

# Deriving Generic Bounds for Time-Series Constraints Based on Regular Expressions Characteristics

Ekaterina Arafaïlova · Nicolas Beldiceanu ·  
Helmut Simonis

**Abstract** We introduce the concept of regular expression characteristics as a unified way to concisely express bounds on time-series constraints. This allows us not only to define time-series constraints in a compositional way, but also to deal with their combinatorial aspect in a compositional way, without developing ad-hoc bounds for each time-series constraint separately.

## 1 Introduction

A *time series* is here a sequence of integers, corresponding to measurements taken over a time interval. Time series are common in many application areas, such as the output of electric power stations over multiple days [9], or the manpower required in a call-centre [5], or the daily capacity of a hospital clinic over a period of years. Time series are constrained by physical or organisational limits, which restrict the evolution of the series.

We showed in [6] that many constraints  $\gamma(\langle X_1, X_2, \dots, X_n \rangle, N)$  on an unknown time series  $X = \langle X_1, X_2, \dots, X_n \rangle$  can be specified in a *compositional way* by a triple  $\langle \sigma, f, g \rangle$ , where  $\sigma$  is a regular expression over the alphabet  $\Sigma = \{ '<', '=', '>' \}$  (we assume the reader is familiar with regular expressions [13]), while  $f \in \{ \text{max}, \text{min}, \text{one}, \text{surf}, \text{width} \}$  is called a *feature function*, and  $g \in \{ \text{Max}, \text{Min}, \text{Sum} \}$  is called an *aggregator function*. Volume II of the global constraint catalogue [4] contains 266 such functional time-series constraints.

This is an extended version of parts of the CP 2016 paper [3] which involves a subset of the original authors. This paper introduces 6 new characteristics of regular expressions to express generic bounds on time-series constraints, which were not discussed in the original paper [3]. Ekaterina Arafaïlova is supported by the EU H2020 programme under grant 640954 for project GRACeFUL. Nicolas Beldiceanu is partially supported by the GRACeFUL project and by the Gaspard Monge Program for Optimization and Operations Research (PGMO). Helmut Simonis is supported by Science Foundation Ireland (SFI) under grant SFI/10/IN.1/I3032; the Insight Centre for Data Analytics is supported by SFI under grant SFI/12/RC/2289.

E. Arafaïlova · N. Beldiceanu  
TASC (LS2N) IMT Atlantique, FR – 44307 Nantes, France  
E-mail: {Ekaterina.Arafaïlova, Nicolas.Beldiceanu}@imt-atlantique.fr

H. Simonis  
Insight Centre for Data Analytics, University College Cork, Cork, Ireland  
E-mail: Helmut.Simonis@insight-centre.org

It is currently unknown in general, how to maintain efficiently domain consistency for such time-series constraints. Computing bounds on the result variable  $N$  of a time-series constraint is a way to handle the combinatorial aspect and thus improve propagation. Since we have too many time-series constraints deriving such bounds needs to be done in a systematic way. Motivated by this, we sketched in [3] a methodology to obtain such bounds and illustrated it only for time-series constraints when  $g = \text{Max}$  and  $f = \text{min}$ .

The contribution of this paper, which makes explicit the approach sketched in [3], is to introduce the notion of *regular expression characteristics* that provides a unified way to concisely express bounds on the result variable  $N$  of a time-series constraint. Six regular expression characteristics are introduced, which allows coming up in a compositional way with bounds when  $\langle g, f \rangle \in \{\langle \text{Sum}, \text{one} \rangle, \langle \text{Max}, \text{width} \rangle, \langle \text{Min}, \text{width} \rangle, \langle \text{Sum}, \text{width} \rangle\}$ : five main theorems (see Theorems 1, 2, 3, 4, and 5) allow obtaining 95 bounds implemented in Volume II of the global constraint catalogue [4]. When the time-series variables  $\langle X_1, \dots, X_n \rangle$  have the same interval integer domain, these bounds are sharp for all the 22 regular expression of [4]. We now put in perspective with existing work the contribution of this paper.

Going back to the work of Schützenberger [15], *regular cost functions* are quantitative extensions of regular languages that correspond to a function mapping a word to an integer value or infinity (QRE for short). Recently there was a rise of interest in this area, both from a theoretical perspective [10] with max-plus automata, and from a practical point of view with the synthesis of cost register automata [1] for data streams [2]. Within constraint programming automata constraints were introduced in [14] and in [8, 12], the later also computing an integer value from a word. More recently, the work on synthesising automata from transducers [6] for identifying all maximal occurrences of a pattern in a time series is part of the QRE line of research. While most previous mentioned works give quantitative extensions of regular languages as their general motivation, to the best of our knowledge none of them introduced the concept of regular expression characteristics, which is the key abstraction proposed here. The paper is structured in the following way:

- ◊ In Section 2, we recall the background both on regular expressions [13], and on the way of describing time-series constraints in a compositional way [6].
- ◊ In Section 3, we first introduce a notation system for denoting regular expression characteristics. Then we present six regular expressions characteristics, which characterise different quantitative aspects of regular expressions useful for capturing some of their combinatorial flavour. Finally, based on two of these characteristics, we provide a necessary condition for the occurrence of a regular expression in a time-series.
- ◊ In Section 4, we show how to obtain generic bounds for time-series constraints whose result variable is constrained to be the number of occurrences of a regular expression in a time-series, i.e., time-series constraints where  $g = \text{Sum}$  and  $f = \text{one}$ .
- ◊ In Section 5, we show how to obtain generic bounds for the result variables of time-series constraints for which the feature  $f$  is `width`, and the aggregator  $g$  is in  $\{\text{Max}, \text{Min}, \text{Sum}\}$ .
- ◊ In Section 6, we synthesise all the results on bounds we have so far from the CP paper [3], and from Sections 4 and 5 of this paper: for each bound we recall (1) the regular expression characteristics it uses, (2) the generic theorem it comes from, and (3) the properties under which the bound is sharp.
- ◊ We evaluate in Section 7 the beneficial propagation impact of the derived bounds.

## 2 Background

First we give the necessary background on word and regular languages. Then we recall the time-series constraints introduced in [6].

## 2.1 Regular Languages

An *alphabet*  $\mathcal{A}$  is a finite set of symbols and a symbol of  $\mathcal{A}$  is called a *letter*. A *word* on  $\mathcal{A}$  is a finite sequence of symbols belonging to  $\mathcal{A}$ . The empty word is denoted by  $\varepsilon$ . The *length* of a word  $w$  is the number of letters in  $w$  and is denoted by  $|w|$ . For  $i \in [1, |w|]$ ,  $w[i]$  denotes the letter  $i$  of a word  $w$ . The concatenation of two words is denoted by putting them side by side, with an implicit infix operator between them. A word  $w$  is a *factor* of a word  $x$  if there exist two words  $v$  and  $z$  such that  $x = vwz$ ; when  $v = \varepsilon$ ,  $w$  is a *prefix* of  $x$ , when  $z = \varepsilon$ ,  $w$  is a *suffix* of  $x$ . If both  $w$  is not empty and different from  $x$ , then it is a *proper factor* of  $x$ . Given a word  $w$  and a positive integer  $k > 0$ ,  $w^k$  denotes the *concatenation of  $k$  occurrences of  $w$* . Given an integer  $k$  and a language  $\mathcal{L}$ ,  $\mathcal{L}^k$  is defined by  $\mathcal{L}^0 = \{\varepsilon\}$ ,  $\mathcal{L}^1 = \mathcal{L}$  and  $\mathcal{L}^k = \mathcal{L} \cdot \mathcal{L}^{k-1}$  where ‘ $\cdot$ ’ is the concatenation operator. Then the Kleene closure of  $\mathcal{L}$  is defined by  $\cup_{n \geq 0} \mathcal{L}^n$  and denoted by  $\mathcal{L}^*$ .

**Definition 1** A regular expression [11]  $r$  on an alphabet  $\mathcal{A}$  and the language  $\mathcal{L}_r$  it describes, the regular language, are recursively defined as follows:

- (1) 0 and 1 are regular expressions that respectively describe  $\emptyset$  (the empty set) and  $\{\varepsilon\}$ .
- (2) For every letter  $\ell$  of  $\mathcal{A}$ ,  $\ell$  is a regular expression that describes the singleton  $\{\ell\}$ .
- (3) If  $r_1$  and  $r_2$  are regular expressions, respectively describing the regular languages  $\mathcal{L}_{r_1}$  and  $\mathcal{L}_{r_2}$ , then  $r_1 + r_2$ ,  $r_1 * r_2$  and  $r_1^*$  are regular expressions that respectively describe the regular languages  $\mathcal{L}_{r_1} \cup \mathcal{L}_{r_2}$ ,  $\mathcal{L}_{r_1} \cap \mathcal{L}_{r_2}$ , and  $\mathcal{L}_{r_1}^*$ .

*Example 1* Consider the alphabet  $\Sigma = \{<, =, >\}$ .

- **Decreasing** = ‘ $>$ ’ is a regular expression on  $\Sigma$ . The word  $v = >$  is a word of length 1 on  $\Sigma$  that belongs to  $\mathcal{L}_{\text{Decreasing}}$ , and it does not have any proper factors. The word ‘ $>>$ ’ is a word of length 2 on  $\Sigma$ , which does not belong to  $\mathcal{L}_{\text{Decreasing}}$ .
- **Inflexion** = ‘ $<(<|=)*> | >(>|=)*<$ ’ is a regular expression on  $\Sigma$ . The word  $v = >=<$  is a word of length 3 on  $\Sigma$  that belongs to  $\mathcal{L}_{\text{Inflexion}}$ . The word  $v$  has multiple proper factors, e.g., ‘ $>$ ’, ‘ $<$ ’. The word ‘ $>=<<$ ’ does not belong to  $\mathcal{L}_{\text{Inflexion}}$ .  $\triangle$

**Definition 2** A regular expression  $r$  is a *non-fixed length regular expression* if not all words of  $\mathcal{L}_r$  have the same length.

*Example 2* We give two examples of regular expressions, a first one with a fixed length and a second one with a non-fixed length.

- The **Decreasing** = ‘ $>$ ’ regular expression has a fixed length since  $\mathcal{L}_{\text{Decreasing}}$  contains a single word.
- The **Inflexion** = ‘ $<(<|=)*> | >(>|=)*<$ ’ regular expression does not have a fixed length since  $\mathcal{L}_{\text{Inflexion}}$  contains words of different length.  $\triangle$

**Definition 3** A regular expression over an alphabet  $\mathcal{A}$  is *disjunction-capsuled* if it is in the form of ‘ $r_1 r_2 \dots r_p$ ’, where every  $r_i$  (with  $i \in [1, p]$ ) is, either a letter of the alphabet  $\mathcal{A}$ , or a regular expression whose regular language contains the empty word.

Note that Definition 3 is a slight extension of a similar notion introduced in [17].

*Example 3* Table 1 recalls the 22 regular expressions used for describing time-series constraints in [4, 6]. Every regular expression  $\sigma$  in column 2 of Table 1 is in the form of  $\sigma = \sigma_1 | \sigma_2 | \dots | \sigma_t$  with  $t \geq 1$ , and every  $\sigma_i$  (with  $i \in [1, t]$ ) is a disjunction-capsuled regular expression. Then  $\mathcal{L}_\sigma$  is the union of the  $\mathcal{L}_{\sigma_i}$  (with  $i \in [1, t]$ ).

The ‘ $> | > (> |=)* > (< | < (< |=)* <$ ’ regular expression has the same regular language as **Gorge**, but is not disjunction-capsuled.  $\triangle$

## 2.2 Time-Series Constraints

A time series here is a sequence of integers corresponding to measurements taken over the time. We showed in [6] that many constraints  $\gamma(\langle X_1, X_2, \dots, X_n \rangle, N)$  on an unknown time series  $\langle X_1, X_2, \dots, X_n \rangle$  of given length  $n$ , where every  $X_i$  is an *integer variable*, can be specified by a triple  $\langle \sigma, f, g \rangle$ , where  $\sigma$  is a regular expression on the alphabet  $\Sigma = \{ '<', '=', '>' \}$  that is characterised by two integer constants  $a_\sigma$  and  $b_\sigma$ , whose role is to trim the left and right borders of the regular expression, while  $f \in \{\text{max, min, one, surf, width}\}$  is called a *feature*, and  $g \in \{\text{Max, Min, Sum}\}$  is called an *aggregator*. Let the sequence  $S = \langle S_1, S_2, \dots, S_{n-1} \rangle$ , called the *signature* and containing *signature variables*, be linked to the sequence  $X$  via the *signature constraints*  $(X_i < X_{i+1} \Leftrightarrow S_i = '<') \wedge (X_i = X_{i+1} \Leftrightarrow S_i = '=') \wedge (X_i > X_{i+1} \Leftrightarrow S_i = '>')$  for all  $i \in [1, n-1]$ . If a sub-signature  $\langle S_i, S_{i+1}, \dots, S_j \rangle$  is a maximal word matching  $\sigma$  in the signature of  $X$ , then the subseries  $\langle X_{i+b_\sigma}, X_{i+b_\sigma+1}, \dots, X_{j+1-a_\sigma} \rangle$  is called a  $\sigma$ -*pattern* and the subseries  $\langle X_i, X_{i+1}, \dots, X_{j+1} \rangle$  is called an *extended  $\sigma$ -pattern*. Integer variable  $N$  is constrained to be the aggregation, computed using  $g$ , of the list of values of feature  $f$  for all  $\sigma$ -patterns in  $\langle X_1, X_2, \dots, X_n \rangle$ . We name a time-series constraint specified by  $\langle \sigma, f, g \rangle$  as  $g\_f\_ \sigma$ . A time series is *maximal* for  $g\_f\_ \sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  if it yields the maximum value of  $N$  among all time series of length  $n$  that have the same initial domains for the time-series variables.

name $\sigma$	regular expression	$a_\sigma$	$b_\sigma$
BumpOnDecreasingSequence	'>><>>'	1	2
Decreasing	'>'	0	0
DecreasingSequence	'(>(> =)*)*>'	0	0
DecreasingTerrace	'>=+>'	1	1
DipOnIncreasingSequence	'<<><<'	1	2
Gorge	'(>(> =)*)*><((< =)*<)*'	1	1
Increasing	'<'	0	0
IncreasingSequence	'(<(< =)*)*<'	0	0
IncreasingTerrace	'<=+<'	1	1
Inflexion	'(<(< =)*>   >(> =)*<'	1	1
Peak	'(<(< =)* (> =)*>'	1	1
Plain	'>=*<'	1	1
Plateau	'<=*>'	1	1
ProperPlain	'>=+<'	1	1
ProperPlateau	'<=+>'	1	1
Steady	'= '	0	0
SteadySequence	'=+ '	0	0
StrictlyDecreasingSequence	'>+ '	0	0
StrictlyIncreasingSequence	'<+ '	0	0
Summit	'(<(< =)*)*<>((> =)*>)*'	1	1
Valley	'>(> =)* (< =)*<'	1	1
Zigzag	'(<>)+ < (>   $\epsilon$ )   (><)+ > (<   $\epsilon$ )'	1	1

Table 1: Regular expression names  $\sigma$ , corresponding regular expressions, values of the parameters  $a_\sigma$  and  $b_\sigma$

We consider the set of 22 regular expressions used in [4], which is given in Table 1. Most of these regular expressions capture topological patterns that one wants to control when generating time-series, while some of them, like **Zigzag**, correspond to abnormal situations one wants to catch from existing time-series. Within a  $\sigma$ -*pattern* the two integer constants  $b_\sigma$  and  $a_\sigma$  trim respectively the left and right borders of the interval  $[i, j+1]$  to the leftmost and rightmost variable of  $\langle X_1, X_2, \dots, X_n \rangle$  used to compute the corresponding feature: for example for **IncreasingTerrace** = ' $<=+<$ ', since  $b_{\text{IncreasingTerrace}} = a_{\text{IncreasingTerrace}} = 1$ , we only consider the  $X_i$  that are involved in an equality, i.e., the  $X_i$  of the flat part of the terrace.

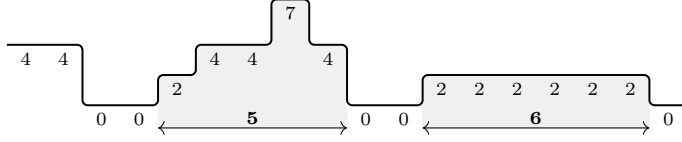


Fig. 1: Illustrating the  $\text{MIN\_WIDTH\_PEAK}(5, \langle 4, 4, 0, 0, 2, 4, 4, 7, 4, 0, 0, 2, 2, 2, 2, 2, 0 \rangle)$  time-series constraint with its two peaks of respective width 5 and 6

*Example 4* Consider the time series  $X = \langle 4, 4, 0, 0, 2, 4, 4, 7, 4, 0, 0, 2, 2, 2, 2, 2, 0 \rangle$  with its signature  $\text{'=>=<=<=<=<=>=<=====>'}$  and the regular expression  $\text{Peak} = \text{'<(<|=)*(>|=)*>'}$  with  $b_{\text{Peak}}$  and  $a_{\text{Peak}}$  being both equal to 1: a **Peak**-pattern, called a *peak*, within a time series corresponds, except for its first and last elements, to a maximal occurrence of **Peak** in the signature, and the **width** feature value of a peak is its number of elements. The time series  $X$  contains two peaks, namely  $\langle 2, 4, 4, 7, 4 \rangle$  and  $\langle 2, 2, 2, 2, 2, 2 \rangle$ , visible the way  $X$  is plotted in Figure 1, of widths 5 and 6 respectively, hence the minimal-width peak, obtained by using the aggregator **Min**, has width  $N = 5$ : the corresponding constraint is named  $\text{MIN\_WIDTH\_PEAK}$ .  $\triangle$

### 3 Regular Expressions Characteristics

To get parametrised bounds, this section introduces regular expressions characteristics used for deriving sharp lower and upper bounds on the result variable of a time series constraint when the feature is in  $\{\text{one}, \text{width}\}$ . For all characteristics we use a notation system, which is described in Section 3.1. We introduce the following characteristics:

- The *width* of a regular expression in Section 3.2.
- The *height* of a regular expression in Section 3.3.
- The *range* of a regular expression wrt a time series length in Section 3.4.
- The *set of inducing words* of a regular expression in Section 3.5.
- The *overlap* of a regular expression wrt an integer interval domain in Section 3.6.
- The *smallest variation of maxima* of a regular expression wrt an integer interval domain in Section 3.7.

Section 3.8 presents a summary example combining all the introduced regular expressions characteristics for the **DecreasingTerrace** regular expression. Section 3.9 introduces a necessary and sufficient condition for the existence of at least one occurrence of a regular expression within the signature of a time series under some hypothesis on the domain of time-series variables. Table 3 provides for each of the 22 regular expressions in Table 1 the corresponding value of each regular expression characteristics.

#### 3.1 A Notation System for Regular Expression Characteristics

We introduce a notation system for naming the characteristics of regular expressions. A regular expression characteristic  $C$  is a function, denoted by  $C_R^P(V)$ , whose arguments are  $R$ ,  $P$ , and  $V$  explained below:

- $R$  is a regular expression over  $\Sigma = \{<, =, >\}$ .
- $P$  is a subset, possible empty, drawn from  $\{\ell, u, n\}$ , where  $[\ell, u]$  is the domain of the variables of a time series, and  $n$  is the length of a given time series.

- $V$  is a vector of additional arguments, which are different from  $R$ ,  $\ell$ ,  $u$ , and  $n$ . If  $V$  is empty, then we simply write  $C_R^P$ . Quite often these additional arguments correspond to words in  $\mathcal{L}_R$  since a characteristic  $C_R^P$  will be defined in terms of an other characteristic  $C_R^P(V)$ : for instance the height of a regular expression  $R$  will be defined in terms of the heights of words in  $\mathcal{L}_R$ .

The domain of the function  $C_R^P(V)$  is the cartesian product of the following elements in the given order:

- The domain of  $R$ , namely  $\mathcal{R}_\Sigma$ , which is the set of regular expressions over  $\Sigma$ .
- The cartesian product of the domains of the elements of  $P$ , if  $P$  is not empty.
- The cartesian product of the domains of the arguments of  $V$ , if  $V$  is not empty.

The font used for the symbol ‘ $C$ ’ depends on the type of values returned by  $C_R^P(V)$ :

- If the codomain of  $C_R^P(V)$  is  $\mathbb{Z}$ , then ‘ $C$ ’ is a lower-case Greek letter, e.g.,  $\delta$ .
- If the codomain of  $C_R^P(V)$  is the power set of some set, then ‘ $C$ ’ is an upper-case Greek letter, e.g.,  $\Delta$ .

Some characteristics are associated with, either the lower or the upper bound on the value of the result variable of a time-series constraint. In this case, the ones associated with the upper (respectively lower) bound are denoted by  $\overline{C}_R^P(V)$  (respectively  $\underline{C}_R^P(V)$ ).

### 3.2 Width

This section introduces the *width* characteristic of a regular expression; it will be used in Theorem 2 for computing the sharp upper bound on the number of occurrences of the regular expression within the signature of a time series. This characteristic is also used for defining a necessary and sufficient condition, see Property 1, for the existence of at least one occurrence of a regular expression within the signature of a time series over an integer interval domain.

**Definition 4** Consider a regular expression  $\sigma$ . The *width* of  $\sigma$ , denoted by  $\omega_\sigma$ , is a function that maps an element of  $\mathcal{R}_\Sigma$  to  $\mathbb{N}$ . It is defined by  $\omega_\sigma = \min_{v \in \mathcal{L}_\sigma} |v|$ .

*Example 5* Consider the  $\sigma = \text{DecreasingTerrace}$  regular expression. There is a single shortest word in  $\mathcal{L}_\sigma$ , namely ‘ $>=>$ ’. Thus the width of  $\sigma$  is 3. Hence, any extended  $\sigma$ -pattern has at least  $3 + 1$  time-series variables.  $\triangle$

### 3.3 Height

We introduce the notion of height of a regular expression, which is used for defining a necessary and sufficient condition, see Property 1, for the existence of at least one occurrence of a regular expression within the signature of a time series. This characteristic is also used in Theorem 2 of Section 4 for computing a sharp upper bound on the number of occurrences of the regular expression within the signature of a time series. Definitions 5 and 6 are only used for introducing Definition 7.

**Definition 5** Consider a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$ . The *set of supporting time series* of a word  $v$  in  $\mathcal{L}_\sigma$  wrt  $\langle \ell, u \rangle$ , denoted by  $\Omega_\sigma^{\langle \ell, u \rangle}(v)$ , is a function that maps an element of  $\mathcal{R}_\Sigma \times \mathbb{Z} \times \mathbb{Z} \times \Sigma^*$  to  $\mathcal{P}(\mathbb{Z}^*)$ , where  $\mathcal{P}(\mathbb{Z}^*)$  is the power set of  $\mathbb{Z}^*$ . Each element of  $\Omega_\sigma^{\langle \ell, u \rangle}(v)$  is a time series over  $[\ell, u]$  whose signature is  $v$ , and is called a *supporting time series* of  $v$  wrt  $\langle \ell, u \rangle$ .

**Definition 6** Consider a regular expression  $\sigma$ . The *height* of a word  $v$  in  $\mathcal{L}_\sigma$ , denoted by  $\eta_\sigma(v)$ , is a function that maps an element of  $\mathcal{R}_\Sigma \times \Sigma^*$  to  $\mathbb{N}$ . It is defined by  $\eta_\sigma(v) = \min_{\Omega_\sigma^{(\ell, u)}(v) \neq \emptyset} (u - \ell)$ , where  $[\ell, u]$  is an integer interval domain.

**Definition 7** Consider a regular expression  $\sigma$ . The *height* of  $\sigma$ , denoted by  $\eta_\sigma$ , is a function that maps an element of  $\mathcal{R}_\Sigma$  to  $\mathbb{N}$ . It is defined by  $\eta_\sigma = \min_{v \in \mathcal{L}_\sigma} \eta_\sigma(v)$ .

