

Lecture 3 Week 4 (March 31)

33459-01 Principles of Knowledge Discovery in Data

Sequential Pattern Analysis

Lecture by: Dr. Osmar R. Zaiane

33459-01: Principles of Knowledge Discovery in Data – March-June, 2006

(Dr. O. Zaiane) 1

Course Content

- Introduction to Data Mining
- Association analysis
- **Sequential Pattern Analysis**
- Classification and prediction
- Contrast Sets
- Data Clustering
- Outlier Detection
- Web Mining



33459-01: Principles of Knowledge Discovery in Data – March-June, 2006

(Dr. O. Zaiane) 2

Sequence of Transactions

- Association rule mining searches for relationships between items in a dataset where time is irrelevant.



- Sequential Pattern Analysis considers time (or order of transactions).



Data: sequences of evidences in time order
Target: sub-sequences that happened frequently

33459-01: Principles of Knowledge Discovery in Data – March-June, 2006

(Dr. O. Zaiane) 3

Sequence Pattern Examples

- Examples 1
 - 60% of customers typically rent "star wars", then "Empire strikes back", and then "Return of Jedi".
 - Note: these rentals need not to be consecutive.
 - <SW>, ..., <ESB>, ..., <RJ>
- Example 2
 - 60% of customers buy "Fitted Sheet and flat sheet and pillow", followed by "comforter", followed by "drapes and ruffles"
 - Note: elements of a sequential pattern need not to be simple items.
 - <FittedSheet, FlatSheet, Pillow>, ..., <Comforter>, ..., <Drapes, Ruffles>

33459-01: Principles of Knowledge Discovery in Data – March-June, 2006

(Dr. O. Zaiane) 4

Lecture Outline

Part I: Concepts (30 minutes)

- Basic concepts

Part II: Apriori-based Approaches (45 minutes)

- Apriori-all
- GSP

Part III: Pattern-Growth-based Approaches (45 minutes)

- Free-Span
- Prefix-Span

33459-01: Principles of Knowledge Discovery in Data – March-June, 2006

(Dr. O. Zaiane) 5

Why sequential pattern mining?

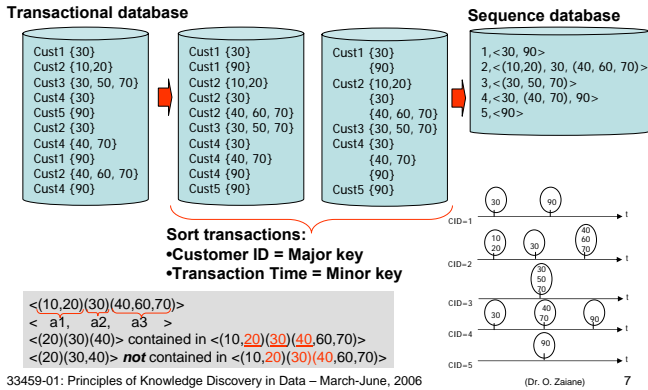
- Time or order in which actions appear or happen can be relevant in decision making.
- Many applications of sequential pattern mining
 - Customer shopping sequences (idem for book/video rental):
 - First buy computer, then CD-ROM, and then digital camera, within 3 months.
 - Medical treatment (e.g., symptoms and diseases)
 - Serial crime solving
 - Natural disasters (e.g., earthquakes),
 - Science & engineering processes,
 - Stocks and markets,
 - Telephone calling patterns,
 - Web access log click streams,
 - DNA sequences and gene structures, etc.

33459-01: Principles of Knowledge Discovery in Data – March-June, 2006

(Dr. O. Zaiane) 6

Sequence Database

Converts the original transaction database into a database of customer sequences.



What Is Sequential Pattern Mining?

- Given a set of sequences, find the complete set of **frequent** subsequences

A **sequence**: $\langle (ef)(ab)(df)cb \rangle$

A **sequence database**

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	$\langle eg(af)cbc \rangle$

An element may contain a set of items. Items within an element are unordered and we list them alphabetically.

