

1



Agenda

- What is Sequence Data?
- Sequential Pattern Mining
 - GSP
 - Prefixspan



2

2

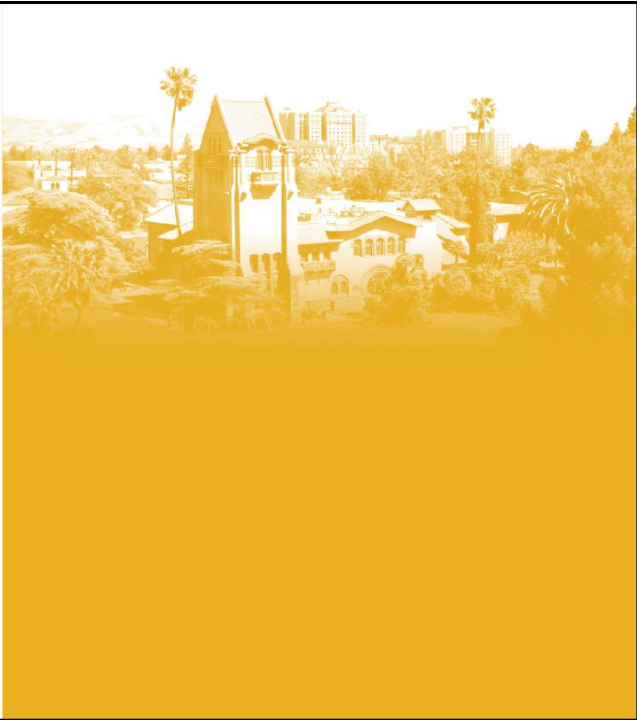
Data Mining / Machine Learning Methods

	Vector Data	Set Data	Sequence Data / Time Series Data	Graph Data
Classification	Logistic Regression			
Clustering	K-means Hierarchical Clustering Gaussian Mixture Models			
Prediction / Regression	Linear Regression Generalized Linear Models			
Frequent Pattern Mining		Apriori	GSP; PrefixSpan	
Similarity Search				

3

3

Basics of Sequence Data



4

What is Sequence Data?

- A sequence dataset consists of ordered elements or events
- Recorded with or without notion of time
- Transaction dataset vs. sequence dataset

Transaction Dataset

TID	itemsets
10	a, b, d
20	a, c, d
30	a, d, e
40	b, e, f

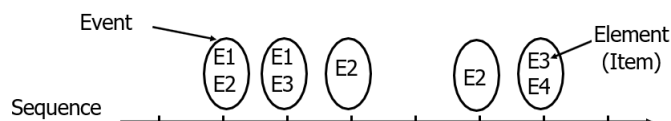
Sequence Dataset

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

5

Sequence Data

Sequence Database	Sequence	Event	Element
Customer	Purchase history of a given customer	A set of items bought by a customer at time t	Books, diary products, CDs, etc
Web Data	Browsing activity of a particular Web visitor	A collection of files viewed by a Web visitor after a single mouse click	Home page, index page, contact info, etc
Event data	History of events generated by a given sensor	Events triggered by a sensor at time t	Types of alarms generated by sensors
Genome sequences	DNA sequence of a particular species	An element of the DNA sequence	Bases A, T, G, C



6

Examples of Sequence Data

Music Sequence

DNA Sequence

```

MD106 ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
NEWC  ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
W501  ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
MD199 ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
C1674 ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
SIM4  ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG

MD106 CTACGGCCTAATGGTGCTAACGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
NEWC  CTACGGCCTAATGGTGCTAACGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
W501  CTACGGCCTAATGGTGCTAACGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
MD199 CTACGGCCTAATGGTGCTAACGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
C1674 CTACGGCCTAATGGTGCTAACGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
SIM4  CTACGGCCTAATGGTGCTAACGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT

MD106 CCGTTTCAAGTACCAAACCTGAGTGC GGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
NEWC  CCGTTTCAAGTACCAAACCTGAGTGC GGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
W501  CCGTTTCAAGTACCAAACCTGAGTGC GGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
MD199 CCGTTTCAAGTACCAAACCTGAGTGC GGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
C1674 CCGTTTCAAGTACCAAACCTGAGTGC GGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
SIM4  CCGTTTCAAGTACCAAACCTGAGTGC GGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG

MD106 CTGCAGGAGGCGTCCACCACCAAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
NEWC  CTGCAGGAGGCGTCCACCACCAAGTGCCCCAATCTACAGGTCATCGGCCGAGAAATAG
W501  CTGCAGGAGGCGTCCACCACCAAGTGCCCCAATCTACAGGTCATCGGCCGAGAAATAG
MD199 CTGCAGGAGGCGTCCACCACCAAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
C1674 CTGCAGGAGGCGTCCACCACCAAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
SIM4  CTGCAGGAGGCGTCCACCACCAAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG

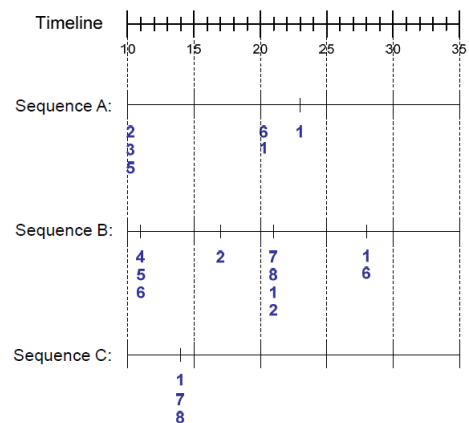
```

