

Sequence Data Mining

Young-Rae Cho, Ph.D.

Associate Professor

Division of Software / Division of Digital Healthcare

Yonsei University – Mirae Campus

❑ Single Sequence Mining

- Frequent string pattern mining
 - Finding substrings that frequently occur in a single sequence

ACTTCGATGGAGCCAGTCGCGAAATTCGACTAGATCG

❑ Sequence Dataset Mining

- Frequent string/sequence pattern mining
 - Finding substrings/sub-sequences that frequently occur among sequences
- Sequence data clustering
 - Grouping similar sequences
- Sequence data classification
 - Classifying a new sequence

id	sequence
1	ACTTCG
2	GGAGC
3	CGTCGAT
4	ATCGATCGC
5	TCGACT
6	AGATCGC

Applications



□ Examples

- Customer shopping sequential patterns
e.g., First buy a computer, then a CD-ROM, and then a printer within 3 months
- Stock market changes
- Web-log patterns
- Medical treatment records
- Gene or protein sequences
- Natural disaster records

□ Challenges

- Finding the complete set satisfying the minimum support threshold
- Developing efficient and scalable algorithms
- Incorporating various kinds of user-specific constraints

Frequent String Pattern Mining



□ Problem Definition

- Finding all substrings that occur frequently among multiple sequences
- Sequence: an ordered list of items
- Substring of a sequence x : a part of consecutive items from x
- Frequent string pattern ?
 - Suppose minimum support of 75%

```
T1= cctgatagacgctatctggctatccacgtacgtaggtcctctgtgCGaatctatgCGtttccaacat  
T2= agtactggtgtacatttgatacgtacgtacaccggcaacctgaaacaaacgctcagaaccagaagtgc  
T3= aaacgtacgtgcaccctctttcttcgTggctctggccaacgagggctgatgtataagacgaaaatttt  
T4= agcctccgatgtaagtcatagctgtaactattacctgccaccctattacatcttacgTgcgtataca  
T5= ctgttatacaacgCGtcatggcggggtatgCGttttggtcGtcgTtacgctcgatcgTtaacgtacgtc
```

- “acgtacgt” occurs in 4 out of 5 sequences (80%) – more frequent than minimum support

Properties of String Patterns



□ Properties

- Anti-monotonic property → Apriori algorithm
- If a string S is not frequent, then none of the super-strings of S are frequent
- Example: If $\langle a b \rangle$ is infrequent, $\langle a b c \rangle$ and $\langle d a b \rangle$ are infrequent too.

□ Process to Solve Frequent String Pattern Mining

- Iteratively increase length of substrings that occur more frequently than minimum support

□ Problems?

- Selective joining?
- Apriori pruning?
- Support computing?

Frequent Sequence Pattern Mining



□ Scope

- Transaction data → Sequential transaction data
- Frequent itemset patterns → (Frequent) Sequential patterns

□ Problem Definitions

- Sequence: an ordered list of items
 - e.g., $x = \langle a \ c \ f \ e \ c \ d \rangle$
- Sub-sequence of x : an ordered sequence of items from x
 - Not necessarily consecutive (different from a substring)
 - $\langle a \ e \ c \rangle \langle c \ f \ d \rangle \langle a \ c \rangle \langle c \ e \ c \rangle \langle a \ f \ c \ d \rangle \langle c \ f \ c \ e \rangle$
- Frequent sequence pattern ?
 - Suppose minimum support of 75%
 - $\langle a \ c \ d \rangle$ is one of the frequent sequence patterns
 - Any others?

SID	sequences
10	$\langle acadcf \rangle$
20	$\langle abcaed \rangle$
30	$\langle efadfcb \rangle$
40	$\langle egafcbd \rangle$

Frequent Sequence Pattern Minding – cont'



Extended Problem Definitions

- Consider “or” at a single position
- Sequence: an ordered list of elements and each element is a set of items,
 - e.g., $x = \langle (ab) c (df) (ef) d \rangle$
- Sub-sequence of x
 - $\langle (ab) d \rangle \langle a c (ef) \rangle \langle a f e \rangle \langle \del{(ab) e f} \rangle \langle b f d \rangle \langle c d f \rangle$
- Frequent sequence pattern
 - Suppose minimum support of 75%

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

- $\langle (ab) c \rangle$ is one of the frequent sequence patterns
- Any others?



