



# Using Answer Set Programming for pattern mining

Thomas Guyet, Yves Moinard, René Quiniou

## ► To cite this version:

Thomas Guyet, Yves Moinard, René Quiniou. Using Answer Set Programming for pattern mining. Intelligence Artificielle Fondamentale, Jun 2014, Angers, France. hal-01069092

**HAL Id: hal-01069092**

**<https://hal.inria.fr/hal-01069092>**

Submitted on 26 Sep 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.

# Using Answer Set Programming for pattern mining

Thomas Guyet<sup>1</sup>, Yves Moinard<sup>2</sup> and René Quiniou<sup>2</sup>

<sup>1</sup> AGROCAMPUS-OUEST/IRISA UMR6074, 35042 Rennes, France

<sup>2</sup> INRIA, Campus de Beaulieu, 35042 Rennes, France

## Abstract

Serial pattern mining consists in extracting the frequent sequential patterns from a unique sequence of itemsets. This paper explores the ability of a declarative language, such as Answer Set Programming (ASP), to solve this issue efficiently. We propose several ASP implementations of the frequent sequential pattern mining task: a non-incremental and an incremental resolution. The results show that the incremental resolution is more efficient than the non-incremental one, but both ASP programs are less efficient than dedicated algorithms. Nonetheless, this approach can be seen as a first step toward a generic framework for sequential pattern mining with constraints.

## 1 Introduction

Sequential pattern mining aims at analysing ordered or timed data to extract interesting patterns the elements. Broadly, a pattern is considered interesting if it occurs frequently in the data, *i.e.* the number of its occurrences is greater than a fixed given threshold.

As non informed mining methods tend to generate massive results, there is more and more interest in pattern mining algorithms able to mine data considering some expert knowledge. Though a generic pattern mining tool that could be tailored to the specific task of a data-scientist is still a holy grail for pattern mining software designers, some recent attempts have proposed generic pattern mining tools [7, 12] for itemset mining tasks.

One way to introduce expert knowledge in pattern mining algorithms is to define constraints on the expected patterns. Motivated by the use of constraints in sequential pattern mining. Pei *et al.* [14, 15] defined seven types of generic constraints that a data scientist would like to express. However, beyond generic

---

\*This work first appeared in the proceedings of the IAF conference (French Fundamental Artificial Intelligence Conference) in June 2014.

constraints, a data scientist would like to express a lot of domain-specific constraints reflecting his own knowledge and preferences.

ASP (Answer Set Programming) is a declarative language that enables an easy expression of constraints as well as domain knowledge and associated formal reasoning tools. To engage in the generic way sketched before, this paper proposes to implement a classical sequential pattern mining algorithm in ASP with the future objective of letting the data scientist define constraints on its patterns and associated domain knowledge.

We show how to encode sequential pattern mining tasks in ASP. The first intuitive and naive encoding is next improved with the incremental resolution facilities of the `clingo` ASP solver [6]. The performance of the method is evaluated on simulated data and compared with the dedicated algorithm MinEpi [10].

## 2 State of the art

Designing pattern mining languages based on pattern constraints has developed recent interest in the literature [7, 12, 2, 3]. The aim of these proposals is to obtain a *declarative* constraint-based language even at the cost of degraded runtime performance compared to a specialized algorithm.

All these approaches have been conducted on itemset mining in transaction databases, which is much simpler than sequential pattern mining in a sequence database, which in turn is simpler than sequential pattern mining in a unique long sequence.

Sequential pattern mining in a sequence database have been addressed by numerous algorithms inspired by algorithms for mining frequent itemsets. The most known algorithms are GSP [17], SPIRIT [5], SPADE [22], PrefixSpan [13], and CloSpan [21] or BIDE [20] for closed sequential patterns. To determine the support of a pattern, these algorithms count the number of sequences in the database that contain the pattern without taking into account repetitions of the pattern in the sequences. Deciding if a sequence contains a pattern is generally simple. It is worth-noting that all these algorithms satisfy the anti-monotonicity property which is essential to obtain good mining performances. The anti-monotonicity property says that if some pattern is frequent then all its sub-patterns are also frequent. And reciprocally, if some pattern is not-frequent then all its super-patterns are non-frequent. This property enables the algorithm to prune efficiently the search space and thus reduces its exploration.

Counting the number of occurrences of a pattern, or of patterns, in a long sequence introduces some complexity compared to previous approaches as two occurrences of a pattern can overlap. Defining how to enumerate the occurrences of a pattern is important to ensure the anti-monotonicity property. The MinEpi and WinEpi algorithms [10] extract frequent *episodes*. An episode is a directed acyclic graph where edges express the temporal relations between events. A *proper edge* from node  $v$  to node  $w$  in the pattern graph implies that the events of  $v$  must occur before the events of  $w$ , while a *weak edge* from  $v$  to  $w$  implies

that the events of  $w$  may occur either at the same time as those of  $v$  or later. Testing whether an episode occurs in the sequence is an NP-complete problem [18].

The frequency measure used by MinEpi relies on minimal occurrences of patterns and is anti-monotonic. Other counting methods have been proposed to solve similar problems while preserving the anti-monotonicity property required for the effectiveness of pattern occurrences search (see [1] for a unified formulation of sequential pattern counting methods).

Integrating constraints in sequential pattern mining is often limited to the use of simple anti-monotonic temporal constraints such as **mingap** constraints. In addition to the classical support constraint, Pei et al. [14] defined seven types of constraints on patterns. SPIRIT [5] is one of the rare algorithm that considers complex constraints on patterns such as regular expressions.