*Example 6* Consider the  $\sigma = \text{DecreasingTerrace}$  regular expression and an integer interval domain  $[\ell, u]$ . When  $u - \ell \leq 1$ , there does not exist a time series over  $[\ell, u]$  whose signature is a word in  $\mathcal{L}_\sigma$ ; but when  $u - \ell = 2$ , there exists a time series over  $[\ell, u]$  whose signature is a word, for example ' $>=>$ ', in  $\mathcal{L}_\sigma$ . Hence, the height of  $\sigma$  equals 2.  $\triangle$

### 3.4 Range

This section introduces a characteristics needed by Theorems 3, 4, and 5 for characterising sharp bounds on the result value of a time-series constraint when the feature is **width**. This characteristics, described in Definition 8, is called the *range* of a regular expression  $\sigma$ , and shows the variation of the minimum height of words of  $\mathcal{L}_\sigma$  for words of increasing length.

**Definition 8** Consider a regular expression  $\sigma$  and a time series length  $n$ . The *range* of  $\sigma$  wrt  $\langle n \rangle$ , denoted by  $\phi_\sigma^{(n)}$ , is a function that maps an element of  $\mathcal{R}_\Sigma \times \mathbb{N}$  to  $\mathbb{N}$ . It is defined by  $\phi_\sigma^{(n)} = \min_{v \in \mathcal{L}_\sigma, |v|=n-1} \eta_\sigma(v)$ , where  $\eta_\sigma(v)$  is the height of the word  $v$ . If  $\mathcal{L}_\sigma$  does not contain any word of length  $n - 1$ , then the value of  $\phi_\sigma^{(n)}$  is undefined.

*Example 7* Consider the  $\sigma = \text{SteadySequence}$  regular expression. For every integer  $n > \omega_\sigma$ , the language  $\mathcal{L}_\sigma$  contains a word with  $n - 1$  equalities. Any word of this type has a height of 0. Hence, the range of  $\sigma$  is a constant function of  $n$ , which equals 0.

### 3.5 Set of Inducing Words

Given a disjunction-capsuled regular expression  $\sigma$ , we first introduce the notion of *inducing word* of  $\mathcal{L}_\sigma$ , which is a maximum sequence of letters that appears in every word of  $\mathcal{L}_\sigma$  in a fixed order. Then we generalise this notion to any disjunction of disjunction-capsuled regular expression.

**Definition 9** Consider a disjunction-capsuled regular expression  $\sigma$ . The (unique) non-empty shortest word of  $\mathcal{L}_\sigma$  is the *inducing word* of  $\mathcal{L}_\sigma$ .

**Definition 10** Consider a regular expression  $\sigma$  that is in the form of  $\sigma = \sigma_1 | \sigma_2 | \dots | \sigma_t$  with  $t \geq 1$ , where every  $\sigma_i$  (with  $i \in [1, t]$ ) is a disjunction-capsuled regular expression. The *set of inducing words* of  $\sigma$ , denoted by  $\Theta_\sigma$ , is a function that maps an element of  $\mathcal{R}_\Sigma$  to  $\mathcal{P}(\Sigma^*)$ , where  $\mathcal{P}(\Sigma^*)$  is the power set of  $\Sigma^*$ . The value of  $\Theta_\sigma$  is the union of inducing words of every  $\sigma_i$ .

*Example 8* Consider the  $\text{Inflexion} = '<(<|=)*> | >(>|=)*<'$  regular expression. It can be represented as  $\text{Inflexion}_1 | \text{Inflexion}_2$ , where  $\text{Inflexion}_1 = '<(<|=)*>'$ ,  $\text{Inflexion}_2 = '>(>|=)*<'$ , and both  $\text{Inflexion}_1$  and  $\text{Inflexion}_2$  are disjunction-capsuled. The word  $v = '<>'$  is the inducing word of  $\mathcal{L}_{\text{Inflexion}_1}$ , the word  $v = '><'$  is the inducing word of  $\mathcal{L}_{\text{Inflexion}_2}$ , and both  $v$  and  $v$  are inducing words of  $\mathcal{L}_{\text{Inflexion}}$ . Hence,  $\Theta_{\text{Inflexion}} = \{'<>', '><'\}$ .  $\triangle$

### 3.6 Overlap

This section introduces the *overlap* characteristic of a regular expression; it will be used in Theorem 2 for computing the sharp upper bound on the number of occurrences of the regular expression within the signature of a time series. To define the overlap of a regular expression  $\sigma$  wrt to an integer interval domain  $[\ell, u]$ , Definition 11 first introduces the notion of *set of superpositions* of two words  $v$  and  $w$  in  $\mathcal{L}_\sigma$  wrt  $\langle \ell, u \rangle$ . Intuitively the superposition of  $v$  and  $w$  wrt  $\langle \ell, u \rangle$  is the signature  $z$  of some ground time series over  $[\ell, u]$  that contains exactly two  $\sigma$ -patterns, i.e.,  $v$  as a prefix and  $w$  as a suffix of  $z$ , and whose length does not exceed the length of  $vw$ .

**Definition 11** Consider a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$ . The *set of superpositions* of two words,  $v$  and  $w$  in  $\mathcal{L}_\sigma$ , wrt  $\langle \ell, u \rangle$ , denoted by  $\Gamma_\sigma^{\langle \ell, u \rangle}(v, w)$ , is a function that maps an element of  $\mathcal{R}_\Sigma \times \mathbb{Z} \times \mathbb{Z} \times \Sigma^* \times \Sigma^*$  to  $\mathcal{P}(\Sigma^*)$ , where  $\mathcal{P}(\Sigma^*)$  is the power set of  $\Sigma^*$ . Each element  $z$  in  $\Gamma_\sigma^{\langle \ell, u \rangle}(v, w)$  is a word over  $\Sigma$ , called a superposition of  $v$  and  $w$  wrt  $\langle \ell, u \rangle$  and satisfying all the following conditions:

- (1)  $z \notin \mathcal{L}_\sigma$ , (2)  $\Omega_\sigma^{\langle \ell, u \rangle}(z) \neq \emptyset$ , (3)  $v$  is a prefix of  $z$ , (4)  $w$  is a suffix of  $z$ , (5)  $|z| \leq |vw|$ .

*Example 9* Consider  $\sigma = \text{DecreasingTerrace}$ , and an integer interval domain  $[\ell, u]$  allowing to have at least one occurrence of  $\sigma$  in the signature of a time series over  $[\ell, u]$ , i.e.,  $u - \ell \geq 2$ . We compute a superposition of the pair  $\langle v, v \rangle$ , where  $v = '>=>>' \in \mathcal{L}_\sigma$ . Let  $z$  denote  $'>=>=>'$ .

- \* First, assume that  $u - \ell = 2$ . The word  $z$  is not a superposition of  $v$  and  $v$ , since the number of consecutive increases in the word  $z$  is 3, which is strictly greater than  $u - \ell = 2$ , and thus  $\Omega_\sigma^{\langle \ell, u \rangle}(z) = \emptyset$ . Indeed, when  $u - \ell = 2$ , there is no superposition of  $v$  and  $v$ , because any word different from  $z$  satisfying the first four conditions of Definition 11 will violate Condition (5) of Definition 11, i.e., will be strictly longer than  $2 \cdot |v|$ , thus  $\Gamma_\sigma^{\langle \ell, u \rangle}(v, v) = \emptyset$ .
- \* Now assume that  $u - \ell = 3$ . Then,  $\Omega_\sigma^{\langle \ell, u \rangle}(z) \neq \emptyset$ , and the word  $z$  is the only superposition of  $v$  and  $v$ , thus  $\Gamma_\sigma^{\langle \ell, u \rangle}(v, v) = \{'>=>=>'\}$ .
- \* Finally, assume that  $u - \ell \geq 4$ . The sets of supporting time series of both words  $'>=>=>'$  and  $'>=>=>=>'$  wrt  $\langle \ell, u \rangle$  are not empty, and these two words are the only superpositions of  $v$  and  $v$  wrt  $\langle \ell, u \rangle$ , thus  $\Gamma_\sigma^{\langle \ell, u \rangle}(v, v) = \{'>=>=>', '>=>=>=>'\}$ .  $\triangle$

For a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$ , we now introduce the *overlap* characteristic of  $\sigma$  wrt  $\langle \ell, u \rangle$ , which is a crucial component in the sharp upper bound formula stated in Theorem 2. It corresponds to the maximum number of time-series variables that can be shared by two consecutive extended  $\sigma$ -patterns: when maximising the number of  $\sigma$ -patterns in a time series, we need to enforce as many consecutive extended  $\sigma$ -patterns as possible to have as many common time-series variables as possible.

**Definition 12** Consider a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$ . The *overlap* of two words,  $v$  and  $w$  in  $\mathcal{L}_\sigma$ , wrt  $\langle \ell, u \rangle$ , denoted by  $o_\sigma^{\langle \ell, u \rangle}(v, w)$ , is a function that maps an element of  $\mathcal{R}_\Sigma \times \mathbb{Z} \times \mathbb{Z} \times \Sigma^* \times \Sigma^*$  to  $\mathbb{N}$ . It is defined by

$$o_\sigma^{\langle \ell, u \rangle}(v, w) = \begin{cases} \left( |vw| - \min_{z \in \Gamma_\sigma^{\langle \ell, u \rangle}(v, w)} |z| \right) + 1 & \text{if } \Gamma_\sigma^{\langle \ell, u \rangle}(v, w) \neq \emptyset \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

Case (1) of Definition 12 corresponding to condition  $\Gamma_\sigma^{\langle \ell, u \rangle}(v, w) \neq \emptyset$  states that the overlap is strictly greater than 0 iff there exists at least one ground time series over  $[\ell, u]$  that is not strictly longer than  $|vw|$  and that has exactly two  $\sigma$ -patterns corresponding to the occurrences of  $v$  and  $w$  in its signature. The



term  $|vw| - \min_{z \in \Gamma_\sigma^{\langle \ell, u \rangle}(v, w)} |z|$  denotes the maximum possible overlap that is actually achieved by the shortest superposition. The increment  $+1$  denotes that we count the number of time-series variables rather than the number of signature variables.

We now generalise in Definition 13 the notion of overlap wrt  $\langle \ell, u \rangle$  upon a regular expression.

**Definition 13** Consider a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$ . The *overlap* of  $\sigma$  wrt  $\langle \ell, u \rangle$ , denoted by  $o_\sigma^{\langle \ell, u \rangle}$ , is a function that maps an element of  $\mathcal{R}_\Sigma \times \mathbb{Z} \times \mathbb{Z}$  to  $\mathbb{N}$ . If there exists a constant  $c$  in  $\mathbb{N}$  such that for any pair of words  $v, w$  in  $\mathcal{L}_\sigma$ , the value of  $o_\sigma^{\langle \ell, u \rangle}(v, w)$  is bounded by  $c$ , then the overlap of  $\sigma$  wrt  $\langle \ell, u \rangle$  is defined by  $o_\sigma^{\langle \ell, u \rangle} = \max_{v, w \in \mathcal{L}_\sigma} o_\sigma^{\langle \ell, u \rangle}(v, w)$ . Otherwise,  $o_\sigma^{\langle \ell, u \rangle}$  is undefined.

By Definition 13, we need to compute the overlap of  $\sigma$  wrt every pair of values  $\langle \ell, u \rangle$ , i.e., every domain  $[\ell, u]$ . Note that it is enough to compute the overlap of  $\sigma$  wrt  $\langle \ell, u \rangle$  once for every value of the difference  $u - \ell$  permitting an occurrence of  $\sigma$  in the signature of a time series, i.e., for a difference that is greater than or equal to the height of the regular expression  $\sigma$ . While in the general case, we do need to consider every value of  $u - \ell$ , this is not required for all the 22 regular expressions in Table 1, where we only need to consider at most two different values of  $u - \ell$ .

*Example 10* We successively consider values of the overlap of two regular expressions.

- Consider the  $\sigma = \text{DecreasingTerrace}$  regular expression, whose height  $\eta_\sigma$  is 2.
  - \* If  $u - \ell = \eta_\sigma = 2$ , then  $o_\sigma^{\langle \ell, u \rangle} = 0$ , because as shown in Example 9, for any pair of words in  $\mathcal{L}_\sigma$ , the set of their superpositions wrt  $\langle \ell, u \rangle$  is empty.
  - \* If  $u - \ell \geq \eta_\sigma + 1 = 3$ , then  $o_\sigma^{\langle \ell, u \rangle} = 2$  and is obtained, for example, for the pair ' $>=>$ ' and ' $>=>$ ', and their superposition ' $>=>=>$ '.
  - \* For any other value of  $u - \ell \geq 4$ , the value of the overlap of  $\sigma$  wrt  $\langle \ell, u \rangle$  equals 2 as well.
- Consider the  $\sigma = '<=*|=*>'$  regular expression and an integer interval domain  $[\ell, u]$  such that  $u > \ell$ . The overlap of  $\sigma$  wrt  $\langle \ell, u \rangle$  is undefined, because for any constant  $c$  in  $\mathbb{N}$ , there always exists a pair of words of length  $c + 1$  whose overlap wrt  $\langle \ell, u \rangle$  equals  $c + 1$ .  $\triangle$

### 3.7 Smallest Variation of Maxima

This section introduces the *smallest variation of maxima* characteristics of a regular expression, which is used in Theorem 2 for computing the sharp upper bound on the number of occurrences of the regular expression within the signature of a time series. To maximise the number of occurrences of a regular expression  $\sigma$  in a time series over an integer interval domain  $[\ell, u]$ , we select extended  $\sigma$ -patterns of minimum length  $\omega_\sigma + 1$  such that two consecutive extended  $\sigma$ -patterns maximise the number of shared time-series variables, i.e., share  $o_\sigma^{\langle \ell, u \rangle}$  variables. Unfortunately, for a few regular expressions like **DecreasingTerrace**, it is not always possible that all  $\sigma$ -patterns share  $o_\sigma^{\langle \ell, u \rangle}$  time-series variables: since we decrease by at least one unit between two consecutive overlapping extended  $\sigma$ -patterns we will be blocked at some point by the lower limit  $\ell$ , even if we start from the upper limit  $u$ . To maximise the number of  $\sigma$ -patterns in a time series, we must decrease as little as possible on two consecutive overlapping extended  $\sigma$ -patterns. Definition 16 formalises the notion of *smallest variation of the maxima* of a regular expression wrt  $\langle \ell, u \rangle$ . First, Definition 14 defines the notion of shift of a proper factor in a word in the language of a regular expression wrt some integer interval domain. Then, using this notion, Definition 15 (respectively Definition 16) introduces the smallest variation of the maxima of two words (respectively a language  $\mathcal{L}_\sigma$ ) wrt  $\langle \ell, u \rangle$ .

**Definition 14** Consider a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$ . The *shift* of a proper factor  $w$  in a word  $v$  in  $\mathcal{L}_\sigma$  wrt  $\langle \ell, u \rangle$ , denoted by  $\bar{\nu}_\sigma^{\langle \ell, u \rangle}(v, w, i)$ , is a function that maps an element of  $\mathcal{R}_\Sigma \times \mathbb{Z} \times \mathbb{Z} \times \Sigma^* \times \Sigma^* \times \mathbb{N}$  to  $\mathbb{N}$ . It is defined by

$$\bar{\nu}_\sigma^{\langle \ell, u \rangle}(v, w, i) = \min_{t \in \Omega_\sigma^{\langle \ell, u \rangle}(v)} \min_{x \in t_{w_i}} (\max(t) - x),$$

where  $\max(t)$  is the maximum value of a time series  $t$ , a supporting time series of  $v$  wrt  $\langle \ell, u \rangle$ , and  $t_{w_i}$  is a subseries of  $t$  corresponding to the  $i^{\text{th}}$  extended  $\sigma$ -pattern whose signature is  $w$ . If  $w$  is not a proper factor of  $v$ , or  $i$  is strictly greater than the number of occurrences of  $w$  in  $v$ , then  $\bar{\nu}_\sigma^{\langle \ell, u \rangle}(v, w, i)$  is undefined.

Consider a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$ . For any  $v$  in  $\mathcal{L}_\sigma$ , if  $u - \ell \geq \eta_\sigma(v)$ , then the value of  $\bar{\nu}_\sigma^{\langle \ell, u \rangle}(v, w, i)$  does not depend on the domain  $[\ell, u]$ , because there always exists a time series in  $\Omega_\sigma^{\langle \ell, u \rangle}(v)$  where each variable has its largest value compared to the other time series of  $\Omega_\sigma^{\langle \ell, u \rangle}(v)$ . Then,  $\bar{\nu}_\sigma^{\langle \ell, u \rangle}(v, w, i)$  does not depend on the values in the domain, but only on the structure of the word  $v$ . Hence, w.l.o.g. the notation for  $\bar{\nu}_\sigma^{\langle \ell, u \rangle}(v, w, i)$  can be simplified to  $\bar{\nu}_\sigma(v, w, i)$ .

*Example 11* Consider  $\sigma = \text{DecreasingTerrace}$  when  $u - \ell \geq 3$ , and two words  $v = \langle = < = < \rangle$  and  $w = \langle = < \rangle$ . The word  $v$  contains two occurrences of  $w$ , thus the value of  $\bar{\nu}_\sigma(v, w, i)$  is defined when  $i \in \{1, 2\}$ :

- \* When  $i$  is 1, the value of  $\bar{\nu}_\sigma(v, w, 1)$  equals 0, since the first extended  $\sigma$ -pattern whose signature is  $w$  necessarily contains the maximum of any time series in  $\Omega_\sigma^{\langle \ell, u \rangle}(v)$ .
- \* When  $i$  is 2, the value of  $\bar{\nu}_\sigma(v, w, 2)$  equals 1, since the maximum of the second extended  $\sigma$ -pattern whose signature is  $w$  has a difference of at least one with the maximum of any time series in  $\Omega_\sigma^{\langle \ell, u \rangle}(v)$ .  $\triangle$

**Definition 15** Consider a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$ . The *smallest variation of maxima* of superpositions of two words  $w$  and  $v$  in  $\mathcal{L}_\sigma$  wrt  $\langle \ell, u \rangle$ , denoted by  $\delta_\sigma^{\langle \ell, u \rangle}(v, w)$ , is a function that maps an element of  $\mathcal{R}_\Sigma \times \mathbb{Z} \times \mathbb{Z} \times \Sigma^* \times \Sigma^*$  to  $\mathbb{N}$ . It is defined by

$$\delta_\sigma^{\langle \ell, u \rangle}(v, w) = \begin{cases} \bar{\nu}_\sigma(z_*, v, 1) - \bar{\nu}_\sigma(z_*, w, 1), & \text{if } v \neq w \text{ and } \Gamma_\sigma^{\langle \ell, u \rangle}(v, w) \neq \emptyset \\ \bar{\nu}_\sigma(z_{**}, v, 1) - \bar{\nu}_\sigma(z_{**}, w, 2), & \text{if } v = w \text{ and } \Gamma_\sigma^{\langle \ell, u \rangle}(v, w) \neq \emptyset \\ 0, & \text{if } \Gamma_\sigma^{\langle \ell, u \rangle}(v, w) = \emptyset \end{cases}$$

where the words  $z_*$  and  $z_{**}$  both belongs to  $\Gamma_\sigma^{\langle \ell, u \rangle}(v, w)$ , and the value  $\min_{z \in \Gamma_\sigma^{\langle \ell, u \rangle}(v, w)} |\bar{\nu}_\sigma(z, v, 1) - \bar{\nu}_\sigma(z, w, 1)|$  (respectively  $\min_{z \in \Gamma_\sigma^{\langle \ell, u \rangle}(v, w)} |\bar{\nu}_\sigma(z, v, 1) - \bar{\nu}_\sigma(z, w, 2)|$ ) is reached when  $z$  is  $z_*$  (respectively  $z_{**}$ ).

In Definition 15, either  $\bar{\nu}_\sigma(z_*, v, 1)$  (respectively  $\bar{\nu}_\sigma(z_{**}, v, 1)$ ) or  $\bar{\nu}_\sigma(z_*, w, 1)$  (respectively  $\bar{\nu}_\sigma(z_{**}, w, 2)$ ) equals zero, since for any time series  $t$  whose signature is  $z_*$  (respectively  $z_{**}$ ), at least one of the two extended  $\sigma$ -patterns contains the maximum of  $t$ .

The next lemma introduces a property of words whose smallest variation of maxima wrt some integer interval domain is not zero.

**Lemma 1** Consider a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$ . If  $\delta_\sigma^{\langle \ell, u \rangle}(v, w)$ , the smallest variation of maxima of two words  $v$  and  $w$  in  $\mathcal{L}_\sigma$  wrt  $\langle \ell, u \rangle$ , is positive (respectively negative), then  $v$  (respectively  $w$ ) does not contain any ' $>$ ' (respectively ' $<$ ').

*Proof* For brevity, we consider only the case of  $\delta_\sigma^{\langle \ell, u \rangle}(v, w)$  being positive, the case of a negative value of  $\delta_\sigma^{\langle \ell, u \rangle}(v, w)$  being symmetric, and w.l.o.g. we assume that  $v \neq w$ .

Since  $\delta_\sigma^{\langle \ell, u \rangle}(v, w) > 0$ , there exists at least one superposition  $z$  of  $v$  and  $w$  wrt  $\langle \ell, u \rangle$  such that  $\bar{v}_\sigma(z, v, 1) = \delta_\sigma^{\langle \ell, u \rangle}(v, w)$ , and  $\bar{v}_\sigma(z, w, 1) = 0$ . Assume that  $v$  contains at least one ' $>$ '. Let  $i$  denote the position of the first ' $>$ ' in  $z$ , which is necessarily within its proper factor  $v$ . Since there exists a time series in  $\Omega_\sigma^{\langle \ell, u \rangle}(z)$  such that the time-series variable at position  $i$  equals  $u$ ,  $\bar{v}_\sigma(z, v, 1)$  equals 0. This contradicts the fact that  $\bar{v}_\sigma(z, v, 1) = \delta_\sigma^{\langle \ell, u \rangle}(v, w) > 0$ , thus the word  $v$  does not contain any ' $>$ '.  $\square$

**Definition 16** Consider a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$ . The *smallest variation of maxima* of  $\sigma$  wrt  $\langle \ell, u \rangle$ , denoted by  $\delta_\sigma^{\langle \ell, u \rangle}$ , is a function that maps an element of  $\mathcal{R}_\Sigma \times \mathbb{Z} \times \mathbb{Z}$  to  $\mathbb{N}$ . It is defined by

$$\delta_\sigma^{\langle \ell, u \rangle} = \begin{cases} \text{undefined,} & \text{if } \exists v_1, v_2, w_1, w_2 \in \mathcal{L}_\sigma \text{ s.t. } \delta_\sigma^{\langle \ell, u \rangle}(v_1, w_1) > 0 \text{ and } \delta_\sigma^{\langle \ell, u \rangle}(v_2, w_2) < 0 \\ 0, & \text{if } \delta_\sigma^{\langle \ell, u \rangle} = 0 \\ \delta_\sigma^{\langle \ell, u \rangle}(v_*, w_*), & \text{otherwise} \end{cases}$$

where the words  $v_*$  and  $w_*$  both belong to  $\mathcal{L}_\sigma$  and the value  $\min_{\substack{v, w \in \mathcal{L}_\sigma \\ \delta_\sigma^{\langle \ell, u \rangle}(v, w) \neq 0}} |\delta_\sigma^{\langle \ell, u \rangle}(v, w)|$  is reached when  $v$

is  $v_*$  and  $w$  is  $w_*$ .

