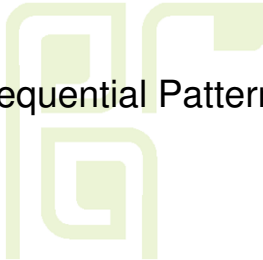


WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME



Sequential Patterns

HUMAN CAPITAL
HUMAN - BEST INVESTMENT

EUROPEAN UNION
EUROPEAN SOCIAL FUND

Project is co-financed by European Union within European Social Fund

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Sequential Patterns - Informally

- Sequential patterns – patterns occurring frequently in data sequences in which the order of elements is important.
- Example:** In the case of a set of events, the order is determined by timestamps, while in the case of a document, the order is determined by positions of paragraphs, sentences or words.

2

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Sequential Patterns

Sample dataset D

CId	Tid	Items
1	10	cd
1	15	abc
1	20	abf
1	25	acdf
2	15	abf
2	20	e
3	10	abf
4	10	dgh
4	20	bf
4	25	agh

- In the context of market basket data, a sequential pattern is a typical purchase behavior of customers.
- Example dataset D contains 4 data (customer) sequences.
- Purchase sequence $\langle(d)(bf)(a)\rangle$ occurs for customers: 1 and 4 in D.

3

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Support of a Sequence

- Support of a sequence S* is denoted as $sup(S)$ and defined as the number of data sequences containing S.
- Property.** Support of a subsequence S of a sequence S' is not less than $sup(S')$.
- Property.** Support of a supersequence S of a sequence S' is not greater than $sup(S')$.

4

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Sequence's Support

Dataset D

CId	Tid	Items
1	10	cd
1	15	abc
1	20	abf
1	25	acdf
2	15	abf
2	20	e
3	10	abf
4	10	dgh
4	20	bf
4	25	agh

- $sup(\langle(d)(bf)(a)\rangle)=2$
- $sup(\langle(b)(a)\rangle)=2 \geq sup(\langle(d)(bf)(a)\rangle)$
- $sup(\langle(cd)(bf)(a)\rangle)=1 \leq sup(\langle(d)(bf)(a)\rangle)$

5

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Sequential Patterns - Formally

- Sequence S is defined as a *sequential pattern* (or alternatively, as a *frequent sequence*) if its support is above a threshold $minSup$.

6

SPADE: Creation of Candidate Sequences

- Candidate sequences of size n are created from pairs of sequential patterns of size $n-1$.

7

SPADE: Creation of Candidates Sequences of Size 2

Sequential pattern	Sequential pattern	Candidate sequences for $x \neq y$	Candidate sequences for $x = y$
$\langle x \rangle$	$\langle y \rangle$	$\langle \textcolor{red}{x}y \rangle$ $\langle x \textcolor{red}{y} \rangle$ $\langle \textcolor{red}{y}x \rangle$	$\langle \textcolor{red}{x}x \rangle$

8

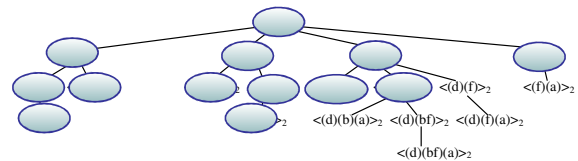
SPADE: Creation of Candidates of Size Greater than 2

Sequential pattern 1	Sequential pattern 2	Candidate sequences for $x \neq y$	Candidate sequences for $x = y$
$\langle G(P)x \rangle$	$\langle G(P)y \rangle$	$\langle G(P)\textcolor{red}{x}y \rangle$ $\langle G(P)\textcolor{red}{x}(y) \rangle$ $\langle G(P)y\textcolor{red}{x} \rangle$	$\langle G(P)\textcolor{red}{x}(x) \rangle$
$\langle G(Px) \rangle$	$\langle G(Py) \rangle$	$\langle G(P\textcolor{red}{x}y) \rangle$	-
$\langle G(Px) \rangle$	$\langle G(P)y \rangle$	$\langle G(P\textcolor{red}{x})(y) \rangle$	$\langle G(P\textcolor{red}{x})(x) \rangle$
$\langle G(P)x \rangle$	$\langle G(Py) \rangle$	$\langle G(P\textcolor{red}{y})(x) \rangle$	$\langle G(P\textcolor{red}{y})(y) \rangle$

9

Example: Result of SPADE

- $\text{minSup} = 1$.

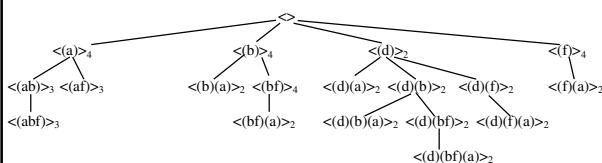


- Infrequent candidate sequences are skipped.

10

Example: Result of SPADE

- $\text{minSup} = 1$.



- Infrequent candidate sequences are skipped.

11

SPADE: Evaluation of Candidate Sequences

- Evaluation of each candidate sequence S is carried out by means of its list of transaction identifiers (shortly, *tidlists*; denoted as $t(S)$); namely,

$\text{sup}(S)$ is equal to the number of candidate sequences registered in $t(S)$.

12

SPADE: Tidlists of Sequences of Size 1

- Determined based on given dataset D.

$t(<(a)>)$		$t(<(b)>)$		$t(<(d)>)$		$t(<(f)>)$	
Cld	Tld	Cld	Tld	Cld	Tld	Cld	Tld
1	15	1	15	1	10	1	20
1	20	1	20	1	25	1	25
1	25	2	15	4	10	2	15
2	15	3	10			3	10
3	10	4	20			4	20
4	25						

- $sup(<(a)>) = 4$, $sup(<(b)>) = 4$, $sup(<(d)>) = 2$, $sup(<(f)>) = 4$.

