

Data Mining

Knowledge Discovery in Databases

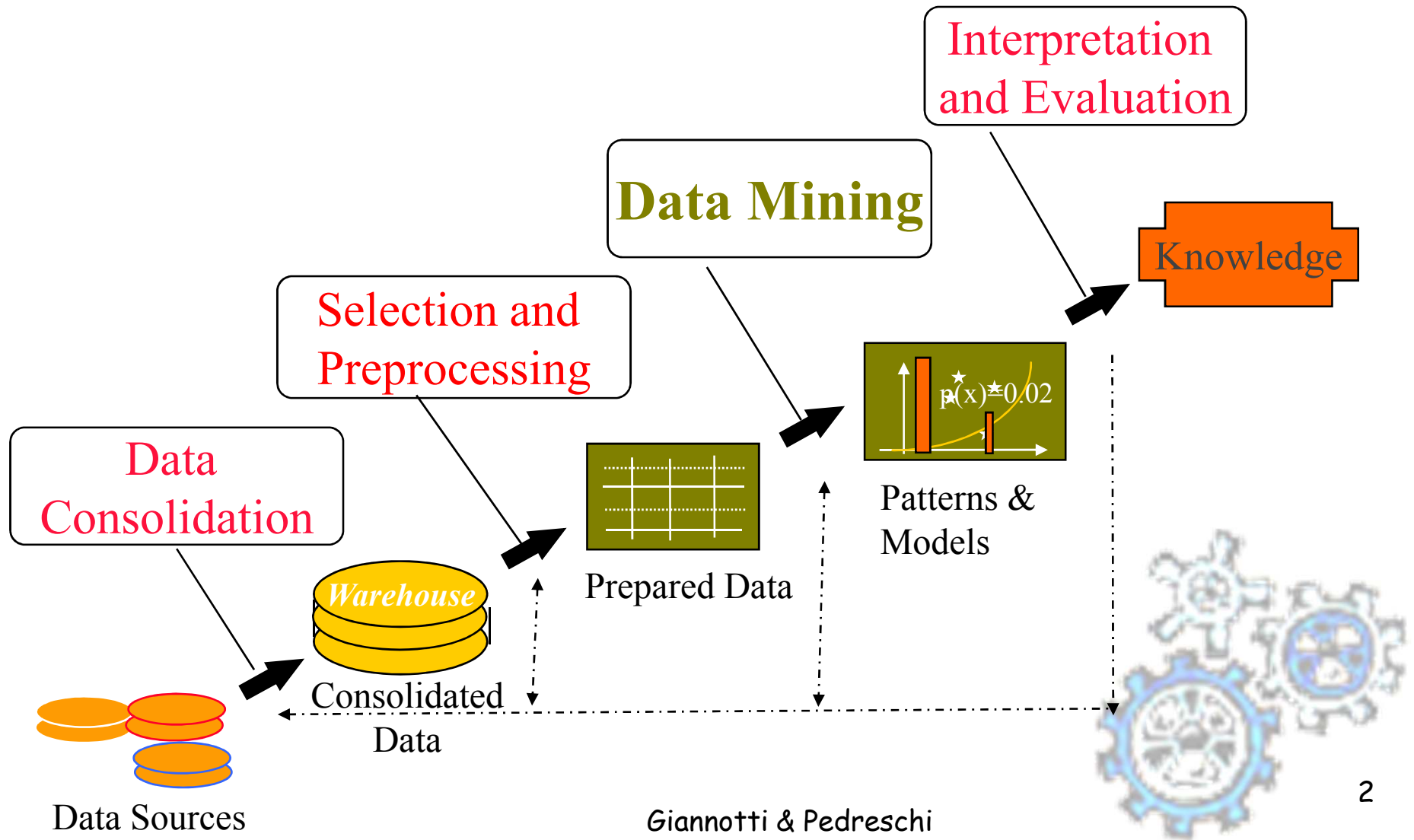
Fosca Giannotti and Dino Pedreschi
Pisa KDD Lab, ISTI-CNR & Univ. Pisa

<http://www-kdd.cnuce.cnr.it/>



MAINS – Master in Management dell’Innovazione
Scuola Superiore S. Anna

KDD Process



Association rules and market basket analysis



Association rules - module outline

1. What are association rules (AR) and what are they used for:

1. The paradigmatic application: Market Basket Analysis
2. The single dimensional AR (intra-attribute)

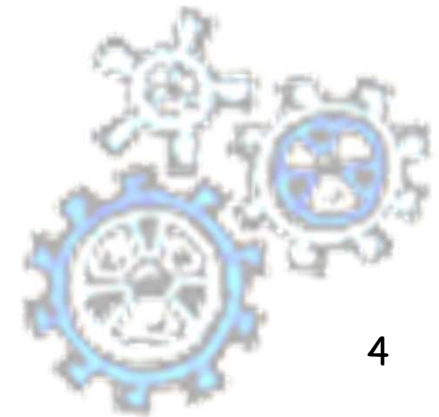


2. How to compute AR

1. Basic Apriori Algorithm and its optimizations
2. Multi-Dimension AR (inter-attribute)
3. Quantitative AR
4. Constrained AR

3. How to reason on AR and how to evaluate their quality

1. Multiple-level AR
2. Interestingness
3. Correlation vs. Association



Market Basket Analysis: the context

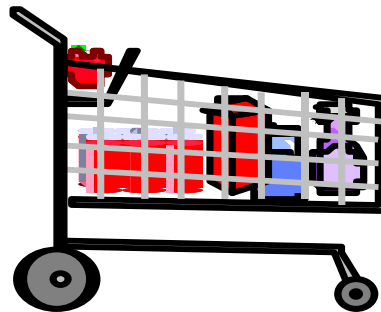
Customer buying habits by finding associations and correlations between the different items that customers place in their "shopping basket"

Milk, eggs, sugar,
bread



Customer1

Milk, eggs, cereal, bread



Customer2

Eggs, sugar

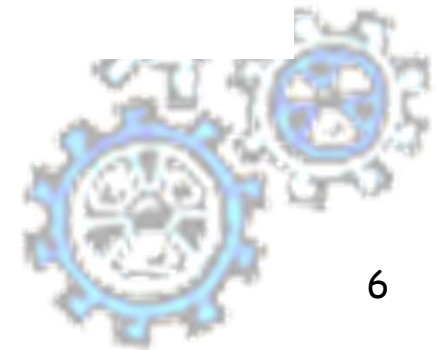
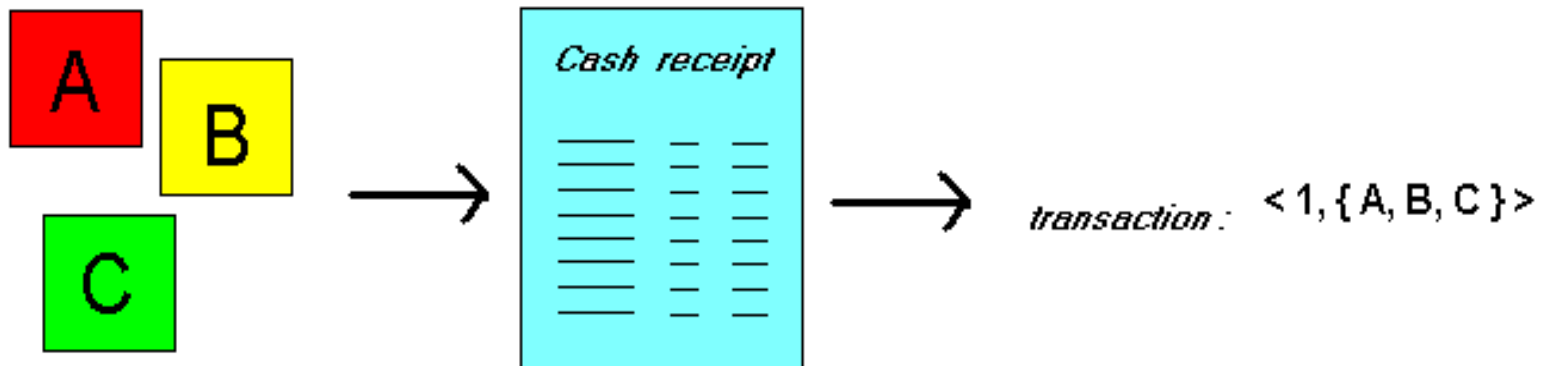


Customer3

Market Basket Analysis: the context

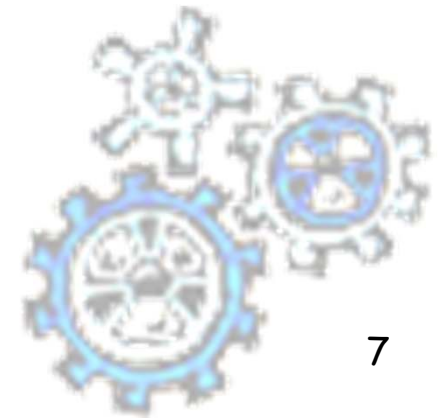
Given: a database of customer **transactions**, where each transaction is a **set of items**

- 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

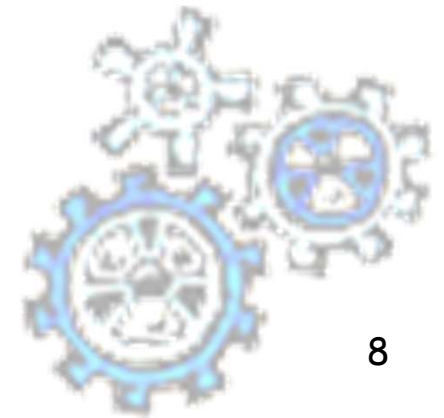
Telecommunication:

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

Atmospheric phenomena:

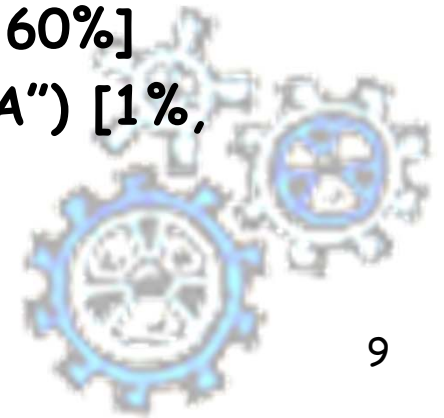
Each time interval (e.g. a day) is a transaction containing the set of observed event (rains, wind, etc.)

Etc.



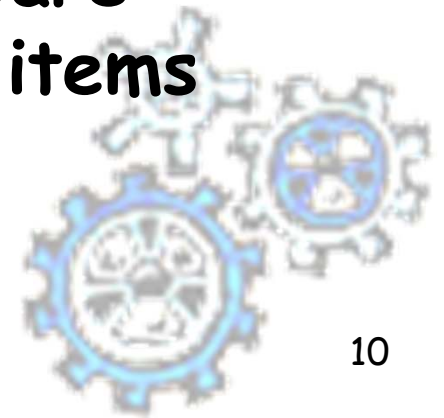
Association Rules

- 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 \rightarrow Head [support, confidence]”.
 - $\text{buys}(x, \text{"diapers"}) \rightarrow \text{buys}(x, \text{"beers"})$ [0.5%, 60%]
 - $\text{major}(x, \text{"CS"}) \wedge \text{takes}(x, \text{"DB"}) \rightarrow \text{grade}(x, \text{"A"})$ [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 hardware store opens, one of the most sold items is toilet rings.”



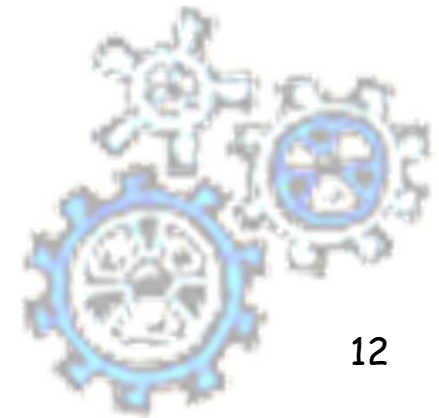
Association Rules Road Map

- **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.



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



Data Mining

Association Analysis: Basic Concepts and Algorithms

Lecture Notes for Chapter 6

Introduction to Data Mining

by

Tan, Steinbach, Kumar

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

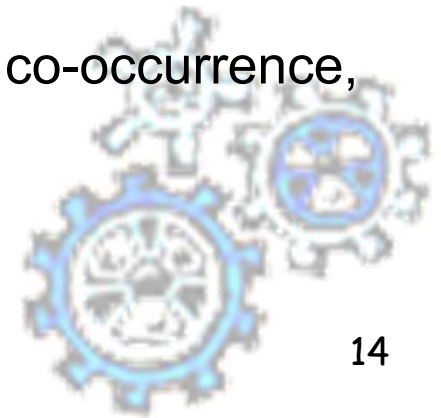
Market-Basket transactions

<i>TID</i>	<i>Items</i>
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

$\{\text{Diaper}\} \rightarrow \{\text{Beer}\},$
 $\{\text{Milk, Bread}\} \rightarrow \{\text{Eggs, Coke}\},$
 $\{\text{Beer, Bread}\} \rightarrow \{\text{Milk}\},$

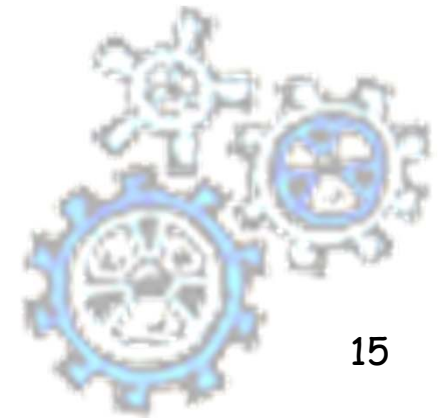
Implication means co-occurrence,
not causality!