Properties of Sequence Patterns

□ Properties

- Anti-monotonic property → Apriori algorithm
- If a sequence S is not frequent, then none of the super-sequences of S are frequent

□ Example

- If $\langle (ab)d \rangle$ is infrequent, so do $\langle (ab)de \rangle$ and $\langle (abc)d \rangle$ and $\langle a(ab)d \rangle$

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

min_sup = 75%



GSP (Generalized Sequential Pattern mining) Algorithm

□ Process

- 1) Initially, find all frequent length-1 sequences (= frequent 1-itemset)
- 2) Generate candidate length-(k+1) sequences from frequent length-k sequences
- 3) Count support for each candidate sequence to select frequent sequences
- 4) Repeat (2) and (3) until no frequent sequence or no candidate is found

□ Example

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

min_sup = 75%

□ Candidate & Frequent Length-1 Sequences

- <a> <c>, <d>, <e> <f> <g>

GSP Algorithm – Length-2 Sequences



Candidate & Frequent Length-2 Sequences

	<a>		<c>	<e>	<f>
<a>	<aa>	<ab> <ba>	<ac> <ca>	<ae> <ea>	<af> <fa>
		<bb>	<bc> <cb>	<be> <eb>	<bf> <fb>
<c>			<cc>	<ce> <ec>	<cf> <fc>
<e>				<ee>	<ef> <fe>
<f>					<ff>

	<a>		<c>	<e>	<f>
<a>		<(ab)>	<(ac)>	<(ae)>	<(af)>
			<(bc)>	<(be)>	<(bf)>
<c>				<(ce)>	<(cf)>
<e>					<(ef)>
<f>					

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

GSP Algorithm – Length-3 Sequences



Candidate & Frequent Length-3 Sequences

	<ab>	<ac>	<bc>	<(ab)>
<ab>	<aab> <abb>	<abc> <acb>	<abc>	<(ab)b> <a(ab)>
<ac>		<aac> <acc>	<abc> <bac>	<(ab)c>
<bc>			<bbc> <bcc>	<(ab)c>
<(ab)>				<(ab)(ab)>

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

Summary of GSP Algorithm



❑ Strength

- Apriori pruning

❑ Weakness

- Generates a huge set of candidate sequences
- Requires multiple scans of database
- Inefficient for mining long sequential patterns

❑ References

- Agrawal, R. and Srikant, R., “Mining sequential patterns.” In Proceedings of ICDE (1995)
- Srikant, R. and Agrawal, R., “Mining sequential patterns: Generalizations and performance improvements.” In Proceedings of EDBT (1996)

□ General Definition

- When a sequence S is a list of items,
- $S[1 .. j]$ is a prefix of a string $S[1 .. n]$ where $j \leq n$
- X is a prefix of Y if $X \cdot Z = Y$ for some string Z

□ Extended Definition

- When a sequence is a list of elements and each element is a set of items,
- Sort the items of each element in an alphabetic order
- Given $x = \langle e_1, e_2, \dots, e_n \rangle$, $y = \langle e'_1, e'_2, \dots, e'_m \rangle$ ($m \leq n$) is a **prefix** of x if and only if
 - 1) $e_i = e'_i$ for $(i \leq m-1)$
 - 2) $e'_m \subseteq e_m$
 - 3) all items in $(e_m - e'_m)$ are alphabetically after those in e'_m
- Examples: $\langle a \rangle$, $\langle aa \rangle$, $\langle a(ab) \rangle$ and $\langle a(abc) \rangle$ are prefixes of $\langle a(abc)(ac)d(cf) \rangle$



PrefixSpan (Prefix-projected Sequential Pattern mining) Algorithm

□ Main Idea

- Keep track of prefixes (instead of all candidate sequences) from the sequence database
- Project their suffixes into projected databases

