Data Mining

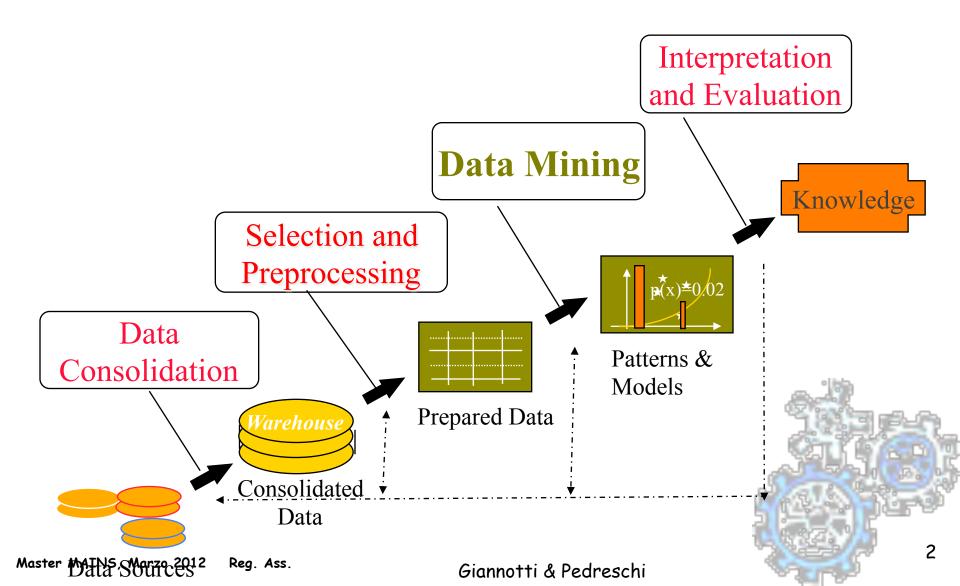
Knowledge Discovery in Databases

Fosca Giannotti and Dino Pedreschi Pisa KDD Lab, ISTI-CNR & Univ. Pisa

Slides available at:

ttp://didawiki.cli.di.unipi.it

KDD Process



Association rules and market basket analysis



Association rules - module outline

- What are association rules (AR) and what are they used for:
 - The paradigmatic application: Market Basket Analysis
 - The single dimensional AR (intra-attribute)

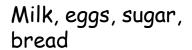


- Basic Apriori Algorithm and its optimizations
- Multi-Dimension AR (inter-attribute)
- Quantitative AR
- Constrained AR
- How to reason on AR and how to evaluate their quality
 - Multiple-level AR
 - Interestingness
 - Correlation vs. Association



Market Basket Analysis: the context

Customer buying habits by finding associations and correlations between the different items that customers place in their "shopping basket"





Milk, eggs, cereal, bread



Customer2

Eggs, sugar

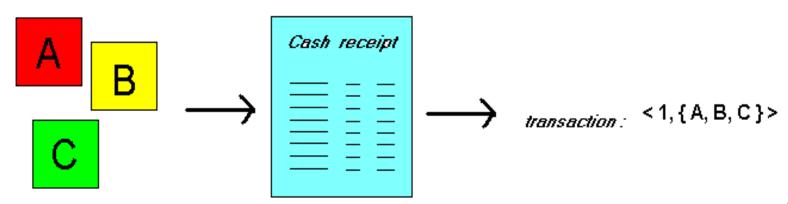


Customer3

Market Basket Analysis: the context

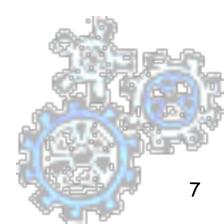
Given: a database of customer transactions, where each transaction is a set of items

I Find groups of items which are frequently purchased together



Goal of MBA

- Extract information on purchasing behavior
- Actionable information: can suggest
 - new store layouts
 - new product assortments
 - which products to put on promotion
- MBA applicable whenever a customer purchases multiple things in proximity
 - credit cards
 - services of telecommunication companies
 - banking services
 - medical treatments



MBA: applicable to many other contexts

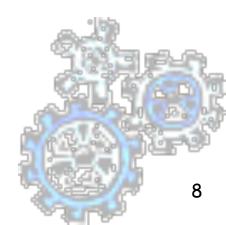
Telecommunication:

Each customer is a transaction containing the set of customer's phone calls

Atmospheric phenomena:

Each time interval (e.g. a day) is a transaction containing the set of observed event (rains, wind, etc.)

Etc.



Association Rules

- Express how product/services relate to each other, and tend to group together
- "if a customer purchases three-way calling, then will also purchase call-waiting"
- simple to understand
- actionable information: bundle three-way calling and call-waiting in a single package
- Examples.
 - Rule form: "Body → Head [support, confidence]".
 - buys(x, "diapers") \rightarrow buys(x, "beers") [0.5%, 60%]
 - major(x, "CS") ^ takes(x, "DB") → grade(x, "A") [1%, 75%]

Useful, trivial, unexplicable

- Useful: "On Thursdays, grocery store consumers often purchase diapers and beer together".
- Trivial: "Customers who purchase maintenance agreements are very likely to purchase large appliances".
- Unexplicable: "When a new hardaware store opens, one of the most sold items is toilet rings."

Association Rules Road Map

- Single dimension vs. multiple dimensional AR
 - E.g., association on items bought vs. linking on different attributes.
 - Intra-Attribute vs. Inter-Attribute
- Qualitative vs. quantitative AR
 - Association on categorical vs. numerical attributes
- Simple vs. constraint-based AR
 - E.g., small sales (sum < 100) trigger big buys (sum > 1,000)?
- Single level vs. multiple-level AR
 - E.g., what brands of beers are associated with what brands of diapers?
- Association vs. correlation analysis.
 - Association does not necessarily imply correlation.

11

Association Rule Mining

Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction

Market-Basket transactions

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

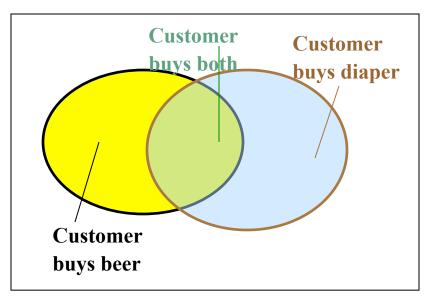
Example of Association Rules

```
{Diaper} → {Beer},
{Milk, Bread} → {Eggs,Coke},
{Beer, Bread} → {Milk},
```

Implication means co-occurrence, not causality!

Basic Concepts: Frequent Patterns and Association Rules

Transaction-id	Items bought
10	A, B, D
20	A, C, D
30	A, D, E
40	B, E, F
50	B, C, D, E, F



Itemset $X = \{x_1, ..., x_k\}$

- Find all the rules $X \rightarrow Y$ with minimum support and confidence
 - support, s, probability that a transaction contains $X \cup Y$
 - confidence, c, conditional probability that a transaction having X also contains Y

Let $sup_{min} = 50\%$, $conf_{min} = 50\%$ Freq. Pat.: {A:3, B:3, D:4, E:3, AD:3}

Association rules:

 $A \rightarrow D \ (60\%, 100\%)$

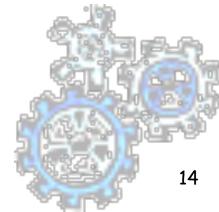
 $D \rightarrow A (60\%, 75\%)$

dicembre 1, 2015

Definition: Frequent Itemset

- Itemset
 - A collection of one or more items
 - ✓ Example: {Milk, Bread, Diaper}
 - k-itemset
 - ✓ An itemset that contains k items
- Support count (σ)
 - Frequency of occurrence of an itemset
 - E.g. $\sigma(\{Milk, Bread, Diaper\}) = 2$
- Support
 - Fraction of transactions that contain an itemset
 - E.g. s({Milk, Bread, Diaper}) = 2/5
- Frequent Itemset
 - An itemset whose support is greater than or equal to a minsup threshold

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke



Definition: Association Rule

- Association Rule
 - An implication expression of the form X → Y, where X and Y are itemsets
 - Example: {Milk, Diaper} → {Beer}

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

- Rule Evaluation Metrics
 - Support (s)
 - Fraction of transactions that contain both X and Y
 - Confidence (c)
 - Measures how often items in Y appear in transactions that contain X

Example:

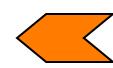
 $\{Milk, Diaper\} \Rightarrow Beer$

$$s = \frac{\sigma(\text{Milk}, \text{Diaper}, \text{Beer})}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

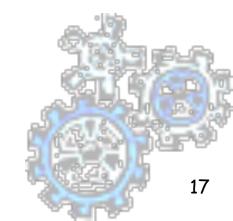
Association rules - module outline

- What are association rules (AR) and what are they used for:
 - The paradigmatic application: Market Basket Analysis
 - The single dimensional AR (intra-attribute)
- How to compute AR
 - Basic Apriori Algorithm and its optimizations
 - Multi-Dimension AR (inter-attribute)
 - Quantitative AR
 - Constrained AR
- How to reason on AR and how to evaluate their quality
 - Multiple-level AR
 - Interestingness
 - Correlation vs. Association



Association Rule Mining Task

- Given a set of transactions T, the goal of association rule mining is to find all rules having
 - support ≥ minsup threshold
 - confidence ≥ minconf threshold
- Brute-force approach:
 - List all possible association rules
 - Compute the support and confidence for each rule
 - Prune rules that fail the minsup and minconf thresholds
 - ⇒ Computationally prohibitive!



