

PrefixSpan: Mining Sequential Patterns Efficiently by Prefix-Projected Pattern Growth

Authors:
Jian Pei, Jiawei Han, Behzad Mortazavi-Asi, Helen Pinto Qiming Chen,
Umeshwar Dayal, Mei-Chun Hsu

Presenter:
Wojciech Stach

Outline

- Mining Sequential Patterns
 - Problem statement
 - Definitions & examples
 - Strategies
- PrefixSpan algorithm
 - Motivation
 - Definitions & examples
 - Algorithm
 - Example
 - Performance study
- Conclusions

Sequential Pattern Mining

- Given
 - a set of sequences, where each sequence consists of a list of elements and each element consists of set of items
 - user-specified min_support threshold

Sequential Pattern Mining

- Find all the frequent subsequences, i.e. the subsequences whose occurrence frequency in the set of sequences is no less than min_support



id	Sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<(ef)(ab)(df)cb>
40	<eg(af)cbc>

<a(abc)(ac)d(cf)> - 5 elements, 9 items

<a(abc)(ac)d(cf)> - 9-sequence

<a(abc)(ac)d(cf)> = <a(cba)(ac)d(cf)>
<a(abc)(ac)d(cf)> ≠ <a(ac)(abc)d(cf)>



id	Sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<(ef)(ab)(df)cb>
40	<eg(af)cbc>

min_support = 2

Solution – 53 frequent subsequences

<a> <aa> <ab> <a(bc)> <a(bc)a> <aba> <abc>
<(ab)> <(ab)c> <(ab)d> <(ab)f> <(ab)dc> <ac>
<aca> <acb> <acc> <ad> <adc> <af>
 <ba> <bc> <(bc)> <(bc)a> <bd> <bdc> <bf>
<c> <ca> <cb> <cc>
<d> <db> <dc> <dcb>
<e> <ea> <eab> <eac> <each> <eb> <ebc> <ec>
<ecb> <ef> <efb> <efc> <efcb>
<f> <fb> <fbc> <fc> <fcb>

Subsequence vs. super sequence

- Given two sequences $\alpha = \langle a_1 a_2 \dots a_n \rangle$ and $\beta = \langle b_1 b_2 \dots b_m \rangle$
- α is called a **subsequence** of β , denoted as $\alpha \subseteq \beta$, if there exist integers $1 \leq j_1 < j_2 < \dots < j_n \leq m$ such that $a_1 \subseteq b_{j_1}, a_2 \subseteq b_{j_2}, \dots, a_n \subseteq b_{j_n}$
- β is a **super sequence** of α



$\beta = \langle a(abc)(ac)d(cf) \rangle$

$\alpha_1 = \langle aa(ac)d(c) \rangle$

$\alpha_2 = \langle (ac)(ac)d(cf) \rangle$

$\alpha_3 = \langle ac \rangle$

~~$\beta = \langle a(abc)(ac)d(cf) \rangle$~~

~~$\alpha_4 = \langle df(cf) \rangle$~~

~~$\alpha_5 = \langle (cf)d \rangle$~~

~~$\alpha_6 = \langle (abc)d(cf) \rangle$~~

5

Sequence Support Count

- A **sequence database** is a set of tuples $\langle sid, s \rangle$
- A tuple $\langle sid, s \rangle$ is said to **contain** a sequence α , if α is a subsequence of s , i.e., $\alpha \subseteq s$
- The **support of a sequence** α is the number of tuples containing α



id	Sequence
10	$\langle a(abc)(ac)d(cf) \rangle$
20	$\langle (ad)c(bc)(ae) \rangle$
30	$\langle (ef)(ab)(df)cb \rangle$
40	$\langle eg(af)cbc \rangle$

$\alpha_1 = \langle a \rangle$

$\text{support}(\alpha_1) = 4$

$\alpha_2 = \langle ac \rangle$

$\text{support}(\alpha_2) = 4$

$\alpha_3 = \langle (ab)c \rangle$

$\text{support}(\alpha_3) = 2$

6

Strategies

- Apriori-property based
 - AprioriSome (1995)
 - AprioriAll (1995)
 - DynamicSome (1995)
 - GSP (1996)
- Regular expression constraints
 - SPiRiT (1999)
- Data projection based
 - FreeSpan (2000)