$\langle a(bc)dc \rangle$ is a **subsequence** of $\langle a(abc)(ac)d(cf) \rangle$

Given **support threshold** $min_sup = 2$, $\langle (ab)c \rangle$ is a **sequential pattern** (cf. SID 10 & 30)

Sequential Patterns

- 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

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	$\langle eg(af)cbc \rangle$

$\langle a(abc)(ac)d(cf) \rangle$ - 5 elements, 9 items
 $\langle a(abc)(ac)d(cf) \rangle$ - 9-sequence
 $\langle a(abc)(ac)d(cf) \rangle = \langle a(cba)(ac)d(cf) \rangle$
 $\langle a(abc)(ac)d(cf) \rangle \neq \langle a(ac)(abc)d(cf) \rangle$

Order doesn't Matter (list)
 Order matters (sequence)

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**

Solution – 53 frequent subsequences

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	$\langle eg(af)cbc \rangle$

$min_support = 2$

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

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
- A sequence s is maximal if it is not contained in any other sequence.

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	$\langle eg(af)cbc \rangle$

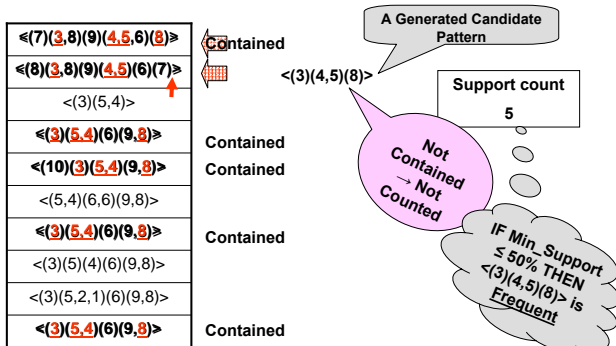
$\beta = \langle a(abc)(ac)d(cf) \rangle$
 $\alpha_1 = \langle aa(ac)d(c) \rangle$
 $\alpha_2 = \langle (ac)(ac)d(cf) \rangle$
 $\alpha_3 = \langle ac \rangle$
 $\beta = \langle a(abc)(ac)d(cf) \rangle$
 $\alpha_4 = \langle d(cf) \rangle$
 $\alpha_5 = \langle (cf)d \rangle$
 $\alpha_6 = \langle (abc)dcf \rangle$

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	Support
10	$\langle a(abc)(ac)d(cf) \rangle$	$\alpha_1 = \langle a \rangle$ support(α_1) = 4
20	$\langle (ad)c(bc)(ae) \rangle$	$\alpha_2 = \langle ac \rangle$ support(α_2) = 4
30	$\langle (ef)(ab)(df)cb \rangle$	$\alpha_3 = \langle (ab)c \rangle$ support(α_3) = 2
40	$\langle eg(af)cbc \rangle$	

Counting Sequences (An example)



Challenges on Sequential Pattern Mining

- A huge number of possible sequential patterns are hidden in databases
- 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
- Comparison of association rules and sequence mining
 - Mining for association rules
 - Purpose: Discovery of frequent unordered itemsets.
 - Complexity: With n items there are $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ k -itemsets (sets with k items).
 - Mining for Sequential Patterns
 - Purpose: Discovery of frequent sequences of (unordered) itemsets.
 - Complexity: With n items there are n^k k -sequences (sequences with k items).

Association mining algorithms discover isolated item sets (intra-event patterns). Sequence mining algorithms discover series of item sets (inter-event patterns).

Studies on Sequential Pattern Mining

- Concept introduction and an initial Apriori-like algorithm
 - R. Agrawal & R. Srikant. "Mining sequential patterns," ICDE'95
- GSP—An Apriori-based, influential mining method (developed at IBM Almaden)
 - R. Srikant & R. Agrawal. "Mining sequential patterns: Generalizations and performance improvements," EDBT'96
- From sequential patterns to episodes (Apriori-like + constraints)
 - H. Mannila, H. Toivonen & A.I. Verkamo. "Discovery of frequent episodes in event sequences," Data Mining and Knowledge Discovery, 1997
- Data Projection-based approaches
 - FreeSpan (Han et al. "frequent pattern-projected sequential pattern mining" SIGKDD 2000)
 - PrefixSpan (Pei et al. "Prefix-projected pattern growth" ICDE 2001)
- Mining sequential patterns with constraints
 - M.N. Garofalakis, R. Rastogi, K. Shim: SPIRIT: Sequential Pattern Mining with Regular Expression Constraints. VLDB 1999

Lecture Outline

Part I: Concepts (30 minutes)

- Basic concepts

Part II: Apriori-based Approaches (45 minutes)

- Apriori-All
- GSP

Part III: Pattern-Growth-based Approaches (45 minutes)

- Free-Span
- Prefix-Span

AprioriAll: The idea

- Basic method to mine sequential patterns
 - Based on the Apriori algorithm
 - Count all the large sequences, including non-maximal sequences
 - Use Apriori-generate function to generate candidate sequences: get candidates for a pass using only the large sequences found in the previous pass and make a scan over the data to find their support

AprioriAll: The big picture

Five-phase algorithm

1. Sort phase:
 - Create the sequence database from transactions.
2. Large itemset phase
 - Find all frequent itemsets using Apriori
3. Transformation phase:
 - Do integer mapping for large itemsets
4. Sequence phase:
 - Find all frequent sequential patterns using Apriori.
5. Maximal phase:
 - Eliminate non maximal sequences.

