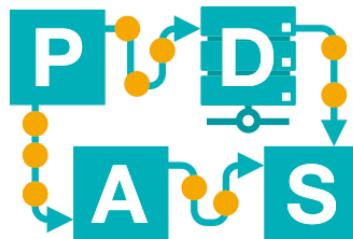


Sequence Mining

Lecture 13

IDS-L13



Chair of Process
and Data Science

RWTH AACHEN
UNIVERSITY

Outline of Today's Lecture

- Short recap: Apriori algorithm
- Event Data
- Mining sequential patterns
- AprioriAll Algorithm
- Extensions



Material

Mining Sequential Patterns

Rakesh Agrawal Ramakrishnan Srikant
IBM Almaden Research Center
650 Harry Road, San Jose, CA 95120

Abstract

We are given a large database of customer transactions, where each transaction consists of customer, transaction time, and the items bought in the transaction. We introduce the problem of mining sequential patterns over such databases. We present three algorithms to solve this problem, and empirically evaluate the performance using synthetic data. Evaluation of the proposed algorithms, AprioriSane and AprioriSome, All have comparable performance, albeit AprioriSome performs a little better when the minimum number of customers that must support a sequential pattern is low. Self-supply experiments show that both AprioriSome and AprioriSane scale well with the number of customer transactions. They also have excellent scalability properties with respect to the number of items in a transaction.

1 Introduction

Database mining is motivated by the decision support problem faced by most large retail organizations. Progress in barcode technology has made it possible for retail stores to automatically capture quantitative amounts of sales data, referred to as the *facts* of sales. A record in such data typically consists of the transaction date and the items bought in the transaction. Very often, data records also contain customer-id, particularly when the purchase has been made using a credit card or a membership card. Calling the sequence of items sold and data using the orders they receive.

We introduce the problem of mining *sequential patterns* over this data. An example of such a pattern is

^{*}Also Department of Computer Science, University of Wisconsin, Madison.

Rakesh Agrawal, Ramakrishnan Srikant: Mining Sequential Patterns. ICDE 1995: 3-14

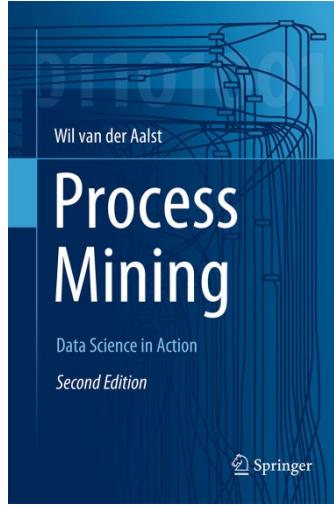
(that customer typically can “Buy Water” then “Enter Supermarket”, and then “Leave” at the end). Note that these events need not be consecutive. Customers who rent some other video in between also support this sequential pattern. Elements of a sequential pattern need not be simple items: “Fitted Sheet and flat sheet and pillow cases”, followed by “comforter”, followed by “linens and pillows” is an example of a sequential pattern in which the elements are sets of items.

Problem Statement: We are given a database D of customer transactions. Each transaction consists of the following fields: customer-id, transaction-time, and a list of items bought in the transaction. A customer has more than one transaction with the same transaction-time. We do not consider quantities of items bought in a transaction: each item is a binary variable representing whether the item was bought or not.

An *itemset* is a non-empty set of items. A *sequence* is an ordered list of itemsets. Without loss of generality, we assume that the set of items is mapped to a set of contiguous integers. We denote an itemset i by $\{i_1, i_2, \dots, i_n\}$, where i_j is an item and a sequence $s = (s_1, s_2, \dots, s_m)$, where s_i is an itemset.

A sequence (a_1, a_2, \dots, a_n) is contained in another sequence (b_1, b_2, \dots, b_m) if there exist integers $i_1 < i_2 < \dots < i_n$ such that $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$. For example, the sequence $(1, 5, 8)$ is contained in $(1, 7, 3, 2, 6, 4, 8, 9)$, since $(1, 5) \subseteq (1, 7, 3, 2, 6, 4, 8, 9)$ and $(8) \subseteq (8)$. However, the sequence $(3, 5)$ is not contained in $(1, 5)$, (and vice versa). The former represents items 3 and 5 being bought one after the other, while the latter represents items 3 and 5 being bought simultaneously. A sequence s is supported if s is not contained in any other itemset.

All the transactions of a customer can together be viewed as a sequence, where each transaction corresponds to a set of items, and the list of



"W. van der Aalst. Process Mining: Data Science in Action. Springer-Verlag, Berlin, 2016"
[\(<http://springer.com/9783662498507>\)](http://springer.com/9783662498507)

The book “Data Mining: Concepts and Techniques (3rd Edition) by Jiawei Han, Micheline Kamber, and Jian Pei. Morgan Kaufmann, June 2011” is rather brief about sequence mining, but many of the general principles are valid also for sequence mining.

Short recap: Apriori algorithm



Apriori algorithm: Basic ingredients

- Let $D \in \mathfrak{B}(\mathcal{P}(\mathcal{I}))$ be a dataset and min_sup a chosen threshold.
- If $A \subseteq B$ and $\text{support}(A) < \text{min_sup}$, then $\text{support}(B) < \text{min_sup}$.
- **Leveling:** $L_k = \{A \subseteq \mathcal{I} | \text{support}(A) \geq \text{min_sup} \wedge |A| = k\}$.

Apriori algorithm: Basic ingredients

- L_k are all frequent itemsets of length k .
- **Construction of candidates:** Any $A = \{a_1, a_2, \dots, a_{k-1}, a_k\} \in L_k$ can be obtained by joining $A' = \{a_1, a_2, \dots, a_{k-2}, a_{k-1}\} \in L_{k-1}$ and $A'' = \{a_1, a_2, \dots, a_{k-2}, a_k\} \in L_{k-1}$.
- Let C_k be the set of all **candidates** created from L_{k-1} in this manner.
- **Prune** the set C_k by removing itemsets that have infrequent subsets.



Algorithm

Algorithm: Apriori. Find frequent itemsets using an iterative level-wise approach based on candidate generation.

Input:

- D , a database of transactions;
- min_sup , the minimum support count threshold.

Output: L , frequent itemsets in D .

Method:

```
(1)    $L_1 = \text{find.frequent.1-itemsets}(D);$ 
(2)   for ( $k = 2; L_{k-1} \neq \emptyset; k++$ ) {
(3)      $C_k = \text{apriori.gen}(L_{k-1});$ 
(4)     for each transaction  $t \in D$  { // scan  $D$  for counts
(5)        $C_t = \text{subset}(C_k, t);$  // get the subsets of  $t$  that are candidates
(6)       for each candidate  $c \in C_t$ 
(7)          $c.\text{count}++;$ 
(8)     }
(9)      $L_k = \{c \in C_k | c.\text{count} \geq min\_sup\}$ 
(10)
(11)    return  $L = \bigcup_k L_k;$ 
```

```
procedure apriori_gen( $L_{k-1}$ :frequent  $(k - 1)$ -itemsets)
(1)   for each itemset  $l_1 \in L_{k-1}$ 
(2)     for each itemset  $l_2 \in L_{k-1}$ 
(3)       if  $(l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2]) \wedge \dots \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$  then {
(4)          $c = l_1 \bowtie l_2;$  // join step: generate candidates
(5)         if  $\text{has\_infrequent\_subset}(c, L_{k-1})$  then
(6)           delete  $c;$  // prune step: remove unfruitful candidate
(7)         else add  $c$  to  $C_k;$ 
(8)       }
(9)   return  $C_k;$ 
```

```
procedure has_infrequent_subset( $c$ : candidate  $k$ -itemset;
 $L_{k-1}$ : frequent  $(k - 1)$ -itemsets); // use prior knowledge
(1)   for each  $(k - 1)$ -subset  $s$  of  $c$ 
(2)     if  $s \notin L_{k-1}$  then
(3)       return TRUE;
(4)   return FALSE;
```

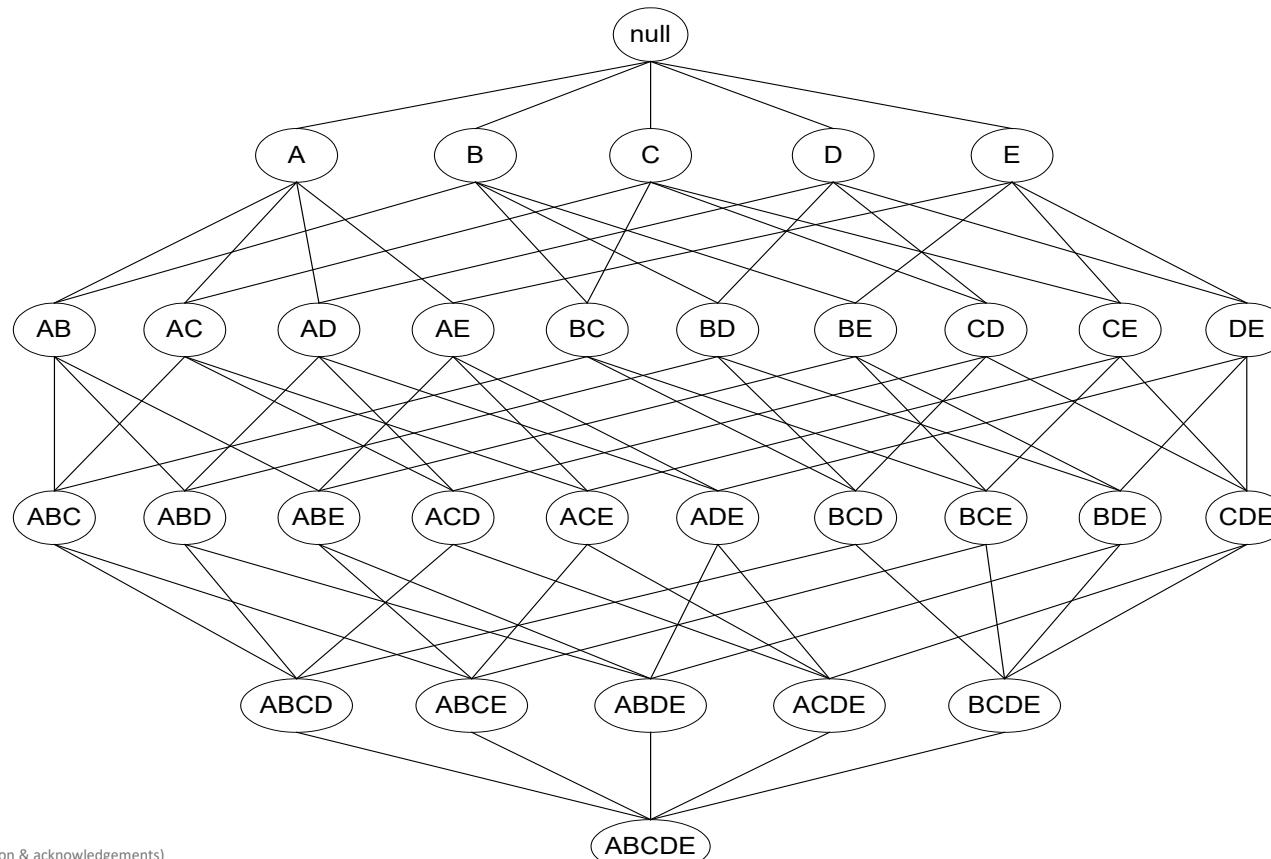
creating candidates
 C_k based on L_{k-1}

testing candidates to
create L_k based on C_k

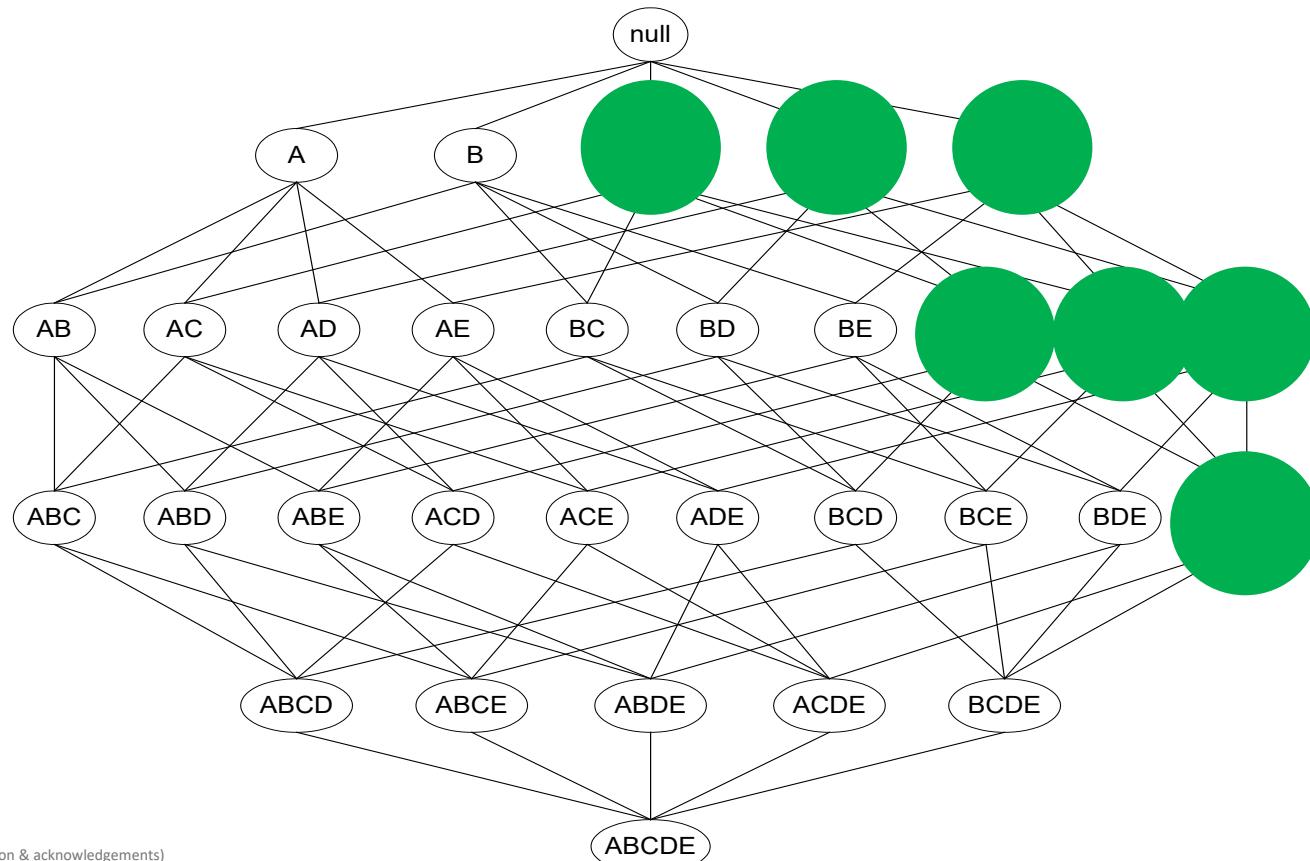
create candidates from
 L_{k-1} (avoiding
duplicates)

prune set of candidates
by checking whether
subsets were frequent
before

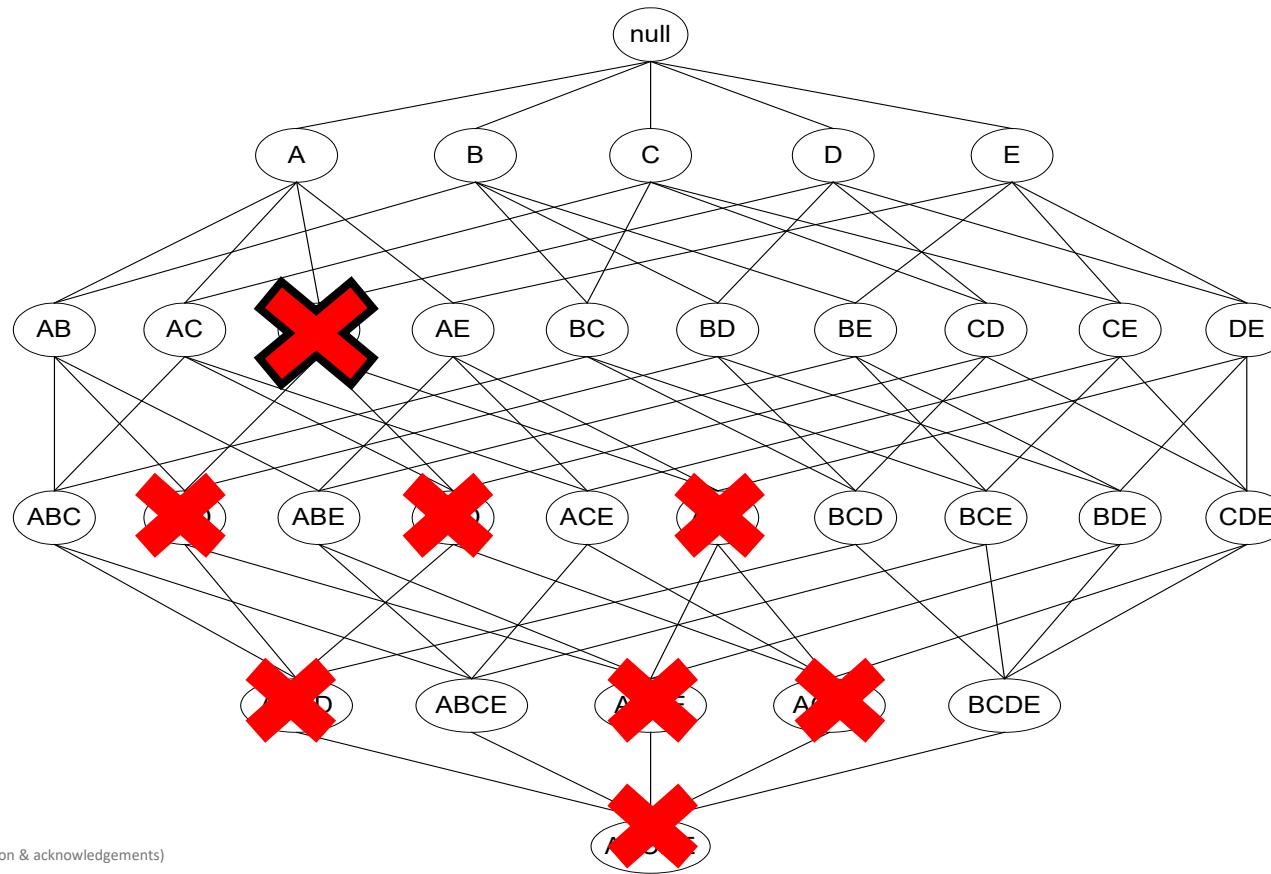
Consider a data set with 5 items



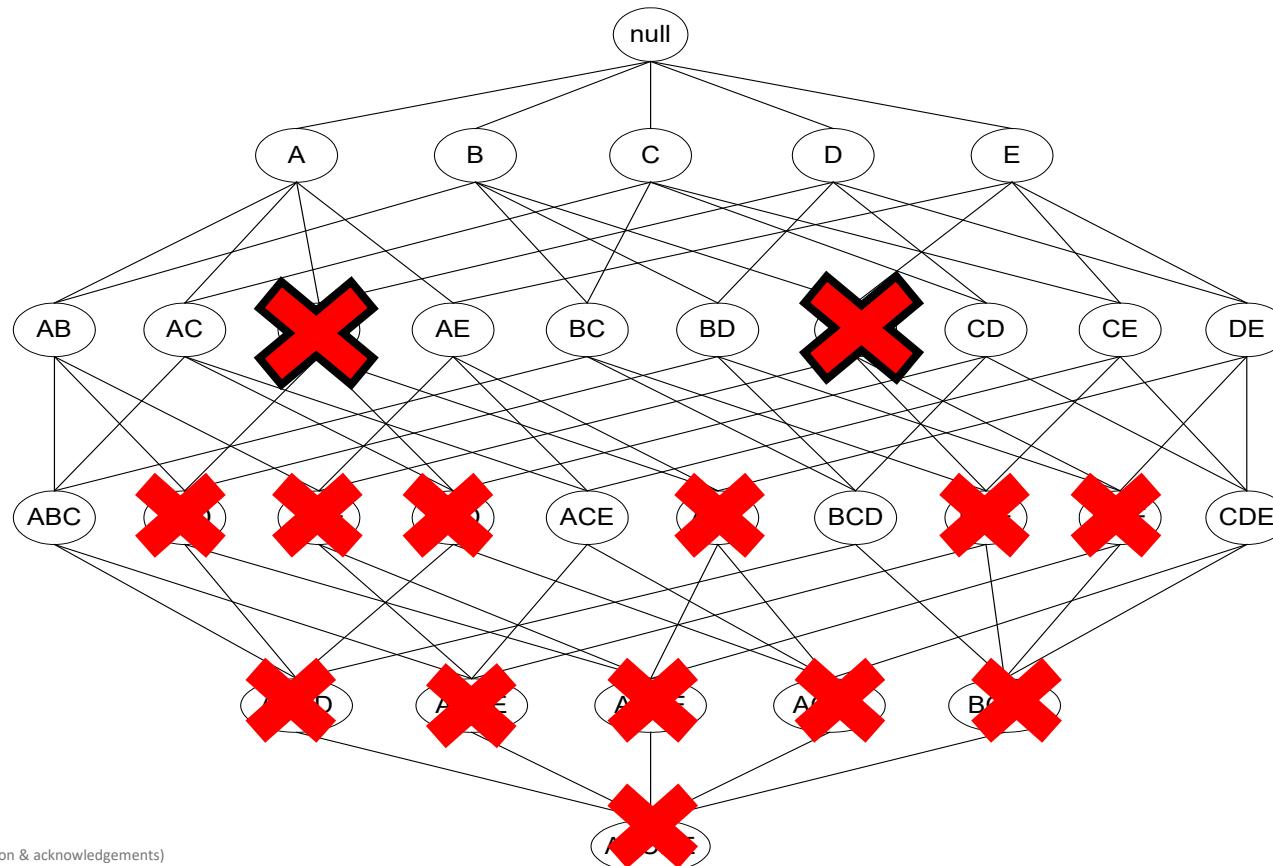
Consider a data set with 5 items



Consider a data set with 5 items



Consider a data set with 5 items



Apriori principle can be used when searching for other “patterns”

- Frequent sequences
- Frequent subgraphs
- Frequent partial orders
- Etc.

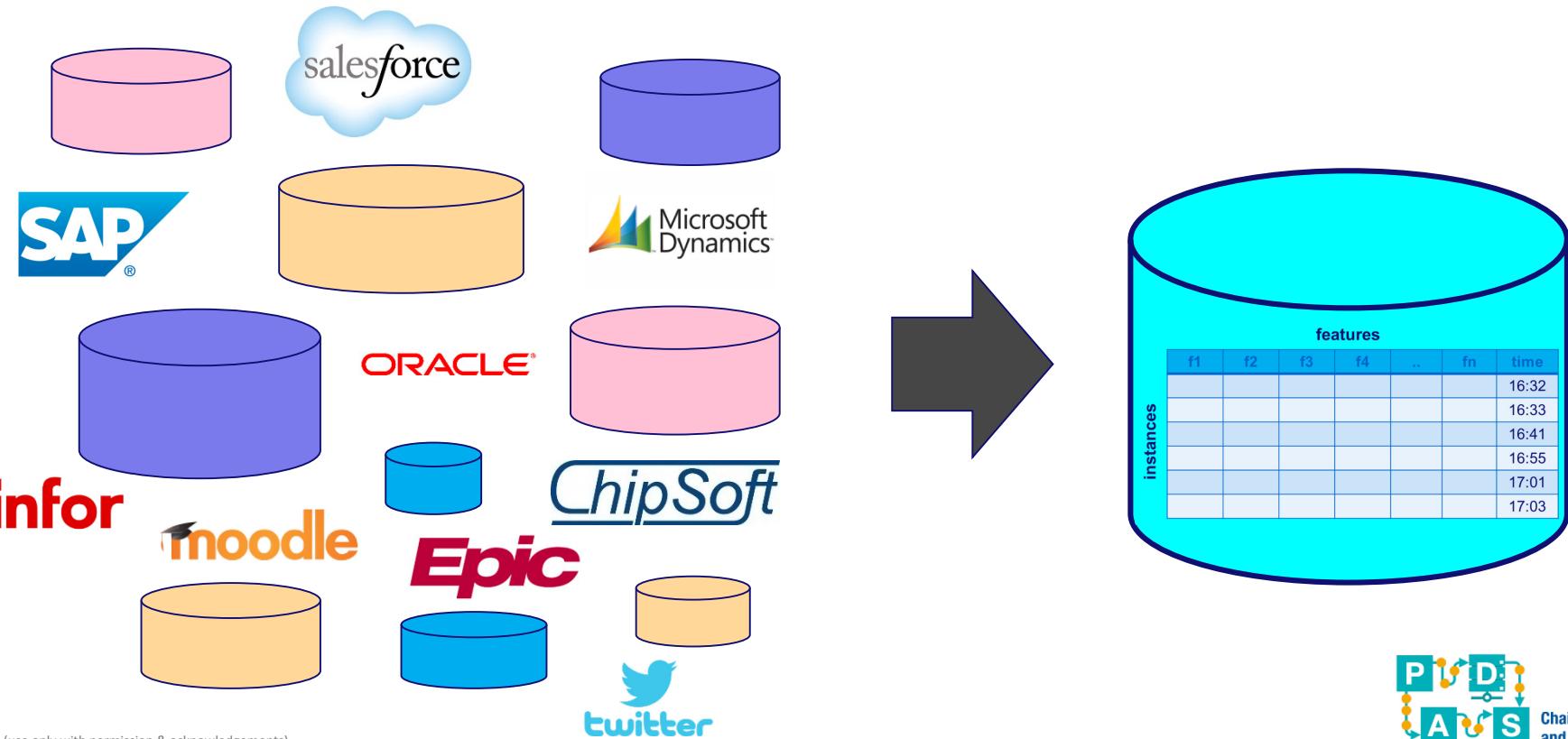
Event Data



Event data: Discrete timestamped events

features							
instances	f1	f2	f3	f4	..	fn	time
							16:32
							16:33
							16:41
							16:55
							17:01
							17:03

Event data are everywhere



Example 1

every row is an event
(here: an exam attempt)

student name	course name	exam date	mark
Peter Jones	Business Information systems	16-1-2014	8
Sandy Scott	Business Information systems	16-1-2014	5
Bridget White	Business Information systems	16-1-2014	9
John Anderson	Business Information systems	16-1-2014	8
Sandy Scott	BPM Systems	17-1-2014	7
Bridget White	BPM Systems	17-1-2014	8
Sandy Scott	Process Mining	20-1-2014	5
Bridget White	Process Mining	20-1-2014	9
John Anderson	Process Mining	20-1-2014	8
...

case id

activity name

timestamp

other data

Example 2

order number	activity	timestamp	user	product	quantity
9901	register order	22-1-2014@09.15	Sara Jones	iPhone5S	1
9902	register order	22-1-2014@09.18	Sara Jones	iPhone5S	2
9903	register order	22-1-2014@09.27	Sara Jones	iPhone4S	1
9901	check stock	22-1-2014@09.49	Pete Scott	iPhone5S	1
9901	ship order	22-1-2014@10.11	Sue Fox	iPhone5S	1
9903	check stock	22-1-2014@10.34	Pete Scott	iPhone4S	1
9901	handle payment	22-1-2014@10.41	Carol Hope	iPhone5S	1
9902	check stock	22-1-2014@10.57	Pete Scott	iPhone5S	2
9902	cancel order	22-1-2014@11.08	Carol Hope	iPhone5S	2
...

case id

activity name

timestamp

resource

other data

Example 3

patient	activity	timestamp	doctor	age	cost
5781	make X-ray	23-1-2014@10.30	Dr. Jones	45	70.00
5541	blood test	23-1-2014@10.18	Dr. Scott	61	40.00
5833	blood test	23-1-2014@10.27	Dr. Scott	24	40.00
5781	blood test	23-1-2014@10.49	Dr. Scott	45	40.00
5781	CT scan	23-1-2014@11.10	Dr. Fox	45	1200.00
5833	surgery	23-1-2014@12.34	Dr. Scott	24	2300.00
5781	handle payment	23-1-2014@12.41	Carol Hope	45	0.00
5541	radiation therapy	23-1-2014@13.57	Dr. Jones	61	140.00
5541	radiation therapy	23-1-2014@13.08	Dr. Jones	61	140.00

...

...

...

...

...

case id

activity name

timestamp

resource

other data

Mandatory attributes for process mining

- **Case id:** the “thing” used to group events (customer, order, student, session, device, suitcase, etc.).
- **Activity name:** description of the event, typically an activity, but it could also be an itemset.
- **Timestamp:** the time at which the event takes place

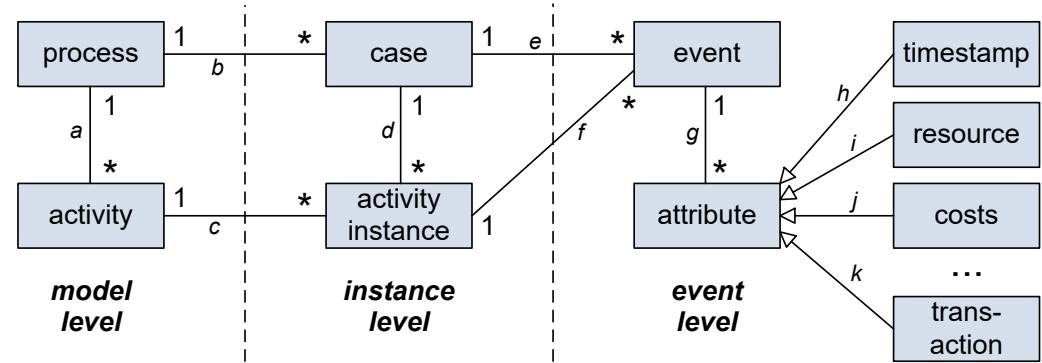
There may be various attributes with a well understood meaning: resource, lifecycle (start, complete, etc.), costs, role, etc. However, events are atomic.



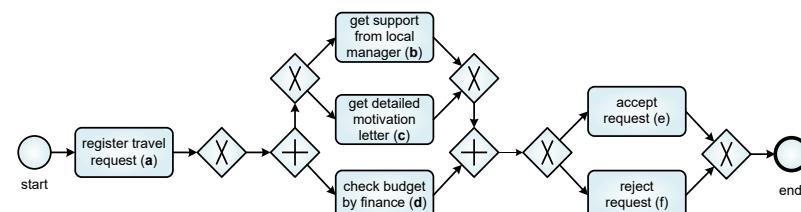
Process mining starts from event data



The screenshot shows the IEEE Xplore Digital Library interface. At the top, it displays 'Access provided by: Universitätsbibliothek der RWTH Aachen' and the IEEE logo. Below the header, there are navigation links for 'Browse', 'My Settings', and 'Get Help'. The main search bar contains the placeholder 'Enter keywords or short phrases (searches metadata only by default)'. Underneath the search bar, there are links for 'Advanced Search' and 'Other Search Options'. The main content area displays the IEEE Standard for eXtensible Event Stream (XES) document. The title is '1849-2016 - IEEE Standard for eXtensible Event Stream (XES) for Achieving Interoperability in Event Logs and Event Streams'. It includes sections for 'Abstract', 'Figures', 'References', 'Citations', 'Keywords', 'Definitions', 'Metrics', 'Versions', and 'Media'. A sidebar on the right lists 'Related Articles' and 'View All'. At the bottom, there is a 'Scope' section with detailed text about the standard's purpose and structure.

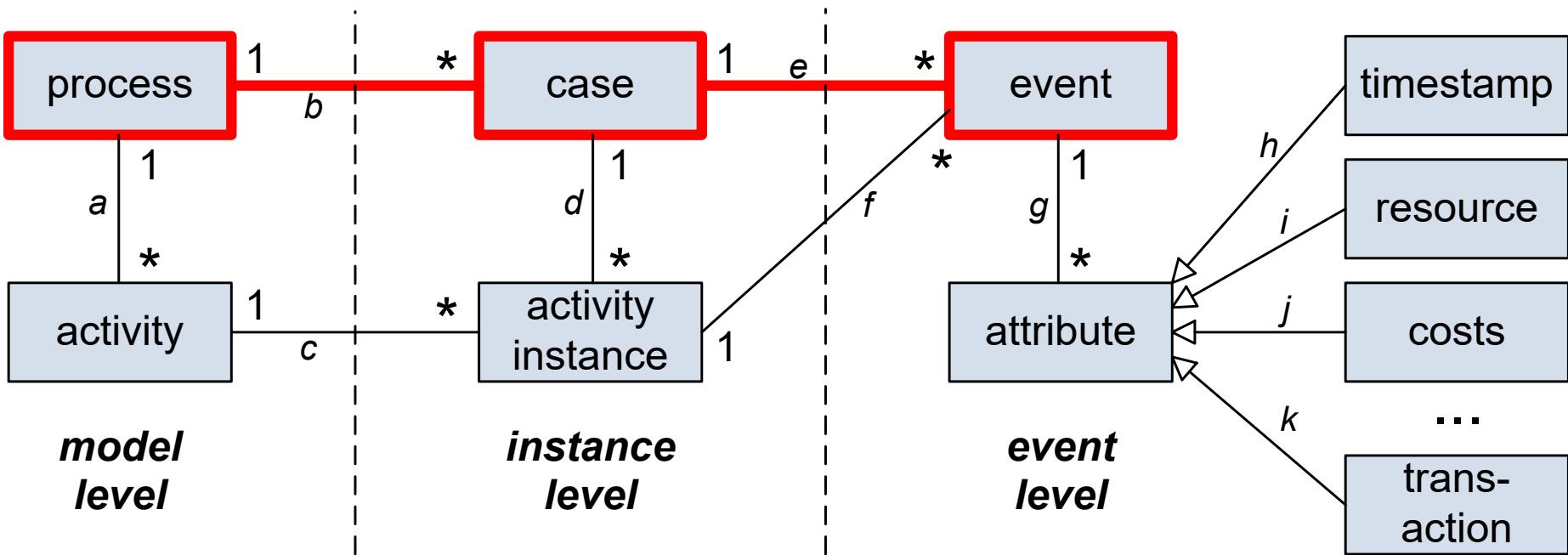


and learns end-to-end process models



Chair of Process
and Data Science

Process - case - event



Events have attributes

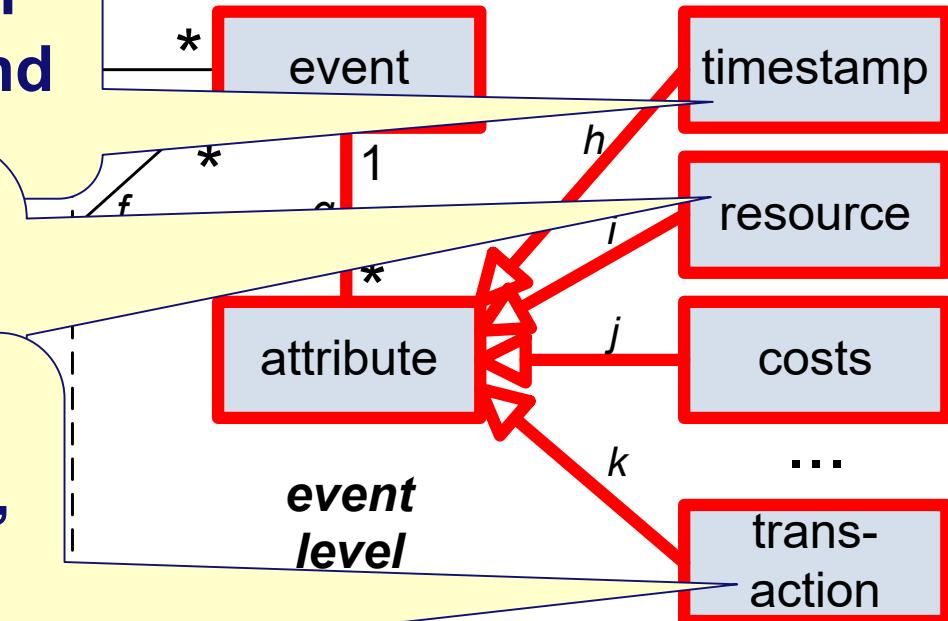
pr

time of event: crucial for
the ordering of events and

~~performance analysis~~

resource associated to
event (person, machine,
~~department, organization~~)

transactional information
(start, suspend, complete,
etc.): activities are often
non-atomic



Typical abstraction used by process mining algorithms: A multiset of activity sequences

order number	activity	timestamp	user	product	quantity
9901	register order	22-1-2014@09.15	Sara Jones	iPhone5S	1
9902	register order	22-1-2014@09.18	Sara Jones	iPhone5S	2
9903	register order	22-1-2014@09.27	Sara Jones	iPhone4S	1
9901	check stock	22-1-2014@09.49	Pete Scott	iPhone5S	1
9901	ship order	22-1-2014@10.11	Sue Fox	iPhone5S	1
9903	check stock	22-1-2014@10.34	Pete Scott	iPhone4S	1
9901	handle payment	22-1-2014@10.41	Carol Hope	iPhone5S	1
9902	check stock	22-1-2014@10.57	Pete Scott	iPhone5S	2
9902	cancel order	22-1-2014@11.08	Carol Hope	iPhone5S	2
...

one activity sequence per order



9901
[<register order, check stock, ship order, handle payment>, <register
9904
9902]

same sequence may appear many times

Typical abstraction used by process mining algorithms: A multiset of activity sequences

patient	activity	timestamp	doctor	age	cost
5781	make X-ray	23-1-2014@10.30	Dr. Jones	45	70.00
5541	blood test	23-1-2014@10.18	Dr. Scott	61	40.00
5833	blood test	23-1-2014@10.27	Dr. Scott	24	40.00
5781	blood test	23-1-2014@10.49	Dr. Scott	45	40.00
5781	CT scan	23-1-2014@11.10	Dr. Fox	45	1200.00
5833	surgery	23-1-2014@12.34	Dr. Scott	24	2300.00
5781	handle payment	23-1-2014@12.41	Carol Hope	45	0.00
5541	radiation therapy	23-1-2014@13.57	Dr. Jones	61	140.00
5541	radiation therapy	23-1-2014@13.08	Dr. Jones	61	140.00
...

one activity sequence
per patient



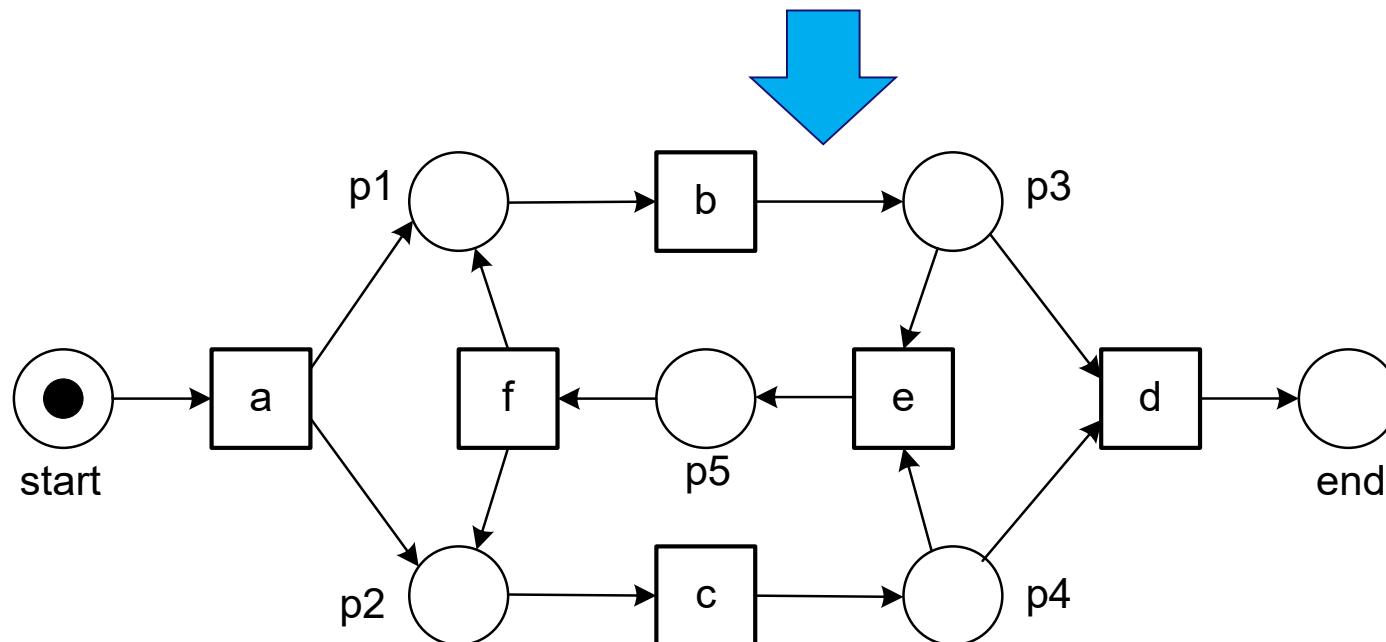
⁵⁷⁸¹ [⟨make X-ray, blood test, handle payment⟩, ⁵⁵⁴¹ ⟨blood test, radiation

same sequence may appear many times

Process discovery

(presented in detail in later lectures)

$$L_2 = [\langle a, b, c, d \rangle^3, \langle a, c, b, d \rangle^4, \langle a, b, c, e, f, b, c, d \rangle^2, \langle a, b, c, e, f, c, b, d \rangle, \\ \langle a, c, b, e, f, b, c, d \rangle^2, \langle a, c, b, e, f, b, c, e, f, c, b, d \rangle]$$



Sequential pattern mining

- **Uses a specific type of event data as input.**
- **For example,**
 - Informal: [ab(cd)e, ab(cd)e, a(cd)e, a(bc)(cde)f].
 - Formal:
 $\langle \{\{a\}, \{b\}, \{c, d\}, \{e\}\}, \{\{a\}, \{b\}, \{c, d\}, \{e\}\}, \{\{a\}, \{c, d\}, \{e\}\}, \{\{a\}, \{b, c\}, \{c, d, e\}, \{f\}\} \rangle$
- **Formally, $D \in \mathfrak{B}((\mathcal{P}(\mathcal{I}))^*)$.**

Mining sequential patterns



Two beautiful/influential papers

(especially considering that they are almost 25 years old)

Fast Algorithms for Mining Association Rules

Rakesh Agrawal

Ramakrishnan Srikant*

IBM Almaden Research Center
650 Harry Road, San Jose, CA 95120

Abstract

We consider the problem of discovering association rules between items in a large database of sales transactions. We present two new algorithms for solving this problem that are much faster than the known algorithms. Other applications include catalog design, add-on sales, and market basket analysis, all based on buying patterns. Empirical evaluation shows that these algorithms outperform the known algorithms by factors ranging from three for small problems to more than an order of magnitude for large problems. We also show that the best fast algorithm presented here can be easily combined into a hybrid algorithm, called AprioriHybrid. Scale-up experiments show that AprioriHybrid scales linearly with the number of transactions. AprioriHybrid also has excellent scale-up properties with respect to the transaction size and the number of items in the database.

1 Introduction

Progress in bar-code technology has made it possible for retail organizations to collect and store massive amounts of sales data, referred to as the *basket* data. A record in such data typically consists of the transaction date and the items bought in the transaction. Such an ordered list of items bought in a transaction is an important piece of the marketing infrastructure. They are interested in instituting information-driven marketing processes, managed by database technology, that enable marketers to develop and implement customized marketing programs and strategies [6].

The problem of mining association rules over basket data was introduced in [4]. An example of such a rule might be that 98% of customers that purchase

*Visiting from the Department of Computer Science, University of Wisconsin-Madison.
Permission to copy without fee for all or part of this material is granted provided that the copies are not made or distributed for commercial advantage and that copies bear this notice and notice of the title of the publication and its date, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.

Proceedings of the 20th VLDB Conference
Santiago, Chile, 1994

tires and auto accessories also get automotive services done. Finding all such rules is valuable for cross-marketing and targeted mailing applications. Other applications include catalog design, add-on sales, and market basket analysis, all based on buying patterns. The databases involved in these applications are very large. It is imperative, therefore, to have fast algorithms for this task.

The following is a formal statement of the problem [4]: Let $T = \{t_1, t_2, \dots, t_n\}$ be a set of literals, called items. Let \mathcal{D} be a set of transactions, where each transaction T is a set of items such that $T \subseteq T$. A transaction T is a set of items such that $T \subseteq T$, associated with each item $t_i \in T$ a unique identifier, called its *TID*. We say that a transaction T contains X , a set of some items in T , if $X \subseteq T$. An association rule is an implication of the form $X \Rightarrow Y$, where $X \subseteq T$, $Y \subseteq T$, and $X \cap Y = \emptyset$. The rule $X \Rightarrow Y$ holds in the transaction set \mathcal{D} with confidence c if $c\%$ of transactions in \mathcal{D} that contain X also contain Y . The rule $X \Rightarrow Y$ supports s in the transaction set \mathcal{D} if $s\%$ of transactions in \mathcal{D} contain $X \cup Y$. Our rules are somewhat more general than in [4] in that we allow a consequent to have more than one item.

Given a set of transactions \mathcal{D} , the problem of mining association rules is to generate all association rules that have support and confidence greater than the user-specified minimums (*called support and minimum confidence*) (called *Apriori*, respectively). Our discussion is now with respect to the representation of \mathcal{D} . For example, \mathcal{D} could be a data file, a relational table, or the result of a relational expression.

An algorithm for finding all association rules, henceforth referred to as the *AIS* algorithm, was presented in [4]. Another algorithm for this task, called the *SETM* algorithm, has been proposed in [13]. In this paper, we present two new algorithms, *Apriori* and *AprioriTid*, that differ fundamentally from these algorithms. We present experimental results showing

1 year, both
using the same
Apriori principle

Rakesh Agrawal, Ramakrishnan Srikant.
Fast Algorithms for Mining Association Rules in Large
Databases. VLDB 1994: 487-499

Mining Sequential Patterns

Rakesh Agrawal Ramakrishnan Srikant*

IBM Almaden Research Center
650 Harry Road, San Jose, CA 95120

Abstract

We are given a large database of customer transactions, where each transaction consists of customer, transaction time, and the items bought in the transaction. We introduce the problem of mining sequential patterns over such databases. We present three algorithms to solve this problem, namely, Apriori, AprioriAll, and AprioriSome, using synthetic data. Two of the proposed algorithms, AprioriSome and AprioriAll, have comparable performance, albeit AprioriSome performs a little better when the minimum number of customers that must support a sequential pattern is low. Synthetic experiments also show that the Apriori, AprioriAll, and AprioriSome scale linearly with the number of customer transactions. They also have excellent scale-up properties with respect to the number of transactions per customer and the number of items in a transaction.

1 Introduction

Database mining is motivated by the decision support problem faced by most large retail organizations. Progress in barcode technology has made it possible to store massive amounts of sales data, referred to as the *basket* data. A record in such data typically consists of the transaction date and the items bought in the transaction. Very often, data records also contain customer-id, purchase-order, and purchase-line information, making a record an *item-fingerprint* [6]. Customers may also collect such data using the orders they receive.

We introduce the problem of mining *sequential* patterns over this data. An example of such a pattern is

*Also Department of Computer Science, University of Wisconsin-Madison.

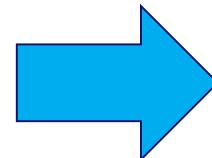
Rakesh Agrawal, Ramakrishnan
Srikant: Mining Sequential
Patterns. ICDE 1995: 3-14



Chair of Process
and Data Science

Input: Multiset of sequences of itemsets

Customer Id	Transaction Time	Items Bought
1	June 25 '93	30
1	June 30 '93	90
2	June 10 '93	10, 20
2	June 15 '93	30
2	June 20 '93	40, 60, 70
3	June 25 '93	30, 50, 70
4	June 25 '93	30
4	June 30 '93	40, 70
4	July 25 '93	90
5	June 12 '93	90



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

Multiset of sequences of itemsets

Taken from Rakesh Agrawal,
Ramakrishnan Srikant: Mining
Sequential Patterns. ICDE 1995: 3-14

Sequential Patterns with support > 25%
$\langle (30) (90) \rangle$
$\langle (30) (40 70) \rangle$

Input: Multiset of sequences of itemsets

- \mathcal{I} is the **set of all items**. An **itemset** is a nonempty set of items, e.g., $i = \{i_1, i_2, \dots, i_m\} \subseteq \mathcal{I}$ with $m \geq 1$.
- A **sequence** is a nonempty sequence of itemsets, e.g., $s = \langle s_1, s_2, \dots, s_n \rangle \in (\mathcal{P}(\mathcal{I}))^*$ with $n \geq 1$.
- A **dataset D** is a multiset of sequences.
- Technically, $D \in \mathfrak{B}\left((\mathcal{P}(\mathcal{I}))^*\right)$.
(\mathfrak{B} is the multiset and \mathcal{P} is the powerset operator)

Example input: $D \in \mathfrak{B} \left((\mathcal{P}(\mathcal{I}))^* \right)$

```
[⟨{a}, {a, b}, {b, c}, {c}⟩, ⟨{a}, {a}, {a, b}, {b, c}, {c}⟩,  
⟨{a}, {a}, {a}, {a}⟩, ⟨{a, b}, {a, b}, {a, b}, {a, b}⟩,  
⟨{a, b}, {b, c}, {c, d, e}⟩, ⟨{a}, {a, b}, {b, c}, {c}⟩]
```

Informal short-hand notation:

```
[a(ab)(bc)c, aa(ab)(bc)c, aaaa, (ab)(ab)(ab)(ab),  
(ab)(bc)(cde), a(ab)(bc)c]
```



5€ geschenkt beim Kauf eines 30€ Geschenkgutscheins

Jetzt entdecken ▾

Einkaufswagen

Preis

Menge

1 ▾

**Philips Hue White E27 LED Lampe Erweiterung, dimmbar, warmweißes Licht, steuerbar via App, kompatibel mit Amazon Alexa (Echo, Echo Dot)**

Auf Lager.

✓prime

 Dies ist ein Geschenk Erfahren Sie mehr[Löschen](#) | [Auf die Merkliste](#)

EUR 16,99

**Philips Hue Bridge, zentrales, intelligentes Steuerelement des Hue Systems**

EUR 46,95

1 ▾

Auf Lager.

✓prime

 Dies ist ein Geschenk Erfahren Sie mehr[Löschen](#) | [Auf die Merkliste](#)**Philips Hue White E27 LED Lampe Doppelpack, dimmbar, warmweißes Licht, steuerbar via App, kompatibel mit Amazon Alexa (Echo, Echo Dot)**

EUR 27,99

1 ▾

Auf Lager.

✓prime

 Dies ist ein Geschenk Erfahren Sie mehr[Löschen](#) | [Auf die Merkliste](#)**Philips Hue White und Color Ambiance E27 LED Lampe Starter Set, drei Lampen**

EUR 134,90

1 ▾

Auf Lager.

✓prime

 Dies ist ein Geschenk Erfahren Sie mehr[Löschen](#) | [Auf die Merkliste](#)
$$D \in \mathfrak{B} \left((\mathcal{P}(\mathcal{I}))^* \right)$$

Summe (4 Artikel): EUR 226,83



Kaffee

Nespresso &
You

Maschinen



Accessories



Geschenke



Our Choices



Nachhaltigkeit



Storefinder



Service | FAQ



Professional

$$D \in \mathfrak{B}\left(\left(\mathcal{P}(\mathcal{I})\right)^*\right)$$

MEIN KONTO

Willkommen Wil van der Aalst

Mitglied seit 21-09-2018

Kundennummer: 3868857



Meine Bestellungen



Meine Adressen



Meine persönlichen Daten



Meine Maschinen

Benachrichtigungen &
Erinnerungen

Marketingpräferenzen



Express Checkout



Mein Kaffee Abo

Meine Bestellungen

BESTELLDATUM	STATUS	QUELLE	Liefermethode	BESTELLNUMMER	BETRAG	
01/11/2018	Geliefert	Internet	Standardlieferung - Lieferung am nächsten Werktag	25250378	97,80 €	
21/09/2018	Geliefert	Internet	Standardlieferung innerhalb von 2 Werktagen	24533528	78,60 €	

[Mehr Bestellungen anzeigen >](#)

Bestellung - 01/11/2018

[Wieder bestellen](#)

Kapseln (250)		Stückpreis	Menge	Gesamt
	Ristretto	0,38 €	x	30
	Roma	0,38 €	x	80
	Vivalto Lungo	0,40 €	x	50
	Linizio Lungo	0,40 €	x	80
	Ristretto Decaffeinato	0,40 €	x	10

$$D \in \mathfrak{B}\left(\left(\mathcal{P}(\mathcal{I})\right)^*\right)$$



Goal: Find frequent sequential patterns

- Given a dataset $D \in \mathfrak{B}\left(\left(\mathcal{P}(\mathcal{I})\right)^*\right)$, i.e., a multiset of sequences of itemsets, find all frequent sequential patterns.
- Sequential patterns are of the same form $p \in \left(\mathcal{P}(\mathcal{I})\right)^*$, i.e., sequences of itemsets.
- The support of a sequential pattern p is the fraction (or absolute number) of sequences in D that supports pattern p (i.e., is contained).



Containment

- Let $a = \langle a_1, a_2, \dots, a_n \rangle \in (\mathcal{P}(I))^*$ and $b = \langle b_1, b_2, \dots, b_m \rangle \in (\mathcal{P}(I))^*$ be two itemset sequences.
- a is contained in b if there exist integers $i_1 < i_2 < \dots < i_n$ such that $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$.
- $a \sqsubseteq b$ if a is contained in b .
- If $a \sqsubseteq b$, then a is a **subsequence** of b and b is a **supersequence** of a .

Examples

- $\langle \{a\}, \{a, b\}, \{b, c\}, \{c\} \rangle \sqsubseteq \langle \{a\}, \{a\}, \{a, b, c\}, \{b, c\}, \{b, c\}, \{a, c\} \rangle$
- $\langle \{a\}, \{a, b\}, \{b, c\}, \{c\} \rangle \not\sqsubseteq \langle \{a\}, \{a, b, c\}, \{b, d\}, \{b, e\}, \{a, c\} \rangle$

Informal notation:

- $a(ab)(bc)c \sqsubseteq aa(abc)(bc)(bc)(ac)$
- $a(ab)(bc)c \not\sqsubseteq a(abc)(bd)(be)(ac)$

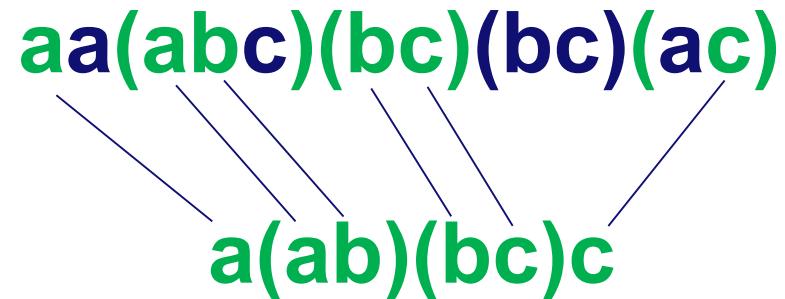
$\langle a_1, a_2, \dots, a_n \rangle \sqsubseteq \langle b_1, b_2, \dots, b_m \rangle$ iff there exist integers $i_1 < i_2 < \dots < i_n$ such that $a_1 \subseteq b_{i_1}$, $a_2 \subseteq b_{i_2}$, ..., $a_n \subseteq b_{i_n}$

Examples

- $\langle \{a\}, \{a, b\}, \{b, c\}, \{c\} \rangle \sqsubseteq \langle \{a\}, \{a\}, \{a, b, c\}, \{b, c\}, \{b, c\}, \{a, c\} \rangle$
- $\langle \{a\}, \{a, b\}, \{b, c\}, \{c\} \rangle \not\sqsubseteq \langle \{a\}, \{a, b, c\}, \{b, d\}, \{b, e\}, \{a, c\} \rangle$

Informal notation:

- $a(ab)(bc)c \sqsubseteq aa(abc)(bc)(bc)(ac)$
- $a(ab)(bc)c \not\sqsubseteq a(abc)(bd)(be)(ac)$



multiple mappings possible

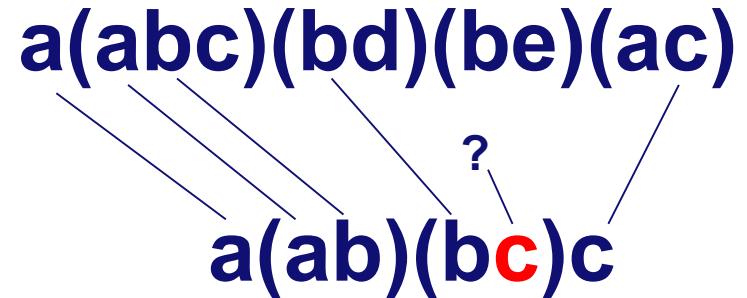
$\langle a_1, a_2, \dots, a_n \rangle \sqsubseteq \langle b_1, b_2, \dots, b_m \rangle$ iff there exist integers $i_1 < i_2 < \dots < i_n$ such that $a_1 \sqsubseteq b_{i_1}$, $a_2 \sqsubseteq b_{i_2}$, ..., $a_n \sqsubseteq b_{i_n}$

Examples

- $\langle \{a\}, \{a, b\}, \{b, c\}, \{c\} \rangle \sqsubseteq \langle \{a\}, \{a\}, \{a, b, c\}, \{b, c\}, \{b, c\}, \{a, c\} \rangle$
- $\langle \{a\}, \{a, b\}, \{b, c\}, \{c\} \rangle \not\sqsubseteq \langle \{a\}, \{a, b, c\}, \{b, d\}, \{b, e\}, \{a, c\} \rangle$

Informal notation:

- $a(ab)(bc)c \sqsubseteq aa(abc)(bc)(bc)(ac)$
- $a(ab)(bc)c \not\sqsubseteq a(abc)(bd)(be)(ac)$



$\langle a_1, a_2, \dots, a_n \rangle \sqsubseteq \langle b_1, b_2, \dots, b_m \rangle$ iff there exist integers $i_1 < i_2 < \dots < i_n$
such that $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$

Questions

- $(ab)(bc) \sqsubseteq (bc)(ab)$?
- $ab \sqsubseteq a(ac)(bc)c$?
- $aa(ab)(bc) \sqsubseteq (ab)(ace)(bce)(ab)$?
- $(abc)ef \sqsubseteq (ab)(bc)(ef)f$?
- $(abc)ef \sqsubseteq (ab)(bc)(abcd)(ef)f$?

$\langle a_1, a_2, \dots, a_n \rangle \sqsubseteq \langle b_1, b_2, \dots, b_m \rangle$ iff there exist integers $i_1 < i_2 < \dots < i_n$
such that $a_1 \sqsubseteq b_{i_1}$, $a_2 \sqsubseteq b_{i_2}$, ..., $a_n \sqsubseteq b_{i_n}$

Answers

- $(ab)(bc) \not\sqsubseteq (bc)(ab)$ (incompatible order)
- $ab \sqsubseteq a(ac)(bc)c$
- $aa(ab)(bc) \not\sqsubseteq (ab)(ace)(bce)(ab)$
((ab) cannot be mapped without also handling aa)
- $(abc)ef \not\sqsubseteq (ab)(bc)(ef)f$ (no match for (abc))
- $(abc)ef \sqsubseteq (ab)(bc)(abcd)(ef)f$

$\langle a_1, a_2, \dots, a_n \rangle \sqsubseteq \langle b_1, b_2, \dots, b_m \rangle$ iff there exist integers $i_1 < i_2 < \dots < i_n$
such that $a_1 \sqsubseteq b_{i_1}$, $a_2 \sqsubseteq b_{i_2}$, ..., $a_n \sqsubseteq b_{i_n}$

Support (relative)

- The **support** of a sequence pattern $p \in (\mathcal{P}(\mathcal{I}))^*$ given a dataset $D \in \mathfrak{B}((\mathcal{P}(\mathcal{I}))^*)$ is the fraction of sequences in D that are supersequences of p .
- $\text{support}(p) = |[s \in D | p \sqsubseteq s]| / |D|$ (relative)
- Minimum support threshold ***min_sup***: lower bound for $\text{support}(p)$.

Support (absolute)

- $\text{support_count}(p) = |[s \in D | p \sqsubseteq s]|$ (absolute, also called frequency or count)
- Minimum support count threshold: lower bound for $\text{support_count}(p)$.

Questions

- $D = [abcd, (abcd), (ab)(cd), (ab)(bc)(cd)]$
- What is the $support_count(p)$ for
 - $p=a$
 - $p=ab$
 - $p=(ab)$
 - $p=(ab)c$
 - $p=(ab)(bd)$
 - $p=ab(cd)$

Questions

- $D=[abcd,(abcd),(ab)(cd),(ab)(bc)(cd)]$
- What is the $support_count(p)$ for
 - $p=a : 4 - D=[abcd,(abcd),(ab)(cd),(ab)(bc)(cd)]$
 - $p=ab : 2 - D=[abcd,(abcd),(ab)(cd),(ab)(bc)(cd)]$
 - $p=(ab) : 3 - D=[abcd,(abcd),(ab)(cd),(ab)(bc)(cd)]$
 - $p=(ab)c : 2 - D=[abcd,(abcd),(ab)(cd),(ab)(bc)(cd)]$
 - $p=(ab)(bd) : 0 - D=[abcd,(abcd),(ab)(cd),(ab)(bc)(cd)]$
 - $p=ab(cd) : 1 - D=[abcd,(abcd),(ab)(cd),(ab)(bc)(cd)]$

AprioriAll Algorithm



Brute force approach

- Let k be the length of the longest sequence in $D \in \mathfrak{B}\left(\left(\mathcal{P}(\mathcal{I})\right)^*\right)$ and q the largest itemset.
 - Generate all sequence patterns $p \in \left(\mathcal{P}(\mathcal{I})\right)^*$ having a length $\leq k$ and itemsets of size $\leq q$. This number is finite.
 - Compute the support of each candidate pattern.
 - Return all that have a support higher than min_sup .
- Obviously this is very expensive.

A smarter approach based on Apriori

- First described in Rakesh Agrawal, Ramakrishnan Srikant: Mining Sequential Patterns. ICDE 1995: 3-14.
- Just like for frequent itemsets: Avoid testing hopeless candidates.
- If $p_1 \sqsubseteq p_2$ (p_1 is a subsequence of p_2), then p_2 cannot be frequent if p_1 is not frequent.
 - $\text{support}(p_1) \geq \text{support}(p_2)$ if $p_1 \sqsubseteq p_2$
 - if $p_1 \sqsubseteq p_2$ and $\text{support}(p_1) < \text{min_sup}$, then $\text{support}(p_2) < \text{min_sup}$

Step 1: Determine all itemsets L

- $L = \{i \subseteq I | support(\langle i \rangle) \geq min_sup\}$ are all itemsets that appear in a sufficient number of sequences.
- These itemsets are called **itemsets**.
- For example, for $D=[abcd,(abcd),(ab)(cd),(ab)(bc)(cd)]$ and $min_sup = 0.7$ the following itemsets are frequent: **a** (support = 4/4), **b** (support = 4/4), **C** (support = 4/4), **d** (support = 4/4), **(ab)** (support = 3/4), **(cd)** (support = 3/4).
- To determine all **itemsets**, we can use a variant of the **original Apriori algorithm** (the only difference is that support is counted per sequence and not per itemset).

Step 2: Preprocess dataset

- We only need to consider itemsets (i.e., L): There cannot be any patterns that involve other itemsets.
- Frequent sequence patterns must be of the form L^* !
- The set $L_1 = \{\langle i \rangle | i \in L\}$ is the set of all frequent sequence patterns of length 1.
- $L_k \subseteq L^*$ is the set of all frequent sequence patterns of length k (to be computed).

Transformation

(only for optimization, using formal notation)

- Transform $D \in \mathfrak{B}\left(\left(\mathcal{P}(\mathcal{I})\right)^*\right)$ into $D_T \in \mathfrak{B}\left(\left(\mathcal{P}(L)\right)^*\right)$.
- Each sequence is now described by a sequence of sets of itemsets.
- Let $L = \{\{a\}, \{b\}, \{c\}, \{a, b\}\}$ and $D = [\langle\{a, c\}, \{a, b, c\}\rangle, \langle\{c\}, \{a, c\}\rangle, \dots]$.
- $D_T = [\langle\{\{a\}\{c\}\}, \{\{a\}, \{b\}, \{c\}, \{a, b\}\}\rangle, \langle\{\{c\}\}, \{\{a\}, \{c\}\}\rangle, \dots]$
- The new representation makes it very easy to test whether a sequence pattern is supported by a sequence in the data set.

Step 3: Generate set of candidate sequences

- Assume we have L_{k-1} , the set of all frequent sequence patterns of length $k - 1$. Recall that $L_1 = \{\langle i \rangle | i \in L\}$.
- C_k is the set of all candidate sequences obtained by taking two sequences from L_{k-1} where the first $k - 1$ are the same.

Step 4: Prune the set of candidate sequences

- For all $c \in C_k$.
 - Consider all subsequences of c of length $k - 1$.
 - If one of these subsequences is not in L_{k-1} , then remove c from C_k .

Step 5: Test all candidate sequences

- For each transformed sequence $s \in D_T$: Increment the count of $c \in C_k$ if c is contained in s .
- Remove all candidates $c \in C_k$ that do no meet the threshold and the result is L_k .
- $L_k = \{c \in C_k | support(c) \geq min_sup\}$.
- Increment k and goto Step 3 until $L_k = \emptyset$.
- $\bigcup_k L_k$ is the set of all frequent sequence patterns

Step 6 (optional): Remove non-maximal patterns

- An sequence s is a maximal sequence in data set D if s is frequent, and there is no other supersequence s' that is also frequent ($s \sqsubset s'$).
- It is possible to keep only the maximal sequences. However, support information for the subsequences will be lost (these may be more frequent).

Other sequential pattern mining approaches

From Carl Mooney, John Roddick, Sequential pattern mining: approaches and algorithms, ACM Comput. Surv., 45(2):1-39, 2013.

Table VII. A summary of Apriori-based algorithms.

Algorithm Name	Author	Year	Notes
Candidate Generation: Horizontal Database Format			
Apriori (All, Some, Dynamic Some)	[Agrawal and Srikant]	1995	
Generalised Sequential Patterns (GSP)	[Srikant and Agrawal]	1996	Max/Min Gap, Window, Taxonomies
PSP	[Masseglia et al.]	1998	Retrieval optimisations
Sequential Pattern mIning with Regular expressIon consTraints (SPIRIT)	[Garofalakis et al.]	1999	Regular Expressions
Maximal Frequent Sequences (MFS)	[Zhang et al.]	2001	Based on GSP, uses Sampling
Regular Expression-Highly Adaptive Constrained Local Extractor (RE-Hacke)	[Albert-Lorincz and Boulicaut]	2003	Regular Expressions, similar to SPIRIT
Maximal Sequential Patterns using Sampling (MSPS)	[Luo and Chung]	2004	Sampling
Candidate Generation: Vertical Database Format			
Sequential PAttern Discovery using Equivalence classes (SPADE)	[Zaki]	2001	Equivalence Classes
Sequential PAttern Mining (SPAM)	[Ayres et al.]	2002	Bitmap representation
LAst Position INduction (LAPIN)	[Yang and Kitsuregawa]	2004	Uses last position
Cache-based Constrained Sequence Miner (CCSM)	[Orlando et al.]	2004	k-way intersections, cache
Indexed Bit Map (IBM)	[Savary and Zeitouni]	2005	Bit Map, Sequence Vector, Index, NB table
LAst Position INduction Sequential PAttern Mining (LAPIN-SPAM)	[Yang and Kitsuregawa]	2005	Uses SPAM, Uses last position



Table IX. A summary of pattern growth algorithms.

Algorithm Name	Author	Year	Notes
Pattern Growth			
FREQuEnt pattern-projected Sequential PAtterN mining (FreeSpan)	[Han et al.]	2000	Projected sequence database
PREFIX-projected Sequential PAtterN mining (PrefixSpan)	[Pei et al.]	2001	Projected prefix database
Sequential pattern mining with Length-decreasing suPPort (SLPMiner)	[Seno and Karypis]	2002	Length-decreasing support

Table XI. A summary of temporal sequence algorithms.

Algorithm Name	Author	Year	Notes
Candidate Generation: Episodes			
WINEPI	[Mannila et al.]	[1995], 1996	State automata, Window
WINEPI, MINEPI	[Mannila et al.]	1997	State automata, Window, Maximal
Pattern Directed Mining	[Guralnik et al.]	1998	Pattern language
Event-Oriented Patterns	[Sun et al.]	2003	Target events
Pattern Growth: Episodes			
PROjected Window Lists (PROWL)	[Huang et al.]	2004	Projected window lists



Extensions



Association rules based on frequent sequences

- Frequent sequence patterns can be split in an “if” and “then” part (similar to “normal” association rules).
- $\langle \{beer\}, \{red, white\} \rangle \Rightarrow \langle \{beer\}, \{red, white\}, \{wodka\} \rangle$
- $\langle \{beer\}, \{beer\} \rangle \Rightarrow \langle \{beer\}, \{beer\}, \{beer\}, \{beer\}, \{beer\} \rangle$
- Many variants possible.

How interesting are frequent sequences?

- For any technique that identifies patterns, it is important to filter out the less interesting ones.
- One can look at things like correlations and base frequencies to decide how surprising sequences or sequence-based rules are.
- It is also possible to add further constraints.

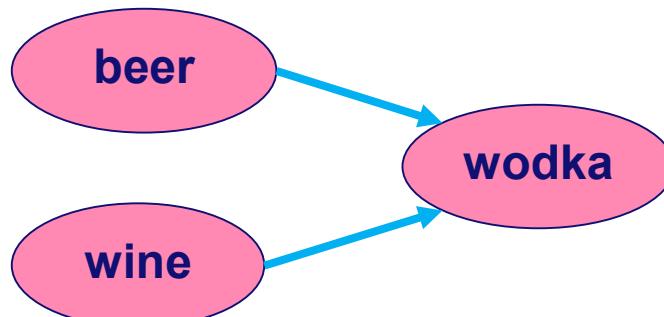
Additional constraints

- **Item constraints** (only consider sequences that include or exclude a set of items).
- **Length constraints** (only consider patterns of a given size).
- **Time constraints** (only consider patterns that occur in a short timeframe). This includes gap and duration constraints.
- **Regular expression constraints** (only consider patterns that satisfy a regular expression or temporal constraint).



Chair of Process
and Data Science

Episode mining



Rather than looking for sequences we look for embedded **partial orders**.

Reminder: Data science is like a making cocktail. Mix the proper ingredients in the right way!



Python
Visualization
Decision trees
Regression
Support vector machines
Neural networks
Evaluation
Clustering
Frequent items sets
Association rules
Sequence mining
Process mining
Process discovery
Conformance checking
Text mining
Preprocessing
Visual analytics
Encryption
Anonymization
Big data infra
Distribution



Python
Visualization
Decision trees
Regression
Support vector machines
Neural networks
Evaluation
Clustering
Frequent items sets
Association rules
Sequence mining
Process mining
Process discovery
Conformance checking
Text mining
Preprocessing
Visual analytics
Encryption
Anonymization
Big data infra
Distribution



Visualization



Support vector machines



Clustering



Sequence mining



Preprocessing

Conclusion

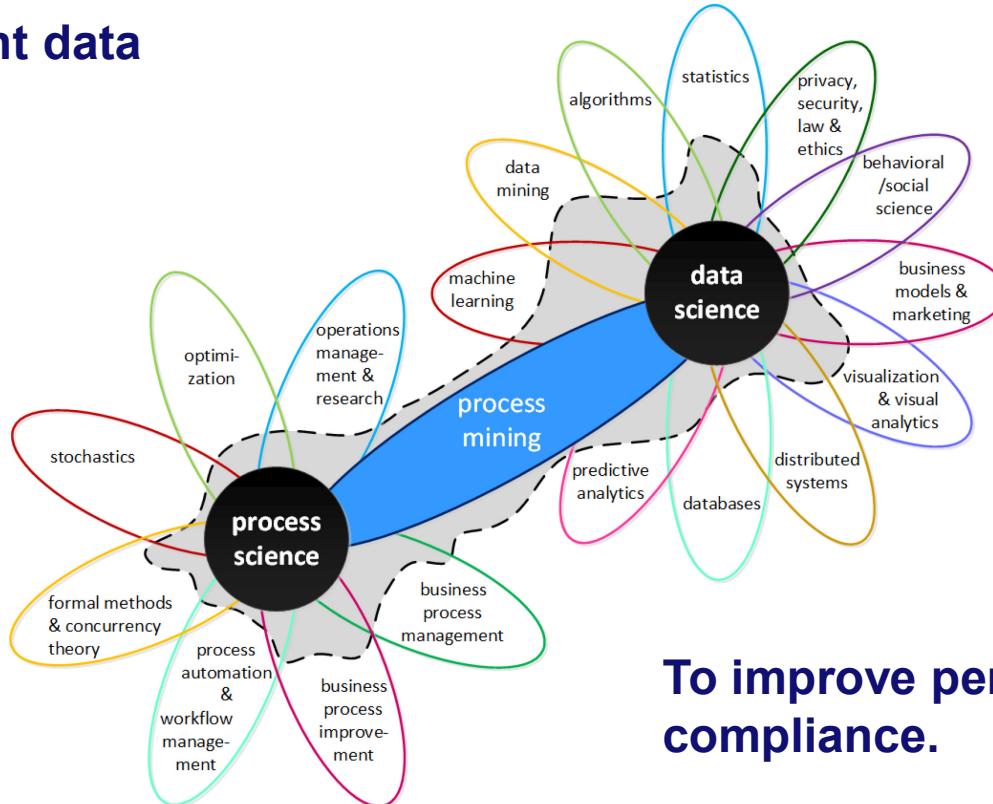


Short summary of lecture

- Event data are everywhere!
- Many questions related to event data!
- Sequence mining looks for local patterns.
- Apriori-style approach can be used (builds on frequent itemsets).

Next process mining

Starting point: Event data



Chair of Process
and Data Science

Sequence mining versus process mining

Sequence mining	Process mining
Focus on local patterns	Focus on end-to-end processes
Focus on itemsets	Focus on activities
Low-level patterns	Higher-level behaviors (e.g., concurrency and choice).
Specific event data (itemset based)	General event data (case, activity, resource, lifecycle, etc.)
...	...

#	Lecture	Date	Day
	Lecture 1	Introduction	10/10/2018 Wednesday
	Lecture 2	Crash Course in Python	
Instruction 1		Python	
	Lecture 3	Basic data visualisation	
	Lecture 4	Decision trees	
Instruction 2		Decision trees and data mining	
	Lecture 5	Regression	
	Lecture 6	Support vector machines	
Instruction 3		Regression and support vector machines	
	Lecture 7	Neural networks (1/2)	
Instruction 4		Neural networks and neural network architectures	
	Lecture 8	Neural networks (2/2)	
	Lecture 9	Evaluation of supervised learning models	
Instruction 5		Neural networks and evaluation of supervised learning models	
	Lecture 10	Clustering	
	Lecture 11	Frequent items sets	15/11/2018 Thursday
	Lecture 12	Association rules	21/11/2018 Wednesday
	Lecture 13	Sequence mining	22/11/2018 Thursday
Instruction 6		Clustering, frequent items sets, association rules	23/11/2018 Friday
	Lecture 14	Process mining (unsupervised)	28/11/2018 Wednesday
	Lecture 15	Process mining (supervised)	29/11/2018 Thursday
Instruction 7		Process mining and sequence mining	30/11/2018 Friday
	Lecture 16	Text mining (1/2)	05/12/2018 Wednesday
Instruction 8		Text mining and process mining	06/12/2018 Thursday !!
	Lecture 17	Text mining (2/2)	12/12/2018 Wednesday
	Lecture 18	Data preprocessing, data quality, binning, etc	13/12/2018 Thursday
	Lecture 19	Visual analytics & information visualization	19/12/2018 Wednesday
	backup		20/12/2018 Thursday
Instruction 9		Text mining, preprocessing and visualization	21/12/2018 Friday
	Lecture 20	Responsible data science (1/2)	09/01/2019 Wednesday
	Lecture 21	Responsible data science (2/2)	10/01/2019 Thursday
Instruction 10		Responsible data science	11/01/2019 Friday
	Lecture 22	Big data (1/2)	16/01/2019 Wednesday
	Lecture 23	Big data (2/2)	17/01/2019 Thursday
Instruction 11		Big data	18/01/2019 Friday
	Lecture 24	Closing	23/01/2019 Wednesday
	backup		24/01/2019 Thursday
Instruction 12		Example exam questions	25/01/2019 Friday
	backup		30/01/2019 Wednesday
	backup		31/01/2019 Thursday
extra		Question hour	01/02/2019 Friday