

Ludwig-Maximilians-Universität München Institut für Informatik Lehr- und Forschungseinheit für Datenbanksysteme



# Knowledge Discovery in Databases SS 2016

# Chapter 3b: Sequential Pattern Mining

Lecture: Prof. Dr. Thomas Seidl

Tutorials: Julian Busch, Evgeniy Faerman, Florian Richter, Klaus Schmid

Knowledge Discovery in Databases I: Sequential Pattern Mining





- 1. Sequential Pattern Mining Basics
- 2. Sequential Pattern Mining Algorithms
- 3. Interval-based Sequential Pattern Mining
- 4. Process Mining





- <u>Task 1:</u> find all subsets of items that occur with a specific sequence in many transactions.
  - E.g.: 97% of transactions contain the sequence {jogging  $\rightarrow$  high ECG  $\rightarrow$  sweating}
- <u>Task 2:</u> find all rules that correlate the **order** of one set of items after that of another set of items in the transaction database.
  - E.g.: 72% of users who perform a web search *then* make a long eye gaze over the ads *follow that* by a successful add-click
- The order of the items matters, thus all possible permutations of items must be considered when checking possible frequent sequences, not only the combinations of items
- Applications: data with temporal order (streams), e.g.: bioinformatics, Web mining, text mining, sensor data mining, process mining etc.



# Sequential Pattern Mining vs. Frequent Itemset Mining



- Both can be applied on similar dataset
  - Each customer has a customer id and aligned with transactions.
  - Each transaction has a transaction id and belongs to one customer.
  - Based on the transaction id, each customer also aligned to a transaction sequence.

Cid	Tid	ltem		
1	1	{butter}		
	2	{milk}		
	3	{sugar}		
2	4	{butter, sugar}		
	5	{milk, sugar}		
	6	{butter, milk, sugar}		
	7	{eggs}		
3	8	{sugar}		
	9	{butter, milk}		
	10	{eggs}		
	11	{milk}		

Cid	Item
1	{butter} ,{milk}, {sugar}
2	{butter, sugar}, {milk, sugar}, {butter, milk, sugar}, {eggs}
3	{sugar}, {butter, milk}, {eggs}, {milk}

#### Sequential pattern mining

• The order of items matters

sequences	frequency		
{butter}	4		
{butter, milk}	2		
{butter},{milk}	4		
{milk},{butter}	1		
{butter},{butter,milk}	1		

Frequent itemset mining

 No temporal importance in the order of items happening together

items	frequency
{butter}	4
{milk}	5
{butter, milk}	2



#### Sequential Pattern Mining → Basics





- *Alphabet* Σ is set symbols or characters (denoting items)
- Sequence S = s<sub>1</sub>s<sub>2</sub> ... s<sub>k</sub> is an ordered list of a length |S| = k items where s<sub>i</sub> ∈ Σ is an item at position i also denoted as S[i]
- *K-sequence* is a sequence of length k
- Consecutive subsequence  $\mathbf{R} = r_1 r_2 \dots r_m$  of  $\mathbf{S} = s_1 s_2 \dots s_n$  is also a sequence in  $\Sigma$  such that  $r_1 r_2 \dots r_m = s_j s_{j+1} \dots s_{j+m-1}$  with  $1 \le j \le n m + 1$ . we say  $\mathbf{S}$  contains  $\mathbf{R}$  and denote this by  $\mathbf{R} \subseteq \mathbf{S}$
- In the more general: subsequence R of S we allow for gaps between the items of R, i.e. the items of the subsequence R ⊆ S must have the same order of the ones in S but there can be some other items between them
- A prefix of a sequence **S** is any consecutive subsequence of the form  $S[1:i] = s_1s_2 \dots s_i$  with  $0 \le i \le n$ , S[1:0] is the empty prefix
- A suffix of a sequence S is any consecutive subsequence of the form  $S[i:n] = s_i s_{i+1} \dots s_n$  with  $1 \le i \le n+1$ , S[n+1:n] is the empty suffix





