

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Discovering Sequential
Patterns with SPADE

Marzena Kryszkiewicz

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 $x = y$
$\langle x \rangle$	$\langle y \rangle$	$\langle xy \rangle$ $\langle x(y) \rangle$ $\langle (y)x \rangle$	$\langle (x)x \rangle$

8

SPADE: Creation of Candidates of Size Greater than 2

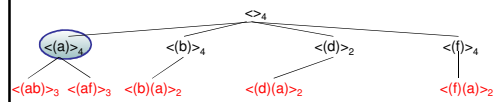
Sequential pattern 1	Sequential pattern 2	Candidate sequences for $x \neq y$	Candidate sequences for $x = y$
$\langle G(P)x \rangle$	$\langle G(P)y \rangle$	$\langle G(P)(xy) \rangle$ $\langle G(P)(x)(y) \rangle$ $\langle G(P)(y)(x) \rangle$	$\langle G(P)(x)(x) \rangle$
$\langle G(Px) \rangle$	$\langle G(Py) \rangle$	$\langle G(Pxy) \rangle$	-
$\langle G(Px) \rangle$	$\langle G(P)y \rangle$	$\langle G(Px)(y) \rangle$	$\langle G(Px)(x) \rangle$
$\langle G(P)x \rangle$	$\langle G(Py) \rangle$	$\langle G(Py)(x) \rangle$	$\langle G(Py)(y) \rangle$

9

Example: Result of SPADE...

- $\text{minSup} = 1$.

$\langle x \rangle$	$\langle y \rangle$	$\langle xy \rangle$ $\langle x(y) \rangle$ $\langle (y)x \rangle$	$\langle (x)x \rangle$
---------------------	---------------------	--	------------------------



Candidates:

- $\langle a \rangle_4$
- $\langle ab \rangle_3, \langle a(b) \rangle_1, \langle (b)a \rangle_2$
- $\langle ad \rangle_1, \langle a(d) \rangle_1, \langle (d)a \rangle_2$
- $\langle af \rangle_3, \langle a(f) \rangle_1, \langle (f)a \rangle_2$

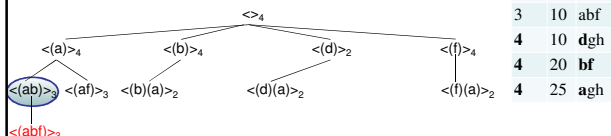
CId	Tid	Items
1	10	cd
1	15	abc
1	20	abf
1	25	acdf
2	15	abf
2	20	e
3	10	abf
4	10	dgh
4	20	bf
4	25	agh

10

Example: Result of SPADE...

$\langle G(P)x \rangle$	$\langle G(P)y \rangle$	$\langle G(P)(xy) \rangle$ $\langle G(P)(x)(y) \rangle$ $\langle G(P)(y)(x) \rangle$	$\langle G(P)(x)(x) \rangle$
$\langle G(Px) \rangle$	$\langle G(Py) \rangle$	$\langle G(Pxy) \rangle$	-
$\langle G(Px) \rangle$	$\langle G(P)y \rangle$	$\langle G(Px)(y) \rangle$	$\langle G(Px)(x) \rangle$
$\langle G(P)x \rangle$	$\langle G(Py) \rangle$	$\langle G(Py)(x) \rangle$	$\langle G(Py)(y) \rangle$

- $\text{minSup} = 1$

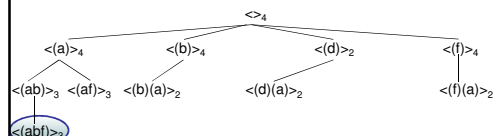


11

Example: Result of SPADE...

$\langle G(P)x \rangle$	$\langle G(P)y \rangle$	$\langle G(P)(xy) \rangle$ $\langle G(P)(x)(y) \rangle$ $\langle G(P)(y)(x) \rangle$	$\langle G(P)(x)(x) \rangle$
$\langle G(Px) \rangle$	$\langle G(Py) \rangle$	$\langle G(Pxy) \rangle$	-
$\langle G(Px) \rangle$	$\langle G(P)y \rangle$	$\langle G(Px)(y) \rangle$	$\langle G(Px)(x) \rangle$
$\langle G(P)x \rangle$	$\langle G(Py) \rangle$	$\langle G(Py)(x) \rangle$	$\langle G(Py)(y) \rangle$

- $\text{minSup} = 1$



12

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

• minSup = 1

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

13

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

• minSup = 1.

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

14

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

minSup = 1.

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

15

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

minSup = 1.

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

16

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

minSup = 1.

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

17

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

• minSup = 1.

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

18

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

minSup = 1.

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

19

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

minSup = 1.

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

20

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

minSup = 1.

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

21

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

minSup = 1.

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

22

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

minSup = 1.

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

23

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Example: Result of SPADE...

minSup = 1.

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

24

Example: Result of SPADE...

minSup = 1.

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

25

Example: Result of SPADE...

minSup = 1.

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

26

Example: Result of SPADE

minSup = 1.

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

27

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,

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

28

SPADE: Tidlists of Sequences of Size 1

- Determined based on given dataset D .

$t(<(a)>)$		$t(<(b)>)$		$t(<(d)>)$		$t(<(f)>)$	
CId	Tid	CId	Tid	CId	Tid	CId	Tid
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						

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

29

SPADE: Tidlists of Sequences of Size Greater than 1

- Determined based on tidlists of parents.

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

Annotations:

- d and f occur simultaneously in a data sequence (pointing to the entry for $t(<(df)>)$).
- f occurs later than d in a data sequence (pointing to the entry for $t(<(d)(f)>)$).

30

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				

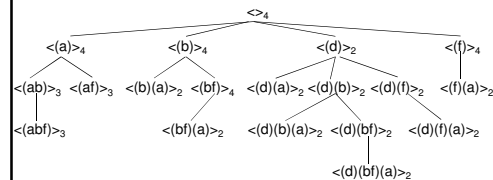
a occurs later than $\langle (d)(bf) \rangle$
in a data sequence

31

Sequential Rules

- Sequential rules created from $\langle(d)(bf)(a)\rangle$:

- $\langle \rangle \rightarrow \langle (d)(bf)(a) \rangle$ [sup.: 2, conf.: 2/4]
- $\langle (d) \rangle \rightarrow \langle (bf)(a) \rangle$ [sup.: 2, conf.: 2/2]
- $\langle (d)(bf) \rangle \rightarrow \langle (a) \rangle$ [sup.: 2, conf.: 2/2]

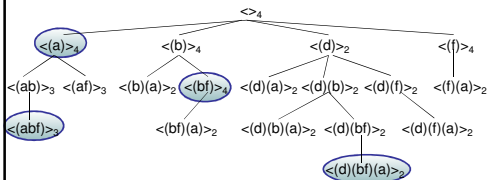


32

<i>Cld</i>	<i>Tid</i>	<i>Items</i>
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

Closed Sequential Patterns

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



33

<i>CId</i>	<i>Tid</i>	<i>Items</i>
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

References

- 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)

34