

Mining Maximal Sequential Patterns without Candidate Maintenance

Artificial Intelligence
Seminar
10 March 2016

Syméon Malengreau
Ayad Aldayeh

Table of contents

Authors presentation

Concepts

Sequence database

Sequential pattern

Closed Sequential pattern

Maximal Sequential pattern

Algorithm

PrefixSpan

MaxSP

Measures

Authors

Vincent S. Tseng



Philippe Fournier-Viger



Cheng-Wei Wu



Authors : Vincent S. Tseng

Ph.D. in Computer Science

Professor, Dept. Computer Science, National Chiao Tung University, Taiwan

H-index : 32

Citations : 3308



Authors : Philippe Fournier-Viger

Ph.D. in Computer Science

**Associate Professor, Harbin Institute of Technology, Shenzhen
Graduate School**

He created the SPMF library

H-index : 17

Citations : 712



Authors : Cheng-Wei Wu

H-index : 12

Citations : 518



About the article

We can trust the author, but

We found errors in some examples,

Some part of the articles were mathematically not clear enough, and some definition were simply wrong,

In the end, the article was a complete mess, we cannot understand who accepted to publish it !

After some search else where we made it possible to understand the content of the articles. But based only on it, it wouldn't have be possible.

It was published in a book : *Advanced Data Mining and Applications*

This particular article has (when the slides where made) 10 citations

Concepts

- **Sequence Database**
- **Sequential pattern**
- **Closed sequential pattern**
- **Maximal sequential pattern**

Sequence Database

A sequence database consist of :

A set of items

$\{1, 2, 3, 4, \dots, N\}$

Itemset (set of item, distinct and unordered)

$\{1, 2, 3, 5\}$ or $\{4, 5\}$ or $\{3, 7\}$ or ...

Sequence (set of itemsets)

$\langle \{1,2\}, \{3\}, \{5\} \rangle$ or $\langle \{4\}, \{6\} \rangle$ or ...

The sequence database is a set of sequences

What do theses concepts represents ?

Sequence Database : Illustration

Let's take as an example a book

Set of item \rightarrow The words

{He, nice, the, is, a, guy, sun, shine, ...}

Itemset \rightarrow A sentence (where words are distinct and unordered)

{He, a, nice, guy, is}

{The, sun, shine, in, the, sky}

Sequence \rightarrow A chapter of the book

Sequence Database \rightarrow The book

Sequential pattern

Synonyms are *sub-sequence* or *frequent sequence*

It is a sequence of item that appears a certain number of time, that number is the *minimum support threshold* (or *minsup*)

Sequence database

$\langle \{1,2\}, \{3\}, \{4\}, \{6\} \rangle$

$\langle \{2\}, \{5\}, \{6\} \rangle$

$\langle \{1,3\}, \{5\}, \{6\} \rangle$

With minsup = 2, some examples of sequential pattern

$\{5\}, \{6\}$

$\{1\}$

$\{3\}, \{6\}$

...

Closed sequential pattern

A closed sequential pattern is a sequential pattern not included in another closed pattern having the same frequency.

$\langle \{1\}, \{1\ 2\ 3\}, \{1\ 3\}, \{4\}, \{3\ 6\} \rangle$

$\langle \{1\ 4\}, \{3\}, \{2\ 3\}, \{1\ 5\} \rangle$

$\langle \{5\ 6\}, \{1\ 2\}, \{4\ 6\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\ 6\}, \{3\}, \{2\}, \{3\} \rangle$

With support 2 (or 2/4 entry \rightarrow 50 %), here are some closed sequential pattern

$\{1\}, \{3\}$ 100 % (4/4)

$\{1\}, \{3\}, \{2\}$ 75 % (3/4)

$\{5\}, \{1\}, \{3\}, \{2\}$ 50% (2/4)

$\{5\}$ 75 % (3/4)

And this one is NOT

$\{1\}$ 100 % (4/4)

Maximal sequential pattern

The same as the closed sequential pattern, but if one sequence is in another one, it is not maximal.

Interesting property :

You can derive every closed sequential patterns from the maximal sequential patterns



Question 1 : Closed and Maximal pattern

Considering the database

1. Which one of these is not a closed sequential pattern ? Why ?

→ $\langle \{b\}, \{f\} \rangle$

→ $\langle \{b\} \rangle$

→ $\langle \{a,b\} \rangle$

→ $\langle \{a\}, \{b\}, \{e\} \rangle$

2. Which one of these is a maximal sequential pattern ? Why ?

→ $\langle \{a\}, \{e\} \rangle$

→ $\langle \{b\}, \{b\} \rangle$

→ $\langle \{b\}, \{f\}, \{e\} \rangle$

→ $\langle \{a\}, \{f\} \rangle$

$\langle \{a,b\}, \{c\}, \{f,g\}, \{g\}, \{e\} \rangle$

$\langle \{a,d\}, \{c\}, \{b\}, \{a,b,e,f\} \rangle$

$\langle \{a\}, \{b\}, \{f,g\}, \{e\} \rangle$

$\langle \{b\}, \{f,g\} \rangle$

MaxSP Algorithm

Find the maximal sequential pattern

It is build uppon the PrefixSpan Algorithm

Why the need for a new algorithm ?

- Less memory usage
- Faster to find sequential pattern

PrefixSpan : Start

First let's explain the PrefixSpan Algorithm

It's the most efficient pattern mining algorithm

We start with a sequence database

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$

PrefixSpan : Pattern-growth

It works by pattern-growth, which does not generate any candidates (saving memory)

1. Scan : *Calculate support for each item and existing itemset*
2. Output : *Output item that have enough support*
3. Projection : *Recursively project the database with every item that have enough support*



PrefixSpan : Scan

MinSup 75 % (3)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$

Item	Support
1	
2	
3	
4	
5	
6	
7	

PrefixSpan : Scan

MinSup 75 % (3)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$

Item	Support
1	100% (4)
2	
3	
4	
5	
6	
7	

PrefixSpan : Scan

MinSup 75 % (3)

$\langle \{ \}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{ \} \rangle$

$\langle \{5\}, \{ \}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{ \}, \{3\}, \{2\}, \{3\} \rangle$

