



**HAL**  
open science

# NegPSpan: efficient extraction of negative sequential patterns with embedding constraints

Thomas Guyet, René Quiniou

► **To cite this version:**

Thomas Guyet, René Quiniou. NegPSpan: efficient extraction of negative sequential patterns with embedding constraints. 2018. hal-01743975v1

**HAL Id: hal-01743975**

**<https://inria.hal.science/hal-01743975v1>**

Preprint submitted on 4 Apr 2018 (v1), last revised 25 Jul 2018 (v2)

**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.

# NegPSpan: efficient extraction of negative sequential patterns with embedding constraints

Thomas Guyet – Agrocampus-Ouest/IRISA UMR6074  
René Quiniou, Univ Rennes, Inria, CNRS, IRISA

April 4, 2018

## Abstract

Mining frequent sequential patterns consists in extracting recurrent behaviors, modeled as patterns, in a big sequence dataset. Such patterns inform about which events are frequently observed in sequences, *i.e.* what does really happen. Sometimes, knowing that some specific event does not happen is more informative than extracting a lot of observed events. Negative sequential patterns (NSP) formulate recurrent behaviors by patterns containing both observed events and absent events. Few approaches have been proposed to mine such NSPs. In addition, the syntax and semantics of NSPs differ in the different methods which makes it difficult to compare them. This article provides a unified framework for the formulation of the syntax and the semantics of NSPs. Then, we introduce a new algorithm, NEGPSpan, that extracts NSPs using a PrefixSpan depth-first scheme and enabling *maxgap* constraints that other approaches do not take into account. The formal framework allows for highlighting the differences between the proposed approach wrt to the methods from the literature, especially wrt the state of the art approach eNSP. Intensive experiments on synthetic and real datasets show that NEGPSpan can extract meaningful NSPs and that it can process bigger datasets than eNSP thanks to significantly lower memory requirements and better computation times.

## 1 Introduction

In many application domains such as diagnosis or marketing, decision makers show a strong interest for rules that associates specific events (a context) to undesirable events to which they are correlated or that are frequently triggered in such a context. Sequential pattern mining algorithms can extract such hidden rules from execution traces or transactions. In the classical setting, sequential patterns contain only positive events, *i.e.* really observed events. However, the absence of a specific action or event can often better explain the occurrence of an undesirable situation [2]. For example in diagnosis, if some maintenance operations have not been performed, *e.g.* damaged parts have not been replaced, then a fault should occur in a short delay while if these operations were performed in time the fault would not occur. In marketing, if some market-place customer has not received special offers or coupons for a long time then she/he has a high probability of churning while if she/he were provided such special offers she/he should remain loyal to her/his market-place. In these two cases, mining specific events, some present and some absent, to discover under which context some undesirable situation occur or not may provide interesting so-called *actionable* information

for determining which action should be performed to avoid the undesirable situation, *i.e.* fault in diagnosis, churn in marketing.

We aim at discovering sequential patterns that take into account the absence of some events called *negative events* [2]. Moreover, we want to take into account some aspect of the temporal dimension as well, maximal pattern span or maximal gap between the occurrences of pattern events. For example, suppose that from a sequence dataset, we want to mine a sequential pattern  $\mathbf{p} = \langle a b \rangle$  with the additional *negative* constraint telling that the event  $c$  should not appear between events  $a$  and  $b$  in  $\mathbf{p}$ . The corresponding negative pattern is represented as  $\mathbf{p} = \langle a \neg c b \rangle$ , where the logical sign  $\neg$  denotes an absent event or set of events. Once the general idea of introducing negative statements in a pattern has been stated, the syntax and semantics of such negative patterns should be clearly formulated since they have a strong impact both on algorithms outcome and their computational efficiency. As we will see, the few algorithms from literature do not use the same syntactical constraints and rely on very different semantics principles (see Section 2). More precisely, the efficiency of eNSP [1], the state-of-the-art algorithm for NSP mining, comes from a negation semantics that enables efficient operations on the sets of supported sequences. The two computational limits of eNSP are memory requirements and the impossibility for eNSP to handle embedding constraints such as the classical *maxgap* and *maxspan* constraints. When mining relatively long sequences (above 20 itemsets), such constraints appear semantically sound to consider short pattern occurrences where events are not too distant. In addition, such constraints can efficiently prune occurrence search.

This article provides two main contributions:

- we clarify the syntactic definition of negative sequential patterns and we provide different negation semantics with their properties.
- we propose NEGPSpan, an algorithm inspired by algorithm PrefixSpan to extract negative sequential patterns with maxgap constraints.

Intensive experiments compare, on synthetic and real datasets, the performance of NEGPSpan and eNSP as well as the pattern sets extracted by each of them. We show that algorithm NEGPSpan is more efficient than eNSP for mining long sequences thanks to the maxgap constraint and that its memory requirement is several orders of magnitude lower, enabling to process much larger datasets. In addition, we highlight that eNSP misses interesting patterns on real datasets due to semantic restrictions.

## 2 Related work

Kamepalli et al. provide a survey of the approaches proposed for mining negative patterns [6]. The three most significant algorithms appear to be PNSP, Neg-GSP and eNSP. We briefly review each of them in the following paragraphs.

PNSP (Positive and Negative Sequential Patterns mining) [5] is the first algorithm proposed for mining full negative sequential patterns where negative itemsets are not only located at the end of the pattern. PNSP extends algorithm GSP to cope with mining negative sequential patterns. PNSP consists of three steps: i) mine frequent positive sequential patterns, by using algorithm GSP, ii) preselect negative sequential itemsets – for PNSP, negative itemsets must not be too infrequent (should have a support less than a threshold *miss\_freq*) – iii) generate candidate negative sequences levelwise and scan the sequence dataset again to compute the support of these candidates and prune

the search when the candidate is infrequent. This algorithm is incomplete: the second parameter reduces the set of potential negative itemsets. Moreover, the pruning strategy of PNSP is not correct [13] and PNSP misses potentially frequent negative patterns.

Zheng et al. [13] also proposed a negative version of algorithm GSP, called Neg-GSP, to extract negative sequential patterns. They showed that traditional Apriori-based negative pattern mining algorithms relying on support anti-monotonicity have two main problems. The first one is that the Apriori principle does not apply to negative sequential patterns. They gave an example of sequence that is frequent even if one of its sub-sequence is not frequent. The second problem has to do with the efficiency and the effectiveness of finding frequent patterns due to a vast candidate space. Their solution was to prune the search space using the support anti-monotonicity over positive parts. This pruning strategy is correct but is not really efficient considering the huge number of remaining candidates whose support has to be evaluated. We will see in Section 3.2 that anti-monotonicity can be defined considering an order relation based on common prefixes.

eNSP (efficient NSP) has been recently proposed by Cao et al. [16]. It identifies NSPs by computing only frequent positive sequential patterns and deducing negative sequential patterns from positive patterns. Precisely, Cao et al. showed that the support of some negative pattern can be computed by arithmetic operations on the support of its positive sub-patterns, thus avoiding additional sequence database scans to compute the support of negative patterns. However, this necessitates to store all the (positive) sequential patterns with their set of covered sequences (tid-lists) which may be impossible in case of big dense datasets and low minimal support thresholds. This approach makes the algorithm more efficient but it hides some restrictive constraints on the extracted patterns. First, a frequent negative pattern whose so-called positive partner (the pattern where all negative events have been switched to positive) is not frequent will not be extracted. Second, every occurrences of a negative pattern in a sequence should satisfy absence constraints. We call this *strong absence semantics* (see Section 3.2). These features lead eNSP to extract less patterns than previous approaches. In some practical applications, eNSP may miss potentially interesting negative patterns from the dataset.

The first constraint has been partly tackled by Dong et al. with algorithm eNSPFI, an extension of eNSP which mines NSPs from frequent and some infrequent positive sequential patterns from the negative border [4]. E-msNSP [12] is another extension of eNSP which uses multiple minimum supports: an NSP is frequent if its support is greater than a local minimal support threshold computed from the content of the pattern and not a global threshold as in classical approaches. A threshold is associated with each item, and the minimal support of a pattern is defined from the most constrained item it contains. Such kind of adaptive support prevents from extracting some useless patterns still keeping the pattern support anti-monotonic. The same authors also proposed high utility negative sequential patterns based on the same principles [11]. It is worth noting that these algorithm relies basically on the same principles of eNSP and so, present the same features, strong memory requirements, strong absence semantics for negation.

### 3 Negative sequential patterns

This section introduces sequential patterns and negative sequential pattern mining. First we remind some basic definitions about sequences of itemsets, and classical sequential pattern mining then we introduce some definitions of negative sequential patterns.