□ Projected Database

- A set of maximal-length suffixes with a given prefix

□ Example

- <a>-projected database has
 - <(abc)(ac)d(cf)>
 - <(_b)a(bc)(ae)>
 - <(_b)(df)cb>
 - <(_f)cbc>

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

PrefixSpan Algorithm – cont'



□ Process

1) Find all frequent length-1 sequences

- $\langle a \rangle, \langle b \rangle, \langle c \rangle, \langle e \rangle, \langle f \rangle$

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

2) Divide search space to

- Maximal-length suffixes of a prefix ending with $\langle a \rangle$
- Maximal-length suffixes of a prefix ending with $\langle b \rangle$
- Maximal-length suffixes of a prefix ending with $\langle c \rangle$
- Maximal-length suffixes of a prefix ending with $\langle e \rangle$
- Maximal-length suffixes of a prefix ending with $\langle f \rangle$

3) Construct a projected database for each search space

4) Find sequential patterns recursively

PrefixSpan Algorithm – cont'



Projected Databases in the 1st Recursion

min_sup = 75%

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

pattern	projected database
<a>	<(abc)(ac)d(cf)>, <(_b)a(bc)(ae)>, <(_b)(df)cb>, <(_f)cbc>
	<(_c)(ac)d(cf)>, <a(bc)(ae)>, <(df)cb>, <c>
<c>	<(ac)d(cf)>, <(ae)>, , <bc>
<e>	<(_f)(ab)(df)cb>, <g(af)cbc>
<f>	<(ab)(df)cb>, <cbc>

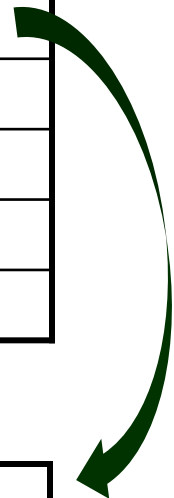
PrefixSpan Algorithm – cont'



Projected Databases in the 2nd Recursion

pattern	projected database
<a>	<(abc)(ac)d(cf)>, <(_b)a(bc)(ae)>, <(_b)(df)cb>, <(_f)cbc>
	<(_c)(ac)d(cf)>, <a(bc)(ae)>, <(df)cb>, <c>
<c>	<(ac)d(cf)>, <(ae)>, , <bc>
<e>	<(_f)(ab)(df)cb>, <g(af)cbc>
<f>	<(ab)(df)cb>, <cbc>

pattern	projected database
<ab>	<(_c)(ac)d(cf)>, <(_c)(ae)>, <c>
<ac>	<(ac)d(cf)>, <(ae)>, , <bc>
<(ab)>	<(_c)(ac)d(cf)>, <a(bc)(ae)>, <(df)cb>



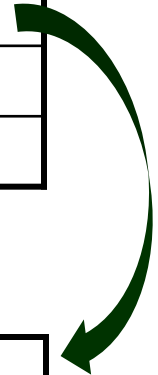
PrefixSpan Algorithm – cont'



❑ Projected Databases in the 2nd Recursion

pattern	projected database
<ab>	<(_c)(ac)d(cf)>, <(_c)(ae)>, <c>
<ac>	<(ac)d(cf)>, <(ae)>, , <bc>
<(ab)>	<(_c)(ac)d(cf)>, <a(bc)(ae)>, <(df)cb>

pattern	projected database
<(ab)c>	<d(cf)>, <(ae)>,



Summary of PrefixSpan Algorithm



❑ Strength

- Efficient (major cost is to construct projected databases)
- Project databases keep shrinking rapidly

❑ Weakness

- Searching frequency redundantly

❑ References

- Pei, J., et al., “PrefixSpan: Mining sequential patterns efficiently by prefix-projected pattern growth.” In Proceedings of ICDE (2001)

Questions?



- ❑ Lecture Slides on the Course Website, “https://ads.yonsei.ac.kr/faculty/data_mining”