Mining Association Rules

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example of Rules:

```
{Milk, Diaper} \rightarrow {Beer} (s=0.4, c=0.67)

{Milk, Beer} \rightarrow {Diaper} (s=0.4, c=1.0)

{Diaper, Beer} \rightarrow {Milk} (s=0.4, c=0.67)

{Beer} \rightarrow {Milk, Diaper} (s=0.4, c=0.67)

{Diaper} \rightarrow {Milk, Beer} (s=0.4, c=0.5)

{Milk} \rightarrow {Diaper, Beer} (s=0.4, c=0.5)
```

Observations:

- All the above rules are binary partitions of the same itemset: {Milk, Diaper, Beer}
- Rules originating from the same itemset have identical support but can have different confidence
- Thus, we may decouple the support and confidence requirements

Mining Association Rules

Two-step approach:

- 1. Frequent Itemset Generation
 - Generate all itemsets whose support ≥ minsup

2. Rule Generation

 Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

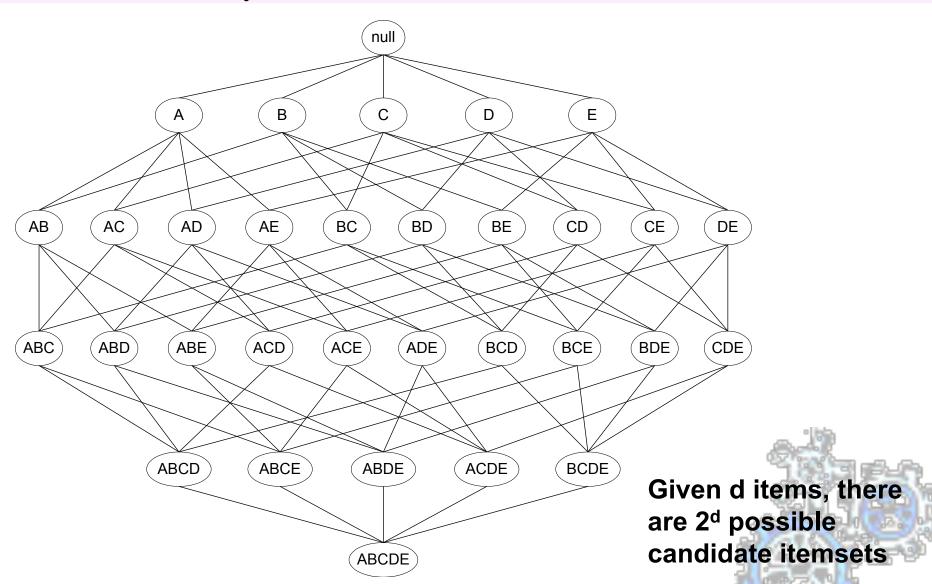
Frequent itemset generation is still computationally expensive

Basic Apriori Algorithm

Problem Decomposition

- ① Find the frequent itemsets: the sets of items that satisfy the support constraint
 - A subset of a frequent itemset is also a frequent itemset,
 i.e., if {A,B} is a frequent itemset, both {A} and {B} should
 be a frequent itemset
 - Iteratively find frequent itemsets with cardinality from 1 to k (k-itemset)
- 2 Use the frequent itemsets to generate association rules.

Frequent Itemset Generation



Frequent Itemset Generation

- Brute-force approach:
 - Each itemset in the lattice is a candidate frequent itemset
 - Count the support of each candidate by scanning the database
 Transactions
 List of

			Candidates
	TID	Items	Canuluates
\blacktriangle	1	Bread, Milk	
Τ	2	Bread, Diaper, Beer, Eggs	
N	3	Milk, Diaper, Beer, Coke	M
1	4	Bread, Milk, Diaper, Beer	
4	5	Bread, Milk, Diaper, Coke	
•			

- Match each transaction against every candidate
- Complexity ~ O(NMw) => Expensive since M = 2d !!!

Frequent Itemset Generation Strategies

- Reduce the number of candidates (M)
 - Complete search: M=2^d
 - Use pruning techniques to reduce M
- Reduce the number of transactions (N)
 - Reduce size of N as the size of itemset increases
 - Used by DHP and vertical-based mining algorithms
- Reduce the number of comparisons (NM)
 - Use efficient data structures to store the candidates or transactions
 - No need to match every candidate against every transaction

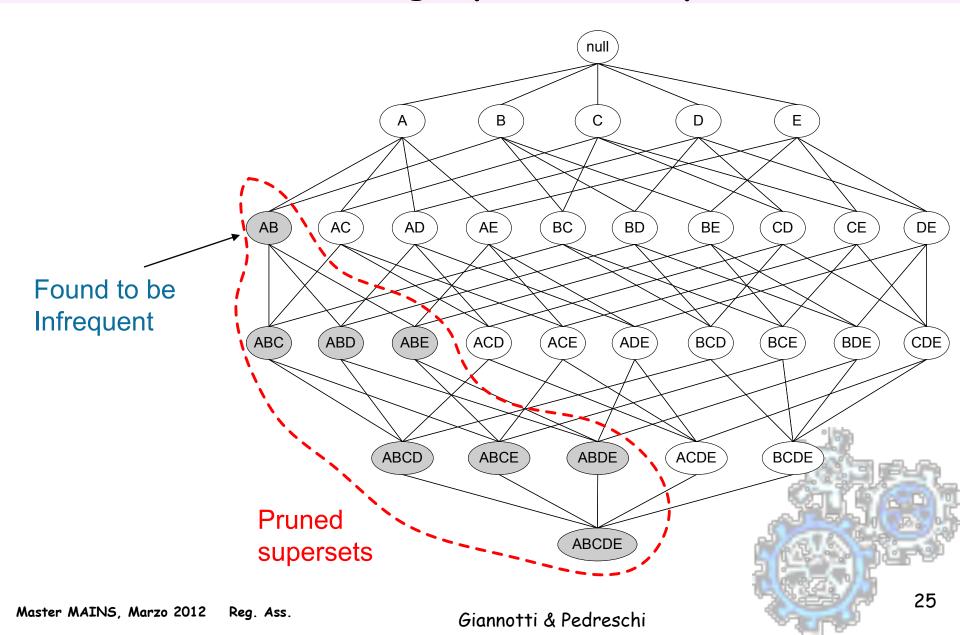
Reducing Number of Candidates

- Apriori principle:
 - If an itemset is frequent, then all of its subsets must also be frequent
- Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \ge s(Y)$$

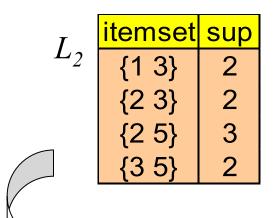
- Support of an itemset never exceeds the support of its subsets
- This is known as the anti-monotone property of support

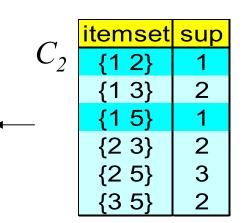
Illustrating Apriori Principle

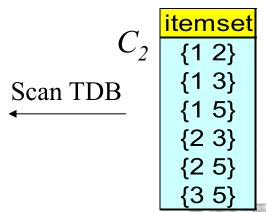


Apriori Execution Example (min_sup = 2)

Data	base TDI	\mathcal{C}	itemset	sup.	, r	itemset	sup.
TID	Items		{1}	2	L_1	{1}	2
100	1 3 4	Scan TDB	{2}	3		{2}	3
200	2 3 5		{3}	3		{3}	3
300	1235		{4 }	1		{ 5 }	3
400	2 5		{5 }	3			







C_{α}	itemset		
<i>C</i> ³	{2 3 5}		

Scan TDB	I_{i2}	item
	— — 3	{23

The Apriori Algorithm

- Join Step: C_k is generated by joining L_{k-1} with itself
- Prune Step: Any (k-1)-itemset that is not frequent cannot be a subset of a frequent k-itemset
- Pseudo-code:

```
C_k: Candidate itemset of size k

L_k: frequent itemset of size k

L_1 = {frequent items};

for (k = 1; L_k != \emptyset; k++) do begin

C_{k+1} = candidates generated from L_k;

for each transaction t in database do

increment the count of all candidates in C_{k+1}

that are contained in t

L_{k+1} = candidates in C_{k+1} with min_support end

return \bigcup_k L_k;
```



How to Generate Candidates?

- Suppose the items in L_{k-1} are listed in an order
- **Step 1:** self-joining L_{k-1}

```
insert into C_k
select p.item_1, p.item_2, ..., p.item_{k-1}, q.item_{k-1}
from L_{k-1} p, L_{k-1} q
where p.item_1 = q.item_1, ..., p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1}
```

Step 2: pruning

```
forall itemsets c in C_k do

forall (k-1)-subsets s of c do

if (s \text{ is not in } L_{k,l}) then delete c from C_k
```

Example of Generating Candidates

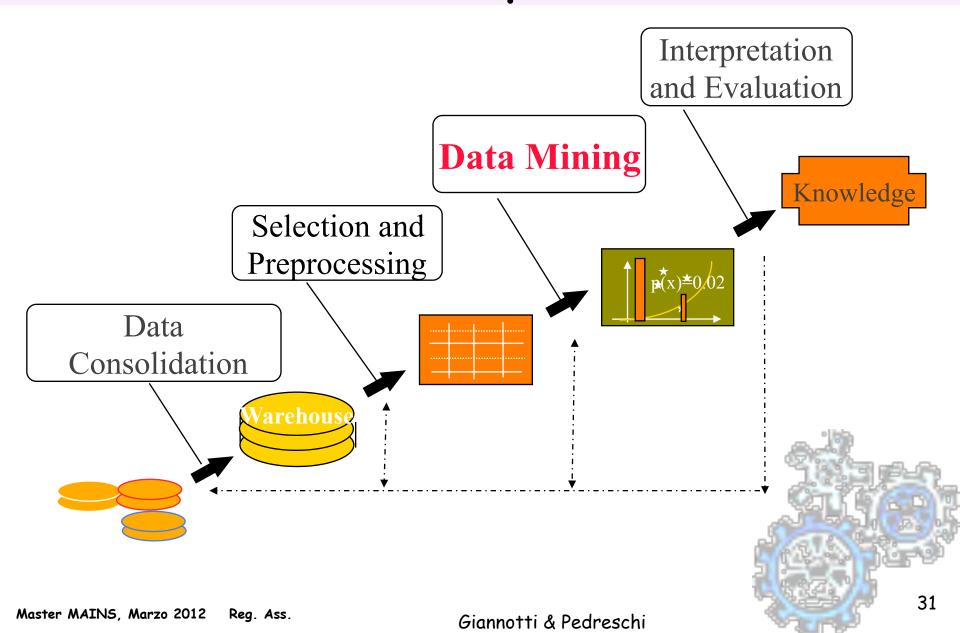
- $L_3=\{abc, abd, acd, ace, bcd\}$
- Self-joining: L,*L,
 - abcd from abc and abd
 - acde from acd and ace
- Pruning:
 - acde is removed because ade is not in L_3
- C₄={abcd}