In the sequel,  $[n] = \{1, \dots, n\}$  denotes the set of the first  $n$  strictly positive integers. Let  $(\mathcal{I}, <)$  be the set of items (alphabet) associated with a total order (e.g. lexicographic order). An *itemset*

$A = \{a_1 a_2 \dots a_n\} \subset \mathcal{I}$  is a set of ordered items. A *sequence*  $\mathbf{s}$  is a set of sequentially ordered itemsets  $\mathbf{s} = \langle s_1 s_2 \dots s_n \rangle: \forall i, j \in [n], i < j$  means that  $s_i$  is located before  $s_j$  in sequence  $\mathbf{s}$  which starts by  $s_1$  and finishes by  $s_n$ . Mining sequential patterns from a dataset of sequences, denoted  $\mathcal{D}$ , consists in extracting the frequent subsequences (patterns) included in database sequences having a support (*i.e.* the number of sequences in which the pattern occurs) greater than a given threshold  $\sigma$ . There is a huge literature about sequential pattern mining. We will not go into details and refer the reader to a survey of the literature, such as Mooney et al. [7].

Negative sequential patterns (NSP) extend classical sequential patterns by enabling the specification of absent itemsets. For example,  $\mathbf{p} = \langle a b \neg c e f \rangle$  is a negative pattern. The symbol  $\neg$  before  $c$  denotes that  $c$  is a negative itemset (here reduced to an item). Semantically,  $\mathbf{p}$  specifies that events  $a$  and  $b$  happen in a row, then events  $e$  and  $f$  occur in a row, but event  $c$  does not occur between  $c$  and  $e$  occurrences.

In the field of string matching, negation is classically defined for regular expression. In this case, a pattern is an expression that can hold any kind of negated *pattern*. The same principle gives the following most generic definition of negative sequential patterns:

**Definition 1** (Most generic negative sequential patterns). Let  $\mathcal{N}$  be the set of negative patterns. A negative pattern  $\mathbf{p} = \langle p_1 \dots p_n \rangle \in \mathcal{N}$  is a sequence where  $\forall i, p_i$  is a positive itemset ( $p_i \subset \mathcal{I}$ ) or a negated pattern ( $p_i = \neg\{q_i\}, q_i \in \mathcal{N}$ ).

Due to its recursive definition,  $\mathcal{N}$  appears to be too huge to be an interesting and tractable search space for pattern mining. For instance, with  $\mathcal{I} = \{a, b, c\}$ , it is possible to express simple patterns like  $\langle a \neg b c \rangle$  but also complex patterns like  $\langle a, \neg\langle b, c \rangle \rangle$ . The combinatorics for such patterns is infinite. Also, recursive patterns, such as the last one, are really difficult to interpret by experts.

As a consequence, some restrictions are required to ease the design of algorithms that could extract a reasonable number of understandable patterns. No semantics is “more” correct or relevant than another one. It depends on the information to be captured. Our objective is to give the opportunity to data scientists to make an informed choice and also to highlight that holding on execution time only is not necessarily relevant.

Moreover, if these syntactic restrictions appear to be important to clearly state the mining task, they must also be taken into account to describe precisely the semantics of negation in a sequence of itemsets. In section 3.2, we formulate and illustrate some alternative semantics and their consequences on pattern support evaluation.

### 3.1 Syntactic constraints for negative sequential patterns

Just above, we have introduced a generic definition for sequential patterns with negation. Now, we provide a more specific definition of negative sequential patterns (NSP) which introduces some syntactic constraints that are broadly used in the literature [6].

**Definition 2** (Negative sequential patterns (NSP)). A negative pattern  $\mathbf{p} = \langle p_1 \dots p_n \rangle$  is a sequence where  $\forall i, p_i$  is a positive itemset ( $p_i = \{p_i^j\}, p_i^j \in \mathcal{I}$ ) or a negated itemset ( $p_i = \neg\{q_i^j\}, q_i^j \in \mathcal{I}$ ). The *positive part*<sup>1</sup> of pattern  $\mathbf{p}$ , denoted  $\mathbf{p}^+$ , is the subsequence of  $\mathbf{p}$  that holds only its positive itemsets.

---

<sup>1</sup>Called the *maximal positive subsequence* in PNSP and Neg-GSP or the *positive element id-set* in eNSP.

According to definition 2, negative sequential patterns consist of positive or negative itemsets but an itemset consists of positive items only. Moreover, the following assumptions on NSP syntax are classically used to reduce the pattern search space:

- (A1) consecutive negative itemsets are forbidden since it is difficult to distinguish successive negations as in  $\langle a \neg b \neg c d \rangle$  from conjunctive negations as in  $\langle a \neg(b c) d \rangle$  (no  $b$  neither  $c$  between  $a$  and  $d$ ).
- (A2) a negated item is not surrounded by itemsets containing this item. This simplification is motivated by the objective to simplify pattern understanding. A pattern  $\langle a \neg b c \rangle$  may be interpreted as “there is exactly one occurrence of  $b$  between  $a$  and  $c$ ”. But, this leads to redundant patterns:  $\langle ab \neg bc \rangle$  matches exactly the same sequences than  $\langle a \neg bbc \rangle$ .
- (A3) negated itemsets are restricted to a subset of possible itemsets, denoted by some language  $\mathcal{L}^-$ . This specification changes from one proposal to another. For example,  $\mathcal{L}^-$  could be the frequent itemsets or the itemsets built from the frequent items or any arbitrary subset of itemsets.

It is worth-noting that these assumptions are not mandatory to state the formal properties of most algorithms. They improve the efficiency by reducing the search space and also prevent the generation of numerous patterns that would be useless or difficult to interpret by the user.

### 3.2 Semantics of negative sequential patterns

The semantics of negative sequential patterns relies on *negative containment*: a sequence  $s$  supports pattern  $p$  if  $s$  contains a sub-sequence  $s'$  such that every positive itemset of  $p$  is included in some itemset of  $s'$  in the same order and for any negative itemset  $\neg i$  of  $p$ ,  $i$  is *not included* in any itemset occurring in the sub-sequence of  $s'$  located between the occurrence of the positive itemset preceding  $\neg i$  in  $p$  and the occurrence of the positive itemset following  $\neg i$  in  $p$ .

So far in the literature, the absence or non-inclusion of negative itemsets has been specified by loose formulations. The authors of PNSP have proposed the set symbol  $\not\subseteq$  to specify non-inclusion. This symbol is misleading since it does not correspond to the associated semantics given in PNSP: an itemset  $I$  is absent from an itemset  $I'$  if all items in  $I$  are absent altogether from  $I'$  which is covered by  $I \cap I' = \emptyset$  and not  $I \not\subseteq I'$ . We will call PNSP interpretation *total non inclusion*. It should be distinguished from *partial non inclusion* which corresponds (correctly) to the set symbol  $\not\subset$ . The symbol  $\not\subseteq$  was further used by the authors of Neg-GSP and eNSP. The semantics of non inclusion is not detailed in Neg-GSP and one cannot determine if it means total or partial non inclusion<sup>2</sup>. eNSP does not define explicitly the semantics of non inclusion but, from the procedure used to compute the support of patterns, one can deduce that it uses total non inclusion.

**Definition 3** (non inclusion). We introduce two operators relating two itemsets  $P$  and  $I$ :

- partial non inclusion:  $P \not\subset I \Leftrightarrow \exists e \in P, e \notin I$
- total non inclusion:  $P \not\subseteq I \Leftrightarrow \forall e \in P, e \notin I$

---

<sup>2</sup>Actually, though not clearly stated, it seems that the negative elements of Neg-GSP patterns consist of items rather than itemsets. In this case, total and partial inclusion are identical.

Table 1: Support in  $\mathcal{D}$  of increasingly extended patterns under the total and partial non inclusion semantics

	total non inclusion	partial non inclusion
$\mathbf{ns}_1 = \langle b \neg ca \rangle$	1, 3, 4	1, 3, 4
$\mathbf{ns}_2 = \langle b \neg (cd) a \rangle$	1, 2, 3, 4	1, 4
$\mathbf{ns}_3 = \langle b \neg (cde) a \rangle$	1, 2, 3, 4	1
$\mathbf{ns}_4 = \langle b \neg (cdeg) a \rangle$	1, 2, 3, 4, 5	1
	monotonic	anti monotonic

Choosing one non inclusion interpretation or the other has consequences on extracted patterns as well as on pattern search. Let's illustrate this on related pattern support in the sequence dataset  $\mathcal{D} = \{s_1 = \langle (bc) f a \rangle, s_2 = \langle (bc) (cf) a \rangle, s_3 = \langle (bc) (df) a \rangle, s_4 = \langle (bc) (ef) a \rangle, s_5 = \langle (bc) (cdef) a \rangle\}$ . Table 1 compares the support of progressively extended patterns under the two semantics to show whether anti-monotonicity is respected or not.