Item	Support
1	100% (4)
2	100% (4)
3	
4	
5	
6	
7	

PrefixSpan : Scan

MinSup 75 % (3)

$\langle \{ \}, \{3\}, \{4\}, \{6\}, \{ \}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{ \}, \{ \} \rangle$

$\langle \{5\}, \{ \}, \{4\}, \{3\}, \{ \} \rangle$

$\langle \{5\}, \{7\}, \{ \}, \{3\}, \{ \}, \{3\} \rangle$

Item	Support
1	100% (4)
2	100% (4)
3	100% (4)
4	
5	
6	
7	

PrefixSpan : Scan

MinSup 75 % (3)

$\langle \{ \}, \{ \}, \{4\}, \{6\}, \{ \}, \{ \} \rangle$

$\langle \{4\}, \{ \}, \{ \}, \{ \} \rangle$

$\langle \{5\}, \{ \}, \{4\}, \{ \}, \{ \} \rangle$

$\langle \{5\}, \{7\}, \{ \}, \{ \}, \{ \}, \{ \} \rangle$

Item	Support
1	100% (4)
2	100% (4)
3	100% (4)
4	75% (3)
5	
6	
7	

PrefixSpan : Scan

MinSup 75 % (3)

$\langle \{\}, \{\}, \{\}, \{6\}, \{\}, \{\} \rangle$

$\langle \{\}, \{\}, \{\}, \{\} \rangle$

$\langle \{5\}, \{\}, \{\}, \{\}, \{\} \rangle$

$\langle \{5\}, \{7\}, \{\}, \{\}, \{\}, \{\} \rangle$

Item	Support
1	100% (4)
2	100% (4)
3	100% (4)
4	75% (3)
5	50% (2)
6	
7	

PrefixSpan : Scan

MinSup 75 % (3)

$\langle \{\}, \{\}, \{\}, \{6\}, \{\}, \{\} \rangle$

$\langle \{\}, \{\}, \{\}, \{\} \rangle$

$\langle \{\}, \{\}, \{\}, \{\}, \{\} \rangle$

$\langle \{\}, \{7\}, \{\}, \{\}, \{\}, \{\} \rangle$

Item	Support
1	100% (4)
2	100% (4)
3	100% (4)
4	75% (3)
5	50% (2)
6	25% (1)
7	

PrefixSpan : Scan

MinSup 75 % (3)

<{ }, { }, { }, { }, { }, { }>

<{ }, { }, { }, { }>

<{ }, { }, { }, { }, { }>

<{ }, { 7 }, { }, { }, { }, { }>

Item	Support
1	100% (4)
2	100% (4)
3	100% (4)
4	75% (3)
5	50% (2)
6	25% (1)
7	25% (1)

PrefixSpan : Scan

MinSup 75 % (3)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$

Item	Support
1	100% (4)
2	100% (4)
3	100% (4)
4	75% (3)
5	50% (2)
6	25% (1)
7	25% (1)

PrefixSpan : Scan

We take each item with the support \geq minsup, and output it as a sequence with one item.

Here the output :

$\langle \{1\} \rangle$

$\langle \{2\} \rangle$

$\langle \{3\} \rangle$

$\langle \{4\} \rangle$



PrefixSpan : Scan

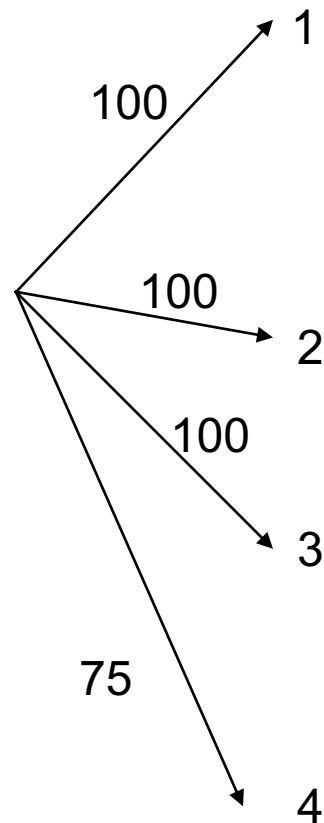
MinSup 75 % (3)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$



New Concept : Projection

We need to define a new concept \rightarrow *Projection*

If we project a sequence $\langle \{1\}, \{2\}, \{3\} \rangle$ by a prefix $\langle \{1\} \rangle$, we take the part of the sequence that follow the prefix. Here $\langle \{2\}, \{3\} \rangle$

Some examples :

$$\langle \{1\}, \{2\}, \{1\}, \{3\} \rangle \text{ by } \langle \{1\} \rangle \rightarrow \langle \{2\}, \{1\}, \{3\} \rangle$$

$$\langle \{3\}, \{4\}, \{5\} \rangle \text{ by } \langle \{3\}, \{4\} \rangle \rightarrow \langle \{5\} \rangle$$

$$\langle \{1\}, \{3, 4\}, \{5\}, \{6\} \rangle \text{ by } \langle \{3\} \rangle \rightarrow \langle \{5\}, \{6\} \rangle$$

$$\langle \{2\}, \{3\}, \{4\}, \{5\}, \{6\} \rangle \text{ by } \langle \{3\}, \{5\} \rangle \rightarrow \langle \{6\} \rangle$$

\rightarrow **Projecting a database, means to project every sequence**

PrefixSpan : Projection

First we will output the result we found with enough support.

Then we will recursively project the database with every of those items.

Lets take the result we have found so far to make it clearer.