In [8], Guns et al. proposed a new perspective on data mining based on constraint programming. Their aim is to separate modelling and solving in data mining in order to gain more flexibility for elaborating new mining models. They chose to illustrate the approach on itemset mining. Only few recent works have proposed to explore constraint programming for sequential pattern mining. Métivier, Loudni et al. [11] have developed a constraint programming method for mining sequential patterns with constraints in a sequence database. The constraints are based on *amongst* and *regular expression* constraints and expressed by automata. Coquery et al. [4] have proposed a SAT based approach for sequential pattern mining. The patterns are of the form  $ab?c$  and an occurrence corresponds to an exact substring (without gap) with joker (the character  $?$  replaces exactly one item different from  $b$  and  $c$ ).

To the best of our knowledge, Järvisalo’s work [9] is the unique application of ASP to pattern mining. Following Guns et al.’s proposal, Järvisalo designed an ASP program to extract frequent itemsets in a transaction database. A main feature of Järvisalo is that each answer set (AS) contains only one frequent itemset associated with the identifiers of the transactions where it occurs. Järvisalo addressed this problem as a new challenge for the ASP solver, but did not highlight the potential benefit of this approach to improve the expressiveness of pattern mining tools.

In this paper, we explore the use of ASP to extract frequent string patterns (*i.e.* patterns without parallel events) in a unique long sequence of itemsets where the occurrences of a string pattern are the minimal occurrences. In addition, we implemented the following simple anti-monotonic constraints to improve the algorithm efficiency:

1. an item constraint states which particular items should not be present in the patterns.
2. a length constraint fixes the maximal length of patterns (number of items).

3. a duration constraint sets the maximal duration of patterns (time elapsed between the occurrence of the first itemset and the occurrence of the last one - difference between timestamps).
4. a gap constraint sets the maximal time delay between the respective timestamps of two successive itemsets.

### 3 Sequential pattern mining

#### 3.1 Items, itemsets and sequences

From now on,  $[n]$  denotes the set of the  $n$  first integers, *i.e.*  $[n] = \{1, \dots, n\}$ .

Let  $\mathcal{E}$  be the set of items and  $\preceq$  a total (reflexive) order on this set (*e.g.* lexicographical). An *itemset*  $A = (a^1, a^2, \dots, a^n)$ ,  $a^i \in \mathcal{E}$  is an ordered set of distinct items, *i.e.*  $\forall i \in [n-1]$ ,  $a^i \preceq a^{i+1}$  and  $i \neq j \Rightarrow a^i \neq a^j$ . The size of an itemset  $\alpha$ , denoted  $|\alpha|$  is the number of items it contains. An itemset  $\beta = (b^1, \dots, b^m)$  is a sub-itemset of  $\alpha = (a^1, \dots, a^n)$ , denoted  $\beta \sqsubseteq \alpha$ , iff  $\beta$  is a subset of  $\alpha$ .

A *sequence*  $S = \langle s_1, s_2, \dots, s_n \rangle$  is an ordered series of itemsets  $s_i$ . The length of a sequence  $S$ , denoted  $|S|$ , is the number of itemsets that make up the sequence. The size of a sequence  $S$ , denoted  $\|S\|$ , is the total number of items it contains:  $\|S\| = \sum_{i=1}^{|S|} |s_i|$ .

$T = \langle t_1, t_2, \dots, t_m \rangle$  is a *sub-sequence* of  $S = \langle s_1, s_2, \dots, s_n \rangle$ , denoted  $T \preceq S$ , iff there exists a sequence of integers  $1 \leq i_1 < i_2 < \dots < i_m \leq n$  such that  $\forall k \in [m]$ ,  $t_k \sqsubseteq s_{i_k}$ .  $T = \langle t_1, t_2, \dots, t_m \rangle$  is a *prefix* of  $S = \langle s_1, s_2, \dots, s_n \rangle$ , denoted  $T \preceq_b S$ <sup>1</sup>, iff  $\forall k \in [m-1]$ ,  $t_k = s_k$  and  $t_m \sqsubseteq s_m$ ,  $t_m \neq s_m$ .

**Example 1.** Let  $\mathcal{E} = \{a, b, c\}$  with the lexicographical order ( $a \preceq b$ ,  $b \preceq c$ ) and the sequence  $S = \langle a(bc)c(abc)cb \rangle$ . To simplify the notation, we omit the parentheses around itemsets containing a single item, and commas inside itemsets. The size of  $S$  is 9 and its length is 6. For instance, sequence  $\langle (bc)(ac) \rangle$  is a sub-sequence of  $S$  and  $\langle a(bc)c(ac) \rangle$  is a prefix of  $S$ .

The relations  $\preceq$  and  $\preceq_b$  define partial orders on sets of sequences.

#### 3.2 Mining string patterns in a long sequence

A *long sequence*  $F = \{f_i\}_{i \in \mathbb{N}}$  is a single sequence of itemsets. A **string pattern**  $P = \langle p_1, \dots, p_n \rangle$  is a sequence of items ( $p_k$ ) (all itemsets are singleton). In the sequel, we use the term **pattern** to refer to our notion of string pattern. Our notion of pattern is more specific than the one defined by Mannila et al. [10].

**Definition 1** (Occurrences of a pattern, minimal occurrences). *An occurrence of the pattern  $P = \langle p_1, \dots, p_n \rangle$  of length  $n$  in a long sequence  $S = (s_1, \dots, s_m)$  is noted by the  $n$ -tuple  $T = (t_1, \dots, t_n)$ , where  $t_k, k \in [n]$  is the position of the occurrence of element  $p_k \in s_{t_k}$  in  $S$ .  $P$  is said to occur at interval  $[t_1, t_n]$  in  $S$ .*

<sup>1</sup>The  $b$  in  $\preceq_b$  stands for backward subsequence.