Obviously, partial non inclusion satisfies anti-monotonicity while total non inclusion does not. In the sequel we will denote the general form of itemset non inclusion by the symbol  $\not\subseteq$ , meaning either  $\not\subseteq$  or  $\not\supseteq$ .

Now, we formulate the notions of sub-sequence, non inclusion and absence by means of the concept of embedding.

**Definition 4** (positive pattern embedding). Let  $\mathbf{s} = \langle s_1 \dots s_n \rangle$  be a sequence and  $\mathbf{p} = \langle p_1 \dots p_m \rangle$  be a (positive) sequential pattern.  $\mathbf{e} = (e_i)_{i \in [m]} \in [n]^m$  is an *embedding* of pattern  $\mathbf{p}$  in sequence  $\mathbf{s}$  iff  $\forall i \in [m], p_i \subseteq s_{e_i}$  and  $\forall i \in [m-1], e_i < e_{i+1}$

**Definition 5** (Strict and soft embeddings of negative patterns). Let  $\mathbf{s} = \langle s_1 \dots s_n \rangle$  be a sequence and  $\mathbf{p} = \langle p_1 \dots p_m \rangle$  be a negative sequential pattern.

$\mathbf{e} = (e_i)_{i \in [m]} \in [n]^m$  is a **soft-embedding** of pattern  $\mathbf{p}$  in sequence  $\mathbf{s}$  iff  $\forall i \in [m]$ :

- $p_i \subseteq s_{e_i}$  if  $p_i$  is positive
- $\forall j \in [e_{i-1} + 1, e_{i+1} - 1], p_i \not\subseteq s_j$  if  $p_i$  is negative

$\mathbf{e} = (e_i)_{i \in [m]} \in [n]^m$  is a **strict-embedding** of pattern  $\mathbf{p}$  in sequence  $\mathbf{s}$  iff for all  $i \in [m]$ :

- $p_i \subseteq s_{e_i}$  if  $p_i$  is positive
- $p_i \not\subseteq \bigcup_{j \in [e_{i-1} + 1, e_{i+1} - 1]} s_j$  if  $p_i$  is negative

**Proposition 1.** *soft-* and *strict-*embeddings are equivalent when  $\not\subseteq \stackrel{\text{def}}{=} \not\supseteq$ .<sup>3</sup>

Let  $\mathbf{p}^+ = \langle p_{k_1} \dots p_{k_l} \rangle$  be the positive part of some pattern  $\mathbf{p}$ , where  $l$  denotes the number of positive itemsets in  $\mathbf{p}$ . If  $\mathbf{e}$  is an embedding of pattern  $\mathbf{p}$  in some sequence  $\mathbf{s}$ , then  $\mathbf{e}^+ = \langle e_{k_1} \dots e_{k_l} \rangle$  is an embedding of the positive sequential pattern  $\mathbf{p}^+$  in  $\mathbf{s}$ .

The following examples illustrate the impact of itemset non-inclusion operator and of embedding type.

<sup>3</sup>Proofs and additional properties are given in extended version available online.

Table 2: Comparison of negative pattern mining proposals. Optional constraints are specified in *Italic*.

	PNSP [5]	NegGSP [13]	eNSP [1]	NEGSPAN
<b>negative elements</b>	itemsets	items?	itemsets	itemsets
<b>itemset non inclusion</b>	strict	strict?	strict	strict/ <i>soft</i>
<b>itemset absence</b>	weak	weak	strong	weak
<b>constraints on negative itemsets</b>	not too infrequent ( <i>supp</i> $\leq$ <i>less_freq</i> )	frequent items	positive partner is frequent	frequent items, <i>bounded size</i>
<b>global constraints on patterns</b>	positive part is frequent	positive part is frequent	positive part is frequent	positive part is frequent, <i>maxspan</i> , <i>maxgap</i>

**Example 1** (Itemset absence semantics). Let  $\mathbf{p} = \langle a \neg(bc) d \rangle$  be a pattern and  $\mathbf{s}_1 = \langle a c b e d \rangle$ ,  $\mathbf{s}_2 = \langle a (bc) e d \rangle$ ,  $\mathbf{s}_3 = \langle a b e d \rangle$  and  $\mathbf{s}_4 = \langle a e d \rangle$  be four sequences. One can notice that each sequence contains a unique occurrence of  $\langle a d \rangle$ , the positive part of pattern  $\mathbf{p}$ . Using soft-embeddings and total non-inclusion ( $\not\subseteq^{\text{def}} \not\subseteq$ ),  $\mathbf{p}$  occurs in  $\mathbf{s}_1$ ,  $\mathbf{s}_3$  and  $\mathbf{s}_4$  but not in  $\mathbf{s}_2$ . Using the strict-embedding semantics and total non-inclusion,  $\mathbf{p}$  occurs in sequence  $\mathbf{s}_3$  and  $\mathbf{s}_4$  considering that items  $b$  and  $c$  occur between occurrences of  $a$  and  $d$  in sequences 1 and 2. Note that the order of  $b$  and  $c$  in  $\mathbf{s}_1$  has no importance for strict-embedding, since the union of intermediate itemsets is considered.

With partial non-inclusion ( $\not\subseteq^{\text{def}} \not\subseteq$ ) and either type of embeddings, the absence of an itemset is satisfied if any of its item is absent. As a consequence,  $\mathbf{p}$  occurs only in sequence  $\mathbf{s}_4$ .

Another point that determines the semantics of negative containment concerns the multiple occurrences of some pattern in a sequence: should every or only one occurrence of the pattern positive part in the sequence satisfy the non inclusion constraints? This point is not discussed in previous propositions for negative sequential pattern mining. Actually, PNSP and Neg-GSP require a weak absence (at least one occurrence should satisfy the non inclusion constraints) while eNSP requires a strong absence (every occurrence should satisfy non inclusion constraints).

**Definition 6** (Negative pattern occurrence). Let  $\mathbf{s}$  be a sequence,  $\mathbf{p}$  be a negative sequential pattern, and  $\mathbf{p}^+$  the positive part of  $\mathbf{p}$ .

- Pattern  $\mathbf{p}$  *softly-occurs* in sequence  $\mathbf{s}$ , denoted  $\mathbf{p} \preceq \mathbf{s}$ , iff there exists at least one (strict/soft)-embedding of  $\mathbf{p}$  in  $\mathbf{s}$ .
- Pattern  $\mathbf{p}$  *strictly-occurs* in sequence  $\mathbf{s}$ , denoted  $\mathbf{p} \sqsubseteq \mathbf{s}$ , iff for any embedding  $\mathbf{e}'$  of  $\mathbf{p}^+$  in  $\mathbf{s}$  there exists an embedding  $\mathbf{e}$  of  $\mathbf{p}$  in  $\mathbf{s}$  such that  $\mathbf{e}' = \mathbf{e}^+$ .

Definition 6 allows for formulating two notions of absence semantics for negative sequential patterns depending on the occurrences of the positive part:

- *strict occurrence/strong absence*: a negative pattern  $\mathbf{p}$  occurs in a sequence  $\mathbf{s}$  iff there exists at least one occurrence of the positive part of pattern  $\mathbf{p}$  in sequence  $\mathbf{s}$  and **every** such occurrence satisfies the negative constraints,



- *soft occurrence/weak absence*: a negative pattern  $\mathbf{p}$  occurs in a sequence  $\mathbf{s}$  iff there exists at least one occurrence of the positive part of pattern  $\mathbf{p}$  in sequence  $\mathbf{s}$  and **one** of these occurrences satisfies the negative constraints.

**Example 2** (Strong vs weak absence semantics). Let  $\mathbf{p} = \langle a b \neg c d \rangle$  be a pattern and  $\mathbf{s}_1 = \langle a b e d \rangle$  and  $\mathbf{s}_2 = \langle a b c a d e b d \rangle$  be two sequences. The positive part of  $\mathbf{p}$  is  $\langle a b d \rangle$ . It occurs once in  $\mathbf{s}_1$  so there is no difference for occurrences under the two semantics. But, it occurs twice in  $\mathbf{s}_2$  with embeddings (1, 2, 5) and (4, 7, 8). The first occurrence does not satisfy the negative constraint ( $\neg c$ ) while the second occurrence does. Under the weak absence semantics, pattern  $\mathbf{p}$  occurs in sequence  $\mathbf{s}_2$  whereas under the strong absence semantics it does not.