13

SPADE: Tidlists of Sequences of Size Greater than 1

- Determined based on tidlists of parents.

$t(<(d)>)$		$t(<(f)>)$		$t(<(df)>)$		$t(<(d)(f)>)$	
Cld	Tld	Cld	Tld	Cld	Tld	Cld	Tld
1	10	1	20	1	25	1	20
1	25	1	25			1	25
4	10	2	15			4	20
		3	10				
		4	20				

d and f occur simultaneously in a data sequence

f occurs later than d in a data sequence

14

SPADE: Tidlists of Sequences of Size Greater than 1

- Determined based on tidlists of parents.

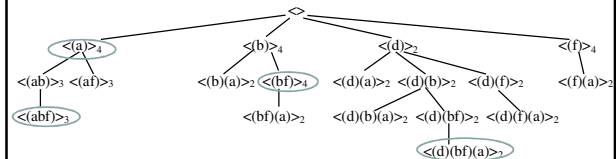
$t(<(d)(b)(a)>)$		$t(<(d)(bf)>)$		$t(<(d)(bf)(a)>)$	
Cld	Tld	Cld	Tld	Cld	Tld
1	20	1	20	1	25
1	25	4	20	4	25
4	25				

a occurs later than <(d)(bf)> in a data sequence

15

Closed Sequential Patterns

- $minSup = 1$.



- A sequential pattern S is closed if all its proper supersequences have supports different from (less than) $sup(S)$.

16

GSP: Constraints on Elements of Generalized Sequential Patterns

Cld	Tld	Items
1	10	1, 2
1	25	4, 6
1	45	3
1	50	1, 2
1	65	3
1	90	2, 4
1	95	6
2
...

- for each element:
 - $t_k - t_s \leq wS$
- for each 2 neighboring elements:
 - $t'_s - t_k > minGap$
 - $t'_k - t_s \leq maxGap$

- $wS = 20$ (0)
- $minGap = 19$ (0)
- $maxGap = 50$ (∞)
- $sup(<(1,4)(2,3)(2,6)>) = ?$

$\begin{matrix} <(1,4) & (2,3) & (2,6)> \\ \longleftrightarrow & \longleftrightarrow & \longleftrightarrow \\ t_s & t_k & t'_s & t'_k & t''_s & t''_k \end{matrix}$

17

GSP: Example...

Cld	Tld	Items
1	10	1, 2
1	25	4, 6
1	45	3
1	50	1, 2
1	65	3
1	90	2, 4
1	95	6
2
...

- constraint for each element:
 - $t_k - t_s \leq wS$
- constraints for each 2 neighboring elements:
 - $t'_s - t_k > minGap$
 - $t'_k - t_s \leq maxGap$

- $wS = 20$ (0)
- $minGap = 19$ (0)
- $maxGap = 50$ (∞)

- $sup(<(1,4)(2,3)(2,6)>) = ?$

$\begin{matrix} <(1,4) & (2,3) & (2,6)> \\ \longleftrightarrow & \longleftrightarrow & \longleftrightarrow \\ t_s & t_k & t'_s & t'_k & t''_s & t''_k \\ 10 & 25 & 50 & 65 & 90 & 95 \end{matrix}$

18

Warsaw University of Technology

Development Programme

GSP: Example

CId	Tid	Items
1	10	1, 2
1	25	4, 6
1	45	3
1	50	1, 2
1	65	3
1	90	2, 4
1	95	6
2
...

- constraint for each element:
 - $t_k - t_s \leq wS$
- constraints for each 2 neighboring elements:
 - $t_s' - t_k > minGap$
 - $t_k' - t_s \leq maxGap$
- $wS = 20$ (0)
- $minGap = 19$ (0)
- $maxGap = 50$ (∞)
- $sup(<(1,4)(2,3)(2,6)>) = ?$

$<(1,4)>$

$<(2,3)>$

$<(2,6)>$

$t_s \quad t_k \quad t_s' \quad t_k' \quad t_s'' \quad t_k''$

$10 \quad 25 \quad 45 \quad 50 \quad 90 \quad 95$

19

Warsaw University of Technology

Development Programme

Creating candidates for GSP

CId	Tid	Items	GSP_3	C_3 after merging	C_3 after pruning
1	10	1, 2			
1	25	4, 6	$<(1,2) (3)>$	$<(1,2) (3,4)>$	$<(1,2) (3,4)>$
1	45	3	$<(1,2) (4)>$	$<(1,2) (3) (5)>$	
1	50	1, 2	$<(1) (3,4)>$		
1	65	3	$<(1,3) (5)>$		
1	90	2, 4	$<(2) (3,4)>$		
1	95	6	$<(2) (3) (5)>$		
2			
...			

20

Warsaw University of Technology

Development Programme

References

- Marzena Kryszkiewicz, Łukasz Skonieczny: Fast Discovery of Generalized Sequential Patterns, Intelligent Methods and Big Data in Industrial Applications, 155-170, online 2018
- Ramakrishnan Srikant, Rakesh Agrawal: Mining Sequential Patterns: Generalizations and Performance Improvements. [EDBT 1996](#): 3-17
- Jianyong Wang, Jiawei Han, Chun Li: Frequent Closed Sequence Mining without Candidate Maintenance. [IEEE Trans. Knowl. Data Eng.](#) 19(8): 1042-1056 (2007)
- Mohammed Javeed Zaki: SPADE: An Efficient Algorithm for Mining Frequent Sequences. [Machine Learning](#) 42(1/2): 31-60 (2001)

21