7

Outline

- Mining Sequential Patterns
 - Problem statement
 - Definitions & examples
 - Strategies
- PrefixSpan algorithm**
 - Motivation
 - Definitions & examples
 - Algorithm
 - Example
 - Performance study
- Conclusions

8

Motivation and Background

- Shortcomings of Apriori-like approaches
 - Potentially huge set of candidate sequences
 - Multiple scans of databases
 - Difficulties at mining long sequential patterns
- FreeSpan (**F**requent pattern-projected **S**equential **p**attern mining) – pattern growth method
 - General idea is to use frequent items to recursively project sequence databases into a smaller projected databases and grow subsequence fragments in each projected database
- PrefixSpan (**P**refix-projected **S**equential **p**attern mining)
 - Less projections and quickly shrinking sequences

9

Prefix

- Given two sequences $\alpha = \langle a_1 a_2 \dots a_n \rangle$ and $\beta = \langle b_1 b_2 \dots b_m \rangle$, $m \leq n$
- Sequence β is called a **prefix** of α if and only if:
 - $b_i = a_i$ for $i \leq m-1$;
 - $b_m \subseteq a_m$;
 - All the items in $(a_m - b_m)$ are alphabetically after those in b_m



$\alpha = \langle a(abc)(ac)d(cf) \rangle$

$\beta = \langle a(abc)a \rangle$

~~$\alpha = \langle a(abc)(ac)d(cf) \rangle$~~

~~$\beta = \langle a(abc)c \rangle$~~

10

Projection

- Given sequences α and β , such that β is a subsequence of α .
- A subsequence α' of sequence α is called a **projection** of α w.r.t. β prefix if and only if
 - α' has prefix β ;
 - There exist no proper super-sequence α'' of α' such that α'' is a subsequence of α and also has prefix β



$\alpha = \langle a(abc)(ac)d(cf) \rangle$

$\beta = \langle (bc)a \rangle$

$\alpha' = \langle (bc)(ac)d(cf) \rangle$

11

Postfix

- Let $\alpha' = \langle a_1 a_2 \dots a_n \rangle$ be the projection of α w.r.t. prefix $\beta = \langle a_1 a_2 \dots a_{m-1} a'_m \rangle$ ($m \leq n$)
- Sequence $\gamma = \langle a''_m a_{m+1} \dots a_n \rangle$ is called the **postfix** of α w.r.t. prefix β , denoted as $\gamma = \alpha / \beta$, where $a''_m = (a_m - a'_m)$
- We also denote $\alpha = \beta \cdot \gamma$



$\alpha' = \langle a(abc)(ac)d(cf) \rangle$

$\beta = \langle a(abc)a \rangle$

$\gamma = \langle _c \rangle d(cf) \rangle$

12

PrefixSpan – Algorithm

- Input:** A sequence database S, and the minimum support threshold min_sup
- Output:** The complete set of sequential patterns
- Method:** Call PrefixSpan(<>, 0, S)
- Subroutine** PrefixSpan($\alpha, l, S|_{\alpha}$)
- Parameters:**
 - α : sequential pattern,
 - l: the length of α ;
 - $S|_{\alpha}$: the α -projected database, if $\alpha \neq \langle \rangle$; otherwise; the sequence database S.

PrefixSpan – Algorithm (2)

- Method**
 - Scan $S|_{\alpha}$ once, find the set of frequent items b such that:
 - b can be assembled to the last element of α to form a sequential pattern; or
 - $\langle b \rangle$ can be appended to α to form a sequential pattern.
 - For each frequent item b, append it to α to form a sequential pattern α' , and output α' ;
 - For each α' , construct α' -projected database $S|_{\alpha'}$, and call PrefixSpan($\alpha', l+1, S|_{\alpha'}$).

PrefixSpan - Example

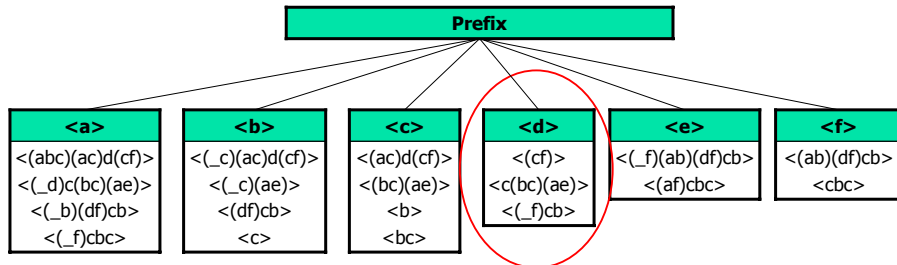
- Find length-1 sequential patterns

<a>		<c>	<d>	<e>	<f>	<g>
4	4	4	3	3	3	1

id	Sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<(ef)(ab)(df)cb>
40	<eg(af)cbc>

