

Data Mining (EECS 4412)

~~Assignment Project Exam Help~~

<https://powcoder.com>
Sequential Pattern Mining
Add WeChat powcoder

Parke Godfrey

EECS

Lassonde School of Engineering

York University

Thanks to

Professor Aijun An

Assignment Project Exam Help

for creation & use of these slides.

<https://powcoder.com>

Add WeChat powcoder

Outline

- ▶ Basic concepts of sequential pattern mining
- ▶ A Simplified Version of GSP Algorithm
Assignment Project Exam Help
<https://powcoder.com>
- ▶ PrefixSpan
Add WeChat powcoder

An Example Sequence Database

Customer Id	TransactionTime	Items Bought
1	June 25 '03	30
1	June 30 '03	90
2	June 10 '03	10, 20
2	June 15 '03	30
2	June 20 '03	40, 60, 70
3	June 25 '03	30, 50, 70
4	June 25 '03	30
4	June 30 '03	40, 70
4	July 25 '03	90
5	June 12 '03	90

Figure 1: Database Sorted by Customer Id and Transaction Time

Customer Id	Customer Sequence
1	{ (30) (90) }
2	{ (10 20) (30) (40 60 70) }
3	{ (30 50 70) }
4	{ (30) (40 70) (90) }
5	{ (90) }

Figure 2: Customer-Sequence Version of the Database

► A *sequence database* consists of a set of

sequences

► A *sequence* is an ordered list of itemsets

Itemset: a set of items

► Items within an itemset are unordered.

► *Element*: an itemset in a sequence

Example of a Sequential Pattern

- ▶ Database of an online book store
 - ▶ Contains data sequences
 - ▶ Each sequence corresponds to a list of transactions by a given customer.
 - ▶ Each transaction contains the books selected by the customer in one order.
- ▶ A sequential pattern
 - ▶ 5% of customers bought “*Foundation*”, then “*Foundation and Empire*” and “*Ringworld*”, then “*Second Foundation*”.
- ▶ Usefulness
 - ▶ Making recommendations.
 - ▶ Knowing what to stock.

Mining Sequential Patterns

- ▶ Objective
 - ▶ Finding (interesting) frequent sequences from a sequence database
- ▶ Applications
 - ▶ Customer shopping sequence analysis
 - ▶ Web log analysis
 - ▶ DNA or protein analysis
 - ▶ Stock market analysis
 - ▶ Medical domain
 - ▶ Predict onset of a disease from a sequence of symptoms, etc.
 - ▶ Earthquake prediction
 - ▶ etc.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Basic Concepts

- ▶ A sequence $\langle a_1, a_2, \dots, a_n \rangle$ *is contained in* (or *is a subsequence of*) another sequence $\langle b_1, b_2, \dots, b_m \rangle$ if there exist integers $i_1 < i_2 < \dots < i_n$ such that $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$.
 - ▶ E.g., $\langle (3) (4\ 5) (8) \rangle$ is contained in $\langle (7) (3\ 8) (9) (4\ 5\ 6) (8) \rangle$
- ▶ *Support of a sequence* is defined as the fraction of total sequences in the database that contains this sequence.
 - ▶ Support of $\langle (30) (90) \rangle$?
 - ▶ Support of $\langle (20) (60\ 70) \rangle$?
 - ▶ Support of $\langle (30) (70) \rangle$?
- ▶ *Frequent sequences* (also called *sequential patterns*)
 - ▶ sequences that satisfy a minimum support (min_sup).

Customer Id	Customer Sequence
1	$\langle (30) (90) \rangle$
2	$\langle (10\ 20) (30) (20\ 60\ 70) \rangle$
3	$\langle (30\ 50\ 70) \rangle$
4	$\langle (30) (40\ 60\ 70) (90) \rangle$
5	$\langle (90) \rangle$

Sequential Pattern Mining

- ▶ What is sequential pattern mining
 - ▶ Given a sequence database, find the set of frequent sequences that satisfy a minimum support (min_sup).
- ▶ Algorithms
 - ▶ Initial Apriori-like algorithms (Agrawal and Srikant, 95)
 - ▶ AprioriAll, AprioriSome, and DynamicSome
 - ▶ GSP - an Apriori-like, influential mining method (Srikant and Agrawal, 96)
 - ▶ PrefixSpan (Pei, et al, 01)
 - ▶ SPADE (Zaki, 01)
 - ▶ Mining sequential patterns with constraints (Pei, et al, 02)
 - ▶ etc.