We keep *minsup* of 75 % (3)



PrefixSpan : Projection

1. Scan Database

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$

Item	Support
1	100% (4)
2	100% (4)
3	100% (4)
4	75% (3)
5	50% (2)
6	25% (1)
7	25% (1)

PrefixSpan : Projection

2. Output first item

$\langle \{1\} \rangle \rightarrow \text{Support : 100 \% (4)}$

3. Project first item (1)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$

Item	Support
1	100% (4)
2	100% (4)
3	100% (4)
4	75% (3)
5	50% (2)
6	25% (1)
7	25% (1)

PrefixSpan : Projection

1. Scan Again

$\langle \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \rangle$

$\langle \{4\}, \{3\}, \{2\} \rangle$

$\langle \{3\}, \{2\}, \{3\} \rangle$

Item	Support
1	0% (0)
2	75% (3)
3	75% (3)
4	50% (2)
5	0% (0)
6	25% (1)
7	0% (0)

PrefixSpan : Projection

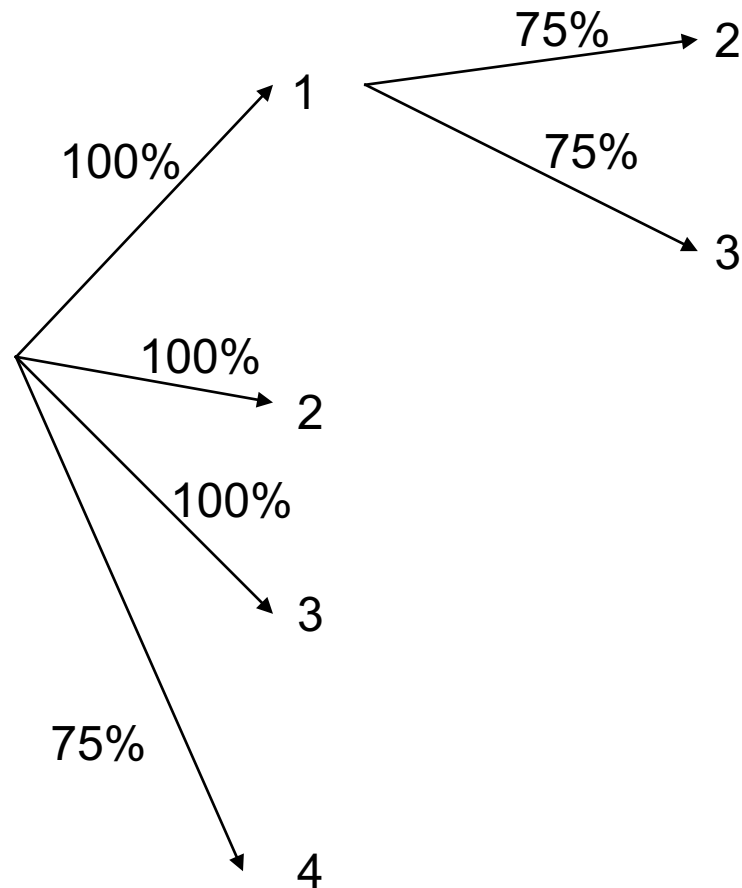
MinSup 75 % (3)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$



PrefixSpan : Projection

2. Output the sequence

$\langle \{1\}, \{2\} \rangle \rightarrow \text{Support : 75 \% (3)}$

3. Project first item (2)

$\langle \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \rangle$

$\langle \{4\}, \{3\}, \{2\} \rangle$

$\langle \{3\}, \{2\}, \{3\} \rangle$

Item	Support
1	0% (0)
2	75% (3)
3	75% (3)
4	50% (2)
5	0% (0)
6	25% (1)
7	0% (0)

PrefixSpan : Projection

1. Scan again

<{3}>

<>

◇

<{3}>

Item	Support
1	0% (0)
2	0% (0)
3	50% (2)
4	0% (0)
5	0 % (0)
6	0% (0)
7	0 % (0)

PrefixSpan : Projection

Operation are over

We continue with other items

Item	Support
1	0% (0)
2	0% (0)
3	50% (2)
4	0% (0)
5	0 % (0)
6	0% (0)
7	0 % (0)

PrefixSpan : Projection

2. Output the sequence

$\langle \{1\}, \{3\} \rangle \rightarrow \text{Support} : 75\% (3)$

3. Project second item (3)

$\langle \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \rangle$

$\langle \{4\}, \{3\}, \{2\} \rangle$

$\langle \{3\}, \{2\}, \{3\} \rangle$

Item	Support
1	0% (0)
2	75% (3)
3	75% (3)
4	50% (2)
5	0% (0)
6	25% (1)
7	0% (0)

PrefixSpan : Projection

1. Scan again

$\langle \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \rangle$

$\langle \{2\} \rangle$

$\langle \{2\}, \{3\} \rangle$

Item	Support
1	0% (0)
2	75% (3)
3	50% (2)
4	25% (1)
5	0% (0)
6	25% (1)
7	0% (0)

PrefixSpan : Projection

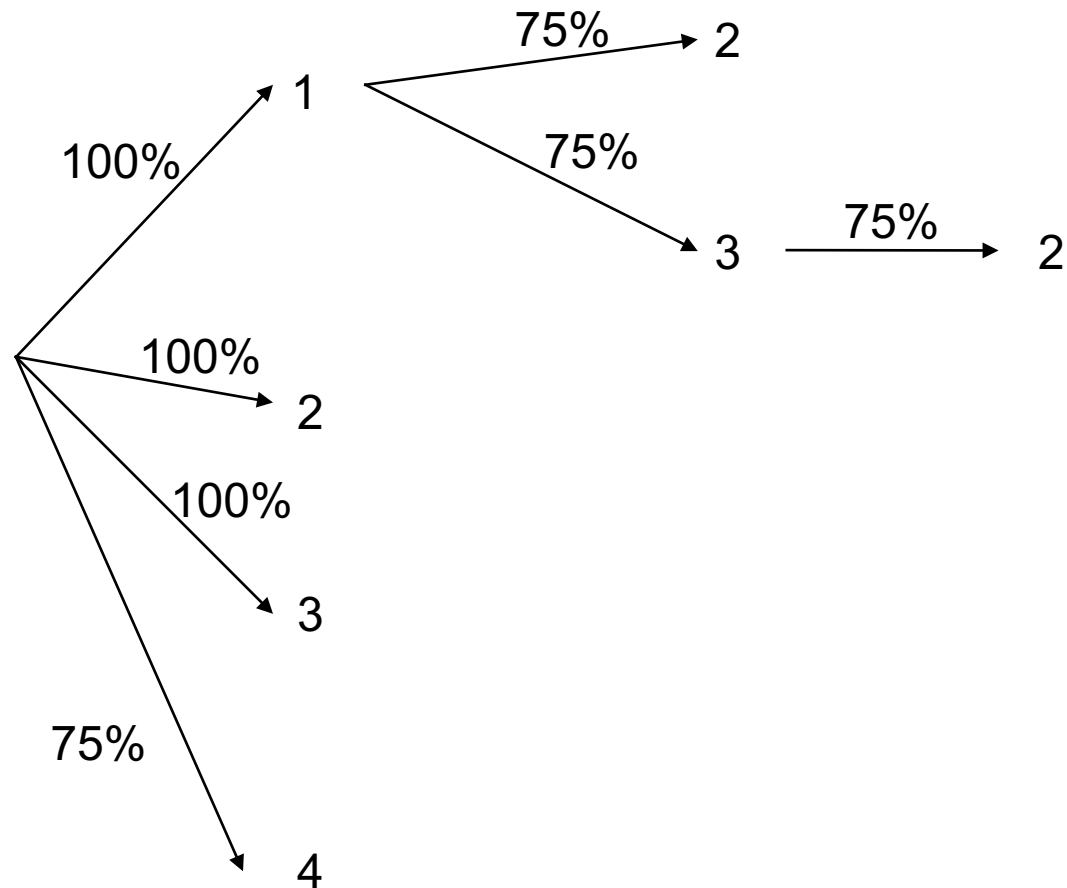
MinSup 75 % (3)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$



