



Sequential Pattern Mining

Lecture Notes for Chapter 7 – Introduction to Data Mining
Tan, Steinbach, Kumar

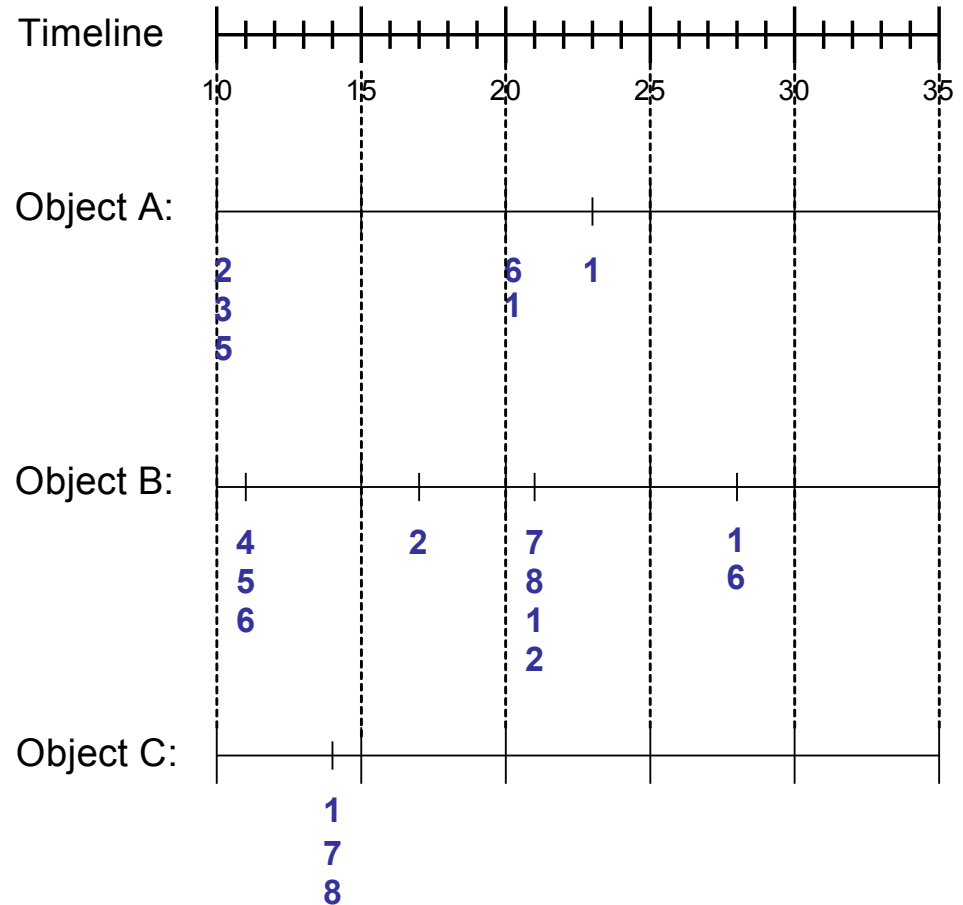
From itemsets to sequences

- Frequent itemsets and association rules focus on transactions and the items that appear there
- Databases of transactions usually have a temporal information
 - Sequential pattern exploit it
- Example data:
 - Market basket transactions
 - Web server logs
 - Tweets
 - Workflow production logs

Sequence Data

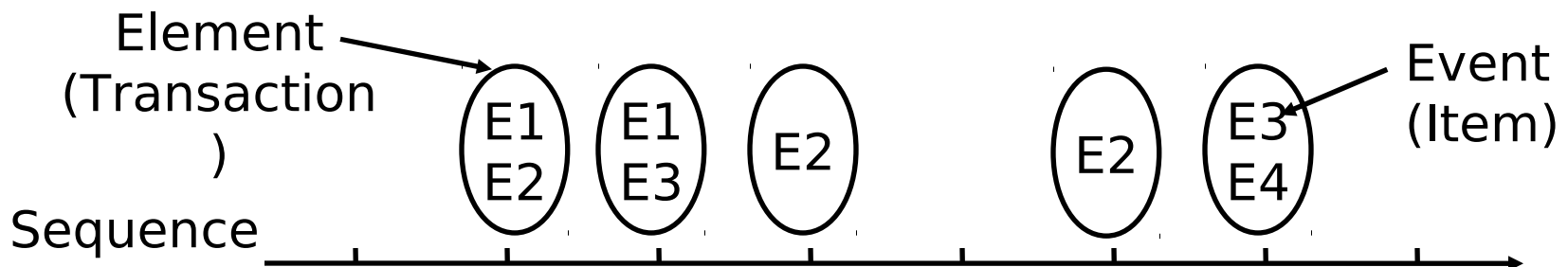
Sequence Database:

Object	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



Examples of Sequence Data

Sequence Database	Sequence	Element (Transaction)	Event (Item)
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



Formal Definition of a Sequence

- A sequence is an ordered list of elements (transactions)

$$S = \langle e_1 e_2 e_3 \dots \rangle$$

- Each element contains a collection of events (items)

$$e_i = \{i_1, i_2, \dots, i_k\}$$

- Each element is attributed to a specific time or location

- Length of a sequence, $|s|$, is given by the number of elements of the sequence
- A k -sequence is a sequence that contains k events (items)

Examples of Sequence

- Web sequence:

< {Homepage} {Electronics} {Digital Cameras} {Canon Digital Camera}
{Shopping Cart} {Order Confirmation} {Return to Shopping} >

- Sequence of initiating events causing the nuclear accident at 3-mile Island:

(http://stellar-one.com/nuclear/staff_reports/summary_SOE_the_initiating_event.htm)

< {clogged resin} {outlet valve closure} {loss of feedwater}
{condenser polisher outlet valve shut} {booster pumps trip}
{main waterpump trips} {main turbine trips} {reactor pressure increases}>

- Sequence of books checked out at a library:

<{Fellowship of the Ring} {The Two Towers} {Return of the King}>

Formal Definition of a Subsequence

- A sequence $\langle a_1 a_2 \dots a_n \rangle$ is contained in another sequence $\langle b_1 b_2 \dots b_m \rangle$ ($m \geq n$) if there exist integers $i_1 < i_2 < \dots < i_n$ such that $a_1 \subseteq b_{i_1}$, $a_2 \subseteq b_{i_2}$, ..., $a_n \subseteq b_{i_n}$

Data sequence	Subsequence	Contain?
$\langle \{2,4\} \{3,5,6\} \{8\} \rangle$	$\langle \{2\} \{3,5\} \rangle$	Yes
$\langle \{1,2\} \{3,4\} \rangle$	$\langle \{1\} \{2\} \rangle$	No
$\langle \{2,4\} \{2,4\} \{2,5\} \rangle$	$\langle \{2\} \{4\} \rangle$	Yes

- The support of a subsequence w is defined as the fraction of data sequences that contain w
- A *sequential pattern* is a frequent subsequence (i.e., a subsequence whose support is $\geq \text{minsup}$)

Sequential Pattern Mining: Definition

- Given:
 - a database of sequences
 - a user-specified minimum support threshold, *minsup*