Apriori Property

- ▶ The **Apriori** property in sequential patterns:
Any nonempty subsequence of a frequent sequence must be frequent
~~Assignment Project Exam Help~~
- ▶ If a sequence is infrequent, then none of its super-sequences is frequent.
~~https://powcoder.com~~
~~Add WeChat powcoder~~
- ▶ i.e., if $\langle (3) (4\ 5) \rangle$ is infrequent, so are $\langle (3) (4\ 5) (8) \rangle$ and $\langle (3\ 6) (4\ 5) \rangle$

GSP

- ▶ GSP: Generalized Sequential Pattern Mining
- ▶ Proposed in
 - ▶ R. Srikant and R. Agrawal. Mining Sequential Patterns: Generalizations and Performance Improvements. In Proc. of EDBT'96.
 - ▶ Can be downloaded from the course web site
- ▶ GSP considers some time constraints and item taxonomy
 - ▶ More general than simply mining sequential patterns

A Simplified Version of GSP

- ▶ GSP without the generalizations
- ▶ Problem statement:
 - ▶ find all frequent sequences from a database of sequences given a min_sup
- ▶ Length of a sequence = number of items in the sequence
 - ▶ Length of $\langle(a)(b)\rangle$ is 2
 - ▶ Length of $\langle(a\ b)\rangle$ is 2
 - ▶ Length of $\langle(a\ b)\ (c)\rangle$ is 3
- ▶ A *length-k sequence* is also called *k-sequence*.

General Description of Simplified GSP

► Method

- Take sequences in form of $\langle(x)\rangle$ as length-1 candidates
- Scan database once, find L_1 , the set of length-1 sequential patterns
- Let $k=1$; while L_k is not empty do
 - Form C_{k+1} , the set of length- $(k+1)$ candidates from L_k ;
 - If C_{k+1} is not empty, scan database once, find L_{k+1} , the set of length- $(k+1)$ sequential patterns
 - Let $k=k+1$;

Finding Length-1 Sequential Patterns

► Initial candidates

- $\langle (a) \rangle, \langle (b) \rangle, \langle (c) \rangle, \langle (d) \rangle, \langle (e) \rangle, \langle (f) \rangle, \langle (g) \rangle, \langle (h) \rangle$

► Scan database once

- count support for candidates

Seq-id	Sequence
10	$\langle (bd)(c)(b)(ac) \rangle$
20	$\langle (bf)(ce)(b)(fg) \rangle$
30	$\langle (ah)(bf)(a)(b)(f) \rangle$
40	$\langle (be)(ce)(d) \rangle$
50	$\langle (a)(bd)(b)(c)(b)(ade) \rangle$

$min_sup = 40\%$

$min_sup_count = 2$

Cand	Sup
$\langle (a) \rangle$	3
$\langle (b) \rangle$	5
$\langle (c) \rangle$	4
$\langle (d) \rangle$	3
$\langle (e) \rangle$	3
$\langle (f) \rangle$	2
$\langle (g) \rangle$	1
$\langle (h) \rangle$	1

Length-2 Candidate Generation

- ▶ How to generate C_2 from L_1
 - ▶ Merge every pair of frequent length-1 sequences.
 - ▶ Merging two frequent length-1 sequences $\langle i_1 \rangle$ and $\langle i_2 \rangle$ will produce three candidate 2-sequences:
 $\langle i_1 i_2 \rangle$, $\langle i_2 i_1 \rangle$ and $\langle i_1 i_1 \rangle$ assuming that i_1 is different from i_2 .
 - ▶ Every frequent length-1 sequence $\langle i \rangle$ can join with itself to produce $\langle i i \rangle$.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Length-2 Candidates

(Cont'd from the example on Slide 13)