PrefixSpan : Projection

2. Output the sequence

$\langle \{1\}, \{3\}, \{2\} \rangle \rightarrow \text{Support : 75 \%}$
(3)

3. Project first item (2)

$\langle \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \rangle$

$\langle \{2\} \rangle$

$\langle \{2\}, \{3\} \rangle$

Item	Support
1	0% (0)
2	75% (3)
3	50% (2)
4	25% (1)
5	0% (0)
6	25% (1)
7	0% (0)

PrefixSpan : Projection

Doing all recursive projection would take time

We will skip to the final result, as the process repeats itself



PrefixSpan : Projection

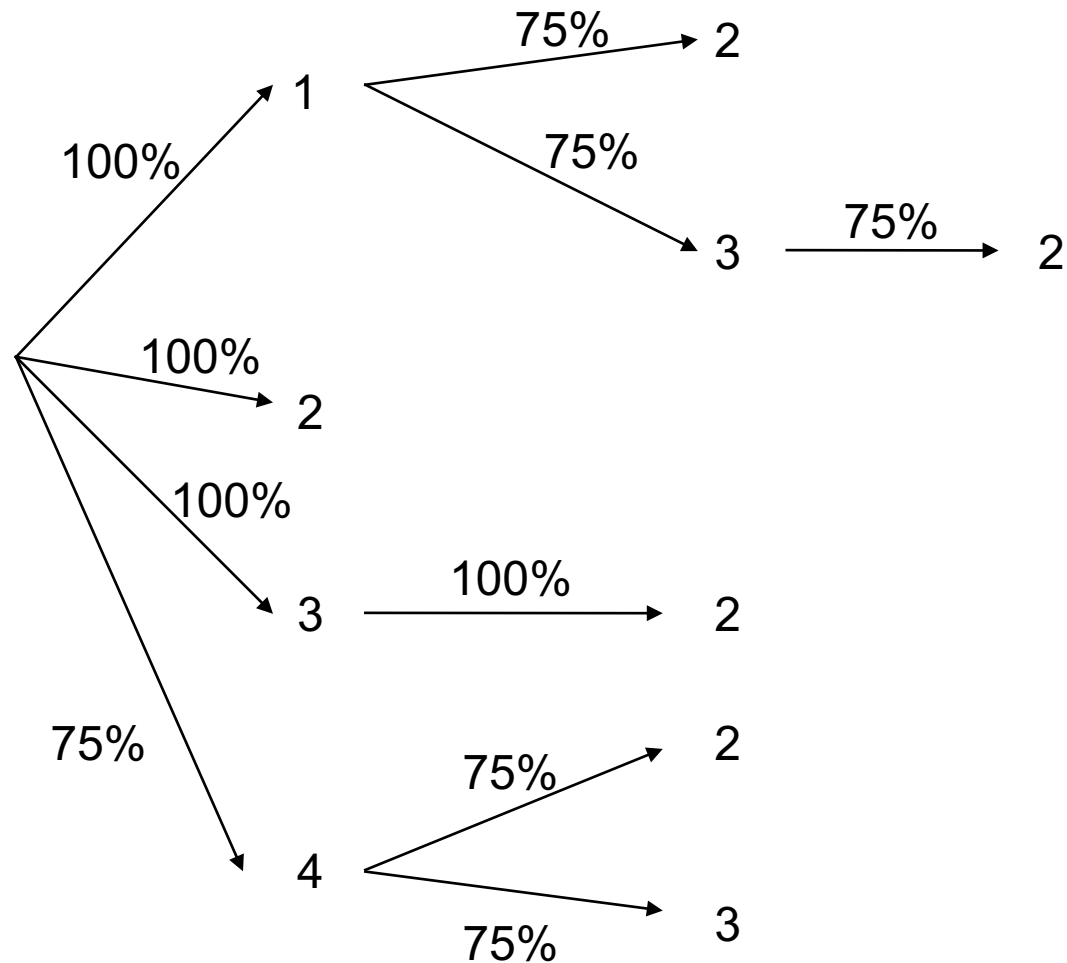
MinSup 75 % (3)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$



PrefixSpan : Projection

At the end of the execution we will get the following result.

Pattern	Support ($\geq 75\%$)	Pattern	Support ($\geq 75\%$)
$\langle\{1\}\rangle$	100 %	$\langle\{3\}\rangle$	100 %
$\langle\{1\},\{2\}\rangle$	75 %	$\langle\{3\},\{2\}\rangle$	100 %
$\langle\{1\},\{3\}\rangle$	75%	$\langle\{4\}\rangle$	75 %
$\langle\{1\},\{3\},\{2\}\rangle$	75 %	$\langle\{4\},\{2\}\rangle$	75 %
$\langle\{2\}\rangle$	100 %	$\langle\{4\},\{3\}\rangle$	75 %

Question 2 : Projection

Considering the database

1. What is the result of the projection of $\langle \{b\}, \{f\} \rangle$ on the database ?
2. In previous sequence, which are not closed and which are maximal ?