We also introduce **constrained negative sequential patterns**. We consider the two most common anti-monotonic constraints on sequential patterns: *maxgap* ( $\theta \in \mathbb{N}$ ) and *maxspan* ( $\tau \in \mathbb{N}$ ) constraints. These constraints impact NSP embeddings. An embedding  $\mathbf{e}$  of a pattern  $\mathbf{p}$  in some sequence  $\mathbf{s}$  satisfies the maxgap (resp. maxspan) constraint iff  $\mathbf{e}^+ = \{e_i, \dots, e_n\}$ , the embedding of the positive part of  $\mathbf{p}$  satisfies the constraint, *i.e.*  $\forall i \in [n-1], e_{i+1} - e_i \leq \theta$  (resp.  $e_n - e_1 \leq \tau$ ).

The definitions of pattern support, frequent pattern and pattern mining task derives naturally from the notion of occurrence of a negative sequential pattern, no matter the choices for embedding (soft or strict), non inclusion (partial or total) and absence (weak or strong). However, these choices concerning the semantics of NSPs impact directly the number of frequent patterns (under the same minimal threshold) and further the computation time. The stronger the negative constraints, the lesser the number of sequences that hold some pattern, and the lesser the number of frequent patterns.

Finally, we introduce two partial orders on NSPs that will be useful for designing an efficient NSP mining algorithm.

**Definition 7** (NSP partial orders).  $\triangleleft$  and  $\blacktriangleleft$  are two partial orders on the set of NSPs. Let  $\mathbf{p} = \langle p_1 \dots p_n \rangle$  and  $\mathbf{q} = \langle q_1 \dots q_m \rangle$  be two NSPs,

- $\mathbf{p} \triangleleft \mathbf{q}$  iff  $n \leq m$ ,  $\forall i \in [n]$ ,  $p_i$  negative  $\Leftrightarrow q_i$  negative and  $\forall i \in [n]$ ,  $p_i$  positive  $\Rightarrow p_i \subseteq q_i$ , and  $\forall i \in [n]$ ,  $p_i \subseteq q_i$ .
- $\mathbf{p} \blacktriangleleft \mathbf{q}$  iff  $n \leq m$ ,  $\forall i \in [n]$ ,  $p_i$  negative  $\Leftrightarrow q_i$  negative and  $\forall i \in [n]$ ,  $p_i$  positive  $\Rightarrow p_i \subseteq q_i$ , and  $p_i$  negative  $\Rightarrow q_i \subseteq p_i$ .

Note that the definition of order  $\blacktriangleleft$  specifies a reverse subset inclusion of negative patterns compared to  $\triangleleft$ . The reason is that depending on itemset non-inclusion (see Definition 3), the support has not the same monotonicity property.

**Proposition 2** (Anti-monotonicity of NSP support). The support of NSP is anti-monotonic with respect to  $\triangleleft$  (resp.  $\blacktriangleleft$ ) when  $\mathcal{Z} \stackrel{\text{def}}{=} \mathcal{Z}$  (resp.  $\mathcal{Z} \stackrel{\text{def}}{=} \mathcal{Z}$ ) is considered.

To conclude this section on formal aspects of negative pattern mining, we provide in Table 2 a comparison of several negative sequential pattern mining approaches wrt several features investigated in this section.

## 4 Algorithm NEGSPAN

In this section, we introduce algorithm NEGSPAN for mining NSPs from a sequence database under maxgap and maxspan constraints and under a weak absence semantics with  $\not\subseteq \stackrel{\text{def}}{=} \not\subseteq$  for itemset inclusion. As stated in proposition 1, no matter the embedding strategy, they are equivalent under strict itemset inclusion. Let  $\mathcal{L}^-$  denote the set of itemsets that can be built from frequent items.

### 4.1 Main algorithm

NEGSPAN is based on algorithm PrefixSpan [9] which implements a depth first search and uses the principle of database projection to reduce the number of sequence scans. NEGSPAN adapts the pseudo-projection principle of PrefixSpan which uses a projection pointer to avoid copying the data. For NEGSPAN, a projection pointer of some pattern  $p$  is a triple  $\langle sid, ppred, pos \rangle$  where  $sid$  is a sequence identifier in the database,  $pos$  is the position in sequence  $sid$  that matches the last itemset of the pattern (necessarily positive) and  $ppred$  is the position of the previous positive pattern.

Algorithm 1 details the main recursive function of NEGSPAN for extending a current pattern  $p$ . The principle of this function is similar to PrefixSpan. Every pattern  $p$  is associated with a pseudo-projected database represented by both the original set of sequences  $\mathcal{S}$  and a set of projection pointers  $occs$ . First, the function evaluates the size of  $occs$  to determine whether pattern  $p$  is frequent or not. If so, it is outputted, otherwise, the recursion is stopped because no larger patterns are possible (anti-monotonicity property).

Then, the function tries three types of pattern extensions:

- the positive sequence composition consists in adding one item to the last itemset of  $p$ ,
- the positive sequence extension consists in adding a new positive singleton itemset at the end of  $p$ ,
- the negative sequence extension consists in inserting a negative itemset between the positive penultimate itemset of  $p$  and the last positive itemset of  $p$ . According to the chosen NSP syntactic restrictions, this extension is possible only when the pattern length is greater than or equal to 2.

The negative pattern extension is specific to our algorithm and is detailed in the next section. The first two extensions are identical to PrefixSpan pattern extensions, including their gap constraints management, *i.e.* maxgap and maxspan constraints between positive patterns. For the sake of space limitation, positive extensions are not detailed in the sequel.

The proposed algorithm is correct and complete. The proposed algorithm is also consistent with order  $\triangleleft$  on NSP and, according to proposition 2, NSP support is anti-monotonic. More specifically, the negative composition by an item is anti-monotonic under the strict-embedding semantics <sup>4</sup>.

### 4.2 Extension of patterns with negated itemsets

Algorithm 2 extends the current pattern  $p$  with negative items. It generates new candidates by inserting an item  $it \in \mathcal{I}^f$ , the set of frequent items. Let  $p[-2]$  and  $p[-1]$  denote respectively the

---

<sup>4</sup>A similar algorithm extracts NSPs with a non occurrence based on operator  $\not\subseteq$ .  $\triangleleft$  order forbids the incremental extension of negative itemsets. Thus, all candidate itemsets, *i.e.*  $\mathcal{L}^-$ , are computed and evaluated without recursions.

---

**Algorithm 1: NEGPSPAN: recursive function for negative sequential pattern extraction**

---

**input:**  $\mathcal{S}$ : set of sequences,  $p$ : current pattern,  $\sigma$ : minimum support threshold,  $occs$ : list of occurrences,  $\mathcal{I}^f$ : set of frequent items,  $\theta$ : maxgap,  $\tau$ : maxspan

```
1 Function NEGPSPAN( $\mathcal{S}, \sigma, p, occs, \mathcal{I}^f, \theta, \tau$ ):
  //Support evaluation of pattern  $p$ 
2  if  $|occs| \geq \sigma$  then
3     $\lfloor$  OutputPattern( $p, occs$ );
4  else
5     $\lfloor$  return;
  //Positive itemset composition
6  PositiveComposition( $\mathcal{S}, \sigma, p, occs, \mathcal{I}^f, \theta, \tau$ );
  //Positive sequential extension
7  PositiveSequence( $\mathcal{S}, \sigma, p, occs, \mathcal{I}^f, \theta, \tau$ );
8  if  $|p| \geq 2$  then
9     $\lfloor$  //Negative sequential extension
      NegativeExtension( $\mathcal{S}, \sigma, p, occs, \mathcal{I}^f, \theta, \tau$ );
```

---

---

**Algorithm 2: NEGPSPAN: negative extensions**

---

**input:**  $\mathcal{S}$ : set of sequences,  $p$ : current pattern,  $\sigma$ : minimum support threshold,  $occs$ : list of occurrences,  $\mathcal{I}^f$ : set of frequent items,  $\theta$ : maxgap,  $\tau$ : maxspan

```
1 Function NegativeExtension( $\mathcal{S}, \sigma, p, occs, \mathcal{I}^f, \theta, \tau$ ):
2   for  $it \in \mathcal{I}^f$  do
3     if  $p[-2]$  is pos then
4        $\lfloor$  //Insert the negative item at the penultimate position
           $p.insert(\neg it)$ ;
5     else
6       if  $it > p[-2].back()$  then
7          $\lfloor$  //Insert an item to the penultimate (negative) itemset
           $p[-2].insert(\neg it)$ ;
8       else
9          $\lfloor$  continue;
10     $newoccs \leftarrow \emptyset$ ;
11    for  $occ \in occs$  do
12       $found \leftarrow false$ ;
13      for  $sp = [occ.pred + 1, occ.pos - 1]$  do
14        if  $it \in s_{occ.sid}[sp]$  then
15           $\lfloor$   $found \leftarrow true$ ;
16           $\lfloor$  break;
17      if  $!found$  then
18         $\lfloor$   $newoccs \leftarrow newoccs \cup \{occ\}$ ;
19      else
20         $\lfloor$  //Look for an alternative occurrence
           $newoccs \leftarrow newoccs \cup Match(s_{sid}, p, \theta, \tau)$ ;
21    NEGPSPAN( $\mathcal{D}, \sigma, p, newoccs, \mathcal{I}^f$ );
22     $p[-2].pop()$ ;
```

---

penultimate itemset and the last itemset of  $p$ . If  $p[-2]$  is positive, then a new negated itemset is inserted between  $p[-2]$  and  $p[-1]$ . Otherwise, if  $p[-2]$  is negative, item  $it$  is added to  $p[-2]$ . To prevent redundant enumeration of negative itemsets, only items  $it$  (lexicographically) greater than the last item of  $p[-2]$  can be added.