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. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$
- **Support**
 - **Fraction of transactions that contain an itemset**
 - E.g. $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$
- **Frequent Itemset**
 - **An itemset whose support is greater than or equal to a *minsup* threshold**

<i>TID</i>	<i>Items</i>
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 $X \rightarrow Y$, where X and Y are itemsets
- Example:
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

<i>TID</i>	<i>Items</i>
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:

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

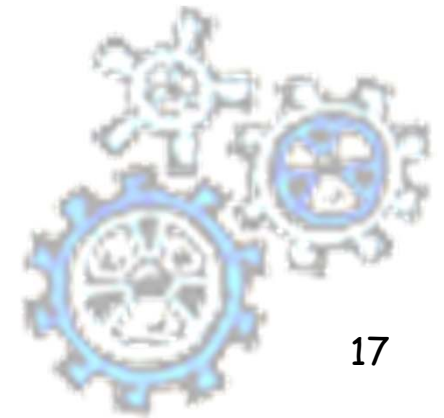
$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4$$

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

Association Rule Mining Task

- Given a set of transactions T , the goal of association rule mining is to find all rules having
 - support \geq *minsup* threshold
 - confidence \geq *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

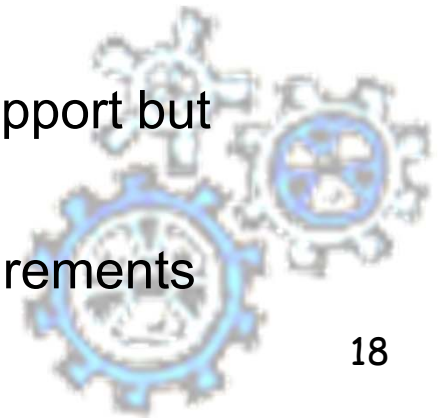
<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example of Rules:

$\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$ ($s=0.4, c=0.67$)
 $\{\text{Milk, Beer}\} \rightarrow \{\text{Diaper}\}$ ($s=0.4, c=1.0$)
 $\{\text{Diaper, Beer}\} \rightarrow \{\text{Milk}\}$ ($s=0.4, c=0.67$)
 $\{\text{Beer}\} \rightarrow \{\text{Milk, Diaper}\}$ ($s=0.4, c=0.67$)
 $\{\text{Diaper}\} \rightarrow \{\text{Milk, Beer}\}$ ($s=0.4, c=0.5$)
 $\{\text{Milk}\} \rightarrow \{\text{Diaper, Beer}\}$ ($s=0.4, c=0.5$)

Observations:

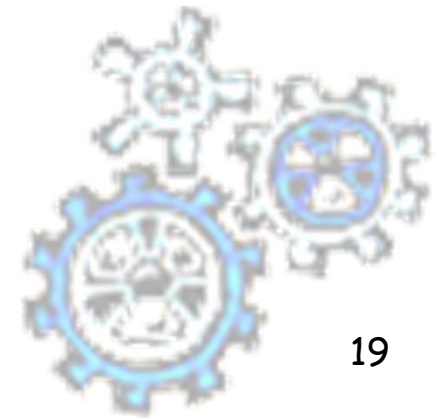
- All the above rules are binary partitions of the same itemset:
 $\{\text{Milk, Diaper, Beer}\}$
- Rules originating from the same itemset have identical support but can have different confidence
- Thus, we may decouple the support and confidence requirements



Mining Association Rules

- **Two-step approach:**
 1. **Frequent Itemset Generation**
 - Generate all itemsets whose support \geq minsup
 2. **Rule Generation**
 - Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

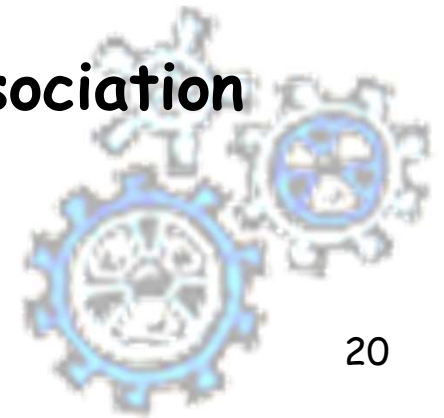
- **Frequent itemset generation is still computationally expensive**



Basic Apriori Algorithm

Problem Decomposition

- 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)
- Use the frequent itemsets to generate association rules.

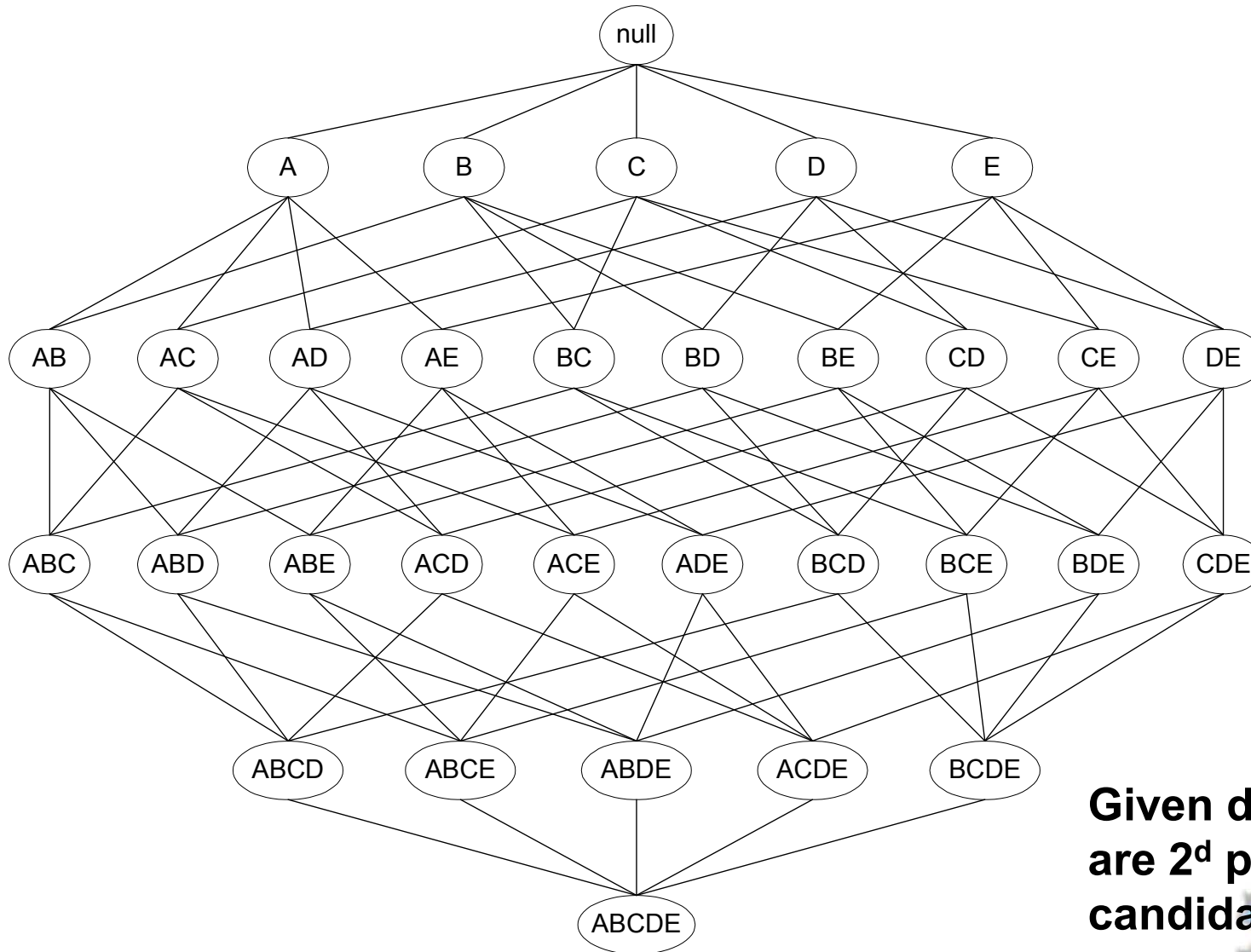


Frequent Itemset Mining Problem

- $I = \{x_1, \dots, x_n\}$ set of distinct literals (called **items**)
- $X \subseteq I, X \neq \emptyset, |X| = k, X$ is called **k-itemset**
- A **transaction** is a couple $\langle tID, X \rangle$ where X is an itemset
- A **transaction database** TDB is a set of transactions
- An itemset X is **contained** in a transaction $\langle tID, Y \rangle$ if $X \subseteq Y$
- Given a TDB the subset of transactions of TDB in which X is contained is named $TDB[X]$.
- The **support** of an itemset X , written $supp_{TDB}(X)$ is the cardinality of $TDB[X]$.
- Given a user-defined **min_sup** threshold an itemset X is **frequent** in TDB if its support is no less than min_sup .

- Given a user-defined min_sup and a transaction database TDB , the **Frequent Itemset Mining Problem** requires to compute all frequent itemsets in TDB w.r.t min_sup .

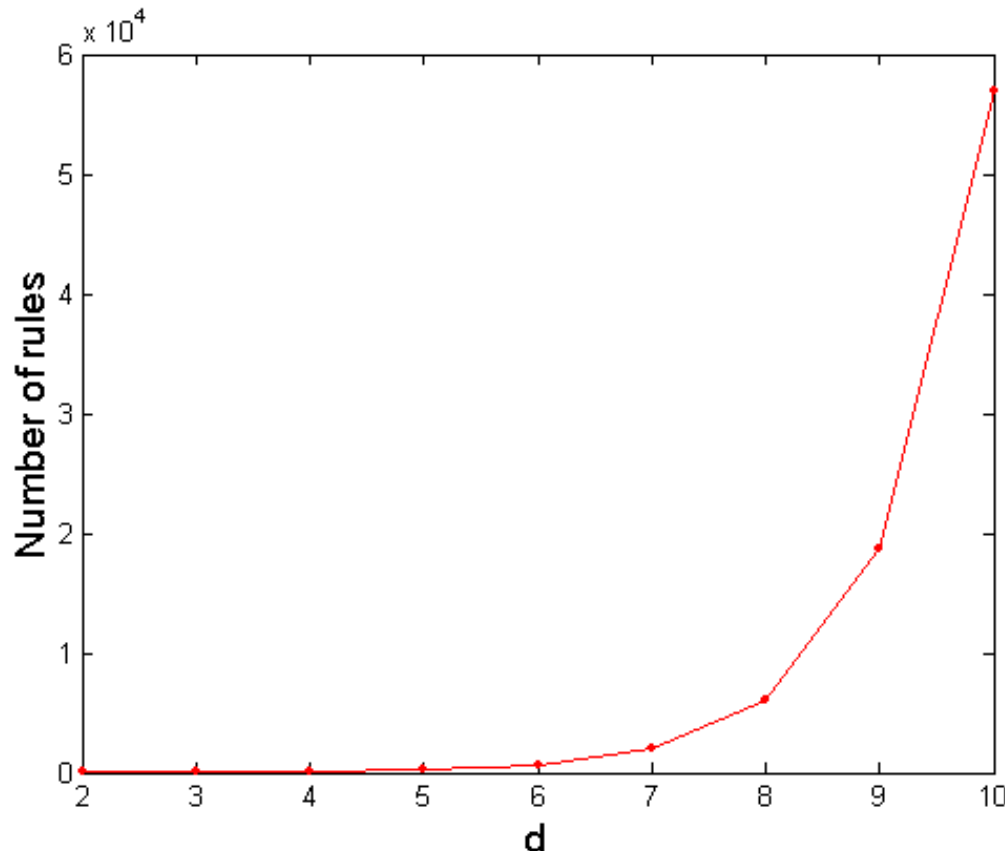
Frequent Itemset Generation



Given d items, there are 2^d possible candidate itemsets

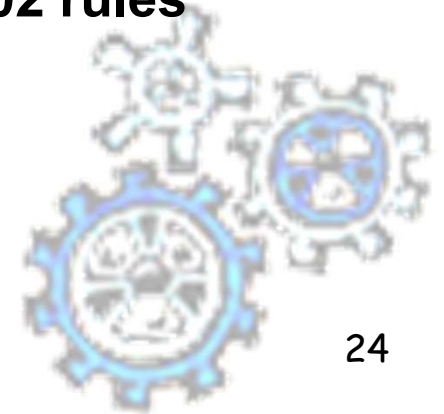
Computational Complexity

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



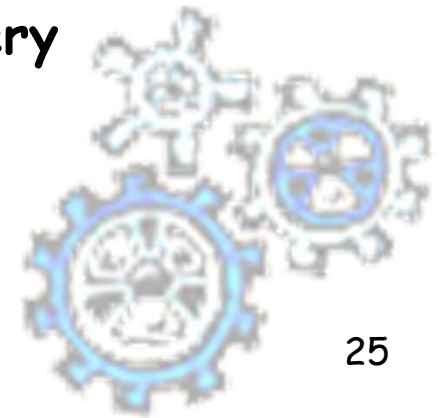
$$R = \sum_{k=1}^{d-1} \left[\binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right]$$
$$= 3^d - 2^{d+1} + 1$$

If $d=6$, $R = 602$ rules



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

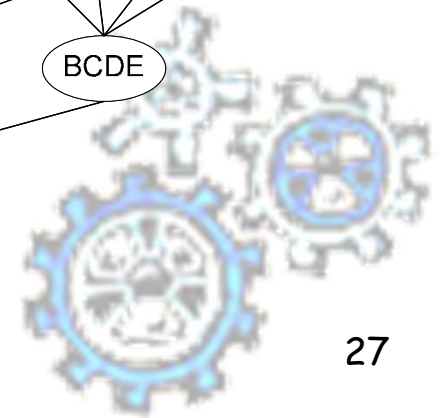
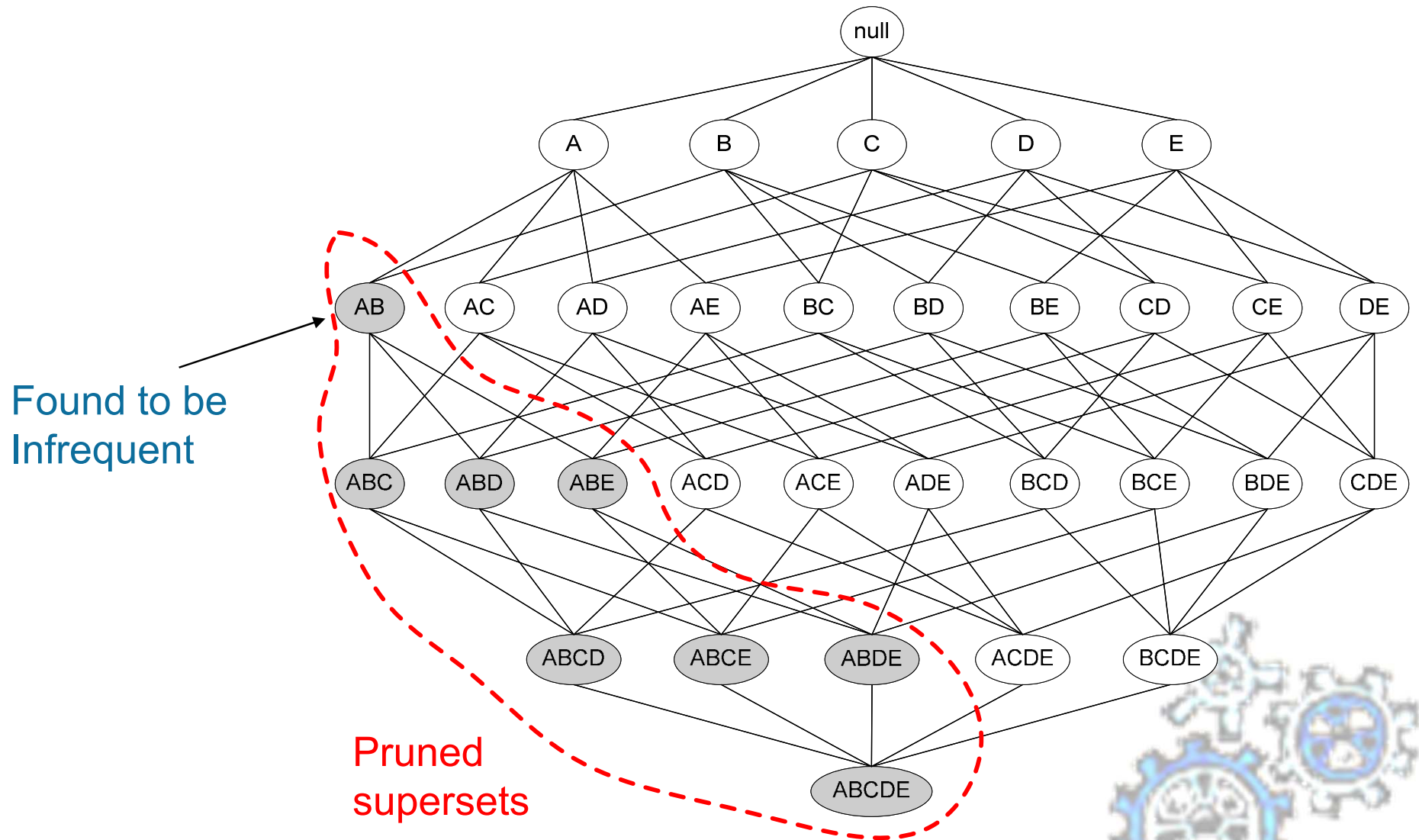
- Apriori principle holds due to the following property of the support measure:

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

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



Illustrating Apriori Principle



Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

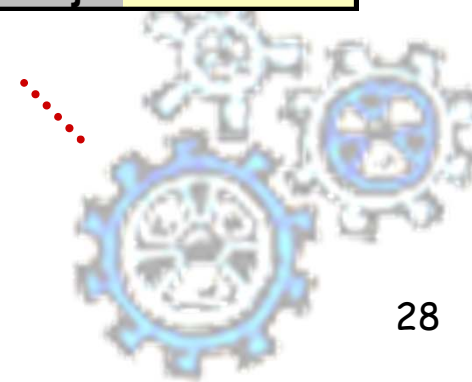
Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3 = 41$
 With support-based pruning,
 $6 + 6 + 1 = 13$

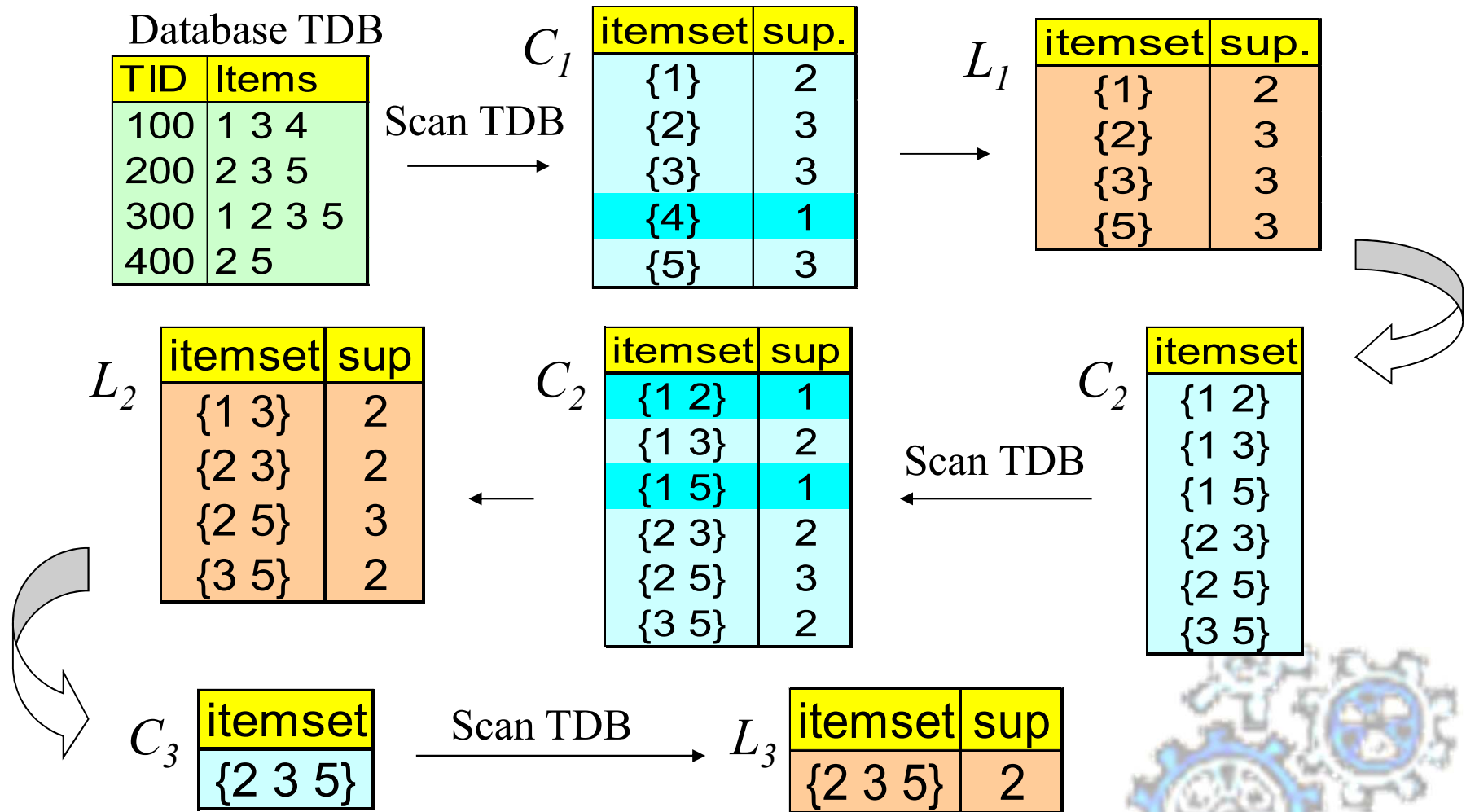


Triplets (3-itemsets)

Itemset	Count
{Bread,Milk,Diaper}	3



Apriori Execution Example ($min_sup = 2$)



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

$L_1 = \{\text{frequent items}\};$

for ($k = 1; L_k \neq \emptyset; k++$) **do begin**

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

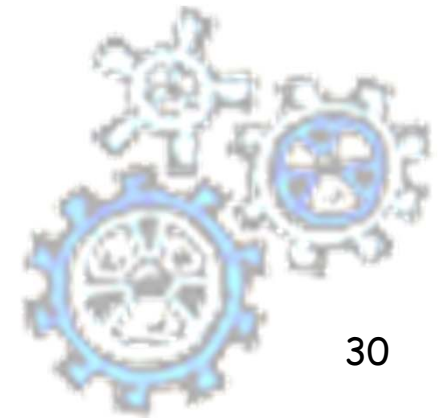
for each transaction t in database **do**

increment the count of all candidates in C_{k+1}
that are contained in t

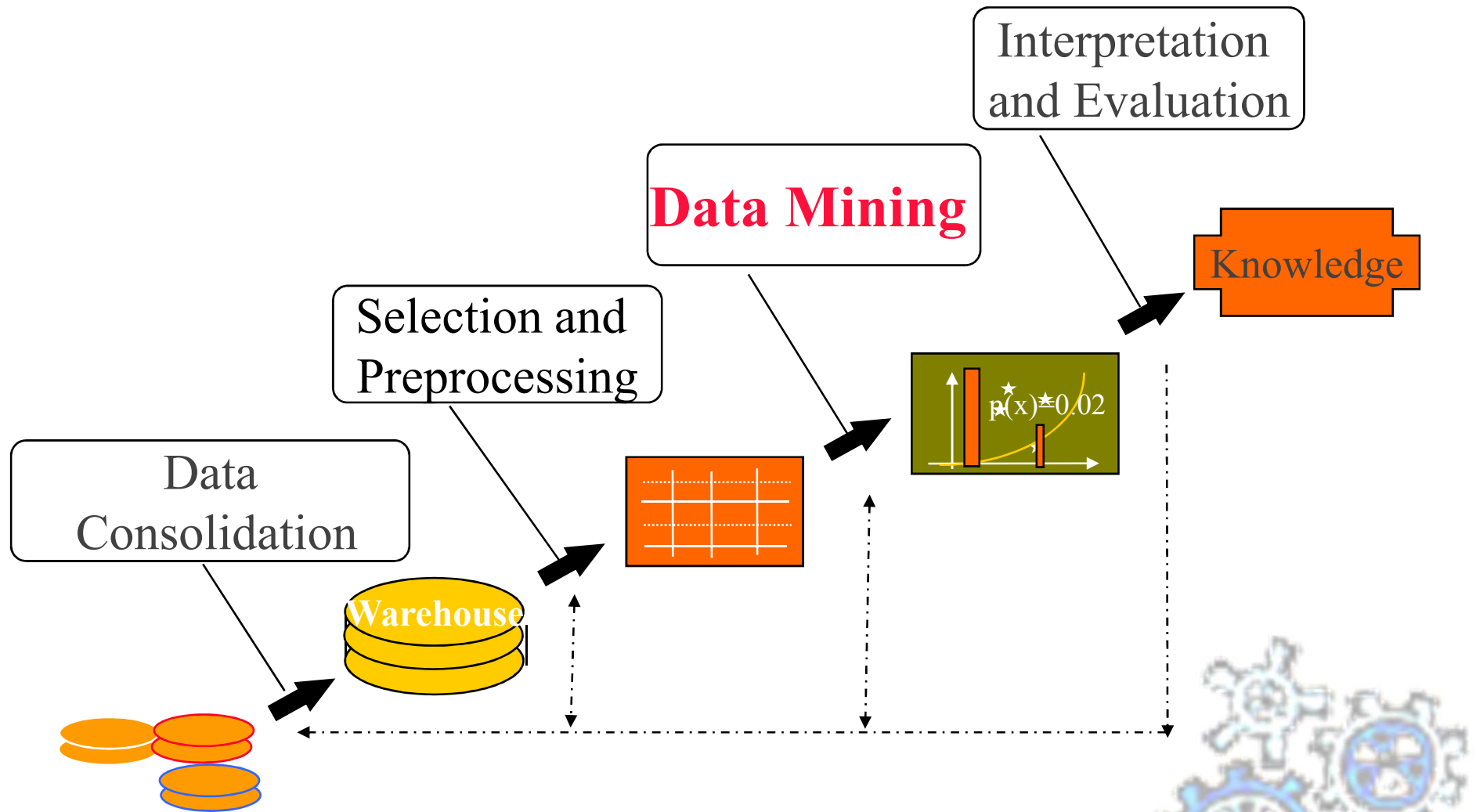
$L_{k+1} =$ candidates in C_{k+1} with `min_support`

end

return $\cup_k L_k$;



The KDD process

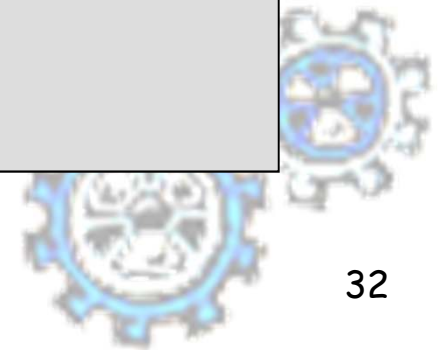


Generating Association Rules from Frequent Itemsets

- Only strong association rules are generated
- Frequent itemsets satisfy minimum support threshold
- Strong rules are those that satisfy minimum confidence threshold

$$\frac{\text{support}(A \cup B)}{\text{support}(A)}$$

- **C** For each frequent itemset, **f**, generate all non-empty subsets of **f**
 - **For every** non-empty subset **s** of **f** **do**
 - **if** $\text{support}(f)/\text{support}(s) \geq \text{min_confidence}$ **then**
 - output rule $s \implies (f-s)$
 - **end**



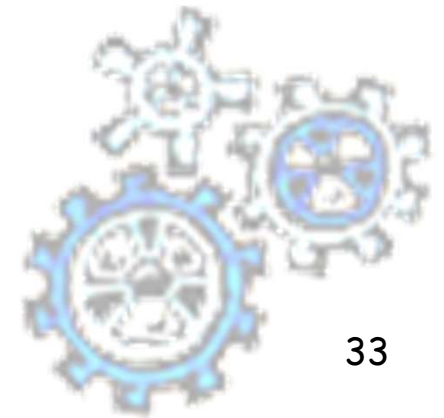
Rule Generation

- Given a frequent itemset L , find all non-empty subsets $f \subset L$ such that $f \rightarrow L - f$ satisfies the minimum confidence requirement

- If $\{A, B, C, D\}$ is a frequent itemset, candidate rules:

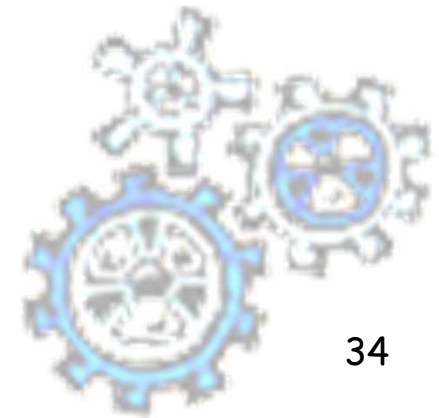
$ABC \rightarrow D,$	$ABD \rightarrow C,$	$ACD \rightarrow B,$	$BCD \rightarrow A,$
$A \rightarrow BCD,$	$B \rightarrow ACD,$	$C \rightarrow ABD,$	$D \rightarrow ABC$
$AB \rightarrow CD,$	$AC \rightarrow BD,$	$AD \rightarrow BC,$	$BC \rightarrow AD,$
$BD \rightarrow AC,$	$CD \rightarrow AB,$		

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



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 :

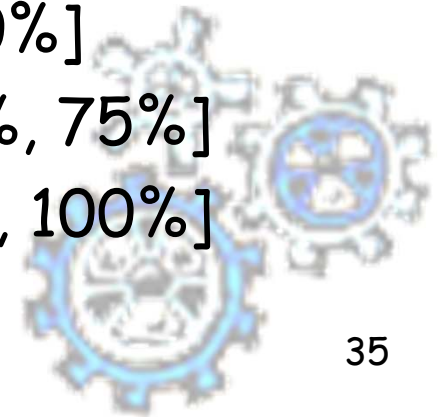
CID	nationality	age	income
1	Italian	50	low
2	French	40	high
3	French	30	high
4	Italian	50	medium
5	Italian	45	high
6	French	35	high

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

Single-dimensional (Intra-attribute)

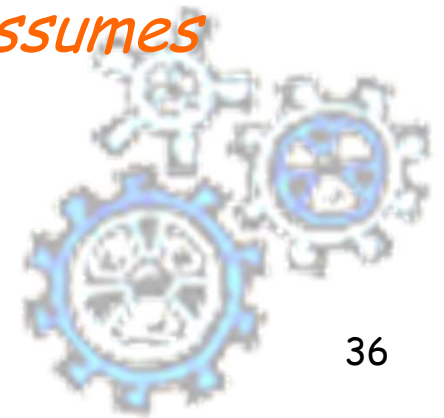
The events are: *items A, B and C belong to the same transaction*

Occurrence of events: *transactions*

Multi-dimensional (Inter-attribute)

The events are : *attribute A assumes value a, attribute B assumes value b and attribute C assumes value c.*

Occurrence of events: *tuples*



Single-dimensional vs Multi-dimensional AR

Multi-dimensional

<1, Italian, 50, low>
<2, French, 45, high>



Single-dimensional

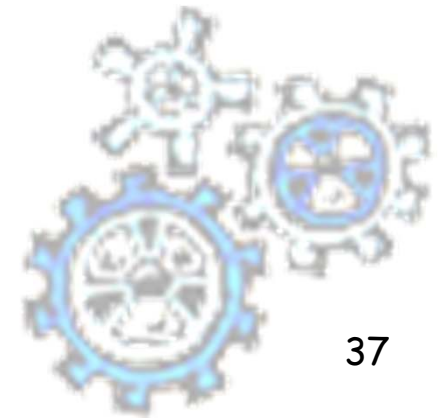
<1, {nat/Ita, age/50, inc/low}>
<2, {nat/Fre, age/45, inc/high}>

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

<1, yes, yes, no, no>
<2, yes, no, yes, no>



<1, {a, b}>
<2, {a, c}>



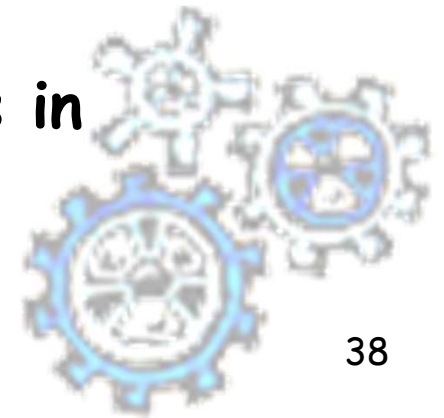
Quantitative Attributes

- 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	70,3	25,8
4	170	65,2	27,0

Problem: too many distinct values

Solution: transform quantitative attributes in categorical ones via **discretization**.