$\langle \{a,b\}, \{c\}, \{f,g\}, \{g\}, \{e\} \rangle$

$\langle \{a,d\}, \{c\}, \{b\}, \{a,b,e,f\} \rangle$

$\langle \{a\}, \{b\}, \{f,g\}, \{e\} \rangle$

$\langle \{b\}, \{f,g\} \rangle$

Pattern	Support ($\geq 75\%$)	Pattern	Support ($\geq 75\%$)
$\langle \{1\} \rangle$	100 %	$\langle \{3\} \rangle$	100 %
$\langle \{1\}, \{2\} \rangle$	75 %	$\langle \{3\}, \{2\} \rangle$	100 %
$\langle \{1\}, \{3\} \rangle$	75%	$\langle \{4\} \rangle$	75 %
$\langle \{1\}, \{3\}, \{2\} \rangle$	75 %	$\langle \{4\}, \{2\} \rangle$	75 %
$\langle \{2\} \rangle$	100 %	$\langle \{4\}, \{3\} \rangle$	75 %

PrefixSpan : Projection

Here are the **closed**.

Pattern	Support ($\geq 75\%$)	Pattern	Support ($\geq 75\%$)
<{1}>	100 %	<{3}>	100 %
<{1},{2}>	75 %	<{3},{2}>	100 %
<{1},{3}>	75%	<{4}>	75 %
<{1},{3},{2}>	75 %	<{4},{2}>	75 %
<{2}>	100 %	<{4},{3}>	75 %

PrefixSpan : Projection

Here are the **closed** and **maximal**.

Pattern	Support ($\geq 75\%$)	Pattern	Support ($\geq 75\%$)
<{1}>	100 %	<{3}>	100 %
<{1},{2}>	75 %	<{3},{2}>	100 %
<{1},{3}>	75%	<{4}>	75 %
<{1},{3},{2}>	75 %	<{4},{2}>	75 %
<{2}>	100 %	<{4},{3}>	75 %

MaxSp : Basic Idea

MaxSP extends the PrefixSpan

A naïve approach would be to keep all sequence in memory and to check every time a new sequence arrives if it is maximal.

That is CloSpan

- Inefficient
- Memory consuming

MaxSp : Basic Idea

The question is

How to know if a pattern P is maximal, without maintaining pattern in memory ?

The solution : Can P be extended by appending items ?

→ YES ? It isn't a maximal sequential pattern

Two check :

1. Maximal backward extension check
2. Maximal forward extension check

MaxSp : Maximal forward extension check

With the Maximal forward extension, we search if we can extends a pattern with upcoming items

→ Concretly it is what PrefixSpan already do



MaxSp : Maximal backward extension check

With maximal backward extension we check if we can extend the pattern with item that we might have passed.

If so we stop looking further because we will find the pattern in another branch.

Let's look at an example.



MaxSp : Maximal backward extension check

To find the maximal backward extension of a prefix P in a sequence S, there is a few steps

1. Find the last in last appearance of all items
2. Find the Maximum period
3. If an item support is \geq than the minsup, there is a backward extension

Lets demonstrate, step by step



MaxSp : Maximal backward extension check

1. Find the last in last appearance of all item.

→ Just find the prefix in reverse order (starting from the end of the Sequence)

P = ABC

S = ABDBCAB

Current

Found Item

Discarded



MaxSp : Maximal backward extension check

1. Find the last in last appearance of all item.

→ Just find the prefix in reverse order (starting from the end of the Sequence)

P = ABC

S = ABDDBCAB

Current

Found Item

Discarded



MaxSp : Maximal backward extension check

1. Find the last in last appearance of all item.

→ Just find the prefix in reverse order (starting from the end of the Sequence)

P = ABC

S = ABD^BC^AB

Current

Found Item

Discarded



MaxSp : Maximal backward extension check

1. Find the last in last appearance of all item.

→ Just find the prefix in reverse order (starting from the end of the Sequence)

P = ABC

S = ABDBCAB

Current

Found Item

Discarded



MaxSp : Maximal backward extension check

1. Find the last in last appearance of all item.

→ Just find the prefix in reverse order (starting from the end of the Sequence)

P = ABC

S = ABDBCAB

Current

Found Item

Discarded



MaxSp : Maximal backward extension check

1. Find the last in last appearance of all item.

→ Just find the prefix in reverse order (starting from the end of the Sequence)

P = ABC

S = ABDBCAB

Current

Found Item

Discarded



MaxSp : Maximal backward extension check

2. Find the maximum period of the prefix in the sequence.

- There is as much maximum period as there is items in the prefix (so here 3)
- It is the space between the last-in-last appearance of item i and the first occurrence of the prefix before the item i

1st Maximum Period

$P = ABC$

$S = [][ABDBCAB]$

$MP = (\text{none})$

2nd Maximum Period

$P = ABC$

$S = [A]BD[BCAB]$

$MP = BD$

3rd Maximum Period

$P = ABC$

$S = [AB]DB[CAB]$

$MP = DB$

MaxSp : Maximal backward extension check

We do that with all sequence in the database

3. We count support for all item in every maximum period.

If for a same i -th maximum period, one item support $\geq \text{minsup} \rightarrow$ There is a Maximal backward extension, the element is not maximal.

Question 3 : Extension

With sequence Database and minsup = 50 %

ABCDBC

CDBDCA

ACBDB

1. Is there any Maximal Forward Extension ?

P = DB

2. Is there any Maximal Backward Extension ?

P = AD

MaxSp Algorithm

Concretly, how does one implement MaxSP based on PrefixSpan ?