7

Sequence Data with Timestamp

Sequence Database

Sequence ID	Timestamp	Events
A	10	2, 3, 5
A	20	6, 1
A	23	1
B	11	4, 5, 6
B	17	2
B	21	7, 8, 1, 2
B	28	1, 6
C	14	1, 8, 7



8

Definitions of Sequences, Events...

- A sequence is an ordered list of events, denoted $\langle e_1 e_2 \dots e_l \rangle$
- An event is a *non-empty* collection of items/elements: $e_1 = (abc)$
- 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

- Examples

$\alpha = \langle (ab)d \rangle$ and $\beta = \langle (abc)(de) \rangle$

$\alpha = \langle a(bc)dc \rangle$ and $\beta = \langle a(abc)(ac)d(cf) \rangle$

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

9

9

Subsequence Exercises

Data Sequence	Subsequence	Contain?
$\langle (2\ 4)(3\ 5\ 6)(8) \rangle$	$\langle (2)(8) \rangle$	
$\langle (1\ 2)(3\ 4) \rangle$	$\langle (1)(2) \rangle$	
$\langle (2\ 4)(2\ 4)(2\ 5) \rangle$	$\langle (2)(4) \rangle$	
$\langle (2\ 4)(2\ 5)(4\ 5) \rangle$	$\langle (2)(4)(5) \rangle$	
$\langle (2\ 4)(2\ 5)(4\ 5) \rangle$	$\langle (2)(5)(5) \rangle$	
$\langle (2\ 4)(2\ 5)(4\ 5) \rangle$	$\langle (2\ 4\ 5) \rangle$	

10

10

Definitions of Sequences, Events...

- Length of a sequence, $|s|$, is given by # of elements in the sequence
e.g. What is the length of $\langle (ef)(ab)(df)cb \rangle$?
- A k -sequence is a sequence that contains k events:
e.g. $\langle (ab)(a) \rangle$ is a 2-sequence and has a length of 3

11

11

Support of a Sequence

- Support of a Sequence α : # of sequences in a sequence dataset D that contains a sequence α on and is denoted as $\text{sup}(\alpha)$:

$$\text{sup}(\alpha) = |\{S \mid S \in D \text{ and } \alpha \sqsubseteq S\}|$$

- α is frequent if $\text{sup}(\alpha) \geq \sigma$ (*minsup* threshold)
- A frequent sequence is called sequential pattern

12

12

Example: Support of a Sequence

Consider the sequence dataset as shown:

- Example 1: $\text{sup}(\langle a \rangle) = ?$
- Example 2: $\text{sup}(\langle b \rangle) = ?$
- Example 3: $\text{sup}(\langle a \rangle \langle b \rangle) = ?$
- Example 4: $\text{sup}(\langle a \ b \rangle) = ?$

SID	Sequence
10	$\langle (ab)(c)(a) \rangle$
20	$\langle (a)(b)(c) \rangle$
30	$\langle (b)(c)(d) \rangle$
40	$\langle (b)(ab)(c) \rangle$

$\text{sup}(\langle a \rangle) = 3$ (75%)
 $\text{sup}(\langle b \rangle) = 4$ (100%)
 $\text{sup}(\langle c \rangle) = 4$ (100%)
 $\text{sup}(\langle a \rangle \langle c \rangle) = 3$ (75%)
 $\text{sup}(\langle ab \rangle) = 2$ (50%)
 $\text{sup}(\langle b \rangle \langle c \rangle) = 4$ (100%)
 $\text{sup}(\langle ab \rangle \langle c \rangle) = 2$ (50%)

SID	Sequence
10	$\langle (ab)(c)(a) \rangle$
20	$\langle (a)(b)(c) \rangle$
30	$\langle (b)(c)(d) \rangle$
40	$\langle (b)(ab)(c) \rangle$

SID	Sequence
10	$\langle (ab)(c)(a) \rangle$
20	$\langle (a)(b)(c) \rangle$
30	$\langle (b)(c)(d) \rangle$
40	$\langle (b)(ab)(c) \rangle$

SID	Sequence
10	$\langle (ab)(c)(a) \rangle$
20	$\langle (a)(b)(c) \rangle$
30	$\langle (b)(c)(d) \rangle$
40	$\langle (b)(ab)(c) \rangle$

SID	Sequence
10	$\langle (ab)(c)(a) \rangle$
20	$\langle (a)(b)(c) \rangle$
30	$\langle (b)(c)(d) \rangle$
40	$\langle (b)(ab)(c) \rangle$

13

13



14

What Is Sequential Pattern Mining?