*Example 12* Consider the  $\sigma = \text{DecreasingTerrace}$  regular expression, an integer interval domain  $[\ell, u]$  such that  $u - \ell \geq 3$ , and the superposition  $z = '>=>=>'$  of the words  $v = '>=>'$  and  $v = '>=>'$  in  $\mathcal{L}_\sigma$ . The value of  $\bar{v}_\sigma(z, v, 1) - \bar{v}_\sigma(z, v, 2)$  is equal to  $0 - 1 = -1$ . For any other pair of words of  $\mathcal{L}_\sigma$  whose set of superpositions wrt  $\langle \ell, u \rangle$  is not empty, we obtain a same or a smaller negative value. Hence, if two extended  $\sigma$ -patterns share some time-series variables, then the maximum of a second extended  $\sigma$ -pattern is at least one unit below, i.e.,  $\delta_\sigma^{\langle \ell, u \rangle} = -1$ , from the maximum of the first extended  $\sigma$ -pattern.  $\triangle$

If  $\delta_\sigma^{\langle \ell, u \rangle}$  is positive (respectively negative), then for any two extended  $\sigma$ -patterns that have at least one common time-series variable, the maximum of the first extended  $\sigma$ -pattern is strictly less (respectively greater) than the maximum of the second extended  $\sigma$ -pattern, e.g., for **DecreasingTerrace**,  $\delta_\sigma^{\langle \ell, u \rangle}$  equals  $-1$ , but for **IncreasingTerrace**,  $\delta_\sigma^{\langle \ell, u \rangle}$  equals  $+1$ .

### 3.8 Summary Example Illustrating All Characteristics

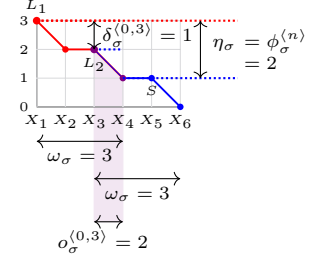
This section illustrates the various regular expression characteristics introduced in the previous sections.

*Example 13* Consider the  $\sigma = \text{DecreasingTerrace}$  regular expression and a time series  $X$  of length 6 over an integer interval domain  $[0, 3]$ . Let us recall the characteristics mentioned in Examples 6, 10, 5, and 12, which are illustrated in Figure 2.

- The *width* of  $\sigma$ , denoted by  $\omega_\sigma$ , equals 3.
- The *height* of  $\sigma$ , denoted by  $\eta_\sigma$ , equals 2. This is the difference between the  $y$ -coordinates of the points  $L_1$  and  $S$  in Figure 2, which are respectively the maximum and the minimum points of the first extended  $\sigma$ -pattern of  $X$ .
- The *range* of  $\sigma$  wrt  $\langle n \rangle$ , denoted by  $\phi_\sigma^{\langle n \rangle}$ , equals 2, with  $n \in \mathbb{N}$  being greater than or equal to  $\omega_\sigma = 3$ .
- The *overlap* of  $\sigma$  wrt  $\langle 0, 3 \rangle$ , denoted by  $o_\sigma^{\langle 0, 3 \rangle}$ , equals 2. It is the number of common points of the first and the second extended  $\sigma$ -patterns in Figure 2, i.e., the number points coloured in violet.

- The *smallest variation of maxima* of  $\sigma$  wrt  $\langle 0, 3 \rangle$ , denoted by  $\delta_\sigma^{(0,3)}$ , equals 1. It is the difference between the  $y$ -coordinates of the  $L_1$  and the  $L_2$  points in Figure 2, which are the maxima points of the first, respectively the second, extended  $\sigma$ -pattern of  $X$ .  $\triangle$

Fig. 2: A time series of length  $n = 6$  over the integer interval domain  $[0, 3]$  containing two extended  $\sigma$ -patterns, where  $\sigma$  is **DecreasingTerrace**. The  $x$ -axis is for time-series variables, the  $y$ -axis is for domain values. The first (respectively second) extended  $\sigma$ -pattern is shown in red (respectively blue). The common time-series variables of the two extended  $\sigma$ -patterns are coloured in violet.  $L_1$  (respectively  $L_2$ ) is the point whose  $y$ -coordinate is maximum among all points of the first (respectively second) extended  $\sigma$ -pattern.  $S$  is the point whose  $y$ -coordinate is minimum among all points of the first extended  $\sigma$ -pattern.



### 3.9 Necessary and Sufficient Condition for the Existence of an Occurrence of a Regular Expression

Consider a regular expression  $\sigma$  and a time series  $X = \langle X_1, X_2, \dots, X_n \rangle$  with every  $X_i$  ranging over the same integer interval domain. There exists a necessary and sufficient condition, based on the domains and the number of time-series variables, for  $\sigma$  to occur at least once within the signature of  $X$ . In order to define this condition we use the *width* of a regular expression, introduced in Definition 4, and the *height* of a regular expression, introduced in Definition 7.

*Property 1* Consider a regular expression  $\sigma$  and a time series  $\langle X_1, X_2, \dots, X_n \rangle$  with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ . The *necessary-and-sufficient condition* is satisfied if the two following conditions hold:

- (i) The value of  $n$  is strictly greater than  $\omega_\sigma$ , the width of  $\sigma$ .
- (ii) The difference between  $u$  and  $\ell$  is greater than or equal to  $\eta_\sigma$ , the height of  $\sigma$ .

*Example 14* Consider the  $\sigma = \text{DecreasingTerrace}$  regular expression and a time series of length  $n$  over an integer interval domain  $[\ell, u]$ . We recall the values computed in Examples 6 and 5, namely the height of  $\sigma$  is 2, and the width of  $\sigma$  is 3. Hence, the necessary-and-sufficient condition is satisfied if  $n > 3$  and  $u - \ell \geq 2$ .  $\triangle$

All formulae presented in all the next sections assume that Property 1 holds.

## 4 Constraints that Restrict the Number of Occurrences of a Regular Expression

The first family of time-series constraints we consider is the  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  family. Given a sequence  $X = \langle X_1, X_2, \dots, X_n \rangle$ , where all  $X_i$  are integer variables, it enforces the number of occurrences of pattern  $\sigma$  in  $X$  to be equal to  $N$ . Within this constraint family the aggregator is **Sum**, and the feature is **one**. The main results of Section 4 are described by Theorems 1 and 2, which respectively provide a sharp lower bound and a sharp upper bound on the number of occurrences of a regular expression  $\sigma$  in the signature of a time series provided all  $X_i$  (with  $i \in [1, n]$ ) have the same integer interval domain  $[\ell, u]$ . Section 4 is structured in the following way:

- ◊ First, Section 4.1 introduces Property 2, and gives a sharp lower bound on  $N$  provided Property 2 holds.
- ◊ Second, Section 4.2 provides an upper bound, not necessarily sharp, on  $N$ . This bound is valid for any regular expression  $\sigma$  for which the overlap characteristics is defined and does not exceed the width of  $\sigma$ .

- ◇ Third, Section 4.3 extends the upper bound on  $N$  of Section 4.2, and shows that the extended formula is sharp under some hypothesis on the regular expression characteristics:
  - \* Section 4.3.1 defines Properties 3 and 4 of regular expressions that must hold to obtain a sharp upper bound.
  - \* Section 4.3.2 describes the structure of a maximal time series for  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  provided either Property 3 or Property 4 holds.
  - \* Based on the structure of a maximal time series for  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$ , Section 4.3.3 provides a sharp upper bound on  $N$ , provided either Property 3 or Property 4 holds.
- ◇ Finally, Section 4.4 gives a sharp upper bound on  $N$  in a special case of  $\sigma$  being **SteadySequence**, where neither Property 3 nor Property 4 is satisfied.

#### 4.1 A Sharp Lower Bound on the Number of Pattern Occurrences

Consider a  $\text{NB\_}\sigma(X, N)$  time-series constraint with  $X = \langle X_1, X_2, \dots, X_n \rangle$ , where every  $X_i$  (with  $i \in [1, n]$ ) is over the same integer interval domain  $[\ell, u]$ . This section focusses on providing the lower bound on  $N$ . For almost every regular expression of Table 1, we can assign the variables of  $X$  to values in  $[\ell, u]$  in a way that the signature of  $X$  does not contain any occurrence of the regular expression  $\sigma$ . The only two exceptions are the **Steady** = '=' and the **SteadySequence** = '=+' regular expressions when  $\ell = u$ . The next theorem, namely Theorem 1, provides a sharp lower bound on  $N$  assuming the property that we now introduce holds.

*Property 2* A regular expression  $\sigma$  has the NB-simple property for an integer interval domain  $[\ell, u]$  if  $\sigma$  is a disjunction of disjunction-capsuled regular expressions and if at least one of the following conditions holds:

- (i) Every inducing word of  $\sigma$  includes at least one letter that is different from '='.
- (ii) Every inducing word of  $\sigma$  includes at least one '=', and  $u > \ell$ .

**Theorem 1** Consider a  $\text{NB\_}\sigma(X, N)$  time-series constraint with  $X = \langle X_1, X_2, \dots, X_n \rangle$ , where every  $X_i$  (with  $i \in [1, n]$ ) is over the same integer interval domain  $[\ell, u]$ , and, where  $\sigma$  is a disjunction of disjunction-capsuled regular expressions. If  $\sigma$  has the NB-simple property for  $[\ell, u]$ , then a sharp lower bound on  $N$  is 0.

*Proof* If Condition (i) of Property 2 is satisfied, then by definition of an inducing word, every word of  $\mathcal{L}_\sigma$  contains at least one letter that is not '='. Hence, the time series  $X$ , where all variables are assigned to the same value, has no occurrences of  $\sigma$  in its signature, and thus a sharp lower bound on  $N$  is 0.

If Condition (ii) of Property 2 is satisfied, then every word in  $\mathcal{L}_\sigma$  contains at least one '='. The ground time series of length  $n$  with alternating  $\ell$  and  $\ell + 1$  has no equalities in its signature, and thus no occurrences of  $\sigma$ . Hence, a sharp lower bound on  $N$  equals 0.  $\square$

Every regular expression in Table 1 has the NB-simple property for any integer interval domain  $[\ell, u]$ , except **Steady** and **SteadySequence** for the domain  $[\ell, u]$  such that  $\ell = u$ . We now consider the cases of **Steady** and **SteadySequence** where neither condition of Property 2 holds, which means that Theorem 1 cannot be applied for computing a sharp lower bound on  $N$ .

**Proposition 1** Consider a  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with  $\sigma$  being the **Steady** regular expression, and with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$  such that  $\ell = u$ . A sharp lower bound on  $N$  equals  $n - 1$ .

*Proof* Since  $\ell = u$ , there exists a single ground time series  $t$  of length  $n$  over  $[\ell, u]$ . All the time-series variables of  $t$  have the same value, namely  $\ell$ , and thus its signature consists of  $n - 1$  equalities. The number of occurrences of  $\sigma$  in the signature of  $t$  equals  $n - 1$ , which is thus a sharp lower bound on  $N$ .  $\square$

**Proposition 2** Consider a  $\text{NB}_\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with  $\sigma$  being the **SteadySequence** regular expression, and with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$  such that  $\ell = u$ . A sharp lower bound on  $N$  equals 1.

*Proof* Since  $\ell = u$ , there exists a single ground time series  $t$  of length  $n$  over  $[\ell, u]$ . All the time-series variables of  $t$  have the same value, namely  $\ell$ , and thus its signature consists of  $n - 1$  equalities. The number of occurrences of  $\sigma$  in the signature of  $t$  equals 1, which is thus a sharp lower bound on  $N$ .  $\square$

#### 4.2 Step 1: A Not Necessarily Sharp Upper Bound

Consider a  $\text{NB}_\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ . Lemma 2 of this section provides an upper bound, not necessarily sharp, on  $N$ . Intuitively, to get a maximal number of  $\sigma$ -patterns, every extended  $\sigma$ -pattern should be as short as possible and every two consecutive extended  $\sigma$ -patterns should have a maximal number of time-series variables in common. Although, it is not sharp in general, it is sharp for all regular expressions in Table 1, except **Decreasing**, **Increasing**, **DecreasingTerrace**, and **IncreasingTerrace**.

We first define the notion of *interval without restart*, in order to identify a subseries such that every two consecutive extended  $\sigma$ -patterns within this subseries have  $o_\sigma^{\langle \ell, u \rangle}$  common time-series variables. This notion will be reused in Section 4.3 for deriving a sharp upper bound on  $N$ .

**Definition 17** Consider a regular expression  $\sigma$  and a ground time series  $X = \langle X_1, X_2, \dots, X_n \rangle$  over  $[\ell, u]$ . An *interval without restart* of  $X$  is any interval  $[\alpha, \beta]$  (with  $1 \leq \alpha \leq \beta \leq n$ ), such that all the following conditions hold:

- (1) Every  $X_k$  (with  $k \in [\alpha, \beta]$ ) belongs to at least one extended  $\sigma$ -pattern for which all time-series variables have indices in  $[\alpha, \beta]$ .
- (2) The width of every extended  $\sigma$ -pattern whose time-series variable indices are in  $[\alpha, \beta]$  is equal to  $\omega_\sigma + 1$ .
- (3) Every pair of consecutive extended  $\sigma$ -patterns, whose time-series variables indices are in  $[\alpha, \beta]$ , share  $o_\sigma^{\langle \ell, u \rangle}$  time-series variables.
- (4) When  $o_\sigma^{\langle \ell, u \rangle} > 0$  every extended  $\sigma$ -pattern, whose time-series variables indices are in  $[\alpha, \beta]$ , has no common time-series variables with any extended  $\sigma$ -pattern that has an index outside  $[\alpha, \beta]$ .

Note that, by Condition (4) of Definition 17, the intervals without restart of a ground time series are always disjoint. Consequently two consecutive extended  $\sigma$ -patterns belonging to distinct intervals without restart do not share any time-series variable.

*Example 15* We consider an example of intervals without restart for the  $\sigma = \text{DecreasingTerrace}$  regular expression. For the time series  $X = \langle 4, 3, 3, 2, 2, 1, 4, 2, 2, 1 \rangle$ , the intervals  $[1, 6]$  and  $[7, 10]$  are intervals without restart corresponding to the subseries  $t_1 = \langle 4, 3, 3, 2, 2, 1 \rangle$  and  $t_2 = \langle 4, 2, 2, 1 \rangle$ , because:

- \* Each  $X_i$  (with  $i \in [1, 6]$  or  $i \in [7, 10]$ ) belongs to at least one extended  $\sigma$ -pattern (Condition (1) of Definition 17).
- \* The subseries  $t_1$  (respectively  $t_2$ ) contains 2 (respectively 1) extended  $\sigma$ -patterns of shortest length 4 (Condition (2) of Definition 17).
- \* The two consecutive extended  $\sigma$ -patterns of  $t_1$  have  $o_\sigma^{\langle 1, 4 \rangle} = 2$  time-series variables in common (Condition (3) of Definition 17).
- \* There is no extended  $\sigma$ -pattern straddling between  $[1, 6]$  and  $[7, 10]$  (Condition (4) of Definition 17).  $\triangle$

**Lemma 2** Consider a regular expression  $\sigma$ , and a time series  $X = \langle X_1, X_2, \dots, X_n \rangle$ , with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$  such that  $o_\sigma^{\langle \ell, u \rangle} \leq \omega_\sigma$ .

- (i) The number of  $\sigma$ -patterns in  $X$  is bounded by  $\left\lfloor \frac{\max(0, n - o_\sigma^{\langle \ell, u \rangle})}{\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}} \right\rfloor$ .
- (ii) In addition, if  $n \leq \omega_\sigma$  or there exists at least one ground time series of length  $n$  over  $[\ell, u]$  that contains a single interval without restart longer than  $n - \omega_\sigma - 1 + o_\sigma^{\langle \ell, u \rangle}$ , then the mentioned upper bound is sharp.

*Proof* Since  $o_\sigma^{\langle \ell, u \rangle} \leq \omega_\sigma$ , the denominator  $\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}$  of the considered bound, is always positive. When  $n \leq \omega_\sigma$  the formula  $\left\lfloor \frac{\max(0, n - o_\sigma^{\langle \ell, u \rangle})}{\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}} \right\rfloor$  gives 0 as the result, which is the upper bound on  $N$ . Consider now the case when  $n > \omega_\sigma$ .

[Proof of (i)] We construct a time series  $t$  that we prove to have the maximum number of  $\sigma$ -patterns among all ground time series of length  $n$  without considering any domain restrictions.

- ◇ We assume that the constructed time series  $t$  has a single interval without restart, which is longer than  $n - \omega_\sigma - 1 + o_\sigma^{\langle \ell, u \rangle}$ . Note that such a time series may not be feasible over  $[\ell, u]$ .
- ◇ By definition of an interval without restart, every pair of consecutive extended  $\sigma$ -patterns of  $t$  has  $o_\sigma^{\langle \ell, u \rangle}$  common time-series variables. In addition, every extended  $\sigma$ -pattern has exactly  $\omega_\sigma + 1$  time-series variables and every time-series variable whose indice is in the interval without restart belongs to at least one extended  $\sigma$ -pattern.
- ◇ We now prove that, for any ground time series, the number of  $\sigma$ -patterns cannot exceed the number of  $\sigma$ -patterns of the constructed time series  $t$ .
  - \* Assume that this is not true, then there exists a ground time series whose extended  $\sigma$ -patterns are either strictly shorter than  $\omega_\sigma + 1$  or have a number common time-series variables strictly greater than  $o_\sigma^{\langle \ell, u \rangle}$ .
  - \* Neither of this statements can be possible by construction of  $t$  and the definitions of  $\omega_\sigma$  and  $o_\sigma^{\langle \ell, u \rangle}$ .
  - \* Since the smallest length of an extended  $\sigma$ -pattern equals  $\omega_\sigma + 1$ , and since the number of time-series variables outside the interval without restart of  $t$  is strictly smaller than  $\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}$ , the time series  $t$  does not have any  $\sigma$ -pattern outside of its single interval without restart.
  - \* Hence,  $t$  has the maximum number of  $\sigma$ -patterns compared to all ground time series of length  $n$ .

Let us now estimate the maximum number  $P$  of potential  $\sigma$ -patterns in the time series  $t$ . From the construction of  $t$  we have

$$\underbrace{\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle} + \omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle} + \dots + \omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}}_{P-1 \text{ times}} + \underbrace{\omega_\sigma + 1}_{1 \text{ time}} + n_r = n, \quad (3)$$

where  $n_r$  is the number of time-series variables outside the interval without restart of  $t$ . From Equality (3) and from  $n_r < \omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}$  we obtain

$$P \cdot (\omega_\sigma + 1) - (P - 1) \cdot o_\sigma^{\langle \ell, u \rangle} + n_r = n \Rightarrow P = \left\lfloor \frac{n - o_\sigma^{\langle \ell, u \rangle}}{\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}} \right\rfloor \quad (4)$$

From the right-hand side of Implication (4) we have that the maximum number of  $\sigma$ -patterns among all time series of length  $n$  over  $[\ell, u]$  is less than or equal to  $\left\lfloor \frac{n - o_\sigma^{\langle \ell, u \rangle}}{\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}} \right\rfloor$ .

[Proof of (ii)] If the time series  $t$  constructed in the first part of this proof is feasible wrt  $[\ell, u]$ , then the obtained bound is sharp.  $\square$



#### 4.3 Step 2: Extending the Upper Bound to Get a Sharp Bound Under Some Hypothesis

Consider a  $\text{NB}_\sigma(X, N)$  time-series constraint with  $X = \langle X_1, X_2, \dots, X_n \rangle$ , where every  $X_i$  (with  $i \in [1, n]$ ) is over the same integer interval domain  $[\ell, u]$ . This section focusses on computing a *sharp* upper bound on  $N$  under some hypothesis on the characteristics of  $\sigma$ .

##### 4.3.1 Required Properties of Regular Expressions

Building in a greedy way a time-series that maximises the number of  $\sigma$ -pattern occurrences requires finding a pair of words in  $\mathcal{L}_\sigma$  such that the superposition of these two words wrt an integer interval domain *simultaneously optimises* several characteristics. Depending on the value of the overlap of  $\sigma$  wrt  $\langle \ell, u \rangle$ , we have either the  $\overline{\text{NB}}\text{-overlap}$  property when  $o_\sigma^{\langle \ell, u \rangle} > 0$ , introduced in Property 3, or the  $\overline{\text{NB}}\text{-no-overlap}$  property when  $o_\sigma^{\langle \ell, u \rangle} = 0$ , introduced in Property 4.

- The  $\overline{\text{NB}}\text{-overlap}$  property holds when there exists a pair of words in  $\mathcal{L}_\sigma$ , whose lengths and heights are minimum, and both the overlap and the smallest variation of maxima are reached for a superposition of these words, which is not a factor of any word in  $\mathcal{L}_\sigma$ .
- The  $\overline{\text{NB}}\text{-no-overlap}$  property holds when there exists a word in  $\mathcal{L}_\sigma$ , whose length and height are minimum, and this word can be a maximal occurrence of  $\sigma$  in the signature of a time series over  $[\ell, u]$ .

*Property 3* A regular expression  $\sigma$  has the  $\overline{\text{NB}}\text{-overlap}$  property for an integer interval domain  $[\ell, u]$ , if there exists a pair of not necessarily distinct words  $v$  and  $w$  in  $\mathcal{L}_\sigma$ , and there exists a superposition  $z_1$  (respectively  $z_2$ ) of  $v$  and  $w$  (respectively  $w$  and  $v$ ) wrt  $\langle \ell, u \rangle$ , i.e.,  $o_\sigma^{\langle \ell, u \rangle} > 0$ , such that the following conditions are all satisfied:

- (i)  $\eta_\sigma(v) = \eta_\sigma(w) = \eta_\sigma$ , i.e.,  $v$  and  $w$  have their heights being equal to the height of  $\sigma$ .
- (ii)  $|v| = |w| = \omega_\sigma$ , i.e.,  $v$  and  $w$  are shortest words in  $\mathcal{L}_\sigma$ .
- (iii)  $|v| + |w| - |z_1| + 1 = |w| + |v| - |z_2| + 1 = o_\sigma^{\langle \ell, u \rangle} \leq \omega_\sigma$ , i.e., the overlap between  $v$  and  $w$  (respectively  $w$  and  $v$ ) wrt  $\langle \ell, u \rangle$  is maximum, and its value is bounded by the width of  $\sigma$ .
- (iv) Both superpositions  $z_1$  and  $z_2$  are not factors of any word in  $\mathcal{L}_\sigma$ .
- (v)

$$\delta_\sigma^{\langle \ell, u \rangle} = \begin{cases} \bar{v}_\sigma(z_1, v, 1) - \bar{v}_\sigma(z_1, w, 1) = \bar{v}_\sigma(z_2, w, 1) - \bar{v}_\sigma(z_2, v, 1), & \text{if } v \neq w \\ \bar{v}_\sigma(z_1, v, 1) - \bar{v}_\sigma(z_1, w, 2), & \text{if } v = w, \end{cases}$$

i.e., the smallest variation of maxima of superpositions of  $v$  and  $w$  (respectively  $w$  and  $v$ ) wrt  $\langle \ell, u \rangle$  is reached for their superposition  $z_1$  (respectively  $z_2$ ), and is equal to the smallest variation of maxima of  $\sigma$  wrt  $\langle \ell, u \rangle$ .

- (vi)  $\eta_\sigma(z_1) = \eta_\sigma(z_2) = \eta_\sigma + |\delta_\sigma^{\langle \ell, u \rangle}|$ , i.e., the height of each of these two superpositions  $z_1$  and  $z_2$  is the height of  $\sigma$  plus the absolute value of the smallest variation of maxima of  $\sigma$  wrt  $\langle \ell, u \rangle$ .
- (vii) If  $\delta_\sigma^{\langle \ell, u \rangle} > 0$  (respectively  $\delta_\sigma^{\langle \ell, u \rangle} < 0$ ), then neither ' $v <$ ' (respectively ' $v >$ ') nor ' $w <$ ' (respectively ' $w >$ ') is a factor of any word in  $\mathcal{L}_\sigma$ .

Every regular expression  $\sigma$  in Table 1 has the  $\overline{\text{NB}}\text{-overlap}$  property for any integer interval domain  $[\ell, u]$  such that  $o_\sigma^{\langle \ell, u \rangle} > 0$ .