Let  $T = (t_1, \dots, t_n)$  and  $T' = (t'_1, \dots, t'_n)$ ,

$$T \triangleleft T' \Leftrightarrow \begin{cases} [t_1, t_n] \subset [t'_1, t'_n], \text{ or} \\ n > 1 \wedge t_1 = t'_1 \wedge t_n = t'_n \wedge \\ (t_1, \dots, t_{n-1}) \triangleleft (t'_1, \dots, t'_{n-1}) \end{cases}$$

An occurrence  $T$  of a pattern  $P$  is minimal if it is minimal for  $\triangleleft$ .  $\mathcal{I}_S(P)$  denotes the set of all minimal occurrences of  $P$  in  $S$ .

**Example 2.** Let  $S = \langle a(bc)(ac)cd \rangle$  be a long sequence and  $P = \langle abc \rangle$  be a pattern. Finding occurrences of pattern  $P$  in  $S$  means to locate items of  $P$  in  $S$ :  $\langle a \rangle$  appears only at positions 1 and 3,  $\langle b \rangle$  appears at position 2,  $\langle c \rangle$  appears at positions 2, 3 and 4. Pattern  $P$  could have two occurrences:  $(1, 2, 3)$  and  $(1, 2, 4)$ . The minimal occurrence condition eliminates occurrence  $(1, 2, 4)$  because  $[1, 2]$  is (strictly) included in  $[1, 4]$ . Thus  $\mathcal{I}_S(\langle abc \rangle) = \{(1, 2, 3)\}$ .

We now consider the pattern  $P = \langle acd \rangle$  in order to illustrate the second (recursive) minimal occurrence condition. The occurrences of the pattern would be  $\{(1, 2, 5), (1, 3, 5)\}$ . The occurrence  $(1, 3, 5)$  is not minimal because the occurrences bounds are equivalents but  $[1, 2]$  is (strictly) included in  $[1, 3]$ , thus  $(1, 2, 5) \triangleleft (1, 3, 5)$ .

Let  $S = \langle aaaa \rangle$  be the long sequence.  $\mathcal{I}_S(\langle aa \rangle) = \{(1, 2), (2, 3), (3, 4)\}$  and  $\mathcal{I}_S(\langle aaa \rangle) = \{(1, 3), (2, 4)\}$ .

Let  $S$  be a long sequence and  $P$  a pattern. The **support** of pattern  $P$ , denoted by  $\text{supp}(P)$ , is the cardinality of  $\mathcal{I}_S(P)$ , i.e.  $\text{supp}(P) = \text{card}(\mathcal{I}_S(P))$ .

Note that the support function  $\text{supp}(\cdot)$  is not anti-monotonic on the set of patterns with associated partial order  $\preceq$  (see [19]). Considering order  $\preceq_b$ ,  $\text{supp}(\cdot)$  is anti-monotonic.

**Definition 2** (Mining a long sequence). Given a threshold  $\sigma$ , a pattern  $P$  is frequent in a long sequence  $S$  iff  $\text{supp}(P) \geq \sigma$ . Mining a long sequence consists in extracting all frequent patterns.

**Example 3.** Let  $\sigma = 2$ . The set of frequent patterns in  $S = \langle a(bc)(abc)c(bc) \rangle$  is  $\{ \langle a \rangle, \langle b \rangle, \langle c \rangle, \langle ab \rangle, \langle ac \rangle, \langle bb \rangle, \langle bc \rangle, \langle cb \rangle, \langle cc \rangle, \langle acb \rangle, \langle acc \rangle, \langle bcc \rangle, \langle ccc \rangle \}$ .

## 4 Mining serial patterns with ASP

This section shows a first way of specifying serial pattern mining in classical ASP programming style. After presenting how to model the data in ASP, we give an intuitive algorithm, then we show how to improve its resolution using `clingo` control abilities.

Our proposal is borrowed from Järvisalo's [9]: a solution of the ASP program is an answer set (AS) that contains a single frequent pattern as well as its occurrences. The resolution relies on the “generate and test principle”: generate combinatorially all the possible patterns, one at a time, associated with their minimal occurrences and test whether they satisfy the specified constraints.

## 4.1 Modelling the long sequence, patterns and occurrences

A long sequence is modelled by the predicate **seq**/2. For example, the atom **seq**(3,5) declares that an item 5 occurs at timestamp 3. Similarly, the current pattern is modelled by the predicate **pattern**/2.

**Example 4.** Let  $S = \langle aca(ac)bc \rangle$  be a long sequence and  $P = \langle acc \rangle$  be a pattern.

```

1 %sequence description : a=1, b=2, c=3
2 seq(1,1). seq(2,3). seq(3,1). seq(4,1). seq(4,3). seq(5,2). seq(6,3).
3 %pattern description
4 pattern(1,1). pattern(2,3). pattern(3,3).
```

An occurrence of some pattern  $P$  is described by a set of atoms **occ**(**I**,**P**,**S**) where **I** is the identifier of the occurrence, (an AS describes all the occurrences of its pattern), **P** is the item position in the pattern and **S** is the timestamp of a matching itemset in the sequence.

**Example 5.** Continuing example above, the occurrences of pattern  $P$  are  $\mathcal{I}_S(P) = \{(1, 2, 4), (3, 4, 6)\}$ . The AS for the program above should contain:

```

%occurrences
occ(1,1,1). occ(1,2,2). occ(1,3,4).
occ(2,1,3). occ(2,2,4). occ(2,3,6).
```

## 4.2 An ASP program for extracting frequent patterns of length $n$

In this section, we detail an ASP program that generates all the frequent patterns of length  $n$ . According to the ASP programming principle, the program contains a generation part and a constraint part. The generation part specifies the pattern structure as a set of **pattern** atoms and their related occurrences as a set of **occ** atoms. The constraints eliminate pattern candidates that are not frequent.