Length of a sequence = the number of items

51 length-2 Candidates

	<a>		<c>	<d>	<e>	<f>
<a>	<(a)(a)>	<(a)(b)>	<(a)(c)>	<(a)(d)>	<(a)(e)>	<(a)(f)>
	<(b)(a)>	<(b)(b)>	<(b)(c)>	<(b)(d)>	<(b)(e)>	<(b)(f)>
<c>	<(c)(a)>	<(c)(b)>	<(c)(c)>	<(c)(d)>	<(c)(e)>	<(c)(f)>
<d>	<(d)(a)>	<(d)(b)>	<(d)(c)>	<(d)(d)>	<(d)(e)>	<(d)(f)>
<e>	<(e)(a)>	<(e)(b)>	<(e)(c)>	<(e)(d)>	<(e)(e)>	<(e)(f)>
<f>	<(f)(a)>	<(f)(b)>	<(f)(c)>	<(f)(d)>	<(f)(e)>	<(f)(f)>

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

	<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>						

Without Apriori

property,

$$8*8 + 8*7/2 = 92$$

candidates

Apriori prunes

44.57% candidates

Finding Length-2 Sequential Patterns

- ▶ Scan database one more time, collect support count for each length-2 candidate
- ▶ There are 19 length-2 candidates which pass the minimum support threshold in our example DB

- ▶ They are length-2 sequential patterns:

<(a)(a)>, <(a)(b)>

<(b)(a)>, <(b)(b)>, <(b)(c)>, <(b)(d)>, <(b)(e)>, <(b)(f)>

<(c)(a)>, <(c)(b)>, <(c)(d)>

<(d)(a)>, <(d)(b)>, <(d)(c)>

<(f)(b)>, <(f)(f)>

<(bd)>, <(bf)>, <(ce)>

Generating C_k from L_{k-1} ($k > 2$)

- ▶ Join step: generate C_k by joining L_{k-1} with L_{k-1}
 - ▶ All items in an itemset are ranked in an order.
 - ▶ A sequence s_1 joins with s_2 if the subsequence obtained by dropping the first item of s_1 is the same as the subsequence obtained by dropping the last item of s_2 .
 - ▶ The joined sequence is s_1 plus the last item of s_2 .
 - ▶ The added item becomes a separate element if it was a separate element in s_2 , or part of the last element of s_1 otherwise.
- ▶ Examples:
 - ▶ Joining $\langle(1)(2\ 3)(4)\rangle$ and $\langle(2\ 3)(4)(5)\rangle$ produces $\langle(1)\ (2\ 3)(4)(5)\rangle$
 - ▶ Joining $\langle(1)(2\ 3)(4)\rangle$ and $\langle(2\ 3)(4\ 5)\rangle$ produces $\langle(1)(2\ 3)(4\ 5)\rangle$

Generating C_k from L_{k-1} ($k > 2$) (Cont'd)

- ▶ Prune step: delete candidates in C_k that have infrequent $(k-1)$ -subsequence
 - ▶ A $(k-1)$ -subsequence of sequence s is derived by dropping an item from s
 - ▶ Check each $(k-1)$ -subsequence against L_{k-1}
 - ▶ Example:
 - ▶ If $\langle (ab)(d) \rangle$, $\langle (b)(ad) \rangle$, $\langle (b)(de) \rangle$ are all length-3 frequent sequences, then
 - ▶ Join result: $\langle (ab)(de) \rangle$
 - ▶ Its length-3 subsequences:
 - ▶ $\langle (b)(de) \rangle$, $\langle (a)(de) \rangle$, $\langle (ab)(e) \rangle$, $\langle (ab)(d) \rangle$
 - ▶ $\langle (a)(de) \rangle$ and $\langle (ab)(e) \rangle$ are infrequent, $\langle (ab)(de) \rangle$ is pruned
 - ▶ As long as one of them is infrequent, can stop checking others.
 - ▶ No need to check $\langle (ab)(d) \rangle$ and $\langle (b)(de) \rangle$ (Why?)
 - ▶ C_4 becomes empty

Generating C_k from L_{k-1} ($k > 2$) (Cont'd)

More Examples:

- ▶ If $\langle (a)(b) \rangle$, $\langle (a)(a) \rangle$ and $\langle (b)(a) \rangle$ are all length-2 sequential patterns, then length-3 candidates are:
 - ▶ Join result:
 - ▶ $\langle (a)(b)(a) \rangle$, $\langle (a)(a)(b) \rangle$, $\langle (a)(a)(a) \rangle$, $\langle (b)(a)(b) \rangle$, and $\langle (b)(a)(a) \rangle$.
 - ▶ After pruning:
 - ▶ $\langle (a)(b)(a) \rangle$, $\langle (a)(a)(b) \rangle$, $\langle (a)(a)(a) \rangle$, $\langle (b)(a)(a) \rangle$.
- ▶ If $\langle (bd) \rangle$, $\langle (b)(b) \rangle$ and $\langle (d)(b) \rangle$ are all length-2 sequential patterns, what are the length-3 candidates?
 - ▶ Join result:
 - ▶ $\langle (bd)b \rangle$, $\langle (b)(bd) \rangle$, $\langle (b)(b)(b) \rangle$, $\langle (d)(bd) \rangle$, $\langle (d)(b)(b) \rangle$
 - ▶ After pruning:
 - ▶ $\langle (bd)b \rangle$, $\langle (b)(b)(b) \rangle$, $\langle (d)(b)(b) \rangle$

Example Continued

- ▶ L_4 : Length-4 sequential patterns: $min_sup_count=2$

- ▶ $\langle (b)(c)(b)(a) \rangle$

- ▶ $\langle (bd)(b)(a) \rangle$

- ▶ $\langle (bd)(b)(c) \rangle$

- ▶ $\langle (bd)(c)(a) \rangle$

- ▶ $\langle (bd)(c)(b) \rangle$

- ▶ $\langle (bf)(b)(f) \rangle$

- ▶ $\langle (d)(c)(b)(a) \rangle$

Seq-id	Sequence
10	$\langle (bd)(c)(b)(ac) \rangle$
20	$\langle (bf)(ce)(b)(fg) \rangle$
30	$\langle (ah)(bf)(a)(b)(f) \rangle$
40	$\langle (be)(ce)(d) \rangle$
50	$\langle (a)(bd)(b)(c)(b)(ade) \rangle$

- ▶ C_5 : Length-5 candidates (after join and prune):

- ▶ $\langle (bd)(c)(b)(a) \rangle$

- ▶ L_5 : Length-5 sequential pattern (found by scanning DB):

- ▶ $\langle (bd)(c)(b)(a) \rangle$: 2

Summary of Simplified GSP

- ▶ Make the first pass over the sequence database D to find all the 1-element (length-1) frequent sequences
- ▶ Repeat until no new candidate or frequent sequences are found
- ▶ **Candidate Generation:**
 - ▶ **Merge** joinable pairs of frequent sequences of length $(k-1)$ to generate candidate sequences that contain k items
 - ▶ **Prune** candidate k -sequences that contain infrequent $(k-1)$ -subsequences
- ▶ **Support Counting and Candidate Elimination:**
 - ▶ Make a new pass over the sequence database D to find the support count for each candidate sequence
 - ▶ Eliminate candidate k -sequences whose actual support is less than min_sup

Bottlenecks of GSP

- ▶ A huge set of candidates
 - ▶ 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 in mining
- ▶ Real challenge: mining long sequential patterns
 - ▶ An exponential number of short candidates
 - ▶ A length-100 sequential pattern needs 10^{30} candidate sequences!

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

Better Method: PrefixSpan

- ▶ Generate all frequent sequences without candidate generation and testing.

<https://powcoder.com>

- ▶ J. Pei et al. PrefixSpan: Mining sequential patterns efficiently by prefix-projected pattern growth. In Proc. of ICDE'01.

- ▶ Will talk about it next.