- Task:
 - Find all subsequences with support \geq *minsup*

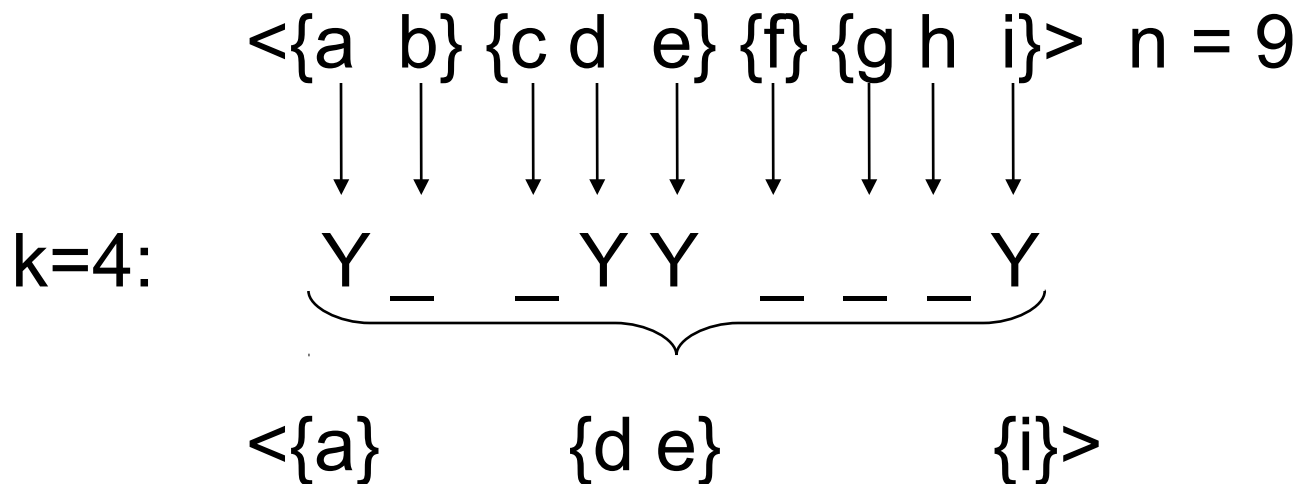
Sequential Pattern Mining: Challenge

- Given a sequence: $\langle \{a\ b\} \{c\ d\ e\} \{f\} \{g\ h\ i\} \rangle$

- Examples of subsequences:

$\langle \{a\} \{c\ d\} \{f\} \{g\} \rangle$, $\langle \{c\ d\ e\} \rangle$, $\langle \{b\} \{g\} \rangle$, etc.

- How many k -subsequences can be extracted from a given n -sequence?



Answer:

$$\binom{n}{k} = \binom{9}{4} = 126$$

Sequential Pattern Mining: Example

Object	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

Minsup = 50%

Examples of Frequent Subsequences:

< {1,2} > s=60%
< {2,3} > s=60%
< {2,4}> s=80%
< {3} {5}> s=80%
< {1} {2} > s=80%
< {2} {2} > s=60%
< {1} {2,3} > s=60%
< {2} {2,3} > s=60%
< {1,2} {2,3} > s=60%

Exercises

- find instances/occurrence of the following patterns

$\langle \{C\} \{H\} \{C\} \rangle$

$\langle \{A\} \{F\} \rangle$

$\langle \{A\} \{A\} \{D\} \rangle$

$\langle \{A\} \{A,B\} \{F\} \rangle$

- in the input sequence below

$\langle \begin{array}{cccccccc} \{A,C\} & \{C,D\} & \{F,H\} & \{A,B\} & \{B,C,D\} & \{E\} & \{A,B,D\} & \{F\} \\ t=0 & t=1 & t=2 & t=3 & t=4 & t=5 & t=6 & t=7 \end{array} \rangle$

Exercises

- find instances/occurrence of the following patterns

$\langle \{C\} \{H\} \{C\} \rangle$

$\langle \{A\} \{B\} \rangle$

$\langle \{C\} \{C\} \{E\} \rangle$

$\langle \{A\} \{E\} \rangle$

- in the input sequence below

$\langle \begin{array}{ccccccc} \{A,C\} & \{C,D,E\} & \{F\} & \{A,H\} & \{B,C,D\} & \{E\} & \{A,B,D\} \\ t=0 & t=1 & t=2 & t=3 & t=4 & t=5 & t=6 \end{array} \rangle$

Extracting Sequential Patterns

- Given n events: $i_1, i_2, i_3, \dots, i_n$
- Candidate 1-subsequences:
 $\langle \{i_1\} \rangle, \langle \{i_2\} \rangle, \langle \{i_3\} \rangle, \dots, \langle \{i_n\} \rangle$
- Candidate 2-subsequences:
 $\langle \{i_1, i_2\} \rangle, \langle \{i_1, i_3\} \rangle, \dots, \langle \{i_1\} \{i_1\} \rangle, \langle \{i_1\} \{i_2\} \rangle, \dots, \langle \{i_{n-1}\} \{i_n\} \rangle$
- Candidate 3-subsequences:
 $\langle \{i_1, i_2, i_3\} \rangle, \langle \{i_1, i_2, i_4\} \rangle, \dots, \langle \{i_1, i_2\} \{i_1\} \rangle, \langle \{i_1, i_2\} \{i_2\} \rangle, \dots,$
 $\langle \{i_1\} \{i_1, i_2\} \rangle, \langle \{i_1\} \{i_1, i_3\} \rangle, \dots, \langle \{i_1\} \{i_1\} \{i_1\} \rangle, \langle \{i_1\} \{i_1\} \{i_2\} \rangle,$
...