*Example 16* We now illustrate the  $\overline{\text{NB}}\text{-overlap}$  property on two regular expressions.

- The  $\sigma = \text{DecreasingTerrace}$  regular expression has the  $\overline{\text{NB}}\text{-overlap}$  property for the integer interval domain  $[\ell, u]$  such that  $u - \ell \geq 3$ , because there exists a pair of words  $v = w = '>=>'$  in  $\mathcal{L}_\sigma$  and their superposition  $z = '>=>=>'$  wrt  $\langle \ell, u \rangle$ , such that all the following conditions are satisfied:



- \*  $\eta_\sigma(v) = \eta_\sigma(w) = \eta_\sigma = 2.$  (Cond. (i) of Prop. 3)
- \*  $|v| = |w| = \omega_\sigma = 3.$  (Cond. (ii) of Prop. 3)
- \*  $|v| + |w| - |z| + 1 = o_\sigma^{\langle \ell, u \rangle} = 2 \leq \omega_\sigma = 3.$  (Cond. (iii) of Prop. 3)
- \* Since any word in  $\mathcal{L}_\sigma$  contains only consecutive equalities, the word  $z$  is not a factor of any word in  $\mathcal{L}_\sigma.$  (Cond. (iv) of Prop. 3)
- \*  $\delta_\sigma^{\langle \ell, u \rangle} = \bar{v}_\sigma^{\langle \ell, u \rangle}(z, v, 1) - \bar{v}_\sigma^{\langle \ell, u \rangle}(z, w, 2) = -1.$  (Cond. (v) of Prop. 3)
- \* The height of  $z$  is 3, which equals  $\eta_\sigma + |\delta_\sigma^{\langle \ell, u \rangle}|.$  (Cond. (vi) of Prop. 3)
- \* No word in  $\mathcal{L}_\sigma$  has ' $<$ ', thus ' $v <$ ' is not a factor of any word in  $\mathcal{L}_\sigma.$  (Cond. (vii) of Prop. 3)
- The  $\sigma = \text{SteadySequence}$  regular expression does not have the  $\overline{\text{NB}}$ -overlap property for any integer interval domain  $[\ell, u]$ , because for any pair of words  $v, w$  in  $\mathcal{L}_\sigma$ , the set of superpositions of  $v$  and  $w$  wrt  $\langle \ell, u \rangle$  is empty, and thus  $o_\sigma^{\langle \ell, u \rangle} = 0.$   $\triangle$

*Property 4* A regular expression  $\sigma$  has the  $\overline{\text{NB}}$ -no-overlap property for an integer interval domain  $[\ell, u]$ , if  $o_\sigma^{\langle \ell, u \rangle} = 0$  and if there exists a word  $v$  in  $\mathcal{L}_\sigma$  such that all the following conditions are satisfied:

- (i)  $|v| = \omega_\sigma$ , i.e.,  $v$  is a shortest word in  $\mathcal{L}_\sigma$ .
- (ii)  $\eta_\sigma(v) = \eta_\sigma$ , i.e.,  $v$  has a minimum height among all words in  $\mathcal{L}_\sigma$ .
- (iii) Either both words ' $v >$ ' and ' $v <$ ' are not factors of any word in  $\mathcal{L}_\sigma$ , or at least one of the three words ' $v > v$ ', ' $v < v$ ', ' $v = v$ ' is not a factor of any word in  $\mathcal{L}_\sigma$ , and its height is equal to  $\eta_\sigma$ .
- (iv) For any integer  $n > \omega_\sigma$ , there exists at least one ground time series of length  $n$  over  $[\ell, u]$ , whose signature contains  $v$  as a maximal occurrence of  $\sigma$ .

Any regular expression  $\sigma$  in Table 1 has the  $\overline{\text{NB}}$ -no-overlap property for any integer interval domain  $[\ell, u]$  such that  $o_\sigma^{\langle \ell, u \rangle} = 0$ , except the **SteadySequence** regular expression for  $[\ell, u]$  such that  $\ell = u$ . The case of **SteadySequence** when  $\ell = u$  is discussed in Example 17.

*Example 17* We illustrate the  $\overline{\text{NB}}$ -no-overlap property on two regular expressions.

- The  $\sigma = \text{DecreasingTerrace}$  regular expression has the  $\overline{\text{NB}}$ -no-overlap property for any integer interval domain  $[\ell, u]$  such that  $u - \ell = 2$  because (1) as shown in Example 9, for any two words of  $\mathcal{L}_\sigma$ , the set of their superpositions wrt  $\langle \ell, u \rangle$  is empty, and (2) there exists a word  $v = '>=>'$  in  $\mathcal{L}_\sigma$  that satisfies all the following conditions:
  - \*  $|v| = \omega_\sigma = 3.$  (Cond. (i) of Prop. 4)
  - \*  $\eta_\sigma(v) = \eta_\sigma = 2.$  (Cond. (ii) of Prop. 4)
  - \* The word ' $v < v$ ' is not a factor of any word in  $\mathcal{L}_\sigma$ , and its height is 2. (Cond. (iii) of Prop. 4)
  - \* For any integer  $n > \omega_\sigma$ , there exists a ground time series of length  $n$  over  $[\ell, u]$  whose signature contains  $v$  as a maximal occurrence of  $\sigma$ . (Cond. (iv) of Prop. 4)
- Consider the  $\sigma = \text{SteadySequence}$  regular expression.
  - \* First,  $\sigma$  does not have the  $\overline{\text{NB}}$ -no-overlap property for an integer interval domain  $[\ell, u]$  such that  $u - \ell = 0$ , since Condition (iv) of Property 4 is violated: the shortest word of  $\mathcal{L}_\sigma$ , namely  $v = '='$  cannot be a maximal occurrence of  $\sigma$  in the signature of any ground time series longer than 2 over  $[\ell, u]$ ; indeed, for any time-series length, there exists a single ground time series with all equal values, thus its signature contains only equalities, which prevents  $v$  to be a maximal occurrence of  $\sigma$ .
  - \* Second,  $\sigma$  has the  $\overline{\text{NB}}$ -no-overlap property for an integer interval domain  $[\ell, u]$  such that  $u - \ell > 0$  because there exists a word  $v = '='$  in  $\mathcal{L}_\sigma$  that satisfies all the following conditions:
    - $|v| = \omega_\sigma = 1.$  (Cond. (i) of Prop. 4)
    - $\eta_\sigma(v) = \eta_\sigma = 0.$  (Cond. (ii) of Prop. 4)
    - No word of  $\mathcal{L}_\sigma$  contains ' $>$ ' or ' $<$ ', hence neither ' $v >$ ', nor ' $v <$ ' are factors of any word in  $\mathcal{L}_\sigma$ . (Cond. (iii) of Prop. 4)

- For any integer  $n > \omega_\sigma$ , there exists a ground time series of length  $n$  over  $[\ell, u]$  whose signature contains  $v$  as a maximal occurrence of  $\sigma$ . (Cond. (iv) of Prop. 4)  $\triangle$

#### 4.3.2 Structure of a Maximal Time Series

Consider a  $\text{NB}_\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with every  $X_i$  having the same integer interval domain  $[\ell, u]$ . Lemma 3 describes the structure of a maximal time series for  $\text{NB}_\sigma(\langle X_1, \dots, X_n \rangle, N)$  under the hypothesis that  $\sigma$  has either the  $\overline{\text{NB}}$ -overlap or the  $\overline{\text{NB}}$ -no-overlap property for  $[\ell, u]$ .

**Lemma 3** *Consider a regular expression  $\sigma$  that has either the  $\overline{\text{NB}}$ -overlap or the  $\overline{\text{NB}}$ -no-overlap property for an integer interval domain  $[\ell, u]$ . Then for any integer number  $n > \omega_\sigma$ , there exists a word  $q_{opt}$  such that any ground time series  $t$  of length  $n$  over  $[\ell, u]$  whose signature contains  $q_{opt}$  has the maximum number of  $\sigma$ -patterns among all ground time series of length  $n$  over  $[\ell, u]$ .*

*Proof* We first construct a word  $q_{opt}$  and we show that there is at least one time series of length  $n$  over  $[\ell, u]$  whose signature contains  $q_{opt}$ . Then, we prove that any time series  $t$  of length  $n$  over  $[\ell, u]$  whose signature contains  $q_{opt}$  is maximal for the  $\text{NB}_\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with every  $X_i$  ranging over  $[\ell, u]$ .

**Case (1):**  $\sigma$  has the  $\overline{\text{NB}}$ -overlap property for  $[\ell, u]$ .

Then there exist two words  $v$  and  $w$  of  $\mathcal{L}_\sigma$  and a superposition  $z_1$  (respectively  $z_2$ ) of  $v$  and  $w$  (respectively  $w$  and  $v$ ) wrt  $\langle \ell, u \rangle$  such that all the six conditions of Property 3 are satisfied. Let  $w_1$  and  $w_2$  be the words such that  $z_1 = vw_2$  and  $z_2 = ww_1$ . The figure on the right shows the relations between the words  $z_1, z_2, v, w, w_1$ , and  $w_2$ .

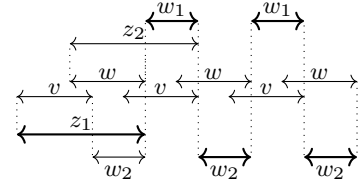


Fig. 3: Illustration of the word  $z_1w_1w_2w_1w_2$  belonging to the language of ' $v | z_1(w_1w_2)^*(w_1 | \varepsilon)$ '

- **Step 1:** Construction of the word  $q_{opt}$ .

When constructing the word  $q_{opt}$  we consider two cases.

- \* **Case (1.1):** The variation of maxima  $\delta_\sigma^{(\ell, u)}$  equals zero.

In this case,  $t$  has a single interval without restart that contains all  $\sigma$ -patterns of  $t$ . We construct the signature  $q_{opt}$  of this interval without restart by imposing the following conditions:

- (a) The word  $q_{opt}$  is in the language of the ' $v | z_1(w_1w_2)^*(w_1 | \varepsilon)$ ' regular expression.
- (b) The length of  $q_{opt}$  is less than  $n$ .
- (c) The length of  $q_{opt}$  is maximum among all words satisfying Conditions (a), and (b).

By condition (i) of Property 3, the heights of both  $v$  and  $w$  equal  $\eta_\sigma$ , the height of  $\sigma$ . Since  $\delta_\sigma^{(\ell, u)} = 0$ , by Conditions (v) and (vi) of Property 3, the height of both words  $z_1$  and  $z_2$  is  $\eta_\sigma$ . Hence, the height of  $q_{opt}$  is also  $\eta_\sigma$ , thus  $q_{opt}$  indeed appears in the signature of some ground time series of length  $n$  over  $[\ell, u]$ , and  $t$  is feasible.

- \* **Case (1.2):** The variation of maxima  $\delta_\sigma^{(\ell, u)}$  does not equal zero.

For brevity, we consider only the case when  $\delta_\sigma^{(\ell, u)} > 0$ , the case of a negative  $\delta_\sigma^{(\ell, u)}$  being symmetric. The time series  $t$  may have  $p \geq 1$  intervals without restart, hence in order to construct  $q_{opt}$  we first construct the signature  $\tilde{q}_{opt}$  of every, except possibly the last one, interval without restart of  $t$  by imposing the following conditions:

- (d) The word  $\tilde{q}_{opt}$  is in the language of the ' $v | z_1(w_1w_2)^*(w_1 | \varepsilon)$ ' regular expression.
- (e) The set of supporting time series of  $\tilde{q}_{opt}$  wrt  $\langle \ell, u \rangle$  is not empty.
- (f) The length of  $\tilde{q}_{opt}$  is less than  $n$ .
- (g) The length of  $\tilde{q}_{opt}$  is maximum among all words satisfying Conditions (d), (e) and (f).

Note that  $\tilde{q}_{opt}$  always exists, since there is at least one word, namely  $v$ , satisfying Conditions (d), (e) and (f). Then, the word  $q_{opt}$  must satisfy the following conditions:

- (h) The word  $q_{opt}$  belongs to the language of the  $(\tilde{q}_{opt} >)^* \tilde{q}_{rest}$  regular expression, where  $\tilde{q}_{rest}$  is a word in the language of the  $v | z_1(w_1 w_2)^*(w_1 | \varepsilon)$  regular expression such that  $|\tilde{q}_{rest}| \leq |\tilde{q}_{opt}|$ .
- (i) The length of  $q_{opt}$  is less than  $n$ .
- (j) The length of  $q_{opt}$  is maximum among all words satisfying Conditions (h) and (i).

Since  $\delta_\sigma^{(\ell, u)} > 0$ , by Lemma 1, and by construction of  $\tilde{q}_{opt}$ , the word  $\tilde{q}_{opt}$  does not contain any ' $>$ '. Then, the concatenation of  $\tilde{q}_{opt}$  and ' $>$ ' has the same height as  $\tilde{q}_{opt}$ . Hence, the height of  $q_{opt}$  equals the height of  $\tilde{q}_{opt}$ , whose set of supporting time series wrt  $\langle \ell, u \rangle$  is not empty, thus  $q_{opt}$  indeed appears in the signature of some ground time series of length  $n$  over  $[\ell, u]$ , and  $t$  is feasible.

◦ **Step 2:** *Maximality of any time series  $t$  whose signature contains  $q_{opt}$ .*

We now prove that  $t$  is a maximal time series for  $\text{NB}_\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$ .

\* First, we show that the number  $p$  of  $\sigma$ -patterns of  $t$  equals the number of occurrences of the words  $v$  and  $w$  in its signature. By Conditions (iv) and (vii) of Property 3, the words  $v$  and  $w$  appearing in  $q_{opt}$  cannot be factors of any other occurrence of  $\sigma$  in  $q_{opt}$ , hence  $p$  is not less than the number of occurrences of the words  $v$  and  $w$  in  $q_{opt}$ . By Conditions (iii) of Property 3, no extended  $\sigma$ -pattern can straddle between two other extended  $\sigma$ -patterns. In addition, by the maximality of the length of  $q_{opt}$  there is no occurrence of  $\sigma$  in the part of the signature of  $t$  that is not  $q_{opt}$ . Hence, neither is  $p$  greater than the number of occurrences of the words  $v$  and  $w$  in  $q_{opt}$ , and thus these values are equal.

\* Second, we prove that  $t$  is maximal for  $\text{NB}_\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$ . Suppose that  $t$  is not maximal for  $\text{NB}_\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  and there exists a time series  $t'$  of length  $n$  over  $[\ell, u]$  that has a number of  $\sigma$ -patterns strictly greater than  $t$ . Then at least one of the following conditions must be satisfied:

- (k) There is a smaller number of intervals without restart of the same total length.
- (l) Some extended  $\sigma$ -patterns of such a time series are of length shorter than  $\omega_\sigma + 1$ .
- (m) Some pairs of consecutive extended  $\sigma$ -patterns have more common time-series variables than  $o_\sigma^{(\ell, u)}$ .
- (n) There is an extended  $\sigma$ -pattern that straddles between two other extended  $\sigma$ -patterns.

Assumption (k) contradicts Condition (v) of Property 3 and the construction of the signature of intervals without restart. Assumptions (l) and (m) contradict Conditions (ii) and (iii) of Property 3. Finally, Assumption (n) is not possible because of the bound imposed on the value of the overlap in Condition (iii) of Property 3. Hence,  $t$  has the maximum number of  $\sigma$ -patterns among all ground time series of the same length over  $[\ell, u]$ .

**Case (2):**  $\sigma$  has the  $\overline{\text{NB}}$ -no-overlap property for  $[\ell, u]$ .

There exists a word  $v$  such that all the conditions of Property 4 are satisfied. The construction of  $q_{opt}$  is similar to Case (1), but the word  $q_{opt}$  will always be the signature of a single interval without restart. The word  $q_{opt}$  is built using the following rules:

- (o) If both words ' $v >$ ' and ' $v <$ ' are not proper factors of any word in  $\mathcal{L}_\sigma$ , then  $q_{opt}$  is in the language of the  $(v > v <)^* v$  regular expression.
- (p) If at least one word  $w$  in  $\{v >, v =, v <\}$  is not a proper factor of any word in  $\mathcal{L}_\sigma$ , and its height equals  $\eta_\sigma$ , then  $q_{opt}$  is in the language of the  $w^* v$  regular expression.
- (q) The length of  $q_{opt}$  is less than  $n$ .
- (r) The length of  $q_{opt}$  is maximum among all words satisfying Conditions (o), (p), and (q).

Since all the conditions of Property 4 are satisfied, it can be shown that the height of  $q_{opt}$  is not greater than  $u - \ell$ , and thus at least one time series of length  $n$  over  $[\ell, u]$  contains  $q_{opt}$  in its signature. Then, in a similar fashion as in Case (1), one can prove that any time series whose signature contains  $q_{opt}$  is maximal for  $\text{NB}_\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$ .  $\square$

#### 4.3.3 A Sharp Upper Bound on the Number of Occurrences of Regular Expression

Consider a  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ . First, Lemma 4 gives an upper bound on the maximum length of an interval without restart in a time series over  $[\ell, u]$ . Second, based on this upper bound and the structure of a maximal time series for  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  showed in Lemma 3, Theorem 2 provides a sharp upper bound on  $N$  under some hypothesis on the characteristics of the regular expression  $\sigma$ .

**Lemma 4** *Consider a regular expression  $\sigma$  and an integer interval domain  $[\ell, u]$  such that one of the following conditions is satisfied:*

- (i) *The value of  $\delta_\sigma^{\langle \ell, u \rangle}$  equals zero.*
- (ii) *The value of  $\delta_\sigma^{\langle \ell, u \rangle}$  does not equal zero and  $\sigma$  has the  $\overline{\text{NB}}$ -overlap property.*

*Then, the maximum length of an interval without restart of any ground time series over  $[\ell, u]$  is bounded by*

$$m_\sigma^{\langle \ell, u \rangle} = \begin{cases} \left\lfloor \frac{u - \ell - \eta_\sigma + |\delta_\sigma^{\langle \ell, u \rangle}|}{|\delta_\sigma^{\langle \ell, u \rangle}|} \right\rfloor \cdot (\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}) + o_\sigma^{\langle \ell, u \rangle}, & \text{if } \delta_\sigma^{\langle \ell, u \rangle} \neq 0, \\ +\infty, & \text{otherwise.} \end{cases}$$

*Proof Case (1): Condition (i) is satisfied.*

Since  $\delta_\sigma^{\langle \ell, u \rangle} = 0$ , the condition that the maximum length of an interval without restart is bounded by  $+\infty$  is trivially satisfied. This upper bound reflects the fact that when  $\delta_\sigma^{\langle \ell, u \rangle} = 0$ , the maximum length of an interval without restart does not depend on the domain  $[\ell, u]$ .

**Case (2): Condition (ii) is satisfied.**

Consider now the case when  $\delta_\sigma^{\langle \ell, u \rangle} \neq 0$  and  $\sigma$  has the  $\overline{\text{NB}}$ -overlap property. Let  $\tilde{q}_{opt}$  be a word such that (1)  $\tilde{q}_{opt}$  is the signature of an interval without restart of maximum length constructed in Lemma 3 for a time series of some length  $n$  over  $[\ell, u]$ ; (2) for any time series of length  $n' > n$  over  $[\ell, u]$ ,  $\tilde{q}_{opt}$  is also the signature of an interval without restart of maximum length. Note that such  $\tilde{q}_{opt}$  necessarily exists by condition that the set of supporting time series of  $\tilde{q}_{opt}$  wrt  $\langle \ell, u \rangle$  must not be empty. Then, there exists a ground time series  $t$  of length  $n$  over  $[\ell, u]$  whose signature is  $\tilde{q}_{opt}$ . By construction of  $\tilde{q}_{opt}$ , the maximum of every extended  $\sigma$ -pattern of  $t$ , except the first one, is  $|\delta_\sigma^{\langle \ell, u \rangle}|$  units smaller or greater, depending on the sign of  $\delta_\sigma^{\langle \ell, u \rangle}$ , compared to the maximum of the preceding extended  $\sigma$ -pattern. Thus, the maxima of these extended  $\sigma$ -patterns form a monotonously decreasing (respectively increasing) sequence of integer numbers. By Conditions (i), (iii), (iv) and (v) of Property 3, the number of elements of such a sequence is bounded by  $\left\lfloor \frac{u - \ell - \eta_\sigma + |\delta_\sigma^{\langle \ell, u \rangle}|}{|\delta_\sigma^{\langle \ell, u \rangle}|} \right\rfloor$ . Since every extended  $\sigma$ -pattern is of length  $\omega_\sigma + 1$ , has a height  $\eta_\sigma$ , and the number of common time-series variable between two extended  $\sigma$ -patterns equals  $o_\sigma^{\langle \ell, u \rangle}$ , the value  $\left\lfloor \frac{u - \ell - \eta_\sigma + |\delta_\sigma^{\langle \ell, u \rangle}|}{|\delta_\sigma^{\langle \ell, u \rangle}|} \right\rfloor \cdot (\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}) + o_\sigma^{\langle \ell, u \rangle}$  is the maximum length of an interval without restart of a ground time series among all ground time series over  $[\ell, u]$ .

□

**Theorem 2** *Consider a  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ . If  $\sigma$  has either the  $\overline{\text{NB}}$ -overlap or the  $\overline{\text{NB}}$ -no-overlap properties for  $[\ell, u]$ , then a sharp upper bound on  $N$  is*

$$\underbrace{\left\lfloor \frac{\max(0, m - o_\sigma^{\langle \ell, u \rangle})}{\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}} \right\rfloor}_A \cdot \underbrace{\left\lfloor \frac{n}{m} \right\rfloor}_B + \underbrace{\left\lfloor \frac{\max(0, (n \bmod m) - o_\sigma^{\langle \ell, u \rangle})}{\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}} \right\rfloor}_C \quad (5)$$

where:

- $m = \min(n, \max(1, m_\sigma^{\langle \ell, u \rangle}))$ , where  $m_\sigma^{\langle \ell, u \rangle}$  is the upper bound on the maximum length of an interval without restart in a time series over  $[\ell, u]$ , introduced by Lemma 4.
- $A$  is the maximum number of  $\sigma$ -patterns in an interval without restart of maximum length.
- $B$  is the number of intervals without restart of maximum length in a maximal time series for the  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint.
- $C$  is the maximum number of  $\sigma$ -patterns in an interval without restart of non-maximum length in a maximal time series for  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$ .

*Proof* Lemma 3 showed the existence of a word  $q_{opt}$  such that any time series  $t$  of length  $n$  over  $[\ell, u]$  whose signature contains  $q_{opt}$  is maximal for  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$ . Hence, a sharp upper bound on  $N$  can be obtained by counting the number of occurrences of  $\sigma$  in  $q_{opt}$ .

Case (a):  $m_\sigma^{\langle \ell, u \rangle} \geq n - \omega_\sigma + o_\sigma^{\langle \ell, u \rangle}$ . Then,  $t$  contains a single interval without restart longer than  $n - \omega_\sigma + o_\sigma^{\langle \ell, u \rangle}$ . Further, the value of  $\min(n, \max(1, m_\sigma^{\langle \ell, u \rangle}))$  equals  $n$ , and the components  $B$  and  $C$  become respectively equal to 1 and 0, thus Formula (5) simplifies to  $A$ . By Lemma 2, the obtained value is a sharp upper bound on  $N$ .

Case (b):  $m_\sigma^{\langle \ell, u \rangle} < n - \omega_\sigma + o_\sigma^{\langle \ell, u \rangle}$ . Then  $t$  may contain multiple intervals without restart. Furthermore, the length of all intervals without restart of  $t$ , except maybe the last one, equals  $m_\sigma^{\langle \ell, u \rangle}$ . By Lemma 2, the maximum number of  $\sigma$ -patterns within every interval without restart of maximum length is  $\left\lfloor \frac{\max(0, m - o_\sigma^{\langle \ell, u \rangle})}{\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}} \right\rfloor$ , i.e., the term  $A$ . The number of intervals without restart of maximum length is  $\left\lfloor \frac{n}{m} \right\rfloor$ , i.e., the term  $B$ . The last interval without restart of  $t$  may be shorter than  $m_\sigma^{\langle \ell, u \rangle}$ , then its length is computed as  $n \bmod m$ , and the number of  $\sigma$ -patterns in the last interval without restart is computed as  $\left\lfloor \frac{\max(0, (n \bmod m) - o_\sigma^{\langle \ell, u \rangle})}{\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}} \right\rfloor$ , which is  $C$ .  $\square$

*Example 18* Consider a  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ . Let  $\sigma$  be the **DecreasingTerrace** regular expression.

- \* First, assume that  $u - \ell = 2$ , and recall some of the computed characteristics, namely  $o_\sigma^{\langle \ell, u \rangle} = 0$ ,  $\omega_\sigma = 3$  and  $\delta_\sigma^{\langle \ell, u \rangle} = 0$ . It was shown in Example 17 that  $\sigma$  has the  $\overline{\text{NB}}$ -no-overlap property for  $[\ell, u]$ , thus Theorem 2 can be applied for computing a sharp upper bound on  $N$ . By Lemma 4, we have that  $m_\sigma^{\langle \ell, u \rangle} = +\infty$ , and thus a sharp upper bound on  $N$  is  $\left\lfloor \frac{\max(0, \min(n, \max(1, m_\sigma^{\langle \ell, u \rangle})) - o_\sigma^{\langle \ell, u \rangle})}{\omega_\sigma + 1 - o_\sigma^{\langle \ell, u \rangle}} \right\rfloor = \left\lfloor \frac{\max(0, \min(n, \max(1, +\infty)) - 0)}{3 + 1 - 0} \right\rfloor = \left\lfloor \frac{n}{4} \right\rfloor$ .
- \* Second, assume  $u - \ell \geq 3$ , then  $o_\sigma^{\langle \ell, u \rangle}$  is now equal to 2, and  $\delta_\sigma^{\langle \ell, u \rangle}$  is equal to  $-1$ . It was shown in Example 17 that  $\sigma$  has the  $\overline{\text{NB}}$ -overlap property for  $[\ell, u]$ , thus Theorem 2 can be applied for computing a sharp upper bound on  $N$ , and a sharp upper bound on  $N$  is equal to  $\left\lfloor \frac{\max(0, m - 2)}{2} \right\rfloor \cdot \left\lfloor \frac{n}{m} \right\rfloor + \left\lfloor \frac{\max(0, (n \bmod m) - 2)}{2} \right\rfloor$  where  $m = \min(n, \max(1, m_\sigma^{\langle \ell, u \rangle})) = \min(n, \max(1, (u - \ell - 1) \cdot 2 + 2))$ , computed by using Lemma 4.  $\triangle$

All the 22 regular expression in Table 1 have either the  $\overline{\text{NB}}$ -overlap or the  $\overline{\text{NB}}$ -no-overlap property for any integer interval domain  $[\ell, u]$ , except the **SteadySequence** regular expression when  $\ell = u$ . A sharp upper bound on the result variable of a time-series constraint in this case is given in Proposition 3.

#### 4.4 A Sharp Upper Bound: Special Case

Proposition 3 provides a sharp upper bound on the number of occurrences of the **SteadySequence** regular expression in the signature of a time series over an integer interval domain  $[\ell, u]$  such that  $\ell = u$ .

**Proposition 3** *Consider a  $\text{NB\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with  $\sigma$  being the **SteadySequence** regular expression and with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$  such that  $\ell = u$ . A sharp upper bound on  $N$  equals 1.*

*Proof* Since  $\ell = u$ , there exists a single time series of length  $n$  over  $[\ell, u]$ , and all its time-series variables have the same value, namely  $\ell$ . The entire signature of this time series is the word in  $\mathcal{L}_\sigma$ , thus a sharp upper bound on  $N$  equals 1.  $\square$

### 5 Time-Series Constraints with Feature **WIDTH**

We now consider the  $\text{G\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  family of time-series constraints with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ , i.e., the case when the feature is **width**,  $\text{G}$  is in the set  $\{\text{Max}, \text{Min}, \text{Sum}\}$  and  $\sigma$  is a non-fixed length regular expression. Section 5.1 defines Properties 5 and 6 of regular expressions that we use to obtain sharp upper bounds on  $N$ . All the regular expressions in Table 1 have both Properties 5 and 6. Based on these properties, Section 5.2 (respectively Section 5.3) provides a sharp upper bound on  $N$  when  $\text{G}$  is **Max** (respectively **Sum**). Finally, Section 5.4 gives a sharp lower bound on  $N$  when  $\text{G}$  is **Sum**. Note that we do not consider a lower (respectively upper) bound for the case when the aggregator is **Max** (respectively **Min**), since when  $\sigma$  has the NB-simple property (see Property 2) for  $[\ell, u]$ , there exists a time series of length  $n$  over  $[\ell, u]$  that has no  $\sigma$ -patterns, and thus yields the default value of  $N$ , namely  $-\infty$  (respectively  $+\infty$ ). Among the 22 regular expressions in Table 1 only the **Steady** and the **SteadySequence** regular expressions do not have the NB-simple property for a domain with a single element, i.e.,  $\ell = u$ .

#### 5.1 Properties of Regular Expressions

Property 5 is used for deriving a sharp upper bound on  $N$  for a  $\text{MAX\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint. Property 5 requires the range of a regular expression be a monotonically increasing linear function of  $n$ .

*Property 5* A regular expression  $\sigma$  has the WIDTH-max property if the following conditions are all satisfied:

- (i) There exists a shortest word in  $\mathcal{L}_\sigma$  whose height equals  $\eta_\sigma$ , the height of  $\sigma$ .
- (ii) For every time-series length  $n > \omega_\sigma + 1$ , the range of  $\sigma$  wrt  $\langle n \rangle, \phi_\sigma^{(n)}$ , is defined and equals  $e_\sigma \cdot (n - 1 - \eta_\sigma) + c_\sigma + \eta_\sigma$  with  $\langle e_\sigma, c_\sigma \rangle \in \{\langle 0, 0 \rangle, \langle 0, 1 \rangle, \langle 1, 0 \rangle\}$ .

Property 6 is used for deriving a sharp upper bound on  $N$  for a  $\text{SUM\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint.

*Property 6* A regular expression  $\sigma$  has the WIDTH-sum property for an integer interval domain  $[\ell, u]$  if the following conditions are all satisfied:

- (i)  $o_\sigma^{\langle \ell, u \rangle} \leq a_\sigma + b_\sigma$ .
- (ii) If for every time-series length  $n > \omega_\sigma + 1$ , the range of  $\sigma$  wrt  $\langle n \rangle, \phi_\sigma^{(n)}$ , equals  $n - 1$ , then  $a_\sigma, b_\sigma$  and  $o_\sigma^{\langle \ell, u \rangle}$  are all equal to 0, and  $\omega_\sigma$ , the width of  $\sigma$ , is equal to 1.

Condition (i) of Property 6 withdraws from consideration a regular expression  $\sigma$  whose  $\sigma$ -patterns overlap, i.e., some time-series variables belong simultaneously to two  $\sigma$ -patterns, which will be formalised in Lemma 5. Condition (ii) of Property 6 restricts further a class of regular expressions whose range depends linearly on  $n$ .

## 5.2 Upper Bound for $\text{MAX\_WIDTH}_\sigma$

We first consider the case when the aggregator is **Max**, i.e., the  $\text{MAX\_WIDTH}_\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  family of time-series constraints with  $\sigma$  being a non-fixed length regular expression and every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ . To compute a sharp upper bound on  $N$ , we maximise the width of a  $\sigma$ -pattern in  $X = \langle X_1, X_2, \dots, X_n \rangle$ . We do so by detecting a longest word in  $\mathcal{L}_\sigma$  that may appear in the signature of  $X$ . The transition from the length of a  $\sigma$ -pattern to the length of the corresponding word in  $\mathcal{L}_\sigma$  is sound because the width of the  $\sigma$ -pattern is the width of the corresponding word plus 1 and minus the sum of  $a_\sigma$  and  $b_\sigma$ , which are constant parameters of  $\sigma$ , introduced in Table 1.

A trivial but, possibly not sharp upper bound on  $N$  is  $n - a_\sigma - b_\sigma$ . Further, for regular expressions that have the  $\overline{\text{WIDTH-max}}$  property, we show that the sharpness of the mentioned upper bound depends only on the difference between  $u$  and  $\ell$ .

The idea for computing a sharp upper bound on  $N$  when  $\sigma$  has the  $\overline{\text{WIDTH-max}}$  property is to identify the minimum value  $d$  of  $u - \ell$  such that the bound  $n - a_\sigma - b_\sigma$  is still sharp. When  $u - \ell$  is smaller than  $d$  we need to find the maximum value of  $k < n$ , such that  $k - a_\sigma - b_\sigma$  is a sharp upper bound on  $N$  for a  $\text{MAX\_WIDTH}_\sigma(\langle X_1, X_2, \dots, X_k \rangle, N)$  time-series constraint with every  $X_i$  ranging over  $[\ell, u]$ .

The next theorem provides a sharp upper bound on  $N$  when the regular expression  $\sigma$  has the  $\overline{\text{WIDTH-max}}$  property.

**Theorem 3** *Consider a  $\text{MAX\_WIDTH}_\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with  $\sigma$  being a non-fixed length regular expression, and all  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ . If  $\sigma$  has the  $\overline{\text{WIDTH-max}}$  property, then a sharp upper bound on  $N$  is*

$$\begin{cases} n - a_\sigma - b_\sigma & \text{if } u - \ell \geq \phi_\sigma^{(n)} \\ e_\sigma \cdot (u - \ell + 1 - a_\sigma - b_\sigma) + c_\sigma \cdot (\omega_\sigma + 1 - a_\sigma - b_\sigma) & \text{if } u - \ell < \phi_\sigma^{(n)} \end{cases} \quad (1)$$

where  $e_\sigma$  and  $c_\sigma$  are parameters of the regular expression  $\sigma$ , introduced in Property 5.

*Proof* When the regular expression  $\sigma$  has the  $\overline{\text{WIDTH-max}}$  property, the range  $\phi_\sigma^{(n)}$  of  $\sigma$  wrt  $\langle n \rangle$  is a monotonically increasing function of  $n$ . It implies that, if the bound  $n - a_\sigma - b_\sigma$  is sharp for some interval integer domain  $[\ell_1, u_1]$ , then it is also sharp for any interval integer domain  $[\ell_2, u_2]$  such that  $u_2 - \ell_2 > u_1 - \ell_1$ . Hence, the sharpness of the upper bound  $n - a_\sigma - b_\sigma$  depends only on  $u - \ell$ .

[Case (1):  $u - \ell \geq \phi_\sigma^{(n)}$ ]. By definition of  $\phi_\sigma^{(n)}$ , we have that if  $u - \ell \geq \phi_\sigma^{(n)}$ , then there exists a word in  $\mathcal{L}_\sigma$  of length  $n - 1$  whose height is not greater than  $u - \ell$ . Hence,  $n - a_\sigma - b_\sigma$  is a sharp upper bound on  $N$ .

[Case (2):  $u - \ell < \phi_\sigma^{(n)}$ ]. This case requires a more detailed analysis than Case (1). Let us consider the three distinct pairs of  $\langle e_\sigma, c_\sigma \rangle$  from Condition (ii) of Property 5:

- (a) The case of  $\langle e_\sigma, c_\sigma \rangle$  being  $\langle 0, 0 \rangle$ . Since  $u - \ell < 0 \cdot (n - 1 - \eta_\sigma) + 0 + \eta_\sigma = \eta_\sigma$ , the necessary-and-sufficient condition, i.e., Property 1, is not satisfied, thus  $N$  is equal to its default value, namely 0.
- (b) The case of  $\langle e_\sigma, c_\sigma \rangle$  being  $\langle 0, 1 \rangle$ . Since  $u - \ell < 0 \cdot (n - 1 - \eta_\sigma) + 1 + \eta_\sigma = \eta_\sigma + 1$ , the only words in  $\mathcal{L}_\sigma$  that can appear in the signature of a ground time series over  $[\ell, u]$  are the ones with the minimum height, namely  $\eta_\sigma$ . For every time-series length  $n > \omega_\sigma + 1$ , we have that  $\phi_\sigma^{(n)} = \eta_\sigma + 1$ , which implies that for every word in  $\mathcal{L}_\sigma$  of length strictly greater than  $\omega_\sigma$ , the height is at least  $\eta_\sigma + 1$ . Hence, only a



word of length  $\omega_\sigma$  and of height  $\eta_\sigma$  can be an occurrence of  $\sigma$  in the signature of a ground time series over  $[\ell, u]$ . By Condition (i) of Property 5, such a word exists in  $\mathcal{L}_\sigma$  and thus, a sharp upper bound on  $N$  is  $\omega_\sigma + 1 - a_\sigma - b_\sigma$ .

- (c) The case of  $\langle e_\sigma, c_\sigma \rangle$  being  $\langle 1, 0 \rangle$ . Since  $u - \ell < 1 \cdot (n - 1 - \eta_\sigma) + 0 + \eta_\sigma = n - 1$ , we have that  $u - \ell < n - 1$ . Hence, we aim at finding the longest time-series length  $k < n$  such that  $u - \ell = k - 1$ , and a sharp upper bound on  $N$  will be  $k - a_\sigma - b_\sigma$ . The largest value of such  $k$  equals  $u - \ell + 1$ , thus a sharp upper bound on  $N$  is  $u - \ell + 1 - a_\sigma - b_\sigma$ .  $\square$

*Example 19* Consider a  $\text{MAX\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with every  $X_i$  having the same integer interval domain  $[\ell, u]$ . The three items of this example cover each value of  $\langle e_\sigma, c_\sigma \rangle$  in the set  $\{\langle 0, 0 \rangle, \langle 0, 1 \rangle, \langle 1, 0 \rangle\}$ .

- Consider the  $\sigma = \text{Inflexion}$  regular expression. Recall that both  $a_\sigma$  and  $b_\sigma$  are equal to 1, the width  $\omega_\sigma$  of  $\sigma$  is equal to 2, the height  $\eta_\sigma$  of  $\sigma$  is equal to 1, and for any time-series length  $n > \omega_\sigma + 1$ , the range  $\phi_\sigma^{(n)}$  of  $\sigma$  wrt  $\langle n \rangle$  is equal to  $e_\sigma \cdot (n - 1 - \eta_\sigma) + c_\sigma + \eta_\sigma = \eta_\sigma = 1$ . Since there exists a word, namely  $v = \langle < > \rangle$ , in  $\mathcal{L}_\sigma$  whose length equals 2 and whose height is equal to 1, and  $\langle e_\sigma, c_\sigma \rangle$  is  $\langle 0, 0 \rangle$ ,  $\sigma$  has the  $\overline{\text{WIDTH}}$ -max property. Hence, we apply Theorem 3 for computing a sharp upper bound on  $N$ .
  - \* If  $u - \ell \geq \phi_\sigma^{(n)} = 1$ , then a sharp upper bound on  $N$  is equal to  $n - a_\sigma - b_\sigma = n - 2$ .
  - \* If  $u - \ell < \phi_\sigma^{(n)} = 1$ , then a sharp upper bound on  $N$  is equal to 0.
- Consider the  $\sigma = \text{Gorge}$  regular expression. Recall that both  $a_\sigma$  and  $b_\sigma$  are equal to 1, the width  $\omega_\sigma$  of  $\sigma$  is equal to 2, the height  $\eta_\sigma$  of  $\sigma$  is equal to 1, and for any time-series length  $n > \omega_\sigma + 1$ , the range  $\phi_\sigma^{(n)}$  of  $\sigma$  wrt  $\langle n \rangle$  is equal to  $e_\sigma \cdot (n - 1 - \eta_\sigma) + c_\sigma + \eta_\sigma = \eta_\sigma + 1 = 2$ . Since there exists a word, namely  $v = \langle > < \rangle$ , in  $\mathcal{L}_\sigma$  whose length equals 2 and whose height is equal to 1, and  $\langle e_\sigma, c_\sigma \rangle$  is  $\langle 0, 1 \rangle$ ,  $\sigma$  has the  $\overline{\text{WIDTH}}$ -max property. Hence, we apply Theorem 3 for computing a sharp upper bound on  $N$ .
  - \* If  $u - \ell \geq 2$ , then a sharp upper bound on  $N$  is equal to  $n - a_\sigma - b_\sigma = n - 2$ .
  - \* If  $u - \ell < 2$ , then a sharp upper bound on  $N$  is equal to  $\omega_\sigma + 1 - a_\sigma - b_\sigma = 1$ .
- Consider the  $\sigma = \text{StrictlyDecreasingSequence}$  regular expression. Recall that both  $a_\sigma$  and  $b_\sigma$  are equal to 0, the width  $\omega_\sigma$  of  $\sigma$  is equal to 1, the height  $\eta_\sigma$  of  $\sigma$  is equal to 1, and for any time-series length  $n > \omega_\sigma + 1$ , the range  $\phi_\sigma^{(n)}$  of  $\sigma$  wrt  $\langle n \rangle$  is equal to  $e_\sigma \cdot (n - 1 - \eta_\sigma) + c_\sigma + \eta_\sigma = n - 1 - \eta_\sigma + \eta_\sigma = n - 1$ . Since there exists a word, namely  $v = \langle > \rangle$ , in  $\mathcal{L}_\sigma$  whose length is equal to 1 and whose height is equal to 1, and  $\langle e_\sigma, c_\sigma \rangle$  is  $\langle 1, 0 \rangle$ ,  $\sigma$  has the  $\overline{\text{WIDTH}}$ -max property. Hence, we apply Theorem 3 for computing a sharp upper bound on  $N$ .
  - \* If  $u - \ell \geq \phi_\sigma^{(n)} = n - 1$ , then a sharp upper bound on  $N$  is equal to  $n - a_\sigma - b_\sigma = n$ .
  - \* If  $u - \ell < \phi_\sigma^{(n)} = n - 1$ , then a sharp upper bound on  $N$  is equal to  $u - \ell + 1$ .  $\triangle$

### 5.3 Upper Bound for $\text{SUM\_WIDTH\_}\sigma$

We now consider the  $\text{SUM\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  family of time-series constraints with  $\sigma$  being a non-fixed length regular expression and with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ . Under some hypothesis on the overlap of  $\sigma$  wrt  $\langle \ell, u \rangle$ , Lemma 5 provides an upper bound on  $N$  and a condition when this bound is sharp. Then, Theorem 4 extends the bound of Lemma 5 and gives a more general condition under which the extended bound on  $N$  is sharp.

**Lemma 5** Consider a  $\text{SUM\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ , and with  $\sigma$  being a non-fixed length regular expression.

- (i) If  $o_\sigma^{(\ell, u)} \leq a_\sigma + b_\sigma$  then  $n - a_\sigma - b_\sigma$  is an upper bound on  $N$ .



(ii) If, in addition,  $u - \ell \geq \phi_\sigma^{(n)}$ , then this bound is sharp.

*Proof* [Proof of (i)] Let us consider a time series  $t$  of length  $n$  over  $[\ell, u]$  that has  $p > 1$   $\sigma$ -patterns. Let  $\omega_i$  be the length of the  $\sigma$ -pattern  $i$  (with  $i$  in  $[1, p]$ ); let  $n_r$  be the number of time-series variables that are not in any extended  $\sigma$ -pattern of  $t$ ; and let  $o_i$  be the number of common time-series variables of the extended  $\sigma$ -patterns  $i$  and  $i + 1$ . Then, the following equality holds

$$n = \omega_1 + a_\sigma + b_\sigma + \sum_{i=1}^{p-1} (\omega_{i+1} + a_\sigma + b_\sigma - o_i) + n_r \quad (1)$$

The time series  $t$  yields  $\sum_{i=1}^p \omega_i$  as the value of  $N$ , thus we express this quantity from Equality (1) and obtain

$$N = n - n_r - p \cdot (a_\sigma + b_\sigma) + \sum_{i=1}^{p-1} o_i \quad (2)$$

In order to prove that  $n - a_\sigma - b_\sigma$  is a valid upper bound on  $N$ , we show that the difference between  $n - a_\sigma - b_\sigma$  and the right-hand side of Equality (2) is always non-negative if  $o_\sigma^{(\ell, u)} \leq a_\sigma + b_\sigma$ .

$$n - (a_\sigma + b_\sigma) - n + n_r + p \cdot (a_\sigma + b_\sigma) - \sum_{i=1}^{p-1} o_i = n_r + (p-1) \cdot (a_\sigma + b_\sigma) - \sum_{i=1}^{p-1} o_i \quad (3)$$

The value of  $n_r$  is non-negative, and by the definition of  $o_\sigma^{(\ell, u)}$ , every  $o_i$  is not greater than  $o_\sigma^{(\ell, u)}$ . In addition, we have the following inequality  $o_\sigma^{(\ell, u)} \leq a_\sigma + b_\sigma$ . Hence, a lower estimate of the right-hand side of Equality (3) is given by the following inequality

$$n_r + (p-1) \cdot (a_\sigma + b_\sigma) - \sum_{i=1}^{p-1} o_i \geq 0 + (p-1) \cdot (a_\sigma + b_\sigma) - (p-1) \cdot (a_\sigma + b_\sigma) = 0 \quad (4)$$

By Inequality (4) we obtain that, when  $o_\sigma^{(\ell, u)} \leq a_\sigma + b_\sigma$ , the difference between  $n - a_\sigma - b_\sigma$  and the value of  $N$  is always non-negative. Hence,  $n - a_\sigma - b_\sigma$  is an upper bound on  $N$ .

[Proof of (ii)] We now show that  $n - a_\sigma - b_\sigma$  is a sharp upper bound on  $N$ , when  $u - \ell \geq \phi_\sigma^{(n)}$ . By definition of  $\phi_\sigma^{(n)}$ , the range of  $\sigma$  wrt  $\langle n \rangle$ , there exists a word  $v$  of length  $n-1$  in  $\mathcal{L}_\sigma$  whose height is at most  $u - \ell$ . Hence, there exists at least one ground time series of length  $n$  over  $[\ell, u]$  whose signature is  $v$ , all its time-series variable belong to a single extended  $\sigma$ -pattern. For such a time series, the value of  $p$  equals 1, and  $n_r$  equals 0. By the right-hand side of Equality (2), we have that  $N$  equals  $n - a_\sigma - b_\sigma - 0 - (1-1)(a_\sigma + b_\sigma) = n - a_\sigma - b_\sigma$ , which was proved to be an upper bound. Hence, in this case  $n - a_\sigma - b_\sigma$  is a sharp upper bound on  $N$ .  $\square$

**Theorem 4** Consider a  $\text{SUM\_WIDTH\_}\sigma((X_1, X_2, \dots, X_n), N)$  time-series constraint with  $\sigma$  being a non-fixed-length regular expression and every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ . If  $\sigma$  has both the  $\overline{\text{WIDTH-max}}$  property and the  $\overline{\text{WIDTH-sum}}$  property for  $[\ell, u]$ , then a sharp upper bound on  $N$  is

$$\begin{cases} n - a_\sigma - b_\sigma & \text{if } u - \ell \geq \phi_\sigma^{(n)} \\ e_\sigma \cdot (n - \rho_\sigma^{(\ell, u, n)}) + c_\sigma \cdot (\omega_\sigma + 1 - a_\sigma - b_\sigma) \cdot \tau_\sigma^{(\ell, u, n)} & \text{if } u - \ell < \phi_\sigma^{(n)} \end{cases} \quad (1)$$

$$\quad (2)$$

where:

$\circ$   $e_\sigma$  and  $c_\sigma$  are parameters of the regular expression  $\sigma$ , introduced in Property 5.