- Given a set of sequences and support threshold, find the complete set of **frequent subsequences**.

Sequence Data

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

Given support threshold $min_sup = 2$, <(ab)c> is a sequential pattern

15

15

Applications

Applications of Sequential Pattern Mining

- Customer shopping sequences:
 - First buy computer, then CD-ROM, and then digital camera, within 3 months.
- Medical treatments, natural disasters (e.g., earthquakes), science & eng. processes, stocks and markets, etc.
- Telephone calling patterns
- Weblog click streams
- DNA sequences and gene structures

16

16

Example: Sequential Pattern Mining

Sequence ID	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1,2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3,4
D	3	4,5
E	1	1,3
E	2	2,4,5

Frequent Subsequences:

$\text{sup}\langle(1,2)\rangle = 60\%$
 $\text{sup}\langle(2,3)\rangle = 60\%$
 $\text{sup}\langle(2,4)\rangle = 80\%$
 $\text{sup}\langle(3)(5)\rangle = 80\%$
 $\text{sup}\langle(1)(2)\rangle = 80\%$
 $\text{sup}\langle(2)(2)\rangle = 60\%$
 $\text{sup}\langle(1)(2,3)\rangle = 60\%$
 $\text{sup}\langle(2)(2,3)\rangle = 60\%$
 $\text{sup}\langle(1,2)(2,3)\rangle = 60\%$

17

17

Extracting Sequential Patterns

- Given n items: $i_1, i_2, i_3, \dots, i_n$

- Candidate 1-subsequences:

$\langle(i_1)\rangle, \langle(i_2)\rangle, \langle(i_3)\rangle, \dots, \langle(i_n)\rangle$

- Candidate 2-subsequences:

$\langle(i_1, i_2)\rangle, \langle(i_1, i_3)\rangle, \dots,$

$\langle(i_1)(i_1)\rangle, \langle(i_1)(i_2)\rangle, \dots, \langle(i_n)(i_n)\rangle$

- Candidate 3-subsequences:

$\langle(i_1, i_2, i_3)\rangle, \langle(i_1, i_2, i_4)\rangle, \dots,$

$\langle(i_1, i_2)(i_1)\rangle, \langle(i_1, i_2)(i_2)\rangle, \dots,$

$\langle(i_1)(i_1, i_2)\rangle, \langle(i_1)(i_1, i_3)\rangle, \dots,$

$\langle(i_1)(i_1)(i_1)\rangle, \langle(i_1)(i_1)(i_2)\rangle, \dots$

18

18

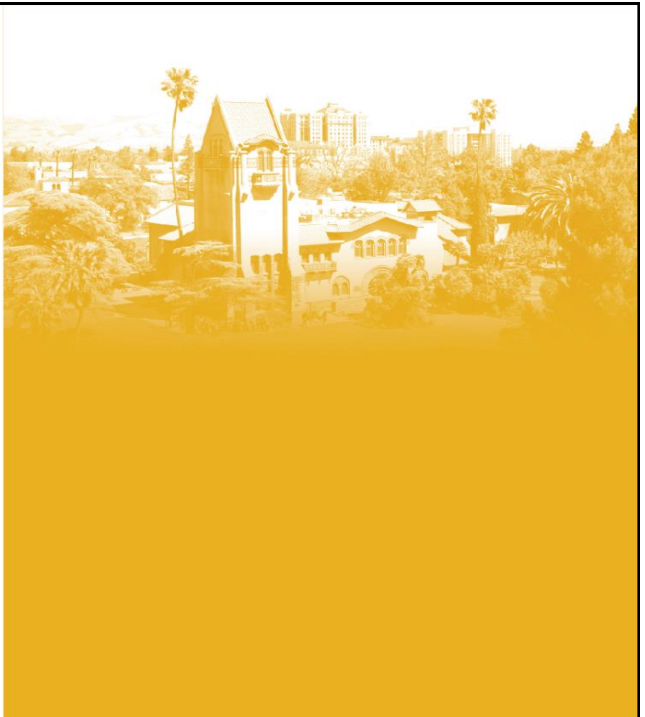
Challenges on Sequential Pattern Mining

- A huge number of possible sequential patterns are hidden in dataset
- A mining algorithm should
 - find the complete set of patterns, when possible, satisfying the minimum support (frequency) threshold
 - be highly efficient, scalable, involving only a small number of database scans
 - be able to incorporate various kinds of user-specific constraints

19

19

Generalized Sequential Pattern (GSP)



20

The Apriori Property of Sequential Patterns

A basic property: **Apriori** (Agrawal & Srikant'94)

- If a sequence S is not frequent, then none of the super-sequences of S is frequent.

Example: $\langle abc \rangle$ is infrequent $\rightarrow \langle abcde \rangle$ and $\langle (abc)(de) \rangle$ are infrequent too!