min_support = 2

- Divide search space

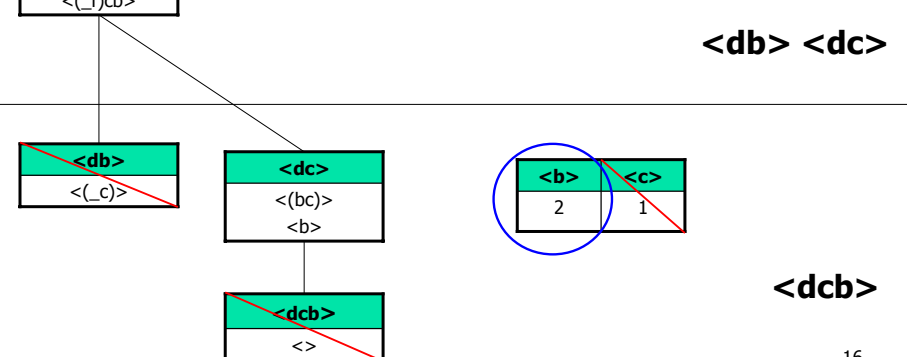


PrefixSpan – Example (2)

- Find subsets of sequential patterns

<d>	<a>		<c>	<d>	<e>	<(_e)>	<f>	<(_f)>
<(cf)> <c(bc)(ae)> <(_f)cb>	1	2	3	0	1	0	1	1

<db> <dc>



<dcb>

PrefixSpan - characteristics

- No candidate sequence needs to be generated by PrefixSpan
- Projected databases keep shrinking
- The major cost of PrefixSpan is the construction of projected databases



How to reduce this cost?



Different projection methods

- Bi-level projection**
 - reduces the number and the size of projected databases
- Pseudo-Projection**
 - reduces the cost of projection when projected database can be held in main memory

17

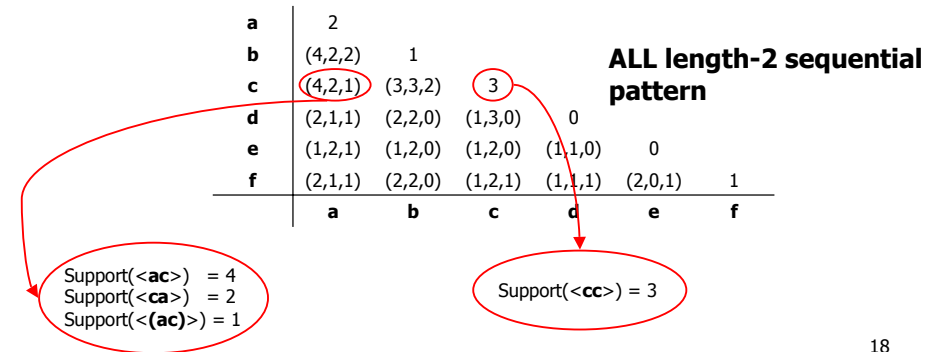
Bi-level Projection



id	Sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<(ef)(ab)(df)cb>
40	<eg(af)cbc>

min_support = 2

- Scan to get 1-length sequences
- Construct a **triangular matrix** instead of **projected databases** for each length-1 patterns



18

Bi-level projection (2)

- For each length-2 sequential pattern α , construct the α -projected database and find the frequent items
- Construct corresponding S-matrix



<ab>	a	b	c	(c)	d	(d)	e	(e)	f	(f)
<(_c)(ac)(cf)>	2	0	2	2	0	1	0	0	1	0
<(_c)a>										
<c>										

<aba> <abc> <a(bc)>

a	0		
c	(1,0,1)	1	
(c)	(0,2,0)	(0,1,0)	0
	a	c	(c)

<a(bc)a>

19

Bi-level projection (3) - optimization

- "Do we need to include every item in a postfix in the projected databases?"
- NO!** Item pruning in projected database by 3-way Apriori checking



<ac> is not frequent

Any super-sequence of it can never be a sequential pattern

c can be excluded from construction of <ab> - projected database

<a(bd)> is not frequent

To construct <a(bc)>-projected database, sequence <a(bcde)df> should be projected to <(_e)df> instead of <(_de)df>