32

Factors Affecting Complexity

- Choice of minimum support threshold
 - lowering support threshold results in more frequent itemsets
 - this may increase number of candidates and max length of frequent itemsets
- Dimensionality (number of items) of the data set
 - more space is needed to store support count of each item
 - if number of frequent items also increases, both computation and I/O costs may also increase
- Size of database
 - since Apriori makes multiple passes, run time of algorithm may increase with number of transactions
- Average transaction width
 - transaction width increases with denser data sets
 - This may increase max length of frequent itemsets and traversals of hash tree (number of subsets in a transaction increases with its width)

The KDD process



Generating Association Rules from Frequent Itemsets

- Only strong association rules are generated
- **■** Frequent itemsets satisfy minimum support threshold
- Strong rules are those that satisfy minimum confidence threshold
- confidence(A ==> B) = $Pr(B \mid A) = \frac{support(A \cup B)}{support(A)}$

For each frequent itemset, f, generate all non-empty subsets of f
For every non-empty subset s of f do

if support(f)/support(s) ≥ min_confidence then

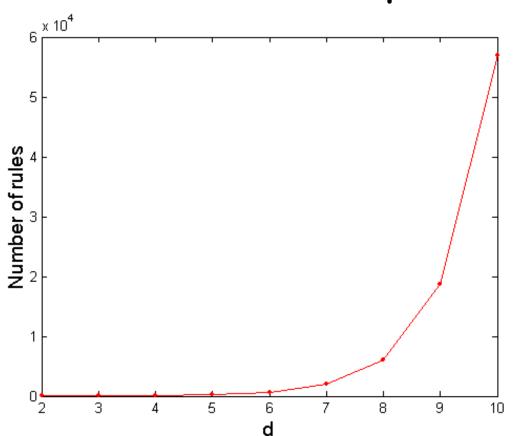
output rule s ==> (f-s)

end

Computational Complexity

Given d unique items:

- Total number of itemsets = 2^d
- Total number of possible association rules:



$$R = \sum_{k=1}^{d-1} \left[\binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right]$$
$$= 3^{d} - 2^{d+1} + 1$$

If d=6, R = 602 rules

33

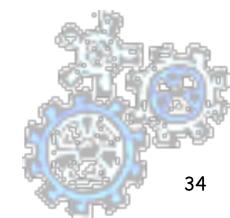
Master MAINS, Marzo 2012

Rule Generation

- Given a frequent itemset L, find all non-empty subsets $f \subset L$ such that $f \to L f$ satisfies the minimum confidence requirement
 - If {A,B,C,D} is a frequent itemset, candidate rules:

$$ABC \rightarrow D$$
, $ABD \rightarrow C$, $ACD \rightarrow B$, $BCD \rightarrow A$, $A \rightarrow BCD$, $B \rightarrow ACD$, $C \rightarrow ABD$, $D \rightarrow ABC$, $AB \rightarrow CD$, $AC \rightarrow BD$, $AD \rightarrow BC$, $BC \rightarrow AD$, $BD \rightarrow AC$, $CD \rightarrow AB$,

If |L| = k, then there are $2^k - 2$ candidate association rules (ignoring $L \to \emptyset$ and $\emptyset \to L$)



Rule Generation

- How to efficiently generate rules from frequent itemsets?
 - In general, confidence does not have an anti-monotone property

 $c(ABC \rightarrow D)$ can be larger or smaller than $c(AB \rightarrow D)$

- But confidence of rules generated from the same itemset has an anti-monotone property
- e.g., L = {A,B,C,D}:

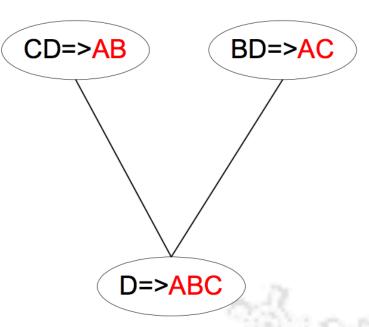
$$c(ABC \rightarrow D) \ge c(AB \rightarrow CD) \ge c(A \rightarrow BCD)$$

✓ Confidence is anti-monotone w.r.t. number of items on the RHS
of the rule

Reg. Ass.

Rule Generation for Apriori Algorithm

- Candidate rule is generated by merging two rules that share the same prefix in the rule consequent
- join(CD=>AB,BD=>AC)
 would produce the candidate
 rule D => ABC
- Prune rule D=>ABC if its subset AD=>BC does not have high confidence



Reg. Ass.

Association rules - module outline

- What are association rules (AR) and what are they used for:
 - The paradigmatic application: Market Basket Analysis
 - The single dimensional AR (intra-attribute)
- How to compute AR
 - Basic Apriori Algorithm and its optimizations
 - Multi-Dimension AR (inter-attribute)
 - Quantitative AR
 - Constrained AR
- How to reason on AR and how to evaluate their quality
 - Multiple-level AR
 - Interestingness
 - Correlation vs. Association



Multidimensional AR

Associations between values of different attributes:

CID	nationality	age	income
1	Italian	50	low
2	French	40	high
3	French	30	high
4	Italian	50	medium
5	Italian	45	high
6	French	35	high

RULES:

nationality = French \Rightarrow income = high [50%, 100%] income = high \Rightarrow nationality = French [50%, 75%] age = 50 \Rightarrow nationality = Italian [33%, 100%]

Single-dimensional vs Multi-dimensional AR

Multi-dimensional

<1, Italian, 50, low>

<2, French, 45, high>

Single-dimensional

<1, {nat/Ita, age/50, inc/low}>

<2, {nat/Fre, age/45, inc/high}>

Schema: <ID, a?, b?, c?, d?>

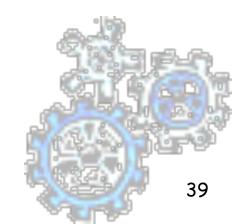
<1, yes, yes, no, no>

<2, yes, no, yes, no>



<1, {a, b}>

<2, {a, c}>



Quantitative Attributes

- Quantitative attributes (e.g. age, income)
- Categorical attributes (e.g. color of car)

CID	height	weight	income
1	168	75,4	30,5
2	175	80,0	20,3
3	174	80,0 70,3 65,2	30,5 20,3 25,8
4	170	65,2	27,0

Problem: too many distinct values

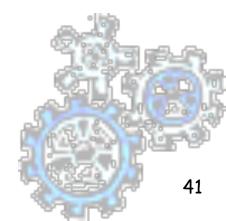
Solution: transform quantitative attributes in categorical ones via discretization.

Quantitative Association Rules

CID	Age	Married	NumCars
1	23	No	1
2	25	Yes	1
3	29	No	0
4	34	Yes	2
5	38	Yes	2

[Age: 30..39] and [Married: Yes] \Rightarrow [NumCars:2]

support = 40% confidence = 100%



Discretization of quantitative attributes

Solution: each value is replaced by the interval to which it belongs.

height: 0-150cm, 151-170cm, 171-180cm, >180cm

weight: 0-40kg, 41-60kg, 60-80kg, >80kg

income: 0-10ML, 11-20ML, 20-25ML, 25-30ML, >30ML

CID	height	weight	income
1	151-171	60-80	>30
2	171-180	60-80	20-25
3	171-180	60-80	25-30
4	151-170	60-80	25-30

Problem: the discretization may be useless (see weight).

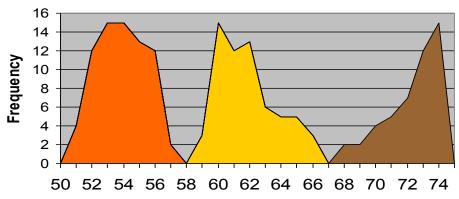
How to choose intervals?

- 1. Interval with a fixed "reasonable" granularity Ex. intervals of 10 cm for height.
- 2. Interval size is defined by some domain dependent criterion

Ex.: 0-20ML, 21-22ML, 23-24ML, 25-26ML, >26ML

3. Interval size determined by analyzing data, studying the distribution or using clustering





weight

50 - 58 kg 59-67 kg

> 68 kg

43

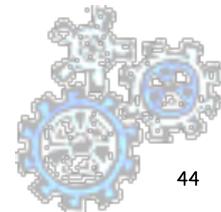
Discretization of quantitative attributes

- 1. Quantitative attributes are statically discretized by using predefined concept hierarchies:
 - elementary use of background knowledge

Loose interaction between Apriori and discretizer

- 2. Quantitative attributes are dynamically discretized
 - into "bins" based on the distribution of the data.
 - considering the distance between data points.

Tighter interaction between Apriori and discretizer



Constraints and AR

- Preprocessing: use constraints to focus on a subset of transactions
 - Example: find association rules where the prices of all items are at most 200 Euro
- Optimizations: use constraints to optimize Apriori algorithm
 - Anti-monotonicity: when a set violates the constraint, so does any of its supersets.
 - Apriori algorithm uses this property for pruning
- Push constraints as deep as possible inside the frequent set computation

Constraint-based AR

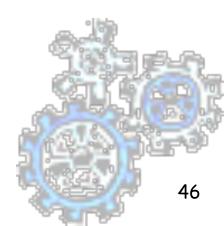
What kinds of constraints can be used in mining?

Data constraints:

- √ SQL-like queries
 - Find product pairs sold together in Vancouver in Dec. '98.
- ✓ OLAP-like queries (Dimension/level)
 - in relevance to region, price, brand, customer category.

Rule constraints:

- ✓ specify the form or property of rules to be mined.
- ✓ Constraint-based AR



Rule Constraints

- Two kind of constraints:
 - Rule form constraints: meta-rule guided mining.
 - ✓ $P(x, y) ^ Q(x, w) \rightarrow takes(x, "database systems").$
 - Rule content constraint: constraint-based query optimization (Ng, et al., SIGMOD'98).
 - √ sum(LHS) < 100 ^ min(LHS) > 20 ^ sum(RHS) > 1000
- 1-variable vs. 2-variable constraints (Lakshmanan, et al. SIGMOD'99):
 - 1-var: A constraint confining only one side (L/R)
 of the rule, e.g., as shown above.
 - 2-var: A constraint confining both sides (L and R).
 - √ sum(LHS) < min(RHS) ^ max(RHS) < 5* sum(LHS)</p>

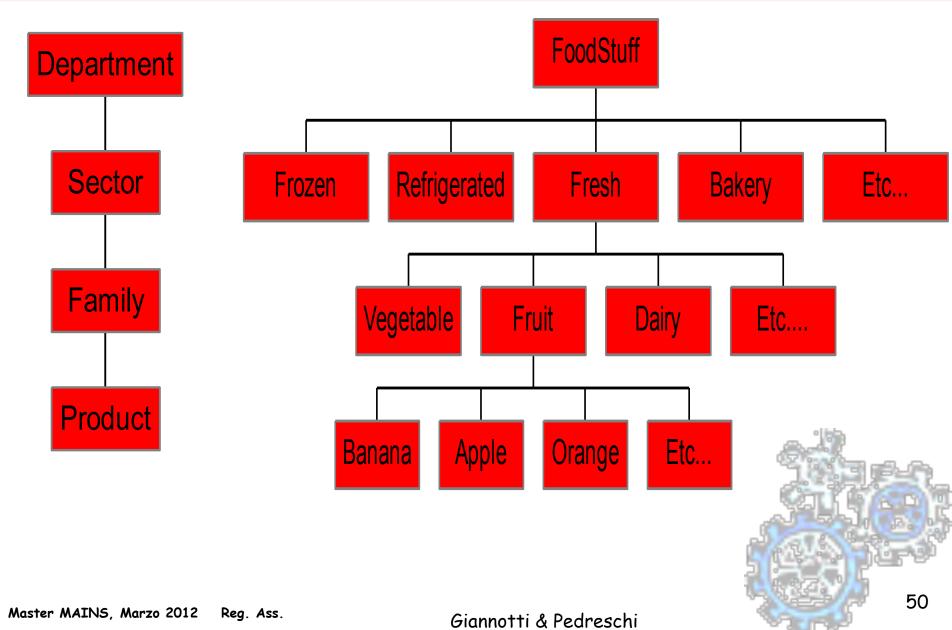
Association rules - module outline

- What are association rules (AR) and what are they used for:
 - The paradigmatic application: Market Basket Analysis
 - The single dimensional AR (intra-attribute)
- How to compute AR
 - Basic Apriori Algorithm and its optimizations
 - Multi-Dimension AR (inter-attribute)
 - Quantitative AR
 - Constrained AR
- How to reason on AR and how to evaluate their quality
 - Multiple-level AR
 - Interestingness
 - Correlation vs. Association

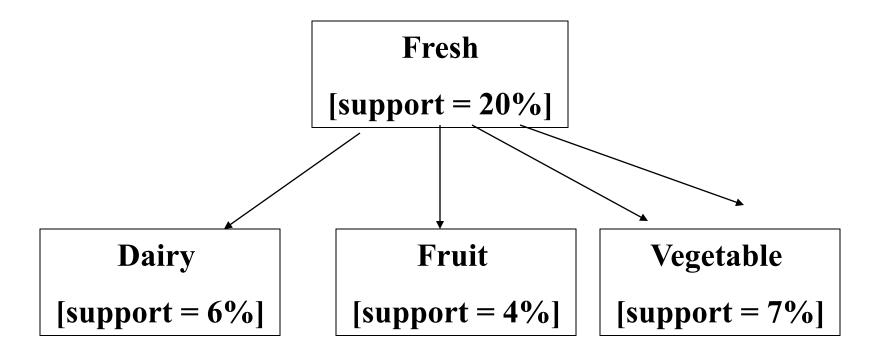
Multilevel AR

- Is difficult to find interesting patterns at a too primitive level
 - high support = too few rules
 - low support = too many rules, most uninteresting
- Approach: reason at suitable level of abstraction
- A common form of background knowledge is that an attribute may be generalized or specialized according to a hierarchy of concepts
- Dimensions and levels can be efficiently encoded in transactions
- Multilevel Association Rules: rules which combine associations with hierarchy of concepts

Hierarchy of concepts



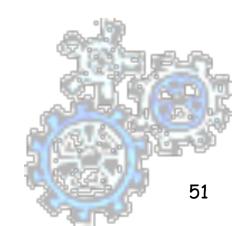
Multilevel AR



Fresh \Rightarrow Bakery [20%, 60%]

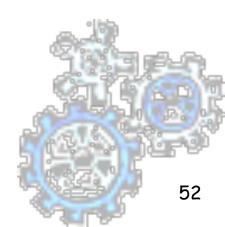
Dairy \Rightarrow Bread [6%, 50%]

Fruit \Rightarrow Bread [1%, 50%] is not valid

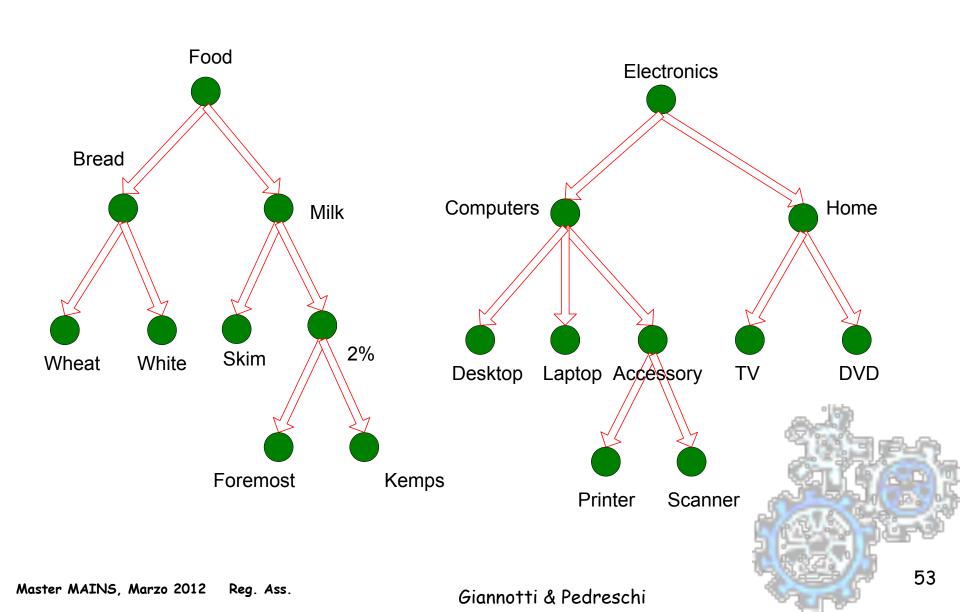


Support and Confidence of Multilevel AR

- from specialized to general: support of rules increases (new rules may become valid)
- from general to specialized: support of rules decreases (rules may become not valid, their support falls under the threshold)
- Confidence is not affected



Multi-level Association Rules



Multi-level Association Rules

- Why should we incorporate concept hierarchy?
 - Rules at lower levels may not have enough support to appear in any frequent itemsets
 - Rules at lower levels of the hierarchy are overly specific
 - ✓ e.g., skim milk → white bread, 2% milk → wheat bread, skim milk → wheat bread, etc.

are indicative of association between milk and bread

54

Reasoning with Multilevel AR

■ Too low level => too many rules and too primitive.

Example: Apple Melinda ⇒ Colgate Tooth-paste

It is a curiosity not a behavior

- Too high level => uninteresting rules
 - Example: Foodstuff ⇒ Varia
- Redundancy => some rules may be redundant due to "ancestor" relationships between items.
 - A rule is redundant if its support is close to the "expected" value, based on the rule's ancestor.
- Example (milk has 4 subclasses)
 - milk ⇒ wheat bread, [support = 8%, confidence = 70%]
 - 2%-milk ⇒ wheat bread, [support = 2%, confidence = 72%]

Mining Multilevel AR

- Calculate frequent itemsets at each concept level, until no more frequent itemsets can be found
- For each level use Apriori
- A top_down, progressive deepening approach:
 - First find high-level strong rules:

fresh \rightarrow bakery [20%, 60%].

Then find their lower-level "weaker" rules:

fruit \rightarrow bread [6%, 50%].

- Variations at mining multiple-level association rules.
 - Level-crossed association rules:

fruit → wheat bread

- Association rules with multiple, alternative hierarchies:

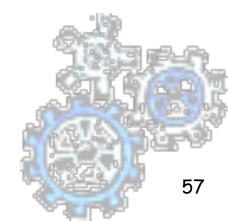
fruit → Wonder bread

Reasoning with AR

Significance:

```
Example: <1, {a, b}>
<2, {a} >
<3, {a, b, c}>
<4, {b, d}>
```

 $\{b\} \Rightarrow \{a\}$ has confidence (66%), but is not significant as support($\{a\}$) = 75%.



Beyond Support and Confidence

Example 1: (Aggarwal & Yu, PODS98)

	coffee	not coffee	sum(row)
tea	20	5	25
not tea	70	5	75
sum(col.)	90	10	100

- {tea} => {coffee} has high support (20%) and confidence (80%)
- However, a priori probability that a customer buys coffee is 90%
 - A customer who is known to buy tea is less likely to buy coffee (by 10%)
 - There is a negative correlation between buying tea and buying coffee
 - {~tea} => {coffee} has higher confidence(93%)

Correlation and Interest

- Two events are independent if $P(A \land B) = P(A)*P(B)$, otherwise are correlated.
- Interest = $P(A \land B) / P(B)*P(A)$
- Interest expresses measure of correlation
 - $= 1 \Rightarrow A$ and B are independent events
 - less than $1 \Rightarrow A$ and B negatively correlated,
 - greater than $1 \Rightarrow A$ and B positively correlated.
 - In our example, I(buy tea ∧ buy coffee)=0.89 i.e. they are negatively correlated.

Computing Interestingness Measure

■ Given a rule X → Y, information needed to compute rule interestingness can be obtained from a contingency table

Contingency table for $X \rightarrow Y$

	У	À	
X	f ₁₁	f ₁₀	f ₁₊
X	f ₀₁	f ₀₀	f _{o+}
	f ₊₁	f ₊₀	ITI

 f_{11} : support of X and Y f_{10} : support of X and Y f_{01} : support of X and Y f_{00} : support of X and Y

Used to define various measures

support, confidence, lift, Gini,
 J-measure, etc.

Statistical-based Measures

Measures that take into account statistical dependence

$$Lift = \frac{P(Y \mid X)}{P(Y)}$$

$$P(X \mid X) = \frac{P(Y \mid X)}{P(X \mid X)}$$

$$Interest = \frac{P(X,Y)}{P(X)P(Y)}$$

$$PS = P(X,Y) - P(X)P(Y)$$

$$\phi - coefficient = \frac{P(X,Y) - P(X)P(Y)}{\sqrt{P(X)[1 - P(X)]P(Y)[1 - P(Y)]}}$$

Example: Lift/Interest

	Coffe e	Coffe e	
Tea	15	5	20
Tea	75	5	80
	90	10	100

Association Rule: Tea → Coffee

Confidence = P(Coffee|Tea) = 0.75

but P(Coffee) = 0.9

 \Rightarrow Lift = 0.75/0.9= 0.8333 (< 1, therefore is negatively associated)

Drawback of Lift & Interest

	У	À	
×	10	0	10
×	0	90	90
	10	90	100

	У	Ā	
X	90	0	90
×	0	10	10
	90	10	100

$$Lift = \frac{0.1}{(0.1)(0.1)} = 10$$

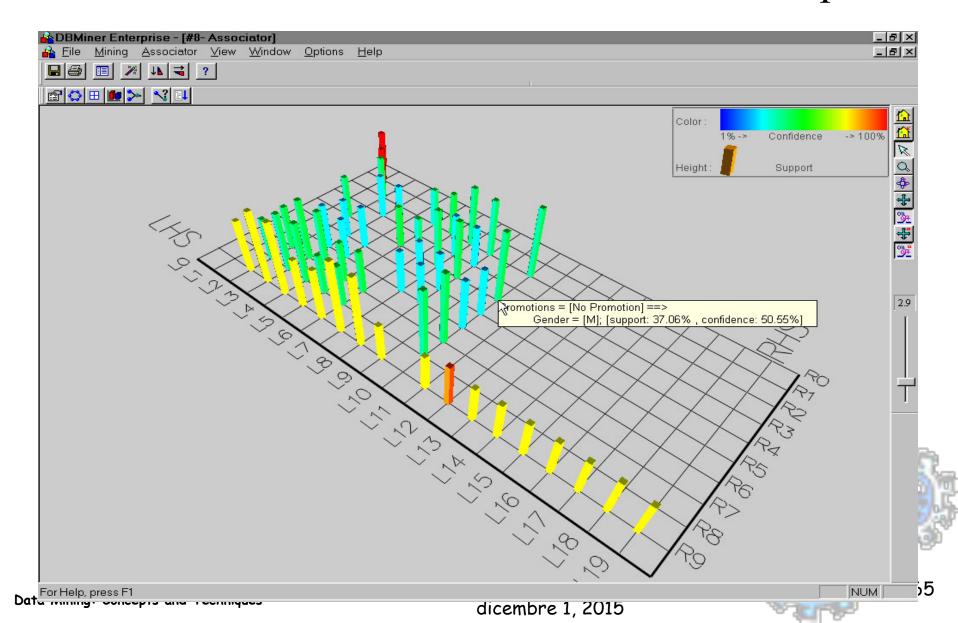
$$Lift = \frac{0.9}{(0.9)(0.9)} = 1.11$$

Statistical independence:

If
$$P(X,Y)=P(X)P(Y) \Rightarrow Lift = 1$$

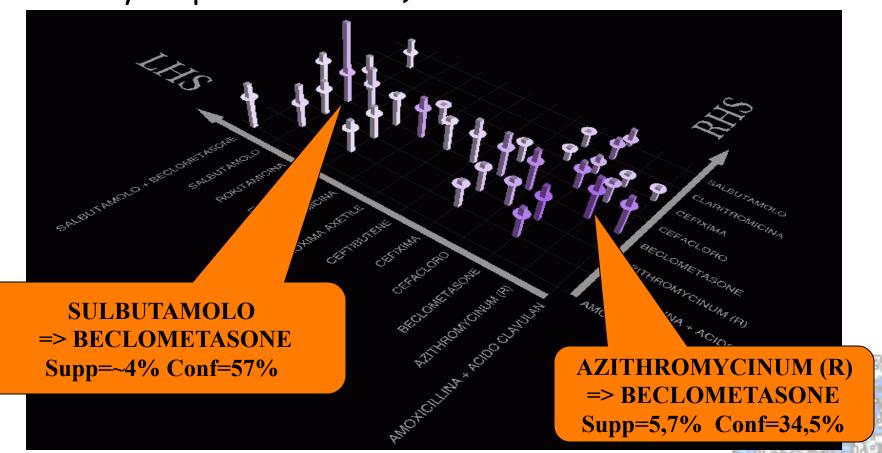
	#	Measure	Formula
The second of the set	1	ϕ -coefficient	P(A,B)-P(A)P(B)
There are lots of	١	·	$\frac{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}{\sum_{j} \max_{k} P(A_{j}, B_{k}) + \sum_{k} \max_{j} P(A_{j}, B_{k}) - \max_{j} P(A_{j}) - \max_{k} P(B_{k})}{2-\max_{j} P(A_{j}) - \max_{k} P(B_{k})}$
measures proposed	2	Goodman-Kruskal's (λ)	$\frac{2-\max_{j}P(A_{j})-\max_{k}P(B_{k})}{P(A,B)P(\overline{A},\overline{B})}$
in the literature	3	$\text{Odds ratio } (\alpha)$	$\overline{P(A,\overline{B})P(\overline{A},B)}$
	4	Yule's Q	$\frac{P(A,B)P(\overline{AB}) - P(A,\overline{B})P(\overline{A},B)}{P(A,B)P(\overline{AB}) + P(A,\overline{B})P(\overline{A},B)} = \frac{\alpha - 1}{\alpha + 1}$
	5	Yule's Y	$\frac{\sqrt{P(A,B)P(\overline{AB})} - \sqrt{P(A,\overline{B})P(\overline{A},B)}}{\sqrt{P(A,B)P(\overline{AB})} + \sqrt{P(A,\overline{B})P(\overline{A},B)}} = \frac{\sqrt{\alpha} - 1}{\sqrt{\alpha} + 1}$
Some measures are	6	Kappa (κ)	$\begin{array}{c} V (A,B) + (A,B)$
good for certain	"	rappa (w)	$\frac{1 - P(A)P(B) - P(\overline{A})P(\overline{B})}{\sum_{i} \sum_{j} P(A_{i}, B_{j}) \log \frac{P(A_{i}, B_{j})}{P(A_{i})P(B_{j})}}$
applications, but not	7	Mutual Information (M)	$\frac{\sum_{i} \sum_{j} P(A_i) P(B_j)}{\min(-\sum_{i} P(A_i) \log P(A_i), -\sum_{j} P(B_j) \log P(B_j))}$
for others	8	J-Measure (J)	$\max\left(P(A,B)\log(\frac{P(B A)}{P(B)}) + P(A\overline{B})\log(\frac{P(\overline{B} A)}{P(\overline{B})}),\right)$
			$P(A,B)\log(\frac{P(A B)}{P(A)}) + P(\overline{A}B)\log(\frac{P(\overline{A} B)}{P(\overline{A})})$
	9	Gini index (G)	$\max \left(P(A)[P(B A)^2 + P(\overline{B} A)^2] + P(\overline{A})[P(B \overline{A})^2 + P(\overline{B} \overline{A})^2] \right)$
What criteria should		(• / • • • • • • • • • • • • • • • • • • •	$-P(B)^2-P(\overline{B})^2$,
we use to determine			$P(B)[P(A B)^{2} + P(\overline{A} B)^{2}] + P(\overline{B})[P(A \overline{B})^{2} + P(\overline{A} \overline{B})^{2}]$
whether a measure			$-P(A)^{2}-P(\overline{A})^{2}$
is good or bad?	10	Support (s)	P(A,B)
io good or bad.	11	Confidence (c)	$\max(P(B A), P(A B))$
	12	Laplace (L)	$\max\left(\frac{NP(A,B)+1}{NP(A)+2},\frac{NP(A,B)+1}{NP(B)+2}\right)$
What about Apriori		-	
What about Apriori-	13	Conviction (V)	$\max\left(rac{P(A)P(\overline{B})}{P(A\overline{B})},rac{P(B)P(\overline{A})}{P(B\overline{A})} ight) \ rac{P(A,B)}{P(A,B)}$
style support based	14	Interest (I)	$ \begin{array}{c} \frac{P(A,B)}{P(A)P(B)} \\ P(A,B) \end{array} $
pruning? How does	15	cosine (IS)	$\frac{1}{\sqrt{P(A)P(B)}}$
it affect these	16	Piatetsky-Shapiro's (PS)	P(A,B) - P(A)P(B)
measures?	17	Certainty factor (F)	$\max\left(rac{P(B A)-P(B)}{1-P(B)},rac{P(A B)-P(A)}{1-P(A)} ight)$
	18	Added Value (AV)	$\max(P(B A) - P(B), P(A B) - P(A))$
	19	Collective strength (S)	$\frac{\frac{P(A,B)+P(\overline{AB})}{P(A)P(B)+P(\overline{A})P(\overline{B})}}{\frac{P(A,B)}{P(A)+P(B)-P(A,B)}} \times \frac{\frac{1-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A,B)-P(\overline{AB})}}{\frac{P(A,B)}{P(A)+P(B)-P(A,B)}}$
	20	Jaccard (ζ)	$\frac{P(A,B)}{P(A)+P(B)-P(A,B)}$
Master MAINS, Marzo 2012 Re	g.2 1 \ss	$\operatorname{Klosgen}(K)$ Gio	P(A,B) mark $P(B A) - P(B), P(A B) - P(A))$

Visualization of Association Rules: Plane Graph

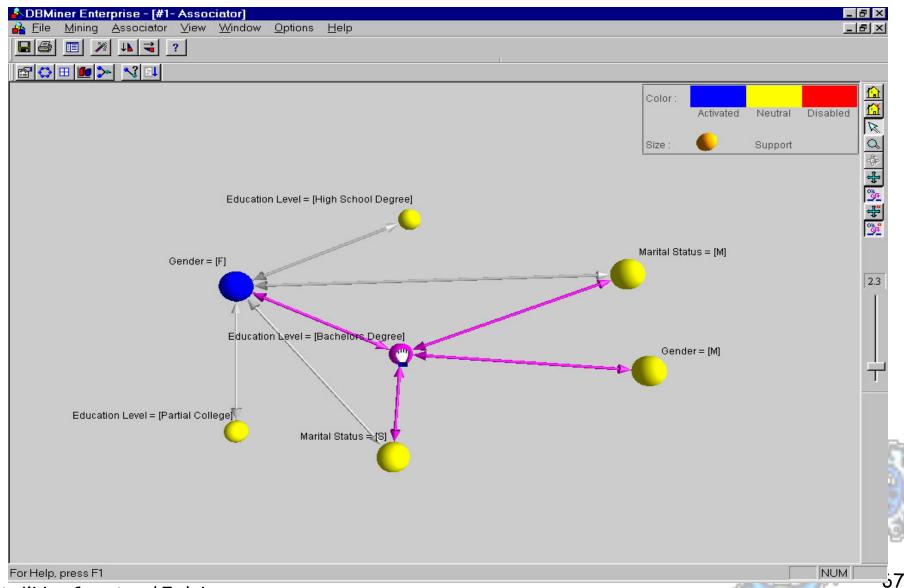


Association Rules – visualization

(Patients <15 old for USL 19 (a unit of Sanitary service), January-September 1997)



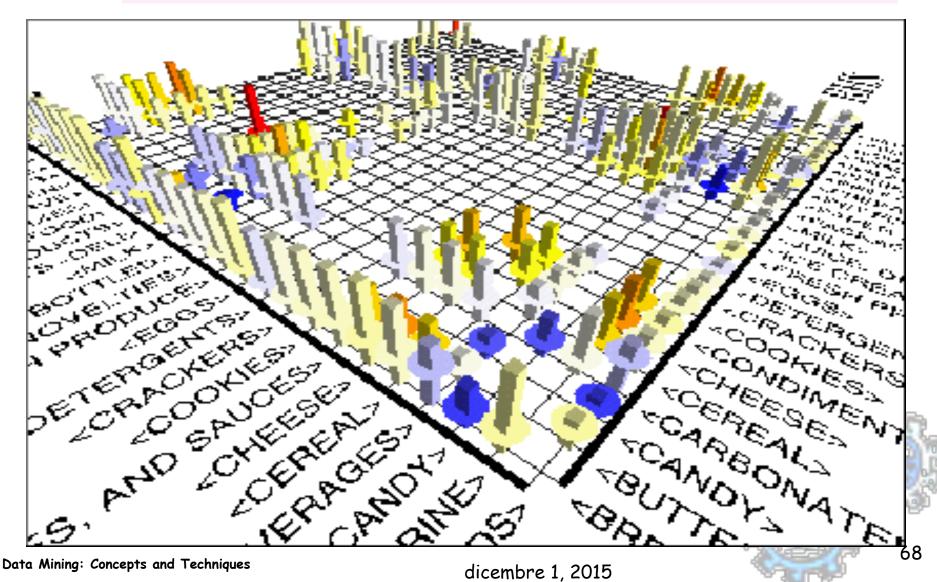
Visualization of Association Rules: Rule Graph



Data Mining: Concepts and Techniques

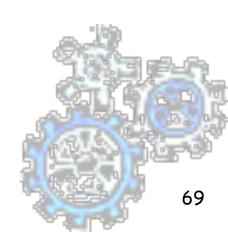
dicembre 1, 2015

Visualization of Association Rules (SGI/MineSet 3.0)



Domain dependent measures

- Together with support, confidence, interest, ..., use also (in post-processing) domain-dependent measures
- E.g., use rule constraints on rules
- Example: take only rules which are significant with respect their economic value
- sum(LHS)+ sum(RHS) > 100



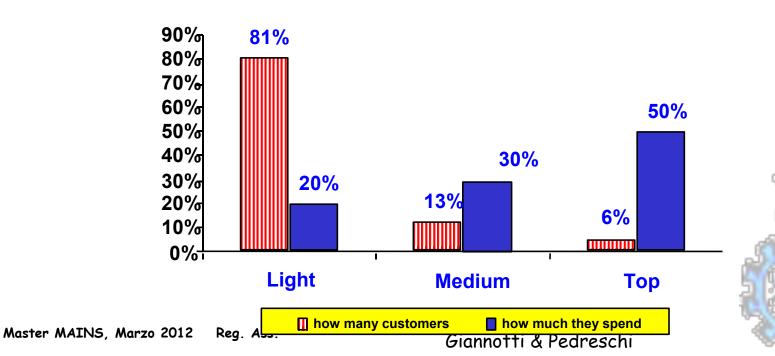
Conclusions

- Association rule mining
 - probably the most significant contribution from the database community to KDD
 - A large number of papers have been published
- Many interesting issues have been explored
- An interesting research direction
 - Association analysis in other types of data: spatial data, multimedia data, time series data, etc.

70

Conclusion (2)

- MBA is a key factor of success in the competition of supermarket retailers.
- Knowledge of customers and their purchasing behavior brings potentially huge added value.

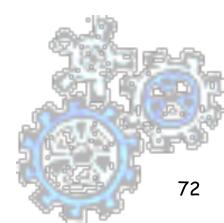


71

Which tools for market basket analysis?

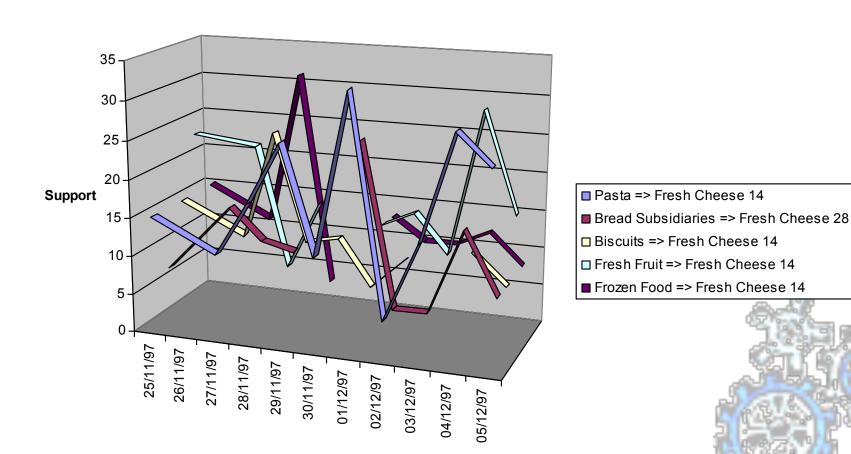
Association rule are needed but insufficient

- Market analysts ask for business rules:
 - Is supermarket assortment adequate for the company's target class of customers?
 - Is a promotional campaign effective in establishing a desired purchasing habit?



Business rules: temporal reasoning on AR

- Which rules are established by a promotion?
- How do rules change along time?



Association rules - module2 Examples

Association Rules in Web Miming AR & Atherosclerosis prevention study Moviegoer Data bases



MBA in Web Usage Mining

- Association Rules in Web Transactions
 - discover affinities among sets of Web page references across user sessions

Examples

- 60% of clients who accessed /products/, also accessed /products/software/webminer.htm
- 30% of clients who accessed /special-offer.html, placed an online order in /products/software/
- Actual Example from IBM official Olympics Site:
 - √ {Badminton, Diving} ==> {Table Tennis} [conf = 69.7%, sup = 0.35%]

Applications

- Use rules to serve dynamic, customized contents to users
- prefetch files that are most likely to be accessed

75

Web Usage Mining: Example

Association Rules From Cray Research Web Site

Conf	supp	Association Rule
82.8	3.17	/PUBLIC/product-info/T3E
		===>
		/PUBLIC/product-info/T3E/CRAY_T3E.html
90	0.14	/PUBLIC/product-info/J90/J90.html,
		/PUBLIC/product-info/T3E
		===>
		/PUBLIC/product-info/T3E/CRAY_T3E.html
97.2	0.15	/PUBLIC/product-info/J90,
		/PUBLIC/product-info/T3E/CRAY_T3E.html,
		/PUBLIC/product-info/T90,
		===>
		/PUBLIC/product-info/T3E,
		/PUBLIC/sc.html

Design "suggestions"

from rules 1 and 2: there is something in J90.html that should be moved to th page /PUBLIC/product-info/T3E (why?)

MBA in Text / Web Content Mining

Documents Associations

- Find (content-based) associations among documents in a collection
- Documents correspond to items and words correspond to transactions
- Frequent itemsets are groups of docs in which many words occur in common

	Doc 1	Doc 2	Doc 3	 Doc n
business	5	5	2	 1
capital	2	4	3	 5
fund	0	0	0	 1
:	:	:	•	 :
invest	6	0	0	 3

Term Associations

- Find associations among words based on their occurrences in documents
- similar to above, but invert the table (terms as items, and docs as transactions)

Atherosclerosis prevention study

2nd Department of Medicine, 1st Faculty of Medicine of Charles University and Charles University Hospital, U nemocnice 2, Prague 2 (head. Prof. M. Aschermann, MD, SDr, FESC)

Atherosclerosis prevention study:

The STULONG 1 data set is a real database that keeps information about the study of the development of atherosclerosis risk factors in a population of middle aged men.

Used for Discovery Challenge at PKDD 00-02-03-04

Atherosclerosis prevention study:

- Study on 1400 middle-aged men at Czech hospitals
 - Measurements concern development of cardiovascular disease and other health data in a series of exams
- The aim of this analysis is to look for associations between medical characteristics of patients and death causes.
- Four tables
 - Entry and subsequent exams, questionnaire responses, deaths

The input data

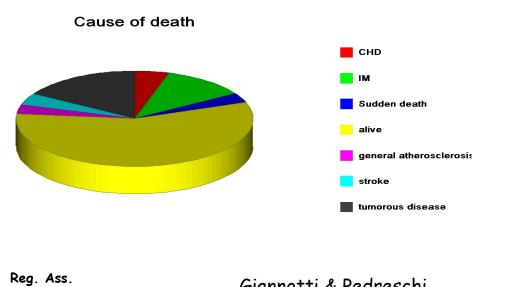
Data f	Data from Entry and Exams				
General characteristics	Examinations	habits			
Marital status	Chest pain	Alcohol			
Transport to a job	Breathlesness	Liquors			
Physical activity in a job	Cholesterol	Beer 10			
Activity after a job	Urine	Beer 12			
Education	Subscapular	Wine			
Responsibility	Triceps	Smoking			
Age		Former smoker			
Weight		Duration of smoking			
Height		Tea			
		Sugar			
		Coffee			

The input data

DEATH CAUSE	PATIENTS	%
myocardial infarction	80	20.6
coronary heart disease	33	8.5
stroke	30	7.7
other causes	79	20.3
sudden death	23	5.9
unknown	8	2.0
tumorous disease	114	29.3
general atherosclerosis	22	5.7
TOTAL	389	100.0

Data selection

- When joining "Entry" and "Death" tables we implicitely create a new attribute "Cause of death", which is set to "alive" for subjects present in the "Entry" table but not in the "Death" table.
- We have only 389 subjects in death table.



The prepared data

Patient	General character	istics	Examination	ons	Habits	Cause of
	Activity after work	Education	Chest pain		Alcohol	 death
1	moderate activity	university	not present		no	Stroke
2	great activity		not ischaemic		occasionally	myocardial infarction
3	he mainly sits		other pains		regularly	tumorous disease
						 alive
389	he mainly sits		other pains		regularly	tumorous disease

Descriptive Analysis/ Subgroup Discovery / Association Rules

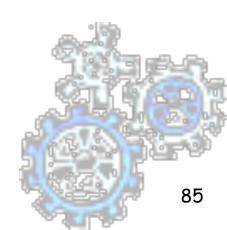
Are there strong relations concerning death cause?

General characteristics $(?) \Rightarrow$ Death cause (?)

Examinations $(?) \Rightarrow Death cause (?)$

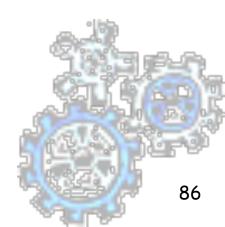
Habits $(?) \Rightarrow$ Death cause (?)

Combinations (?) \Rightarrow Death cause (?)



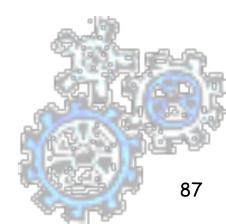
Example of extracted rules

- Education(university) & Height<176-180> ⇒Death cause (tumouros disease), 16; 0.62
- It means that on tumorous disease have died 16, i.e. 62% of patients with university education and with height 176-180 cm.



Example of extracted rules

- Physical activity in work(he mainly sits) & Height<176-180> ⇒ Death cause (tumouros disease), 24; 0.52
- It means that on tumorous disease have died 24 i.e. 52% of patients that mainly sit in the work and whose height is 176-180 cm.

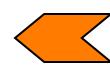


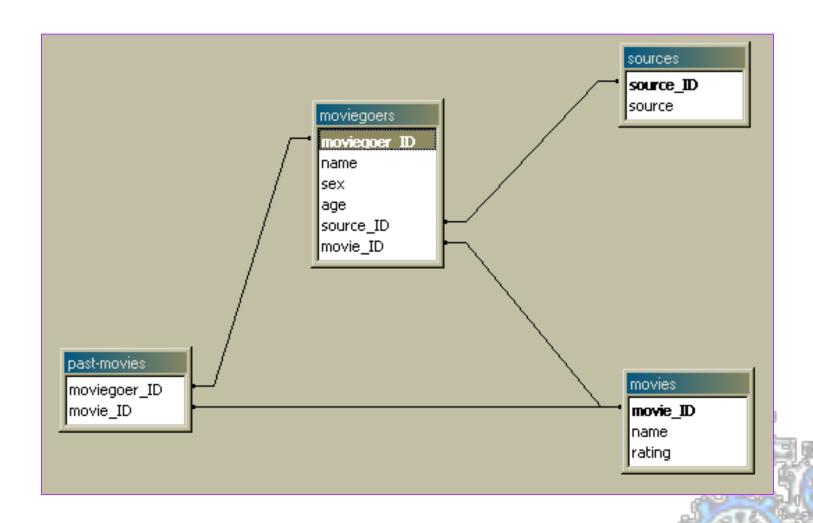
Example of extracted rules

- Education(university) & Height<176-180> ⇒Death cause (tumouros disease), 16; 0.62; +1.1;
- the relative frequency of patients who died on tumorous disease among patients with university education and with height 176-180 cm is 110 per cent higher than the relative frequency of patients who died on tumorous disease among all the 389 observed patients

Association rules - module2 Examples

Association Rules in Web Miming AR & Atherosclerosis prevention study Moviegoer Data bases





moviegoers.name	sex	age	source	movies.name
Amy	f	27	Oberlin	Independence Day
Andrew	m	25	Oberlin	12 Monkeys
Andy	m	34	Oberlin	The Birdcage
Anne	f	30	Oberlin	Trainspotting
Ansje	f	25	Oberlin	I Shot Andy Warhol
Beth	f	30	Oberlin	Chain Reaction
Bob	m	51	Pinewoods	Schindler's List
Brian	m	23	Oberlin	Super Cop
Candy	f	29	Oberlin	Eddie
Cara	f	25	Oberlin	Phenomenon
Cathy	f	39	Mt. Auburn	The Birdcage
Charles	m	25	Oberlin	Kingpin
Curt	m	30	MRJ	T2 Judgment Day
David	m	40	MRJ	Independence Day
Erica	f	23	Mt. Auburn	Trainspotting

Classification

- determine sex based on age, source, and movies seen
- determine source based on sex, age, and movies seen
- determine most recent movie based on past movies, age, sex, and source

Estimation

- for predict, need a continuous variable (e.g., "age")
- predict age as a function of source, sex, and past movies
- if we had a "rating" field for each moviegoer, we could predict the rating a new moviegoer gives to a movie based on age, sex, past movies, etc.

Clustering

- find groupings of movies that are often seen by the same people
- find groupings of people that tend to see the same movies
- clustering might reveal relationships that are not necessarily recorded in the data (e.g., we may find a cluster that is dominated by people with young children; or a cluster of movies that correspond to a particular genre)

Association Rules

- market basket analysis (MBA): "which movies go together?"
- need to create "transactions" for each moviegoer containing movies seen by that moviegoer:

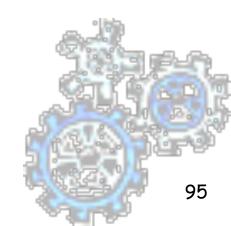
name	TID	Transaction
Amy	001	{Independence Day, Trainspotting}
Andrew	002	{12 Monkeys, The Birdcage, Trainspotting, Phenomenon}
Andy	003	{Super Cop, Independence Day, Kingpin}
Anne	004	{Trainspotting, Schindler's List}
		•••

may result in association rules such as:

```
{"Phenomenon", "The Birdcage"} ==> {"Trainspotting"} {"Trainspotting", "The Birdcage"} ==> {sex = "f"}
```

Sequence Analysis

- similar to MBA, but order in which items appear in the pattern is important
- e.g., people who rent "The Birdcage" during a visit tend to rent "Trainspotting" in the next visit.



Sequential Patterns



Master MAINS, Marzo 2012 Reg. Ass.

Sequential / Navigational Patterns

- Sequential patterns add an extra dimension to frequent itemsets and association rules time.
 - Items can appear before, after, or at the same time as each other.
 - General form: "x% of the time, when A appears in a transaction, B appears within z transactions."
 - ✓ note that other items may appear between A and B, so sequential patterns do not necessarily imply consecutive appearances of items (in terms of time)

Examples

- Renting "Star Wars", then "Empire Strikes Back", then "Return of the Jedi" in that order
- Collection of ordered events within an interval
- Most sequential pattern discovery algorithms are based on extensions of the Apriori algorithm for discovering itemsets

Navigational Patterns

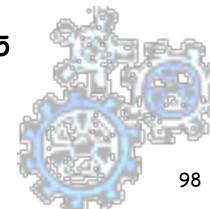
- they can be viewed as a special form of sequential patterns which capture navigational patterns among users of a site
- in this case a session is a consecutive sequence of pageview references for a user over a specified period of time

Mining Sequences - Example

Customer-sequence

CustId	Video sequence
1	$\{(C), (H)\}$
2	$\{(AB), (C), (DFG)\}$
3	{(CEG)}
4	$\{(C), (DG), (H)\}$
5	$\{(H)\}$