First, we introduce some additional predicates and constants:

- constant **th** represents the minimum frequency threshold,
- predicate **patlen**/1 sets the length of patterns,
- predicate **symb**/1 specifies the items that can appear in a pattern.

### 4.2.1 Generating patterns and occurrences

The program below gives the rules for generating occurrences of patterns and their occurrences. The rule in line 6 generates patterns containing exactly  $L$  **pattern** atoms, expressed combinatorially on the vocabulary specified by **symb**/1 atoms. As only minimal models are solutions of ASP programs, every AS will contain **L** **pattern** atoms defining only one pattern.

Next, at line 9, predicate **instid**/1 is used to enumerate exactly **th** occurrences of the specified pattern. All the occurrences of a pattern could be generated, but only **th** are kept in the solution, *i.e.* the minimal number of occurrences for a pattern to be frequent.

Finally, at line 12, for each of the **th** pattern occurrences, a set of exactly **L occ** atoms is generated. In this combinatorial step, every possible association between a pattern item and a sequence timestamp **P** is blindly considered.

```

5 %generating all the patterns
6 L { pattern(1..L, P) : symb(P) } L :- patlen(L).
7
8 %list of occurrences
9 instid(1..th).
10
11 % I: occurrence id, L: pattern length, P: sequence position
12 L{ occ(I, 1..L, P): seq(P, S) }L :- instid(I), patlen(L).

```

#### 4.2.2 Minimal occurrences constraints

The loose specification above generates many solutions, some being correct but redundant, *e.g.* with different occurrence identifiers, some being incorrect, *e.g.* with non minimal occurrences. Now, we add some constraints to keep only the required models.

The first constraint expresses that the item at position **N** in the pattern must correspond to one of the items at timestamp **P** in the given sequence. Constraints are expressed by negations. The following constraint requires that for any occurrence atom **occ(I,N,P)** of a pattern having item **S** at position **N**, it is impossible to not have the same item **S** at timestamp **P** in the sequence.

```

13 :- occ(I, N, P), pattern(N, S), not seq(P, S).

```

The next constraint expresses that the occurrences must respect the ordering of pattern items: it is not possible to have 2 items ordered in some way in the pattern and the 2 items they are mapped with ordered in the other way in the sequence.

```

14 :- occ(I, N, P), occ(I, M, Q), N<M, P>=Q.

```

The next constraint imposes that the *N*-th item of the pattern maps with a single timestamp in the sequence.

```

15 :- occ(I, N, P), occ(I, N, Q), P<Q.

```

Notice that here **P < Q** is equivalent to **P!=Q** (meaning  $P \neq Q$  in ASP) since **P** and **Q** play symmetric roles in the two atoms in **occ**. This “trick” is often used in ASP (*cf.* also lines 16 and 17 in next listing), since it has the advantage of reducing the size of the grounding of the program.

The final constraints are related to minimal occurrences. The first constraint forbids solutions with two occurrences starting at the same timestamp. The second constraint forbids solutions with two occurrences ending at the same timestamp. The third constraint forbids solutions with two occurrences such that one “contains” the other.



```

16 :- occ(I, 1, P), occ(J, 1, Q), P=Q, I<J.
17 :- occ(I, L, P), occ(J, L, Q), patlen(L), P=Q, I<J.
18 :- occ(I, 1, SI), occ(J, 1, SJ), occ(I, L, EI), occ(J, L, EJ), patlen(L), SI<SJ, EJ<EI, I
    !=J.

```

Constraints 16 to 18 could be summed up by the single constraint below. However, the single constraint is significantly less computationally efficient than the three elementary constraints.

```

:- occ(I, 1, SI), occ(J, 1, SJ), occ(I, L, EI), occ(J, L, EJ), patlen(L), SI<=SJ, EJ<=
    EI, I!=J.

```

#### 4.2.3 Reducing the combinatorics of the solution space

At this step, all the ASs are correct answers to the sequence mining problem, but many of them represents the same solution due to the combinatorial nature of the generation step. For instance, extracting patterns of length 4 in a sequence of 20 items yields 2448 answer sets, but there are only 24 different answer sets<sup>2</sup>.

To improve the program efficiency, we add the following constraint on occurrence identifiers imposing that occurrences are ordered by their beginning timestamps. This constraint reduces the number of answer sets to 408. It is better but not perfect.

```

19 :- occ(I, 1, P), occ(J, 1, Q), I<J, P>Q.

```

#### 4.3 Illustration of constraint effectiveness

Let  $W = \langle abab \rangle$  be a sequence and  $th = 2$  the support threshold. We consider only the models generated with **pattern**(1,1) and **pattern**(2,2), *i.e.* that evaluates the occurrences of  $\langle ab \rangle$ . The models that contain at least two occurrences of  $\langle ab \rangle$  contain atoms of the form  $occ(1, w, \alpha)$ ,  $occ(1, x, \beta)$  for occurrence 1 and  $occ(2, y, \gamma)$ ,  $occ(2, z, \delta)$  for occurrence 2, with  $\alpha, \beta, \gamma, \delta \in [4]$  and  $w, x, y, z \in [2]$ .