Generalized Sequential Pattern (GSP)

- **Step 1:**

- Make the first pass over the sequence database D to yield all the 1-element frequent sequences

- **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 k -sequences that contain infrequent $(k-1)$ -subsequences

- **Support Counting:**

- ◆ Make a new pass over the sequence database D to find the support for these candidate sequences

- **Candidate Elimination:**

- ◆ Eliminate candidate k -sequences whose actual support is less than $minsup$

Candidate Generation

- Base case ($k=2$):
 - Merging two frequent 1-sequences $\langle\{i_1\}\rangle$ and $\langle\{i_2\}\rangle$ will produce two candidate 2-sequences: $\langle\{i_1\} \{i_2\}\rangle$ and $\langle\{i_1 i_2\}\rangle$
- General case ($k>2$):
 - A frequent $(k-1)$ -sequence w_1 is merged with another frequent $(k-1)$ -sequence w_2 to produce a candidate k -sequence if the subsequence obtained by removing the first event in w_1 is the same as the subsequence obtained by removing the last event in w_2
 - ◆ The resulting candidate after merging is given by the sequence w_1 extended with the last event of w_2 .
 - If the last two events in w_2 belong to the same element, then the last event in w_2 becomes part of the last element in w_1
 - Otherwise, the last event in w_2 becomes a separate element appended to the end of w_1

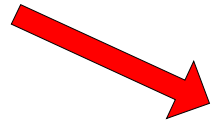
Candidate Generation Examples

- Merging the sequences
 $w_1 = \langle \{1\} \{2\ 3\} \{4\} \rangle$ and $w_2 = \langle \{2\ 3\} \{4\ 5\} \rangle$
will produce the candidate sequence $\langle \{1\} \{2\ 3\} \{4\ 5\} \rangle$ because the last two events in w_2 (4 and 5) belong to the same element
- Merging the sequences
 $w_1 = \langle \{1\} \{2\ 3\} \{4\} \rangle$ and $w_2 = \langle \{2\ 3\} \{4\} \{5\} \rangle$
will produce the candidate sequence $\langle \{1\} \{2\ 3\} \{4\} \{5\} \rangle$ because the last two events in w_2 (4 and 5) do not belong to the same element
- We do not have to merge the sequences
 $w_1 = \langle \{1\} \{2\ 6\} \{4\} \rangle$ and $w_2 = \langle \{1\} \{2\} \{4\ 5\} \rangle$
to produce the candidate $\langle \{1\} \{2\ 6\} \{4\ 5\} \rangle$ because if the latter is a viable candidate, then it can be obtained by merging w_1 with
 $\langle \{1\} \{2\ 6\} \{5\} \rangle$

GSP Example

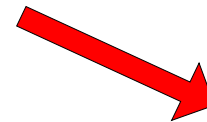
Frequent 3-sequences

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



Candidate Generation

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



Candidate Pruning

< {1} {2 5} {3} >

GSP Exercise

- Given the following dataset of sequences

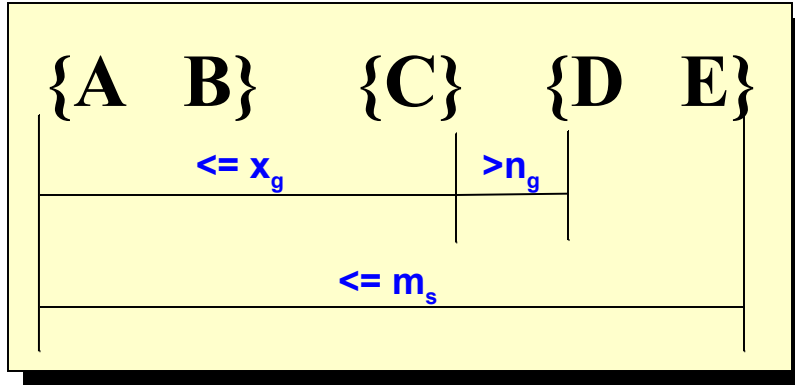
ID	Sequence
1	a b → a → b
2	b → a → c d
3	a → b
4	a → a → b d

- Generate sequential patterns if $\text{min_sup} = 35\%$

GSP Exercise - solution

Sequential pattern	Support
a	100 %
b	100 %
d	50 %
a → a	50 %
a → b	75 %
a → d	50 %
b → a	50 %
a → a → b	50 %

Timing Constraints (I)



x_g : max-gap

n_g : min-gap

m_s : maximum span

$$x_g = 2, n_g = 0, m_s = 4$$

Data sequence	Subsequence	Contain?
$\langle \{2,4\} \{3,5,6\} \{4,7\} \{4,5\} \{8\} \rangle$	$\langle \{6\} \{5\} \rangle$	Yes
$\langle \{1\} \{2\} \{3\} \{4\} \{5\} \rangle$	$\langle \{1\} \{4\} \rangle$	No
$\langle \{1\} \{2,3\} \{3,4\} \{4,5\} \rangle$	$\langle \{2\} \{3\} \{5\} \rangle$	Yes
$\langle \{1,2\} \{3\} \{2,3\} \{3,4\} \{2,4\} \{4,5\} \rangle$	$\langle \{1,2\} \{5\} \rangle$	No

Mining Sequential Patterns with Timing Constraints

- Approach 1:
 - Mine sequential patterns without timing constraints
 - Postprocess the discovered patterns

- Approach 2:
 - Modify GSP to directly prune candidates that violate timing constraints
 - Question:
 - ◆ Does Apriori principle still hold?