AprioriAll Algorithm(1)

AprioriAll Algorithm

C_k : Candidate sequence of size k
 L_k : frequent or large sequence of size k

```

L1 = {large 1-sequence}; //result of itemset phase
for (k = 2; L_{k-1} != ∅; k++) do begin
    C_k = candidates generated from L_{k-1};
    for each customer-sequence c in database do
        Increment the count of all candidates in C_k
        that are contained in c
    L_k = Candidates in C_k with minimum support
end
Answer = Maximal sequences in ∪_k L_k;
    
```

Candidate Generation --Join Step:

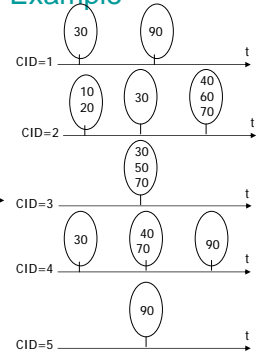
```

C_k is generated by joining L_{k-1} with itself
insert into C_k,
select p.itemset_1, ..., p.itemset_{k-1}, q.itemset_{k-1}
from L_{k-1} p, L_{k-1} q
where p.itemset_1 = q.itemset_1, ...,
p.itemset_{k-2} = q.itemset_{k-2}
    
```

For example:
 $\{1,2,3\} \times \{1,2,4\}$
 $= \{1,2,3,4\}$
 and
 $\{1,2,4,3\}$

Sequence Database Example

Customer ID	TransactionTime	Items
1	1	30
1	2	90
2	1	10,20
2	2	30
2	3	40,60,70
3	1	30,50,70
4	1	30
4	2	40,70
4	3	90
5	1	90



MinSupport = 40%, i.e. 2 customers

Answer: $\langle 30 \rangle \langle 90 \rangle$ (CID1,4) $\langle 30 \rangle \langle 40,70 \rangle$ (CID2,4)

Not Answer: $\langle 30 \rangle \langle 40 \rangle \langle 70 \rangle \langle 90 \rangle \langle 30 \rangle \langle 40 \rangle \langle 30 \rangle \langle 70 \rangle \langle 40 \rangle \langle 70 \rangle$ Why?

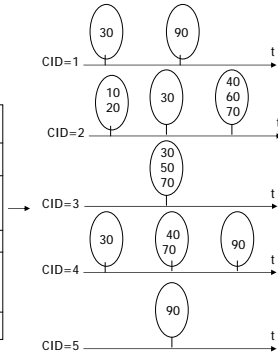


Sort Phases

Sort Phases

CID: major key, TID: secondary key

Customer ID	TransactionTime	Items
1	1	30
1	2	90
2	1	10,20
2	2	30
2	3	40,60,70
3	1	30,50,70
4	1	30
4	2	40,70
4	3	90
5	1	90



Itemset Phase

Itemset Phase:

- Find all large itemset (Why? Because each itemset in a large sequence has to be a large itemset.)
- To get large (frequent) itemsets → Use Apriori algorithm
- Need to modify support counting. (For sequential patterns, support is measured by fraction of customers.)

Customer ID	TransactionTime	Items
1	1	30
1	2	90
2	1	10,20
2	2	30
2	3	40,60,70
3	1	30,50,70
4	1	30
4	2	40,70
4	3	90
5	1	90

MinSupport = 40%, i.e. 2 customers

Itemset Result:

$\{30\} \{40\} \{70\} \{40,70\} \{90\}$

Difference from Apriori:

- the support count should be incremented only once per customer

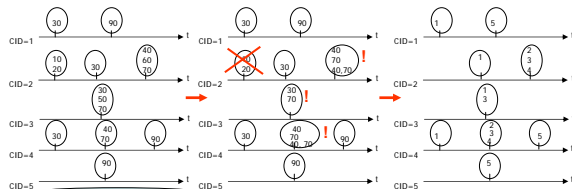
Transformation Phase

Transformation Phase:

- Each large itemset is then mapped to a set of contiguous integers (Why? to be able to compare two frequent itemsets in constant time)
- Represent transactions as sets of large itemsets.

itemset	Map
$\{30\}$	1
$\{40\}$	2
$\{70\}$	3
$\{40,70\}$	4
$\{90\}$	5

Itemsets



Sequence
1, <30, 90>
2, <(10,20), 30, (40, 60, 70)>
3, <(30, 50, 70)>
4, <30, (40, 70), 90>
5, <90>

Sequence database

Transformed database

Sequence
1, <{1}, {5}>
2, <{1}, {2, 3, 4}>
3, <{1}, {3}>
4, <{1}, {2, 3, 4}, {5}>
5, <{5}>

Sequence Phase

Sequence Phase:

- Use set of large itemsets to find the desired sequences.
- Similar structure to Apriori algorithms used to find large itemsets.
 - Use seed set to generate candidate sequences.
 - Count support for each candidate.
 - Eliminate candidate sequences which are not large.

itemset	Map
$\{30\}$	1
$\{40\}$	2
$\{70\}$	3
$\{40,70\}$	4
$\{90\}$	5

Itemsets

MinSupport = 40%, i.e. 2 customers

F.Seq.	Sup.
{1}	4
{2}	2
{3}	3
{4}	2
{5}	3
{1,2}	2
{1,3}	3
{1,4}	2
{1,5}	2

Apriori

Re-mapping

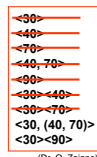
<30>
<40>
<70>
<40,70>
<90>
<30><40>
<30><70>
<30, (40, 70)>
<30><90>

Maximal phase

- **Maximum Phase:**

- Find the maximal sequences among the set of frequent sequences
- delete all sequences that are sub-sequences of other frequent sequences.

```
for (k=n; k>1; k--) do
  for each k-sequence Sk do
    Delete from all subsequences of Sk
```



Summary for AprioriAll

- Algorithm wastes much time in counting non-maximal sequences, which can not be sequential patterns
- There are other variations of AprioriAll that reduce the candidates that are not maximal: AprioriSome and DynamicSome
- Absence of time window constraints
- *AprioriALL* is the basis of many efficient algorithm developed later. GSP is among them.

GSP—A Generalized Sequential Pattern Mining Algorithm

- GSP (Generalized Sequential Pattern) mining algorithm
 - proposed by Agrawal and Srikant, EDBT'96
- Outline of the method
 - Initially, every item in DB is a candidate of length-1
 - for each level (i.e., sequences of length-k) do
 - scan database to collect support count for each candidate sequence
 - generate candidate length-(k+1) sequences from length-k frequent sequences using Apriori
 - repeat until no frequent sequence or no candidate can be found
- Major strength: Candidate pruning by Apriori

A Basic Property of Sequential Patterns: Apriori

- A basic property: Apriori (Agrawal & Srikant'94)
 - If a sequence S is not frequent
 - Then none of the super-sequences of S is frequent
 - E.g, <hb> is infrequent → so do <hab> and <(ah)b>

Seq. ID	Sequence
10	<(bd)cb(ac)>
20	<(bf)(ce)b(fg)>
30	<(ah)(bf)abf>
40	<(be)(ce)d>
50	<a(bd)bc(bade)>

Given support threshold
 $min_sup = 2$

The GSP Algorithm

- Take sequences in form of <x> as length-1 candidates
- Scan database once, find F_1 , the set of length-1 sequential patterns
- Let $k=1$; while F_k is not empty do
 - Form C_{k+1} , the set of length-(k+1) candidates from F_k ;
 - If C_{k+1} is not empty, scan database once, find F_{k+1} , the set of length-(k+1) sequential patterns
 - Let $k=k+1$;

Finding Length-1 Sequential Patterns

- Examine GSP using an example
- Initial candidates: all singleton sequences
 - <a>, , <c>, <d>, <e>, <f>, <g>, <h>
- Scan database once, count support for candidates

Seq. ID	Sequence
10	<(bd)cb(ac)>
20	<(bf)(ce)b(fg)>
30	<(ah)(bf)abf>
40	<(be)(ce)d>
50	<a(bd)bc(bade)>