PrefixSpan

- ▶ Generate all frequent sequences without candidate generation and testing.
- ▶ Strategy
 - ▶ Divide and conquer
 - ▶ Divide the patterns to be mined into subsets and find patterns in each subset recursively.
 - ▶ Projection-based
 - ▶ Recursively project a sequence database into a set of smaller databases based on the frequent “prefix” mined so far
 - ▶ Mine each projected database to find frequent “suffixes”

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Prefix and Suffix

- ▶ Suppose all the items in an element (i.e. an itemset) of a sequence are listed alphabetically.
- ▶ *Prefixes* of sequence $\langle a(abc)(ac)d(cf) \rangle$ are $\langle a \rangle$, $\langle aa \rangle$, $\langle a(ab) \rangle$, $\langle a(abc) \rangle$, $\langle a(abc)a \rangle$, ...
- ▶ *Suffixes* of sequence $\langle a(abc)(ac)d(cf) \rangle$:
 - ▶ $\langle (abc)(ac)d(cf) \rangle$ is the suffix wrt prefix $\langle a \rangle$
 - ▶ $\langle (_bc)(ac)d(cf) \rangle$ is the suffix wrt prefix $\langle aa \rangle$
 - ▶ $\langle (_c)(ac)d(cf) \rangle$ is the suffix wrt prefix $\langle a(ab) \rangle$
 - ▶ ...

Mining Sequential Patterns by Prefix Projections

- ▶ Step 1: find length-1 sequential patterns
 - ▶ $\langle a \rangle, \langle b \rangle, \langle c \rangle, \langle d \rangle, \langle e \rangle, \langle f \rangle$
- ▶ Step 2: divide search space. The complete set of freq. seq. can be partitioned into 6 subsets:
 - ▶ The ones having prefix $\langle a \rangle$;
 - ▶ The ones having prefix $\langle b \rangle$;
 - ▶ ...
 - ▶ The ones having prefix $\langle f \rangle$

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	$\langle eg(af)cbc \rangle$

min_sup count=2

Finding Freq. Seq. with Prefix <a>

- ▶ Only need to consider sequences containing <a>
- ▶ In a sequence containing <a>, only the subsequence (suffixes) prefixed with the first occurrence of <a> should be considered.
- ▶ The collection of such subsequences is called *<a>-projected 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>

min_sup count=2

Next: recursively mine *<a>-projected database* to find frequent sequences in the projected database.

Finding Freq. Seq. with Prefix <a> (Cont'd)

- ▶ Find local frequent length-1 sequences by scanning <a>-projected database
 - ▶ *<a>-projected database*: <(abc)(ac)d(cf)>, min_sup count=2
<(_d)c(bc)(ae)>, <(_b)(df)gh>, <(_f)abe>
 - ▶ *Local frequent length-1 sequences*:
<a>:2, :4, <_b>:2, <c>:4, <d>:2, <f>:2
- ▶ Generate all the length-2 freq. seq. having prefix <a>:
<aa>:2, <ab>:4, <(ab)>:2, <ac>:4, <ad>:2, <af>:2
- ▶ Further partition frequent sequences with prefix <a> into 6 subsets
 - ▶ Having prefix <aa>;
 - ▶ Having prefix <ab>;
 - ▶ ...
 - ▶ Having prefix <af>

Finding Freq. Seq. with Prefix <aa>

- ▶ From <a>-projected database:
<(abc)(ac)d(cf)>, <(_d)c(bc)(ae)>, <(_b)(df)cb>, <(_f)cbc>
- ▶ Generate <aa>-projected database:
<(_bc))(ac)d(cf)> <(_e)>
- ▶ No local frequent items, stop growing prefix <aa>.

min_sup count=2

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

Finding Freq. Seq. with Prefix <ab>

- ▶ From <a>-projected database:
<(abc)(ac)d(cf)>, <(_d)c(bc)(ae)>, <(_b)(df)cb>, <(_f)cbc>
- ▶ Generate <ab>-projected database:
<(_c)(ac)d(cf)>, <(_c)(ae)>, <c>
- ▶ Local frequent length-1 sequences
▶ <a>:2, <c>:2, <(_c)>:2
- ▶ Generate length-3 freq. seq. with prefix <ab>:
▶ <aba>:2, <abc>:2, <a(bc)>:2
- ▶ Further partition the set of freq. Seq. prefixed with <ab>:
 - ▶ Sequences with prefix <aba>
 - ▶ Sequences with prefix <abc>
 - ▶ Sequences with prefix <a(bc)>

Finding Freq. Seq. with Prefix <aba>

- ▶ From <ab>-projected database:
 <(_c) (ac)d(cf)>, <(_c)(ae)>, <c>
- ▶ Generate <aba>-projected database.
 - ▶ <(_c)d(cf)>, <(e)>
- ▶ No local frequent length-1 sequence, stop growing <aba>.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Finding Freq. Seq. with Prefix <abc>

- ▶ From <ab>-projected database:
 <(_c) (ac)d(cf)>, <(_c)(ae)>, <c>
- ▶ Generate <abc>-projected database.
 - ▶ <d(cf)> <https://powcoder.com>
- ▶ No local frequent length-1 sequence, stop growing <abc>.