Sequential patterns with support > 0.25 $\{(C), (H)\}\$

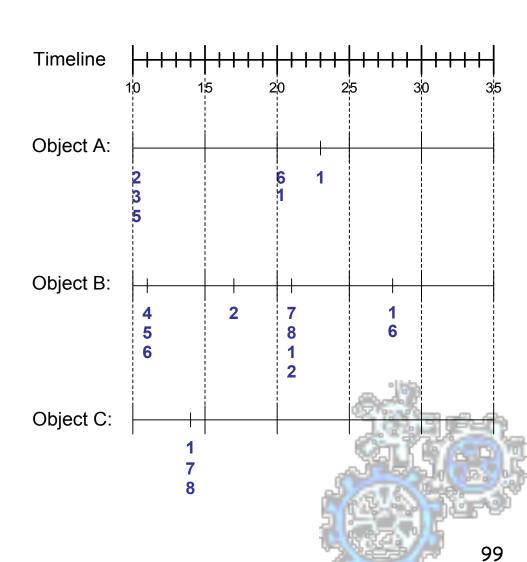


Sequence Data

Giannotti & Pedreschi

Sequence Database:

Object	Timestamp	Events
Α	10	2, 3, 5
Α	20	6, 1
Α	23	1
В	11	4, 5, 6
В	17	2
В	21	7, 8, 1, 2
В	28	1, 6
С	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 Eler	DNA sequence of a particular species	An element of the DNA sequence	Bases A,T,G,C

(Transaction) E1 E1

Sequence

 $\begin{array}{c}
E1 \\
E2
\end{array}$ $\begin{array}{c}
E1 \\
E3
\end{array}$ $\begin{array}{c}
E2
\end{array}$

 $\left(\mathsf{E2}\right)$

(Item)

100

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)

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

- 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 $\langle a_1 a_2 ... a_n \rangle$ is contained in another sequence $\langle b_1 b_2 ... b_m \rangle$ (m \geq n) if there exist integers

 $i_1 < i_2 < ... < i_n$ such that $a_1 \subseteq b_{i1}$, $a_2 \subseteq b_{i1}$, ..., a_n

Daira sequence	Subsequence	Contain?
< {2,4} {3,5,6} {8} >	< {2} {3,5} >	Yes
< {1,2} {3,4} >	< {1} {2} >	No
< {2,4} {2,4} {2,5} >	< {2} {4} >	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 Master MA (Nis, Charzo 2012 Subsequence whose support is > minsup)

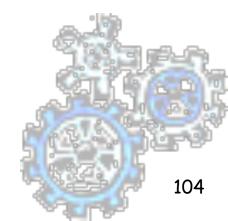
Sequential Pattern Mining: Definition

Given:

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

Task:

Find all subsequences with support ≥ minsup



Sequential Pattern Mining: Challenge

- Given a sequence: <{a b} {c d e} {f} {g h i}>
 - Examples of subsequences:
 <{a} {c d} {f} {g} >, < {c d e} >, < {b} {g} >, etc.
- How many k-subsequences can be extracted from a given n-sequence?

$$\{a b\} \{c d e\} \{f\} \{g h i\} > n = 9$$

105

Sequential Pattern Mining: Example

Object	Timestamp	Events
Α	1	1,2,4
Α	2	2,3
Α	3	5
В	1	1,2
В	2	2,3,4
С	1	1, 2
С	2	2,3,4
С	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
Е	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%

Extracting Sequential Patterns

- Given n events: i_1 , i_2 , i_3 , ..., i_n
- Candidate 1-subsequences:

$$\langle \{i_1\} \rangle, \langle \{i_2\} \rangle, \langle \{i_3\} \rangle, ..., \langle \{i_n\} \rangle$$

Candidate 2-subsequences:

$$\{i_1, i_2\}$$
, $\{i_1, i_3\}$, ..., $\{i_1\} \{i_1\}$, $\{i_1\}$, $\{i_2\}$, ..., $\{i_{n-1}\} \{i_n\}$

Candidate 3-subsequences:

$$\langle \{i_1, i_2, i_3\} \rangle, \langle \{i_1, i_2, i_4\} \rangle, ..., \langle \{i_1, i_2\} \{i_1\} \rangle, \langle \{i_1, i_2\} \{i_2\} \rangle, ..., \langle \{i_1\} \{i_1, i_2\} \rangle, \langle \{i_1\} \{i_1\} \{i_1\} \{i_2\} \rangle, ..., \langle \{i_1\} \{i_1\} \{i_2\} \rangle, ..., \langle \{i_1\} \{i_2\} \rangle, ..., \langle \{i_1\} \{i_2\} \rangle, ..., \langle \{i_2\} \{i_2\} \}, ..., \langle \{i_2\} \{i_2\} \}, ..., \langle \{i_2\} \{i_2\} \rangle, ..., \langle \{i_2\} \{i_2\} \rangle, ..., \langle \{i_2\} \{i_2\} \rangle,$$

Generalized Sequential Pattern (GSP)

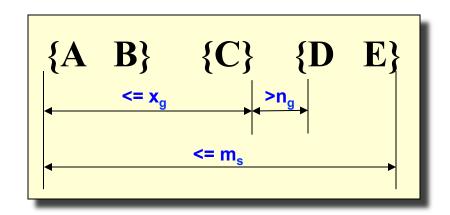
- Step 1:
 - Make the first pass over the sequence database D to yield all the 1element frequent sequences
- Step 2:

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:
 - \checkmark 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:
 - \checkmark Eliminate candidate k-sequences whose actual support is less than minsup

108

Timing Constraints (I)



x_g: max-gap

n_g: min-gap

m_s: maximum span

$x_g = 2$, $n_g = 0$, $m_g = 4$ Data sequence	Subsequence	Contain?
< {2,4} {3,5,6} {4,7} {4,5} {8} >	< {6} {5} >	Yes
< {1} {2} {3} {4} {5}>	< {1} {4} >	No
< {1} {2,3} {3,4} {4,5}>	< {2} {3} {5} >	Yes
< {1,2} {3} {2,3} {3,4} {2,4} {4,5}>	< {1,2} {5} >	No

References - Association rules

- R. Agrawal, T. Imielinski, and A. Swami. Mining association rules between sets of items in large databases. SIGMOD'93, 207-216, Washington, D.C.
- R. Agrawal and R. Srikant. Fast algorithms for mining association rules. VLDB'94 487-499, Santiago, Chile.
- R. Agrawal and R. Srikant. Mining sequential patterns. ICDE'95, 3-14, Taipei, Taiwan.
- R. J. Bayardo. Efficiently mining long patterns from databases. SIGMOD'98, 85-93, Seattle, Washington.
- 5. Brin, R. Motwani, and C. Silverstein. Beyond market basket: Generalizing association rules to correlations. SIGMOD'97, 265-276, Tucson, Arizona..
- D.W. Cheung, J. Han, V. Ng, and C.Y. Wong. Maintenance of discovered association rules in large databases: An incremental updating technique. ICDE'96, 106-114, New Orleans, LA..
- T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama. Data mining using two-dimensional optimized association rules: Scheme, algorithms, and visualization. SIGMOD'96, 13-23, Montreal, Canada.
- E.-H. Han, G. Karypis, and V. Kumar. Scalable parallel data mining for association rules. SIGMOD'97, 277-288, Tucson, Arizona.
- J. Han and Y. Fu. Discovery of multiple-level association rules from large databases. VLDB'95, 420-431, Zurich, Switzerland.
- M. Kamber, J. Han, and J. Y. Chiang. Metarule-guided mining of multi-dimensional association rules using data cubes. KDD'97, 207-210, Newport Beach, California.
- M. Klemettinen, H. Mannila, P. Ronkainen, H. Toivonen, and A.I. Verkamo. Finding interesting rules from large sets of discovered association rules. CIKM'94, 401-408, Gaithersburg, Maryland.
- R. Ng, L. V. S. Lakshmanan, J. Han, and A. Pang. Exploratory mining and pruning optimizations of constrained associations rules. SIGMOD'98, 13-24, Seattle, Washington.
- B. Ozden, S. Ramaswamy, and A. Silberschatz. Cyclic association rules. ICDE'98, 412-421, Orlando, FL.
- J.S. Park, M.S. Chen, and P.S. Yu. An effective hash-based algorithm for mining association rules. SIGMOD'95, 175-186, San Jose, CA.
- S. Ramaswamy, S. Mahajan, and A. Silberschatz. On the discovery of interesting patterns in association rules. VLDB'98, 368-379, New York, NY.
- S. Sarawagi, S. Thomas, and R. Agrawal. Integrating association rule mining with relational database systems: Alternatives and implications. SIGMOD'98, 343-354, Seattle, WA.

References - Association rules

- A. Savasere, E. Omiecinski, and S. Navathe. An efficient algorithm for mining association rules in large databases. VLDB'95, 432-443, Zurich, Switzerland.
- C. Silverstein, S. Brin, R. Motwani, and J. Ullman. Scalable techniques for mining causal structures.
 VLDB'98, 594-605, New York, NY.
- R. Srikant and R. Agrawal. Mining generalized association rules. VLDB'95, 407-419, Zurich, Switzerland.
- R. Srikant and R. Agrawal. Mining quantitative association rules in large relational tables. SIGMOD'96, 1-12, Montreal, Canada.
- R. Srikant, Q. Vu, and R. Agrawal. Mining association rules with item constraints. KDD'97, 67-73, Newport Beach, California.
- D. Tsur, J. D. Ullman, S. Abitboul, C. Clifton, R. Motwani, and S. Nestorov. Query flocks: A generalization of association-rule mining. SIGMOD'98, 1-12, Seattle, Washington.
- B. Ozden, S. Ramaswamy, and A. Silberschatz. Cyclic association rules. ICDE'98, 412-421, Orlando, FL.
- R.J. Miller and Y. Yang. Association rules over interval data. SIGMOD'97, 452-461, Tucson, Arizona.
- J. Han, G. Dong, and Y. Yin. Efficient mining of partial periodic patterns in time series database. ICDE'99, Sydney, Australia.
- F. Giannotti, G. Manco, D. Pedreschi and F. Turini. Experiences with a logic-based knowledge discovery support environment. In Proc. 1999 ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery (SIGMOD'99 DMKD). Philadelphia, May 1999.
- F. Giannotti, M. Nanni, G. Manco, D. Pedreschi and F. Turini. Integration of Deduction and Induction for Mining Supermarket Sales Data. In Proc. PADD'99, Practical Application of Data Discovery, Int. Conference, London, April 1999.