$min_sup = 2$

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

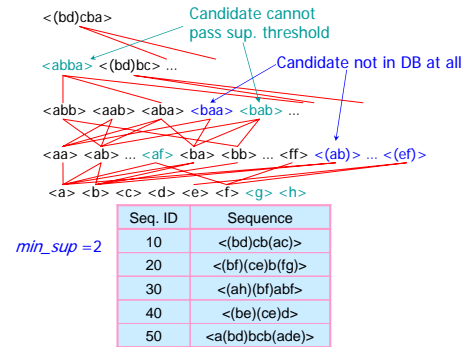
Generating Length-3 Candidates

<a>		<c>	<d>	<e>	<f>
<aa>:2	<ab>:2	<ac>:1	<ad>:1	<ae>:1	<af>:1
<ba>:3	<bb>:4	<bc>:4	<bd>:2	<be>:3	<bf>:2
<ca>:2	<cb>:3	<cc>:1	<cd>:2	<ce>:1	<cf>:1
<da>:2	<db>:2	<dc>:2	<dd>:0	<de>:1	<df>:0
<ea>:0	<eb>:1	<ec>:0	<ed>:1	<ee>:1	<ef>:1
<fa>:1	<fb>:2	<fc>:1	<fd>:0	<fe>:1	<ff>:2

Example of generating <(xy)z> pattern for <(bd)>:

- Need to concatenate another Length-2 frequent itemset
- Concatenating those patterns that **end with b or d** to form something like <a(bd)>, <b(bd)>, <c(bd)>, <d(bd)>, <f(bd)>
- Concatenating those patterns that **starts with b or d** to form something like <(bd)a>, <(bd)b>, <(bd)c>, <(bd)d>, <(bd)e>, <(bd)f>

The GSP Mining Process



Bottlenecks of GSP

- A huge set of candidates could be generated
 - 1,000 frequent length-1 sequences generate $1000 \times 1000 + \frac{1000 \times 999}{2} = 1,499,500$ length-2 candidates!
- Multiple scans of database
- Real challenge: mining long sequential patterns
 - An exponential number of short candidates
 - A length-100 sequential pattern needs 10^{30} candidate sequences!

Lecture Outline

Part I: Concepts (30 minutes)

- Basic concepts

Part II: Apriori-based Approaches (45 minutes)

- Apriori-All
- GSP

Part III: Pattern-Growth-based Approaches (45 minutes)

- Free-Span (Frequent Pattern-Projected Sequential Pattern Mining)
- Prefix-Span (Prefix-Projected Sequential Pattern)

FreeSpan (Generalities)

- J. Han, J. Pei, B. Mortazavi-Asl, Q. Chen, U. Dayal, and M.-C. Hsu. **FreeSpan: Frequent pattern-projected sequential pattern mining**. KDD'00, pages 355-359.
- A divide-and-conquer approach
 - Recursively **project** a sequence database into a set of smaller databases
 - Mine each projected database to find the subset of patterns

FreeSpan (example)

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

- Given a sequence database S and min_support = 2
 - Step 1: find length-1 sequential patterns and list them in support descending order
 - f_list = a:4, b:4, c:4, d:3, e:3, f:3, ~~g:1~~
 - Step 2: divide search space. The complete set of seq. patterns can be partitioned into 6 disjoint subsets (move down the f_list):
 - ones only contain item **a**
 - ones contain item **b** but no items after b in f_list
 - ones contain item **c** but no items after c in f_list
 - ones contain item **d** but no items after d in f_list
 - ones contain item **e** but no items after e in f_list
 - ones contain item **f**
- find subsets of sequential patterns. They can be mined by constructing projected databases and mining each recursively

FreeSpan (con't)

- Finding Seq. Patterns containing item b but no items after b in f_list
 - -projected database:
 - 10:<a(ab)a>, 20:<aba>, 30:<(ab)b>, 40:<ab>
 - Find all the length-2 seq. patterns containing item b but no items after b in f_list :
 - <ab>:4, <ba>:2, <(ab)>:2
 - Further partition and mining

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

f_list = a:4,b:4, c:4,d:3,e:3,f:3

From FreeSpan to PrefixSpan

- Freespan:
 - Projection-based: No candidate sequence needs to be generated
 - But, projection can be performed at any point in the sequence, and the projected sequences may not shrink much. For example, the size of f-projected database is the same as the original sequence database
- PrefixSpan
 - Projection-based
 - But only prefix-based projection: less projections and quickly shrinking sequences

• J. Pei, J. Han, H. Pinto, Q. Chen, U. Dayal, and M.-C. Hsu, "PrefixSpan: Mining Sequential Patterns Efficiently by Prefix-Projected Pattern Growth", Proc. 2001 Int. Conf. on Data Engineering (ICDE'01), Heidelberg, Germany, April 2001.