Then, lines 10 to 20, evaluate the candidate by computing the pseudo-projection of the current database. According to the selected semantics associated with  $\sqsubseteq$ , *i.e.* total non inclusion (see Definition 5), it is sufficient to check the absence of  $it$  in the subsequence included between the occurrences of positive itemsets surrounding  $it$ . To achieve this, the algorithm checks the sequence positions in the interval  $[occ.ppred+1, occ.pos-1]$ . If  $it$  does not occur in itemsets from this interval, then the extended pattern occurs in the sequence  $occ.sid$ . Otherwise, to ensure the completeness of the algorithm, another occurrence of the pattern has to be searched in the sequence (*cf.* `Match` function that takes into account gap constraints).

For example, the first occurrence of pattern  $\mathbf{p} = \langle abc \rangle$  in sequence  $\langle abecabc \rangle$  is  $occ_p = \langle sid, 2, 4 \rangle$ . Let's now consider  $\mathbf{p}' = \langle ab-ec \rangle$ , a negative extension of  $p$ . The extension of the projection-pointer  $occ_p$  does not satisfy the absence of  $e$ . So a new occurrence of  $p$  has to be searched for.  $\langle sid, 6, 7 \rangle$ , the next occurrence of  $\mathbf{p}$ , satisfies the negative constraint. Then, `NEGPSPAN` is called recursively for extending the new current pattern  $\langle ab-ec \rangle$ .

## 5 Experiments

This section presents experiments on synthetic and real data. Experiments on synthetic data aims at exploring and comparing `NEGPSPAN` and `eNSP` for negative sequential pattern mining. The other experiments were conducted on medical care pathways and illustrates results for negative patterns. `NEGPSPAN` and `eNSP` have been implemented in C++. We focus our presentation on the most significant results. More detailed results can be found in a companion website<sup>5</sup>.

### 5.1 Benchmark

This section presents experiments on synthetically generated data. The principle of our sequence generator is the following: generate random negative patterns and hide or not some of their occurrences inside randomly generated sequences. The main parameters are the total number of sequences ( $n$ , default value is  $n = 500$ ), the mean length of sequences ( $l = 20$ ), the number of different items ( $d = 20$ ), the total number of patterns to hide (3), their mean length (4) and the minimum occurrence frequency of patterns in the dataset (10%).

Generated sequences are sequences of items (not itemsets). For such kind of sequences, patterns extracted by `eNSP` hold only items because positive partners have to be frequent. For a fair evaluation and preventing `NEGPSPAN` from generating more patterns, we restricted  $\mathcal{L}^-$  to the set of frequent items. For both approaches, we limit the pattern length to 5 items.

Figure 1 illustrates the computation time and number of patterns extracted by `eNSP` and `NEGPSPAN` on sequences of length 20 and 30, under three minimal thresholds ( $\sigma = 10\%$ ,  $15\%$  and  $20\%$ ) and with different values for the maxgap constraint ( $\tau = 4, 7, 10$  and  $\infty$ ). For `eNSP`, the minimal support of positive partners is set as 80% of the minimal threshold  $f$ . Each boxplot has been obtained with a 20 different sequence datasets. Each run has a timeout of 5 minutes.

---

<sup>5</sup>Code, data generator and synthetic benchmark datasets can be downloaded here: <http://people.irisa.fr/Thomas.Guyet/negativepatterns/>.

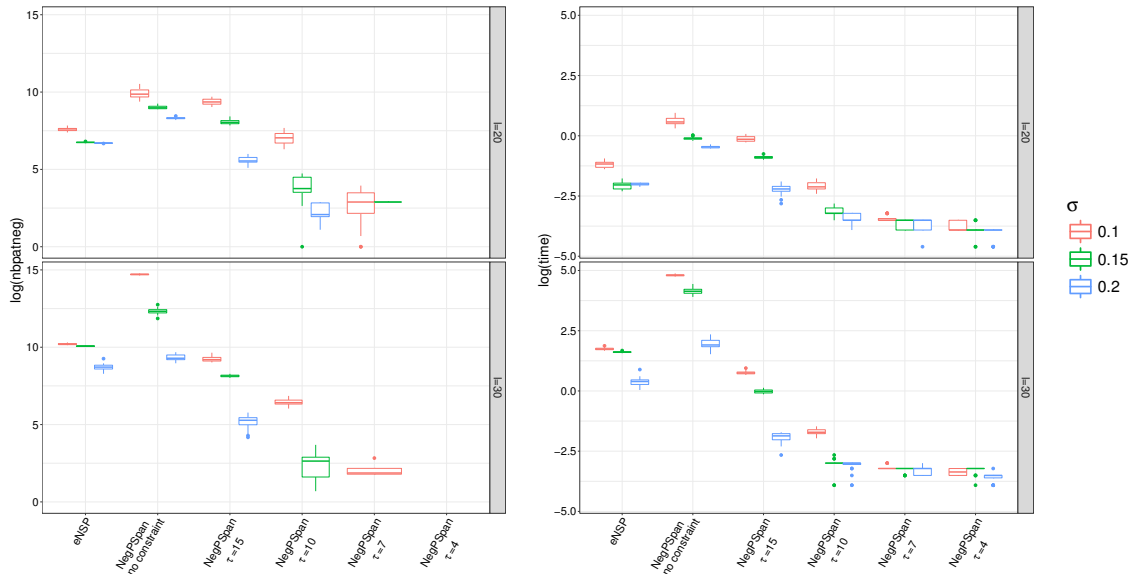


Figure 1: Comparison of number of patterns (left) and computing time (right) between eNSP and NEGSPAN, with different values for maxgap ( $\tau$ ). Top (resp. bottom) figures correspond to database with mean sequence length equal to 20 (resp. 30). Boxplot colors correspond to different values of  $\sigma$  (10%, 15% and 20%).

The main conclusion from Figure 1 is that NEGSPAN is more efficient than eNSP when maxgap constraints are used. As expected, eNSP is more efficient than NEGSPAN without any maxgap constraint. This is mainly due to the number of extracted patterns. NEGSPAN extracts significantly more patterns than eNSP because of different choices for the semantics of NSPs. First, eNSP uses a stronger negation semantics. It can easily be proved that, without maxgap constraints, the set of patterns extracted by NEGSPAN is a superset of those extracted by eNSP<sup>6</sup>. Second, eNSP potentially misses a lot of interesting patterns due to the minimal support imposed on the positive partners, that is a strong additional constraint. Indeed, we set  $\zeta = 0.7\sigma$  which is a high value for specifying the set of positive partners on which negative patterns are explored. It is an advantageous setting for the efficiency of eNSP.

An interesting result is that, for reasonably long sequences (20 or 30), even a weak maxgap constraint ( $\tau = 10$ ) significantly reduces the number of patterns and makes NEGSPAN more efficient. This is of particular interest because the maxgap is a quite natural constraint when mining long sequences. It prevents from taking into account long distance correlations that are more likely irrelevant. Another interesting question raised by these results is the real meaning of extracted patterns by eNSP. In fact, under low frequency thresholds, it extracts numerous patterns that are

<sup>6</sup>Proof is given in an extended version of the paper available online.

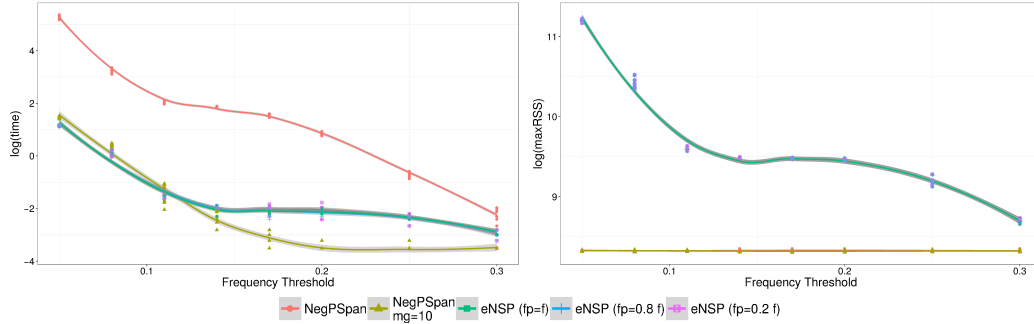


Figure 2: Comparison of computing time (left) and memory consumption (right) between eNSP and NEGPSpan wrt minimal support.

not frequent when weak maxgap constraints are considered. As a consequence, the significance of most of the patterns extracted by eNSP seems poor while processing “long” sequences datasets.

