

# Hybrid ASP-based Approach to Pattern Mining<sup>\*</sup>

Sergey Paramonov<sup>a</sup>, Daria Stepanova<sup>b</sup>, Pauli Miettinen<sup>b</sup>

<sup>a</sup>Department of Computer Science, KU Leuven, Belgium

<sup>b</sup>Max Planck Institute of Informatics, Saarbrücken, Germany

sergey.paramonov@kuleuven.be, {dstepano, pmiettinen}@mpi-inf.mpg.de

**Abstract.** Detecting small sets of relevant patterns from a given dataset is a central challenge in data mining. The relevance of a pattern is based on user-provided criteria; typically, all patterns that satisfy certain criteria are considered relevant. Rule-based languages like Answer Set Programming (ASP) seem well-suited for specifying such criteria in a form of constraints. Although progress has been made, on the one hand, on solving individual mining problems and, on the other hand, developing generic mining systems, the existing methods either focus on scalability or on generality. In this paper we make steps towards combining local (frequency, size, cost) and global (various condensed representations like maximal, closed, skyline) constraints in a generic and efficient way. We present a hybrid approach for itemset and sequence mining which exploits dedicated highly optimized mining systems to detect frequent patterns and then filters the results using declarative ASP. Experiments on real-world datasets show the effectiveness of the proposed method and computational gains both for itemset and sequence mining.

## 1 Introduction

**Motivation.** Availability of vast amounts of data from different domains has led to an increasing interest in the development of scalable and flexible methods for data analysis. A key feature of flexible data analysis methods is their ability to incorporate users' background knowledge and different criteria of interest. They are often provided in the form of *constraints* to the valid set of answers, the most common of which is the frequency threshold: a pattern is only considered interesting if it appears often enough. Mining all frequent (and otherwise interesting) patterns is a very general problem in data analysis, with applications in medical treatments, customer shopping sequences, Weblog click streams and text analysis, to name but a few examples.

Most data analysis methods consider only one (or few) types of constraints, limiting their applicability. Constraint Programming (CP) has been proposed as a general approach for (sequential) mining of frequent patterns [1], and Answer Set Programming (ASP) [6] has been proved to be well-suited for defining the constraints conveniently thanks to its expressive and intuitive modelling language and the availability of optimized ASP solvers (see e.g., [13,5,9] for existing approaches).

In general, all constraints can be classified into *local constraints*, that can be validated by the pattern candidate alone, and *global constraints*, that can only be validated

---

<sup>\*</sup> This work is partially funded by the ERC AdG SYNTH (Synthesising inductive data models).

via an exhaustive comparison of the pattern candidate against all other candidates. Combining local and global constraints in a generic way is an important and challenging problem, which has been widely acknowledged in the constraint-based mining community. Although progress has been made, on the one hand, on solving individual mining problems and, on the other hand, on developing generic mining systems, the existing methods either focus on scalability or on generality, but rarely address both of these aspects. This naturally limits the practical applicability of the existing approaches.

**State of the art and its limitations.** Purely declarative ASP encodings for frequent and maximal itemset mining were proposed in [13]. In this approach, first every item's inclusion into the candidate itemset is guessed, and the guessed candidate pattern is checked against frequency and maximality constraints. While natural, this encoding is not truly generic, as adding extra local constraints requires significant changes in it. Indeed, for a database, where all available items form a frequent (and hence maximal) itemset, the maximal ASP encoding has a single model. The latter is, however, eliminated once restriction on the length of allowed itemsets is added to the program. This is undesired, as being maximal is not a property of an itemset on its own, but rather in the context of a collection of other itemsets [3]. Thus, ideally one would be willing to first apply all local constraints and only afterwards construct a condensed representation of them, which is not possible in [13].

This shortcoming has been addressed in the recent work [5] on ASP-based sequential pattern mining, which exploits ASP preference-handling capacities to extract patterns of interest and supports the combination of local and global constraints. However, both [5] and [13] present purely declarative encodings, which suffer from scalability issues caused by the exhaustive exploration of the huge search space of candidate patterns. The subsequence check amounts to testing whether an embedding exists (matching of the individual symbols) between sequences. In sequence mining, a pattern of size  $m$  can be embedded into a sequence of size  $n$  in  $O(n^m)$  different ways, therefore, clearly a direct pattern enumeration is unfeasible in practice.

While a number of individual methods tackling selective constraint-based mining tasks exist (see Tab. 1 for comparison) there is no uniform ASP-based framework that is capable of effectively combining constraints both on the global and local level and is suitable for itemsets and sequences alike.

**Contributions.** The goal of our work is to make steps towards building a generic framework that supports mining of condensed (sequential) patterns, which (1) effectively combines dedicated algorithms and declarative means for pattern mining and (2) is easily extendable to incorporation of various constraints. More specifically, the salient contributions of our work can be summarized as follows:

- We present a general extensible pattern mining framework for mining patterns of different types using ASP.
- We introduce a feature comparison, such as closedness under solutions, between different ASP mining models and dominance programming, which is a generic itemset mining language and solver.
- We demonstrate the feasibility of our approach with an experimental evaluation across multiple itemset and sequence datasets.

Datatype	Task	[13]	[5]	[16]	Our work
Itemset	frequent pattern mining	✓	–	✓	✓
	condensed (closed, max, etc)	✓*	–	✓	✓
	condensed under constraints	–	–	✓	✓
Sequence	frequent pattern mining	–	✓	–	✓
	condensed (closed, max, etc)	–	✓	–	✓
	condensed under constraints	–	✓	–	✓

Table 1: Feature comparison between various ASP mining models and dominance programming (“–”: “not designed for this datatype”, ✓\*: only maximal is supported)

## 2 Preliminaries

In this section we briefly recap the necessary background both from the fields of pattern mining and Answer Set Programming (ASP).