Finding Freq. Seq. with Prefix $\langle a(bc) \rangle$

- ▶ From $\langle ab \rangle$ -projected database:
 - $\langle (_c) (ac)d(cf) \rangle, \langle (_c)(ae) \rangle, \langle c \rangle$
- ▶ Generate $\langle a(bc) \rangle$ -projected database:
 - ▶ $\langle (ac)d(cf) \rangle, \langle (ae) \rangle$
- ▶ Local frequent length-1 sequences:
 - ▶ $\langle a \rangle:2$
- ▶ Generate length-4 freq. seq. with prefix $\langle a(bc) \rangle$:
 - ▶ $\langle a(bc)a \rangle:2$
- ▶ Further mining $\langle a(bc)a \rangle$ -projected database returns no frequent sequence prefixed with $\langle a(bc)a \rangle$.
- ▶ This finishes generating all the patterns prefixed with $\langle ab \rangle$.

Completeness of PrefixSpan

SDB

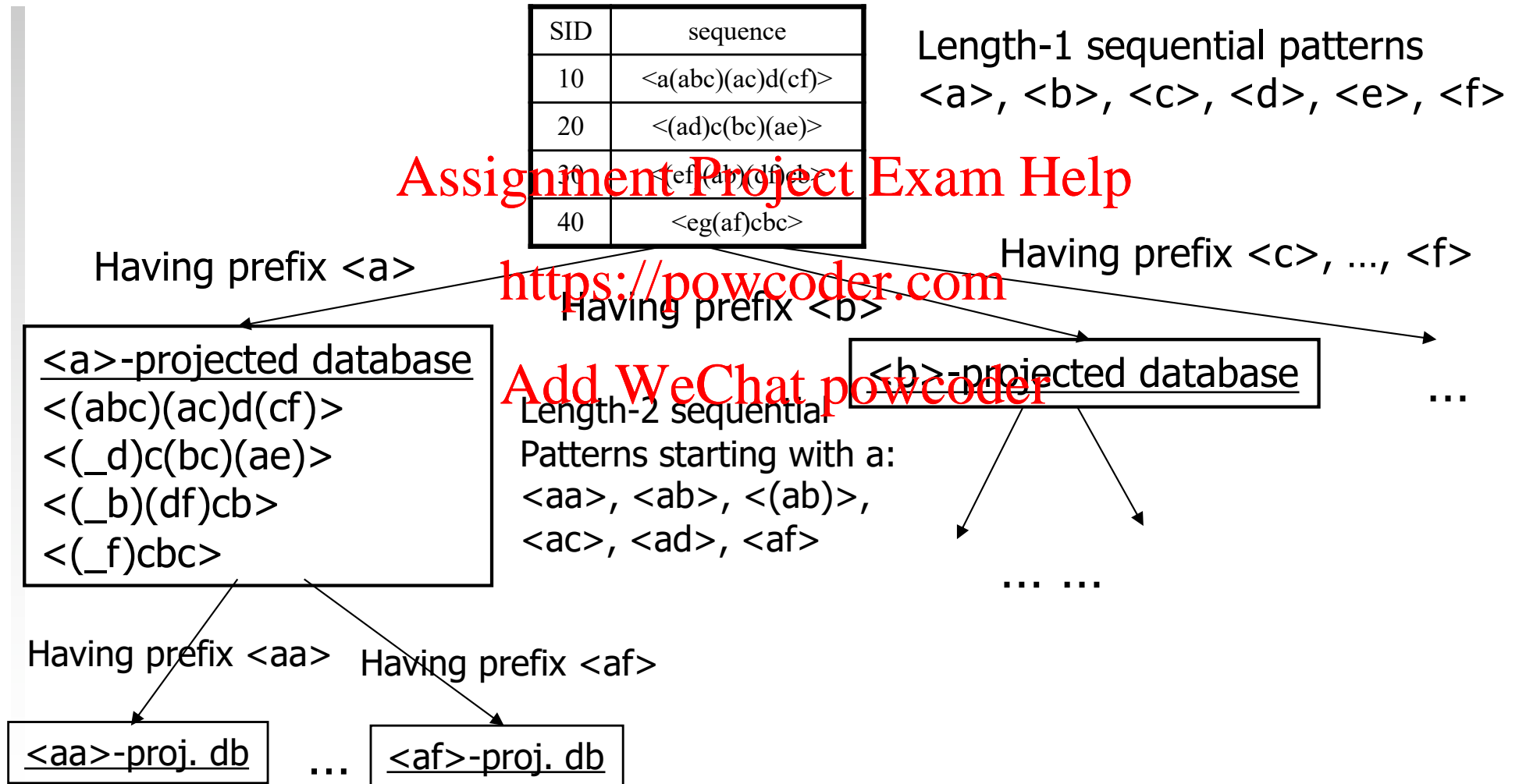
SID	sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<ef(ab)(df)cb>
40	<eg(af)cbc>