21

21

Methods for Sequential Pattern Mining

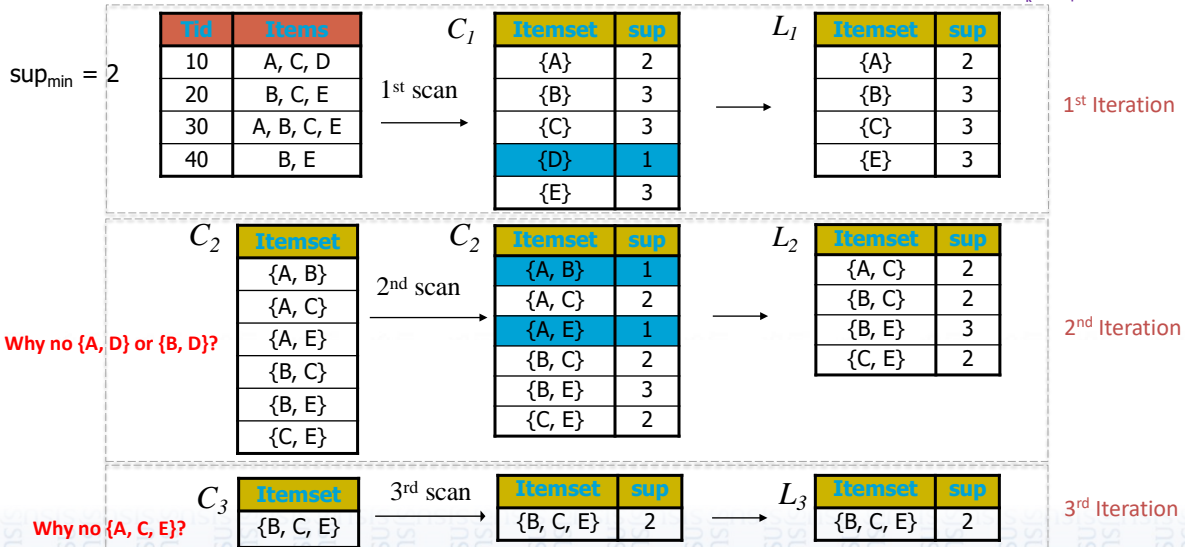
- Apriori-based Approaches
 - GSP
 - SPADE
- Pattern-Growth-based Approaches
 - FreeSpan
 - PrefixSpan

22

22

Recall The Apriori Algorithm

C_k : Candidate itemsets of size k
 L_k : Frequent itemsets of size k



23

23

Generalized Sequential Pattern (GSP) Mining

- Step 1
 - Make the first pass over the sequence database D to yield all candidate of length-1
- Step 2 (Repeat until no new frequent sequences are found)
 - Candidate Generation:
 - Merge pairs of frequent subsequences found in the (k-1)th pass to generate candidate sequences that contain k items
 - Candidate Pruning:
 - Prune candidate length k sequences that contain infrequent length (k-1) subsequences
 - Support Counting:
 - Make a new pass over the sequence database D to find the support for these candidate sequences
 - Candidate Elimination:
 - Eliminate candidate length k sequences whose actual support is less than minsup

24

24

Example: Finding Length-1 Sequential Patterns

- Initial candidates: all singleton sequences: <a>, , <c>, <d>, <e>, <f>, <g>, <h>
- Scan dataset once, count support for candidate.

$min_sup = 2$

SID	sequence
10	<(bd)cb(ac)>
20	<(bf)(ce)b(fg)>
30	<(ah)(bf)abf>
40	<(be)(ce)d>
50	<a(bd)bcb(ade)>



Cand	Sup
<a>	3
	5
<c>	4
<d>	3
<e>	3
<f>	2
<g>	1
<h>	1

25

25

Example: Finding Length-2 Sequential Patterns

- Now, generate length-2 candidates.

	<a>		<c>	<d>	<e>	<f>
<a>	<aa>	<ab>	<ac>	<ad>	<ae>	<af>
	<ba>	<bb>	<bc>	<bd>	<be>	<bf>
<c>	<ca>	<cb>	<cc>	<cd>	<ce>	<cf>
<d>	<da>	<db>	<dc>	<dd>	<de>	<df>
<e>	<ea>	<eb>	<ec>	<ed>	<ee>	<ef>
<f>	<fa>	<fb>	<fc>	<fd>	<fe>	<ff>

	<a>		<c>	<d>	<e>	<f>
<a>		<(ab)>	<(ac)>	<(ad)>	<(ae)>	<(af)>
			<(bc)>	<(bd)>	<(be)>	<(bf)>
<c>				<(cd)>	<(ce)>	<(cf)>
<d>					<(de)>	<(df)>
<e>						<(ef)>
<f>						

- 51 length-2 candidates!!!
- Without Apriori property, $8*8+8*7/2 = 92$ candidates

➔ Apriori prunes 44.6% candidates!