Let  $D$  be a dataset,  $\mathcal{L}$  a language for expressing pattern properties or defining subgroups of the data, and  $q$  a selection predicate. The task of pattern mining is to find  $Th(\mathcal{L}, D, q) = \{\phi \in \mathcal{L} \mid q(D, \phi) \text{ is true}\}$  (see the seminal work [14]).

Pattern mining has been mainly studied for itemsets, sequences and graphs. These settings are determined by the language of  $\mathcal{L}$ . We focus on the first two categories.

### 2.1 Patterns

**Itemsets.** *Itemsets* represent the most simple setting of frequent pattern mining. Let  $\mathcal{I}$  be a set of items  $\{o_1, o_2, \dots, o_n\}$ . Then a nonempty subset of  $\mathcal{I}$  is called an *itemset*. A *transaction dataset*  $D$  is a collection of itemsets,  $D = \{t_1, \dots, t_m\}$ , where  $t_i \subseteq \mathcal{I}$ . For any itemset  $\alpha$ , we denote the set of transactions that contain  $\alpha$  as  $D_\alpha = \{i \mid \alpha \subseteq t_i, t_i \in D\}$  and we refer to  $|D_\alpha|$  as the *support (frequency)* of  $\alpha$  in  $D$ , written  $sup(\alpha)$ . The *relative frequency* of  $\alpha$  in  $D$  refers to the ratio between  $sup(\alpha)$  and  $|D|$ . The *cardinality (or size)*  $|\alpha|$  of an itemset  $\alpha$  is the number of items contained in it.

**Definition 1 (Frequent Itemset).** *For a transaction dataset  $D$  and a frequency threshold  $\sigma \geq 0$ , an itemset  $\alpha$  is frequent in  $D$  if  $sup(\alpha) \geq \sigma$ .<sup>1</sup>*

*Example 1.* Consider a transaction dataset  $D$  from Tab. 2. We have  $\mathcal{I} = \{a, b, c, d, e\}$  and  $|D| = 3$ . For  $\sigma = 2$ , the following itemsets are frequent:  $\alpha_1 = \{a\}$ ,  $\alpha_2 = \{b\}$ ,  $\alpha_3 = \{e\}$ ,  $\alpha_4 = \{a, e\}$  and  $\alpha_5 = \{b, e\}$ .  $\square$

**Sequences.** A *sequence* is an ordered set of items  $\langle s_1, \dots, s_n \rangle$ . The setting of *sequence mining* includes two related yet different cases: frequent substrings and frequent subsequences. In this work we focus on the latter.

<sup>1</sup> In *frequent pattern mining*, often, a *relative threshold*, i.e.,  $\frac{\sigma}{|D|}$  is specified by the user.

ID	a	b	c	d	e
1	✓	✓		✓	✓
2		✓	✓		✓
3	✓				✓

Table 2: Transaction database

ID	Sequence
1	$\langle a b c d a e b \rangle$
2	$\langle b c e b \rangle$
3	$\langle a a e \rangle$

Table 3: Sequence database

**Definition 2 (Embedding in a Sequence).** Let  $S = \langle s_1, \dots, s_m \rangle$  and  $S' = \langle s'_1, \dots, s'_n \rangle$  be two sequences of size  $m$  and  $n$  respectively with  $m \leq n$ . The tuple of integers  $e = (e_1, \dots, e_m)$  is an embedding of  $S$  in  $S'$  (denoted  $S \sqsubseteq_e S'$ ) if and only if  $e_1 < \dots < e_m$  and for any  $i \in \{1, \dots, m\}$  it holds that  $s_i = s'_{e_i}$ .

*Example 2.* For a dataset in Tab. 3 we have that  $\langle b c e b \rangle \sqsubseteq_{e_1} \langle a b c d a e b \rangle$  for  $e_1 = (2, 3, 6, 7)$  and analogously,  $\langle a a e \rangle \sqsubseteq_{e_2} \langle a b c d a e b \rangle$  with  $e_2 = (1, 5, 6)$ .

We are now ready to define an inclusion relation for sequences.

**Definition 3 (Sequence Inclusion).** Given two sequences  $S = \langle s_1, \dots, s_n \rangle$  and  $S' = \langle s'_1, \dots, s'_m \rangle$ , of size  $m$  and  $n$  resp. with  $n \leq m$ , we say that  $S$  is included in  $S'$  or  $S$  is a subsequence of  $S'$  denoted by  $S \sqsubseteq S'$  iff an embedding  $e$  of  $S$  in  $S'$  exists, i.e.

$$S \sqsubseteq S' \leftrightarrow \exists e_1 < \dots < e_m \text{ and } \forall i \in 1 \dots m : s_i = s'_{e_i}. \quad (1)$$

*Example 3.* In Ex. 2 we have  $\langle b c e b \rangle \sqsubseteq \langle a b c d a e b \rangle$  but  $\langle a a e \rangle \not\sqsubseteq \langle b c e b \rangle$ .  $\square$

For a given sequence  $S$  and a sequential dataset  $D = \{S_1, \dots, S_n\}$  we denote by  $D_S$  the subset of  $D$  s.t.  $S \sqsubseteq S'$  for all  $S' \in D_S$ . The support of  $S$  is  $\text{sup}(S) = |D_S|$ .

**Definition 4 (Frequent Sequence).** For a sequential dataset  $D = \{S_1, \dots, S_n\}$  and a frequency threshold  $\sigma \geq 0$ , a sequence  $S$  is frequent in  $D$  if  $\text{sup}(S) \geq \sigma$ .

*Example 4.* For a dataset in Tab. 3 and  $\sigma = 2$ , it holds that  $\langle b c e b \rangle$  and  $\langle a a e \rangle$  are frequent, while  $\langle b d b \rangle$  is not.  $\square$

Note that  $\sqsubseteq$  and  $\subseteq$  are incomparable relations. Indeed, consider two sequences  $s_1 = \langle a b \rangle$  and  $s_2 = \langle b a a \rangle$ . While  $s_1 \subset s_2$ , we clearly have that  $s_1 \not\sqsubseteq s_2$ .