- $\rho_{\sigma}^{\langle \ell, u, n \rangle}$  equals  $\min(1, \max(0, \eta_{\sigma} + 1 - (u - \ell))) \cdot (n \bmod 2)$ .
- $\tau_{\sigma}^{\langle \ell, u, n \rangle}$  is the maximum number of  $\sigma$ -patterns of shortest length in a time series among all ground time series of length  $n$  over  $[\ell, u]$ .

*Proof* When a regular expression  $\sigma$  has the  $\overline{\text{WIDTH}}$ -sum property for  $[\ell, u]$ , Condition (i) of Lemma 5 is satisfied and thus,  $n - a_{\sigma} - b_{\sigma}$  is an upper bound on  $N$ .

[Case (1):  $u - \ell \geq \phi_{\sigma}^{\langle n \rangle}$ ]. Since Condition (ii) of Lemma 5 is also satisfied, by Lemma 5,  $u - \ell \geq \phi_{\sigma}^{\langle n \rangle}$  is a sharp upper bound on  $N$ .

[Case (2):  $u - \ell < \phi_{\sigma}^{\langle n \rangle}$ ]. Let us consider the three potential values of  $\langle e_{\sigma}, c_{\sigma} \rangle$  from Condition (ii) of Property 5:

- The case of  $\langle e_{\sigma}, c_{\sigma} \rangle$  being  $\langle 0, 0 \rangle$ . Since  $u - \ell < \eta_{\sigma}$ , the necessary-sufficient condition, i.e., Property 1, is not satisfied, and thus no word of  $\mathcal{L}_{\sigma}$  can occur in the signature of  $\langle X_1, X_2, \dots, X_n \rangle$ . Hence,  $N$  is equal to its default value, namely 0.
- The case of  $\langle e_{\sigma}, c_{\sigma} \rangle$  being  $\langle 0, 1 \rangle$ . Since  $u - \ell \leq \eta_{\sigma}$ , only a shortest word with a height being  $\eta_{\sigma}$  may occur in a signature of  $\langle X_1, X_2, \dots, X_n \rangle$ , as it was shown in the proof of Theorem 3. By Condition (i) of Property 5, such a word exists, and thus a sharp upper bound on  $N$  is equal to  $\omega_{\sigma} + 1 - a_{\sigma} - b_{\sigma}$ . Hence, any  $\sigma$ -pattern of any ground time series of length  $n$  over  $[\ell, u]$  is of length  $\omega_{\sigma} + 1 - a_{\sigma} - b_{\sigma}$ . Since it is not possible to increase the length of any  $\sigma$ -patterns, in order to maximise  $N$ , it is necessary to maximise the number of  $\sigma$ -patterns of shortest length in a time series of length  $n$  over  $[\ell, u]$ . Since  $\tau_{\sigma}^{\langle \ell, u, n \rangle}$  is the maximum number of  $\sigma$ -patterns of minimum length, a sharp upper bound on  $N$  equals  $(\omega_{\sigma} + 1 - a_{\sigma} - b_{\sigma}) \cdot \tau_{\sigma}^{\langle \ell, u, n \rangle}$ .
- The case of  $\langle e_{\sigma}, c_{\sigma} \rangle$  being  $\langle 1, 0 \rangle$ . When  $\sigma$  has the  $\overline{\text{WIDTH}}$ -sum-property for  $[\ell, u]$ , it belongs to the following class of regular expressions:  $a_{\sigma}, b_{\sigma}, o_{\sigma}^{\langle \ell, u \rangle}$  are all equal to 0, and  $\omega_{\sigma}$  is equal to 1. Consider a time series  $t$  of length  $n$  over  $[\ell, u]$  with  $p \geq 1$   $\sigma$ -patterns, where  $\omega_i$  is the length of the  $\sigma$ -pattern  $i$ ,  $o_i$  is the overlap of the extended  $\sigma$ -patterns  $i$  and  $i + 1$ , and  $\rho_{\sigma}^{\langle \ell, u, n \rangle}$  is the number of time-series variables of  $t$  that do not belong to any extended  $\sigma$ -pattern of  $t$ . Then, the following equality holds

$$N = n - \rho_{\sigma}^{\langle \ell, u, n \rangle} - p \cdot (a_{\sigma} + b_{\sigma}) + \sum_{i=1}^{p-1} o_i$$

In this equality we replace  $a_{\sigma}$ , and  $b_{\sigma}$  with their actual values, namely 0, which gives a simplified equality  $N = n - \rho_{\sigma}^{\langle \ell, u, n \rangle}$ . Since the smaller  $\rho_{\sigma}^{\langle \ell, u, n \rangle}$ , the larger is  $N$ , the aim is to find a time series for which  $\rho_{\sigma}^{\langle \ell, u, n \rangle}$  is minimum. Assume that in such a time series  $p$  equals the maximum number of  $\sigma$ -patterns in a time series among all ground time series of length  $n$  over  $[\ell, u]$ . Then,  $\rho_{\sigma}^{\langle \ell, u, n \rangle}$  is strictly less than  $\omega_{\sigma} + 1 = 2$ , otherwise there would be a contradiction with the maximality of  $p$ . Hence,  $t$  has at most one time-series variable that is outside of any extended  $\sigma$ -pattern of  $t$ . By definition of  $\phi_{\sigma}^{\langle n \rangle}$ , the number of time-series variables in any extended  $\sigma$ -pattern is at most  $u - \ell + 1$ , thus if  $t$  contains at least one  $\sigma$ -pattern shorter than  $u - \ell + 1$  the value of  $\rho_{\sigma}^{\langle \ell, u, n \rangle}$  can be decreased by extending this  $\sigma$ -pattern with one time-series variable. Furthermore, if  $u - \ell \geq \eta_{\sigma} + 1$ , then  $\rho_{\sigma}^{\langle \ell, u, n \rangle} = 0$ , otherwise  $\rho_{\sigma}^{\langle \ell, u, n \rangle} = n \bmod 2$ . Hence, the minimum value of  $\rho_{\sigma}^{\langle \ell, u, n \rangle}$  equals  $\min(1, \max(0, \eta_{\sigma} + 1 - (u - \ell))) \cdot (n \bmod 2)$ .  $\square$

Note that for the 22 regular expressions in Table 1, the maximum number of  $\sigma$ -patterns of shortest length in a time series coincides with the maximum number of  $\sigma$ -patterns in the same time series. Although, in the general case it may not be true.

*Example 20* Consider a  $\text{SUM\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ , and each value of  $\langle e_{\sigma}, c_{\sigma} \rangle$  in  $\{\langle 0, 0 \rangle, \langle 0, 1 \rangle, \langle 1, 0 \rangle\}$ .

- Consider the  $\sigma = \text{Inflexion}$  regular expression. In Example 19, we showed that the regular expression  $\sigma$  has the  $\overline{\text{WIDTH}}$ -max property. Recall that  $o_\sigma^{\langle \ell, u \rangle}$  is equal to 2 and both  $a_\sigma$  and  $b_\sigma$  are equal to 1. Hence, Condition (i) of Property 6 is also satisfied. Since for any time-series length greater than  $\omega_\sigma + 1$ , the value of  $\phi_\sigma^{\langle n \rangle}$  equals  $\eta_\sigma$ , Condition (ii) of Property 6 is trivially satisfied. Hence,  $\sigma$  has also the  $\overline{\text{WIDTH}}$ -sum property, and Theorem 4 can be used for computing a sharp upper bound on  $N$ :
  - \* If  $u - \ell \geq \eta_\sigma = 1$ , then a sharp upper bound on  $N$  is equal to  $n - a_\sigma - b_\sigma = n - 2$ .
  - \* If  $u - \ell < \eta_\sigma = 1$ , then a sharp upper bound on  $N$  is equal to 0.
- Consider the  $\sigma = \text{Gorge}$  regular expression. In Example 19, we showed that the regular expression  $\sigma$  has the  $\overline{\text{WIDTH}}$ -sum property. Recall that  $o_\sigma^{\langle \ell, u \rangle}$  is equal to 1 and both  $a_\sigma$  and  $b_\sigma$  are equal to 1. Hence, Condition (i) of Property 6 is also satisfied. Since for any time-series length greater than  $\omega_\sigma + 1$ , the value of  $\phi_\sigma^{\langle n \rangle}$  equals  $\eta_\sigma + 1$ , Condition (ii) of Property 6 is trivially satisfied. Hence,  $\sigma$  has also the  $\overline{\text{WIDTH}}$ -sum property, and Theorem 4 can be used for computing a sharp upper bound on  $N$ :
  - \* If  $u - \ell \geq 2$ , then a sharp upper bound on  $N$  equals  $n - a_\sigma - b_\sigma = n - 2$ .
  - \* If  $u - \ell < 2$ , then a sharp upper bound on  $N$  is equal to  $\tau_\sigma^{\langle \ell, u, n \rangle} \cdot (\omega_\sigma + 1 - a_\sigma - b_\sigma) = \tau_\sigma^{\langle \ell, u, n \rangle} \cdot (2 + 1 - 1 - 1) = \tau_\sigma^{\langle \ell, u, n \rangle}$ .

For this particular regular expression,  $\tau_\sigma^{\langle \ell, u, n \rangle}$  equals the maximum number of  $\sigma$ -patterns in a time series among all ground time series of length  $n$  over  $[\ell, u]$ , namely  $\lfloor \frac{n-1}{2} \rfloor$ , which is the upper bound obtained in Section 4.
- Consider the  $\sigma = \text{StrictlyDecreasingSequence}$  regular expression. It was shown in Example 19 that  $\sigma$  has the  $\overline{\text{WIDTH}}$ -max-property. Recall that  $o_\sigma^{\langle \ell, u \rangle}$  is equal to 0, and both  $a_\sigma$  and  $b_\sigma$  are equal to 0, thus Condition (i) of Property 6 is also satisfied. Since  $o_\sigma^{\langle \ell, u \rangle}$ ,  $a_\sigma$  and  $b_\sigma$  are all equal to 0, and  $\omega_\sigma$  is equal to 1, Condition (ii) of Property 4 is also satisfied. Hence,  $\sigma$  has the  $\overline{\text{WIDTH}}$ -sum property, and Theorem 4 can be used for computing a sharp upper bound on  $N$ :
  - \* If  $u - \ell \geq n - 1$ , then a sharp upper bound on  $N$  is equal to  $n$ .
  - \* If  $u - \ell < n - 1$ , then a sharp upper bound on  $N$  is equal to  $n - \rho_\sigma^{\langle \ell, u, n \rangle} = n - \min(1, \max(0, (2 - (u - \ell)) \cdot (n \bmod 2)))$ .  $\triangle$

#### 5.4 Lower Bound for $\text{MIN\_WIDTH\_}\sigma$

Finally, consider the  $\text{MIN\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  family of time-series constraints with  $\sigma$  being a non-fixed-length regular expression and with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ . The next theorem, Theorem 5, provides a sharp lower bound on  $N$  assuming the property that we now introduce holds.

*Property 7* A non-fixed-length regular expression  $\sigma$  has the  $\overline{\text{WIDTH}}$ -occurrence property for an integer interval domain  $[\ell, u]$ , if there exist a shortest word  $v$  in  $\mathcal{L}_\sigma$ , i.e.,  $|v| = \omega_\sigma$ , and a word  $w$  in  $\{v <, v =, v >\}$  such that the following conditions are all satisfied:

- (i) The height of  $v$  equals  $\eta_\sigma$ , the height of  $\sigma$ .
- (ii) The height of  $w$  is less than or equal to  $u - \ell$ .
- (iii) The word  $w$  is not a factor of any word in  $\mathcal{L}_\sigma$ .

**Theorem 5** Consider a  $\text{MIN\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with  $\sigma$  being a non-fixed-length regular expression, and with every  $X_i$  having the same integer interval domain  $[\ell, u]$ . If  $\sigma$  has the  $\overline{\text{WIDTH}}$ -occurrence property for  $[\ell, u]$ , then a sharp lower bound on  $N$  equals  $\omega_\sigma + 1 - a_\sigma - b_\sigma$ .

*Proof* Since  $\omega_\sigma$  is the length of a shortest word in  $\mathcal{L}_\sigma$ , the length of any  $\sigma$ -pattern is at least  $\omega_\sigma + 1 - a_\sigma - b_\sigma$ , and thus it is a lower bound on  $N$ . When  $\sigma$  has the WIDTH-occurrence property, there exists a shortest word  $v$  in  $\mathcal{L}_\sigma$  and a word  $w$  in  $\{v <, v =, v >\}$  such that the three conditions of Property 7 are all satisfied. We now show that in this case, the bound is sharp.

Case (a):  $n = \omega_\sigma + 1$ . When Condition (i) of Property 7 is satisfied, there exists a ground time series of length  $n = \omega_\sigma + 1$  over  $[\ell, u]$  whose signature is  $v$ . Hence,  $\omega_\sigma + 1 - a_\sigma - b_\sigma$  is a sharp lower bound on  $N$ .

Case (b):  $n > \omega_\sigma + 1$ . When Condition (ii) of Property 7 is satisfied, there exists a ground time series  $t$  of length  $n$  over  $[\ell, u]$  whose signature is a word in the language of the ' $w = *$ ' regular expression. If Condition (iii) of Property 7 is also satisfied, then the  $v$  in the signature of  $t$  is a maximal occurrence of  $\sigma$ , because  $w$  is not a factor of any word in  $\mathcal{L}_\sigma$ . The length of the corresponding  $\sigma$ -pattern is  $\omega_\sigma + 1 - a_\sigma - b_\sigma$ , thus this value is a sharp lower bound on  $N$ .  $\square$

*Example 21* Consider a  $\text{MIN\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraints with  $\sigma$  being the **Inflexion** regular expression and with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$  such that  $u - \ell \geq \eta_\sigma = 1$ . It was shown in Example 2 that  $\sigma$  is a non-fixed-length regular expression. Furthermore, there exists a word  $v = '< >'$  and a word  $w = '< > ='$  in  $\{v <, v =, v >\}$  such that the following conditions are all satisfied:

- \* The height of  $v$  equals  $\eta_\sigma = 1$ . (Cond. (i) of Prop. 7)
- \* The height of  $w$  equals 1, and thus is less than or equal to  $u - \ell$ . (Cond. (ii) of Prop. 7)
- \* The word  $w$  is not a factor of any word in  $\mathcal{L}_\sigma$ . (Cond. (iii) of Prop. 7)

Hence,  $\sigma$  has the WIDTH-occurrence property for  $[\ell, u]$ , and by Theorem 5, a sharp lower bound on  $N$  equals  $\omega_\sigma + 1 - a_\sigma - b_\sigma = 2 + 1 - 1 - 1 = 1$ .  $\triangle$

All the 22 regular expressions in Table 1 have the WIDTH-occurrence-property for any integer interval domain  $[\ell, u]$ , except the **SteadySequence** regular expression when  $\ell = u$ . This special case is considered in Proposition 4.

**Proposition 4** Consider a  $\text{MIN\_WIDTH\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraints with  $\sigma$  being the **SteadySequence** regular expression and with every  $X_i$  being over an integer interval domain  $[\ell, u]$  such that  $\ell = u$ . A sharp lower bound on  $N$  equals  $n$ .

*Proof* When  $\ell$  equals  $u$ , there exists a single ground time series  $t$  of length  $n$  over  $[\ell, u]$  with all time-series variables having the same value, namely  $\ell$ . The signature of  $t$  is a sequence of  $n - 1$  equalities, which is a word in  $\mathcal{L}_\sigma$ . Hence, every time-series variable of  $t$  belongs to a single extended  $\sigma$ -pattern of  $t$ , and thus a sharp lower bound on  $N$  equals  $n - a_\sigma - b_\sigma = n$ .  $\square$

## 6 Synthesis

Consider a  $\text{G\_F\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraints with every  $X_i$  being over the same integer interval domain  $[\ell, u]$ . Table 2 provides a synthesis of the bounds on  $N$  obtained in Sections 4, 5 and [3], when  $\langle g, f \rangle$  is in  $\{\langle \text{Max}, \text{min} \rangle, \langle \text{Max}, \text{width} \rangle, \langle \text{Min}, \text{width} \rangle, \langle \text{Sum}, \text{one} \rangle, \langle \text{Sum}, \text{width} \rangle\}$ . The theorems and the propositions mentioned in Table 2 were applied for computing sharp bounds on  $N$  for 93 time-series constraints of Volume II of the global constraint catalogue [4]. An entry of Table 2 corresponds to an upper (respectively lower) bound on  $N$  for a  $\text{G\_F\_}\sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  time-series constraint with every  $X_i$  ranging over the same integer interval domain  $[\ell, u]$ , if the corresponding ‘‘Type’’ column contains  $\bar{N}$  (respectively  $\underline{N}$ ). The ‘‘Theorem’’ column contains the theorem or the proposition providing the corresponding sharp bound under the hypothesis that  $\sigma$  has the properties mentioned in the corresponding ‘‘Properties’’

column. The “Theorem” (respectively “Property”) column recalls also the set of characteristics used in the bound of the corresponding theorem or proposition (respectively property).

Note that when the aggregator is **Max** (respectively **Min**) we do not consider a lower (respectively upper) bound on  $N$ . When  $\sigma$  has the **NB-simple** property for  $[\ell, u]$ , there exists a time series of length  $n$  over  $[\ell, u]$  whose signature contains no  $\sigma$ -patterns, and thus such a time series yields the default value of **Max** (respectively **Min**), which is  $-\infty$  (respectively  $+\infty$ ).

$\langle G, F \rangle$	Type	Theorem	Properties
$\langle \text{Sum}, \text{one} \rangle$	$\underline{N}$	Theorem 1	<b>NB-simple</b> ( $\Theta_\sigma$ )
	$\underline{N}$	Proposition 1	$\sigma = \text{Steady}, u = \ell$
	$\underline{N}$	Proposition 2	$\sigma = \text{SteadySequence}, u = \ell$
	$\overline{N}$	Theorem 2 ( $\omega_\sigma, \eta_\sigma, o_\sigma^{(\ell, u)}, \delta_\sigma^{(\ell, u)}$ )	<b>NB-overlap</b> or <b>NB-no-overlap</b> ( $\omega_\sigma, \eta_\sigma, o_\sigma^{(\ell, u)}, \delta_\sigma^{(\ell, u)}$ )
	$\overline{N}$	Proposition 3	$\sigma = \text{SteadySequence}, u = \ell$
$\langle \text{Max}, \text{width} \rangle$	$\overline{N}$	Theorem 3 ( $\omega_\sigma, \eta_\sigma, \phi_\sigma^{(\ell, u)}$ )	<b>WIDTH-max</b> ( $\omega_\sigma, \eta_\sigma, \phi_\sigma^{(\ell, u)}$ )
$\langle \text{Sum}, \text{width} \rangle$	$\overline{N}$	Theorem 4 ( $\omega_\sigma, \eta_\sigma, \phi_\sigma^{(\ell, u)}$ )	<b>WIDTH-max</b> and <b>WIDTH-sum</b> ( $\omega_\sigma, \eta_\sigma, \phi_\sigma^{(\ell, u)}, o_\sigma^{(\ell, u)}$ )
$\langle \text{Min}, \text{width} \rangle$	$\underline{N}$	Theorem 5 ( $\omega_\sigma$ )	<b>WIDTH-occurrence</b> ( $\omega_\sigma, \eta_\sigma$ )
	$\underline{N}$	Proposition 4	$\sigma = \text{SteadySequence}, u = \ell$
$\langle \text{Max}, \text{min} \rangle$	$\overline{N}$	Theorem 1 in [3]	The Condition of Theorem 1 in [3]

Table 2: A synthesis of bounds presented in Sections 4, 5, and in [3].

Table 3 provides for each of the regular expressions in Table 1 the corresponding value of each regular expression characteristics. The 22 regular expressions in Table 1 have the **NB-simple** property for any domain, except the **Steady** and the **SteadySequence** regular expressions when  $\ell = u$ . Table 4 classifies the 22 regular expressions according to the set of properties they share. There are three main groups, and two special ones, namely for the **Steady** and for the **SteadySequence** regular expressions. The partitioning into the three main groups is related to the fact that the entry of Table 2 with Theorem 2, contains a disjunction between the **NB-overlap** and the **NB-no-overlap** properties. Furthermore, a regular expression  $\sigma$  cannot have both properties for the same integer interval domain  $[\ell, u]$ . This allows to partition the 22 regular expressions into three classes, namely:

1. The regular expressions that have the **NB-overlap** property for any  $[\ell, u]$ , i.e., the first group in Table 4.
2. The regular expressions that have the **NB-no-overlap** property for any  $[\ell, u]$ , i.e., the second group in Table 4.
3. The regular expressions that have the **NB-no-overlap** property for any  $[\ell, u]$  such that  $u - \ell = \eta_\sigma$ , and have the **NB-overlap** property for any other  $[\ell, u]$ , i.e., the third group in Table 4.

The **SteadySequence** represents a special case, because when  $u - \ell = \eta_\sigma$ ,  $\sigma$  has neither property for  $[\ell, u]$ , and when  $u - \ell > \eta_\sigma$ ,  $\sigma$  has the **NB-no-overlap** property for  $[\ell, u]$ .

name $\sigma$	$\omega_\sigma$	$\eta_\sigma$	$\langle e_\sigma, c_\sigma \rangle$	$\phi_\sigma^{(n)}$	$\Theta_\sigma$	$o_\sigma^{(\ell, u)}$	$\delta_\sigma^{(\ell, u)}$
Bump	5	2	undefined	$\begin{cases} 2 & \text{if } n = 6 \\ \text{undefined} & \text{otherwise} \end{cases}$	$\{ ' > > < < > > ' \}$	3	0
Dec	1	1	undefined	$\begin{cases} 1 & \text{if } n = 2 \\ \text{undefined} & \text{otherwise} \end{cases}$	$\{ ' > ' \}$	$\begin{cases} 0 & \text{if } u - \ell \leq 1 \\ 1 & \text{otherwise} \end{cases}$	$\begin{cases} 0 & \text{if } u - \ell \leq 1 \\ -1 & \text{otherwise} \end{cases}$
DecSeq	1	1	$\langle 0, 1 \rangle$	$\begin{cases} 1 & \text{if } n = 2 \\ 2 & \text{if } n > 2 \end{cases}$	$\{ ' > ' \}$	0	0
DecTer	3	2	$\langle 0, 0 \rangle$	2	$\{ ' > > < < > > ' \}$	$\begin{cases} 0 & \text{if } u - \ell \leq 2 \\ 2 & \text{otherwise} \end{cases}$	$\begin{cases} 0 & \text{if } u - \ell \leq 2 \\ -1 & \text{otherwise} \end{cases}$
Dip	5	2	undefined	$\begin{cases} 2 & \text{if } n = 6 \\ \text{undefined} & \text{otherwise} \end{cases}$	$\{ ' < < > > < < > ' \}$	3	0
Gorge	2	1	$\langle 0, 1 \rangle$	$\begin{cases} 1 & \text{if } n = 3 \\ 2 & \text{if } n > 3 \end{cases}$	$\{ ' > < ' \}$	1	0
Inc	1	1	undefined	$\begin{cases} 1 & \text{if } n = 2 \\ \text{undefined} & \text{otherwise} \end{cases}$	$\{ ' < ' \}$	$\begin{cases} 0 & \text{if } u - \ell \leq 1 \\ 1 & \text{otherwise} \end{cases}$	$\begin{cases} 0 & \text{if } u - \ell \leq 1 \\ 1 & \text{otherwise} \end{cases}$
IncSeq	1	1	$\langle 0, 1 \rangle$	$\begin{cases} 1 & \text{if } n = 2 \\ 2 & \text{if } n > 2 \end{cases}$	$\{ ' < ' \}$	0	0
IncTer	3	2	$\langle 0, 0 \rangle$	2	$\{ ' < = < ' \}$	$\begin{cases} 0 & \text{if } u - \ell \leq 2 \\ 2 & \text{otherwise} \end{cases}$	$\begin{cases} 0 & \text{if } u - \ell \leq 2 \\ 1 & \text{otherwise} \end{cases}$
Inflexion	2	1	$\langle 0, 0 \rangle$	1	$\{ ' < > ' , ' > < ' \}$	2	0
Peak	2	1	$\langle 0, 0 \rangle$	1	$\{ ' < > ' \}$	1	0
Plain	2	1	$\langle 0, 0 \rangle$	1	$\{ ' > < ' \}$	1	0
Plateau	2	1	$\langle 0, 0 \rangle$	1	$\{ ' < > ' \}$	1	0
PropPlain	3	1	$\langle 0, 0 \rangle$	1	$\{ ' > = < ' \}$	1	0
PropPlateau	3	1	$\langle 0, 0 \rangle$	1	$\{ ' < = > ' \}$	1	0
Steady	1	0	undefined	$\begin{cases} 0 & \text{if } n = 2 \\ \text{undefined} & \text{otherwise} \end{cases}$	$\{ ' = ' \}$	1	0
SteadySeq	1	0	$\langle 0, 0 \rangle$	0	$\{ ' = ' \}$	0	0
SDecSeq	1	1	$\langle 1, 0 \rangle$	$n - 1$	$\{ ' > ' \}$	0	0
SIncSeq	1	1	$\langle 1, 0 \rangle$	$n - 1$	$\{ ' < ' \}$	0	0
Summit	2	1	$\langle 0, 1 \rangle$	$\begin{cases} 1 & \text{if } n = 3 \\ 2 & \text{if } n > 3 \end{cases}$	$\{ ' < > ' \}$	1	0
Valley	2	1	$\langle 0, 0 \rangle$	1	$\{ ' > < ' \}$	1	0
Zigzag	3	1	$\langle 0, 0 \rangle$	1	$\{ ' < > < < ' , ' > < > > ' \}$	$\begin{cases} 0 & \text{if } u - \ell \leq 1 \\ 1 & \text{otherwise} \end{cases}$	0

Table 3: Regular expression names  $\sigma$  and corresponding *width*  $\omega_\sigma$ , *height*  $\eta_\sigma$ , *range*  $\phi_\sigma^{(n)}$  (for a non-fixed-length regular expression  $\sigma$  and for any  $n > \omega_\sigma + 1$ ,  $\phi_\sigma^{(n)} = e_\sigma \cdot (n - 1 - \eta_\sigma) + c_\sigma + \eta_\sigma$ ), *inducing words*  $\Theta_\sigma$ , *overlap*  $o_\sigma^{(\ell, u)}$ , and *smallest variation of maxima*  $\delta_\sigma^{(\ell, u)}$ , where Bump, Dec, DecSeq, DecTer, Dip, Inc, IncSeq, IncTer, PropPlain, PropPlateau, SteadySeq, SDecSeq, SIncSeq are respectively shortcuts for BumpOnDecreasingSequence, Decreasing, DecreasingSequence, DecreasingTerrace, DipOnIncreasingSequence, Increasing, IncreasingSequence, IncreasingTerrace, ProperPlain, ProperPlateau, SteadySequence, StrictlyDecreasingSequence, StrictlyIncreasingSequence.