Given a database D = {S<sub>1</sub>, S<sub>2</sub>, ..., S<sub>N</sub>} of N sequences, the support of a sequence
 R in D is defined as the number of sequences in D that contain:

 $sup(\mathbf{R}) = |\{\mathbf{S}_i \in D | \mathbf{R} \subseteq \mathbf{S}_i\}|$ 

- The *relative support* of **R** is the fraction of sequences that contain **R**:  $rsup(\mathbf{R}) = \sup(\mathbf{R}) / N$
- Given a user-defined minSup threshold, a sequence R is called frequent (or sequential) in database D iff : sup(R) ≥ minSup
- A frequent sequence is *maximal* if it is not a subsequence of any other frequent sequence
- A frequent sequence is *closed* if it is not a subsequence of any other frequent sequence with the same support





- Breadth-first search based
  - GSP (Generalized Sequential Pattern) algorithm<sup>1</sup>
  - SPADE<sup>2</sup>

- ...

- Deepth-first search based
  - PrefixSpan<sup>3</sup>
  - SPAM<sup>4</sup>

- ...

<sup>1</sup>Sirkant & Aggarwal: Mining sequential patterns: Generalizations and performance improvements. EDBT 1996 <sup>2</sup>Zaki M J. SPADE: An efficient algorithm for mining frequent sequences[J]. Machine learning, 2001, 42(1-2): 31-60. <sup>3</sup>Pei at. al.: Mining sequential patterns by pattern-growth: PrefixSpan approach. TKDE 2004 <sup>4</sup>Ayres, Jay, et al: Sequential pattern mining using a bitmap representation. SIGKDD 2002.





- GSP (Generalized Sequential Pattern) algorithm<sup>1</sup> performs a level-wise (breadth-first) search within the prefix tree
- Given the set of frequent sequences at level k, generate all possible sequence extensions or candidates at level k + 1
- Uses the Apriori principle (anti-monotonicity)
- Next compute the support of each candidate and prune the ones with sup < minSup</li>
- Stop the search when no more frequent extensions are possible

<sup>1</sup>Sirkant & Aggarwal: *Mining sequential patterns: Generalizations and performance improvements.* EDBT 1996





- The sequence search space can be organized in a prefix search tree
- The root (Level 0) contains the empty sequence with each item  $x \in \Sigma$  as one of its children
- A node labeled with sequence:  $S = s_1 s_2 \dots s_k$  at Level k has children of the form  $S' = s_1 s_2 \dots s_k s_{k+1}$  at Level k + 1
  - Thus S is a prefix of S' (or S' is an extension of S)
- Example: let  $\Sigma = \{A, C, G, T\}$  and let *D* consist of the following three sequences:

• The prefix search tree is:

IdSequence $S_1$ CAGAAGT $S_2$ TGACAG $S_3$ GAAGT



# **Projection-based sequence mining: Representation**





- Sequence (sup)
- Shaded sequences are infrequent (*minSup* = 3)
- Sequences without sup can be pruned based on infrequent subsequences





- For a database D and an item  $s \in \Sigma$ , the *projected database* w.r.t. s is denoted  $D_s$  and is found as follows:
- For each sequence  $S_i \in D$  do
  - Find the first occurrence of *s* in *S<sub>i</sub>*, say at position *p*
  - $suff_{S_i,s} \leftarrow suffix(S_i)$  starting at position p + 1
  - Remove infrequent items from  $suff_{S_i,s}$
  - $D_s = D_s \cup \text{suff}_{S_i,s}$
- Example for the following *D*:
  - $D_G = \{AAGT, AAG, AAGT\}$



## Sequential Pattern Mining $\rightarrow$ Algorithms





- The *PrefixSpan* algorithm computes the support for only the individual items in the projected databased  $D_s$
- Then performs recursive projections on the frequent items in a depth-first manner
- Initialization:  $D_R \leftarrow D$ ,  $\mathbf{R} \leftarrow \emptyset$ ,  $\mathcal{F} \leftarrow \emptyset$

*PrefixSpan*( $D_R$ , **R**, *minSup*, *F*) For each *s* ∈ Σ such that sup(*s*,  $D_R$ ) ≥ *minSup* do

- $R_s = R + s$  //append s to the end of R
- $\mathcal{F} \leftarrow \mathcal{F} \cup \{(\mathbf{R}_s, \sup(s, D_R))\}$  //calculate the support of s for each  $\mathbf{R}_s$  within  $D_R$
- $D_s \leftarrow \emptyset$  //create projected data for s
- For each  $S_i \in D_R$  do
  - $S'_i \leftarrow \text{projection of } S_i \text{ w.r.t. item } s$
  - Remove an infrequent symbols from  $S_i'$
  - If  $S'_i \neq \emptyset$  then  $D_s = D_s \cup S'_i$
- If  $D_s \neq \emptyset$  then  $PrefixSpan(D_s, \mathbf{R}_s, minSup, \mathcal{F})$



## **PrefixSpan: Example**







## PrefixSpan: Example (cont'd)







## PrefixSpan: Example (cont'd)









- Deals with the more common interval-based items *s* (or events).
- Each event has a starting and an ending time point, where  $t_{s_{start}} < t_{s_{end}}$ .
- Application: Health data analysis, Stock market data analysis, etc.
- Predefined relationships between items are more complex.
  - Point-based relationships: before, after, same time.
  - Interval-based relationships: Allen's relations<sup>1</sup>, End point representation<sup>2</sup>, etc.



<sup>1</sup> Allen: *Maintaining knowledge about temporal intervals.* In Communications of the ACM 1983 <sup>2</sup> Wu, Shin-Yi, and Yen-Liang Chen: *Mining nonambiguous temporal patterns for interval-based events.* TKDE 2007

Sequential Pattern Mining → Interval-based SPM





• *Allen's relationships* is ambiguous when describing interval patterns with more than two events.



TPrefixSpan algorithm<sup>1</sup> converts interval-based sequences into point-based sequences:



- Similar prefix projection mining approach as *PrefixSpan* algorithm.
- Validation checking is necessary in each expanding iteration to make sure that the appended time point can form an interval with a time point in the prefix.

<sup>1</sup> Wu, Shin-Yi, and Yen-Liang Chen: *Mining nonambiguous temporal patterns for interval-based events*. TKDE 2007





• Predefined patterns quantize the relationship of items into limited number of categories:

before | | | | | | | | after

- Timing information ignored
- Accuracy of support affected by noise
- Idea: learn pattern from data set by clustering
  - QTempIntMiner<sup>1</sup>
  - Event Space Miner<sup>2</sup>
  - PIVOTMiner<sup>3</sup>



<sup>1</sup> Guyet, T., & Quiniou, R.: Mining temporal patterns with quantitative intervals. ICDMW 2008

<sup>2</sup> Ruan, G., Zhang, H., & Plale, B.: *Parallel and quantitative sequential pattern mining for large-scale interval-based temporal data*. IEEE Big Data 2014

<sup>3</sup> Hassani M., Lu Y. & Seidl T.: A Geometric Approach for Mining Sequential Patterns in Interval-Based Data Streams. FUZZ-IEEE 2016



## **Process Mining**



- *Processes* consist of (business) actions (sign document, customer call, financial transaction, delivery of goods,...)
- A process instance (case) is a certain sequence of actions in a process e.g. customer call → registration → sending package → invoice
- Processes can be modelled as graphs
- Process Mining creates and enriches models by observing real-world workflows







#### **Process Mining**







"Loring Park desire path" by Jake Krohn https://www.flickr.com/photos/jakekrohn/4143809432 is licensed under a Creative Commons license: http://creativecommons.org/licenses/by/3.0/

- Experience often bends the theoretical assumptions
- Many views, but what is the most useful/appropriate one?
- Manager: Good overview, but too far from the details
- Workers: Great detail knowledge with recent information, but missing links to other departments
- Real world information: Not naturally given, hard to acquire





- Mine multiple sequences of actions to derive a workflow pattern
- Raw material: *event logs*, e.g. produced by ERP systems (enterprise resource planning)
- Each event belongs to a certain case
- Events are labeled with an *activity* and ordered using a *timestamp* at least
- Aims:
  - Improve understanding of "business" workflow
  - Model a simple but precise representation
  - Reveal organizational and social insights





Process Log

Case	Activity	Time	
John	Counter A	11:13	
Mary	Counter B	11:15	
Mary	Receives meal	11:21	
Bob	Enters facility	11:22	
John	Counter S	11:22	
Bob	Leaves facility	11:24	
Mary	Sits down	11:25	
John	Receives meal	11.28	

Sequential Pattern Mining → Process Mining



## **Process Mining – Other Tasks**



- Task 2: Conformance Checking
  - Use previously mined model to judge the validity of a new case
  - Classification task is based on different types of models: procedural models, organizational models, business rules, laws,...
  - Aims: Model reasoning, auditing, security (fraud detection)
- Task 3: Enhancement
  - Business processes change over time, so adapt to new circumstances
  - Evolve models with new data, find deviations
  - Extend/enrich models, using a wider spectrum of data attributes
  - Aims: highlight bottlenecks, combine different models, analyze performance







## **Basis Discovery:** *α***-Algorithm**<sup>1</sup>



- Scan the log for all activities A
- For each pair of activities *a* and *b*, we define the relations
  - a > b iff for some case a is immediately followed by b (direct succession)
  - $a \rightarrow b$  iff a > b and not b > a (causality)
  - a||b iff a > b and b > a (parallelism)
  - a#b iff not a > b and not b > a
- All activities, having only # or  $\rightarrow$  in their row are starting activities. They are collected in  $T_{in}$ .
- Analogously, # or  $\leftarrow$  determine  $T_{out}$ .

Example  $L = \{abcd, acbd, aed\}$ 

	а	b	С	d	е
а	#	$\rightarrow$	$\rightarrow$	#	$\rightarrow$
b	$\leftarrow$	#		$\rightarrow$	#
С	$\leftarrow$		#	$\rightarrow$	#
d	#	$\leftarrow$	$\leftarrow$	#	$\leftarrow$
е	$\leftarrow$	#	#	$\rightarrow$	#

$$T_{in} = \{a\}, T_{out} = \{d\}$$

<sup>1</sup> van der Aalst, W M P and Weijters, A J M M and Maruster, L (2003). "Workflow Mining: Discovering process models from event logs", *IEEE Transactions on Knowledge and Data Engineering*, vol 16





- Prepare a Petri net. The set of transitions is equal to *A*, so each activity represents a transition
- A starting place is created and connected to each node in T<sub>in</sub>
- Also, a final place is created and each node in T<sub>out</sub> is connected to it
- Determine all pairs of sets A and B, such that
  - $\forall a_1, a_2 \in A: a_1 # a_2$
  - $\forall b_1, b_2 \in B: b_1 \# b_2$
  - $\forall a_1 \in A, \forall b_1 \in B: a_1 \rightarrow b_1$
- A place is added in between A and B and connected accordingly



<sup>1</sup> van der Aalst, W M P and Weijters, A J M M and Maruster, L (2003). "Workflow Mining: Discovering process models from event logs", *IEEE Transactions on Knowledge and Data Engineering*, vol 16

## Sequential Pattern Mining → Process Mining







- In this example, valid set pairs are:

   ({a}, {be}), ({a}, {ce}), ({be}, {d}), ({ce}, {d})
- This yields the following Petri net



<sup>1</sup> van der Aalst, W M P and Weijters, A J M M and Maruster, L (2003). "Workflow Mining: Discovering process models from event logs", *IEEE Transactions on Knowledge and Data Engineering*, vol 16

## Sequential Pattern Mining $\rightarrow$ Process Mining