Figure 1 also illustrates classical results encountered with sequential pattern mining algorithms. We can note that, for both algorithms, the number of patterns and runtime increase exponentially as the minimum support decreases. Also, the number of patterns and the runtime increase notably with sequence length.

Figure 2 illustrates computation time and memory consumption with respect to minimum threshold for different settings: eNSP is ran with different values for the minimal frequency of the positive partner of negative patterns (100%, 80% and 20% of the minimal frequency threshold) and NEGPSpan is ran with a maxgap of 10 or without. Computation times show similar results as in previous experiments: NEGPSpan becomes as efficient as eNSP with a (weak) maxgap constraint. We can also notice that the minimal frequency of the positive partners does not impact eNSP computing times neither memory requirements.

The main result illustrated by this Figure is that NEGPSpan consumes significantly less memory than eNSP. This comes from the depth-first search strategy which prevents from memorizing many patterns. On the opposite, eNSP requires to keep in memory all frequent positive patterns and their occurrence list. The lower the threshold is, the more memory is required.

## 5.2 Experiments on real datasets

This section presents experiments on the real datasets from the SPMF repository<sup>7</sup>. These datasets consist of click-streams or texts represented as sequences of items. Datasets features and results are reported in Table 3. For every dataset, we have computed the negative sequential patterns with a maximum length of  $l = 5$  items and a minimal frequency threshold set to  $\sigma = 5\%$ . NEGPSpan is set with a maxgap  $\tau = 10$  and eNSP is set with  $\zeta = .7\sigma$ . For each dataset, we provide the computation time, the memory consumption and the numbers of positive and negative extracted patterns. Note that the numbers of positive patterns for eNSP are given for  $\zeta$  threshold, *i.e.* the support threshold for positive partners used to generate negative patterns.

For the *sign* dataset, the execution has been stopped after 10 mn to avoid running out of

<sup>7</sup><http://www.philippe-fournier-viger.com/spmf/index.php?link=datasets.php>

Table 3: Results on real datasets with setting  $\sigma = 5\%$ ,  $l = 5$ ,  $\tau = 10$ ,  $\zeta = .7\sigma$ . Bold faces highlight lowest computation times or memory consumptions.

	Dataset			NEGPSPAN				eNSP			
	$ \mathcal{D} $	$ \mathcal{I} $	length	time (s)	mem (kb)	#pos	#neg	time (s)	mem (kb)	#pos	#neg
<i>Sign</i>	730	267	51.99	<b>15.51</b>	<b>6,220</b>	348	1,357,278	349.84 (!)	13,901,600	1,190,642	1,257,177
<i>Leviathan</i>	5,834	9,025	33.81	<b>6.07</b>	<b>19,932</b>	110	39797	28.43	428,916	7,691	17,220
<i>Bible</i>	36,369	13,905	21.64	38.82	<b>68,944</b>	102	43,701	<b>27.38</b>	552,288	1,364	2,621
<i>BMS1</i>	59,601	497	2.51	<b>0.16</b>	<b>22,676</b>	5	0	0.18	34,272	8	7
<i>BMS2</i>	77,512	3,340	4.62	0.37	<b>39,704</b>	1	0	<b>0.35</b>	53,608	3	2
<i>kosarak25k</i>	25,000	14804	8.04	0.92	<b>24,424</b>	23	409	<b>0.53</b>	43,124	50	51
<i>MSNBC</i>	31,790	17	13.33	<b>40.97</b>	<b>41,560</b>	613	56,418	41.44	808,744	2,441	5,439

memory. The number of positive patterns extracted by eNSP considering the  $\sigma$  threshold is not equal to NEGPSPAN simply because of the maxgap constraint.

The results presented in Table 3 confirm the results from experiments on synthetic datasets. First, it highlights that NEGPSPAN requires significant less memory for mining every dataset. Second, NEGPSPAN outperforms eNSP for datasets the mean sequence length of which is long (*Sign*, *Leviathan*, and *MSNBC*). In case of the *Bible* dataset, the number of extracted patterns by eNSP is very low compared to NEGPSPAN due to the constraint on minimal frequency of positive partners.

### 5.3 Case study: care pathway analysis

This section presents the use of NSPs for analyzing epileptic patient care pathways. Recent studies suggest that medication changes may be associated with epileptic seizures for patients with long term treatment with anti-epileptic (AE) medication [10]. NSP mining algorithms are used to extract patterns of drugs deliveries that may inform about the suppression of a drug from a patient treatment. In [3], we studied discriminant temporal patterns but it does not explicitly extract the information about medication absence as a possible explanation of epileptic seizures.

Our dataset was obtained from the french insurance database [8] called SNIIRAM. 8,379 Epileptic patients were identified by their hospitalization related to an epileptic event. For each patient, the sequence of drugs deliveries within the 90 days before the epileptic event was obtained from the SNIIRAM. For each drug delivery, the event id is a tuple  $\langle m, grp, g \rangle$  where  $m$  is the ATC code of the active molecule,  $g \in \{0, 1\}$  and  $grp$  is the speciality group. The speciality group identifies the drug presentation (international non-proprietary name, strength per unit, number of units per pack and dosage form). The dataset contains 251,872 events over 7,180 different drugs. The mean length of a sequence is  $7.89 \pm 8.44$  itemsets. Length variance is high due to the heterogenous nature of care pathways. Some of them represent complex therapies involving the consumption of many different drugs while others are simple case consisting of few deliveries of anti-epileptic drugs.

Let first compare results obtained by eNSP and NEGPSPAN to illustrate the differences in the patterns sets extracted by each algorithm. To this end, we set up the algorithms with  $\sigma = 14.3\%$  (1,200 sequences), a maximum pattern length of  $l = 3$ ,  $\tau = 3$  for NEGPSPAN and  $\zeta = .1 \times \sigma$  the minimal support for positive partners for eNSP. eNSP extracts 1,120 patterns and NEGPSPAN only 10 patterns (including positive and negative patterns). Due to a very low  $\zeta$  threshold, many positive patterns are extracted by eNSP leading to generate a lot of singleton negative patterns (*i.e.* a pattern that hold a single negated item).

Table 4: Patterns involving *valproic acid* switches with their supports computed by eNSP and NEGPSPAN.

pattern	support eNSP	support NEGPSPAN
$\mathbf{p}_1 = \langle 383 \neg(86, 383) 383 \rangle$	1,579	
$\mathbf{p}_2 = \langle 383 \neg 86 383 \rangle$	1,251	1,243
$\mathbf{p}_3 = \langle 383 \neg 112 383 \rangle$	1,610	
$\mathbf{p}_4 = \langle 383 \neg 114 383 \rangle$	1,543	1,232
$\mathbf{p}_5 = \langle 383 \neg 115 383 \rangle$	1,568	1,236
$\mathbf{p}_6 = \langle 383 \neg 151 383 \rangle$	1,611	
$\mathbf{p}_7 = \langle 383 \neg 158 383 \rangle$	1,605	
$\mathbf{p}_8 = \langle 383 \neg 7 383 \rangle$		1,243

Precisely, we pay attention to the specific specialty of *valproic acid* which exists in generic form (event 383) or brand-named form (event 114) by selecting patterns that start and finish with event 383. The complete list of these patterns is given in Table 4. Other events correspond to other anti-epileptic drugs (7: *levetiracetam*, 158: *phenobarbital*) or psycholeptic drugs (112: *zolpidem*, 115: *clobazam*, 151: *zopiclone*) except 86 which is *paracetamol*.

First, it is interesting to note that with this setting, the two algorithms share only 3 patterns  $\mathbf{p}_2$ ,  $\mathbf{p}_4$  and  $\mathbf{p}_5$ , which have lower support with NEGPSPAN because of the maxgap constraint. This constraint also explains that pattern  $\mathbf{p}_3$  and  $\mathbf{p}_6$  are not extracted by NEGPSPAN. These patterns illustrate that in some cases, the patterns extracted by eNSP may not be really interesting because they involve distant events in the sequence. Pattern  $\mathbf{p}_1$  is not extracted by NEGPSPAN due to the strict-embedding pattern semantics. With eNSP semantics,  $\mathbf{p}_1$  means that there is no delivery of *paracetamol* and *valproic acid* at the same time. With NEGPSPAN semantics,  $\mathbf{p}_1$  means that there is no delivery of *paracetamol* neither *valproic acid* between two deliveries of *valproic acid*. The latter is stronger and the pattern support is lower. On the opposite, NEGPSPAN can extract patterns that are missed by eNSP. For instance, pattern  $\mathbf{p}_8$  is not extracted by eNSP because its positive partner,  $\langle 383, 7, 383 \rangle$ , is not frequent. In this case, it leads eNSP to miss a potentially interesting pattern involving two anti-epileptic drugs.