20

Pseudo-Projection

- **Observation:** postfixes of a sequence often appear repeatedly in recursive projected databases
- **Method:** instead of constructing *physical* projection by collecting all the postfixes, we can use pointers referring to the sequences in the database as a pseudo-projection
- Every projection consists of two pieces of information: **pointer** to the sequence in database and **offset** to the postfix in the sequence

s1=<a(abc)(ac)d(cf)>

Pointer	Offset	Postfix
s1	2	<(abc)(ac)d(cf)>
s1	5	<(ac)d(cf)>
s1	6	<(_c)d(cf)>

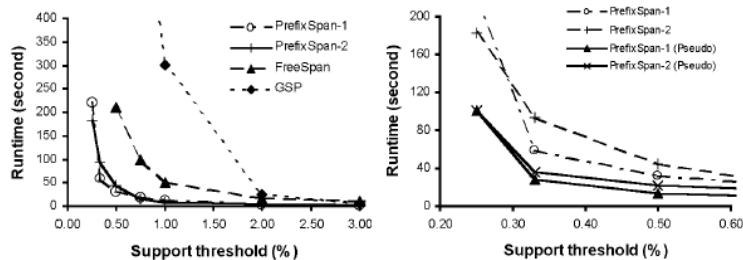
21

Experimental Results

- Environment: 233MHz Pentium PC, 128 MB RAM, Windows NT, Visual C++ 6.0
- Reported test on synthetic data set: C10T8S8I8:
 - 1000 items
 - 10000 sequences
 - Average number of items within elements: 8
 - Average number of elements in a sequence: 8
- Competitors:
 - GSP
 - FreeSpan
 - PrefixSpan-1 (level-by-level projection)
 - PrefixSpan-2 (bi-level projection)

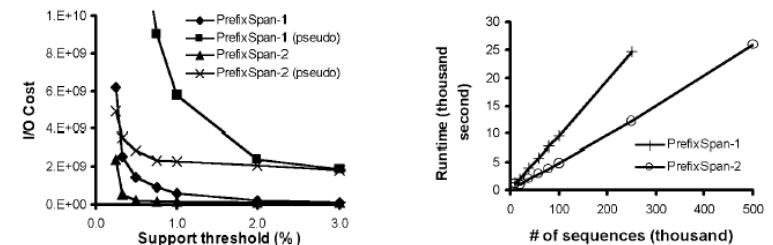
22

Runtime vs. support threshold



23

I/O costs vs. threshold and scalability



24



Outline

- Mining Sequential Patterns
 - Problem statement
 - Definitions & examples
 - Strategies
- PrefixSpan algorithm
 - Motivation
 - Definitions & examples
 - Algorithm
 - Example
 - Performance study
- **Conclusions**

25



Conclusions

- PrefixSpan
 - Efficient pattern growth method
 - Outperforms both GSP and FreeSpan
 - Explores prefix-projection in sequential pattern mining
 - Mines the complete set of patterns but reduces the effort of candidate subsequence generation
 - Prefix-projection reduces the size of projected database and leads to efficient processing
 - Bi-level projection and pseudo-projection may improve mining efficiency

26



References

- Pei J., Han J., Mortazavi-Asl J., Pinto H., Chen Q., Dayal U., Hsu M., PrefixSpan: Mining Sequential Patterns Efficiently by Prefix-Projected Pattern Growth, 17th International Conference on Data Engineering (ICDE), April 2001
- Agrawal R., Srikant R., Mining sequential patterns, Proceedings 1995 Int. Conf. Very Large Data Bases (VLDB'94), pp. 487-499, 1995
- Han J., Dong G., Mortazavi-Asl B., Chen Q., Dayal U., Hsu M.-C., Freespan: Frequent pattern-projected sequential pattern mining, Proceedings 2000 Int. Conf. Knowledge Discovery and Data Mining (KDD'00), pp. 355-359, 2000
- Srikant R., Agrawal R., Mining sequential pattern: Generalizations and performance improvements, Proceedings 5th Int. /conf. Extending Database Technology (EDBT'96), pp. 3-17, 1996
- Zhao Q., Bhowmick S. S., Sequential Pattern Mining: A Survey. Technical Report Center for Advanced Information Systems, School of Computer Engineering, Nanyang Technological University, Singapore, 2003

27



THANK YOU !!!

Any Questions?

28