the area. In our approach (cf. Section 4), we show the interest of the set constraint programming in this general issue of combining local patterns.

¹ <http://www.info.univ-tours.fr/~soulet/music-dfs/music-dfs.html>

Constraint programming is a powerful declarative paradigm for solving difficult combinatorial problems. In a constraint programming approach, one specifies constraints on acceptable solutions and search is used to find a solution that satisfies the constraints. A first approach using Constraint Programming for itemset mining has been proposed in [7]. In this work, constraints such as frequency, closedness, maximality, and constraints that are monotonic or anti-monotonic or variations of these constraints are modeled using 0/1 Linear Programming. Then patterns satisfying these constraints are obtained by using the constraint solver Gecode [11]. This work presents in a unified framework a large set of patterns but does not address patterns modeled by relationships between several local patterns as those described in Section 2. Recently, this work has been extended in order to find correlated patterns (i.e., patterns having the highest score w.r.t. a correlation measure) [24].

3.2 Set CSP

Formally a Constraint Satisfaction Problem (CSP) is a 3-uple $(\mathcal{X}, \mathcal{D}, \mathcal{C})$ where \mathcal{X} is a set of variables, \mathcal{D} is a set of finite domains and \mathcal{C} is a set of constraints that restrict certain simultaneous variables assignments. There are several types of CSPs such as numerical CSPs, boolean CSPs, set CSPs, etc. They differ fundamentally from the domain types and filtering techniques. We present here more precisely set CSPs that are used in our modeling. First, we define Set Intervals. Then we introduce set CSPs, and give an example. Finally we present some filtering rules for set CSPs.

Definition 2 (Set Interval). *let lb and ub be two sets such that $lb \subset ub$, the set interval $[lb..ub]$ is defined as follows: $[lb..ub] = \{E \text{ such that } lb \subseteq E \text{ and } E \subseteq ub\}$.*

Set intervals avoid data storage problems due to the size of domains: they model the domain and encapsulate all the possible values of the variables. For example: $[\{1\}.. \{1, 2, 3\}]$ summarizes $\{\{1\}, \{1, 2\}, \{1, 3\}, \{1, 2, 3\}\}$ and $[\{\}\{1, 2, 3\}]$ summarizes $2^{\{1, 2, 3\}}$.

Definition 3 (Set CSP). *A set constraint satisfaction problem (set CSP) is a 3-uple $(\mathcal{X}, \mathcal{D}, \mathcal{C})$ where $\mathcal{C} = \{c_1, \dots, c_m\}$ is a set of constraints associated to a set $\mathcal{X} = \{X_1, \dots, X_n\}$ of variables. For each variable X_i , an initial domain of set intervals (or union of set intervals) D_{X_i} is given and $D = \{D_{X_1}, \dots, D_{X_n}\}$.*

In order to illustrate the declarative feature and the expressiveness of set CSPs, we give the following example.

Example. [29] Two transmitters have to be assigned to two radio frequencies each. Available frequencies are $\{1, 2, 3, 4\}$ for the first transmitter and $\{3, 4, 5, 6\}$ for the second one. The distance between these two frequencies is equal to the absolute value of the difference between these frequencies. The constraints are:

- two radio frequencies have to be assigned to each transmitter: $c_1 \wedge c_2$.
- both transmitters do not share frequencies: c_3
- two frequencies within a transmitter must have at least a distance equals to 2: c_4
- the first transmitter requires the frequency 3: c_5
- the second transmitter requires the frequency 4: c_6

It can be expressed as a set CSP $(\mathcal{X}, \mathcal{D}, \mathcal{C})$, where:

- $\mathcal{X} = \{t_1, t_2\}$ where t_1 and t_2 are the two transmitters.
- $D(t_1) = [\{ \} .. \{1, 2, 3, 4\}]$ and $D(t_2) = [\{ \} .. \{3, 4, 5, 6\}]$.
- $\mathcal{C} = \{c_1, c_2, c_3, c_4, c_5, c_6\}$ where:
 - $c_1 \quad |t_1| = 2$
 - $c_2 \quad |t_2| = 2$
 - $c_3 \quad t_1 \cap t_2 = \emptyset$
 - $c_4 \quad \forall v_1, v_2 \in t_i, \quad |v_1 - v_2| \geq 2 \quad i = 1, 2$
 - $c_5 \quad 3 \in t_1$
 - $c_6 \quad 4 \in t_2$

This problem has a unique solution where the first transmitter is assigned to the frequencies $\{1, 3\}$ and the second to $\{4, 6\}$.