26

26

Generating Candidates in General

From L_{k-1} to C_k

- Step 1: Join

- Sequences α and β can **join**, if dropping first item in α is the same as dropping the last item in β

Examples:

$\langle(12)3\rangle$ join $\langle(2)34\rangle \rightarrow \langle(12)34\rangle$

$\langle(12)3\rangle$ join $\langle(2)(34)\rangle \rightarrow \langle(12)(34)\rangle$

- Step 2: Pruning

- Check whether all length $k-1$ subsequences of a candidate is contained in L_{k-1}

27

27

The GSP Mining Process

5th scan: 1 cand. 1 length-5 seq. pat.

$\langle(bd)cba\rangle$

Cand. cannot pass sup. threshold

4th scan: 8 cand. 7 length-4 seq. pat.

$\langle abba \rangle \langle (bd)bc \rangle \dots$

Cand. not in DB at all

3rd scan: 46 cand. 20 length-3 seq. pat.
(20 cand. not in DB at all)

$\langle abb \rangle \langle aab \rangle \langle aba \rangle \langle baa \rangle \langle bab \rangle \dots$

2nd scan: 51 cand. 19 length-2 seq. pat.
(10 cand. not in DB at all)

$\langle aa \rangle \langle ab \rangle \dots \langle af \rangle \langle ba \rangle \langle bb \rangle \dots \langle ff \rangle \langle (ab) \rangle \dots \langle (ef) \rangle$

1st scan: 8 cand. 6 length-1 seq. pat.

$\langle a \rangle \langle b \rangle \langle c \rangle \langle d \rangle \langle e \rangle \langle f \rangle \langle g \rangle \langle h \rangle$

28

28

Drawbacks of GSP

- A huge set of candidate sequences generated (especially 2-item candidate sequence)
- Multiple scans of database needed. The length of each candidate grows by one at each database scan.
- Inefficient for mining long sequential patterns.
 - A long pattern grow up from short patterns
 - # of short patterns is exponential to the length of mined patterns.

29

29

Bottlenecks of GSP

- A huge set of candidates could be generated
 - 1,000 frequent length-1 sequences generates huge number of length-2 candidates!
- Multiple scans of database in mining
- Breadth-first search
- Mining long sequential patterns
 - Needs an exponential number of short candidates
 - A length-100 sequential pattern needs $\sim 10^{30}$ candidate sequences!

$$1000 \times 1000 + \frac{1000 \times 999}{2} = 1,499,500$$

$$\sum_{i=1}^{100} \binom{100}{i} = 2^{100} - 1 \approx 10^{30}$$

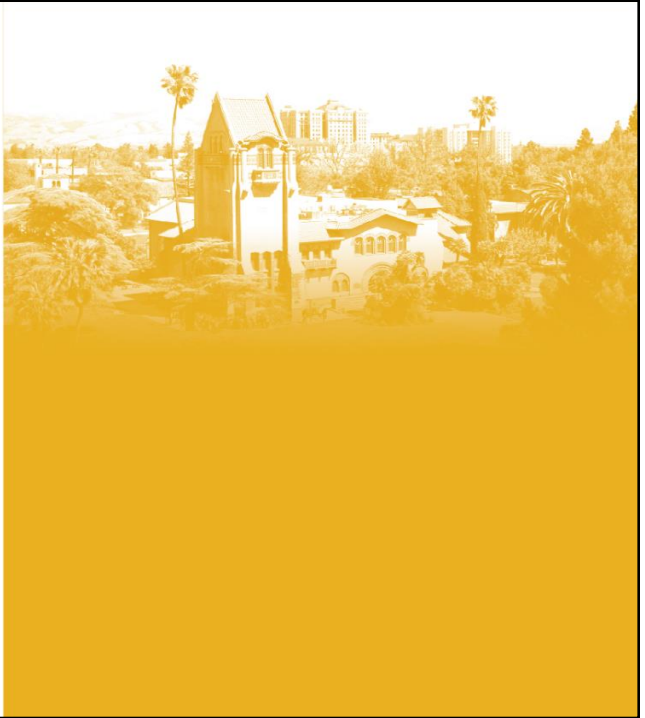
- Some of these can be addressed with PrefixSpan

30

30



Sequential Pattern Mining with PrefixSpan



31



Prefix and Suffix

- Assume a pre-specified order on items, e.g., alphabetical order
 $\langle a \rangle$, $\langle aa \rangle$, $\langle a(ab) \rangle$, $\langle a(abc) \rangle$ are prefixes of sequence $\langle a(abc)(ac)d(cf) \rangle$
 Note: $\langle a(ac) \rangle$ is not a prefix of $\langle a(abc)(ac)d(cf) \rangle$

- Given sequence $\langle a(abc)(ac)d(cf) \rangle$:

Prefix	Suffix
$\langle a \rangle$	$\langle (abc)(ac)d(cf) \rangle$
$\langle aa \rangle$	$\langle (_bc)(ac)d(cf) \rangle$
$\langle a(ab) \rangle$	$\langle (_c)(ac)d(cf) \rangle$

- $(_bc)$ means: the last element in the prefix together with (bc) form one element

32

32

Prefix-based Projection

Given a sequence, α , let α' be subsequence of α :

- α' is called a **projection** of α with respect to prefix β , if and only if
 - α' has prefix β , and
 - α' is the maximum subsequence of α with prefix β