Taking into account the unicity constraint, there are 256 models of the form  $occ(1, 1, \alpha)$ ,  $occ(1, 2, \beta)$ ,  $occ(2, 1, \gamma)$ ,  $occ(2, 2, \delta)$ , with  $\alpha, \beta, \gamma, \delta \in [4]$ . Adding the mapping constraint, the valid models are the following (the 4-tuples below give the related values for  $\alpha, \beta, \gamma, \delta$ ):

(1, 2, 1, 2)	(1, 4, 1, 2)	(1, 2, 1, 4)	(1, 4, 1, 4)
(3, 2, 1, 2)	(3, 4, 1, 2)	(3, 2, 1, 4)	(3, 4, 1, 4)
(1, 2, 3, 2)	(1, 4, 3, 2)	(1, 2, 3, 4)	(1, 4, 3, 4)
(3, 2, 3, 2)	(3, 4, 3, 2)	(3, 2, 3, 4)	(3, 4, 3, 4)

The constraint stating that two occurrences cannot start on the same timestamp, eliminates 8 models:

<del>(1, 2, 1, 2)</del>	<del>(1, 4, 1, 2)</del>	<del>(1, 2, 1, 4)</del>	<del>(1, 4, 1, 4)</del>
(3, 2, 1, 2)	(3, 4, 1, 2)	(3, 2, 1, 4)	(3, 4, 1, 4)
(1, 2, 3, 2)	(1, 4, 3, 2)	(1, 2, 3, 4)	(1, 4, 3, 4)
<del>(3, 2, 3, 2)</del>	<del>(3, 4, 3, 2)</del>	<del>(3, 2, 3, 4)</del>	<del>(3, 4, 3, 4)</del>

---

<sup>2</sup>This can be easily tested using the `project` option of `clingo`.

The constraint about the same end of occurrences eliminates 4 additional models:

<del>(1, 2, 1, 2)</del>	<del>(1, 4, 1, 2)</del>	<del>(1, 2, 1, 4)</del>	<del>(1, 4, 1, 4)</del>
<del>(3, 2, 1, 2)</del>	(3, 4, 1, 2)	(3, 2, 1, 4)	<del>(3, 4, 1, 4)</del>
<del>(1, 2, 3, 2)</del>	(1, 4, 3, 2)	(1, 2, 3, 4)	<del>(1, 4, 3, 4)</del>
<del>(3, 2, 3, 2)</del>	<del>(3, 4, 3, 2)</del>	<del>(3, 2, 3, 4)</del>	<del>(3, 4, 3, 4)</del>

The constraint about the sequentiality of occurrences eliminates 2 additional models:

<del>(1, 2, 1, 2)</del>	<del>(1, 4, 1, 2)</del>	<del>(1, 2, 1, 4)</del>	<del>(1, 4, 1, 4)</del>
<del>(3, 2, 1, 2)</del>	(3, 4, 1, 2)	<del>(3, 2, 1, 4)</del>	<del>(3, 4, 1, 4)</del>
<del>(1, 2, 3, 2)</del>	<del>(1, 4, 3, 2)</del>	(1, 2, 3, 4)	<del>(1, 4, 3, 4)</del>
<del>(3, 2, 3, 2)</del>	<del>(3, 4, 3, 2)</del>	<del>(3, 2, 3, 4)</del>	<del>(3, 4, 3, 4)</del>

The final constraint (line 19) eliminates the model (3, 4, 1, 2) which is a symmetric solution of the remaining model (1, 2, 3, 4).

In this case, the resolution constructs a single answer set:  $\{\text{pattern}(1,1) . \text{pattern}(2,2) . \text{occ}(1,1,1) . \text{occ}(1,2,2) . \text{occ}(2,1,3) . \text{occ}(2,2,4) . \}$ .

This example gives an illustration of successive constraint propagations. This illustration does not correspond to the actual resolution method that uses constraint satisfaction on the grounded program (an intermediary boolean representation of logical rules and facts in the program).

#### 4.4 Extracting all frequent patterns

So far, we have seen how to generate patterns of a given length. To extract all the patterns, the program will proceed level-wise: generate patterns of length 1, then patterns of length 2, ..., then patterns of length  $n$ , then patterns of length  $n+1$ , etc. To this end, we use the control facilities introduced in recent releases of `clingo`. The control part of an ASP program is specified by a python program whose main parts are given below. The program begins with an initial grounding step which is followed by successive resolution steps. At each resolution step, the argument value of the `patlen(k)` atom is incremented to solve the problem for patterns of size  $k$ .

```
prg.ground("base",[])
for k in range(1,maxsize):
    prg.assignExternal(Fun("patlen", [k]), True)
    prg.solve()
    prg.releaseExternal(Fun("patlen", [k]))
```

#### 4.5 Discussion

One of the interest of this program is that it generates only a part of all the occurrences of a pattern in a single AS. The program can be seen as an efficient approach for extracting very frequent patterns, *i.e.*  $\text{supp}(P) \gg \sigma$  (small

patterns for instance). We break down (a part of) the algorithmic complexity by extracting only the number of pattern occurrences required by the support threshold: browsing the whole dataset is not required. Nonetheless, practically, all the occurrences are actually generated but in different answer sets. If a pattern is very frequent, it would be frequent several times with different subsets of its occurrences. The efficiency to compute an AS is lost by the large number of AS generated within this solution.

In addition, the solution does not enforce the anti-monotonicity property since patterns of size  $n$  are computed independently of patterns of size  $m < n$ .

To improve this first solution, we explore the incremental resolution facilities of `clingo`.

## 5 An incremental ASP program for extracting frequent patterns

### 5.1 General principle of incremental extraction of frequent patterns

We use the incremental control facilities of ASP to implement a level-wise extraction of the frequent patterns: patterns of size  $n$  are computed from the patterns of size  $n - 1$ .