Examples of Filtering Rules for Set CSPs. For CSPs, filtering consists on reducing the variable domains in order to remove values that cannot occur in any solution. As soon as a domain D_{X_i} becomes empty (i.e., there is no available value for X_i), a failure is generated for the search. Filtering rules for integer intervals and set intervals are presented in [22][19][12]. We now present two examples of filtering rules for set intervals, the *inclusion* and the *intersection* constraints:

Let $D_x = [a_x .. b_x]$, $D_y = [a_y .. b_y]$ and $D_z = [a_z .. b_z]$ three domains represented by set intervals and D'_x, D'_y and D'_z the filtered domains.

- **Constraint:** $X \subset Y$
Filtering rule: **if** $a_x \subset b_y$ **then**
 $D'_x = [a_x .. b_x \cap b_y]$
 $D'_y = [a_x \cup a_y .. b_y]$
else
 $D'_x = \emptyset, D'_y = \emptyset$
- **Constraint:** $Z = X \cap Y$
Filtering rule: **if** $(b_x \cap b_y) \subset b_z$ and $(b_x \cap b_y) \neq \emptyset$ **then**
 $D'_x = [a_x \cup a_z .. b_x \setminus ((b_x \cap a_y) \setminus b_z)]$
 $D'_y = [a_y \cup a_z .. b_y \setminus ((b_y \cap a_x) \setminus b_z)]$
 $D'_z = [a_z \cup (a_x \cap a_y) .. b_z \cap b_x \cap b_y]$
else
 $D'_x = D'_y = D'_z = \emptyset$

Programming Tool: *ECLⁱPS^e*. [10] is a Constraint Programming Tool supporting the most common techniques used in solving constraints satisfaction (or optimization) problems: Constraint Satisfaction Problems, Mathematical Programming, Local Search and combinations of those. *ECLⁱPS^e* is built around the Constraint Logic Programming paradigm [1]. Different domains of constraints as numeric CSP and Set CSPs can be used together. Finally, libraries for solving set CSPs, as *ic-sets* or *conjunto* [12], are available in *ECLⁱPS^e*.

4 Set Constraint Programming for Pattern Discovery

Our approach is based on two major points. First, we use the wide possibilities of modelization and resolution given by the CSPs, in particular the set CSPs and numeric CSPs. Second, we take benefit from the recent progress on mining local patterns. The last choice is also strengthened by the fact that local constraints can be solved before and regardless n-ary constraints.

In this section, we start by giving an overview of our approach. Then we describe each of the three steps of our method by considering the example of the exception rules described in Section 2.2

4.1 General Overview

Figure 1 provides a general overview of the three steps of our approach:

1. Modeling the query as CSPs, then splitting constraints into local ones and n-ary ones.
2. Solving local constraints using a local pattern extractor (MUSIC-DFS, introduced in Section 3.1) which produces an interval condensed representation of all patterns satisfying the local constraints.
3. Solving n-ary constraints of the CSPs by using *ECLⁱPS^e* (introduced in Section 3.2) where the domain of each variable results from the interval condensed representation (computed in the Step-2).