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

Example:

<(a)bc)(ac)d(cf)> is a projection of <a(abc)(ac)d(cf)>
w.r.t. prefix <a>

<a(d)cf)> is a projection of <a(abc)(a)c)d(cf)>
w.r.t. prefix <ad>

33

33

Projected (Suffix) Database

- Let β be a sequential pattern, β -projected database is the collection of suffixes of projections of sequences in the database w.r.t. prefix β

- Examples:

- <a>-projected database

<(abc)(ac)d(cf)>

<(_d)c(bc)(ae)>

<(_b)(df)cb>

<(_f)cbc>

- <ab>-projected database

<(_c)(ac)d(cf)>

<(_c)(ae)>

<c>

(<a(b)c)(ac)d(cf)> is the projection of <a(abc)(ac)d(cf)> w.r.t. prefix <ab>)

(<a(b)c)(ae)> is the projection of <(ad)c(bc)(ae)> w.r.t. prefix <ab>)

(<abc> is the projection of <eg(af)cbc> w.r.t. prefix <ab>)

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

34

34

Examples

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

Starting with <a>-projected database

- <(abc)(ac)d(cf)>
- <(_d)c(bc)(ae)>
- <(_b)(df)cb>
- <(_f)cbc>

<aa>-projected database

- <(_bc)(ac)d(cf)>
- <(_e)>

<ab>-projected database

- <(_c)(ac)d(cf)>
- <(_c)(ae)>
- <c>

<(ab)>-projected database

- <(_c)(ac)d(cf)>
- <(df)cb>

35

35

Mining Sequential Patterns by Prefix Projections

- Step 1: Find length-1 sequential patterns.
<a>, , <c>, <d>, <e>, <f>
- Step 2: Divide search space. The complete set of sequential patterns can be partitioned into 6 subsets:
The ones having prefix <a>
The ones having prefix
...
The ones having prefix <f>
- Step 3: Mine each subset recursively via corresponding projected databases.

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

36

36

Finding Sequential Patterns with Prefix <a>

- Consider only projections w.r.t. <a>:

- <a>-projected (suffix) database:

- <(abc)(ac)d(cf)>

- <(_d)c(bc)(ae)>

- <(_b)(df)cb>

- <(_f)cbc>

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

- Find all the length-2 sequential patterns having prefix <a>:

<aa>, <ab>, <(ab)>, <ac>, <ad>, <af>

- Further partition into 6 subsets:

- Having prefix <aa>

...

Why are those 6 subsets?

- Having prefix <af>

37

37

Finding Sequential Patterns with Prefix <a>

- By scanning the <a>-projected database once, its locally frequent items are identified as

a : 2, b : 4, _b : 2, c : 4, d : 2, and f : 2.

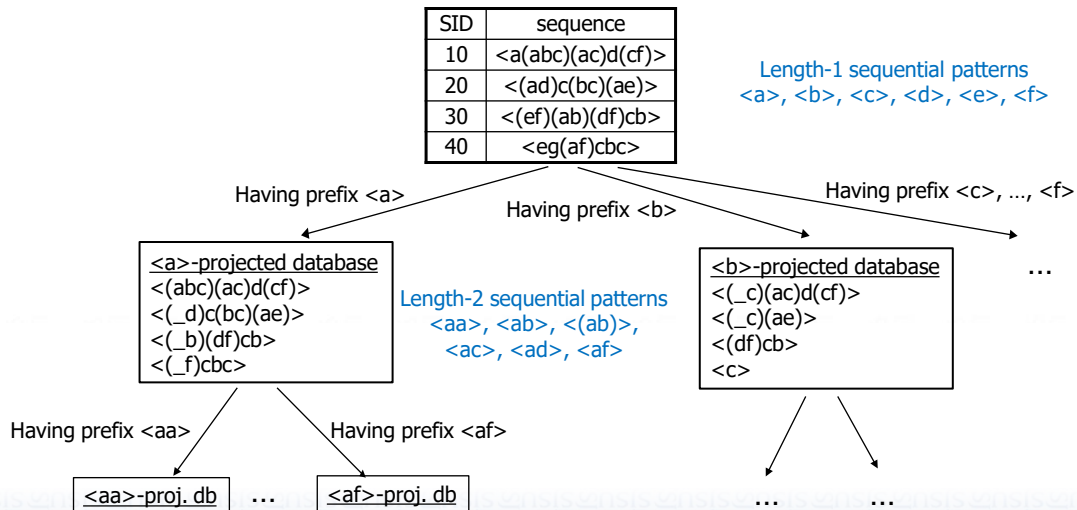
- Thus, all the length-2 sequential patterns prefixed with <a> are found:

<aa>: 2, <ab>: 4, <(ab)>: 2, <ac> : 4, <ad>: 2, and <af>: 2.