Apriori Principle for Sequence Data

Object	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

Suppose:

$$x_g = 1 \text{ (max-gap)}$$

$$n_g = 0 \text{ (min-gap)}$$

$$m_s = 5 \text{ (maximum span)}$$

$$\text{minsup} = 60\%$$

$$\langle \{2\} \{5\} \rangle \text{ support} = 40\%$$

but

$$\langle \{2\} \{3\} \{5\} \rangle \text{ support} = 60\%$$

Problem exists because of max-gap constraint

No such problem if max-gap is infinite

Contiguous Subsequences

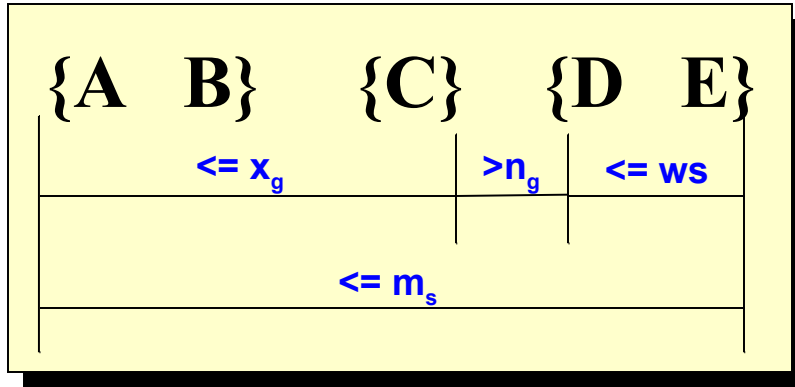
- s is a contiguous subsequence of $w = \langle e_1 \rangle \langle e_2 \rangle \dots \langle e_k \rangle$ if any of the following conditions hold:
 1. s is obtained from w by deleting an item from either e_i or e_k
 2. s is obtained from w by deleting an item from any element e_i that contains more than 2 items
 3. s is a contiguous subsequence of s' and s' is a contiguous subsequence of w (recursive definition)
- Examples: $s = \langle \{1\} \{2\} \rangle$
 - is a contiguous subsequence of $\langle \{1\} \{2\} \{3\} \rangle$, $\langle \{1\} \{2\} \{3\} \{4\} \rangle$, and $\langle \{3\} \{4\} \{1\} \{2\} \{2\} \{3\} \{4\} \rangle$
 - is not a contiguous subsequence of $\langle \{1\} \{3\} \{2\} \rangle$ and $\langle \{2\} \{1\} \{3\} \{2\} \rangle$

Modified Candidate Pruning Step

- Without maxgap constraint:
 - A candidate k -sequence is pruned if at least one of its $(k-1)$ -subsequences is infrequent

- With maxgap constraint:
 - A candidate k -sequence is pruned if at least one of its **contiguous** $(k-1)$ -subsequences is infrequent

Timing Constraints (II)



x_g : max-gap

n_g : min-gap

ws: window size

m_s : maximum span

$x_g = 2$, $n_g = 0$, **ws = 1**, $m_s = 5$

Data sequence	Subsequence	Contain?
$\langle \{2,4\} \{3,5,6\} \{4,7\} \{4,6\} \{8\} \rangle$	$\langle \{3\} \{5\} \rangle$	No
$\langle \{1\} \{2\} \{3\} \{4\} \{5\} \rangle$	$\langle \{1,2\} \{3\} \rangle$	Yes
$\langle \{1,2\} \{2,3\} \{3,4\} \{4,5\} \rangle$	$\langle \{1,2\} \{3,4\} \rangle$	Yes