## 2.2 Condensed Pattern Representations under Constraints

In data mining, constraints are typically specified by the user to encode domain background knowledge. In [17] four types of constraints are distinguished: constraints 1) over the pattern (e.g., restriction on its size), 2) over the cover set (e.g., minimal frequency), 3) over the inclusion relation (e.g., maximal allowed gap in sequential patterns) and 4) over the solution set (e.g., condensed representations).

Orthogonally, constraints can be classified into *local* and *global* ones. A constraint is *local* if deciding whether a given pattern satisfies it is possible without looking at other patterns. For example, minimal frequency or maximal pattern size are local constraints. On the contrary, deciding whether a pattern satisfies a *global* constraint requires

comparing it to other patterns. All constraints from the 4th group are global ones. We are interested in global constraints related to condensed representations.

As argued in Sec. 1, the order in which constraints are applied influences the solution set [3]. As in [3] in this work we apply global constraints only after local ones.

We now present the notions required in our pattern mining framework. Here, the definitions are given for itemsets; for sequences they are identical up to substitution of  $\subset$  with  $\sqsubset$  (subsequence relation). First, to rule out patterns that do not satisfy some of the local constraints, we introduce the notion of validity.

**Definition 5 (Valid pattern under constraints).** *Let  $C$  be a constraint function from  $\mathcal{L}$  to  $\{\top, \perp\}$  and let  $p$  be a pattern in  $\mathcal{L}$ , then the pattern  $p$  is called valid iff  $C(p) = \top$ , otherwise it is referred as invalid.*

*Example 5.* Let  $C$  be a constraint function checking whether a given pattern is of size at least 2. Then for Ex. 1, we have  $C(\alpha_i) = \perp, i = 1..3$  and  $C(\alpha_j) = \top, j = 4..5$ .  $\square$

For detecting patterns that satisfy a given global constraint, the notion of *dominance* is of crucial importance. Intuitively, a dominance relation reflects pairwise preference ( $<^*$ ) between patterns and it is specific for each mining setting. In this work we primarily focus on global constraints related to maximal, closed, free and skyline condensed representations, for which  $<^*$  is defined as follows:

- (i) **Maximal.** For itemsets  $p, q$ ,  $p <^* q$  holds iff  $p \subset q$
- (ii) **Closed.** For itemsets  $p, q$ ,  $p <^* q$  holds iff  $p \subset q \wedge \text{sup}(p) = \text{sup}(q)$
- (iii) **Free.** For itemsets  $p, q$ ,  $p <^* q$  holds iff  $q \subset p \wedge \text{sup}(p) = \text{sup}(q)$
- (iv) **Skyline.** For itemsets  $p, q$ ,  $p <^* q$  holds iff
  - (a)  $\text{sup}(p) \leq \text{sup}(q)$  and  $\text{size}(p) < \text{size}(q)$  or
  - (b)  $\text{sup}(p) < \text{sup}(q)$  and  $\text{size}(p) \leq \text{size}(q)$

Dominated patterns under constraints are now formally defined.

**Definition 6 (Dominated pattern under constraints).** *Let  $C$  be a constraint function, and let  $p$  be a pattern, then  $p$  is called dominated iff there exists a pattern  $p' \in \mathcal{L}$  such that  $p <^* p'$  and  $p'$  is valid under  $C$ .*

*Example 6.* In Ex. 1 for the maximality constraint we have that  $\alpha_1$  is dominated by  $\alpha_4$ ,  $\alpha_2$  by  $\alpha_5$ , while  $\alpha_3$  both by  $\alpha_4$  and  $\alpha_5$ .  $\square$

Exploiting the above definitions we obtain condensed patterns under constraints.

**Definition 7 (Condensed pattern under constraints).** *Let  $p$  be a pattern from  $\mathcal{L}$ , and let  $C$  be a constraint function, then a pattern  $p$  is called condensed under constraints iff it is valid and not dominated under  $C$ .*

*Example 7.* For the constraint function selecting maximal itemsets of size at most 2 and size at least 2,  $\alpha_4$  and  $\alpha_5$  from Ex. 1 are condensed patterns.  $\square$

### 2.3 Answer Set Programming

Answer Set Programming (ASP) [6] is a declarative problem solving paradigm oriented towards difficult search problems. ASP has its roots in Logic Programming and Non-monotonic Reasoning. An ASP program  $\Pi$  is a set of rules of the form

$$a_0 \text{ :- } b_1, \dots, b_k, \text{ not } b_{k+1}, \dots, \text{ not } b_m, \quad (2)$$

where  $1 \leq k \leq m$ , and  $a_0, b_1, \dots, b_m$  are classical literals, and *not* is default negation. The right-hand side of  $r$  is its body,  $Body(r)$ , while the left-hand side is the head,  $Head(r)$ .  $Body^+(r)$  and  $Body^-(r)$  stand for the positive and negative parts of  $Body(r)$  respectively. A rule  $r$  of the form (2) is a *fact* if  $m = 0$ . We omit the symbol  $:-$  when referring to facts. A rule without head literals is a *constraint*. A rule is *positive* if  $k = m$ .

An ASP program  $\Pi$  is *ground* if it consists of only ground rules, i.e. rules without variables. Ground instantiation  $Gr(\Pi)$  of a nonground program  $\Pi$  is obtained by substituting variables with constants in all possible ways. The *Herbrand universe*  $HU(\Pi)$  (resp. *Herbrand base*  $HB(\Pi)$ ) of  $\Pi$ , is the set of all constants occurring in  $\Pi$ , (resp. the set of all possible ground atoms that can be formed with predicates and constants occurring in  $\Pi$ ). Any subset of  $HB(P)$  is a *Herbrand interpretation*.  $MM(\Pi)$  denotes the subset-minimal Herbrand interpretation (*model*) of a ground positive program  $\Pi$ .