38

38

Completeness of PrefixSpan



39

39

Projected Databases & Sequential Patterns

SID	Sequence		prefix	projected (suffix) database	sequential patterns
10	<a(abc)(ac)d(cf)>	➔	<a>	<(abc)(ac)d(cf)>, <(_d)c(bc)(ae)>, <(_b)(df)cb>, <(_f)cbc>	<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>
20	<(ad)c(bc)(ae)>			<(_c)(ac)d(cf)>, <(_c)(ae)>, <(df)cb>, <c>	, <ba>, <bc>, <(bc)>, <(bc)a>, <bd>, <bdc>, <bf>
30	<(ef)(ab)(df)cb>		<c>	<(ac)d(cf)>, <(bc)(ae)>, , <bc>	<c>, <ca>, <cb>, <cc>
40	<eg(af)cbc>		<d>	<(cf)>, <c(bc)(ae)>, <(_f)cb>	<d>, <db>, <dc>, <dcb>
			<e>	<(_f)(ab)(df)cb>, <(af)cbc>	<e>, <ea>, <eab>, <eac>, <eacb>, <eb>, <ebc>, <ec>, <ecb>, <ef>, <efb>, <efc>, <efcb>
			<f>	<(ab)(df)cb>, <cbc>	<f>, <fb>, <fbc>, <fc>, <fcb>

40

40

Efficiency of PrefixSpan

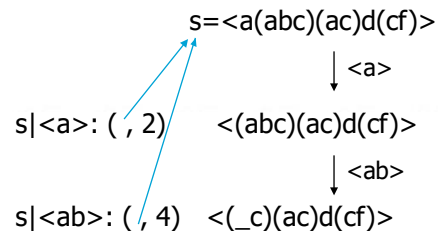
- No candidate sequence needs to be generated
- Projected databases keep shrinking
- Major cost of PrefixSpan: Constructing projected databases
- Can be improved by **pseudo-projections** (virtual projection in memory)

41

41

Speed-up by Pseudo-Projection

- Major cost of PrefixSpan: projection
 - Postfixes of sequences often appear repeatedly in recursive projected databases
- When (projected) database can be held in main memory, use pointers to form projections:
 - Pointer to the sequence
 - Offset of the postfix



42

42

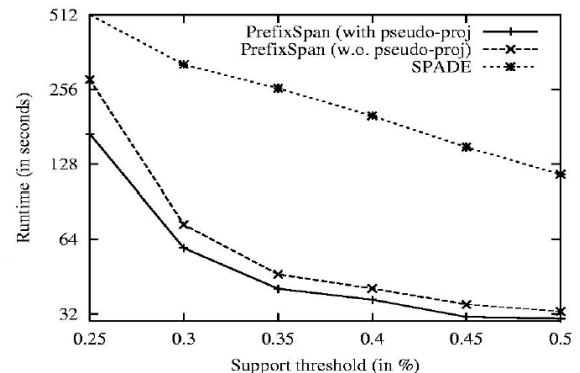
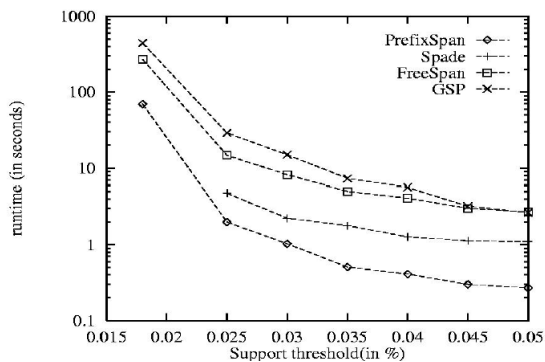
Pseudo-Projection vs. Physical Projection

- Pseudo-projection avoids physically copying postfixes
 - Efficient in running time and space when database can be held in main memory
- However, it is not efficient when database cannot fit in main memory
 - Disk-based random accessing is very costly
- Suggested Approach:
 - Integration of physical and pseudo-projection
 - Swapping to pseudo-projection when the data set fits in memory

43

43

Performance Comparison

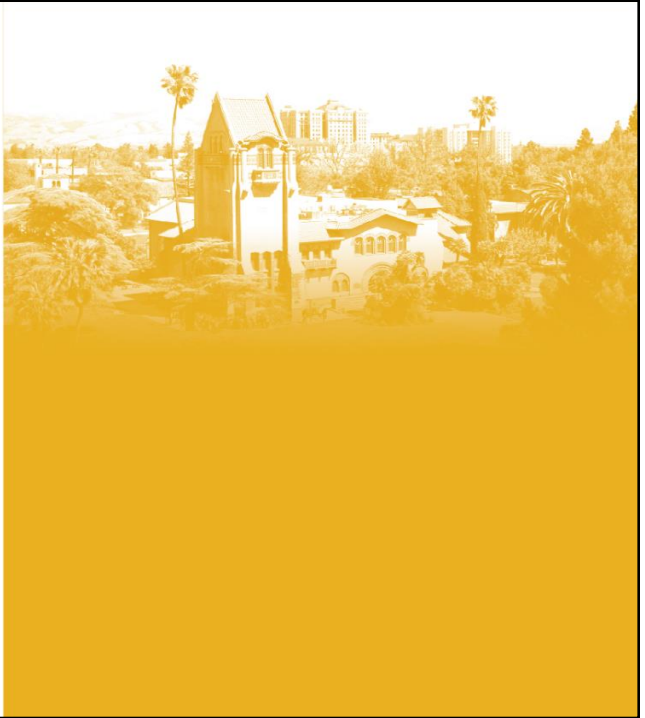


44

44



Sequential Rules



45



Sequential Pattern Implications

- $\langle af \rangle$ is a sequential pattern in the dataset:

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

- Does it imply that if someone buys $\langle a \rangle$, he'll also buy $\langle f \rangle$ later?
- Not really, $\langle a \rangle$ implies $\langle f \rangle$ only 50% of the time! ➔ Need sequential rules!

46

46

Sequential Rules

Two main types of sequential rules

- Standard Sequential rules
- Partially-ordered Sequential rules

47

47

Standard Sequential Rules

Standard Sequential Rules: Rules of the form $X \rightarrow Y$, where X and Y are sequential patterns.

Example: $\langle a(bc) \rangle \rightarrow \langle de \rangle$

- Two thresholds must be set by the user:
 - minimum support > 0
 - minimum confidence > 0
- **Support:** how many sequences contain the rule
- **Confidence:** how many sequences contain the rule divided by how many sequences contain its antecedent

48

48

Example: Standard Sequential Rule

- Consider the standard sequential rule: $\langle ab \rangle \rightarrow \langle f \rangle$

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

Support: 1 sequence (25%)

Confidence: $1 / 4 = 0.25$ (25%)

49

49

Order Matters?

- Notice that some sequential rules that are very similar but have only some small ordering variations.
- Example: $\langle ab \rangle \rightarrow \langle f \rangle$ and $\langle ba \rangle \rightarrow \langle f \rangle$

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



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

Rule	Support	Confidence
$\langle ab \rangle \rightarrow \langle f \rangle$	25%	25%
$\langle ba \rangle \rightarrow \langle f \rangle$	25%	50%

50

50

Partially-Ordered Sequential Rules

Partially-Ordered Sequential Rules: Rules of the form $X \rightarrow Y$, where X and Y are itemsets that are unordered.

Example: $\{a, b\} \rightarrow \{f\}$

- If we observe a and b (in any order), they will be followed by f

51

51

Partially-Ordered Sequential Rules

- This type of rule is often more interesting because it can summarize many standard sequential rules

- Example:

$\{Vivaldi\}, \{Mozart\}, \{Handel\} \Rightarrow \{Berlioz\}$
 $\{Mozart\}, \{Vivaldi\}, \{Handel\} \Rightarrow \{Berlioz\}$
 $\{Handel\}, \{Vivaldi\}, \{Mozart\} \Rightarrow \{Berlioz\}$
 $\{Handel, Vivaldi\}, \{Mozart\} \Rightarrow \{Berlioz\}$
 $\{Handel\}, \{Vivaldi, Mozart\} \Rightarrow \{Berlioz\}$
 $\{Handel, Vivaldi, Mozart\} \Rightarrow \{Berlioz\}$

Standard sequential rules



$\{Mozart, Vivaldi, Handel\} \Rightarrow \{Berlioz\}$

Partially-ordered sequential rules

52

52

Partially-Ordered Sequential Rules

- A **partially-ordered sequential rule** $X \rightarrow Y$ is a relationship between two disjoint and non-empty itemsets X, Y .
- A sequential rule $X \rightarrow Y$ has two properties:
 - **Support**: #of sequences where X occurs before Y , divided by the number of sequences.
 - **Confidence**: # of sequences where X occurs before Y , divided by the number of sequences where X occurs.
- Need to find all valid rules, rules with a support and confidence not less than user-defined thresholds **minSup** and **minConf**

53

53

Example: Mining Partially-Ordered Sequential Rules

- Consider the following sequence database with **minSup= 0.5** and **minConf= 0.5**

ID	Sequences
seq1	{a, b}, {c}, {f}, {g}, {e}
seq2	{a, d}, {c}, {b}, {a, b, e, f}
seq3	{a}, {b}, {f}, {e}
seq4	{b}, {f, g}



ID	Rule	Support	Confidence
r1	$\{a, b, c\} \Rightarrow \{e\}$	0.5	1.0
r2	$\{a\} \rightarrow \{c, e, f\}$	0.5	0.66
r3	$\{a, b\} \rightarrow \{e, f\}$	0.75	1.0
r4	$\{b\} \rightarrow \{e, f\}$	0.75	0.75
r5	$\{a\} \rightarrow \{e, f\}$	0.75	1.0
r6	$\{c\} \rightarrow \{f\}$	0.5	1.0
r7	$\{a\} \rightarrow \{b\}$	0.5	0.66
...

54

54

Summary

- Basic concepts: sequential patterns etc
- Scalable sequential pattern mining methods
 - GSP
 - PrefixSpan
- Sequential Rules: Standard and Partially-Ordered