Length-1 sequential patterns
<a>, , <c>, <d>, <e>, <f>

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



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*

Speed-up by Pseudo-projection

- ▶ Major cost of PrefixSpan: projection
 - ▶ Suffixes of a sequence often appear repeatedly in recursively projected databases

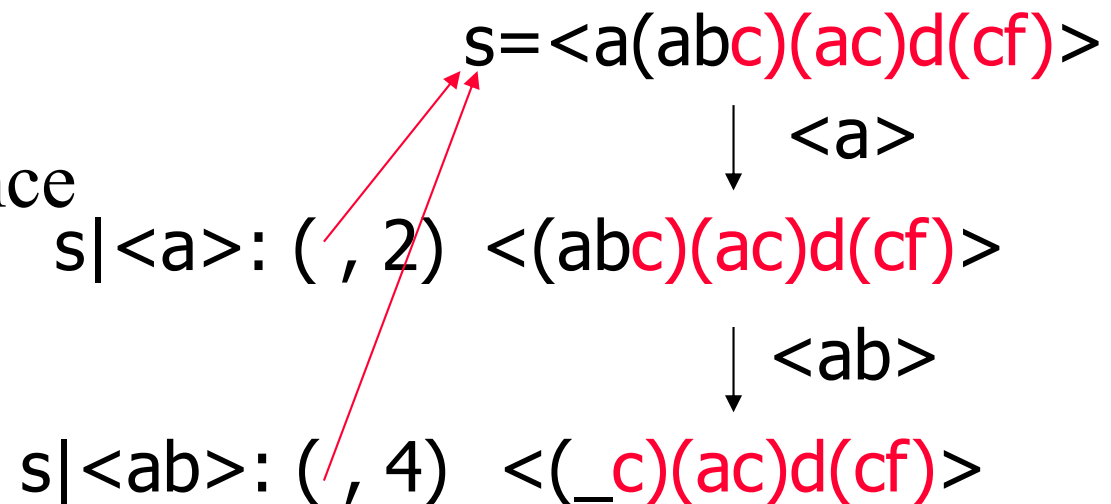
▶ $\langle (abc)(ac)d(cf) \rangle$ appears in $\langle a \rangle$ -projected database

▶ $\langle (_c)(ac)d(cf) \rangle$ appears in $\langle ab \rangle$ -projected database

- ▶ When (projected) database can be held in main memory, use pointers to form projections

▶ Pointer to the sequence

▶ Offset of the suffix



Pseudo-Projection vs. Physical Projection

- ▶ Pseudo-projection avoids physically copying suffixes
 - ▶ 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 projected database fits in memory

Experiments and Performance Analysis

- ▶ Comparing PrefixSpan with GSP, FreeSpan and SPADE in large databases
 - ▶ GSP (IBM Almaden, Srikant & Agrawal EDBT'96)
<https://powcoder.com>
 - ▶ FreeSpan (J. Han J. Pei, B. Mortazavi-Asi, Q. Chen, U. Dayal, M.C. Hsu, KDD '00)
Add WeChat powcoder
 - ▶ SPADE (Zaki, 01)
- ▶ PrefixSpan is fastest and scalable.