The same as PrefixSpan

BUT

1. Before adding a prefix, check that there is no maximal backward extension
2. Add a pattern at the end of the recursive call (at a leaf on the DFS)

MaxSp Algorithm : Steps

Here are the different steps :

1. Scan
2. Project
3. Is there any Maximal Forward Extension in non modified DB (PrefixSpan) ?
Yes, not maximal
4. Else Is there any Maximal Backward Extension in non modified DB ?
Yes, not maximal
5. If previous questions are No ?
Output the sequence

MaxSP: Lets Start From the end

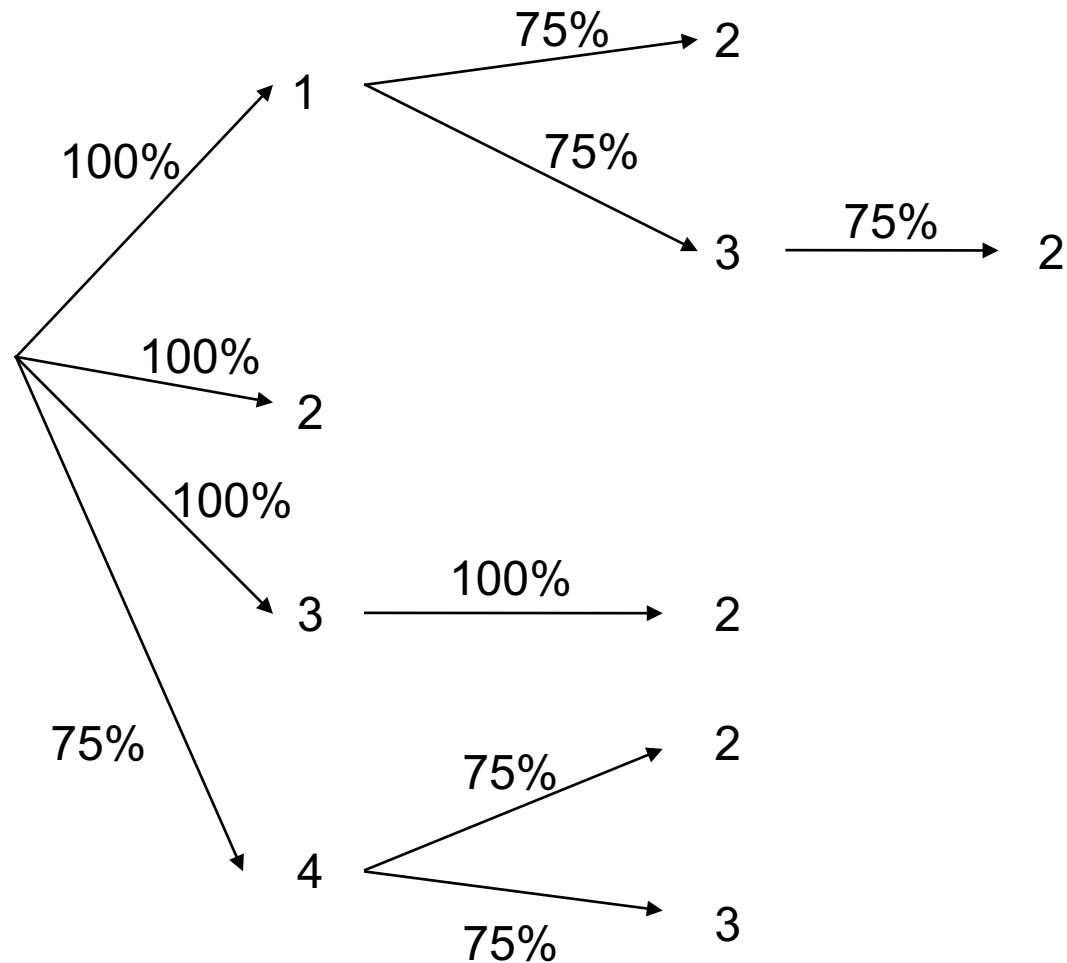
MinSup 75 % (3)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$



MaxSP: Lets Start From the end

MinSup 75 % (3)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$

Here are all possibilities without Maximal Forward Extension

$\langle \{1\}, \{2\} \rangle$

$\langle \{1\}, \{3\}, \{2\} \rangle$

$\langle \{2\} \rangle$

$\langle \{3\}, \{2\} \rangle$

$\langle \{4\}, \{2\} \rangle$

$\langle \{4\}, \{3\} \rangle$

To know if we must add them as maximal, you just check for each of the if there is a Maximal Backward Extension

MaxSP: Lets Start From the end

MinSup 75 % (3)