#### 5.1.1 Frequent pattern incremental discovery

Let  $P$  be a frequent pattern of length  $n$  ( $n \geq 2$ ) and  $\mathcal{I}(P)$  the set of its occurrences. The “incremental” discovery of frequent patterns is based on browsing the pattern space level-wise: for each item  $s$  in the vocabulary, we would like to know if  $Q = P \oplus r$  is a frequent pattern ( $Q$  is a right extension of  $P$ ). The most difficult part of this problem is to build the occurrences of  $Q$  incrementally from the occurrences of  $P$ . The following property will be useful to solve this problem: any valid (minimal) occurrence of the extension  $Q$  of a pattern  $P$  can be obtained by extending some occurrence of pattern  $P$  and only one such occurrence can be extended (other extensions will not be minimal).

**Property 1.** *Let  $n > 1$  and  $Q = P \oplus r = \langle p_1, \dots, p_n, r \rangle$  be a pattern of size  $n + 1$  extending  $P$  and let  $S = (s_k)$  be a long sequence., we have:*

$$\forall I = (i_1, \dots, i_n) \in \mathcal{I}_S(P), (\exists i', s_{i'} = r \wedge \nexists J, J = (j_1, \dots, j_{n+1}) \in \mathcal{I}_S(Q), [j_1, j_{n+1}] \subset [i_1, i'] \Rightarrow I' = (i_1, \dots, i_n, i') \in \mathcal{I}_S(Q)).$$

#### 5.1.2 Incremental ASP program

The program will be organized in three main parts:

1. the **base** part of the program generates all frequent singleton-patterns (*i.e.* patterns of length 1),

2. the **incr**( $n$ ) part of the program generates all frequent patterns of length  $n$  from patterns of length  $n - 1$  (extension),
3. the control part defines the global strategy.

To benefit from the backward anti-monotonicity of frequent pattern mining, items are added at the end of a pattern. The set of occurrences of some pattern can be computed incrementally thanks to property 1. Contrasting with the previous approach, the complete list of occurrences of the pattern is computed because any occurrence of length  $n$  could be useful to compute an occurrence of length  $n + 1$ .

```
prg.ground("base",[])
prg.solve()

for n in range(2,maxsize):
    prg.ground("incr",[n])
    prg.solve()
```

This program instantiates the constraints incrementally and solves a partial problem at each step. Moreover, the resolution of **incr**( $n+1$ ) depends on the AS solved from **incr**( $n$ ).

## 5.2 Generating patterns of length 1

The **base** part of the incremental program is given below.

```
1 %generate all frequent symbols between 1 and nbs
2 symb(S) :- th { seq(P,S) }, S=1..nbs.
3
4 % generate patterns of length 1
5 1{ pattern(1, P) : symb(P) } 1.
6
7 % generate all first elements of occurrences
8 occ(P, 1, P) :- pattern(1,S), seq(P, S).
```

Line 2 generates the list of the frequent symbols. The symbols identifiers are supposed to be between 1 and **nbs**. **nbs** is a program constant. Line 5 generates the singleton-patterns, consisting of a single symbol, and line 8 generates the occurrences. This rule enforces to have all the occurrences in the pattern.

Note that the identifier of an occurrence is the timestamp of the first item (see line 8).

There is no frequency constraints on **occ** atoms because this kind of constraint is satisfied by rule in line 2 which selects frequent symbols only.

## 5.3 Generating patterns of length $n$ from patterns of length $n - 1$

Extending patterns means to add an item at the end of a (frequent) pattern of length  $n - 1$  previously extracted. The generation of occurrences is based on Property 1. To obtain the occurrences of an extended pattern, the rule in line 15 attempts to complete each occurrence of the sub-pattern by an **occ**( $l, n, p$ )

atom. All combinations of **occ** atoms that associate the last item of the pattern with a timestamp of the sequence (containing the item) are generated.

We assess, line 18 to 20, that generated occurrences are minimals, and thus we eliminate a large number of AS.

```

9 #program incr(n).
10
11 % pattern extension
12 1{ pattern(n, P) : symb(P) } 1.
13
14 % occurrences extension
15 0{ occ(I, n, Q) : seq(Q, S), pattern(n, S), Q>P } 1 :- occ(I, n-1, P).
16
17 % minimal occurrences constraints
18 :- occ(I, n-1, P), occ(I, n, Q), seq(PP, S), pattern(n, S), P<PP, PP<Q.
19 :- occ(I, n, P), occ(J, n, P), I<J.
20 :- occ(I, 1, P), occ(J, 1, PP), occ(I, n, Q), occ(J, n, QQ), P<PP, QQ<Q, I!=J.
21
22 % frequency constraint
23 :- { occ(I, n, _) } th-1.

```

Listing 1: Incremental part of the ASP pattern discovery program

## 5.4 Discussion

The main interest of this solution is the ability to use the anti-monotonicity property to prune the search space and to a priori avoid the evaluation of many models. The incremental program is correct and complete (due to space limitation, the proof is omitted).

This second program shows that, in ASP, it is simple to generate incrementally the candidate patterns of size  $n$  from frequent patterns of size  $n - 1$ . The main difficulty was to rely on pattern occurrences of size  $n - 1$  to efficiently compute the pattern occurrences of size  $n$ . This is achieved thanks to Property 1. This ASP program remains simple (only 9 useful lines) and intuitive: the incremental part contains 2 rules, one for extending patterns, one for extending occurrences and 4 constraints for verifying the minimality of occurrences and minimum frequency.

Contrarily to the first program, all the occurrences of a pattern are collected in a single model. The evaluation of a model is a bit more complex due to more occurrences, but there is much less models to evaluate.

All possible extensions of an occurrence are generated but only the lowest position will generate a minimal occurrence. A more efficient way would be to use directive **#min** to select the correct occurrence extension or to order possible extensions and select the minimal ones.