Prefix of a Sequence

- Given two sequences $\alpha = \langle a_1 a_2 \dots a_n \rangle$ and $\beta = \langle b_1 b_2 \dots b_m \rangle$, $m \leq n$
- 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

Given an alphabetical order of items in each itemset (element)



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

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

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



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

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 $\alpha = \beta \cdot \gamma$



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

PrefixSpan – Algorithm

- Input:** A sequence database S, and the minimum support threshold min_sup
- Output:** The complete set of sequential patterns
- Method:** Call PrefixSpan($\langle \rangle, 0, S$)
- 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 $\text{PrefixSpan}(\alpha', l+1, S|_{\alpha'})$.

PrefixSpan - Example

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

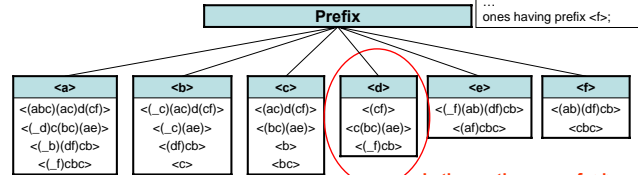
min_support = 2

Partition search space into 6 subsets:
 ones having prefix <a>;
 ones having prefix ;
 ...
 ones having prefix <f>;

- Find length-1 sequential patterns

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

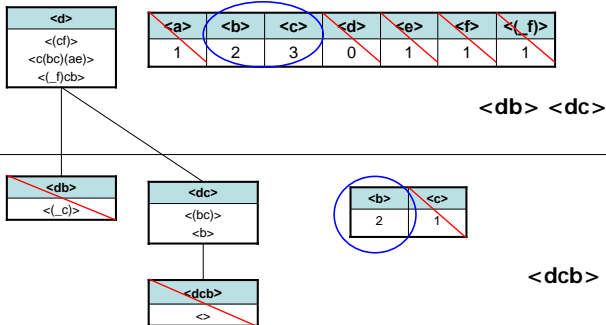
- Divide search space



PrefixSpan – Example (2)

Projected database for <d>

- Find subsets of sequential patterns



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

Bi-level Projection

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

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

	a	b	c	d	e	f
a	2					
b	(4,2,2)	1				
c	(4,2,1)	(3,3,2)	3			
d	(2,1,1)	(2,2,0)	(1,3,0)	0		
e	(1,2,1)	(1,2,0)	(1,2,0)	(1,1,0)	0	
f	(2,1,1)	(2,2,0)	(1,2,1)	(1,1,1)	(2,0,1)	1

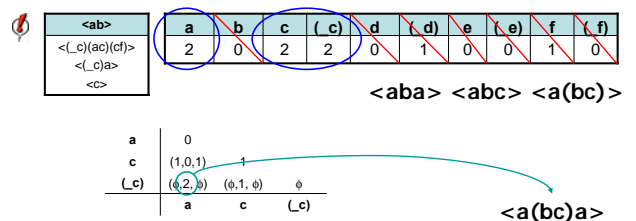
Support(<ac>) = 4
 Support(<ca>) = 2
 Support(<ac>) = 1

Support(<cc>) = 3

ALL length-2 sequential pattern


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



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>

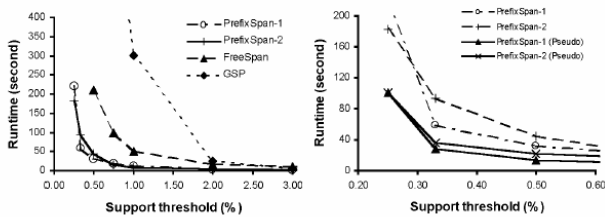
Pseudo-Projection

- **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)>

Efficiency of Prefix-Span and Effect of Pseudo-Projection



PrefixSpan-1 (level-by-level projection)
 PrefixSpan-2 (bi-level projection)

Summary

- Sequential Pattern Mining is useful in many application, e.g. weblog analysis, financial market prediction, Bioinformatics, etc.
- It is similar to the frequent itemsets mining, *but* with consideration of ordering.
- We have looked at different approaches that are descendants from two popular algorithms in mining frequent itemsets
 - Candidates Generation: AprioriAll and GSP
 - Pattern Growth: FreeSpan and PrefixSpan