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$

Lets show with $\langle \{1\}, \{2\} \rangle$

MaxSP: Lets Start From the end

First Maximal Period: MinSup 75 % (3) With $\langle \{1\}, \{2\} \rangle$

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

No item has support \geq MinSup

$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$

Second Maximal Period: MinSup 75 % (3) With $\langle \{1\}, \{2\} \rangle$

$\langle \{1\}, \{3\}, \{4\}, \{6\}, \{2\}, \{3\} \rangle$

$\langle \{4\}, \{3\}, \{2\}, \{1\} \rangle$

Item $\{3\}$ has support of 75%

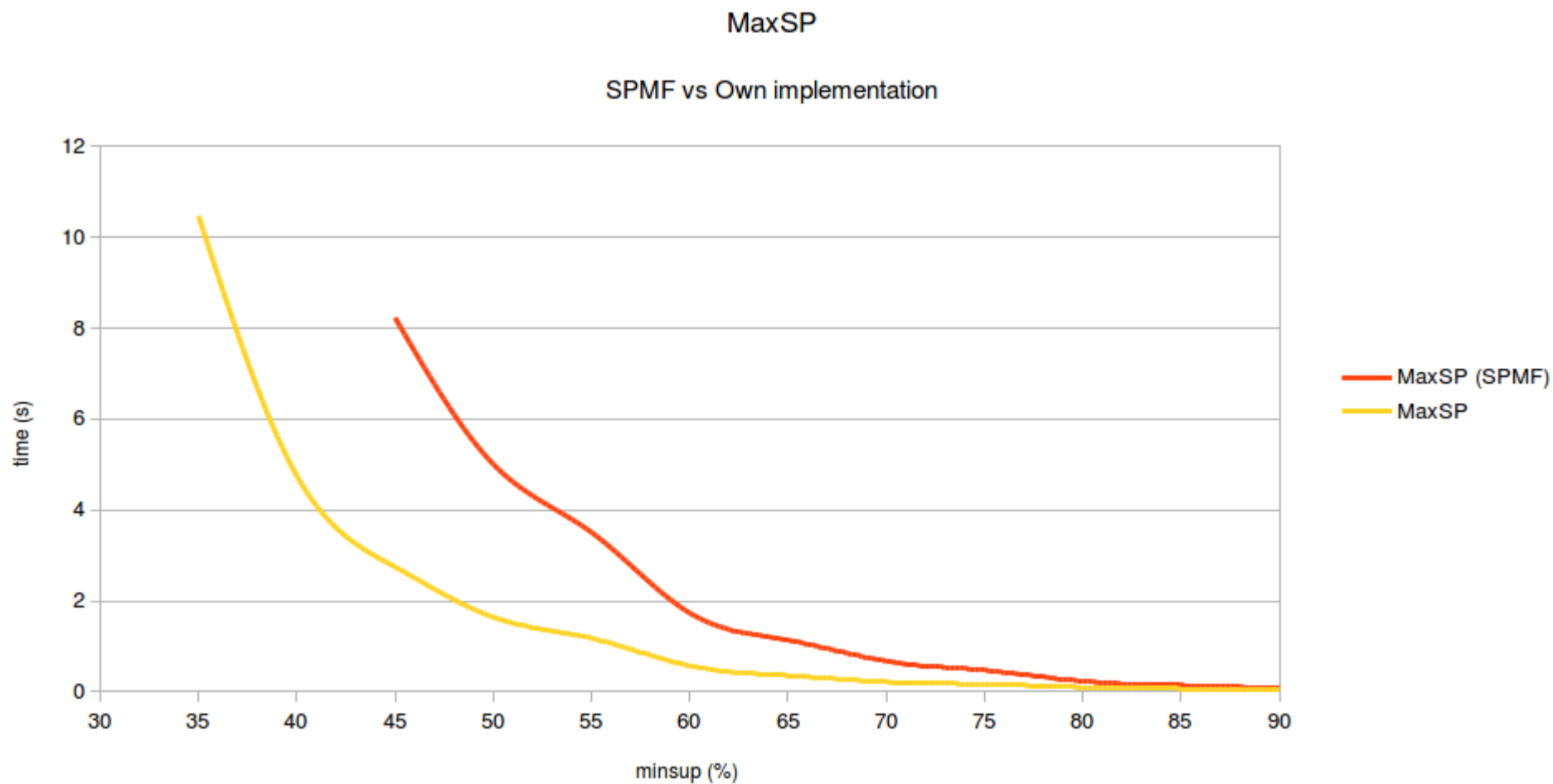
$\langle \{5\}, \{1\}, \{4\}, \{3\}, \{2\} \rangle$

there is a MBE, so $\langle \{1\}, \{2\} \rangle$

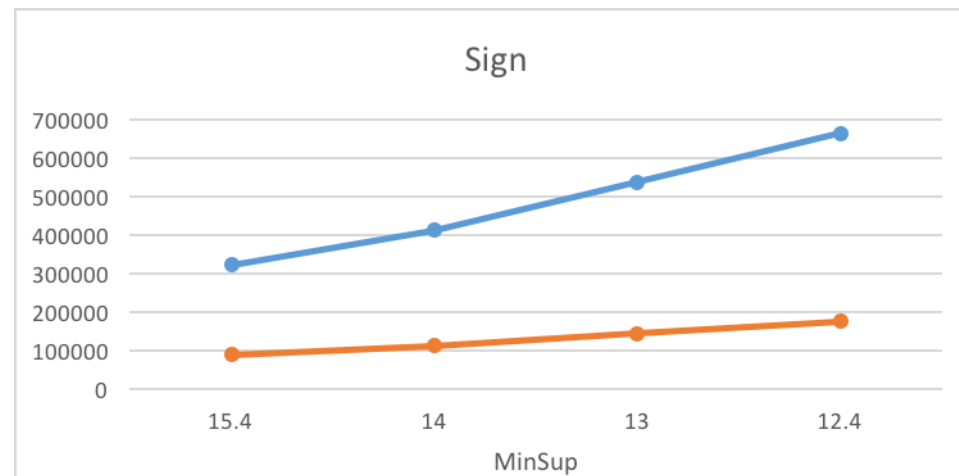
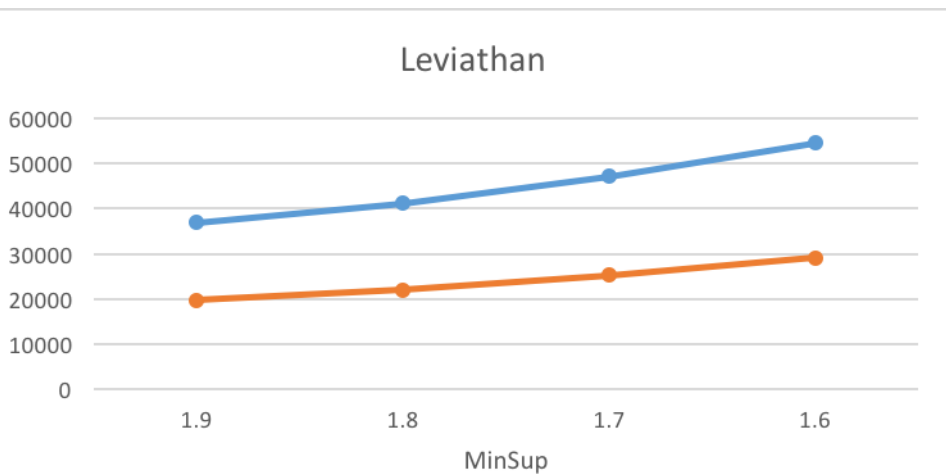
$\langle \{5\}, \{7\}, \{1\}, \{3\}, \{2\}, \{3\} \rangle$