Modified Support Counting Step

- Given a candidate pattern: $\langle \{a, c\} \rangle$

- Any data sequences that contain

$\langle \dots \{a\ c\} \dots \rangle$,

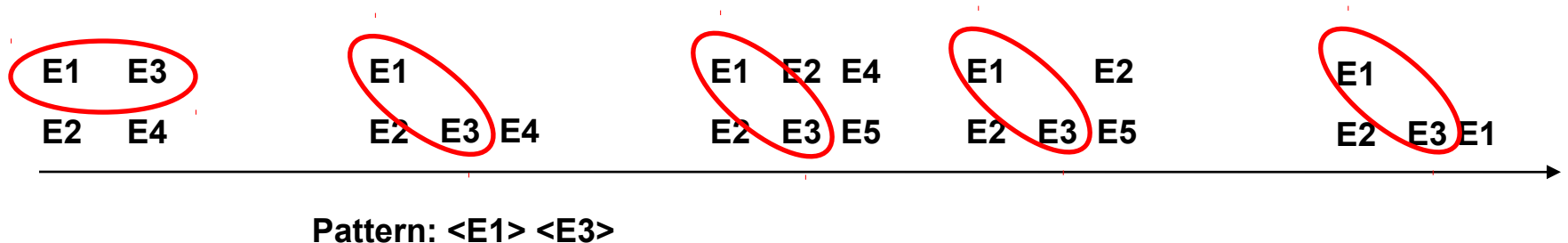
$\langle \dots \{a\} \dots \{c\} \dots \rangle$ (where $\text{time}(\{c\}) - \text{time}(\{a\}) \leq ws$)

$\langle \dots \{c\} \dots \{a\} \dots \rangle$ (where $\text{time}(\{a\}) - \text{time}(\{c\}) \leq ws$)

will contribute to the support count of candidate pattern

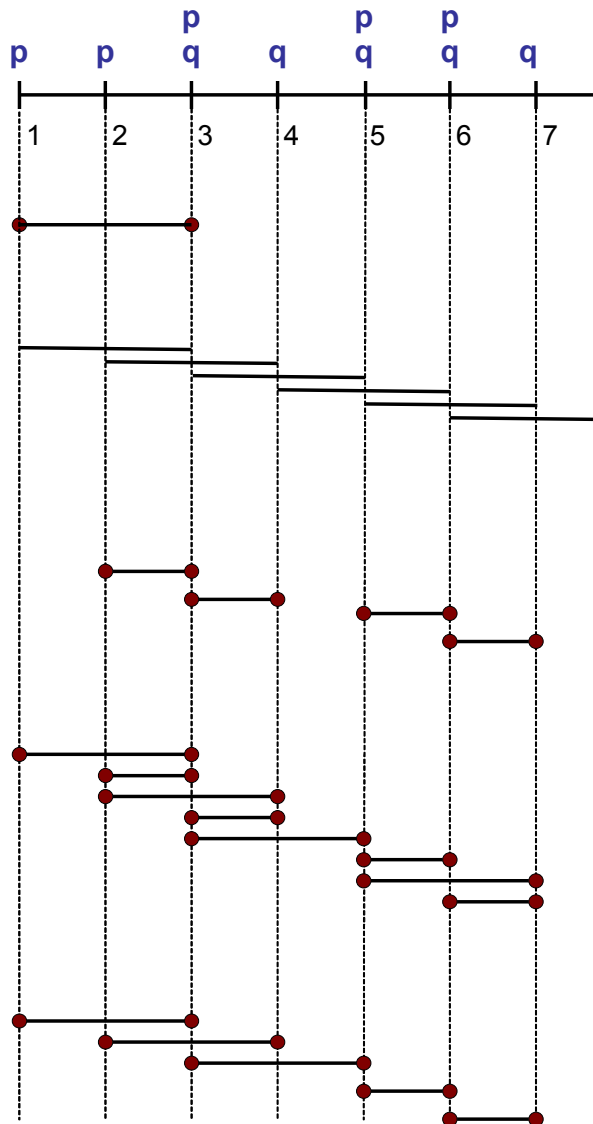
Other Formulation

- In some domains, we may have only one very long time series
 - Example:
 - ◆ monitoring network traffic events for attacks
 - ◆ monitoring telecommunication alarm signals
- Goal is to find frequent sequences of events in the time series
 - This problem is also known as frequent episode mining



General Support Counting Schemes

Object's Timeline



Sequence: (p) (q)

Method Support
 Count

COBJ 1

CWIN 6

CMINWIN 4

CDIST_O 8

CDIST 5

Assume:

$x_g = 2$ (max-gap)

$n_g = 0$ (min-gap)

$ws = 0$ (window size)

$m_s = 2$ (maximum span)