	Regular Expressions	Set of Properties
Overlapping Class	BumpOnDecreasingSequence	$\overline{\text{NB}}$ -simple
	DipOnIncreasingSequence	$\overline{\text{NB}}$ -overlap
	Gorge	$\overline{\text{WIDTH}}$ -max
	Inflexion	$\overline{\text{WIDTH}}$ -sum
	Peak	$\overline{\text{WIDTH}}$ -occurrence
	Plain	Condition of Theorem 1 in [3]
	Plateau	
	ProperPlain	
	ProperPlateau	
	Summit	
	Valley	
Non-Overlapping Class	DecreasingSequence	$\overline{\text{NB}}$ -simple
	IncreasingSequence	$\overline{\text{NB}}$ -no-overlap
	StrictlyDecreasingSequence	$\overline{\text{WIDTH}}$ -max
	StrictlyIncreasingSequence	$\overline{\text{WIDTH}}$ -sum
		$\overline{\text{WIDTH}}$ -occurrence
		Condition of Theorem 1 in [3]
Overlapping Class	Decreasing	$\overline{\text{NB}}$ -simple
	Increasing	$\overline{\text{NB}}$ -no-overlap when $u - \ell = \eta_\sigma$
	DecreasingTerrace	$\overline{\text{NB}}$ -overlap when $u - \ell \geq \eta_\sigma + 1$
	IncreasingTerrace	$\overline{\text{WIDTH}}$ -max
	Zigzag	$\overline{\text{WIDTH}}$ -sum
		$\overline{\text{WIDTH}}$ -occurrence
		Condition of Theorem 1 in [3]
Special Case	Steady	$\overline{\text{NB}}$ -simple when $u - \ell > \eta_\sigma$
		$\overline{\text{NB}}$ -no-overlap when $u - \ell = \eta_\sigma$
		$\overline{\text{NB}}$ -overlap
		$\overline{\text{WIDTH}}$ -max
		$\overline{\text{WIDTH}}$ -sum
		$\overline{\text{WIDTH}}$ -occurrence
		Condition of Theorem 1 in [3]
	SteadySequence	$\overline{\text{NB}}$ -simple when $u - \ell > \eta_\sigma$
		$\overline{\text{NB}}$ -no-overlap when $u - \ell = \eta_\sigma$
		$\overline{\text{NB}}$ -overlap
		$\overline{\text{WIDTH}}$ -max
		$\overline{\text{WIDTH}}$ -sum
		$\overline{\text{WIDTH}}$ -occurrence when $u - \ell > \eta_\sigma$
		Condition of Theorem 1 in [3]

Table 4: Classification of regular expressions: regular expression names  $\sigma$ , their properties and conditions on domain when they hold.

## 7 Evaluation

We evaluate the impact of the methods introduced in the previous sections on both execution time and the number of backtracks (failures) for all the 200 time-series constraints for which the glue constraint exists. Given the time-series constraints  $\gamma(\langle X_1, X_2, \dots, X_n \rangle, N)$ ,  $\gamma(\langle X_1, X_2, \dots, X_i \rangle, N_p)$  and  $\gamma(\langle X_n, X_{n-1}, \dots, X_i \rangle, N_s)$  with  $i \in [1, n]$ , the *glue constraint* links the overall result  $N$  with the two results  $N_p$  and  $N_s$  [7].

In our first experiment, we consider a single  $G\_F\_ \sigma(\langle X_1, X_2, \dots, X_n \rangle, N)$  constraint for which we first enumerate  $N$  and then either find solutions by assigning the  $X_i$  or prove infeasibility of the chosen  $N$ . For each constraint, we compare four variants of *Automaton*, which just states the constraint, using the automaton of [5]: *Glue* adds to *Automaton* the glue constraints [3], [7] for all prefixes and corresponding reversed suffixes; *Bounds* adds to *Automaton* the bound restrictions; *Bounds+Glue* uses both the glue constraints and the bounds; and *Combined* adds to *Bounds+Glue* the bounds for each prefix and corresponding reversed suffix.

In Figure 4, we show results for two problems that are small enough to perform all computations for *Automaton* and all variants within a reasonable time. In the first problem (first row of plots), we use time series of length 10 over the domain  $[1, 5]$ , and find, for each value of  $N$ , the first solution or prove infeasibility. This would be typical for satisfaction or optimisation problems, where one has to detect infeasibility quickly. Our static search routine enumerates the time-series variables  $X_i$  from left to right, starting with the smallest value in the domain. In the case of the initial domains being of the same size, this heuristic typically works best. In the second problem (second row of plots), we consider time series of length 8 over the domain  $[1, 5]$ , and find all solutions for each value of  $N$ . This allows us to verify that no solutions are incorrectly eliminated by any of the variants, and provides a worst-case scenario exploring the complete search tree. Results for the backtrack count are on the left, results for the execution time on the right. We use log scales on both axes, replacing a zero value by one in order to allow plotting. All experiments were run with SICStus Prolog 4.2.3 on a 2011 MacBook Pro 2.2 GHz quadcore Intel Core i7-950 machine with 6 MB cache and 16 GB memory using a single core.

We see that *Bounds* and *Glue* on their own bring good reductions of the search space, but their combinations *Bounds+Glue* and *Combined* in many cases reduce the number of backtracks by more than three orders of magnitude. Indeed, for many constraints, finding the first solution requires no backtracks. On the other hand, there are a few constraints for which the number of backtracks is not reduced significantly. These are constraints for which values of  $N$  in the middle of the domain are infeasible, but this is not detected by any of our variants.

The time for finding the first solution or proving infeasibility is also significantly reduced by the combinations *Bounds+Glue* and *Combined*, even though the glue constraints require two time-series constraints. When finding all solutions, this overhead shows in the total time taken for the three variants using the glue constraints. The bounds on their own reduce the time for many constraints, but rarely by more than a factor of ten.

In our second experiment, shown in Figure 5, we want to see whether the *Combined* variant is scalable. For this, we increase the length of the time series from 10 to 120 over the domain  $[1, 5]$ . We enumerate all possible values of  $N$  and find a first solution or prove infeasibility. For each time-series constraint and value of  $N$ , we impose a timeout of 20 seconds, and we do not consider the constraint if there is a timeout on some value of  $N$ . We plot the percentage of all constraints for which the average runtime is less than or equal to the value on the horizontal axis. For small time values, there are some quantisation effects due to the SICStus time resolution of 10 milliseconds.

For length 10, we find solutions for all values of  $N$  within the timeout, and our plots for *Automaton* (dashed) and *Combined* (solid) reach 100%, but the average time of *Combined* is much smaller. For *Automaton*, the percentage of constraints that are solved within the timeout drops to less than 20% for



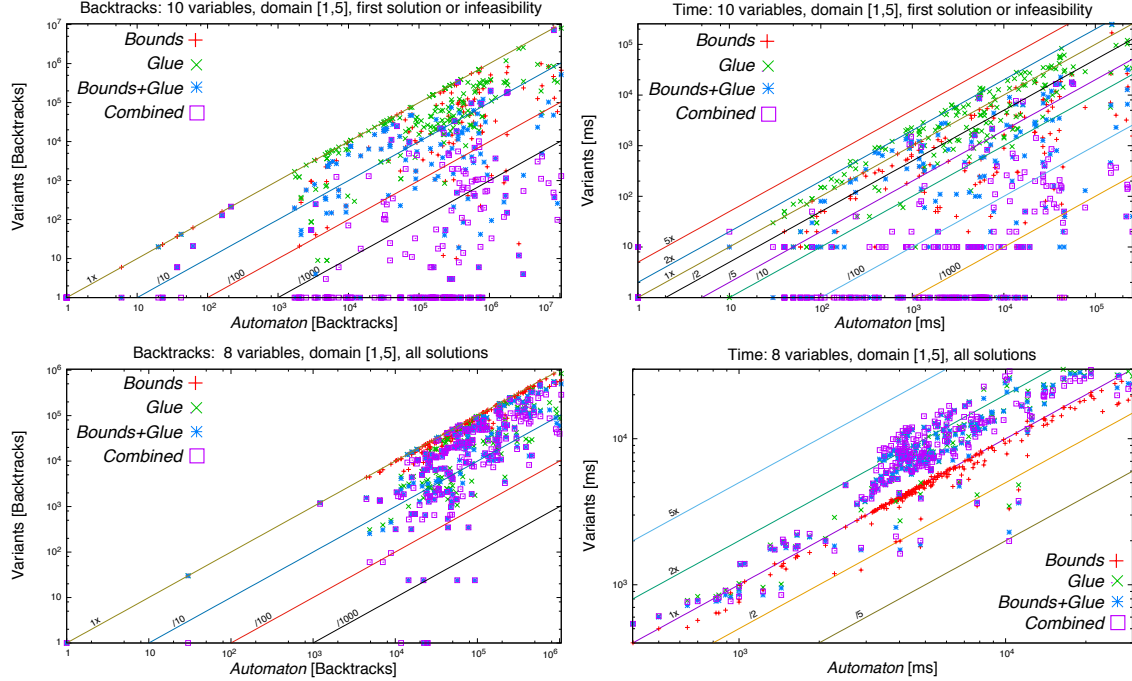


Fig. 4: Comparing backtrack count and runtime for *Automaton* and its variants for the first solution (length 10) and all solutions (length 8).

length 20, and less than 10% for length 40. For *Combined*, we solve over 75% of all constraints within the time limit, even for lengths 100 and 120.

The constraints that are not solved by *Combined* use the feature **surf** or the aggregator **Sum**. The worst performance is observed for constraints combining both **surf** and **Sum**. This is not surprising, as we know that achieving domain consistency for many of those constraints is NP-hard (encoding of *subset-sum*).

As a final experiment, we look at the search trees generated by four solution variants for a single constraint `MAX_SURF_INCREASING_TERRACE`. We only display some of the values for the parameter  $N$ , to make the trees more legible. Figure 6 shows the search tree produced with the help of CP-Viz [16]. Each tree shows the branches explored to find a first solution or proving infeasibility for each parameter value, with the initial choice of the value  $N$  at the top, and then the assignment of ten variables with a standard left-to-right labeling. Failed subtrees are abstracted as red triangles containing two numbers, the one above is the number of internal nodes in the tree, the one below the number of failed leaf nodes. Success nodes are colored in green, while failure nodes are colored red. Internal nodes are labeled by the variable name currently being assigned, and a superscript indicating the number of values in the domain of that variable. Edges indicate choices that are explored, the number indicates the value that is assigned to the selected variable, while a yellow edge color indicates that the value had been fixed by propagation.

In all trees, a first solution for parameter value 4, the smallest feasible value, is found without backtracking. The solution chooses value 1 for  $X_1$  to  $X_7$ , then value 2 for  $X_8$  and  $X_9$ , and finally value 3 for variable  $X_{10}$ . On the other hand, in the initial automaton, a very large failed subtree is shown for the left-most parameter value 3, and a much smaller failed tree for the right-most value 33. Both of those values are infeasible, and are removed by the bounds for this constraint. The *Bounds* version therefore avoids these

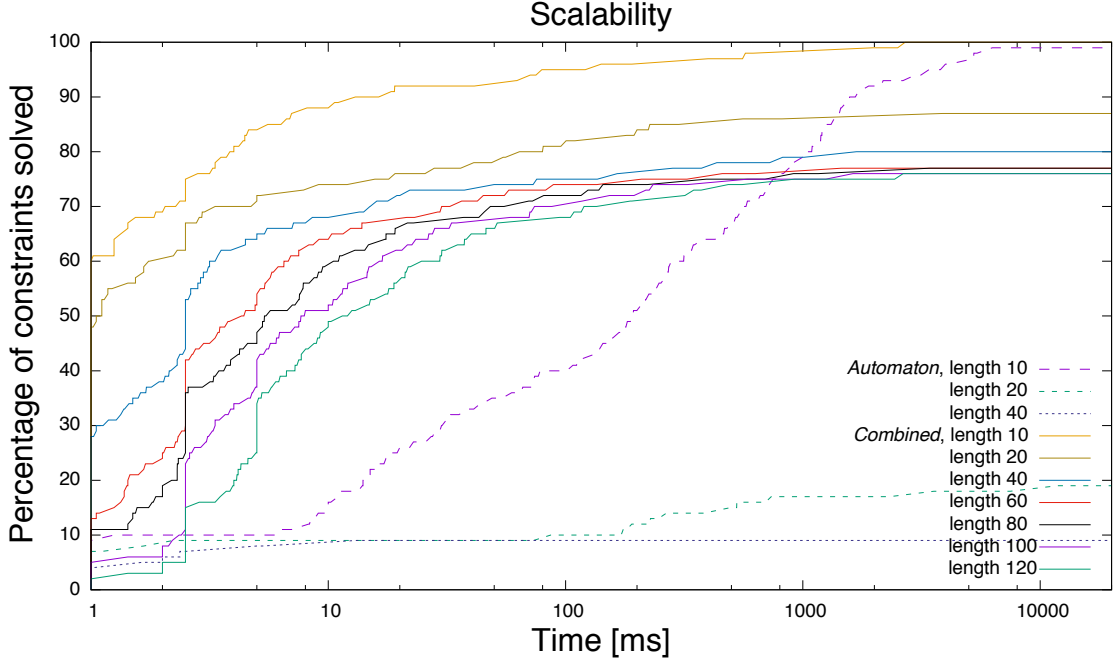


Fig. 5: Scalability results comparing time for *Automaton* and *Combined* on problems of increasing length.

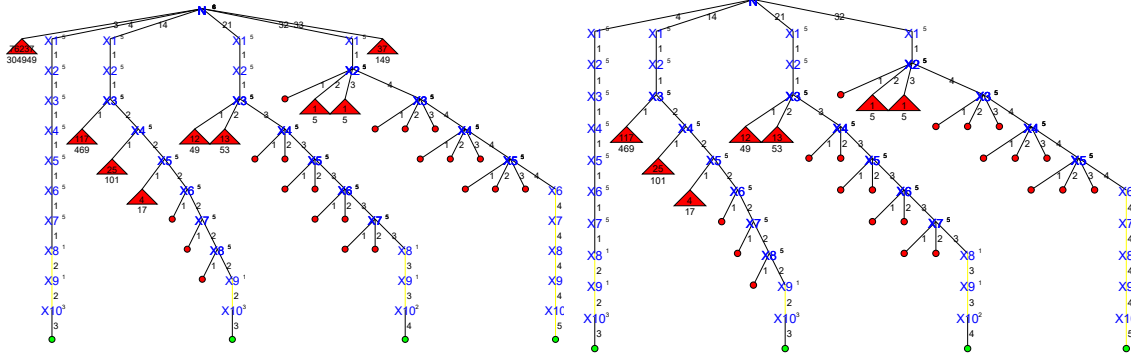
failed sub-trees, but there are no changes for the other, feasible values. When we consider the *Bounds+Glue* version, the search for feasible solutions is reduced, with a further reduction for the *Combined* variant. But we still need search to find the initial solution for some of the parameter values. This occurs since the bounds and the glue matrix reasoning only consider lower and upper bounds, and we don't detect holes in the domain of variable  $N$ . To get the best use of the generated bounds, we have to use the incremental combination of *Bounds* with the *Glue* constraint, as the bounds are then applied for each suffix of unassigned variables to maximise the information extracted.

## 8 Conclusion

Within the context of quantitative extensions of regular languages (QRE) we introduce the concept of regular expression characteristics as a way to unify combinatorial aspects of quantitative extensions of regular languages. We illustrate that approach for time-series constraints where, introducing six regular expression characteristics, allows coming up with generic bounds for families of time-series constraints. We believe the introduction of regular expression characteristics is important for the area of QRE.

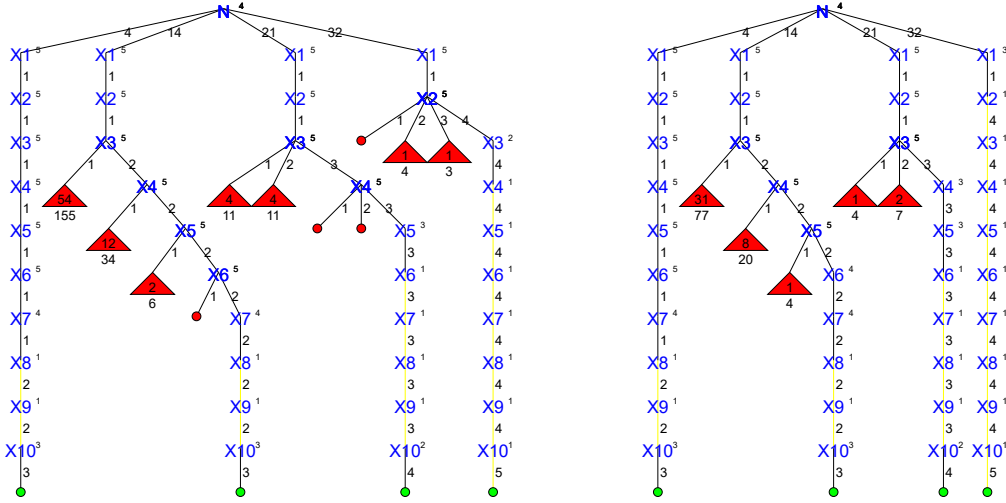
*Acknowledgements.*

We thank Pierre Flener for his feedback on the notation system for regular expression characteristics in Section 3.1.



(a) Automaton

(b) Bounds



(c) Bounds+Glue

(d) Combined

Fig. 6: Comparing parts of the search tree for `max_surf_increasing_terrace`, finding the first solution or proving infeasibility for the manually selected values 3, 4, 14, 21, 32, and 33 of variable  $N$  and 10 variables  $X_1, X_2, \dots, X_{10}$ , each with domain  $[1, 5]$ .