## 6 Adding sequential patterns constraints

In this section, we show how to include some popular constraints on sequential patterns, *e.g.* *max occurrence duration*, *mingap*, *maxgap*, referring to global

constant values. Such constant values may be defined by `#const` directives or be specified in the control part of the program (python part).

```
#const maxlength = 20. % maximum duration of an occurrence
#const maxgap = 7. % maximum delay between two items in an occurrence
#const mingap = 0. % minimum delay between two items in an occurrence
```

The generation rules of program 1 are modified accordingly. It is better to modify the generation rules than to modify the constraint rules. In the first case, unsatisfiable models are not generated a priori. In the second case, the search space related to the AS is pruned according to the result of evaluating constraint rules.

In the incremental part of the program, the generation of `occ(I, n, Q)` atoms is modified as follows:

```
0{ occ(I, n, Q) : seq(Q, S), pattern(n, S), Q <= maxgap + P, Q > mingap + P, Q <= maxlength + P
  } 1 :- occ(I, n-1, P).
```

$Q$  is a valid timestamp for a pattern extension occurrence if it stands within the limits defined by the constraints:

- gap constraints:  $Q \in [mingap + P, maxgap + P]$ , where  $P$  is the timestamp of the first item
- duration constraint:  $Q \leq maxlength + P$ , where  $P$  is the timestamp of the first item

## 7 Evaluation

### 7.1 Experimentations details

The ASP programs presented above have been evaluated on computation time. All ASP programs were run using the `clingo` solver (version 4.3). The `clingo` processes were limited to a memory space of up to 6 GB. The ASP programs were also compared to the dedicated sequential pattern mining algorithm MinEpi [10] through the implementation provided by the tool DMT4SP [16]. As mentioned in the state of the art, we could not find another generic sequential pattern approach to which we could compare the proposed ASP programs.

The dataset were generated by a random itemset sequence generator. The size of itemsets in the sequence was reduced to 1 and the items were equiprobable.  $ql$  denotes the number of different items in the generated sequence.

### 7.2 Incremental *vs* non-incremental version

We briefly compare the non-incremental version with the incremental version on short sequence (20 items). In fact, the non-incremental version quickly overflows the memory.

Method	Computation time (sec.)	Max Memory (bytes)
Non-Incremental	$27.05 \pm 5.45$	$292996 \pm 916$
Incremental	$0.84 \pm 0.45$	$81336 \pm 10221$

### 7.3 ASP resolution *vs* dedicated algorithm

In this section, we compare the incremental ASP program and the DMT4SP implementation of MinEpi. For both programs, the frequency threshold is set to 10%, the number of different items is set to  $ql = 10$  and the maximal pattern length is set to 10. For each configuration, the resolution was repeated several times with different sequences. If an execution requires more than the maximal amount of memory allowed then it is not taken into account in the final evaluation.

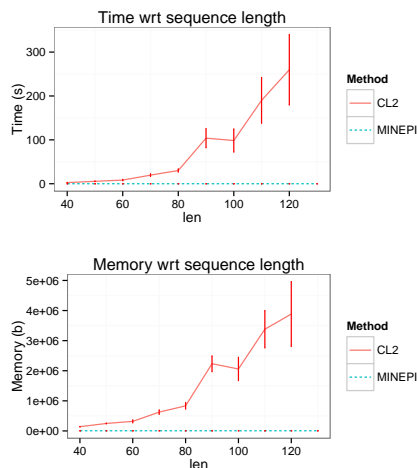


Figure 1: Computation time (in seconds) and memory usage (in bytes) of the incremental ASP program wrt to the length of the sequence.

Figure 1 illustrates the computation time and memory usage for the ASP program. The computation time increases exponentially with length. We can note that the memory allocation is strongly related to the computation time. Both time and memory of ASP are several orders of magnitude higher than MinEpi. The computation time by MinEpi is constantly less than 0.01 second for all the sequences of length fewer than 130 items.

## 7.4 Computing efficiency improvement “with constraints”

In this experiment, we first compared the two incremental ASP programs, without constraints (see section 5) and with sequential pattern constraints (see section 6). The frequency threshold was still set to 10% but the size of the vocabulary was set to  $ql = 7$  (to augment the number of frequent patterns). The constraints were those defined in the program of section 6.

Adding constraints reduces the combinatorics of the search space and, especially, the combinatorics of occurrences extensions. As a consequence, the number of generated models is considerably reduced as well as the computation time. For instance, the mean computation time for a sequence of 70 items is  $166s \pm 34s$  without constraints while with constraints it is  $37s \pm 27s$ . Using constraints, the ASP program becomes computationally more competitive with the dedicated algorithm.

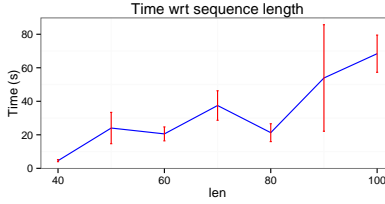


Figure 2: Computation time (in seconds) of the incremental ASP program “with constraints” wrt to the length of the sequence.

Figure 2 illustrates the computation time for the ASP program. The computation times are lower than those of the figure 1. Nonetheless, computation time still increases exponentially with the length of the sequence.

## 8 Conclusion

To the best of our knowledge, this is the first proposed ASP program for frequent pattern mining that implements the anti-monotonicity constraint to prune the search space. It takes advantage of the control facilities of a recent release of `clingo` [6] to design an incremental ASP program. The results show that our proposal is very slow compared to dedicated algorithms such as MinEpi, but our main objective was to demonstrate the feasibility of a declarative and generic approach. ASP solutions are known to be competitive on hard problems. Maybe, enumerating occurrences of a string pattern is not sufficiently hard to be competitive with dedicated algorithms. To assess this assumption, we are working on mining frequent episodes (sequential patterns as sequences of itemsets).