Now, we look at patterns involving a switch from generic form to brand-named form of *valproic acid* with the following settings  $\sigma = 1.2\%$ ,  $l = 3$  and  $\tau = 5$ . Mining only positive patterns extracts the frequent patterns  $\langle 114, 383, 114 \rangle$  and  $\langle 114, 114 \rangle$ . It is impossible to conclude about the possible impact of a switch from 114 to 383 as a possible event triggering an epileptic crisis. From negative patterns extracted by NEGPSPAN, we can observe that the absence of switch  $\langle 114 \neg 383 114 \rangle$  is also frequent in this dataset. Contrary to eNSP semantics which does bring a new information (that can be deduced from frequent patterns), this pattern concerns embeddings corresponding to real interesting cases thanks to gap constraints.

## 6 Conclusion and perspectives

This article have investigated negative sequential pattern mining (NSP). It highlights that state of the art algorithms do not extract the same patterns, not only depending on their syntax and algorithms specificities, but also depending on the semantical choices. This article have proposed definitions that clarify the negation semantics encountered in the literature. We have showed



that NSP support depends on the semantics of itemset non-inclusion, two possible alternatives for considering negation of itemsets and two manners for considering multiple embeddings in a sequences. This let us point out the limits of the state of the art algorithm eNSP that imposes a minimum support for positive partner and that is not able to deal with embedding constraints, and more especially maxgap constraints.

We have proposed NEGPSpan a new algorithm for mining negative sequential patterns that overcomes these limitations. Our experiments show that NEGPSpan is more efficient than eNSP on dataset with medium long sequences (more than 20 itemsets) even when weak maxgap constraints are applied and that it prevents from missing possibly interesting patterns.

In addition, NEGPSpan is based on theoretical foundations that enables to extend it to the extraction of closed or maximal patterns to reduce the number of extracted patterns even more.

## Acknowledgments

The authors would like to thank REPERES Team from Rennes University Hospital for spending time to discuss our case study results. This work is supported by the ANSM/PEPS project.

## References

- [1] Longbing Cao, Xiangjun Dong, and Zhigang Zheng. e-NSP: Efficient negative sequential pattern mining. *Artificial Intelligence*, 235:156–182, 2016.
- [2] Longbing Cao, Philip S. Yu, and Vipin Kumar. Nonoccurring behavior analytics: A new area. *Intelligent Systems*, 30(6):4–11, 2015.
- [3] Yann Dauxais, Thomas Guyet, David Gross-Amblard, and André Happe. Discriminant chronicles mining - application to care pathways analytics. In *Proceedings of 16th Conference on Artificial Intelligence in Medicine*, volume 10259 of *Lecture Notes in Computer Science*, pages 234–244. Springer, 2017.
- [4] Yongshun Gong, Tiantian Xu, Xiangjun Dong, and Guohua Lv. e-nspf: Efficient mining negative sequential pattern from both frequent and infrequent positive sequential patterns. *International Journal of Pattern Recognition and Artificial Intelligence*, 31(02):1750002, 2017.
- [5] Sue-Chen Hsueh, Ming-Yen Lin, and Chien-Liang Chen. Mining negative sequential patterns for e-commerce recommendations. In *Proceedings of Asia-Pacific Services Computing Conference*, pages 1213–1218. IEEE, 2008.
- [6] Sujatha Kamepalli, Raja Sekhara, and Rao Kurra. Frequent Negative Sequential Patterns – a Survey. *International Journal of Computer Engineering and Technology*, 5, 3:115–121, 2014.
- [7] Carl H. Mooney and John F. Roddick. Sequential pattern mining – approaches and algorithms. *ACM Computing Survey*, 45(2):1–39, 2013.
- [8] G. Moulis, M. Lapeyre-Mestre, A. Palmaro, G. Pugno, J.-L. Montastruc, and L. Sailler. French health insurance databases: What interest for medical research? *La Revue de Médecine Interne*, 36:411–417, 2015.

- [9] 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.
- [10] Elisabeth Polard, Emmanuel Nowak, André Happe, Arnaud Biraben, and Emmanuel Oger. Brand name to generic substitution of antiepileptic drugs does not lead to seizure-related hospitalization: a population-based case-crossover study. *Pharmacoepidemiology and drug safety*, 24:1161–1169, 2015.
- [11] Tiantian Xu, Xiangjun Dong, Jianliang Xu, and Xue Dong. Mining high utility sequential patterns with negative item values. *International Journal of Pattern Recognition and Artificial Intelligence*, 31(10):1750035, 2017.
- [12] Tiantian Xu, Xiangjun Dong, Jianliang Xu, and Yongshun Gong. E-msnsp: Efficient negative sequential patterns mining based on multiple minimum supports. *International Journal of Pattern Recognition and Artificial Intelligence*, 31(02):1750003, 2017.
- [13] Zhigang Zheng, Yanchang Zhao, Ziyue Zuo, and Longbing Cao. Negative-GSP: An efficient method for mining negative sequential patterns. In *Proceedings of the Australasian Data Mining Conference*, pages 63–67, 2009.

## A Proofs and additional propositions

Proof of Proposition 1.

*Proof.* Let  $\mathbf{s} = \langle s_1, \dots, s_n \rangle$  be a sequence and  $\mathbf{p} = \langle p_1, \dots, p_m \rangle$  be a negative sequential pattern. Let  $\mathbf{e} = (e_i)_{i \in [m]} \in [n]^m$  be a soft-embedding of pattern  $\mathbf{p}$  in sequence  $\mathbf{s}$ . Then, the definition matches the one for strict-embedding if  $p_i$  is positive. If  $p_i$  is negative then  $\forall j \in [e_{i-1} + 1, e_{i+1} - 1]$ ,  $p_i \not\sqsubseteq s_j$ , i.e.  $\forall j \in [e_{i-1} + 1, e_{i+1} - 1]$ ,  $\forall \alpha \in p_i$ ,  $\alpha \notin s_j$  and then  $\forall \alpha \in p_i$ ,  $\forall j \in [e_{i-1} + 1, e_{i+1} - 1]$ ,  $\alpha \notin s_j$ . It thus implies that  $\forall \alpha \in p_i$ ,  $\alpha \notin \bigcup_{j \in [e_{i-1} + 1, e_{i+1} - 1]} s_j$ , i.e. by definition,  $p_i \not\sqsubseteq \bigcup_{j \in [e_{i-1} + 1, e_{i+1} - 1]} s_j$ .

The exact same reasoning is done in reverse way to prove the equivalence.  $\square$

**Proposition 3.** Soft-embedding  $\implies$  strict-embedding for patterns consisting of items.

*Proof.* Let  $\mathbf{s} = \langle s_1, \dots, s_n \rangle$  be a sequence and  $\mathbf{p} = \langle p_1, \dots, p_m \rangle$  be a NSP s.t. each  $\forall i$ ,  $|p_i| = 1$  and  $\mathbf{p}$  occurs in  $\mathbf{s}$  according to the soft-embedding semantic.

There exists  $\epsilon = (e_i)_{i \in [m]} \in [n]^m$  s.t. for all  $i \in [n]$ ,  $p_i$  is positive implies  $p_i \in s_{e_i}$  and  $p_i$  is negative implies that for all  $j \in [e_{i-1} + 1, e_{e+1} - 1]$ ,  $p_i \notin s_j$  (items only) then  $p_i \notin \bigcup_{j \in [e_{i-1} + 1, e_{e+1} - 1]} s_j$  i.e.  $p_i \not\sqsubseteq \bigcup_{j \in [e_{i-1} + 1, e_{e+1} - 1]} s_j$  (no matter  $\not\sqsubseteq$  or  $\not\sqsubset$ ). As a consequence  $\epsilon$  is a strict-embedding of  $\mathbf{p}$ .  $\square$

**Proposition 4.** Let  $\mathcal{D}$  be a dataset of sequences of items and  $\mathbf{p} = \langle p_1, \dots, p_m \rangle$  be a sequential pattern extracted by eNSP, then without embedding constraints  $\mathbf{p}$  is extracted by NEGPSpan with the same minimum support.

*Proof.* If  $\mathbf{p}$  is extracted by eNSP, it implies that its positive partner is frequent in the dataset  $\mathcal{D}$ . As a consequence, each  $p_i$ ,  $i \in [m]$  is a singleton itemset.