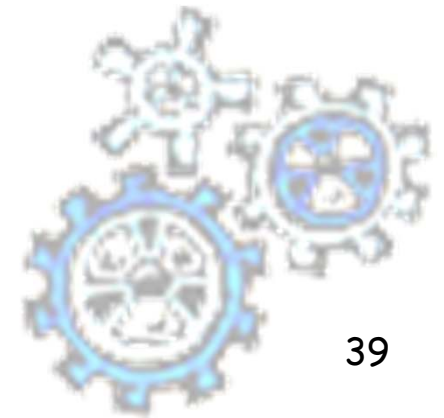
Quantitative Association Rules

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.

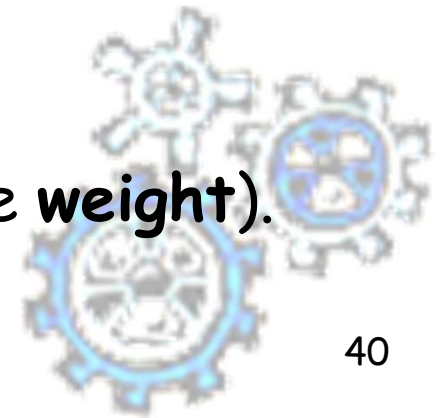
height: 0-150cm, 151-170cm, 171-180cm, >180cm

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

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

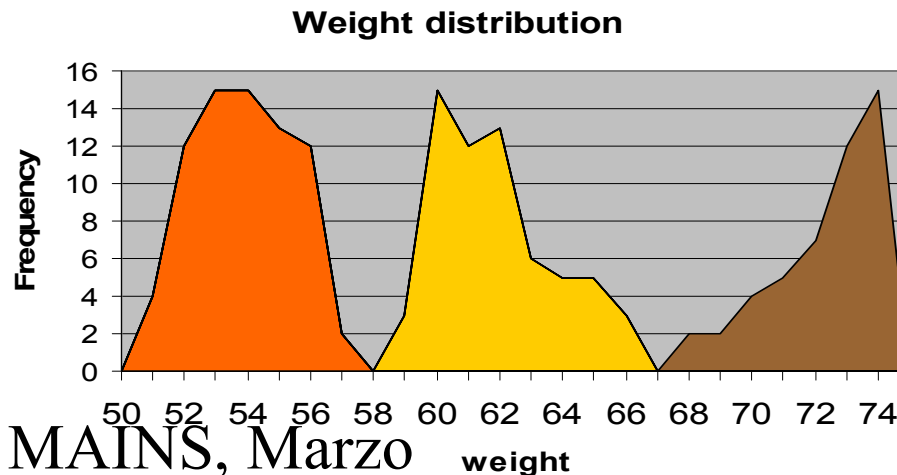
CID	height	weight	income
1	151-171	60-80	>30
2	171-180	60-80	20-25
3	171-180	60-80	25-30
4	151-170	60-80	25-30

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

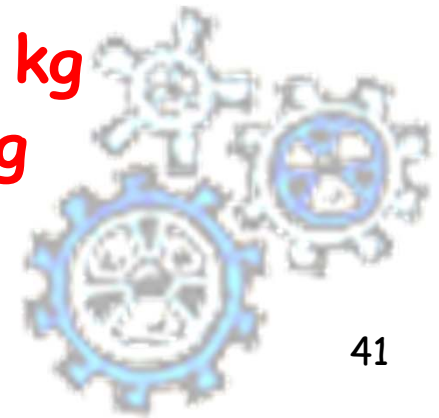


How to choose intervals?

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



50 - 58 kg
59-67 kg
> 68 kg



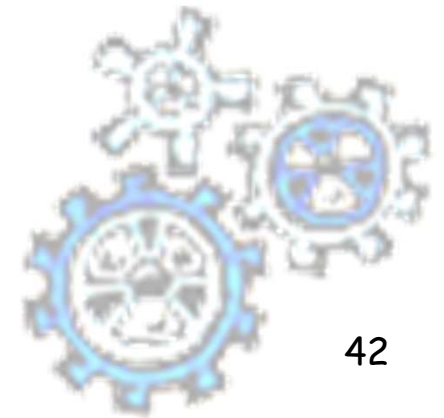
Discretization of quantitative attributes

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



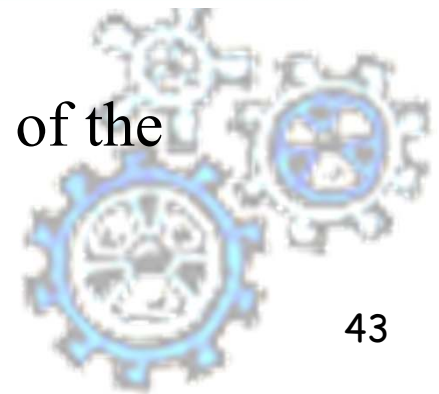
Quantitative Association Rules

RecordID	Age	Married	NumCars
100	23	No	1
200	25	Yes	1
300	29	No	0
400	34	Yes	2
500	38	Yes	2



Sample Rules	Support	Confidence
<age:30..39> and <married: yes> ==> <numCars:2>	40%	100%
<NumCars: 0..1> ==> <Married: No>	40%	66.70%

Handling quantitative rules may require mapping of the **continuous** variables into **Boolean**



Mapping Quantitative to Boolean

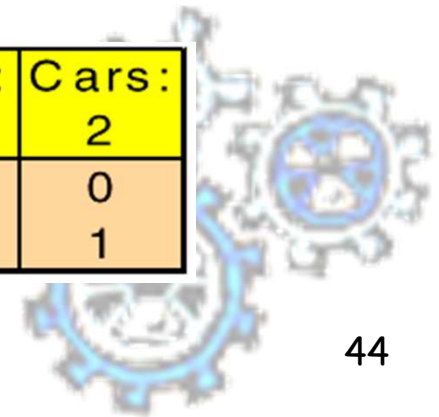
- One possible solution is to map the problem to the Boolean association rules:
 - discretize a non-categorical attribute to intervals, e.g., Age [20,29], [30,39],...
 - categorical attributes: each value becomes one item
 - non-categorical attributes: each interval becomes one item

- Problems with the mapping

- too few intervals: lost information
- too low support: too many rules

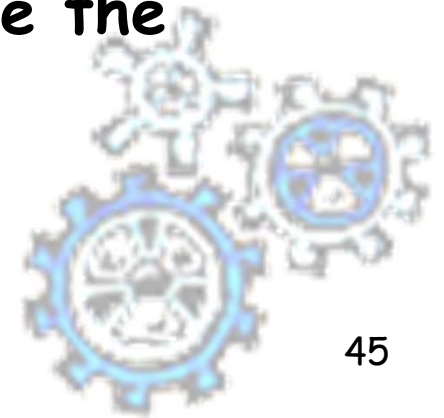
RecordID	Age	Married	NoCars
100	23	No	1
500	38	Yes	2

RecID	Age: 20..29	Age: 30..39	Married: Yes	Married: No	Cars: 0	Cars: 1	Cars: 2
100	1	0	0	1	0	1	0
500	0	1	1	0	0	0	1



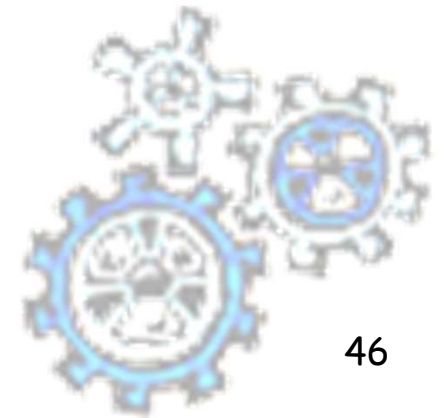
Constraints and AR

- **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.

✓ $P(x, y) \wedge Q(x, w) \rightarrow \text{takes}(x, \text{"database systems"})$.

■ Rule content constraint: constraint-based query optimization (Ng, et al., SIGMOD'98).

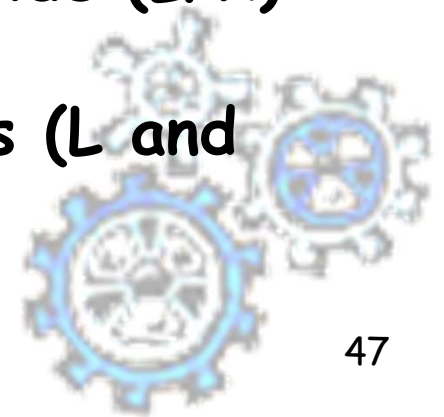
✓ $\text{sum(LHS)} < 100 \wedge \text{min(LHS)} > 20 \wedge \text{sum(RHS)} > 1000$

■ 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).

✓ $\text{sum(LHS)} < \text{min(RHS)} \wedge \text{max(RHS)} < 5 * \text{sum(LHS)}$



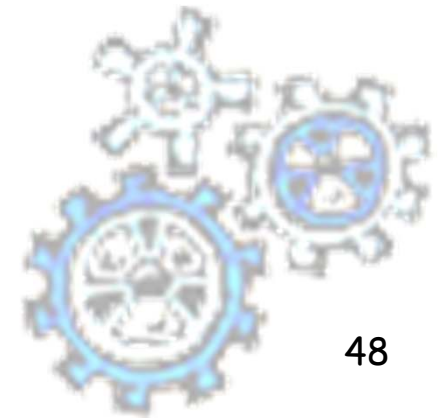
Mining Association Rules with Constraints

■ Postprocessing

- A naïve solution: apply Apriori for finding all frequent sets, and **then** to test them for constraint satisfaction one by one.

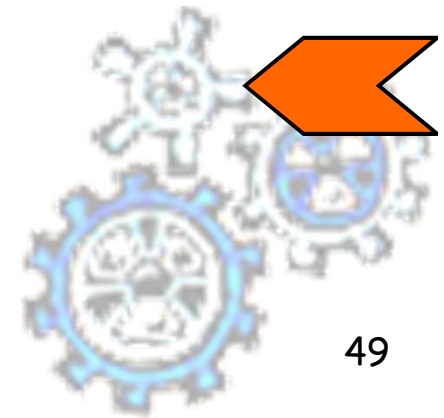
■ Optimization

- Han approach: comprehensive analysis of the properties of constraints and try to **push them as deeply as possible** inside the frequent set computation.



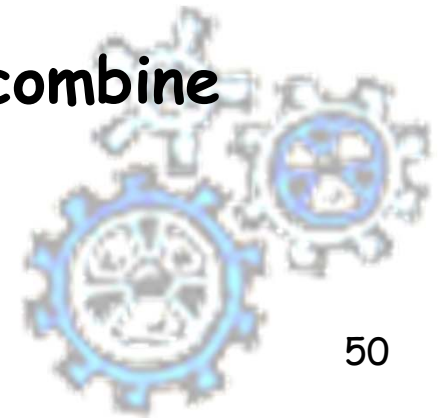
Association rules - module outline

- **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

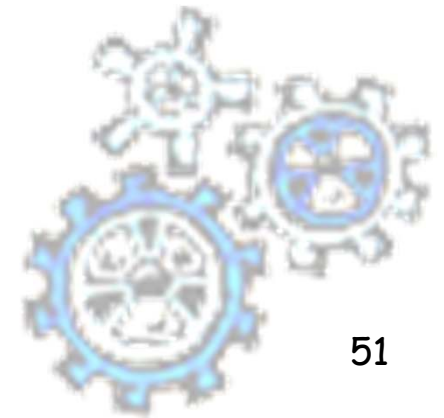
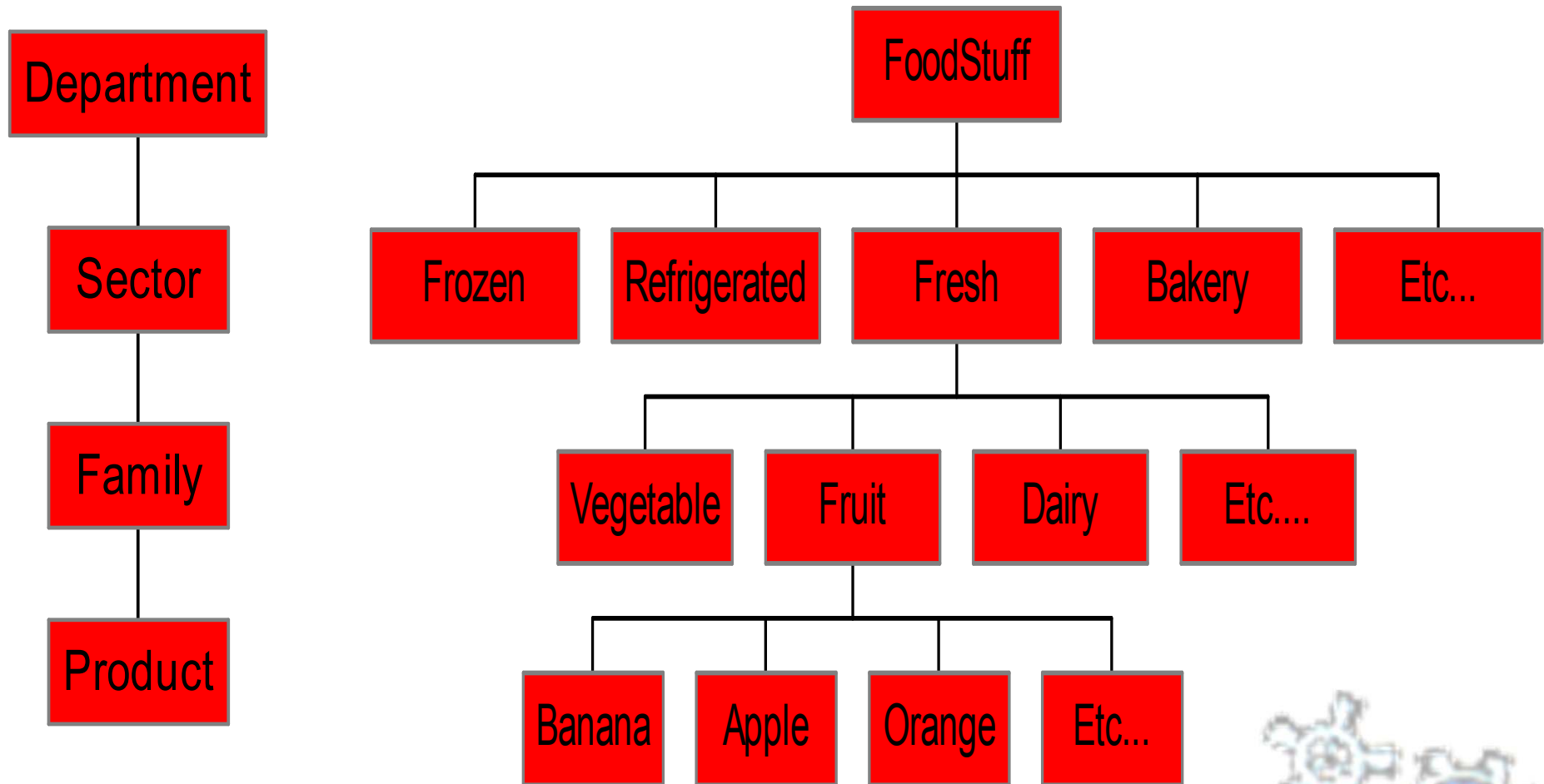


Multilevel AR

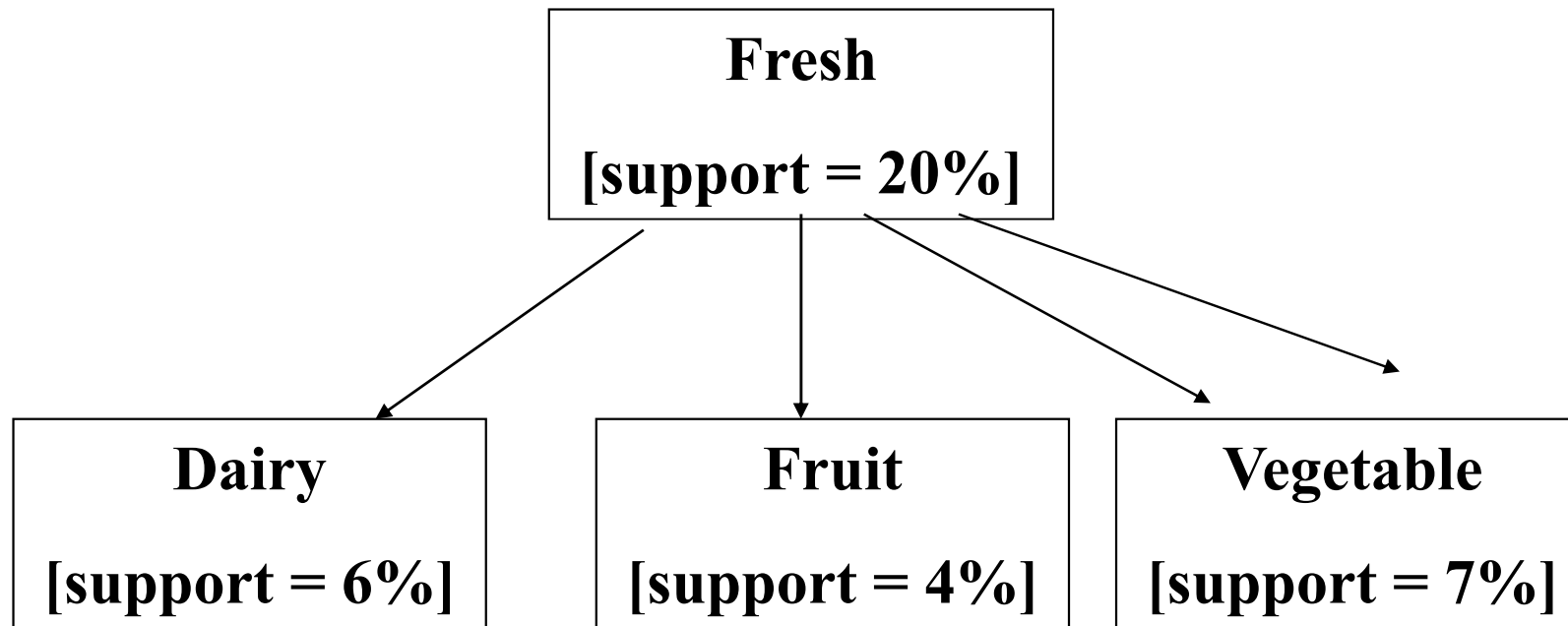
- Is difficult to find interesting patterns at a **too primitive level**
 - high support = too few rules
 - low support = too many rules, most uninteresting
- Approach: reason at suitable level of abstraction
- A common form of background knowledge is that an attribute may be generalized or specialized 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



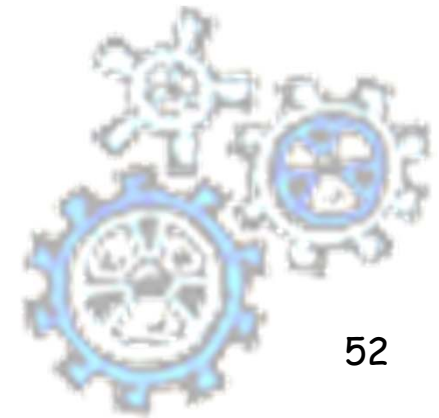
Multilevel AR



Fresh \Rightarrow Bakery [20%, 60%]

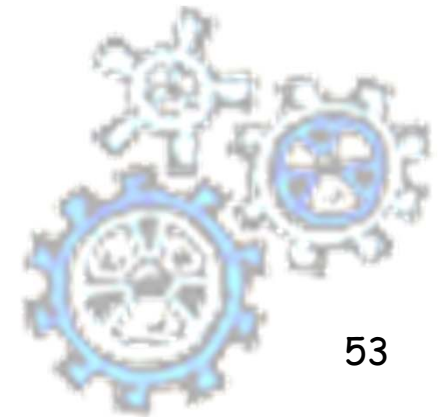
Dairy \Rightarrow Bread [6%, 50%]

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

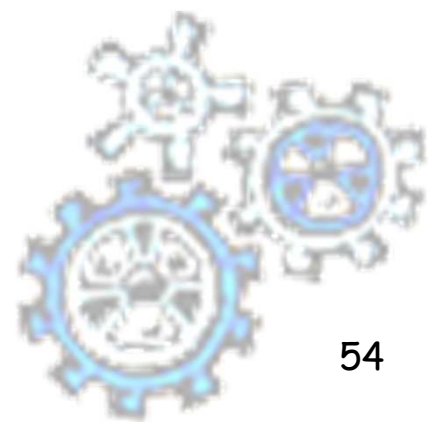
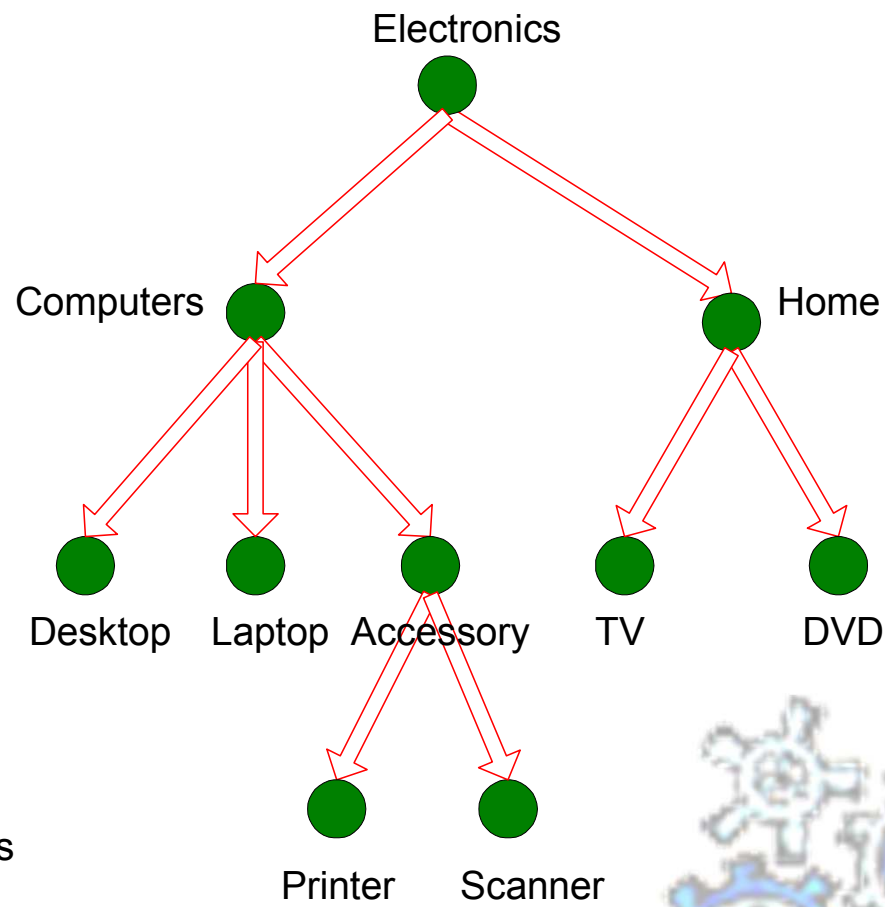
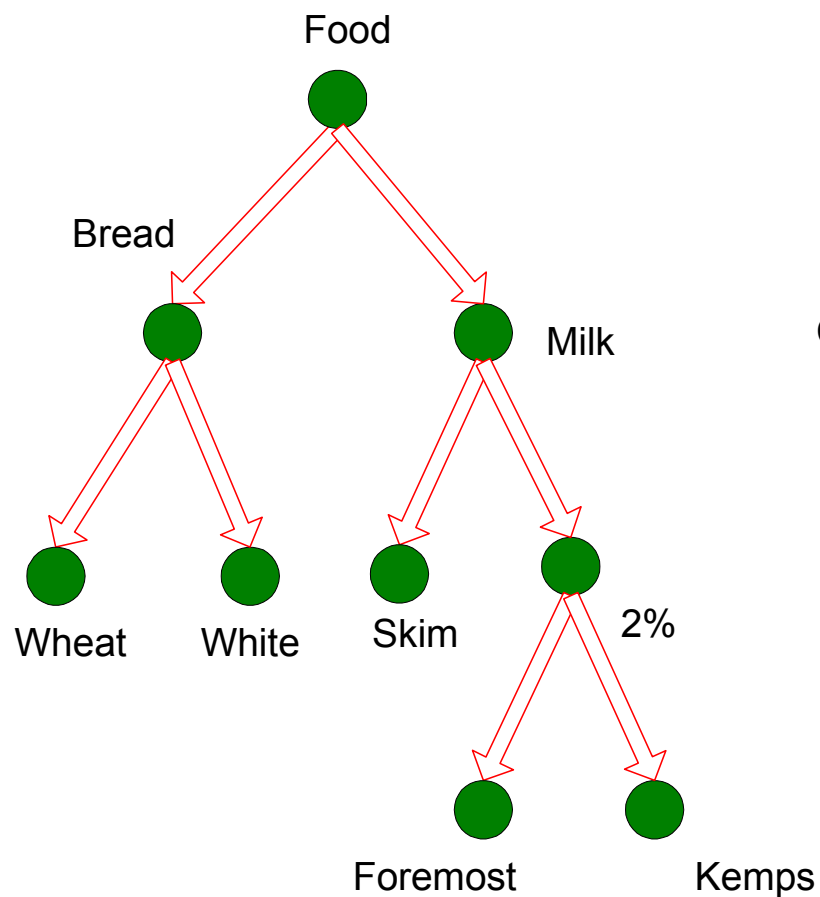


Support and Confidence of Multilevel AR

- **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**

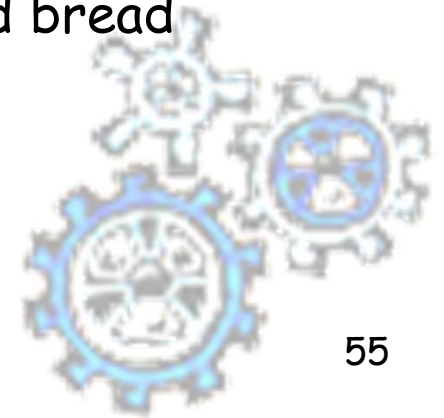


Multi-level Association Rules



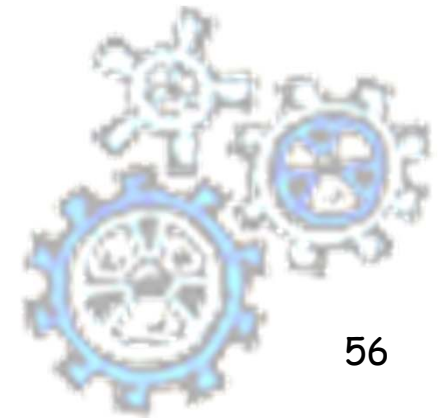
Multi-level Association Rules

- 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



Multi-level Association Rules

- How do support and confidence vary as we traverse the concept hierarchy?
 - If X is the parent item for both $X1$ and $X2$, then $\sigma(X) \leq \sigma(X1) + \sigma(X2)$
 - If $\sigma(X1 \cup Y1) \geq \text{minsup}$,
and X is parent of $X1$, Y is parent of $Y1$
then $\sigma(X \cup Y1) \geq \text{minsup}$, $\sigma(X1 \cup Y) \geq \text{minsup}$
 $\sigma(X \cup Y) \geq \text{minsup}$
 - If $\text{conf}(X1 \Rightarrow Y1) \geq \text{minconf}$,
then $\text{conf}(X1 \Rightarrow Y) \geq \text{minconf}$



Reasoning with Multilevel AR

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

Example: *Apple Melinda* \Rightarrow *Colgate Tooth-paste*

It is a curiosity not a behavior

- Too high level => uninteresting rules

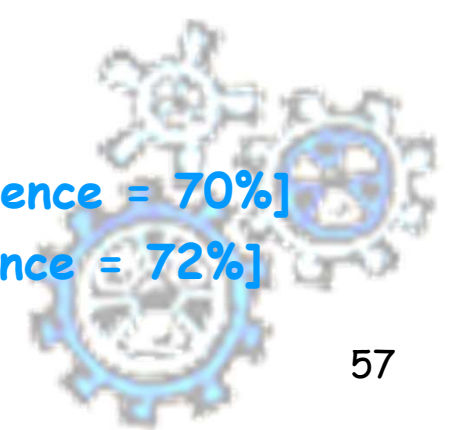
Example: *Foodstuff* \Rightarrow *Varia*

- 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.

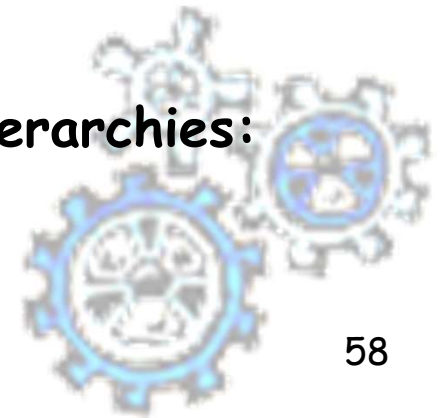
- Example (milk has 4 subclasses)

- *milk* \Rightarrow *wheat bread*, [support = 8%, confidence = 70%]
- *2%-milk* \Rightarrow *wheat bread*, [support = 2%, confidence = 72%]



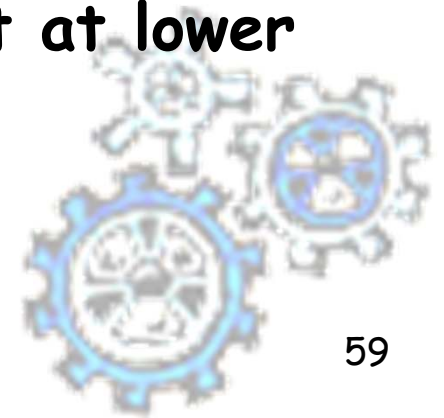
Mining Multilevel AR

- 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 → bakery [20%, 60%].
 - Then find their lower-level “weaker” rules:
fruit → 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*



Multi-level Association: Uniform Support vs. Reduced Support

- **Uniform Support: the same minimum support for all levels**
 - + One minimum support threshold. No need to examine itemsets containing any item whose ancestors do not have minimum support.
 - - If support threshold
 - too high \Rightarrow miss low level associations.
 - too low \Rightarrow generate too many high level associations.
- **Reduced Support: reduced minimum support at lower levels - different strategies possible**



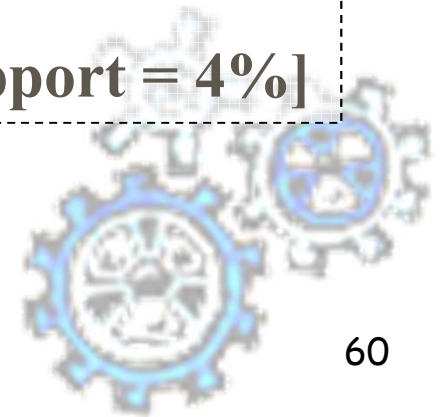
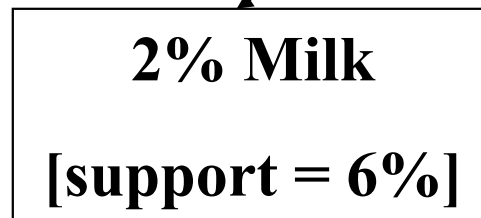
Uniform Support

Multi-level mining with uniform support

Level 1
min_sup = 5%



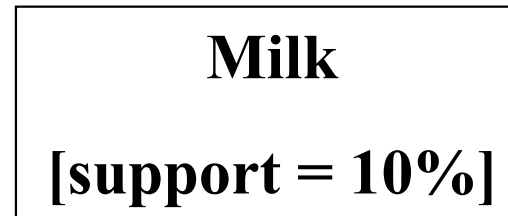
Level 2
min_sup = 5%



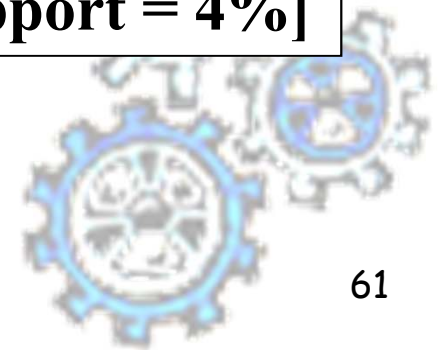
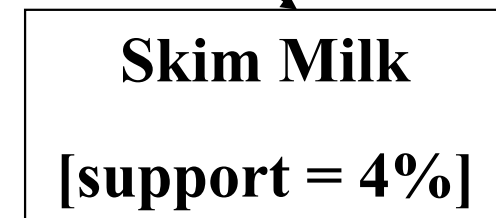
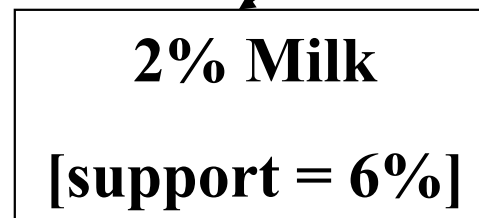
Reduced Support

Multi-level mining with reduced support

Level 1
min_sup = 5%



Level 2
min_sup = 3%

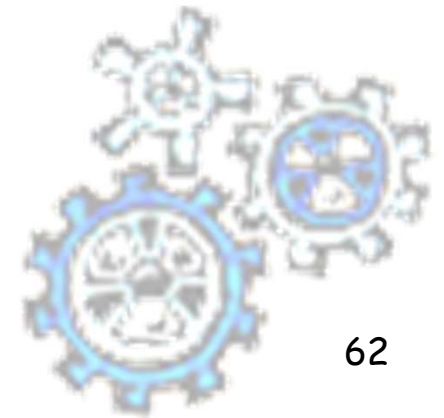


Reasoning with AR

■ Significance:

Example: $\langle 1, \{a, b\} \rangle$
 $\langle 2, \{a\} \rangle$
 $\langle 3, \{a, b, c\} \rangle$
 $\langle 4, \{b, d\} \rangle$

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



Beyond Support and Confidence

■ Example 1: (Aggarwal & Yu, PODS98)

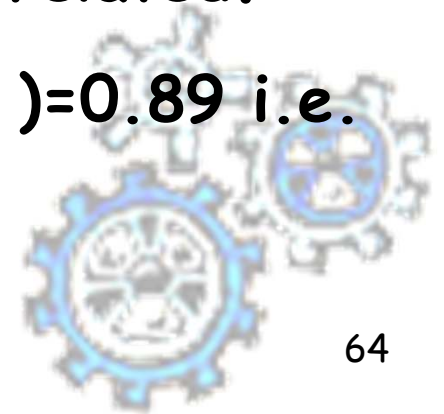
	coffee	not coffee	sum(row)
tea	20	5	25
not tea	70	5	75
sum(col.)	90	10	100

- $\{tea\} \Rightarrow \{coffee\}$ 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
 - $\{\sim tea\} \Rightarrow \{coffee\}$ has higher confidence(93%)



Correlation and Interest

- Two events are independent if $P(A \wedge B) = P(A) * P(B)$, otherwise are correlated.
- Interest = $P(A \wedge B) / P(B) * P(A)$
- 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, $I(\text{buy tea} \wedge \text{buy coffee}) = 0.89$ i.e. they are negatively correlated.



Computing Interestingness Measure

- Given a rule $X \rightarrow Y$, information needed to compute rule interestingness can be obtained from a contingency table

Contingency table for $X \rightarrow Y$

	Y	\bar{Y}	
X	f_{11}	f_{10}	f_{1+}
\bar{X}	f_{01}	f_{00}	f_{0+}
	f_{+1}	f_{+0}	$ T $

f_{11} : support of X and Y

f_{10} : support of X and \bar{Y}

f_{01} : support of \bar{X} and Y

f_{00} : support of \bar{X} and \bar{Y}

Used to define various measures

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

Statistical-based Measures

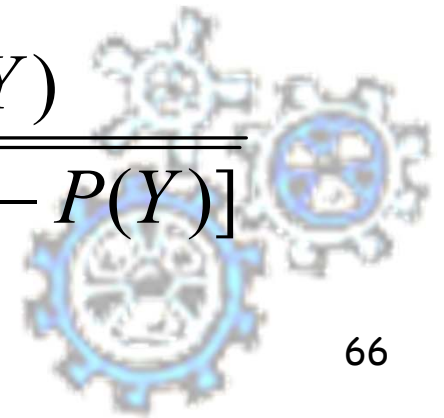
- Measures that take into account statistical dependence

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

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

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

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



Example: Lift/Interest

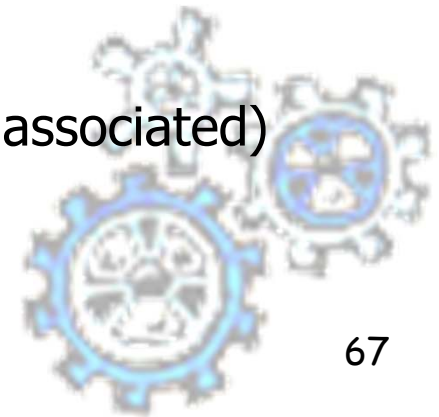
	Coffee	<u>Coffee</u>	
<u>Tea</u>	15	5	20
Tea	75	5	80
	90	10	100

Association Rule: Tea \rightarrow Coffee

Confidence = $P(\text{Coffee}|\text{Tea}) = 0.75$

but $P(\text{Coffee}) = 0.9$

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



Drawback of Lift & Interest

	y	\bar{y}	
x	10	0	10
\bar{x}	0	90	90
	10	90	100

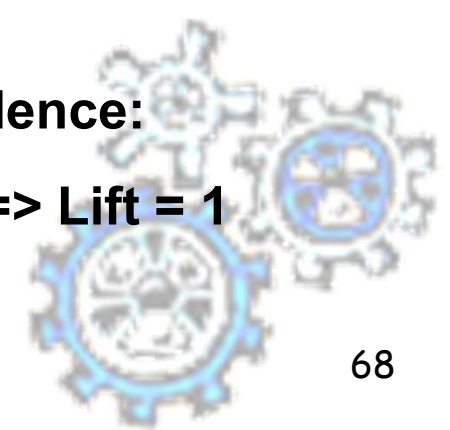
$$Lift = \frac{0.1}{(0.1)(0.1)} = 10$$

	y	\bar{y}	
x	90	0	90
\bar{x}	0	10	10
	90	10	100

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

Statistical independence:

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



There are lots of measures proposed in the literature

Some measures are good for certain applications, but not for others

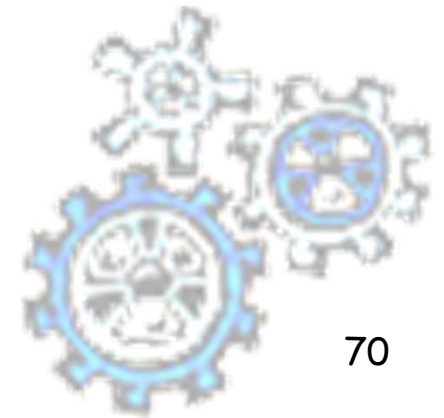
What criteria should we use to determine whether a measure is good or bad?

What about Apriori-style support based pruning? How does it affect these measures?

#	Measure	Formula
1	ϕ -coefficient	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
2	Goodman-Kruskal's (λ)	$\frac{\sum_j \max_k P(A_j, B_k) + \sum_k \max_j P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
3	Odds ratio (α)	$\frac{P(A,B)P(\bar{A},\bar{B})}{P(A,\bar{B})P(\bar{A},B)}$
4	Yule's Q	$\frac{P(A,B)P(\bar{A}\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A}\bar{B}) + P(A,\bar{B})P(\bar{A},B)} = \frac{\alpha - 1}{\alpha + 1}$
5	Yule's Y	$\frac{\sqrt{P(A,B)P(\bar{A}\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A}\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}} = \frac{\sqrt{\alpha - 1}}{\sqrt{\alpha + 1}}$
6	Kappa (κ)	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
7	Mutual Information (M)	$\frac{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\min(-\sum_i P(A_i) \log P(A_i), -\sum_j P(B_j) \log P(B_j))}$
8	J-Measure (J)	$\max \left(P(A, B) \log \left(\frac{P(B A)}{P(B)} \right) + P(\bar{A}\bar{B}) \log \left(\frac{P(\bar{B} \bar{A})}{P(\bar{B})} \right), \right. \\ \left. P(A, B) \log \left(\frac{P(A B)}{P(A)} \right) + P(\bar{A}B) \log \left(\frac{P(\bar{A} B)}{P(\bar{A})} \right) \right)$
9	Gini index (G)	$\max \left(P(A)[P(B A)^2 + P(\bar{B} A)^2] + P(\bar{A})[P(B \bar{A})^2 + P(\bar{B} \bar{A})^2] \right. \\ \left. - P(B)^2 - P(\bar{B})^2, \right. \\ \left. P(B)[P(A B)^2 + P(\bar{A} B)^2] + P(\bar{B})[P(A \bar{B})^2 + P(\bar{A} \bar{B})^2] \right. \\ \left. - P(A)^2 - P(\bar{A})^2 \right)$
10	Support (s)	$P(A, B)$
11	Confidence (c)	$\max(P(B A), P(A B))$
12	Laplace (L)	$\max \left(\frac{NP(A,B)+1}{NP(A)+2}, \frac{NP(A,B)+1}{NP(B)+2} \right)$
13	Conviction (V)	$\max \left(\frac{P(A)P(\bar{B})}{P(\bar{A}B)}, \frac{P(B)P(\bar{A})}{P(\bar{B}A)} \right)$
14	Interest (I)	$\frac{P(A,B)}{P(A)P(B)}$
15	cosine (IS)	$\frac{P(A,B)}{\sqrt{P(A)P(B)}}$
16	Piatetsky-Shapiro's (PS)	$P(A, B) - P(A)P(B)$
17	Certainty factor (F)	$\max \left(\frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)} \right)$
18	Added Value (AV)	$\max(P(B A) - P(B), P(A B) - P(A))$
19	Collective strength (S)	$\frac{P(A,B) + P(\bar{A}\bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A,B) - P(\bar{A}\bar{B})}$
20	Jaccard (ζ)	$\frac{P(A,B)}{P(A) + P(B) - P(A,B)}$
21	Kloggen (K)	$\sqrt{\frac{P(A,B)}{P(A)P(B)}} \max(P(B A) - P(B), P(A B) - P(A))$

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
- $\text{sum(LHS)} + \text{sum(RHS)} > 100$



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

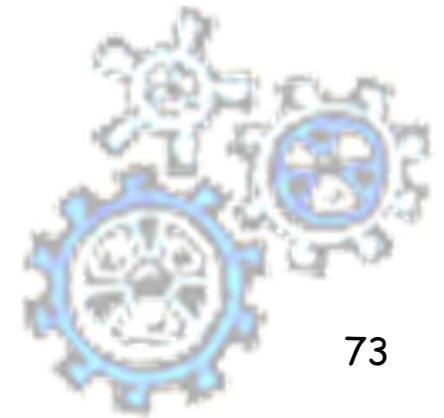
- Use rules to serve dynamic, customized contents to users
- prefetch files that are most likely to be accessed
- determine the best way to structure the Web site (site optimization)

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)**

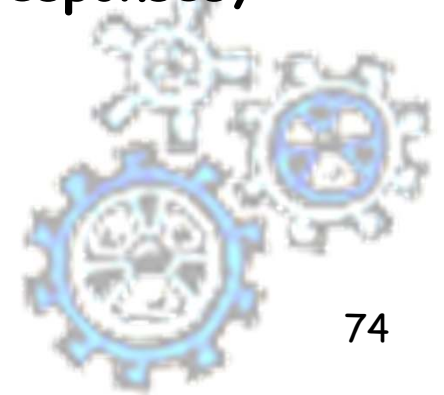
Atherosclerosis prevention study:

- 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



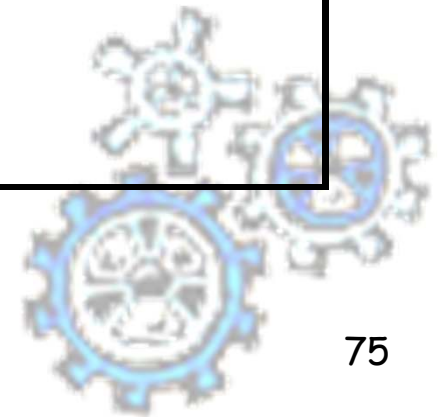
Atherosclerosis prevention study:

- 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 from Entry and Exams		
General characteristics	Examinations	habits
Marital status	Chest pain	Alcohol
Transport to a job	Breathlessness	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

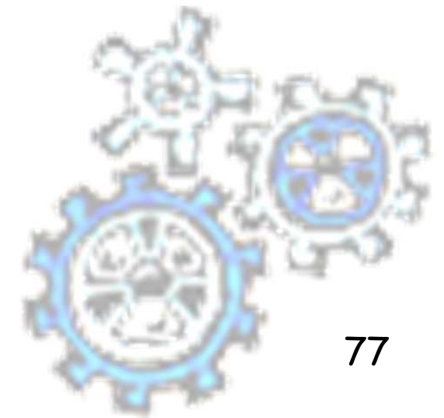
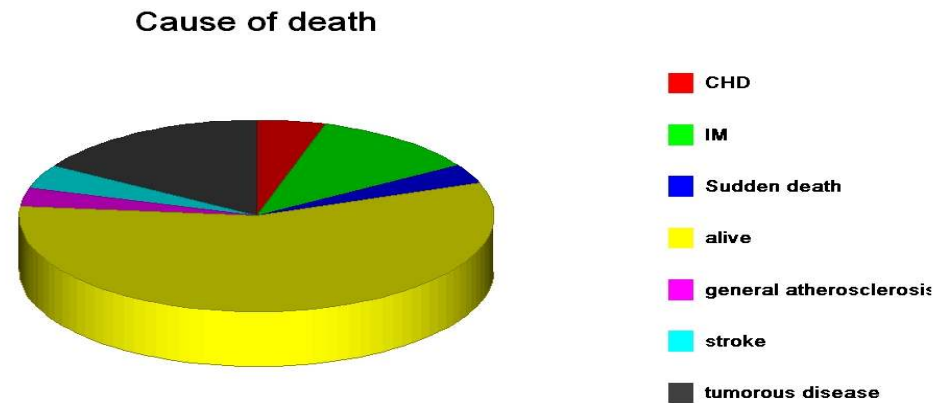


The input data

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 implicitly 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 characteristics		Examinations		Habits		Cause of death
	Activity after work	Education	Chest pain	...	Alcohol	
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

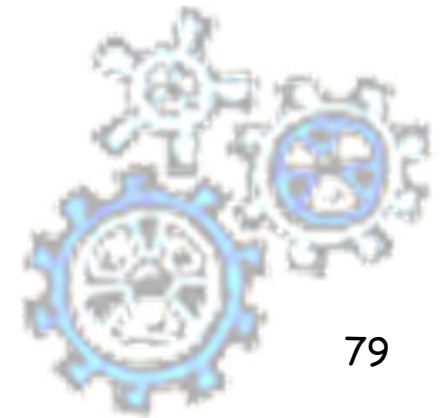
Are there strong relations concerning death cause?

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

Examinations (?) \Rightarrow Death cause (?)

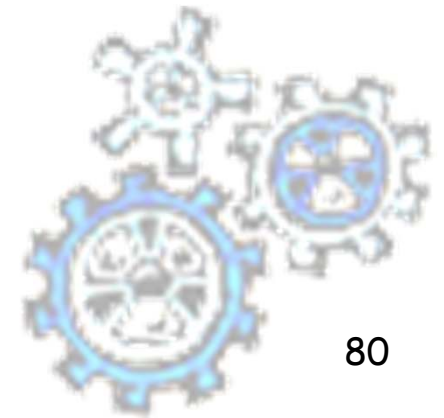
Habits (?) \Rightarrow Death cause (?)

Combinations (?) \Rightarrow Death cause (?)



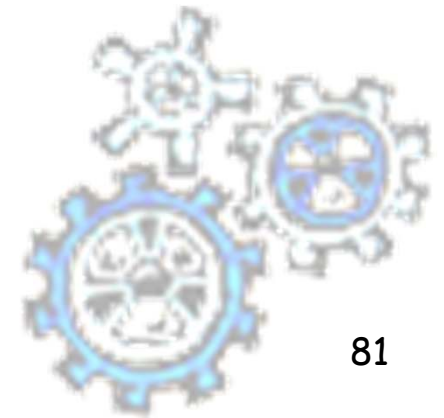
Example of extracted rules

- **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.**



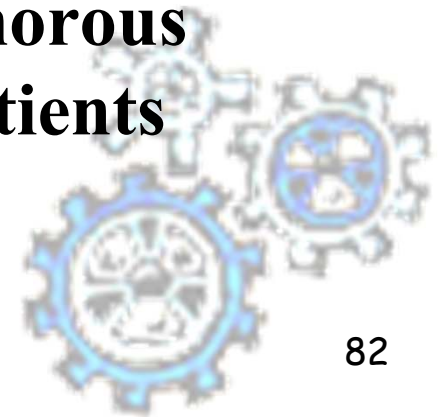
Example of extracted rules

- **Physical activity in work(he mainly sits) & Height<176-180> \Rightarrow 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.**



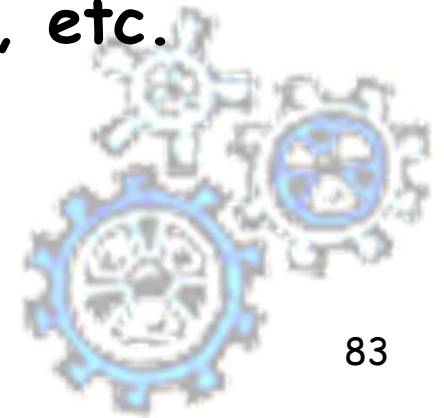
Example of extracted rules

- **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**



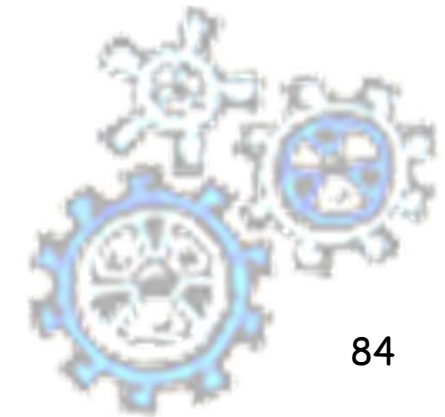
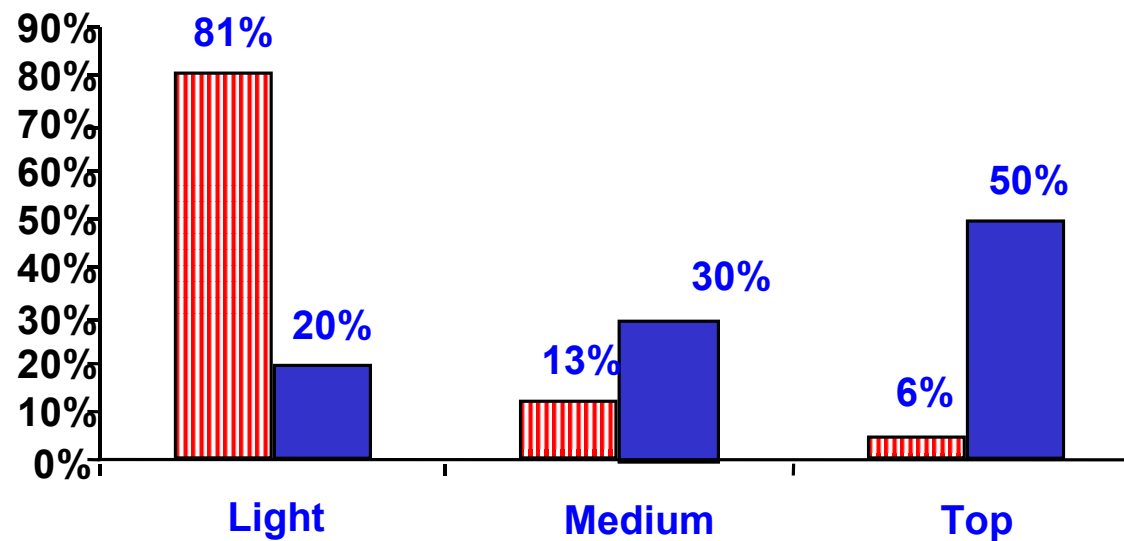
Conclusions

- **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.



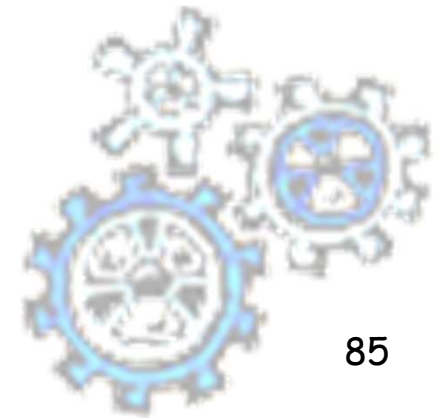
Conclusion (2)

- 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.



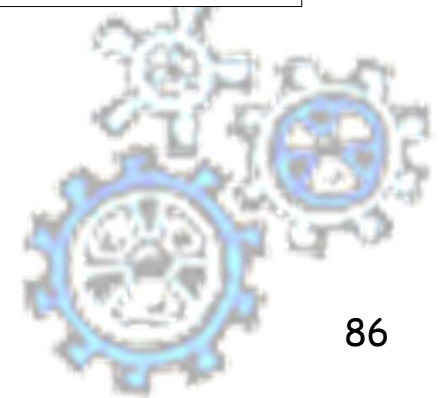
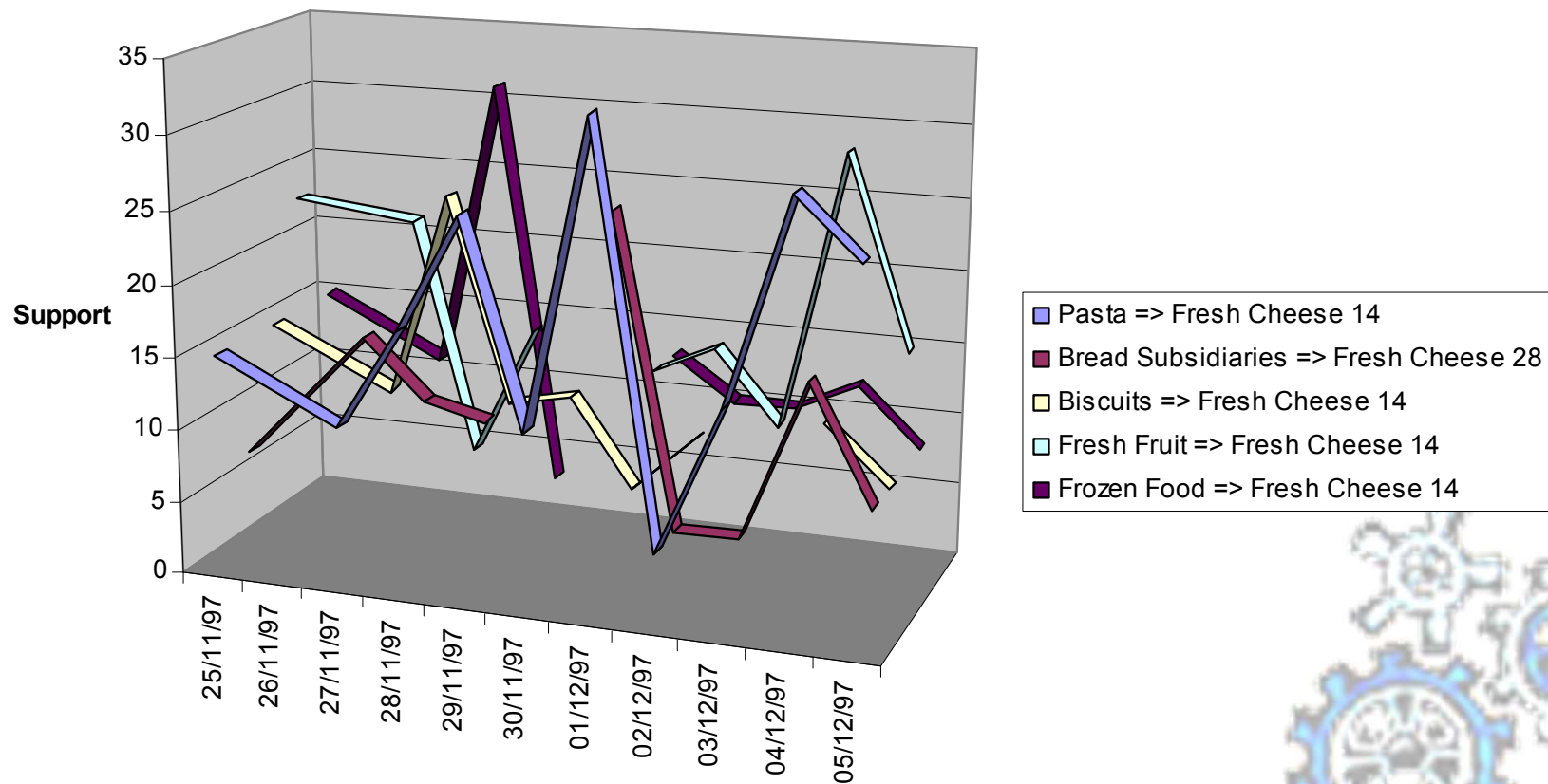
Which tools for market basket analysis?

- Association rules 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?

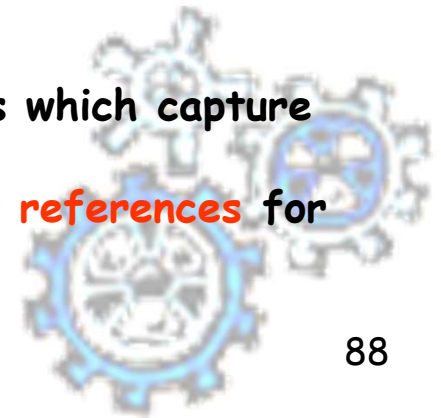


Sequential Patterns



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



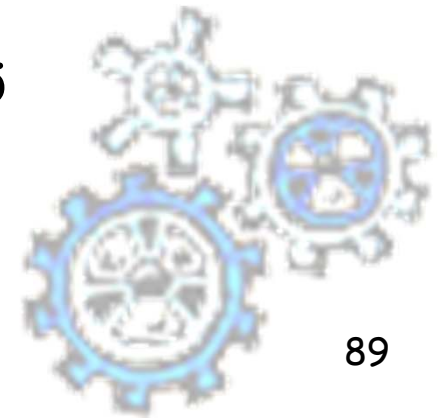
Mining Sequences - Example

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)}
{(C), (DG)}



Intuizione

Obiettivo: personalizzare ed ottimizzare le offerte di vendita ai clienti in base agli acquisti fatti da ciascun cliente in precedenza.

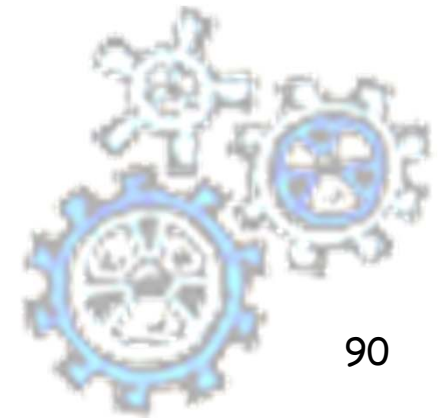
Analisi: studiare il comportamento nel tempo degli acquisti dei clienti !

Metodo: “il 5% dei clienti ha acquistato prima X, poi Y e poi Z”

Requisiti: mantenere traccia degli acquisti dei singoli clienti (nome, fidelity cards, carte di credito, bancomat, e-mail, codice fiscale)

Dominii: vendite al dettaglio, vendite per corrispondenza, vendite su internet, vendite di prodotti finanziari/bancari, analisi mediche

**Intra-Transaction (Regole di Associazione) ...
e Inter-Transaction (Patterns Sequenziali)**



Sequenze

Insieme di transazioni cliente

$$T = \{ (data_1, c_1, t_1), \dots, (data_n, c_n, t_n) \}$$

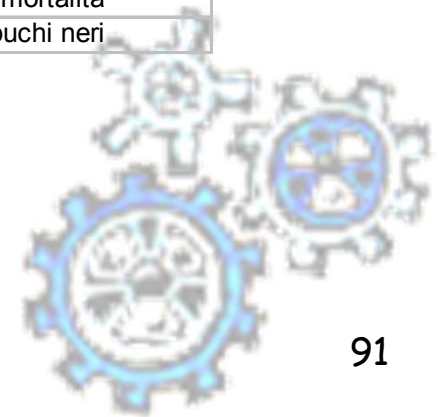
Sequenza di transazioni per cliente c

$$seq(c) = \langle t_1, \dots, t_i, \dots, t_n \rangle$$

ordinate per data

Cliente	Sequenza
1	$\langle \{30\}, \{90\} \rangle$
2	$\langle \{10, 20\}, \{30\}, \{40, 60, 70\} \rangle$
3	$\langle \{10\}, \{30, 50, 70\} \rangle$
4	$\langle \{30\}, \{40, 70\}, \{90\} \rangle$
5	$\langle \{90\} \rangle$

