Introduction to Pattern Discovery (1/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Analysis of Patterns

Introduction to Pattern Discovery

Florent Nicart

University of Rouen

September 29th, 2009

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Outline

Introduction to Pattern Discovery (2/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching v Pattern Discovery

Association Patterns Frequent sets Association Rule: Mining

Sequential Patterns Sequences Languages

Introduction

- Why Pattern Discovery?
- Data sets
- Pattern Matching vs Pattern Discovery

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

Association Patterns

- Frequent Set Mining
- Association Rules Mining



2

Sequential Patterns

- Sequences
- Languages

Why Pattern Discovery?

Introduction to Pattern Discovery (3/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching v Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages





- Increase of data storage capacity,
- increase of computational power,
- inter-disciplinary techniques,
- objective : to understand the underlying process generating the dataset.

Availability of huge datasets

Introduction to Pattern Discovery (4/59)

Cagliari'09

Introduction

Why Pattern Discovery?

Data sets

Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules

Sequential Patterns Sequences Languages

Databases are everywhere

- Supermarket transactions,
- credit card records,
- telephone call details,
- weblogs, ISP logs,
- genetic databases.

What is a pattern?

Pattern matching : a collection of similar values

Pattern discovery : a particuliar combination of values

Pattern Matching vs Pattern Discovery



Pattern Matching vs Pattern Discovery

Introduction to Pattern Discovery (6/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules

Sequential Patterns Sequences Languages



(日)

э

Pattern Discovery

- a dataset,
- a metric and a distance
- "a pattern space"



◆□▶ ◆□▶ ◆三▶ ◆三▶ ・三 の々で

Introduction to Pattern Discovery (8/59)

Cagliari'09

ntroduction

Why Pattern Discovery? Data sets Pattern Matching vs Pattern Discovery

Association Patterns

Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

The pattern :

Finding (subsets) events that occur together [Agrawal93].

Rakesh Agrawal and Ramakrishnan Srikant.

Fast Algorithms for Mining Association Rules in Large Databases,

Proceedings of 20th International Conference on Very Large Data Bases, 1994.

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

Introduction to Pattern Discovery (9/59)

Cagliari'09

ntroduction

Why Pattern Discovery? Data sets Pattern Matching v Pattern Discovery

Association Patterns

Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Typical application (origin) :

Market basket analysis :

- "10% of customers are buying wine and cheese",
- "15% are buying crisps and beer",

and the legendary^a :

 "People buying beer on saturday are very likely to buy nappies" (Wal-Mart).

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

^aBeer and Nappies – A Data Mining Urban Legend http://web.onetel.net.uk/~hibou/Beer and Nappies.html



Introduction to Pattern Discovery (11/59)

Cagliari'09

Introduction

Why Pattern Discovery ? Data sets Pattern Matching v Pattern Discovery

Association Patterns

Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Other typical applications :

- Financial-services/telecommunications companies : to which combination of services customers most often subscribe,
- "Quality-assurance : which combinations of components are most likely to fail at the same time,

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

- Web log analysis,
- Finding Association Rules,
- Classification, ...

Frequent Set Mining Defining the problem

Introduction to Pattern Discovery (12/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets

Association Rules Mining

Sequential Patterns Sequences Languages

The input data

- $I = \{i_1, i_2, ..., i_n\}$, set of items,
- itemset S : any subset^a of \mathcal{I} ,
- transaction : $t = \langle tid, S_{tid} \rangle$, $tid \in \mathbb{N}$ and $S_{tid} \subseteq \mathcal{I}$
- $T = \{t_1, t_2, \dots, t_m\}$, set of transactions,

^aUsually excluding the empty set

Common itemsets

$$t \cap t' = S \cap S'$$
 with $t = \langle tid, S \rangle$ and $t' = \langle tid', S' \rangle$
common $(t, t') = 2^{t \cap t'} \setminus \emptyset$

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

Frequent Set Mining

Example (Transactions and itemsets)

Introduction to Pattern Discovery (13/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets

Association Hules Mining

Sequential Patterns Sequences Languages

	_	_	_	_	_	· · ·
basket	P_1	P_2	P_3	P_4	P_5	I _{tid}
<i>t</i> ₁	1	0	0	0	0	{ <i>P</i> ₁ }
t ₂	1	1	1	1	0	$\{P_1, P_2, P_3, P_4\}$
t ₃	1	0	1	0	1	$\{P_1, P_3, P_5\}$
t_4	0	0	1	0	0	{ P ₃ }
t5	0	1	1	1	0	$\{P_2, P_3, P_4\}$
t ₆	1	1	1	0	0	$\{P_1, P_2, P_3\}$
t ₇	1	0	1	1	0	$\{P_1, P_3, P_4\}$
t ₈	0	1	1	0	1	$\{P_2, P_3, P_5\}$
t ₉	1	0	0	1	0	$\{P_1, P_4\}$
t ₁₀	0	1	1	0	1	$\{P_2, P_3, P_5\}$
$common(t_2, t_3) = \{\{P_1\}, \{P_3\}, \{P_1, P_3\}, \}$						

Frequent Set Mining Defining the problem

Introduction to Pattern Discovery (14/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets

Association Rules Mining

Sequential Patterns Sequences Languages

Frequent Itemset Mining (FIM)

• Find the most frequent itemsets that occur amongst all the transactions of *T*.

• The pattern space is $2^{\mathcal{I}}$

Complexity

• Size of the space of patterns : $2^{|\mathcal{I}|}$

example	#items sold	2 ^I
Small shops	5000	1,4 * 10 ¹⁵⁰⁵
ebay.com	30.147.410	5 * 10 ⁹⁰⁷⁵²⁷⁴
amazon(+third)	+5.000.000	$9,5 * 10^{1505149}$
ebay.com(cat)	20.000	4 * 10 ⁶⁰²⁰

Frequent Set Mining Definitions

Introduction to Pattern Discovery (15/59)

Cagliari'09

Introduction

Discovery ? Data sets Pattern Matching

Association Patterns

Frequent sets

Association Rules Mining

Sequential Patterns Sequences Languages

Definition (Support Count of an itemset)

The *Support count* (or simply count) of an itemset is the number of transactions where it occurs. $Count(S) = |\{ < tid, S_{tid} > | S \subseteq S_{tid} \}|$

Definition (Support of an itemset)

The *Support* of an itemset is the proportion of transactions where it occurs. Support(S) = Count(S)/n

Definition (Minimal support)

minsup is the minimal support of itemset to be considered.

Definition (Supported itemsets)

S is a supported itemset iff Support(S) >= minsup.

Frequent Set Mining Definitions

Introduction to Pattern Discovery (16/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns

Frequent sets

Association Rules Mining

Sequential Patterns Sequences Languages

Example (Supported itemsets)

bas	ket	<i>P</i> ₁	P ₂	P ₃	P ₄	P_5	I _{tid}
t ₁		1	0	0	0	0	{ <i>P</i> ₁ }
t ₂		1	1	1	1	0	$\{P_1, P_2, P_3, P_4\}$
t ₃		1	0	1	0	1	$\{P_1, P_3, P_5\}$
t ₄		0	0	1	0	0	$\{P_3\}$
t ₅		0	1	1	1	0	$\{P_2, P_3, P_4\}$
t ₆		1	1	1	0	0	$\{P_1, P_2, P_3\}$
t7		1	0	1	1	0	$\{P_1, P_3, P_4\}$
t ₈		0	1	1	0	1	$\{P_2, P_3, P_5\}$
t ₉		1	0	0	1	0	$\{P_1, P_4\}$
t ₁₀		0	1	1	0	1	$\{P_2, P_3, P_5\}$

$$\begin{split} &Support(\{P_1,P_3\}) = |\{t_2,t_3,t_6,t_7\}|/10 = 0.4\\ &Support(\{P_2,P_4\}) = |\{t_2,t_5\}|/10 = 0.2\\ &\text{with $minsup = 0.25$, $\{P_1,P_3\}$ is supported, $\{P_2,P_4\}$ is not.} \end{split}$$

Frequent Set Mining Definitions

Introduction to Pattern Discovery (17/59)

Cagliari'09

Introduction Why Pattern Discovery ? Data sets Pattern Matching v Pattern Discovery

Association Patterns Frequent sets

Association Rules Mining

Sequential Patterns Sequences

Theorem (Downward closure property/Support monotonicity.)

$$S \subseteq S' \equiv Support(S) >= Support(S')$$

By corollary, if L_k is the set of supported sets with cardinality k:

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Theorem

$$L_{k} = \emptyset \Rightarrow \forall i > k, L_{i} = \emptyset$$

Frequent Set Mining The Apriori algorithm

Introduction to Pattern Discovery (18/59)

Cagliari'09

The Apriori algorithm

while $(L_{k-1} = \emptyset)$ {

```
L_1 \leftarrow supported itemsets of cardinality one k \leftarrow 2
```

▲□▶▲□▶▲□▶▲□▶ □ のQ@

```
Data sets
Pattern Matching vs
Pattern Discovery
```

Association Patterns Frequent sets

```
Association Rules
Mining
```

Sequential Patterns Sequences Languages $egin{aligned} & C_k \leftarrow \textit{create_candidates}(L_{k-1}) \ & L_k \leftarrow \textit{prune}(C_k) \ & k \leftarrow k+1 \end{aligned}$

return $L_1 \cup L_2 \cup \ldots L_k$

Frequent Set Mining The Apriori algorithm

Introduction to Pattern Discovery (19/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets

Association Rules Mining

Sequential Patterns Sequences Languages

Generating C_k from L_{k-1}

Join Step

Compare each member of L_{k-1} , say A, with every other member, say *B*, in turn. If the first k - 2 items in A and B (i.e. all but the rightmost elements of the two itemsets) are identical, place set $A \cup B$ into C_k .

Prune Step

 L_{k-1}

For each member c of C_k in turn {

Examine all subsets of c with k - 1 elements

Delete c from C_k if any of the subsets is not a member of

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

Frequent Set Mining The Apriori algorithm : an example

Introduction to Pattern Discovery (20/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns

Frequent sets

Association Rule: Mining

Sequential Patterns Sequences

Example (Joint step for C_5)

 $\begin{array}{l} L_4 = \{\{p,q,r,s\},\{p,q,r,t\},\{p,q,r,z\},\{p,q,s,z\},\\ \{p,r,s,z\},\{q,r,s,z\},\{r,s,w,x\},\{r,s,w,z\},\{r,t,v,x\},\\ \{r,t,v,z\},\{r,t,x,z\},\{r,v,x,y\},\{r,v,x,z\},\{r,v,y,z\},\\ \{r,x,y,z\},\{t,v,x,z\},\{v,x,y,z\}\} \end{array}$

<i>I</i> ₁	<i>I</i> ₂	Contribution to C'_5
{ <i>p</i> , <i>q</i> , <i>r</i> , <i>s</i> }	$\{p, q, r, t\}$	$\{p, q, r, s, t\}$
{ p , q , r , s }	$\{p, q, r, z\}$	$\{p, q, r, s, z\}$
$\{p, q, r, t\}$	$\{p, q, r, z\}$	$\{p, q, r, t, z\}$
$\{r, s, w, x\}$	$\{r, s, w, z\}$	$\{r, s, w, x, z\}$
${r, t, v, x}$	$\{r, t, v, z\}$	$\{r, t, v, x, z\}$
$\{r, v, x, y\}$	$\{r, v, x, z\}$	$\{r, v, x, y, z\}$

 $\rightarrow C'_{5} = \{\{p,q,r,s,t\},\{p,q,r,s,z\},\{p,q,r,t,z\},\{r,s,w,x,z\},\{r,t,v,x,z\},\{r,v,x,y,z\}\}$

Frequent Set Mining The Apriori algorithm : an example

Introduction to Pattern Discovery (21/59)

Cagliari'09

Introduction Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns

Frequent sets

Association Rules Mining

Sequential Patterns Sequences Languages

Example (Pruning step for C_5)

 $\begin{array}{l} L_4 = \{\{p,q,r,s\},\{p,q,r,t\},\{p,q,r,z\},\{p,q,s,z\},\\ \{p,r,s,z\},\{q,r,s,z\},\{r,s,w,x\},\{r,s,w,z\},\{r,t,v,x\},\\ \{r,t,v,z\},\{r,t,x,z\},\{r,v,x,y\},\{r,v,x,z\},\{r,v,y,z\},\\ \{r,x,y,z\},\{t,v,x,z\},\{v,x,y,z\}\}\\ C_5' = \{\{p,q,r,s,t\},\{p,q,r,s,z\},\{p,q,r,t,z\},\{r,s,w,x,z\},\\ \{r,t,v,x,z\},\{r,v,x,y,z\}\} \end{array}$

Itemset in C'_5	Subsets all in L ₄ ?
{ <i>p</i> , <i>q</i> , <i>r</i> , <i>s</i> , <i>t</i> }	No : { <i>p</i> , <i>q</i> , <i>s</i> , <i>t</i> } ∉ <i>L</i> ₄
$\{p, q, r, s, z\}$	Yes
$\{p, q, r, t, z\}$	No : $\{p, q, t, z\} \notin L_4$
$\{r, s, w, x, z\}$	No : $\{r, s, x, z\} \notin L_4$
$\{r, t, v, x, z\}$	Yes
$\{r, v, x, y, z\}$	Yes

 $C_5 = \{\{p, q, r, s, z\}, \{r, t, v, x, z\}, \{r, v, x, y, z\}\}$

Association Rules Mining

Introduction to Pattern Discovery (22/59)

Cagliari'09

Introduction

Why Pattern Discovery ? Data sets Pattern Matching Pattern Discovern

Association Patterns

Association Rules Mining

Sequential Patterns Sequences Languages • We would like to use the itemset found to express association rules...

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Association Rules Mining

Introduction to Pattern Discovery (23/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets

Association Rules Mining

Sequential Patterns Sequences Languages

Definition (Rule)

 $L \rightarrow R$ with L (antecedant/condition) and R (consequence) being predicates.

- Reads L implies R,
- ex : "If it rains then the ground will be wet."
- $\{ab\} \rightarrow \{cd\}$ means $\{a, b\} \subseteq S \rightarrow \{c, d\} \subseteq S, \forall S \in T$

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

Definition (Association Rule)

 $L \rightarrow R$ with $L \neq \emptyset$ and $R \neq \emptyset$ being disjoint itemsets.

Association Rules Mining

Introduction to Pattern Discovery (24/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching v Pattern Discovery

Association Patterns Frequent sets

Association Rules Mining

Sequential Patterns Sequences Languages

- Probabilistic AR : "When beer and crisps are bought together, cheese is bought in 45% of cases",
- $L \rightarrow R : L \cup R$ is an itemset :

Definition (Support of an Association Rule)

 $Support(L \rightarrow R) = Support(L \cup R) = Count(L \cup R)/n$

Definition (Confidence of an Association Rule)

 $Confidence(L \rightarrow R) = Support(L \cup R)/Support(L) = Count(L \cup R)/Count(L)$

Association Rules Mining Generating rules

Introduction to Pattern Discovery (25/59)

Cagliari'09

Introduction

Discovery ? Data sets Pattern Matchin

Pattern Discovery

Association Patterns Frequent sets

Association Rules Mining

Sequential Patterns Sequences Languages

- Only supported rules : $Support(L \rightarrow R) \ge minsup$,
- given | L ∪ R |= k, then number of rules that can be generated :

$$\sum_{i=1}^{k-1} C_i^k = 2^k - 2^k$$

• Filtering rules : Confidence $(L \rightarrow R) \ge minconf$

Definition (Confident Right-hand side)

The right-hand side of a rule is said *confident* iff *Confidence*($L \rightarrow R$) \geq *minconf* and *unconfident* otherwise.

Association Rules Mining Generating rules

Introduction to Pattern Discovery (26/59)

Cagliari'09

Introduction Why Pattern Discovery ? Data sets Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets

Association Rules Mining

Sequential Patterns Sequences Languages

theorem

 $Confidence(L \cup \{x\} \rightarrow R) \geq Confidence(L \rightarrow R \cup \{x\})$

Corrolary

Any superset of an unconfident right-hand itemset is unconfident.

Any (non-empty) subset of a confident right-hand itemset is confident.

• Rules can be generated with an Apriori-like algorithm.

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

Sequential Patterns

Introduction to Pattern Discovery (27/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matchir

Pattern Discovery

Patterns Frequent sets

Association Rules

Sequential Patterns

Sequences Languages Temporal sequence data sets :

• Supermaket customer transactions,

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

- Weblogs : user navigation,
- Alarm logs : fault analysis,
- Text sequence data sets :
 - Genomics,
 - documents,
 - etc.

Sequences

Introduction to Pattern Discovery (28/59)

Cagliari'09

Introduction Why Pattern Discovery? Data sets

Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Finding sequential patterns in sequences of events :

• Data set : supermaket customer transactions,

Rakesh Agrawal and Ramakrishnan Srikant.

Extending Database Technology, 1995.

Proceedings of 5th International Conference on

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

Mining Sequential Patterns,

 Pattern : customers renting "The Fellowship of the Ring" rent "The Two Towers" and then "The Return of the King"

Sequences Definitions

Introduction to Pattern Discovery (29/59)

Definition (Sequence)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules

Sequential Patterns Sequences Languages

$s = \langle s_1 \dots s_n \rangle$ where s_k is an itemset

Definition (Contained sequence)

A sequence $\langle a_1 \dots a_n \rangle$ is contained in a sequence $\langle b_1 \dots b_m \rangle$ if there exists integers $i_1 < i_2 < \dots < i_n$ such that $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots a_n \subseteq b_{i_n}$.

Examples

- $\langle \{3\}\{4,5\}\{8\} \rangle$ is contained in $\langle \{7\}\{3,8\}\{9\}\{4,5,6\}\{8\} \rangle$
- $\langle \{3\}\{5\}\rangle$ is not contained in $\langle \{3,5\}\rangle$

Sequences Examples

Introduction to Pattern Discovery (30/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching

Association

Frequent sets

Association Rules Mining

Sequential Patterns Sequences

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	<pre>{{30, 50, 70}{30}{40, 70}{90}}</pre>
5	$\langle \{90\} \rangle$

Sequence Support

Introduction to Pattern Discovery (31/59)

Cagliari'09

Introduction Why Pattern Discovery? Data sets Pattern Matching vs

Association Patterns Frequent sets Association Rules

Sequential Patterns Sequences Languages

Definition (Sequence Support)

Let T_i be the set of transactions of customer *i*, $Support(s) = |\{i \mid s \subseteq T_i\}|$

Definition (Maximal sequence)

In a set of sequences, a sequence is *maximal* if it is not contained in any other sequence.

Pattern discovery task :

Find all maximal sequences (*Sequential Patterns*) in a set of customer transactions having a support greater than *minseqsup*.

Sequence Support



Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets Association Rule

Sequential Patterns Sequences

Customer Id	Customer Sequence
1	$\langle \{ 30 \} \{ 90 \} angle$
2	$\langle \{10, 20\} \{30\} \{40, 60, 70\} angle$
3	$\langle \{30, 50, 70\} \{30\} \{40, 70\} \{90\} \rangle$
5	$\langle \{90\} \rangle$

Sequential Patterns (<i>minseqsup</i> = 0.25)
$\langle \{ 30 \} \{ 90 \} angle$
$\langle \{ 30 \} \{ 40,70 \} angle$

Sequences $\langle \{30\} \rangle$, $\langle \{40\} \rangle$, $\langle \{70\} \rangle$, $\langle \{90\} \rangle$, $\langle \{30\} \{40\} \rangle$, $\langle \{30\} \{70\} \rangle$, $\langle \{40, 70\} \rangle$ are not maximal.

Sequences Algorithm

Introduction to Pattern Discovery (33/59)

Cagliari'09

ntroduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets Association Rules

Sequential Patterns Sequences Languages

Algorithm Agrawal et al

- Sorting phase : (customer id, transaction time)
- Customer-supported itemset phase : customer count based Apriori
- Transformation phase : transforms the customer transaction to accelerate the containing tests
- Sequence phase : find supported sequences from (customer)supported itemsets (AprioriAll, AprioriSome/DynamicSome)
- Maximal phase : keep only the maximal supported sequences.

Sequences Episode Mining

Introduction to Pattern Discovery (34/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets Association Rule Mining

Sequential Patterns Sequences Languages Finding episodes in sequences of events :

- Data set : set of timed events,
- Pattern : groups of events occuring frequently *close* to each other

Heikki Mannila, Hannu Toivonen and A. Inkeri Verkamo, Discovery of Frequent Episodes in Event Sequences, Data Mining and Knowledge Discovery, 1997.

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

Sequences Event sequences

Introduction to Pattern Discovery (35/59)

Cagliari'09

ntroduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets Association Rule Mining

Sequential Patterns Sequences Languages

Definition (event)

Let *E* be a set of event types, an event is a pair (*A*, *t*) where $A \in E$ and *t* is an integer (timestamp).



 $s = \langle (E, 31), (D, 32), (F, 33), (A, 35), (B, 37), \dots, (D, 67) \rangle$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

Sequences Episodes

Introduction to Pattern Discovery (36/59)

Cagliari'09

Introduction Why Pattern

Discovery? Data sets

Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences



 $s = \langle (E, 31), (D, 32), (F, 33), (A, 35), (B, 37), \dots, (D, 67) \rangle$

Definition (episode)

An episode is a partially ordered subset of events.



Subepisodes

Introduction to Pattern Discovery (37/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Definition (Subepisode)

An episode $e = (V', \leq', g')$ is a *subepisode* of $e = (V, \leq, g)$ if $V' \subset V$, $\forall v \in V', g'(v) = g(v)$, and $\forall v, w \in V', v \leq' w \rightarrow v \leq w$



▲□▶▲□▶▲□▶▲□▶ □ のQ@

Sequences Mining episodes



Note : frequency of an episode = amount of windows in which it occurs.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ● ● ●

Sequences Algorithm

Introduction to Pattern Discovery (39/59)

Cagliari'09

Introduction Why Pattern Discovery? Data sets Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Algorithm (Mannila et al)

```
C_{1} \leftarrow \{\{e\} \mid e \in E\}
k \leftarrow 1
while (C_{k} = \emptyset) \{
L_{k} \leftarrow frequent(C_{k}, s, minfr, winsize)
C_{k+1} \leftarrow generate(L_{k}, C)
k \leftarrow k + 1
\}
return L_{1} \cup L_{2} \cup \ldots L_{k}
```

▲□▶▲□▶▲□▶▲□▶ □ のQ@

l an		nee
Lan	yuu	UQUU.
	U	U

Introduction	to
Pattern	
Discovery	
(40/59)	

Cagliari'09

Introduction

Why Pattern Discovery ? Data sets Pattern Matching Pattern Discover

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages Discovering (rational) languages



Languages What is Gramatical Inference?

Introduction to Pattern Discovery (41/59)

Cagliari'09

- Introductior
- Why Pattern Discovery?
- Data sets Pattern Matching
- Pattern Discovery
- Association Patterns Frequent sets Association Rules

Sequential Patterns Sequences Languages • A grammar describes a language/set of words (possibly infinite).

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

- Is it possible, from a finite sample (or a positive and negative sample) to guess the grammar?
- Grammatical inference : find a description of the generating process.

Languages Grammars and their hierarchy

Introduction to Pattern Discovery (42/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets Association Rule

Sequential Patterns Sequences Languages Chomsky hierarchy :



Туре	Grammar	Language	Machine
type-0	Unrestricted	Recurs. enumer.	Turing machine
type-1	Context-sensitive	Context-sensitive	Linear bounded
type-2	Context-free	Context-free	Pushdown aut.
type-3	Regular	Regular	Finite aut.

Languages Definitions : regular language

Introduction to Pattern Discovery (43/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Definition (Regular language)

Let Σ be a finite set of symbols (ex : $\Sigma = \{a, b, c\}$), a regular language $L \subseteq \Sigma^*$ recursively defined by mean of *concatenations, unions* and *stars*

Example (Regular language)

 $L = a \cdot b^* \cdot c = \{ac, abc, abbc, abbbc, \ldots\}$

Theorem (Kleene)

$$Rat = Rec$$

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Languages Definitions : finite state automaton

Introduction to Pattern Discovery (44/59)

Cagliari'09

ntroduction

Why Pattern Discovery? Data sets Pattern Matchin Pattern Discove

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences

Definition (Finite state automaton)

Let Σ be an alphabet, a finite automaton is a 5-tuple $\mathcal{A} = \langle \Sigma, Q, I, F, E \rangle$ where Q is a finite set of states, $I \subseteq Q$ is the set of initial states, $F \subseteq Q$ is the set of final states, $E \subseteq Q \times \Sigma \times Q$ is the set of transitions.

Example (Finite state automaton)



Languages Definitions : finite state automaton

Introduction to Pattern

> Discovery (45/59)

Languages

Example (Finite state automaton)



< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

- Path, accepted path
- accepted word
- Language L(A) recognized by A

Languages Maximal automaton of a sample

Introduction to Pattern Discovery (46/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Definition (Maximal automaton)

Given a sample I_+ , the maximal automaton recognizing I_+ is the biggest automaton recognizing I_+ .

Example (Maximal automaton)

 $\textit{I}_{+} = \{\textit{a},\textit{ab},\textit{bab}\}$:



Languages Quotient of an automaton

Introduction to Pattern Discovery (47/59)

Cagliari'09

Why Pattern Discovery? Data sets Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Definition (Partition of a set)

For any set *S*, a partition π is defined by $\pi = \{s \mid s \subseteq S\}$ $\forall s \in \pi, s \neq \emptyset$, such that $\forall s, s' \in \pi, s \cap s' = \emptyset$, $\cup_{s \in \pi} = S$

Definition (Quotient of an automaton)

 $\begin{array}{ll} \mathcal{A}/\pi = \langle \Sigma, Q', I', F', E' \rangle \text{ is a quotient of } \mathcal{A} = \langle \Sigma, Q, I, F, E \rangle \text{ iff} \\ Q' = & Q/\pi \\ I = & \{B \in Q' \mid \exists q \in Bs.t.q \in I\} \\ F = & \{B \in Q' \mid \exists q \in Bs.t.q \in F\} \\ E = & \{\langle B, a, B' \rangle \in Q' \times \Sigma \times Q' \mid \exists q \in B, q' \in B', a \in \Sigma \\ & s.t.\langle q, a, q' \rangle \in E\} \end{array}$

Languages Quotient of an automaton : property

Introduction to Pattern Discovery (48/59)

Cagliari'09

Introduction Why Pattern Discovery? Data sets Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Property

if \mathcal{A}/π is the quotient of \mathcal{A} w.r.t. relation π , then $L(\mathcal{A}) \subseteq L(\mathcal{A}/\pi)$

- any derivation produces a generalisation,
- the set of automata obtained by subsequent derivation forms a lattice.

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Languages Example : partition π_1



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─のへで

Languages Example : partition π_2



 $L(\mathcal{A}/\pi_2) = a^*b(ba^*b)^*$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Languages Example : partition π_3

Introduction to Pattern Discovery (51/59)

Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages



$$\pi_3 = \{\{0\}, \{1,3\}, \{2,4\}\}$$



Languages Example partition π_4



Cagliari'09

Introduction

Why Pattern Discovery? Data sets Pattern Matching

Association Patterns

Association Rules Mining

Sequential Patterns Sequences Languages



 $L(\mathcal{A}) = a(aa)^*b + ab(bb)^*$

 $\pi_{4} = \{\{0, 1, 3, 2, 4\}\}$



Languages The mining task

Introduction to Pattern Discovery (53/59)

Cagliari'09

Introduction

Why Pattern Discovery ? Data sets Pattern Matching Pattern Discover

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

- Find an automaton that generalizes, *not too much*, the input sample : characterizable, heuristic methods ;
- the space of possibilities is the number of partitions of the set of states of the sample maximal automaton;

(日) (日) (日) (日) (日) (日) (日)

• the size of the exploration space is the number of partitions of the maximal automaton coding the

sample : $B_{n+1} = \sum_{k=0}^{n} \begin{pmatrix} n \\ k \end{pmatrix} B_k$

 $\begin{array}{l} B_5 = 52, \ldots, B_{10} = 115975, B_{20} \approx 5.10^{13}, B_{30} \approx 8.10^{23}, \\ B_{50} \approx 2.10^{49} \ldots \end{array}$

Introduction to Pattern Discovery (54/59)

Cagliari'09

Introduction Why Pattern Discovery? Data sets Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Definition (*k*-deterministic automaton)

 \mathcal{A} is *k*-deterministic if, from any state, there is at most one path for any given word of length $\geq k$.



(日) (日) (日) (日) (日) (日) (日)

Deterministic automaton are 0-deterministic.

Introduction to Pattern Discovery (55/59)

Cagliari'09

ntroduction

Why Pattern Discovery? Data sets Pattern Matching Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Definition (Reverse automaton)

The reverse automaton of $\mathcal{A} = \langle \Sigma, Q, I, F, E \rangle$ is defined by $\mathcal{A}^R = \langle \Sigma, Q, F, I, E' \rangle$ with $E' = \{ \langle q, a, q' \rangle \mid \langle q', a, q \rangle \in E \}.$

Example (Reverse automaton)



Introduction to Pattern Discovery (56/59)

Cagliari'09

Introduction Why Pattern Discovery? Data sets Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

Definition (k-reversible automaton)

An automaton is k-reversible if it is deterministic and its reverse automaton is k-deterministic.

Dana Angluin [Ang82] gave a characterization of reversible languages and an inference algorithm : k - RI.

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

Dana Angluin.

Inference of Reversible Languages, Journal of the ACM, Volume 29, 1982.

Introduction to Pattern Discovery (57/59)

Cagliari'09

The k - RI algorithm

```
Why Pattern
Discovery?
Data sets
Pattern Matching
Pattern Discovery
```

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages

$\begin{array}{l} k \leftarrow \text{ order of the model, } I_+ \leftarrow \text{ positive sample} \\ \mathcal{A} \leftarrow \mathsf{PTA}(I_+) \\ \pi \leftarrow Q \\ \text{while } \neg k \text{-reversible}(\mathcal{A}/\pi) \\ (B_1, B_2) \leftarrow \text{non-reversible}(\mathcal{A}/\pi, \pi) \\ \pi \leftarrow \pi \setminus \{B_1, B_2\} \cup \{B_1 \cup B_2\} \\ \\ \} \\ \text{return } \mathcal{A}/\pi \end{array}$

▲□▶▲□▶▲□▶▲□▶ □ のQ@



Cagliari'09

Introductior Why Pattern

Discovery ? Data sets Pattern Matching v Pattern Discovery

Association Patterns Frequent sets Association Rule: Mining

Sequential Patterns Sequences



Going further ...

Some references to go further ...

Introduction to Pattern Discovery (59/59)

Cagliari'09

Introduction Why Pattern Discovery ? Data sets Pattern Matching vs Pattern Discovery

Association Patterns Frequent sets Association Rules Mining

Sequential Patterns Sequences Languages



Dissister of Data Mising Da

Principles of Data Mining, *David J. Hand, Heikki Mannila and Padhraic Smyth*, ISBN-13 : 978-0-262-08290-7

Principles of Data Mining, Undergraduate Topic in Computer Science, *Max Bramer*, ISBN-13 : 978-1-846-28765-7

WEKA : associations.apriori