There is room for improving our solution. An immediate improvement would be to use an ordered list of possible extensions. Another classical challenge in pattern mining would be for ASP to tackle the extraction of closed patterns



which can considerably improve the performance of data mining algorithms. The closure definition involves evaluating several patterns together, so the current approach (one model = one pattern) does not seem to be relevant to cope with such constraints. The use of the `#external` directive could be explored to achieve this goal.

Another perspective is to investigate how this framework can be extended to deal with other classical pattern mining problems such as tree mining, subgraph mining, etc. Finally, we would like to explore how to let the user express his own constraints on patterns. A meta-language that could automatically generate constraints in ASP is a midterm objective.

## References

- [1] Avinash Achar, Srivatsan Laxman, and P. S. Sastry. A unified view of automata-based algorithms for frequent episode discovery. *CoRR*, abs/1007.0690, 2010.
- [2] Francesco Bonchi and Claudio Lucchese. Extending the state-of-the-art of constraint-based pattern discovery. *Data Knowl. Eng.*, 60(2):377–399, 2007.
- [3] Jean-Francois Boulicaut and Baptiste Jeudy. Constraint-based data mining. In Oded Maimon and Lior Rokach, editors, *Data Mining and Knowledge Discovery Handbook*, pages 399–416. Springer US, 2005.
- [4] Emmanuel Coquery, Said Jabbour, Lakhdar Saïs, and Yakoub Salhi. A SAT-Based approach for discovering frequent, closed and maximal patterns in a sequence. In *20th European Conference on Artificial Intelligence (ECAI’12)*, pages 258–263, 2012.
- [5] Minos N Garofalakis, Rajeev Rastogi, and Kyuseok Shim. SPIRIT: Sequential pattern mining with regular expression constraints. In *Proceedings of the 25th International Conference on Very Large Data Bases*, pages 223–234. Morgan Kaufmann Publishers Inc., 1999.
- [6] Martin Gebser, Roland Kaminski, Benjamin Kaufmann, Max Ostrowski, Torsten Schaub, and Marius Schneider. Potassco: The Potsdam answer set solving collection. *AI Communications*, 24(2):107–124, 2011.
- [7] Tias Guns, Anton Dries, Guido Tack, Siegfried Nijssen, and Luc De Raedt. MiningZinc: A modeling language for constraint-based mining. In *Proceedings of the Twenty-Third international joint conference on Artificial Intelligence*, pages 1365–1372. AAAI Press, 2013.
- [8] Tias Guns, Siegfried Nijssen, and Luc De Raedt. Itemset mining: A constraint programming perspective. *Artificial Intelligence*, 175(12-13):1951–1983, 2011.

- [9] Matti Järvisalo. Itemset mining as a challenge application for answer set enumeration. In *Logic Programming and Nonmonotonic Reasoning*, pages 304–310. Springer, 2011.
- [10] H. Mannila, H. Toivonen, and A. I. Verkamo. Discovering frequent episodes in event sequences. *Journal of Data Mining and Knowledge Discovery*, 1(3):210–215, 1997.
- [11] Jean-Philippe Métivier, Samir Loudni, and Thierry Charnois. A constraint programming approach for mining sequential patterns in a sequence database. *CoRR*, abs/1311.6907, 2013.
- [12] Benjamin Negrevergne, Anton Dries, Tias Guns, and Siegfried Nijssen. Dominance programming for itemset mining. In *International Conference on Data Mining (ICDM)*, 2013.
- [13] Jian Pei, Jiawei Han, Behzad Mortazavi-Asl, Jianyong Wang, Helen Pinto, Qiming Chen, Umeshwar Dayal, and Mei-Chun Hsu. Mining sequential patterns by pattern-growth: the PrefixSpan approach. *IEEE Transactions on Knowledge and Data Engineering*, 16(11):1424–1440, 2004.
- [14] Jian Pei, Jiawei Han, and Wei Wang. Mining sequential patterns with constraints in large databases. In *Proceedings of the eleventh international conference on Information and knowledge management*, pages 18–25. ACM, 2002.
- [15] Jian Pei, Jiawei Han, and Wei Wang. Constraint-based sequential pattern mining: The pattern-growth methods. *J. Intell. Inf. Syst.*, 28(2):133–160, 2007.
- [16] Christophe Rigotti. <http://liris.cnrs.fr/crigotti/dmt4sp.html>.
- [17] Ramakrishnan Srikant and Rakesh Agrawal. Mining sequential patterns: Generalizations and performance improvements. In *Proceedings of the 5th International Conference on Extending Database Technology*, pages 3–17, 1996.
- [18] Nikolaj Tatti and Boris Cule. Mining closed episodes with simultaneous events. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1172–1180, 2011.
- [19] Nikolaj Tatti and Boris Cule. Mining closed strict episodes. *Data Mining and Knowledge Discovery*, 25(1):34–66, 2012.
- [20] Jianyong Wang and Jiawei Han. Bide: Efficient mining of frequent closed sequences. In *Proceedings of the 20th International Conference on Data Engineering*, pages 79–90. IEEE Computer Society, 2004.
- [21] Xifeng Yan, Jiawei Han, and Ramin Afshar. Clospan: Mining closed sequential patterns in large datasets. In *In SDM*, pages 166–177, 2003.

- [22] Mohammed J. Zaki. SPADE: An efficient algorithm for mining frequent sequences. *Journal of Machine Learning*, 42(1/2):31–60, 2001.