The semantics of an ASP program is given in terms of its answer sets. An interpretation  $A$  of  $\Pi$  is an *answer set* (or *stable model*) of  $\Pi$  iff  $A = MM(\Pi^A)$ , where  $\Pi^A$  is the *Gelfond–Lifschitz (GL) reduct* [6] of  $\Pi$ , obtained from  $Gr(\Pi)$  by removing (i) each rule  $r$  such that  $Body^-(r) \cap A \neq \emptyset$ , and (ii) all the negative atoms from the remaining rules. The set of answer sets of a program  $\Pi$  is denoted by  $AS(\Pi)$ .

*Example 8.* Consider the program  $\Pi$  given as follows:

```
(1) pattern(1); (2) pattern(2); (3) item(1,a);
(4) item(1,b); (5) item(2,a);
(6) not_subset(I,J):-pattern(I), item(I,V), I != J,
    pattern(J), not item(J,V).
```

The grounding  $Gr(\Pi)$  of  $\Pi$  is obtained from  $\Pi$  by substituting  $I, J$  with  $1, 2$  and  $V$  with  $b$  resp. The GL-reduct  $\Pi^{A'}$  ( $\Pi$ ) for the interpretation  $A'$  containing facts of  $\Pi$  and  $not\_subset(1, 2)$  differs from  $Gr(\Pi)$  only in that  $not\_item(2, b)$  is not in the body of the rule.  $A'$  is the minimal model of  $\Pi^{A'}$  ( $\Pi$ ), and thus it is in  $AS(\Pi)$ .  $\square$

Other relevant language constructs include conditional literals and cardinality constraints [22]. The former are of the form  $a : b_1, \dots, b_m$ , the latter can be written as  $l\{c_1, \dots, c_n\}t$ , where  $a$  and  $b_i$  are possibly default negated literals and each  $c_j$  is a conditional literal;  $l$  and  $t$  provide lower and upper bounds on the number of satisfied literals in a cardinality constraint. For instance,  $1\{a(X) : b(X)\}3$  holds, whenever between 1 and 3 instances of  $a(X)$  (subject to  $b(X)$ ) are true. Furthermore, aggregates are of the form  $\#sum\{K : cost(I, K)\} > N$ . This atom is true, whenever the sum of all  $K$ , such that  $cost(I, K)$  is true, exceeds  $N$ .

### 3 Hybrid ASP-based Mining Approach

In this section we present our hybrid method for frequent pattern mining. Unlike previous ASP-based mining methods, our approach combines highly optimized algorithms for frequent pattern discovery with the declarative ASP means for their convenient post-processing. Here, we focus on itemset and sequence mining; however our approach can be also applied to subgraph discovery (details are left for future work).

Given a frequency threshold  $\sigma$ , a (sequential) dataset  $D$  and a set of constraints  $\mathcal{C} = \mathcal{C}_l \cup \mathcal{C}_g$ , where  $\mathcal{C}_l$  and  $\mathcal{C}_g$  are respectively local and global constraints, we proceed in two steps as follows.

**Step 1.** First, we launch a dedicated optimized algorithm to extract all (sequential) frequent patterns from a given dataset, satisfying the minimal frequency threshold  $\sigma$ . Here, any frequent pattern mining algorithm can be invoked. We use Eclat [24] for itemsets and PPIC [2] for sequences.

**Step 2.** Second, the computed patterns are post-processed using the declarative means to select a set of *valid* patterns (i.e., those satisfying constraints in  $\mathcal{C}_l$ ). For that the frequent patterns obtained in Step 1 are encoded as facts `item(i, j)` for itemsets and `seq(i, j, p)` for sequences where  $i$  is the pattern's ID,  $j$  is an item contained in it and  $p$  is its position. The local constraints in  $\mathcal{C}_l$  are represented as ASP rules, which collect IDs of patterns satisfying constraints from  $\mathcal{C}_l$  into the dedicated predicate `valid`, while the rest of the IDs are put into the `not_valid` predicate.

Finally, from all valid patterns a desired condensed representation is constructed by storing patterns  $i$  in the `selected` predicate if they are not dominated by other valid patterns based on constraints from  $\mathcal{C}_g$ . Following the principle of [13], in our work every answer set represents a single desired pattern, which satisfies both local and global constraints. The set of all such patterns forms a condensed representation. In what follows we discuss our encodings of local and global constraints in details.

#### 3.1 Encoding Local Constraints

In our declarative program we specify local constraints by the predicate `valid`, which reflects the conditions given in Def. 5. For every constraint in  $\mathcal{C}_l$  we have a set of dedicated rules, stating when a pattern is not valid. For instance, a constraint checking whether the cost of items in a pattern exceeds a given threshold  $N$  is encoded as

```
not_valid(I) :- #sum{C:item(I,J),cost(J,C)} > N, pattern(I).
```

A similar rule for sequences can be defined as follows:

```
not_valid(I) :- #sum{C:seq(I,J,P),cost(J,C)} > N, pattern(I).
```

Analogously, one can specify arbitrary domain constraints on patterns.

```

1 time(1..5).
2 % people born in Germany or France are Europeans
3 eu(I) :- seq(I,bG,P).
4 eu(I) :- seq(I,bF,P).
5 % collect those who moved to France before P
6 moved_before(X,P) :- seq(X,mF,P1), P>P1, time(P), time(P1).
7 % collect those who moved to France after P and before masters
8 moved_after(X,P) :- seq(X,mF,P1), seq(X,ma,P2), P<P1,
9                     p1<P2, time(P), time(P1), time(P2).
10 % keep Europeans who moved to Germany straight before masters
11 keep(X) :- seq(X,ma,P+1), seq(X,mG,P), eu(X).
12 % keep Germans who did not move before masters
13 keep(X) :- seq(X,bG,P1), seq(X,ma,P), not moved_before(X,P).
14 % keep Europeans whose last move before masters was to Germany
15 keep(X) :- seq(X,mG,P1), seq(X,ma,P2), P1<P2,
16             eu(X), not moved_after(X,P1).
17 % a pattern is not valid, if it should not be kept
18 not_valid(X) :- pattern(X), not keep(X)
19