is not maximal sequential pattern

Measures : Own MaxSP vs SPMF MaxSP

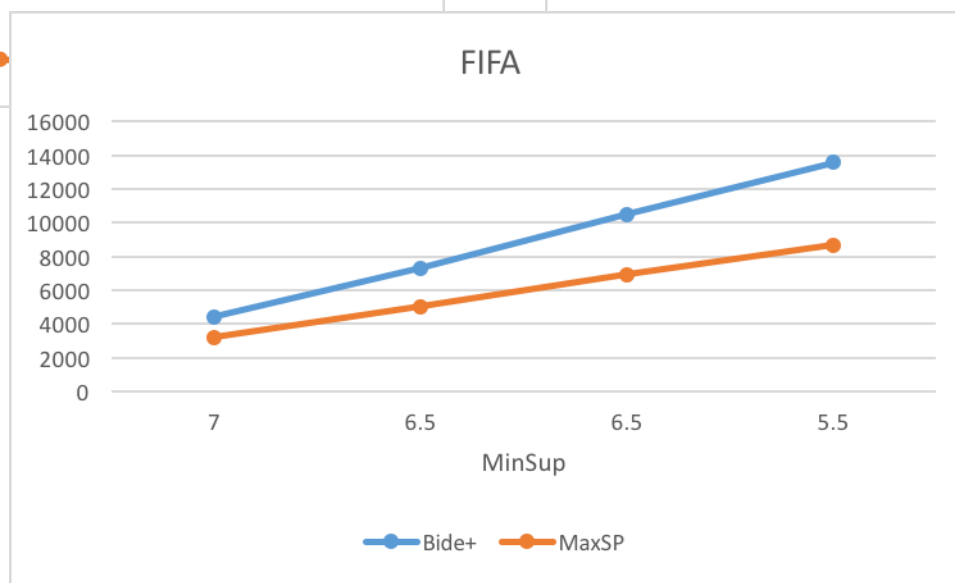


Measures : MaxSP vs Bide+ (Frequent sequences count)



—●— Bide+ —●—

—●— Bide+ —●— MaxSP

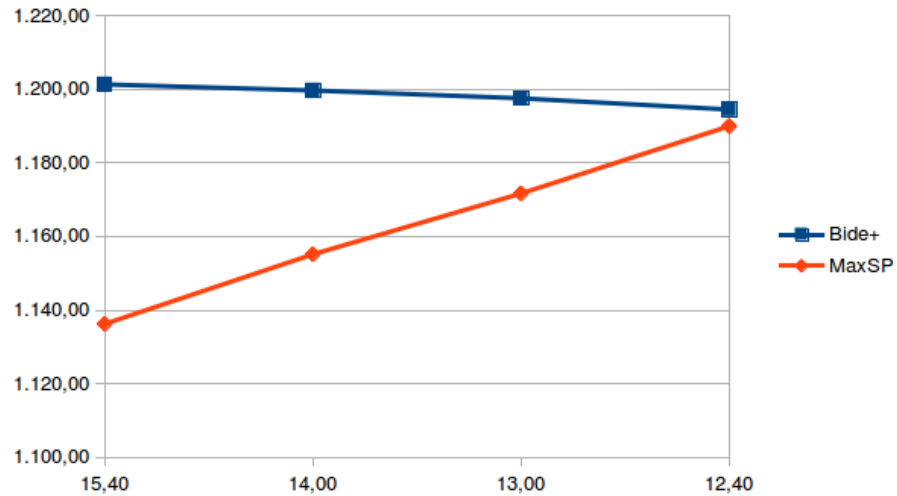


—●— Bide+ —●— MaxSP

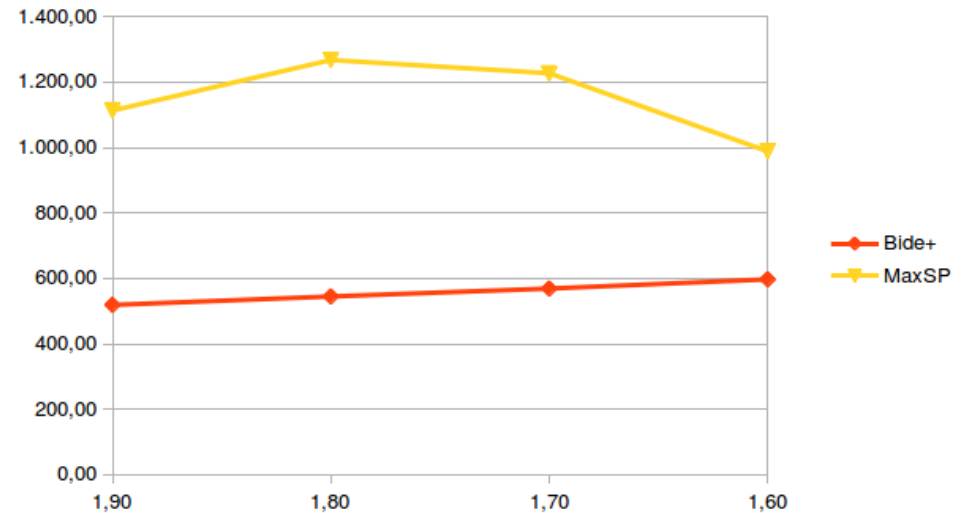
Measures : MaxSP vs Bide+ (Time)



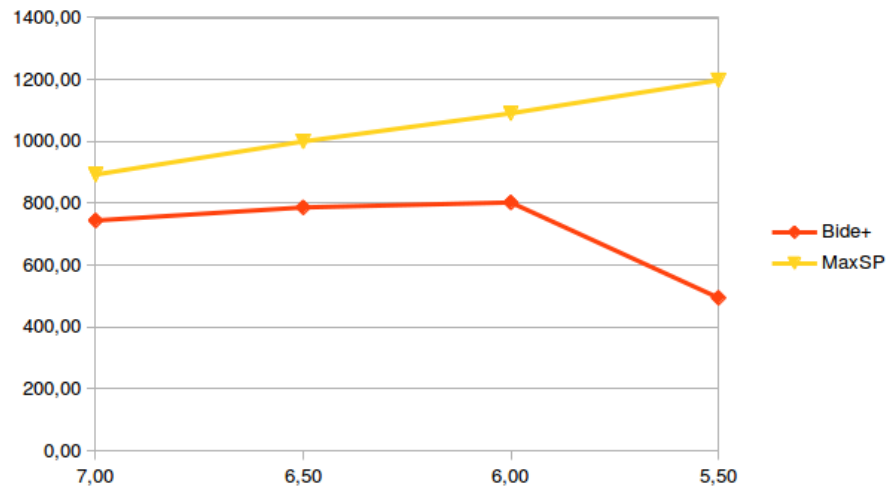
Measures : MaxSP vs Bide+ (Memory)



Sign



Leviathan



FIFA

Last Considerations

We did not considered examples with itemsets (as in the articles), it is not really different, you can apply the same rules

We implemented the Algorithm, the presentation and algorithm are available on our git repository. The algorithm can be optimized as stated in the paper

The End

Thanks for listening !