According to the search space of NEGPSpan defined by  $\triangleleft^8$  is  $\mathbf{p}$  is frequent then it will be reached by the depth-first search. Then it is sufficient to prove that for any sequence  $\mathbf{s} = \langle s_1, \dots, s_n \rangle \in \mathcal{D}$  such that  $\mathbf{p}$  occurs in  $\mathbf{s}$  according to eNSP semantic (strict-embedding, strong absence), then  $\mathbf{p}$  also occurs in  $\mathbf{s}$  according to the NEGPSpan semantics (soft-embedding, weak absence). With that and considering the same minimum support threshold,  $\mathbf{p}$  is frequent according to NEGPSpan. Proposition 3 gives this result.  $\square$

**Example 3** (Pattern set comparisons). NEGPSpan extracts more patterns than eNSP on sequences of items. In fact, NEGPSpan can extract patterns with negative itemsets larger than 2.

eNSP extract patterns that are not extracted by NEGPSpan on sequences of itemsets. Practically, NEGPSpan uses a size limit for negative itemsets  $\nu \geq 1$ . eNSP extracts patterns whose positive partners are frequent. The positive partner, extracted by PrefixSpan may hold itemsets larger than  $\nu$ , and if the pattern with negated itemset is also frequent, then this pattern will be extract by eNSP, but not by NEGPSpan.

**Proposition 5** (Anti-monotonicity of NSP). The support of NSP is anti-monotonic with respect to  $\triangleleft$  (resp.  $\blacktriangleleft$ ) when  $\not\sqsubseteq \stackrel{\text{def}}{=} \not\sqsubseteq$  (resp.  $\not\sqsubset \stackrel{\text{def}}{=} \not\sqsubset$ ) is considered.

*Proof.* Proof of anti-monotonicity of the support wrt  $\triangleleft$ , considering  $\not\sqsubseteq \stackrel{\text{def}}{=} \not\sqsubseteq$ .

Let  $\mathbf{p} = \langle p_1, \dots, p_n \rangle$  and  $\mathbf{q} = \langle q_1, \dots, q_n \rangle$  be two NSP. Let  $|\mathbf{p}| = \sum_i |p_i|$  denote the total number of items in pattern  $\mathbf{p}$ , and we introduce  $\triangleleft_1$  a order relation between immediate ‘‘extensions’’:

---

<sup>8</sup> $\triangleleft$  or  $\blacktriangleleft$  ... no matter with sequences of items for which  $\not\sqsubseteq$  and  $\not\sqsubseteq$  are equivalent.

$\mathbf{p} \triangleleft_1 \mathbf{q} \Leftrightarrow \mathbf{p} \triangleleft \mathbf{q} \wedge |\mathbf{p}| = |\mathbf{q}| + 1$ . We can easily prove recursively that  $\mathbf{p} \triangleleft \mathbf{q} \Leftrightarrow \exists \{\rho_i\}_{i \in [r]}$ , a sequence of patterns<sup>9</sup>  $\rho_0 = p$ ,  $\rho_r = q$ ,  $\forall j, \rho_j \triangleleft_1 \rho_{j+1}$ . Then, it is sufficient to prove that the support is anti-monotonic for relation  $\triangleleft_1$ .

By definition of  $\triangleleft$  and  $\triangleleft_1$ , we have that:  $\mathbf{p} \triangleleft_1 \mathbf{q}$  iff  $m = n$ ,  $\forall i \in [n]$ ,  $p_i$  negative  $\Leftrightarrow q_i$  negative,  $\exists! k \in [n]$ ,  $q_k = p_k \cup \{e\}$  where  $e \in \mathcal{I}$  and  $\forall i \neq k$ ,  $q_i = p_i$ . Note that this definition considers that if  $q_k$  is a singleton, then  $p_k$  is an empty itemset (with same sign as  $q_k$ ) that has been inserted to have the same lengths.

Let  $\mathbf{p}$  and  $\mathbf{q}$  s.t.  $\mathbf{p} \triangleleft_1 \mathbf{q}$ , then we have to show that for any sequence  $\mathbf{s} \in \mathcal{D}$ ,  $\mathbf{q}$  occurs in  $\mathbf{s}$  implies that  $\mathbf{p}$  occurs in  $\mathbf{s}$ .

Let assume that  $\mathbf{q}$  occurs in  $\mathbf{s}$  considering  $\not\subseteq \stackrel{\text{def}}{=} \not\subseteq$  (and thus no matter the embedding strategy, according to proposition 1, we use strict-embedding in the following). Then, for all embedding  $\epsilon = (e_i)_i$ ,  $q_i \subseteq s_{e_i}$  if  $q_i$  is positive and  $\forall j \in [e_{i-1} + 1, e_{i+1} - 1]$ ,  $q_i \not\subseteq s_j$  if  $q_i$  is negative. Then for all  $i \neq k$ , we have immediately that  $p_i \subseteq s_{e_i}$  if  $p_i$  is positive and  $\forall j \in [e_{i-1} + 1, e_{i+1} - 1]$ ,  $p_i \not\subseteq s_j$  if  $p_i$  is negative. The remaining case is for the  $k$ -th itemset of  $\mathbf{q}$ . If  $q_k$  is positive, then  $p_k \subseteq q_k \subseteq s_{e_k}$  and thus  $\epsilon$  is an embedding of  $\mathbf{p}$  in  $\mathbf{s}$ . If  $q_k$  is negative,  $\forall j \in [e_{k-1} + 1, e_{k+1} - 1]$ ,  $q_k \not\subseteq s_j$  i.e.  $\forall j \in [e_{k-1} + 1, e_{k+1} - 1]$ ,  $\forall e \in q_k$ ,  $e \notin s_j$ . Then  $\forall e \in q_k$ ,  $\forall j \in [e_{k-1} + 1, e_{k+1} - 1]$ ,  $e \notin s_j$  and thus  $\forall e \in p_k$ ,  $\forall j \in [e_{k-1} + 1, e_{k+1} - 1]$ ,  $e \notin s_j$  because  $p_k \subset q_k$ <sup>10</sup>, i.e. in short  $\forall j \in [e_{k-1} + 1, e_{k+1} - 1]$ ,  $p_k \not\subseteq s_j$ .

Note any embedding of  $\mathbf{q}$  yields an embedding for  $\mathbf{p}$ , then the property holds for weak and strong absence semantics.

Proof of anti-monotonicity of the support wrt  $\blacktriangleleft$ , considering  $\not\subseteq \stackrel{\text{def}}{=} \not\subseteq$ . Similarly to  $\triangleleft_1$ , we define  $\blacktriangleleft_1$  such that  $\mathbf{p} \blacktriangleleft_1 \mathbf{q}$  iff  $|\mathbf{p}| = |\mathbf{q}|$ <sup>11</sup> there exists  $k \in [n]$  s.t.  $\forall i \neq k$ ,  $p_i = q_i$  and if  $q_k$  is negative,  $q_k \subset p_k$  otherwise  $p_k \subset q_k$ . The exact same reasoning can be done except for the last case of the negative  $k$ -th itemset for  $\mathbf{q}$ . We have here to distinguish strict and soft embeddings. If  $q_k$  is negative and considering strict-embeddings,  $\forall j \in [e_{k-1} + 1, e_{k+1} - 1]$ ,  $q_k \not\subseteq s_j$  i.e.  $\forall j \in [e_{k-1} + 1, e_{k+1} - 1]$ ,  $\exists e \in q_k$ ,  $e \notin s_j$  when using definition of  $\not\subseteq$  (see Definition 3). Considering definition of  $\blacktriangleleft_1$ , we have that  $q_k \subset p_k$  s.t.  $e$  is also an element of  $p_k$  and then  $\forall j \in [e_{k-1} + 1, e_{k+1} - 1]$ ,  $p_k \not\subseteq s_j$ .

If  $q_k$  is negative and considering soft-embeddings,  $q_k \not\subseteq \bigcup_{j \in [e_{k-1} + 1, e_{k+1} - 1]} s_j$ , i.e.  $\exists e \in q_k$ ,  $e \notin \bigcup_{j \in [e_{k-1} + 1, e_{k+1} - 1]} s_j$ . Again, we have  $q_k \subset p_k$  and then  $e$  is also an element of  $p_k$  s.t.  $p_k \not\subseteq \bigcup_{j \in [e_{k-1} + 1, e_{k+1} - 1]} s_j$ .  $\square$

---

<sup>9</sup>This sequence is not unique.

<sup>10</sup>Here, we also have that  $\emptyset \subset q_k$  in case of singleton itemset  $q_k$ . Thus justify the use of the simplified version of

$\triangleleft_1$ .

<sup>11</sup>again, we allow empty sets.