```

Listing 1.1: Moving habits of people during studies

*Example 9.* Consider a dataset storing moving habits of young people during their studies. Let the dedicated frequent sequence mining algorithm return the following patterns:  $S_1 = \langle bG\ mF\ ba\ mG\ ma \rangle$ ;  $S_2 = \langle bF\ mG\ ba\ mF\ ma \rangle$ ;  $S_3 = \langle bA\ ba\ ma \rangle$ , where  $bG$ ,  $bF$ ,  $bA$  stand for born in Germany, France and America,  $ba$ ,  $ma$  stand for bachelors and masters and the predicates  $mG$ ,  $mF$  reflect that a person moved to Germany and France, respectively. Suppose, we are only interested in moving habits of Europeans, who got their masters degree from a German university. The local domain constraint expressing this would state that (1)  $bA$  should not be in the pattern, while (2) either both  $bG$  and  $ma$  should be in it without any  $mF$  in between or  $mG$  should precede  $ma$ . These constraints are encoded in the program in List. 1.1. From the answer set of this program we get that both  $S_2$  and  $S_3$  are not valid, while  $S_1$  is.  $\square$

To combine all local constraints from  $\mathcal{C}_l$  we add to a program a generic rule specifying that a pattern  $I$  is valid whenever `not_valid(I)` cannot be inferred.

```
valid(I) :- pattern(I), not not_valid(I)
```

Patterns  $i$ , for which `valid(i)` is deduced are then further analyzed to construct a condensed representation based on global constraints from  $\mathcal{C}_g$ .

### 3.2 Encoding Global Constraints

The key for encoding global constraints is the declarative formalization of the dominance relation (Def. 6). For example, for itemsets the maximality constraint boils down



```

1 % I is not a subset of J if I has items that are not in J
2 not_subset(J) :- selected(I), item(I,V), not item(J,V),
3                 valid(J), I != J.
4 % derive dominated whenever I is subset of J
5 dominated :- selected(I), valid(J),
6              I != J, not not_subset(J).

```

Listing 1.2: Maximal itemsets encoding

to pairwise checking of subset inclusion between patterns. For sequences this requires a check of embedding existence between sequences.

Regardless of a pattern type from  $\mathcal{L}$  and a constraint from  $\mathcal{C}_g$  every encoding presented in this section is supplied with a rule, which guesses (`selected/1` predicate) a single `valid` pattern to be a candidate for inclusion in the condensed representation, and a constraint that rules out `dominated` patterns thus enforcing a different guess.

```

1 {selected(I) : valid(I)} 1.
   :- dominated.

```

In what follows, we discuss concrete realizations of the dominance relation both for itemsets and sequences for various global constraints, i.e., we present specific rules related to the derivation of the `dominated/0` predicate.

**Itemset Mining.** We first provide an encoding for maximal itemset mining in List 1.2. To recall, a pattern is *maximal* if none of its supersets is frequent. An itemset  $I$  is included in  $J$  iff for every item  $i \in I$  we have  $i \in J$ . We encode the violation of this condition in lines (1)–(3). The second rule presents the dominance criteria.

For closed itemset mining a simple modification of List. 1.2 is required. An itemset is *closed* if none of its supersets has the same support. Thus to both of the rules from List. 1.2 we need to add atoms `support(I,X)`, `support(J,X)`, which store the support sets of  $I$  and  $J$  respectively (extracted from the output of Step 1).

For free itemset mining the rules of the maximal encoding are changed as follows:

```

4 not_superset(J) :- selected(I), item(J,V), not item(I,V),
5                 valid(J), I != J.
6 dominated :- selected(I), valid(J), support(I,X),
7              I != J, not not_superset(J), support(J,X).

```

Finally, the skyline itemset/sequence encoding is given in List. 1.3, where the first two rules specify the conditions (a) and (b) for skyline itemsets as specified in Sec. 2.

**Sequence Mining.** The subpattern relation for sequences is slightly more involved, than for itemsets, as it preserves the order of elements in a pattern. To recall, a sequence  $S$  is included in  $S'$  iff an embedding  $e$  exists, such that  $S \sqsubseteq_e S'$ .

In List. 1.4 we present the encoding for maximal sequence mining. A selected pattern is not maximal if it has at least one valid superpattern. We rule out patterns that are

```

1 % support and size comparison among patterns
2 g_size_geq_fr(J) :- selected(I), valid(J), support(I,X),
3                   support(J,Y), size(I,Si), size(J, Sj),
4                   Si < Sj, X <= Y.
5 geq_size_g_fr(J) :- selected(I), valid(J), support(I,X),
6                   support(J,Y), size(I,Si), size(J, Sj),
7                   Si <= Sj, X < Y.
8 % derivation of the domination condition
9 dominated :- valid(J), g_size_geq_fr(J).
10 dominated :- valid(J), geq_size_g_fr(J).

```

Listing 1.3: Skyline itemsets encoding

for sure not superpatterns of a selected sequence. First, obviously  $J$  is not a superpattern of  $I$  if  $I$  is not a subset of  $J$  (lines (4)–(5)), i.e., if `not_subset(J)` is derived, then  $J$  does not dominate  $I$ . If  $J$  is a superset of  $I$  then to ensure that  $I$  is not dominated by  $J$ , the embedding existence has to be checked (lines (6)–(9)).  $I$  is not dominated by  $J$  if an item exists in  $I$ , which together with its sequential neighbor cannot be embedded in  $J$ . This condition is checked in lines (10)–(13), where `domcand(V, J, P)` is derived if for an item  $V$  at position  $P$  and its follower, embedding in  $J$  can be found.

The encoding for closed sequence mining is obtained from the maximal sequence encoding analogously as it is done for itemsets. The rules for free sequence mining are constructed by substituting lines (4)–(13) of List. 1.4 with the following ones:

```

4 not_superset(J) :- selected(I), in(V, J),
5                   not in(V, I), I != J.
6 domcand(V, J) :- selected(I), seq(J, V, P), item(J, P+1, W),
7                 item(I, V, Q), seq(I, W, Q'), Q' > Q, I != J.
8 not_dominated_by(J) :- selected(I), valid(J), I != J,
9                       seq(J, V, P), seq(J, W, P+1),
10                      not domcand(V, J).

```

Finally, the encoding for mining skyline sequences coincides with the skyline itemsets encoding, which is provided in List. 1.3.

## 4 Evaluation

In this section we evaluate the proposed hybrid approach by comparing it to the existing declarative pattern mining methods: ASP model for sequences from [5] and Dominance Programming (DP) from [16]. We do not consider the itemset mining ASP model [13], since it focuses only on frequent itemset mining and is not applicable to the construction of condensed representations under constraints as in [16]. Moreover, we do not perform comparison with dedicated algorithms designed for a specific problem type; these are known to be more efficient than declarative mining approaches [17].

More specifically, we investigate the following experimental questions.

```

1 % if V appears in a valid pattern I, derive in(V,I)
2 in(V,I) :- seq(I,V,P), valid(I).
3 % I is not a subset of J if I has V that J does not have
4 not_subset(J) :- selected(I), valid(J), I != J,
5                   seq(I,V,P), not in(V,J).
6 % if for a subseq <V,W> in I there is V followed
7 % by W in J then deduce domcand(V,J)
8 domcand(V,J,P) :- selected(I), seq(I,V,P), seq(I,W,P+1), I != J
9                   valid(J), seq(J,V,Q), seq(J,W,Q'), Q'>Q.
10 % if domcand(V,J) does not hold for some V in I
11 % and a pattern J then derive not_dominated_by(J)
12 not_dominated_by(J) :- selected(I), seq(I,V,P), seq(I,W,P+1),
13                        I != J, valid(J), not domcand(V,P,J).
14 % if neither not_dominated_by(J) nor not_subset_of(J)
15 % are derived for some J, then I is dominated
16 dominated :- selected(I), valid(J), I != J,
17              not not_subset_of(J), not not_dominated_by(J).

```

Listing 1.4: Maximal sequence encoding

- **Q<sub>1</sub>**: how does the runtime of our method compare to the existing ASP-based sequence mining models?
- **Q<sub>2</sub>**: what is the runtime gap between the specialized mining languages such as dominance programming and our method?
- **Q<sub>3</sub>**: what is the influence of local constraints on the runtime of our method?

In **Q<sub>1</sub>** we compare our work with the ASP-based model from [5]. In **Q<sub>2</sub>** we measure the runtime difference between specialized itemset mining languages [16] and our ASP-based model. Finally, in **Q<sub>3</sub>** we estimate the runtime effect of adding local constraints.

We report evaluation on 2 transaction datasets,<sup>2</sup> *Mushrooms* (8124 transactions/119 items) and *Vote* (435/48), and on 3 sequence datasets (full),<sup>3</sup> *JMLR* (788 sequences/3847 distinct symbols), *Unix Users* (9099/4093), and *iPRG* (8628/21). All experiments have been performed on a desktop with Ubuntu 14.04, 64-bit environment, Intel Core i5-3570 4xCPU 3.40GHz and 8GB memory using clingo 4.5.4<sup>4</sup> and C++14 for the wrapper. The timeout was set to one hour. Free pattern mining demonstrates the same runtime behavior as closed, due to the symmetric encoding, and is thus omitted.

To investigate **Q<sub>1</sub>**, in Fig. 1a, we compare the ASP model [5] with our method on the default 200 sequence sample, generated by the tool<sup>5</sup> from [5]. We performed the comparison on the synthetic data, as the sequence-mining model [5] failed to compute condensed representations on any of the standard sequence datasets for any support threshold value within the timeout. One can observe that our method consistently out-

<sup>2</sup> From <https://dtai.cs.kuleuven.be/CP4IM/datasets/>.

<sup>3</sup> From <https://dtai.cs.kuleuven.be/CP4IM/cpsm/datasets.html>.

<sup>4</sup> <http://potassco.sourceforge.net>

<sup>5</sup> <https://sites.google.com/site/aspseqmining>

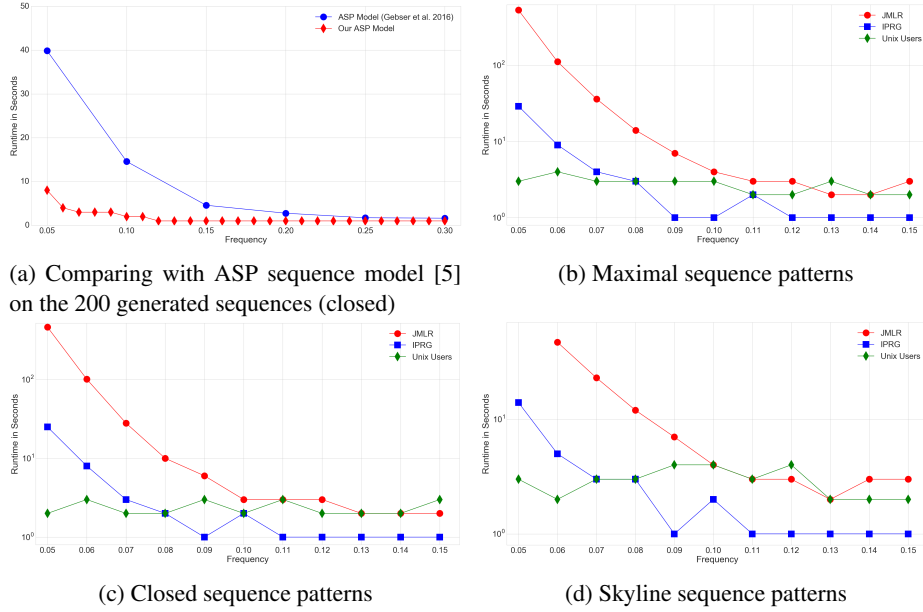


Fig. 1: Investigating  $\mathbf{Q}_1$ : comparison with pure ASP model (1a) and maximal (1b), closed (1c), and skyline (1d) sequence mining on JMLR, Unix Users, and iPRG datasets.

performs the purely declarative approach [5] and the advantage naturally becomes more apparent for smaller frequency threshold values.