4.2 Step-1: Modelling the Query as CSPs

Let \mathbf{r} be a dataset having nb transactions, and \mathcal{I} the set of all its items. We model the problem by using two CSPs \mathcal{P} and \mathcal{P}' that are inter-related:

1. Set CSP $\mathcal{P} = (\mathcal{X}, \mathcal{D}, \mathcal{C})$ where:
 - $\mathcal{X} = \{X_1, \dots, X_n\}$. Each variable X_i represents an unknown itemset.
 - $\mathcal{D} = \{D_{X_1}, \dots, D_{X_n}\}$. The initial domain of each variable X_i is the set interval $\{\} .. \mathcal{I}$.
 - \mathcal{C} is a conjunction of set constraints by using set operators ($\cup, \cap, \setminus, \in, \notin, \dots$)
2. Numeric CSP $\mathcal{P}' = (\mathcal{F}, \mathcal{D}', \mathcal{C}')$ where:
 - $\mathcal{F} = \{F_1, \dots, F_n\}$. Each variable F_i is the frequency of the itemset X_i .

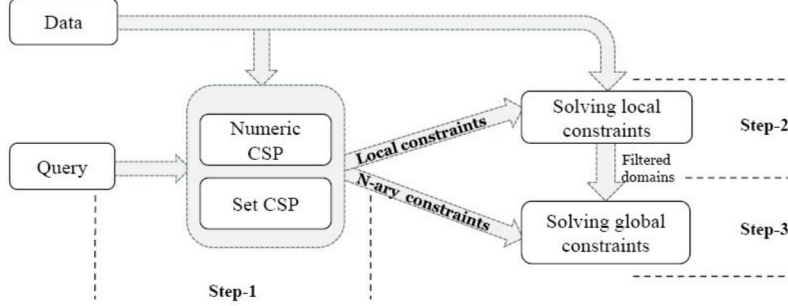


Fig. 1. General overview of our 3-steps method

- $\mathcal{D}' = \{D_{F_1}, \dots, D_{F_n}\}$. The initial domain of each variable F_i is the integer interval $[1 .. nb]$.
- \mathcal{C}' is a conjunction of arithmetic constraints.

Then, the whole set of constraints ($\mathcal{C} \cup \mathcal{C}'$) is divided into two subsets as follows:

- \mathcal{C}_{loc} is the set of local constraints to be solved (by MUSIC-DFS). Solutions are given in the form of an interval condensed representation.
- \mathcal{C}_{n-ary} is the set of n-ary constraints to be solved (by ECL^iPS^e), where the domain of the variables X_i and F_i will be deduced from the interval condensed representation computed in the previous step.

Local (unary) constraints can be solved before and regardless n-ary constraints. The search space of the n-ary constraints is reduced by the space of solutions satisfying local constraints. This ensures that every solution verifies both local and n-ary constraints.

4.3 Example: Modeling the Exception Rules as CSPs

Recall that the definition of the pairs of exception rules is given in Section [2.2](#)

Reformulation: Let $freq(X)$ be the frequency value of the itemset X . Let I and $\neg I \in \mathcal{I}$ (in this example, I and $\neg I$ represent the two class values of the data set). Let $\gamma_1, \gamma_2, \delta_1, \delta_2 \in \mathbb{N}$. The exception rules constraint can be formulated as it follows:

- $X \setminus Y \rightarrow I$ can be expressed by the conjunction: $freq((X \setminus Y) \sqcup I) \geq \gamma_1 \wedge (freq(X \setminus Y) - freq((X \setminus Y) \sqcup I)) \leq \delta_1$ which means that $X \setminus Y \rightarrow I$ must be a frequent rule having a high confidence value.

² The symbol \sqcup denotes the disjoint union operator.

- $X \rightarrow \neg I$ can be expressed by the conjunction: $freq(X \sqcup \neg I) \leq \gamma_2 \wedge (freq(X) - freq(X \sqcup \neg I)) \leq \delta_2$ which means that $X \rightarrow \neg I$ must be a rare rule having a high confidence value.