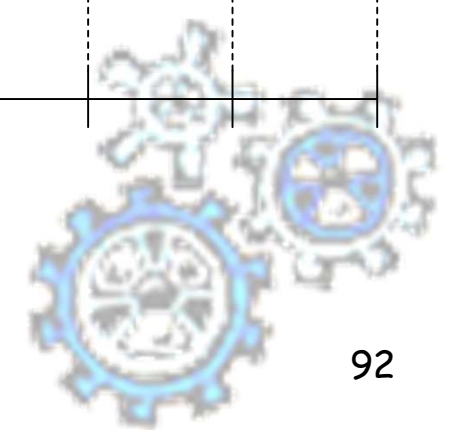
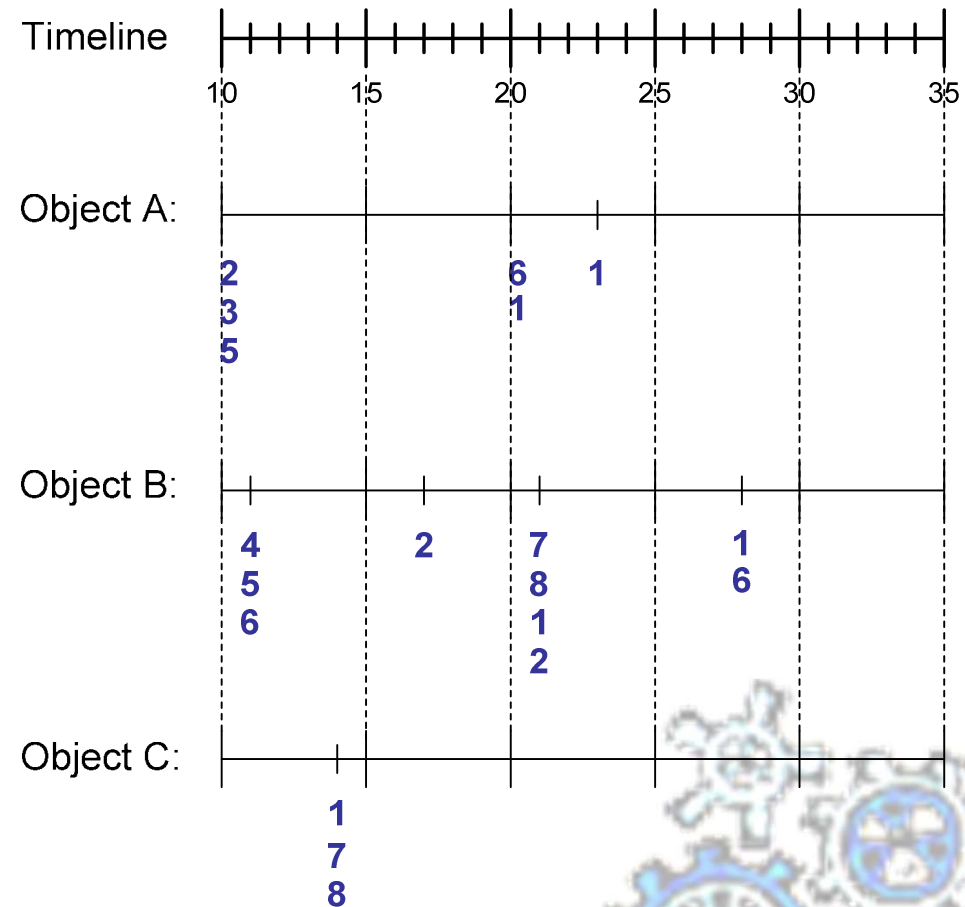
Libro	Titolo
10	Star Wars Episode I
20	La fondazione e l'impero
30	La seconda fondazione
40	Database systems
50	Algoritmi + Strutture Dati =
60	L'insostenibile leggerezza
70	Immortalita'
90	I buchi neri



Sequence Data

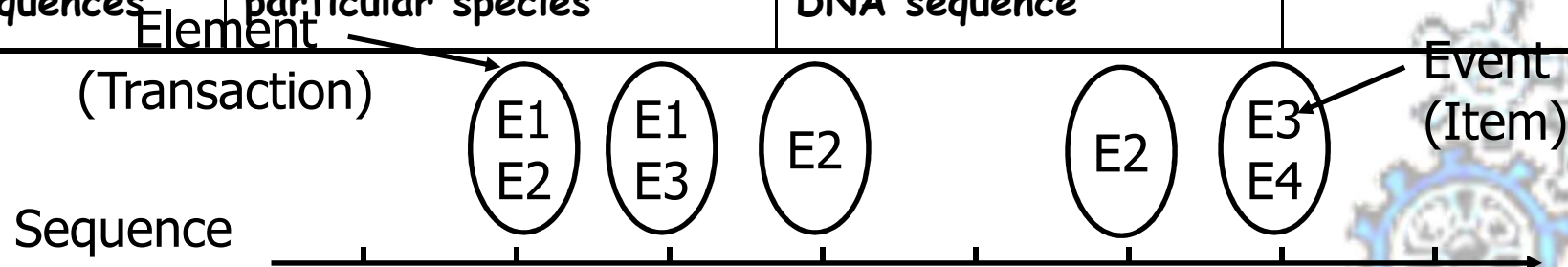
Sequence Database:

Object	Timestamp	Events
A	10	2, 3, 5
A	20	6, 1
A	23	1
B	11	4, 5, 6
B	17	2
B	21	7, 8, 1, 2
B	28	1, 6
C	14	1, 8, 7



Examples of Sequence Data

Sequence Database	Sequence	Element (Transaction)	Event (Item)
Customer	Purchase history of a given customer	A set of items bought by a customer at time t	Books, diary products, CDs, etc
Web Data	Browsing activity of a particular Web visitor	A collection of files viewed by a Web visitor after a single mouse click	Home page, index page, contact info, etc
Event data	History of events generated by a given sensor	Events triggered by a sensor at time t	Types of alarms generated by sensors
Genome sequences	DNA sequence of a particular species	An element of the DNA sequence	Bases A, T, G, C



Formal Definition of a Sequence

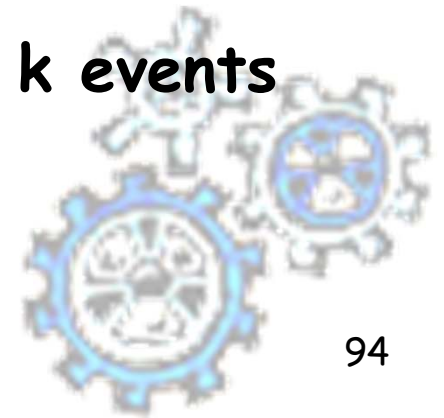
- A sequence is an ordered list of elements (transactions)

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

- Each element contains a collection of events (items)

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

- Each element is attributed to a specific time or location
- Length of a sequence, $|s|$, is given by the number of elements of the sequence
- A k -sequence is a sequence that contains k events (items)



Examples of Sequence

■ Web sequence:

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

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

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

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

■ Sequence of books checked out at a library:

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



Formal Definition of a Subsequence

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

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

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

- The support of a subsequence w is defined as the fraction of data sequences that contain w

- A *sequential pattern* is a frequent subsequence (i.e., a subsequence whose support is $\geq \text{minsup}$)

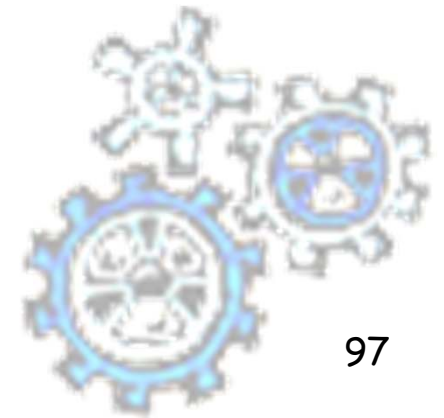
Sequential Pattern Mining: Definition

■ Given:

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

■ Task:

- Find all subsequences with support $\geq \textit{minsup}$

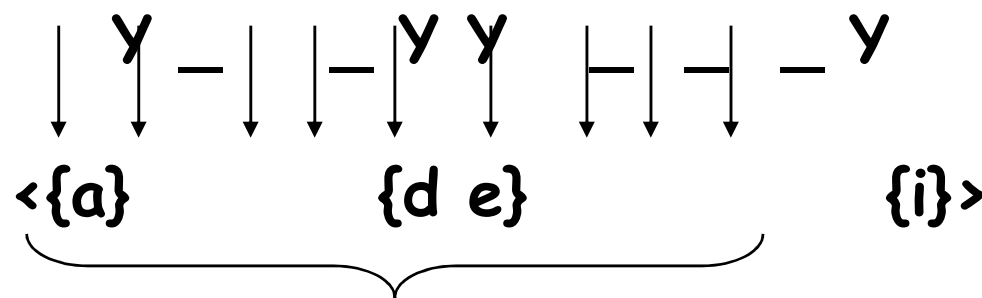


Sequential Pattern Mining: Challenge

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

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

$k=4$:



Answer:

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

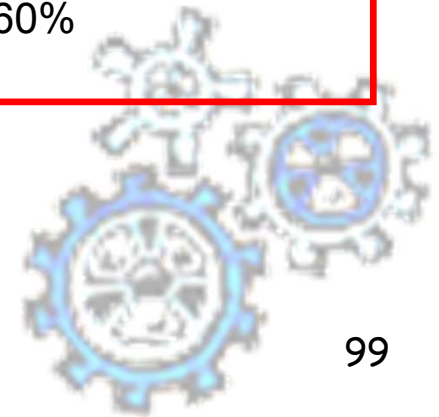
Sequential Pattern Mining: Example

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

Minsup = 50%

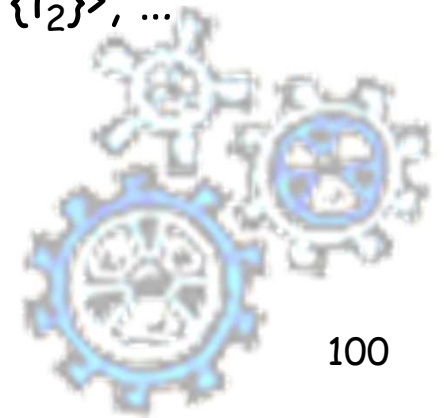
Examples of Frequent Subsequences:

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



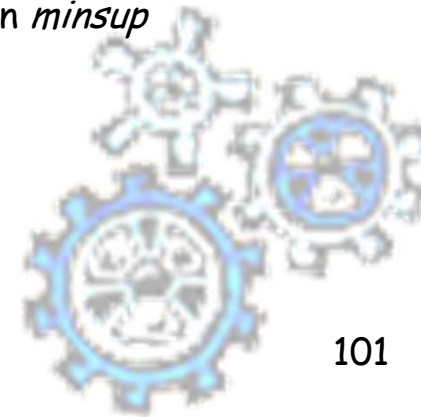
Extracting Sequential Patterns

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

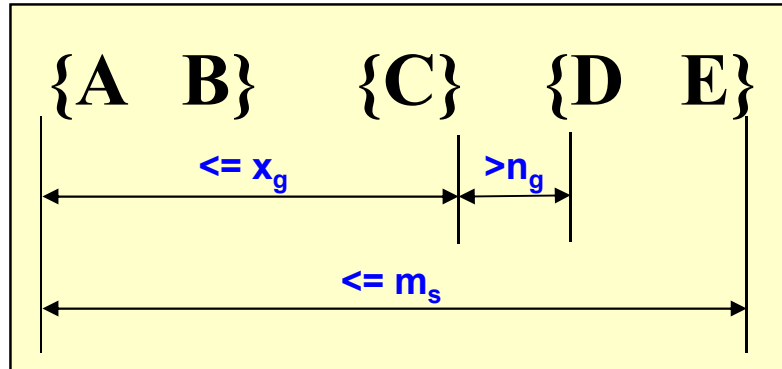


Generalized Sequential Pattern (GSP)

- Step 1:
 - Make the first pass over the sequence database D to yield all the 1-element frequent sequences
- Step 2:
Repeat until no new frequent sequences are found
 - Candidate Generation:
 - ✓ Merge pairs of frequent subsequences found in the $(k-1)$ th pass to generate candidate sequences that contain k items
 - Candidate Pruning:
 - ✓ Prune candidate k -sequences that contain infrequent $(k-1)$ -subsequences
 - Support Counting:
 - ✓ Make a new pass over the sequence database D to find the support for these candidate sequences
 - Candidate Elimination:
 - ✓ Eliminate candidate k -sequences whose actual support is less than *minsup*



Timing Constraints (I)



x_g : max-gap

n_g : min-gap

m_s : maximum span

$x_g = 2, n_g = 0, m_s = 4$ Data sequence	Subsequence	Contain?
$\langle \{2,4\} \{3,5,6\} \{4,7\} \{4,5\} \{8\} \rangle$	$\langle \{6\} \{5\} \rangle$	Yes
$\langle \{1\} \{2\} \{3\} \{4\} \{5\} \rangle$	$\langle \{1\} \{4\} \rangle$	No
$\langle \{1\} \{2,3\} \{3,4\} \{4,5\} \rangle$	$\langle \{2\} \{3\} \{5\} \rangle$	Yes
$\langle \{1,2\} \{3\} \{2,3\} \{3,4\} \{2,4\} \{4,5\} \rangle$	$\langle \{1,2\} \{5\} \rangle$	No

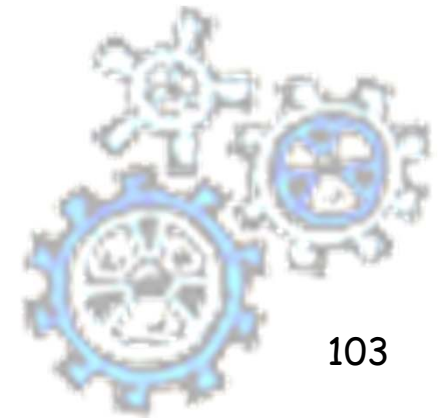
Mining Sequential Patterns with Timing Constraints

■ Approach 1:

- Mine sequential patterns without timing constraints
- Postprocess the discovered patterns

■ Approach 2:

- Modify *GSP* to directly prune candidates that violate timing constraints
- Question:
 - ✓ Does Apriori principle still hold?



Apriori Principle for Sequence Data

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

Suppose:

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

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

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

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

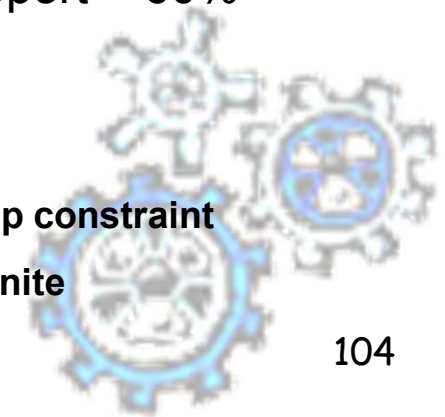
$$\langle \{2\} \{5\} \rangle \text{ support} = 40\%$$

but

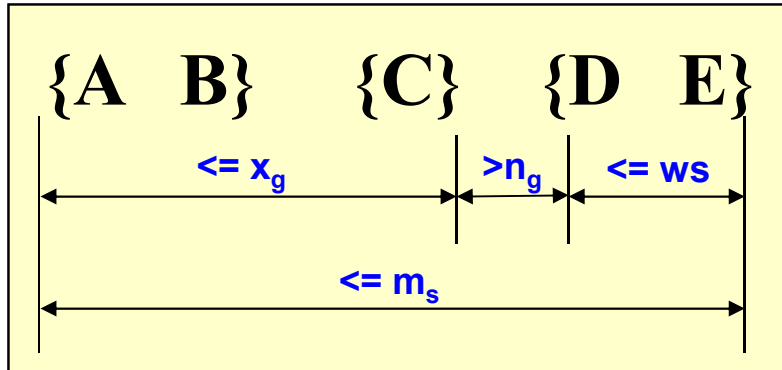
$$\langle \{2\} \{3\} \{5\} \rangle \text{ support} = 60\%$$

Problem exists because of max-gap constraint

No such problem if max-gap is infinite



Timing Constraints



x_g : max-gap

n_g : min-gap

ws: window size

m_s : maximum span

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

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

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.
- S. 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.