In Figs. 1b, 1c and 1d (the point 0.05 for JMLR is a timeout), we present the runtimes of our method for *maximal*, *closed* and *skyline* sequential pattern mining settings on JMRL, Unix Users and iPRG datasets. In contrast to [5], our method managed to produce results on all of these datasets for reasonable threshold values within a couple of minutes.

To investigate  $\mathbf{Q}_2$ , we compare out-of-the-box performance of DP [16] with our approach on maximal, closed and skyline itemset mining problems using standard datasets Vote and Mushrooms. As we see in Figs. 2a and 2b, on average, DP is one-to-two orders of magnitude faster; this gap is, however, diminishing as the minimum frequency increases. Surprisingly, our approach is significantly faster than DP out-of-the-box for skyline patterns (Fig. 2c); this holds also for the Mushrooms dataset, not presented here.

Fine-tuning parameters of DP by changing the order in which operators are applied within the system (skyline+ option) allowed to close this gap. With this adaptation DP demonstrates one-to-two orders of magnitude better performance, as can be seen in Fig. 2c. However, fine-tuning such a system requires the understanding of its inner mechanisms or exhaustive application of all available options.

To address  $\mathbf{Q}_3$  we introduced three simple local constraints for the itemset mining setting from  $\mathbf{Q}_2$ : two size constraints  $size(I) > 2$  and  $size(I) < 7$  and a cost constraint:

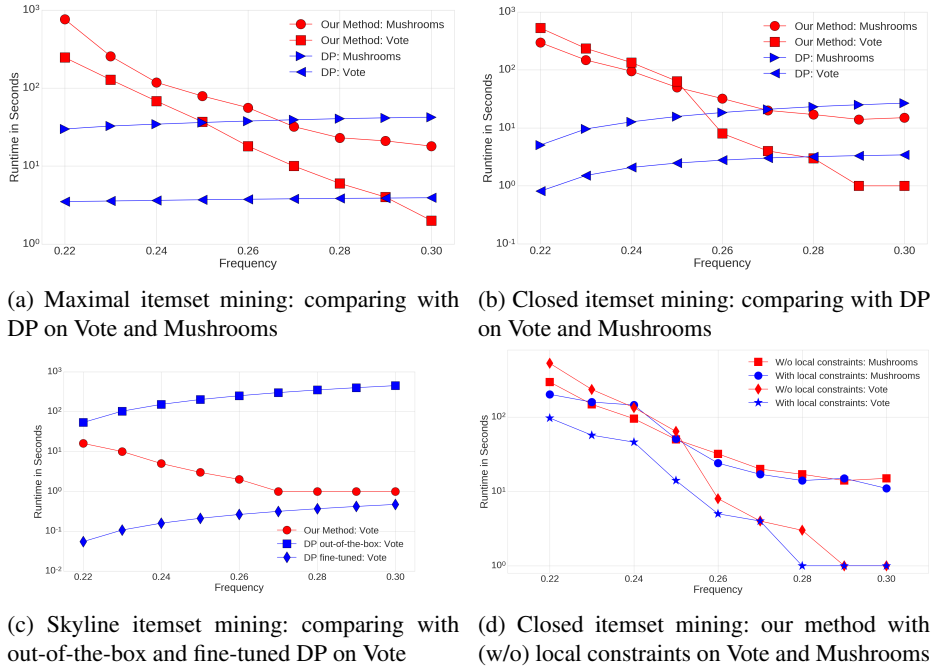


Fig. 2: Investigating  $Q_2$ : comparison with DP [16] (2b, 2a, 2c); and  $Q_3$ : the effect of local constraints on runtime (2d)

each item gets weight equal to its value with the maximal budget of  $n$ , which is set to 20 in the experiments.

In Fig. 2d, we present the results for closed itemset mining with and without local constraints (experiments with other global constraints demonstrate a similar runtime pattern and are not depicted here for space reasons). Local constraints ensure better propagation and speed up the search. One of the key design features of our encoding is the filtering technique used to select candidate patterns among only valid patterns. Its effect can be clearly seen, e.g., for the Vote dataset in Fig. 2d, where for certain frequencies the runtime gap is close to an order of magnitude.

In all experiments, Step 1 of our method contributes to less than 5% of runtime. Overall, our approach can handle real world datasets for sequential pattern mining as demonstrated in  $Q_1$ . In many cases its performance is close to the specialized mining languages, as shown in  $Q_2$ . Finally, as demonstrated in  $Q_3$  various local constraints can be effectively incorporated into our encoding bringing additional performance benefits.

## 5 Related Work

The problem of enhancing pattern mining by injecting various user-specified constraints has recently gained increasing attention. On the one hand, optimized dedicated approaches exist, in which some of the constraints are deeply integrated into the min-

ing algorithm, e.g., [19]. On the other hand, declarative methods based on Constraint Programming [21,17,15], SAT solving [12,11] and ASP [13,5,9] have been proposed.

Techniques from the last group are the closest to our work. However, in contrast to our method, they typically focus only on one particular pattern type and consider local constraints and condensed representations in isolation [20,23]. The works [16,7] focused on CP-based rather than ASP-based itemset mining and did not take into account sequences unlike we do. The authors of [5] studied declarative sequence mining with ASP, but in contrast to our approach, optimized algorithms for frequent pattern discovery are not exploited in their method. A theoretical framework for structured pattern mining was proposed in [8], whose main goal was to formally define the core components of the main mining tasks and compare dedicated mining algorithms to their declarative versions. While generic, this work did not take into account local and global constraints and neither has it been implemented.