To sum up:

$$exception(X \rightarrow \neg I) \equiv \begin{cases} \exists Y \subset X \text{ such that:} \\ freq((X \setminus Y) \sqcup I) \geq \gamma_1 \wedge \\ (freq(X \setminus Y) - freq((X \setminus Y) \sqcup I)) \leq \delta_1 \wedge \\ freq(X \sqcup \neg I) \leq \gamma_2 \wedge \\ (freq(X) - freq(X \sqcup \neg I)) \leq \delta_2 \end{cases}$$

CSP Modelisation: The CSP variables are defined as follows:

- Set variables $\{X_1, X_2, X_3, X_4\}$ representing unknown itemsets:
 - $X_1 : X \setminus Y$,
 - $X_2 : (X \setminus Y) \sqcup I$ (common sense rule),
 - $X_3 : X$,
 - $X_4 : X \sqcup \neg I$ (exception rule).
- Integer variables $\{F_1, F_2, F_3, F_4\}$ representing their frequency values (variable F_i denotes the frequency of the itemset X_i).

Table 2 provides the constraints modeling the exception rules.

Table 2. Exception rules modeled as CSP constraints

Constraints	CSP formulation	Local	N-ary
$freq((X \setminus Y) \sqcup I) \geq \gamma_1$	$F_2 \geq \gamma_1$	×	
	$\wedge I \in X_2$	×	
$freq(X \setminus Y) - freq((X \setminus Y) \sqcup I) \leq \delta_1$	$\wedge X_1 \subsetneq X_3$		×
	$F_1 - F_2 \leq \delta_1$		×
$freq(X \sqcup \neg I) \leq \gamma_2$	$\wedge X_2 = X_1 \sqcup I$		×
	$F_4 \leq \gamma_2$	×	
$freq(X) - freq(X \sqcup \neg I) \leq \delta_2$	$\wedge \neg I \in X_4$	×	
	$F_3 - F_4 \leq \delta_2$		×
	$\wedge X_4 = X_3 \sqcup \neg I$		×

Summary:

- Set CSP
 - $\mathcal{X} = \{X_1, \dots, X_4\}$
 - $\mathcal{C} = \{(I \in X_2), (X_2 = X_1 \sqcup I), (\neg I \in X_4), (X_4 = X_3 \sqcup \neg I), (X_1 \subsetneq X_3)\}$
- Numeric CSP
 - $\mathcal{F} = \{F_1, \dots, F_4\}$
 - $\mathcal{C}' = \{(F_2 \geq \gamma_1), (F_1 - F_2 \leq \delta_1), (F_4 \leq \gamma_2), (F_3 - F_4 \leq \delta_2)\}$
- $\mathcal{C}_{loc} = \{(I \in X_2), (F_2 \geq \gamma_1), (F_4 \leq \gamma_2), (\neg I \in X_4)\}$
- $\mathcal{C}_{n-ary} = \{(F_1 - F_2 \leq \delta_1), (X_2 = X_1 \sqcup I), (F_3 - F_4 \leq \delta_2), (X_4 = X_3 \sqcup \neg I), (X_1 \subsetneq X_3)\}$

4.4 Step-2: Solving Local Constraints

As already said, we use for this task MUSIC-DFS (see Section 3.1) which mines soundly and completely local patterns. In order to fully benefit from the efficiency of the local pattern mining, the set of local constraints \mathcal{C}_{loc} is split into a disjoint union of \mathcal{C}_i (for $i \in [1..n]$) where each \mathcal{C}_i is the set of local constraints related to X_i and F_i . Each \mathcal{C}_i can be separately solved. Let CR_i be the interval condensed representation of all the solutions of \mathcal{C}_i . $CR_i = \bigcup_p (f_p, I_p)$ where I_p is a set interval verifying: $\forall x \in I_p, freq(x) = f_p$. Then the filtered domains (see Section 4.3) for variable X_i and variable F_i are:

- D_{F_i} : the set of all f_p in CR_i
- D_{X_i} : $\bigcup_{I_p \in CR_i} I_p$

Example. Let us consider the dataset r (see Table 1) and the local constraints for the exception rules $\mathcal{C}_{loc} = \{(I \in X_2), (F_2 \geq \gamma_1), (F_4 \leq \gamma_2), (\neg I \in X_4)\}$ (see Section 4.3). The respective values for $(I, \neg I, \gamma_1, \delta_1, \gamma_2, \delta_2)$ are $(c_1, c_2, 2, 1, 1, 0)$. The local constraints set related to X_2 is $\mathcal{C}_2 = \{c_1 \in X_2, F_2 \geq 2\}$ is solved by MUSIC-DFS with the following query showing that the parameters given to MUSIC-DFS are straightforwardly deduced from \mathcal{C}_{loc} .

```
-----
./music-dfs -i donn.bin -q "{c1} subset X2 and freq(X2)>=2;"
X2 in [A, c1]..[A, c1, B ] U [B, c1] -- F2 = 2 ;
X2 in [C, c1] -- F2 = 3
-----
```

4.5 Step-3: Solving n-ary Constraints

Then, from the condensed representation of all patterns satisfying local constraints, domains of the variables X_i and F_i (for $i \in \{1, 2, 3, 4\}$) are updated.

Given the parameters $I = c_1, \neg I = c_2, \delta_1 = 1$ and $\delta_2 = 0$ ($\gamma_1 = 2$ and $\gamma_2 = 1$ are already used in Step-2) and the data set in Table 1 the following *ECLⁱPS^e* session illustrates how all pairs of exception rules can be obtained by using backtracking:

```
-----
[eclipse 1]:
?- exceptions(X1, X2, X3, X4).
Sol1 : X1 = [A,B], X2=[A,B,c1], X3=[A,B,C], X4=[A,B,C,c2];
Sol2 : X1 = [A,B], X2=[A,B,c1], X3=[A,B,D], X4=[A,B,D,c2];
.../...
-----
```

5 Experiments

This section shows the practical usage and the feasibility of our approach. This experimental study is conducted on the *postoperative-patient-data* coming from

the UCI machine learning repository³. This data set gathers 90 objects described by 23 items and characterized by two classes (two objects of a third class value were put aside). We test our approach by using the exception rules as a n -ary constraint (in the following, we use a class value for the item I given in the definition of an exception rule). As previously said, we use MUSIC-DFS (see Section 3.1) and ECL^iPS^e (see Section 3.2). All the tests were performed on a 2 GHz Intel Centrino Duo processor with Linux operating system and 2GB of RAM memory.

These experiments show the feasibility of our approach. Given $(I, \gamma_1, \delta_1, \gamma_2, \delta_2)$ a set of values, our method is able to mine the correct and complete set of all pairs of exception rules.

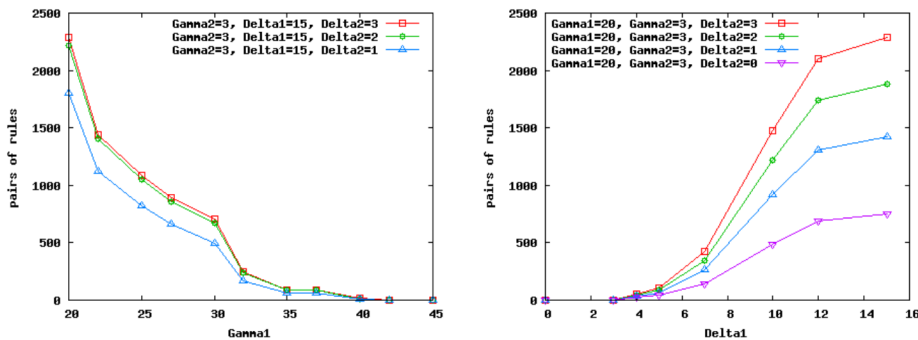


Fig. 2. Number of rules according to γ_1 (left) and δ_1 (right)

