



**OPPA European Social Fund  
Prague & EU: We invest in your future.**

---

# Frequent subsequences, episodal rules

---

**Jiří Kléma**

Department of Cybernetics,  
Czech Technical University in Prague



<http://ida.felk.cvut.cz>















# Canonical form for sequences

---

- a **canonical** (standard) **code word**

- a unique sequence representation, based on the symbol alphabet ordering,
- a usual (not necessary) choice:
  - \* the lexicographical symbol alphabet ordering  $a < b < c < \dots$ ,
  - \* the lexicographically smallest (smaller) code word taken as canonical ( $bac < cab$ ),

- a directed sequence

- the only interpretation (way of reading), each (sub)sequence is a canonical code word,

- an undirected sequence

- two possible ways of reading = two alternative code words,
- the routine application of lexicographical ordering is not possible,
- **prefix property** in a space of canonical code words does not hold:
  - \* every prefix of a canonical word is a canonical word itself,

sequence	canonical form	prefix	canonical form
<i>bab</i>	<i>bab</i>	<i>ba</i>	<i>ab</i>
<i>cabd</i>	<i>cabd</i>	<i>cab</i>	<i>bac</i>

- we have to find a different way of forming code words.

## Canonical form for undirected sequences

- The canonical code words with the prefix property will be formed as follows
  - even and odd length words will be handled separately,
  - code words are started in the middle of sequence,

	even length	odd length
sequence	$a_m a_{m-1} \dots a_2 a_1 b_1 b_2 \dots b_{m-1} b_m$	$a_m a_{m-1} \dots a_2 a_1 a_0 b_1 b_2 \dots b_{m-1} b_m$
code word	$a_1 b_1 a_2 b_2 \dots a_{m-1} b_{m-1} a_m b_m$	$a_0 a_1 b_1 a_2 b_2 \dots a_{m-1} b_{m-1} a_m b_m$
code word	$b_1 a_1 b_2 a_2 \dots b_{m-1} a_{m-1} b_m a_m$	$a_0 b_1 a_1 b_2 a_2 \dots b_{m-1} a_{m-1} b_m a_m$

- canonical is the lexicographically smaller code word in the table,
- the sequence is **extended** by adding
  - a pair  $a_{m+1} b_{m+1}$  or  $b_{m+1} a_{m+1}$ ,
  - one item at the front and one item at the end.
- an example

even length			odd length		
sequence	code words		sequence	code words	
<i>at</i>	<i>at</i>	<i>ta</i>	<i>ule</i>	<i>lue</i>	<i>leu</i>
<i>data</i>	<i>atda</i>	<i>taad</i>	<i>rules</i>	<i>luers</i>	<i>leusr</i>























# PrefixSpan – algorithm, example ( $s_{min} = 2$ )

- PrefixSpan: the input is  $S$  and  $s_{min}$ 
  1.  $i = 1$ , the init projection prefix database  $S|_{\alpha_0} = S|_{\emptyset} = S$ ,
  2. repeat for all the projection prefix databases  $S|_{\alpha_{i-1}}$ 
    - (a) find frequent i-patterns (sufficient support in  $\alpha_{i-1} \cdot S|_{\alpha_{i-1}}$ ),
    - (b) until the set of i-patterns is not empty
      - i. split the state space having the i-patterns ( $\alpha_i$ ) as prefixes
        - a projection database set originates  $S|_{\alpha_i} = (\alpha_{i-1} \cdot S|_{\alpha_{i-1}})|_{\alpha_i}$ ,
        - ii.  $i=i+1$  and go to the step (2).

Id	Sequence	Prefix	Projection database (postfixes) or patterns
		$\langle a \rangle$	$\langle (abc)(ac)d(cf) \rangle, \langle (-d)c(bc)(ae) \rangle, \langle (-b)(df)cb \rangle, \langle (-f)cbc \rangle$
			2-patterns: $\langle aa \rangle:2, \langle ab \rangle:4, \langle ac \rangle:4, \langle ad \rangle:2, \langle af \rangle:2, \langle (ab) \rangle:2$
10	$\langle a(abc)(ac)d(cf) \rangle$	$\langle b \rangle$	$\langle (-c)(ac)d(cf) \rangle, \langle (-c)(ae) \rangle, \langle (df)cb \rangle, \langle c \rangle$
20	$\langle (ad)c(bc)(ae) \rangle$		2-patterns: $\langle ba \rangle:2, \langle bc \rangle:3, \langle (bc) \rangle:2, \langle bd \rangle:2, \langle bf \rangle:2$
30	$\langle (ef)(ab)(df)cb \rangle$	$\langle aa \rangle$	$\langle (-bc)(ac)d(cf) \rangle, \langle (-e) \rangle$
40	$\langle eg(af)cbc \rangle$		<b>STOP (no 3-patterns)</b>
		$\langle (ab) \rangle$	$\langle (-c)(ac)d(cf) \rangle, \langle (df)cb \rangle$
			3-patterns: $\langle (ab)c \rangle:2, \langle (ab)d \rangle:2, \langle (ab)f \rangle:2$













**OPPA European Social Fund  
Prague & EU: We invest in your future.**

---