## References

1. Alur, R., D’Antoni, L., Deshmukh, J.V., Raghothaman, M., Yuan, Y.: Regular functions and cost register automata. In: 28th Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2013, New Orleans, LA, USA, June 25-28, 2013. pp. 13–22. IEEE Computer Society (2013)
2. Alur, R., Fisman, D., Raghothaman, M.: Regular programming for quantitative properties of data streams. In: Thiemann, P. (ed.) Programming Languages and Systems - 25th European Symposium on Programming, ESOP 2016, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2016, Eindhoven, The Netherlands, April 2-8, 2016, Proceedings. Lecture Notes in Computer Science, vol. 9632, pp. 15–40. Springer (2016)
3. Arafailova, E., Beldiceanu, N., Carlsson, M., Flener, P., Francisco Rodríguez, M.A., Pearson, J., Simonis, H.: Systematic derivation of bounds and glue constraints for time-series constraints. In: Rueher, M. (ed.) CP 2016. LNCS, vol. 9892, pp. 13–29. Springer (2016)
4. Arafailova, E., Beldiceanu, N., Douence, R., Carlsson, M., Flener, P., Rodríguez, M.A.F., Pearson, J., Simonis, H.: Global constraint catalog, volume ii, time-series constraints. CoRR abs/1609.08925 (2016), <http://arxiv.org/abs/1609.08925>
5. Arafailova, E., Beldiceanu, N., Douence, R., Flener, P., Francisco Rodríguez, M.A., Pearson, J., Simonis, H.: Time-series constraints: Improvements and application in CP and MIP contexts. In: Quimper, C.G. (ed.) CP-AI-OR 2016. LNCS, vol. 9676, pp. 18–34. Springer (2016)
6. Beldiceanu, N., Carlsson, M., Douence, R., Simonis, H.: Using finite transducers for describing and synthesising structural time-series constraints. Constraints 21(1), 22–40 (January 2016), journal fast track of CP 2015: summary on p. 723 of LNCS 9255, Springer, 2015
7. Beldiceanu, N., Carlsson, M., Flener, P., Francisco Rodríguez, M.A., Pearson, J.: Linking prefixes and suffixes for constraints encoded using automata with accumulators. In: O’Sullivan, B. (ed.) CP 2014. LNCS, vol. 8656, pp. 142–157. Springer (2014)
8. Beldiceanu, N., Carlsson, M., Petit, T.: Deriving filtering algorithms from constraint checkers. In: Wallace, M. (ed.) CP 2004. LNCS, vol. 3258, pp. 107–122. Springer (2004)
9. Beldiceanu, N., Ifrim, G., Lenoir, A., Simonis, H.: Describing and generating solutions for the EDF unit commitment problem with the ModelSeeker. In: Schulte, C. (ed.) CP 2013. LNCS, vol. 8124, pp. 733–748. Springer (2013)
10. Colcombet, T., Daviaud, L.: Approximate comparison of functions computed by distance automata. Theory Comput. Syst. 58(4), 579–613 (2016)
11. Crochemore, M., Hancart, C., Lecroq, T.: Algorithms on Strings. Cambridge University Press (2007)
12. Demasse, S., Pesant, G., Rousseau, L.M.: A **Cost-Regular** based hybrid column generation approach. Constraints 11(4), 315–333 (2006)
13. Hopcroft, J.E., Motwani, R., Ullman, J.D.: Introduction to Automata Theory, Languages, and Computation. Addison-Wesley, 3rd edn. (2007)
14. Pesant, G.: A regular language membership constraint for finite sequences of variables. In: Wallace, M. (ed.) CP 2004. LNCS, vol. 3258, pp. 482–495. Springer (2004)
15. Schützenberger, M.P.: On the definition of a family of automata. Information and Control 4, 245–270 (1961)
16. Simonis, H., Davern, P., Feldman, J., Mehta, D., Quesada, L., Carlsson, M.: A generic visualization platform for CP. In: Cohen, D. (ed.) Principles and Practice of Constraint Programming - CP 2010 - 16th International Conference, CP 2010, St. Andrews, Scotland, UK, September 6-10, 2010. Proceedings. Lecture Notes in Computer Science, vol. 6308, pp. 460–474. Springer (2010), [http://dx.doi.org/10.1007/978-3-642-15396-9\\_37](http://dx.doi.org/10.1007/978-3-642-15396-9_37)
17. Yasunori, I., Takuji, M., Shougo, S., Kenji, H., Toru, F.: A tractable subclass of dtlds for xpath satisfiability with sibling axes. In: Philippa, G., Floris, G. (eds.) Database Programming Languages: 12th International Symposium, DBPL 2009s. LNCS, vol. 5708, pp. 68–83. Springer (2009)

## A Appendix: Tables of Regular Expressions Characteristics

Table 5 gives the *width* characteristics for each regular expression in Table 1. Within Table 5, smallest words are obtained by (1) first discarding from a regular expression all sub-expressions containing the empty word, and then by (2) keeping within each disjunction the smallest length words. For instance, within the **Zigzag** regular expression ' $(\langle \rangle)^* \langle \rangle \langle \rangle (\rangle \mid \varepsilon) \mid (\rangle \langle \rangle)^* \rangle \langle \rangle (\langle \mid \varepsilon)'$ ', we remove the sub-expressions ' $(\langle \rangle)^*$ ', ' $(\rangle \mid \varepsilon)$ ', ' $(\rangle \langle \rangle)^*$ ', ' $(\langle \mid \varepsilon)$ ' and obtain the disjunction ' $\langle \rangle \langle \rangle \mid \rangle \langle \rangle$ ' containing two words of length 3.

name $\sigma$	regular expression	$\omega_\sigma$
BumpOnDecreasingSequence	' $\rangle \rangle \langle \rangle \rangle$ '	5
Decreasing	' $\rangle$ '	1
DecreasingSequence	' $(\rangle (\rangle \mid =)^* \rangle$ '	1
DecreasingTerrace	' $\rangle = =^* \rangle$ '	3
DipOnIncreasingSequence	' $\langle \langle \rangle \langle \rangle \langle \rangle$ '	5
Gorge	' $(\rangle \mid \rangle (\rangle \mid =)^* \rangle)(\langle \mid \langle (\langle \mid =)^* \langle)$ '	2
Increasing	' $\langle$ '	1
IncreasingSequence	' $(\langle (\langle \mid =)^* \langle$ '	1
IncreasingTerrace	' $\langle = =^* \langle$ '	3
Inflexion	' $\langle (\langle \mid =)^* \rangle \mid \rangle (\rangle \mid =)^* \langle$ '	2
Peak	' $\langle (\langle \mid =)^* (\rangle \mid =)^* \rangle$ '	2
Plain	' $\rangle =^* \langle$ '	2
Plateau	' $\langle =^* \rangle$ '	2
ProperPlain	' $\rangle = =^* \langle$ '	3
ProperPlateau	' $\langle = =^* \rangle$ '	3
Steady	' $=$ '	1
SteadySequence	' $= =^*$ '	1
StrictlyDecreasingSequence	' $\rangle \rangle^*$ '	1
StrictlyIncreasingSequence	' $\langle \langle^*$ '	1
Summit	' $(\langle \mid \langle (\langle \mid =)^* \langle)(\rangle \mid \rangle (\rangle \mid =)^* \rangle)$ '	2
Valley	' $\rangle (\rangle \mid =)^* (\langle \mid =)^* \langle$ '	2
Zigzag	' $(\langle \rangle)^* \langle \rangle \langle \rangle (\rangle \mid \varepsilon) \mid (\rangle \langle \rangle)^* \rangle \langle \rangle (\langle \mid \varepsilon)'$ '	3

Table 5: Regular expression names  $\sigma$  and corresponding *width* (see Definition 4); within each regular expression subparts corresponding to a smallest length word are highlighted in yellow.

Table 6 gives the *height* characteristics for each regular expression in Table 1. Within Table 6, the ‘illustration’ column provides for each regular expression  $\sigma$  a word achieving the smallest height among all words of  $\mathcal{L}_\sigma$ . For a regular expression  $\sigma$ , a word  $w$  achieving the smallest height is a word of  $\mathcal{L}_\sigma$  that minimises the number of occurrences of ‘>’ (respectively ‘<’) over all maximal occurrences of ‘> (= | >)\*’ (respectively ‘< (= | <)\*’) in  $w$ . We illustrate this for two regular expressions.

- For the fixed length regular expression  $\sigma = \text{BumpOnDecreasingSequence}$ ,  $\mathcal{L}_\sigma$  contains a single word  $w = '>><>>'$ . Since  $w$  is the concatenation of three proper factors ‘>>’, ‘<’ and ‘>>’ of respective length 2, 1 and 2 we obtain a height of 2.
- For the non-fixed length regular expression  $\sigma = \text{DecreasingTerrace}$ , the word ‘>=>’  $\in \mathcal{L}_\sigma$  has a height of 2. No word in  $\mathcal{L}_\sigma$  can have a smaller height, since any word  $w$  in the language of ‘>=>’ contains two occurrences of ‘>’, one at both extremities of  $w$ , separated by a single stretch of ‘=’.

name $\sigma$	illustration	$\eta_\sigma$
BumpOnDecreasingSequence		2
Decreasing		1
DecreasingSequence		1
DecreasingTerrace		2
DipOnIncreasingSequence		2
Gorge		1
Increasing		1
IncreasingSequence		1
IncreasingTerrace		2
Inflexion		1
Peak		1
Plain		1
Plateau		1
ProperPlain		1
ProperPlateau		1
Steady		0
SteadySequence		0
StrictlyDecreasingSequence		1
StrictlyIncreasingSequence		1
Summit		1
Valley		1
Zigzag		1

Table 6: Regular expression names  $\sigma$  and corresponding *height* shown as thick orange vertical line segments (see Definition 7).

Table 7 gives the *range* characteristics for each regular expression in Table 1. Within Table 7, the ‘illustration’ column provides for each regular expression  $\sigma$  a time series whose signature is a word of the smallest height among all words of the same length  $n - 1$  in  $\mathcal{L}_\sigma$ , i.e. the range of  $\sigma$  wrt  $\langle n \rangle$ . We distinguish three cases:

- For a fixed length regular expression  $\sigma$  (e.g. **BumpOnDecreasingSequence**), the range  $\phi_\sigma^{(n)}$  is only defined for one plus the length of the single word in  $\mathcal{L}_\sigma$ , and is equal to the height  $\eta_\sigma$  of  $\sigma$ .
- For a non-fixed length regular expression  $\sigma$ , if we can find a word of length  $n - 1$  in  $\mathcal{L}_\sigma$  which has the same height as the height  $\eta_\sigma$ , we cannot have a smaller height by definition. This is the case for many of our non-fixed length regular expressions, for example **Peak**, **Inflexion** or **Zigzag**.
- For some non-fixed length regular expressions  $\sigma$  like **DecreasingSequence**, **IncreasingSequence**, **Gorge** or **Summit**, only the corresponding shortest word has a height of  $\eta_\sigma$ . Then, any longer word in  $\mathcal{L}_\sigma$ , has a height of at least  $\eta_\sigma + 1$ .
- For  $\sigma = \mathbf{StrictlyDecreasingSequence}$ ,  $\mathcal{L}_\sigma$  contains a single word of length  $n - 1$ , namely a stretch of  $n - 1$  consecutive ‘>’. Hence, the range of  $\sigma$  wrt  $\langle n \rangle$  is reached for this word and equals  $n - 1$ . The same reasoning applies for **StrictlyIncreasingSequence**.

name $\sigma$	$\langle e_\sigma, c_\sigma \rangle$	illustration	$\phi_\sigma^{(n)}$
BumpOnDecreasingSequence	undefined		$\begin{cases} 2 & \text{if } n = 6 \\ \text{undefined} & \text{otherwise} \end{cases}$
Decreasing	undefined		$\begin{cases} 1 & \text{if } n = 2 \\ \text{undefined} & \text{otherwise} \end{cases}$
DecreasingSequence	$\langle 0, 1 \rangle$		$\begin{cases} 1 & \text{if } n = 2 \\ 2 & \text{if } n > 2 \end{cases}$
DecreasingTerrace	$\langle 0, 0 \rangle$		2
DipOnIncreasingSequence	undefined		$\begin{cases} 2 & \text{if } n = 6 \\ \text{undefined} & \text{otherwise} \end{cases}$
Gorge	$\langle 0, 1 \rangle$		$\begin{cases} 1 & \text{if } n = 3 \\ 2 & \text{if } n > 3 \end{cases}$
Increasing	undefined		$\begin{cases} 1 & \text{if } n = 2 \\ \text{undefined} & \text{otherwise} \end{cases}$
IncreasingSequence	$\langle 0, 1 \rangle$		$\begin{cases} 1 & \text{if } n = 2 \\ 2 & \text{if } n > 2 \end{cases}$
IncreasingTerrace	$\langle 0, 0 \rangle$		2
Inflexion	$\langle 0, 0 \rangle$		1
Peak	$\langle 0, 0 \rangle$		1
Plain	$\langle 0, 0 \rangle$		1
Plateau	$\langle 0, 0 \rangle$		1
ProperPlain	$\langle 0, 0 \rangle$		1
ProperPlateau	$\langle 0, 0 \rangle$		1
Steady	undefined		$\begin{cases} 0 & \text{if } n = 2 \\ \text{undefined} & \text{otherwise} \end{cases}$
SteadySequence	$\langle 0, 0 \rangle$		0
StrictlyDecreasingSequence	$\langle 1, 0 \rangle$		$n - 1$
StrictlyIncreasingSequence	$\langle 1, 0 \rangle$		$n - 1$
Summit	$\langle 0, 1 \rangle$		$\begin{cases} 1 & \text{if } n = 3 \\ 2 & \text{if } n > 3 \end{cases}$
Valley	$\langle 0, 0 \rangle$		1
Zigzag	$\langle 0, 0 \rangle$		1

Table 7: Regular expression names  $\sigma$  and corresponding *range* shown as thick orange vertical line segments (see Definition 8); for a non-fixed-length regular expression  $\sigma$  and for any  $n > \omega_\sigma + 1$ ,  $\phi_\sigma^{(n)} = e_\sigma \cdot (n - 1 - \eta_\sigma) + c_\sigma + \eta_\sigma$ .



Table 8 gives the *inducing words* characteristics for each regular expression in Table 1. Within Table 8, the inducing words characteristics is derived from the corresponding regular expression by removing all sub-expressions containing the empty word and by keeping the rest, i.e. the part highlighted in yellow.

name $\sigma$	regular expression	$\Theta_\sigma$
BumpOnDecreasingSequence	'>><>>'	{ '>><>>' }
Decreasing	'>'	{ '>' }
DecreasingSequence	'(> (>   =)*)*>'	{ '>' }
DecreasingTerrace	'>= * >'	{ '>=>' }
DipOnIncreasingSequence	'<<><<'	{ '<<><<' }
Gorge	'(> (>   =)*)*>< ((<   =)*)*<'	{ '><' }
Increasing	'<'	{ '<' }
IncreasingSequence	'(< (<   =)*)*<'	{ '<' }
IncreasingTerrace	'<= * <'	{ '<=<' }
Inflexion	'< (<   =)*>   > (>   =)*<'	{ '<>', '><' }
Peak	'< (<   =)*(>   =)*>'	{ '<>' }
Plain	'>= * <'	{ '><' }
Plateau	'<= * >'	{ '<>' }
ProperPlain	'>= * <'	{ '>=<' }
ProperPlateau	'<= * >'	{ '<=>' }
Steady	'='	{ '=' }
SteadySequence	'= *'	{ '=' }
StrictlyDecreasingSequence	'> > *'	{ '>' }
StrictlyIncreasingSequence	'< < *'	{ '<' }
Summit	'(< (<   =)*)*<> ((>   =)*)*>'	{ '<>' }
Valley	'> (>   =)*(<   =)*<'	{ '><' }
Zigzag	'(<>)*<>< (>   $\varepsilon$ )   (><)*><> (<   $\varepsilon$ )'	{ '<><', '><>' }

Table 8: Regular expression names  $\sigma$  and corresponding *inducing words* (see Definition 10).

Table 9 gives the *overlap* characteristics for each regular expression in Table 1. Within Table 9 we distinguish the following cases for computing the overlap characteristics:

- Consider a fixed length regular expression  $\sigma$  whose regular language contains a single word  $w$ . Then, we compute the length of the maximum overlap  $o$  between  $w$  and itself for which  $o < |w|$ .
  - \* If such overlap exists the corresponding overlap characteristics  $o_{\sigma}^{(\ell, u)}$  is equal to  $o + 1$ , e.g. for  $\sigma = \text{BumpOnDecreasingSequence}$  the maximum overlap of ' $>><>>$ ' with itself is 2 assuming the two words do not completely overlap, leading to  $o_{\sigma}^{(\ell, u)} = 3$ .
  - \* If such overlap does not exist, depending whether the difference  $u - \ell$  is big enough or not, we can concatenate  $w$  with itself or not, leading to an overlap of 1 (one time-series variable is shared) or to an overlap of 0. This is the case for **Decreasing** and **Increasing** where, depending whether the difference  $u - \ell$  is strictly greater than 1 or not, we get an overlap  $o_{\sigma}^{(\ell, u)}$  of 1 or 0.
- Consider a regular expression  $\sigma$  for which the set of superpositions of any pair of words of  $\mathcal{L}_{\sigma}$  is empty; then the corresponding overlap  $o_{\sigma}^{(\ell, u)}$  is equal to 0. This is the case for **DecreasingSequence**, **IncreasingSequence**, **SteadySequence**, **StrictlyDecreasingSequence**, and **StrictlyIncreasingSequence**, because Condition (1) of Definition 11 is always violated.
- Given a regular expression  $\sigma$  for which (1) the set of superpositions of any pair of words of  $\mathcal{L}_{\sigma}$  is limited to the concatenation of the pair of corresponding words, and (2) any pair of word of  $\mathcal{L}_{\sigma}$  starts with a ' $<$ ' and ends up with a ' $>$ ' (or conversely starts with a ' $>$ ' and ends up with a ' $<$ '), then we can concatenate them so that they share one time-series variable regardless the value of  $u - \ell$ ; we get an overlap  $o_{\sigma}^{(\ell, u)}$  of 1. This is the case for **Gorge**, **Peak**, **Plain**, **Plateau**, **ProperPlain**, **ProperPlateau**, **Summit**, and **Valley**.
- Consider the  $\sigma = \text{DecreasingTerrace}$  regular expression. For any two words  $v = '>=^{(k)}>'$  and  $w = '>=^{(l)}>'$  in  $\mathcal{L}_{\sigma}$  with  $k, l$  being positive integers, the set of their superpositions wrt  $\langle \ell, u \rangle$  contains at most two words, namely  $z_1 = '>=^{(k)}>=^{(l)}>'$  and  $z_2 = '>=^{(k)}>=^{(l)}>'$ . The value of overlap achieved for  $z_1$  and for  $z_2$  is 2 and 1, respectively.
  - \* When  $u - \ell = \eta_{\sigma} = 2$ , neither  $z_1$  nor  $z_2$  can appear in the signature of a ground time series over  $[\ell, u]$ , thus the set of superpositions of  $\sigma$  wrt  $\langle \ell, u \rangle$  is empty, and  $o_{\sigma}^{(\ell, u)} = 0$ .
  - \* When  $u - \ell > \eta_{\sigma} = 2$ , the word  $z_1$  is always in the set of superpositions of  $\sigma$  wrt  $\langle \ell, u \rangle$ , and thus  $o_{\sigma}^{(\ell, u)} = 2$ .

The same reasoning applies for the **IncreasingTerrace** regular expression.

- Consider the  $\sigma = \text{Inflexion}$  regular expression. Any word in  $\mathcal{L}_{\sigma}$  belongs to the language of either  $\sigma_1 = '< (< | =)^* >'$  or  $\sigma_2 = '> (> | =)^* <'$ .
    - \* For any two words  $v, w \in \mathcal{L}_{\sigma_1}$  (respectively  $v, w \in \mathcal{L}_{\sigma_2}$ ), their overlap wrt  $\langle \ell, u \rangle$  is at most 1, since their only possible superposition is  $vw$ .
    - \* For any two words  $v \in \mathcal{L}_{\sigma_1}$  and  $w \in \mathcal{L}_{\sigma_2}$ , their overlap wrt  $\langle \ell, u \rangle$  is at most 2, since the maximum length of a suffix of  $v$  that is also a prefix of  $w$  is 1. Hence,  $o_{\sigma}^{(\ell, u)} \leq 2$ .  
The overlap of the words ' $><$ ' and ' $<>$ ' wrt  $\langle \ell, u \rangle$  such that  $u - \ell \geq \eta_{\sigma}$  is 2, which is maximum. Hence,  $o_{\sigma}^{(\ell, u)} = 2$ .
  - Consider the  $\sigma = \text{Zigzag}$  regular expression. Any word in  $\mathcal{L}_{\sigma}$  belongs to the regular language either of  $\sigma_1 = '(<>)^+ < (> | \varepsilon)'$  or of  $\sigma_2 = '(><)^+ > (< | \varepsilon)'$ .
    - \* For any two words  $v \in \mathcal{L}_{\sigma_1}$  and  $w \in \mathcal{L}_{\sigma_2}$ , their overlap wrt  $\langle \ell, u \rangle$  is always 0, since their set of superpositions wrt  $\langle \ell, u \rangle$  is empty, because Condition (1) of Definition 11 is violated.
    - \* For any two words  $v, w \in \mathcal{L}_{\sigma_1}$  (respectively  $v, w \in \mathcal{L}_{\sigma_2}$ ), their overlap wrt  $\langle \ell, u \rangle$  is at most 1, since their only possible superposition is  $vw$ . Note that the height of every word in  $\mathcal{L}_{\sigma}$  is  $\eta_{\sigma} = 1$ , then the height of  $vw$  is 2, since  $v$  last symbol coincides with the  $w$  first symbol, and it is not '='.
- Hence, when  $u - \ell = \eta_{\sigma} = 1$ , the overlap of  $\sigma$  wrt  $\langle \ell, u \rangle$  is 0, and when  $u - \ell \geq \eta_{\sigma} + 1$ , the overlap of  $\sigma$  wrt  $\langle \ell, u \rangle$  is 1.

name $\sigma$	illustration	$o_{\sigma}^{(\ell, u)}$
BumpOnDecreasingSequence		3
Decreasing		$\begin{cases} 0 & \text{if } u - \ell \leq 1 \\ 1 & \text{otherwise} \end{cases}$
DecreasingSequence		0
DecreasingTerrace		$\begin{cases} 0 & \text{if } u - \ell \leq 2 \\ 2 & \text{otherwise} \end{cases}$
DipOnIncreasingSequence		3
Gorge		1
Increasing		$\begin{cases} 0 & \text{if } u - \ell \leq 1 \\ 1 & \text{otherwise} \end{cases}$
IncreasingSequence		0
IncreasingTerrace		$\begin{cases} 0 & \text{if } u - \ell \leq 2 \\ 2 & \text{otherwise} \end{cases}$
Inflexion		2
Peak		1
Plain		1
Plateau		1
ProperPlain		1
ProperPlateau		1
Steady		1
SteadySequence		0
StrictlyDecreasingSequence		0
StrictlyIncreasingSequence		0
Summit		1
Valley		1
Zigzag		$\begin{cases} 0 & \text{if } u - \ell \leq 1 \\ 1 & \text{otherwise} \end{cases}$

Table 9: Regular expression names  $\sigma$  and corresponding *overlap* between two consecutive pattern occurrences ① and ② illustrated in red, i.e.,  $\bullet$  or  $\circ$  (see Definition 13).

Table 10 gives the *smallest variation of maxima* characteristics for each regular expression in Table 1. Within Table 10 we distinguish the following cases for computing the smallest variation of maxima characteristics:

- When  $o_{\sigma}^{\langle \ell, u \rangle}$  is 0, the quantity  $\delta_{\sigma}^{\langle \ell, u \rangle}$  is also 0 by definition. This is for example the case of **DecreasingSequence** and **Zigzag** when  $u - \ell = 1$ .
- When  $o_{\sigma}^{\langle \ell, u \rangle}$  is not 0, and we can give a pair of words  $v, w$  in  $\mathcal{L}_{\sigma}$  such that their set of superpositions wrt  $\langle \ell, u \rangle$  is not empty and  $\delta_{\sigma}^{\langle \ell, u \rangle}(v, w)$  is 0, the value of  $\delta_{\sigma}^{\langle \ell, u \rangle}$  is also 0. Note that by definition  $\delta_{\sigma}^{\langle \ell, u \rangle}$  has the minimum absolute value. Hence, if the value of zero is reached for at least one pair of words, then  $\delta_{\sigma}^{\langle \ell, u \rangle}$  is zero.
- When  $o_{\sigma}^{\langle \ell, u \rangle}$  is not 0, and when the regular language of  $\sigma$  contains a single word,  $\delta_{\sigma}^{\langle \ell, u \rangle}$  is reached for a superposition of this word with itself. See, for example **Decreasing**, when  $u - \ell \geq 2$ .
- Consider the  $\sigma = \mathbf{DecreasingTerrace}$  regular expression when  $u - \ell \geq 3$ , i.e. the overlap  $o_{\sigma}^{\langle \ell, u \rangle}$  is not 0. For any two words  $v = '>=^{(k)}>'$  and  $w = '>=^{(l)}>'$  in  $\mathcal{L}_{\sigma}$  with  $k, l$  being positive integers, the set of their superpositions wrt  $\langle \ell, u \rangle$  contains at most two words, namely  $'>=^{(k)}>=^{(l)}>'$  and  $'>=^{(k)}>=^{(l)}>'$ . Then, the value of  $\delta_{\sigma}^{\langle \ell, u \rangle}(v, w)$  equals  $-1$  and is reached for the superposition  $'>=^{(k)}>=^{(l)}>'$ . The same reasoning applies for **IncreasingTerrace**.

name $\sigma$	illustration	$\delta_{\sigma}^{(\ell, u)}$
BumpOnDecreasingSequence		0
Decreasing		$\begin{cases} 0 & \text{if } u - \ell \leq 1 \\ -1 & \text{otherwise} \end{cases}$
DecreasingSequence		0
DecreasingTerrace		$\begin{cases} 0 & \text{if } u - \ell \leq 2 \\ -1 & \text{otherwise} \end{cases}$
DipOnIncreasingSequence		0
Gorge		0
Increasing		$\begin{cases} 0 & \text{if } u - \ell \leq 1 \\ 1 & \text{otherwise} \end{cases}$
IncreasingSequence		0
IncreasingTerrace		$\begin{cases} 0 & \text{if } u - \ell \leq 2 \\ 1 & \text{otherwise} \end{cases}$
Inflexion		0
Peak		0
Plain		0
Plateau		0
ProperPlain		0
ProperPlateau		0
Steady		0
SteadySequence		0
StrictlyDecreasingSequence		0
StrictlyIncreasingSequence		0
Summit		0
Valley		0
Zigzag		0

Table 10: Regular expression names  $\sigma$  and corresponding *smallest variation of maxima* (see Definition 16); maxima of two consecutive pattern occurrences ① and ② are shown in red, i.e.,  $\bullet$  or  $\circ$ .