Figure 2 depicts the number of pairs of rules according to γ_1 (left part of the figure) and δ_1 (right part of the figure). We tested several combinations of the parameters. As expected, the lower γ_1 is, the larger the number of pairs of exception rules. Note that the decreasing of the curves is approximately the same for all the combinations of parameters. The result is similar when δ_1 varies (right part of Figure 2): the higher δ_1 is, the larger the number of pairs of exception rules (when δ_1 increases, the confidence decreases so that there are more common sense rules). Interestingly, these curves quantify the number of pairs of exception rules according to the sets of parameters. Some cases seem to point out pairs of rules of good quality. For instance, with $(\gamma_1 = 20, \delta_1 = 5, \gamma_2 = 1, \delta_2 = 0)$, we obtain 25 pairs of rules with a common sense rule having a confidence value greater than or equal to 83% and an exact exception rule (i.e., confidence value equals 100%). Moreover, our approach enables us in a natural way to add new properties such as the control of the sizes of rules. If the user wants that the number of items added to an exception rule remains small with regards to the size of the common sense rule, it can be easily modeled by a new

³ www.ics.uci.edu/~mlearn/MLRepository.html

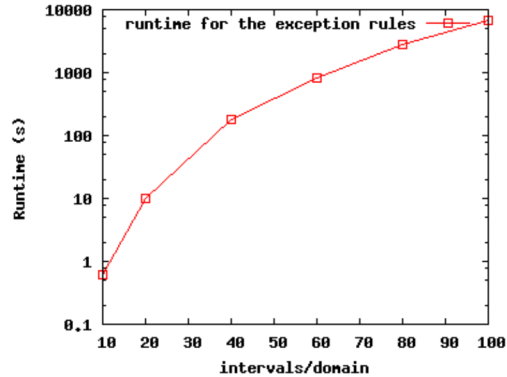


Fig. 3. Runtime according to the number of intervals of the condensed representations

constraint: for instance, the number of added items to an exception rule must be lower than the minimum of a number (e.g., 3) and the size of the common sense rule. It highlights the flexibility of our approach.

Figure 3 details the runtime of our method according to the number of intervals of the condensed representation, i.e., the size of the condensed representation. In this experiment, for each dot of the curve, the four variables have the same domain and thus the same number of intervals. Obviously, the larger the number of intervals is, the higher the runtime (note that we use a logarithmic scale on the Y axis). In the case of exception rules, it is interesting to note that the runtime decreases when the quality of the exception rule pairs increases. Indeed, looking for common sense rules with high frequency and reliable exception rules leads to infer local constraints giving more powerful pruning conditions and thus less intervals.

Table 3 indicates the number of intervals of the variable X_2 in the condensed representation (see Section 4.3) according to several local constraints. It shows the interest of an approach based on local constraint mining.

Discussion. We briefly discuss the set union operator for set CSPs which is a key point in our approach. In order to perform bound consistency filtering, set CSP solvers approximate the union of two set intervals by their convex closure. The convex closure of $[lb_1 .. ub_1]$ and $[lb_2 .. ub_2]$ is defined as $[lb_1 \cap lb_2 .. ub_1 \cup ub_2]$. So, if filtering is applied a lot of times on a same variable domain, this domain may reach the whole set $[\emptyset .. \mathcal{I}]$ and specific information gathered during the search would be lost whereas this information is useful to limit the size of intervals. To circumvent this problem, for each variable X_i with the condensed representation $CR_i = \bigcup_p (f_p, I_p)$, a search is successively performed upon each I_p . This approach is sound and complete and we use it in our experiments. Nevertheless, with this method, we do not fully profit from filtering because removing a value is propagated only in the treated intervals and not in the whole domains. It explains the results of Section 5 showing that the runtime strongly increases when the number

Table 3. Number of intervals according to several local constraints (case of D_{X_2})

Local constraint	Number of intervals in D_{X_2}
-	3002
$I \in X_2$	1029
$I \in X_2 \wedge freq(X_2) \geq 20$	52
$I \in X_2 \wedge freq(X_2) \geq 25$	32

of intervals increases. Alternative solutions consist of implementing a set interval union operator in the kernel of the solver or using non-exact condensed representations to reduce the number of produced intervals (e.g., a condensed representation based on maximal frequent itemsets). In this case, the number of intervals representing the domains will be smaller, but, due to the approximations, it should be necessary to memorize forbidden values.

6 Conclusion and Future Work

In this paper we have presented a new approach for pattern discovery. Its great interest is to model in a flexible way any set of constraints combining several local patterns. The complete and sound set of patterns satisfying the constraints is mined thanks to a joint cooperation between a solver on set constraint programming which copes with n-ary constraints and a solver on local patterns to take benefit on the well-mastered methods on local constraint mining. We think that it is this combination between the local and n-ary levels which enables us the discovery of such patterns. Experiments show the feasibility of our approach.

In classic CSPs, all variables are existentially quantified. Further work is to introduce the universal quantification (\forall): this quantifier would be precious to model important constraints such as the peak constraint (the peak constraint compares neighbor patterns and a peak pattern is a pattern whose all neighbors have a value for a measure lower than a threshold). For that purpose, we think that recent works as Quantified Constraints Satisfaction Problems (QCSP) [220] could be useful.

Acknowledgments. The authors would like to thank Arnaud Soulet for very fruitful discussions and MUSIC-DFS prototype. This work is partly supported by the ANR (French Research National Agency) funded project Bingo2 ANR-07-MDCO-014.

References

1. Apt, K.R., Wallace, M.: Constraint Logic Programming using Eclipse. Cambridge University Press, New York (2007)
2. Benhamou, F., Goualard, F.: Universally quantified interval constraints. In: Dechter, R. (ed.) CP 2000. LNCS, vol. 1894, pp. 67–82. Springer, Heidelberg (2000)

3. Besson, J., Robardet, C., Boulicaut, J.-F.: Mining a new fault-tolerant pattern type as an alternative to formal concept discovery. In: 14th International Conference on Conceptual Structures (ICCS 2006), Aalborg, Denmark, pp. 144–157. Springer, Heidelberg (2006)
4. Bringmann, B., Zimmermann, A.: The chosen few: On identifying valuable patterns. In: Proceedings of the 12th IEEE International Conference on Data Mining (ICDM 2007), Omaha, NE, pp. 63–72 (2007)
5. Calders, T., Rigotti, C., Boulicaut, J.-F.: A survey on condensed representations for frequent sets. In: Boulicaut, J.-F., De Raedt, L., Mannila, H. (eds.) Constraint-Based Mining and Inductive Databases. LNCS (LNAI), vol. 3848, pp. 64–80. Springer, Heidelberg (2006)
6. Crémilleux, B., Soulet, A.: Discovering knowledge from local patterns with global constraints. In: Gervasi, O., Murgante, B., Laganà, A., Tanar, D., Mun, Y., Gavrilova, M.L. (eds.) ICCSA 2008, Part II. LNCS, vol. 5073, pp. 1242–1257. Springer, Heidelberg (2008)
7. De Raedt, L., Guns, T., Nijssen, S.: Constraint Programming for Itemset Mining. In: ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 14th edn., Las Vegas, Nevada, USA (2008)
8. De Raedt, L.: A theory of inductive query answering. In: Proceedings of the IEEE Conference on Data Mining (ICDM 2002), Maebashi, Japan, 2002, pp. 123–130 (2002)
9. De Raedt, L., Zimmermann, A.: Constraint-based pattern set mining. In: Proceedings of the Seventh SIAM International Conference on Data Mining, Minneapolis, Minnesota, USA, April 2007, SIAM (2007)
10. ECLIPSe. Eclipse documentation, <http://www.eclipse-clp.org>
11. Gecode Team. Gecode: Generic constraint development environment (2006), <http://www.gecode.org>
12. Gervet, C.: Interval Propagation to Reason about Sets: Definition and Implementation of a Practical Language. *Constraints* 1(3), 191–244 (1997)
13. Giacometti, A., Miyaneh, E.K., Marcel, P., Soulet, A.: A framework for pattern-based global models. In: 10th Int. Conf. on Intelligent Data Engineering and Automated Learning, Burgos, Spain, pp. 433–440 (2009)
14. Hand, D.J.: ESF exploratory workshop on Pattern Detection and Discovery in Data Mining. In: Hand, D.J., Adams, N.M., Bolton, R.J. (eds.) *Pattern Detection and Discovery*. LNCS (LNAI), vol. 2447, pp. 1–12. Springer, Heidelberg (2002)
15. Kléma, J., Blachon, S., Soulet, A., Crémilleux, B., Gandrillon, O.: Constraint-based knowledge discovery from sage data. *Silico Biology* 8(0014) (2008)
16. Knobbe, A.: From local patterns to global models: The lego approach to data mining. In: International Workshop From Local Patterns to Global Models co-located with ECML/PKDD 2008, Antwerp, Belgium, September 2008, pp. 1–16 (2008)
17. Knobbe, A., Ho, E.: Pattern teams. In: Fürnkranz, J., Scheffer, T., Spiliopoulou, M. (eds.) *PKDD 2006*. LNCS (LNAI), vol. 4213, pp. 577–584. Springer, Heidelberg (2006)
18. Lakshmanan, L.V., Ng, R., Hah, J., Pang, A.: Optimization of constrained frequent set queries with 2-variable constraints (1998)
19. Lhomme, O.: Consistency techniques for numeric csp. In: Proc. of the 13th IJCAI, Chambéry, France, pp. 232–238 (1993)
20. Mamoulis, N., Stergiou, K.: Algorithms for quantified constraint satisfaction problems. In: Wallace, M. (ed.) *CP 2004*. LNCS, vol. 3258, pp. 752–756. Springer, Heidelberg (2004)

21. Mannila, H., Toivonen, H.: Levelwise search and borders of theories in knowledge discovery. *Data Mining and Knowledge Discovery* 1(3), 241–258 (1997)
22. Moore, R.E.: *Interval analysis*. Prentice-Hall, Englewood Cliffs (1966)
23. Ng, R.T., Lakshmanan, V.S., Han, J., Pang, A.: Exploratory mining and pruning optimizations of constrained associations rules. In: *Proceedings of ACM SIGMOD 1998*, pp. 13–24. ACM Press, New York (1998)
24. Nijssen, S., Guns, T., De Raedt, L.: Correlated itemset mining in roc space: a constraint programming approach. In: *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD 2009)*, Paris, France, June 2009, pp. 647–655 (2009)
25. Siebes, A., Vreeken, J., van Leeuwen, M.: Item sets that compress. In: *Proceedings of the Sixth SIAM International Conference on Data Mining*, Bethesda, MD, USA, April 2006, SIAM, Philadelphia (2006)
26. Soulet, A., Crémilleux, B.: An efficient framework for mining flexible constraints. In: Ho, T.-B., Cheung, D., Liu, H. (eds.) *PAKDD 2005*. LNCS (LNAI), vol. 3518, pp. 661–671. Springer, Heidelberg (2005)
27. Soulet, A., Klema, J., Crémilleux, B.: Efficient Mining Under Rich Constraints Derived from Various Datasets. In: Dzeroski, S., Struyf, J. (eds.) *KDID 2006*. LNCS, vol. 4747, pp. 223–239. Springer, Heidelberg (2007)
28. Suzuki, E.: Undirected Discovery of Interesting Exception Rules. *International Journal of Pattern Recognition and Artificial Intelligence* 16(8), 1065–1086 (2002)
29. Thornary, V., Gensel, J., Sherpa, P.: An hybrid representation for set constraint satisfaction problems. In: *Workshop on Set Constraints co-located with the fourth Int. Conf. on Principles and Practice of Constraint Programming*, Pisa, Italy (1998)
30. Yin, X., Han, J.: CPAR: classification based on predictive association rules. In: *proceedings of the 2003 SIAM Int. Conf. on Data Mining (SDM 2003)*, San Fransisco, CA (May 2003)