In [13,5], purely declarative ASP methods have been considered; unlike our approach, they do not admit integration of optimized mining algorithms and thus lack practicality. In fact, the need for such an integration in the context of complex structured mining was even explicitly stated in [18] and in [10], which study formalizations of graph mining problems using logical means.

## 6 Conclusion

We have presented a hybrid approach for condensed itemset and sequence mining, which uses the optimized dedicated algorithms to determine the frequent patterns and post-filters them using a declarative ASP program. The idea of exploiting ASP for pattern mining is not new; it was studied for both itemsets and sequences. However, unlike previous methods we made steps towards optimizing the declarative techniques by making use of the existing specialized methods and also integrated the dominance programming machinery in our implementation to allow for combining local and global constraints on a generic level.

One of the possible future directions is to generalize the proposed approach into an iterative technique, where dedicated data mining and declarative methods are inter-linked and applied in an iterative fashion. More specifically, all constraints can be split into two parts: those that can be effectively handled using declarative means and those for which specialized algorithms are much more scalable. Answer set programs with external computations [4] possibly could be exploited in this mining context.

Another promising but challenging research stream concerns the integration of data *decomposition* techniques into our approach. Here, one can divide a given dataset into several parts, such that the frequent patterns are identified in these parts separately, and then the results are combined.

Orthogonal to this, materialization of the presented ideas on other pattern types including graphs and sequences of itemsets instead of sequences of individual symbols is an interesting future direction.

## References

1. Agrawal, R., Mannila, H., Srikant, R., Toivonen, H., Verkamo, A.I.: Fast discovery of association rules. In: *Advances in Knowledge Disc. and Data Mining*, pp. 307–328, 1996
2. Aoga, J.O.R., Guns, T., Schaus, P.: An efficient algorithm for mining frequent sequence with constraint programming. In: *ECML PKDD*. pp. 315–330 (2016)
3. Bonchi, F., Lucchese, C.: On condensed representations of constrained frequent patterns. *Knowl. Inf. Syst.* 9(2), 180–201 (2006)
4. Eiter, T., Brewka, G., Dao-Tran, M., Fink, M., Ianni, G., Krennwallner, T.: Combining non-monotonic knowledge bases with external sources. In: *Frontiers of Combining Systems, 7th International Symposium, FroCoS Italy, September 16-18*. pp. 18–42 (2009)
5. Gebser, M., Guyet, T., Quiniou, R., Romero, J., Schaub, T.: Knowledge-based sequence mining with ASP. In: *Proc. of 25th Int. Joint Conf. on AI, IJCAI* (2016)
6. Gelfond, M., Lifschitz, V.: The stable model semantics for logic programming. In: *Proc. of ICLP/SLP*. pp. 1070–1080, 1988
7. Guns, T., Dries, A., Nijssen, S., Tack, G., De Raedt, L.: Miningzinc: A declarative framework for constraint-based mining. *Artif. Intell.* 244, 6–29 (2017)
8. Guns, T., Paramonov, S., Négrevergne, B.: On declarative modeling of structured pattern mining. In: *Declarative Learning Based Programming, AAAI Workshop*, 2016
9. Guyet, T., Moinard, Y., Quiniou, R.: Using answer set programming for pattern mining. *CoRR abs/1409.7777* (2014)
10. van der Hallen, M., Paramonov, S., Leuschel, M., Janssens, G.: Knowledge representation analysis of graph mining. *CoRR abs/1608.08956* (2016)
11. Jabbour, S., Sais, L., Salhi, Y.: Boolean satisfiability for sequence mining. In: *Proc. of 22nd ACM Int. Conf. on Information and Knowledge Man., (CIKM)*, 649–658 (2013)
12. Jabbour, S., Sais, L., Salhi, Y.: Decomposition based SAT encodings for itemset mining problems. In: *PAKDD*. 662–674 (2015)
13. Jarvisalo, M.: Itemset mining as a challenge application for answer set enumeration. In: *Logic Prog. and Nonmon. Reas, LPNMR*, 304–310 (2011)
14. Mannila, H., Toivonen, H.: Levelwise search and borders of theories in knowledge discovery. *Data Min. Knowl. Discov.* 1(3), 241–258 (1997)
15. Métivier, J., Loudni, S., Charnois, T.: A constraint programming approach for mining sequential patterns in a sequence database. *CoRR abs/1311.6907* (2013)
16. Négrevergne, B., Dries, A., Guns, T., Nijssen, S.: Dominance programming for itemset mining. In: *2013 IEEE 13th Int. Conf. on Data Mining*, 557–566 (2013)
17. Négrevergne, B., Guns, T.: Constraint-based sequence mining using constraint programming. In: *CPAIOR*. pp. 288–305 (2015)
18. Paramonov, S., van Leeuwen, M., Denecker, M., De Raedt, L.: An exercise in declarative modeling for relational query mining. In: *ILP*. pp. 166–182 (2015)
19. Pei, J., Han, J.: Can we push more constraints into frequent pattern mining? In: *ACM SIGKDD*, 350–354 (2000)
20. Pei, J., Han, J., Mao, R.: CLOSET: an efficient algorithm for mining frequent closed itemsets. In: *ACM SIGMOD Workshop on Res. Issues in Data Min. Knowl. Discov*, pp. 21–30 (2000)
21. Rojas, W.U., Boizumault, P., Loudni, S., Crémilleux, B., Lepailleur, A.: Mining (soft-) sky-patterns using dynamic CSP. In: *CPAIOR*. pp. 71–87 (2014)
22. Simons, P., Niemelä, I., Soinen, T.: Extending and implementing the stable model semantics. *Artif. Intell.* 138(1-2), 181–234 (2002)
23. Yan, X., Han, J., Afshar, R.: Clospan: Mining closed sequential patterns in large datasets. In: *In SDM*. pp. 166–177 (2003)
24. Zaki, M.J., Parthasarathy, S., Ogihara, M., Li, W.: New algorithms for fast discovery of association rules. *Tech. rep., Rochester, NY, USA* (1997)