Problem of Sequential Pattern Mining

- ▶ GSP and PrefixSpan finds all the sequences that satisfy the support threshold.
- ▶ A long frequent sequence contains a combinatorial number of frequent subsequences:
 - ▶ For a length-100 sequential pattern, there exist $2^{100} - 1$ nonempty subsequences.
- ▶ Problem: too many patterns are generated.

Solutions

- ▶ Mining maximal or closed sequential patterns
 - ▶ Maximal sequential pattern:
 - ▶ A frequent sequence s is *maximal* if there exists no frequent super-sequence of s .
 - ▶ Closed sequential pattern:
 - ▶ A sequence s is *closed* if there exists no super-sequence of s with the same support as s .
 - ▶ *CloSpan* (Yan, Han and Afshar, 2003): an efficient closed sequential pattern mining method.

Solutions (Cont'd)

- ▶ Constraint-based mining of sequential patterns
 - ▶ Constraints
 - ▶ Item constraint
 - ▶ Find patterns containing a, b, c, but no d.
 - ▶ Length constraint
 - ▶ Find patterns having at least (or at most) 10 items
 - ▶ Super-pattern constraint
 - ▶ Find patterns that contain <(PC)(Digital-camera)>
 - ▶ Aggregate constraint
 - ▶ Find patterns the average price of whose items is over \$100.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Solutions (Cont'd)

- ▶ More constraints
 - ▶ Regular expression constraint
 - ▶ `<a*{bb|(bc)d|dd}>`
 - ▶ Duration constraint
 - ▶ Find patterns of events about +/- 24 hours of a shooting
 - ▶ Gap constraint
 - ▶ Find purchasing patterns such that the gap between each consecutive purchases is less than 1 month.

Solutions (Cont'd)

► Properties of constraints

► Anti-monotonic constraint

- If a sequence s satisfies constraint C , so does every non-empty subsequence of s

Assignment Project Exam Help

- Examples: $\text{support}(s) \geq 5\%$, $\text{length}(s) < 10$

<https://powcoder.com>

► Monotonic constraint

- If a sequence s satisfies constraint C , so does every super sequence of s .

Add WeChat powcoder

- Examples: $\text{length}(s) \geq 10$, super pattern constraints.

► A systematic study on constraint-based sequential pattern mining:

- J. Pei, J. Han, and W. Wang. "Mining Sequential Patterns with Constraints in Large Databases". *In Proceedings of CIKM'02*.

More Recent Research on Frequent Pattern Mining

- ▶ Mining other kinds of frequent patterns
 - ▶ Frequent tree/graph mining
 - ▶ Find common structural components
 - ▶ Examples of applications
 - ▶ XML documents can be modeled as trees
 - ▶ Molecule or biochemical structures can be modeled as graphs
 - ▶ Web connection structures can be modeled as graphs
- ▶ Finding frequent patterns from data streams
 - ▶ Only one scan of DB is allowed.

More Recent Research on Frequent Pattern Mining (Cont'd)

- ▶ Finding high-utility patterns (either itemsets or sequences) **Assignment Project Exam Help**
 - ▶ Consider the quantity of the item inside a transaction. <https://powcoder.com>
 - ▶ Consider the importance (e.g., price) of an item **Add WeChat powcoder**
 - ▶ Find patterns whose revenue is at least, e.g., \$1000

Frequent Pattern Mining Resources

- ▶ Web site: <http://fimi.cs.helsinki.fi/> contains some frequent itemset mining implementations, datasets and papers.

<https://powcoder.com>

Add WeChat powcoder

- ▶ SPMF: open-source data mining library containing frequent sequence and itemset mining:

<http://www.philippe-fournier-viger.com/spmf/>

Next Class

- ▶ Decision tree learning (Sections 8.1 and 8.2 in Chapter 8)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder