

Data Mining: Concepts and Techniques

— Slides for Textbook —

— Chapter 6 —

Jiawei Han and Micheline Kamber

Intelligent Database Systems Research Lab Simon Fraser
University,

Ari Visa, , Institute of Signal Processing

Tampere University of Technology

Mining Association Rules in Large Databases

- [Association rule mining](#)
- Mining single-dimensional Boolean association rules from transactional databases
- Mining multilevel association rules from transactional databases
- Mining multidimensional association rules from transactional databases and data warehouse
- From association mining to correlation analysis
- Constraint-based association mining
- Summary

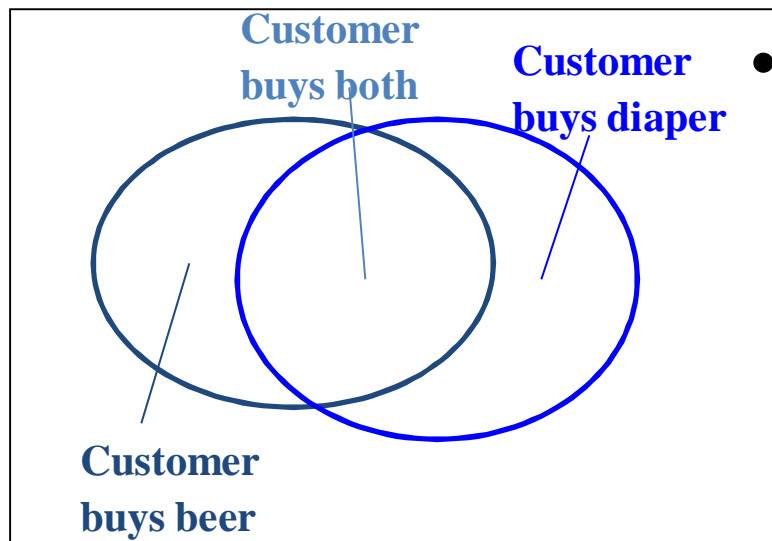
What Is Association Mining?

- Association rule mining:
 - Finding frequent patterns, associations, correlations, or causal structures among sets of items or objects in transaction databases, relational databases, and other information repositories.
- Applications:
 - Basket data analysis, cross-marketing, catalog design, loss-leader analysis, clustering, classification, etc.
- Examples.
 - Rule form: “Body \rightarrow Head [support, confidence]”.
 - $\text{buys}(x, \text{“diapers”}) \rightarrow \text{buys}(x, \text{“beers”}) [0.5\%, 60\%]$
 - $\text{major}(x, \text{“CS”}) \wedge \text{takes}(x, \text{“DB”}) \rightarrow \text{grade}(x, \text{“A”}) [1\%, 75\%]$

Association Rule: Basic Concepts

- Given: (1) database of transactions, (2) each transaction is a list of items (purchased by a customer in a visit)
- Find: all rules that correlate the presence of one set of items with that of another set of items
 - E.g., *98% of people who purchase tires and auto accessories also get automotive services done*
- Applications
 - * \Rightarrow *Maintenance Agreement* (What the store should do to boost Maintenance Agreement sales)
 - *Home Electronics* \Rightarrow * (What other products should the store stocks up?)
 - *Attached mailing* in direct marketing
 - Detecting “ping-pong”ing of patients, faulty “collisions”

Rule Measures: Support and Confidence



- Find all the rules $X \& Y \Rightarrow Z$ with minimum confidence and support
 - support, s , probability that a transaction contains $\{X \wedge Y \wedge Z\}$
 - confidence, c , conditional probability that a transaction having $\{X \wedge Y\}$ also contains Z

Transaction ID	Items Bought
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F

Let minimum support 50%, and minimum confidence 50%, we have

- $A \Rightarrow C$ (50%, 66.6%)
- $C \Rightarrow A$ (50%, 100%)

Association Rule Mining: A Road Map

- Boolean vs. quantitative associations (Based on the types of values handled)
 - $\text{buys}(x, \text{"SQLServer"}) \wedge \text{buys}(x, \text{"DMBook"}) \rightarrow \text{buys}(x, \text{"DBMiner"})$ [0.2%, 60%]
 - $\text{age}(x, \text{"30..39"}) \wedge \text{income}(x, \text{"42..48K"}) \rightarrow \text{buys}(x, \text{"PC"})$ [1%, 75%]
- Single dimension vs. multiple dimensional associations (see ex. Above)
- Single level vs. multiple-level analysis
 - What brands of beers are associated with what brands of diapers?
- Various extensions
 - Correlation, causality analysis
 - Association does not necessarily imply correlation or causality
 - Maxpatterns and closed itemsets
 - Constraints enforced
 - E.g., small sales (sum < 100) trigger big buys (sum > 1,000)?

Mining Association Rules in Large Databases

- Association rule mining
- Mining single-dimensional Boolean association rules from transactional databases
- Mining multilevel association rules from transactional databases
- Mining multidimensional association rules from transactional databases and data warehouse
- From association mining to correlation analysis
- Constraint-based association mining
- Summary

Mining Association Rules—An Example

Transaction ID	Items Bought
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F

Min. support 50%
Min. confidence 50%

Frequent Itemset	Support
{A}	75%
{B}	50%
{C}	50%
{A,C}	50%

For rule $A \Rightarrow C$:

$$\text{support} = \text{support}(\{A \wedge C\}) = 50\%$$

$$\text{confidence} = \text{support}(\{A \wedge C\}) / \text{support}(\{A\}) = 66.6\%$$

The **Apriori** principle:

Any subset of a frequent itemset must be frequent

Mining Frequent Itemsets: the Key Step

- Find the *frequent itemsets*: the sets of items that have minimum support
 - A subset of a frequent itemset must also be a frequent itemset
 - i.e., if $\{AB\}$ 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.

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 Apriori Algorithm — Example

Database D

TID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

Scan D

itemset	sup.
{1}	2
{2}	3
{3}	3
{4}	1
{5}	3

L_1

itemset	sup.
{1}	2
{2}	3
{3}	3
{5}	3

C_2

itemset	sup
{1 2}	1
{1 3}	2
{1 5}	1
{2 3}	2
{2 5}	3
{3 5}	2

C_2

itemset
{1 2}
{1 3}
{1 5}
{2 3}
{2 5}
{3 5}

L_2

itemset	sup
{1 3}	2
{2 3}	2
{2 5}	3
{3 5}	2

C_3

itemset
{2 3 5}

Scan D

L_3

itemset	sup
{2 3 5}	2

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 $p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$
 - from $L_{k-1} p, L_{k-1} q$
 - where $p.item_1=q.item_1, \dots, p.item_{k-2}=q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$
- Step 2: pruning
 - forall *itemsets* c in C_k do
 - forall $(k-1)$ -subsets s of c do
 - if (s is not in L_{k-1}) then delete c from C_k

How to Count Supports of Candidates?

- Why counting supports of candidates a problem?
 - The total number of candidates can be very huge
 - One transaction may contain many candidates
- Method:
 - Candidate itemsets are stored in a *hash-tree*
 - *Leaf node* of hash-tree contains a list of itemsets and counts
 - *Interior node* contains a hash table
 - *Subset function*: finds all the candidates contained in a transaction

Example of Generating Candidates

- $L_3 = \{abc, abd, acd, ace, bcd\}$
- Self-joining: $L_3 * L_3$
 - $abcd$ from abc and abd
 - $acde$ from acd and ace
- Pruning:
 - $acde$ is removed because ade is not in L_3
- $C_4 = \{abcd\}$

Methods to Improve Apriori's Efficiency

- **Hash-based itemset counting**: A k -itemset whose corresponding hashing bucket count is below the threshold cannot be frequent
- **Transaction reduction**: A transaction that does not contain any frequent k -itemset is useless in subsequent scans
- **Partitioning**: Any itemset that is potentially frequent in DB must be frequent in at least one of the partitions of DB
- **Sampling**: mining on a subset of given data, lower support threshold + a method to determine the completeness
- **Dynamic itemset counting**: add new candidate itemsets only when all of their subsets are estimated to be frequent

Is Apriori Fast Enough? — Performance Bottlenecks

- The core of the Apriori algorithm:
 - Use frequent $(k - 1)$ -itemsets to generate [candidate](#) frequent k -itemsets
 - Use database scan and pattern matching to collect counts for the candidate itemsets
- The bottleneck of *Apriori*: [candidate generation](#)
 - Huge candidate sets:
 - 10^4 frequent 1-itemset will generate 10^7 candidate 2-itemsets
 - To discover a frequent pattern of size 100, e.g., $\{a_1, a_2, \dots, a_{100}\}$, one needs to generate $2^{100} \approx 10^{30}$ candidates.
 - Multiple scans of database:
 - Needs $(n + 1)$ scans, n is the length of the longest pattern

Mining Frequent Patterns Without Candidate Generation

- Compress a large database into a compact, Frequent-Pattern tree (FP-tree) structure
 - highly condensed, but complete for frequent pattern mining
 - avoid costly database scans
- Develop an efficient, FP-tree-based frequent pattern mining method
 - A divide-and-conquer methodology: decompose mining tasks into smaller ones
 - Avoid candidate generation: sub-database test only!

Construct FP-tree from a Transaction DB

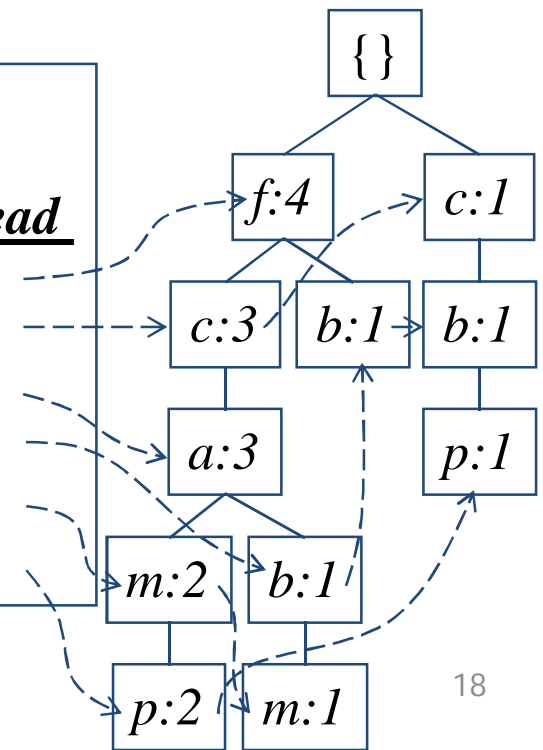
<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

min_support = 0.5

Steps:

1. Scan DB once, find frequent 1-itemset (single item pattern)
2. Order frequent items in frequency descending order
3. Scan DB again, construct FP-tree

Header Table	
<u><i>Item frequency head</i></u>	
<i>f</i>	4
<i>c</i>	4
<i>a</i>	3
<i>b</i>	3
<i>m</i>	3
<i>p</i>	3



Benefits of the FP-tree Structure

- Completeness:
 - never breaks a long pattern of any transaction
 - preserves complete information for frequent pattern mining
- Compactness
 - reduce irrelevant information—infrequent items are gone
 - frequency descending ordering: more frequent items are more likely to be shared
 - never be larger than the original database (if not count node-links and counts)
 - Example: For Connect-4 DB, compression ratio could be over 100

Mining Frequent Patterns Using FP-tree

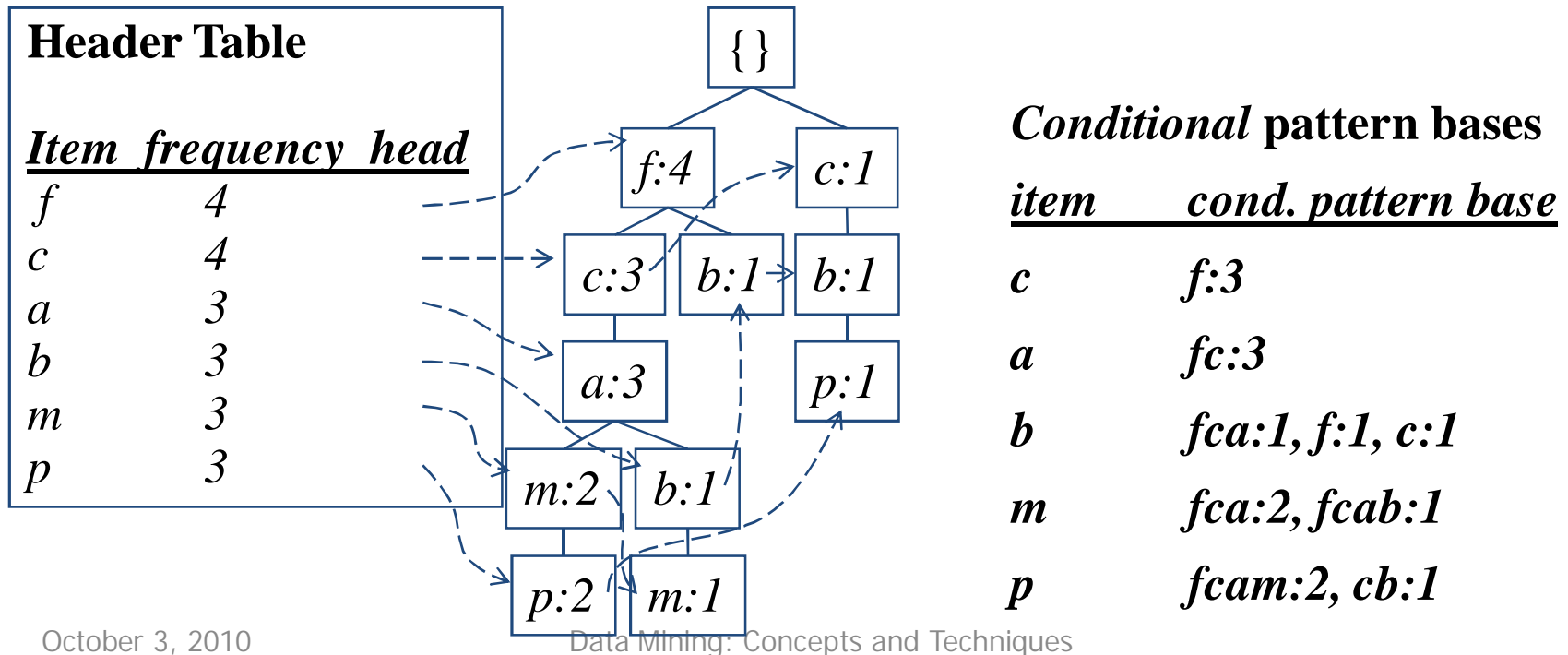
- General idea (divide-and-conquer)
 - Recursively grow frequent pattern path using the FP-tree
- Method
 - For each item, construct its **conditional pattern-base**, and then its **conditional FP-tree**
 - Repeat the process on each newly created conditional FP-tree
 - Until the resulting FP-tree is **empty**, or it contains **only one path** (single path will generate all the combinations of its sub-paths, each of which is a frequent pattern)

Major Steps to Mine FP-tree

- 1) Construct conditional pattern base for each node in the FP-tree
- 2) Construct conditional FP-tree from each conditional pattern-base
- 3) Recursively mine conditional FP-trees and grow frequent patterns obtained so far
 - If the conditional FP-tree contains a single path, simply enumerate all the patterns

Step 1: From FP-tree to Conditional Pattern Base

- Starting at the frequent header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item
- Accumulate all of transformed prefix paths of that item to form a conditional pattern base

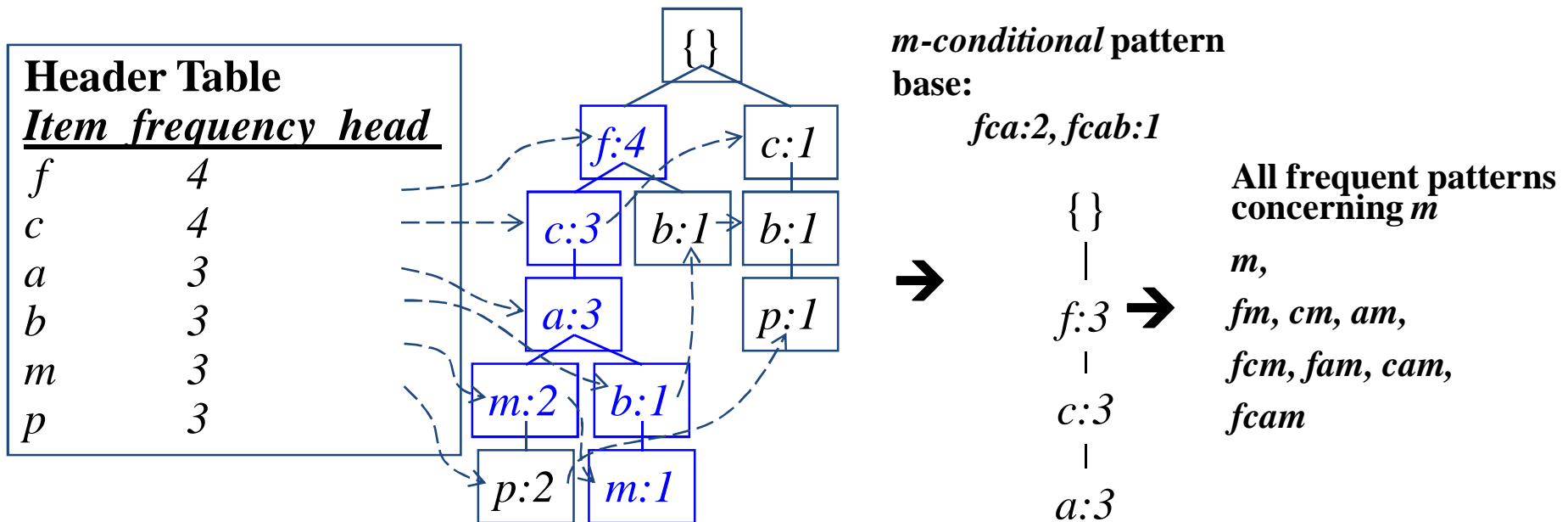


Properties of FP-tree for Conditional Pattern Base Construction

- Node-link property
 - For any frequent item a_i , all the possible frequent patterns that contain a_i can be obtained by following a_i 's node-links, starting from a_i 's head in the FP-tree header
- Prefix path property
 - To calculate the frequent patterns for a node a_i in a path P , only the prefix sub-path of a_i in P need to be accumulated, and its frequency count should carry the same count as node a_i .

Step 2: Construct Conditional FP-tree

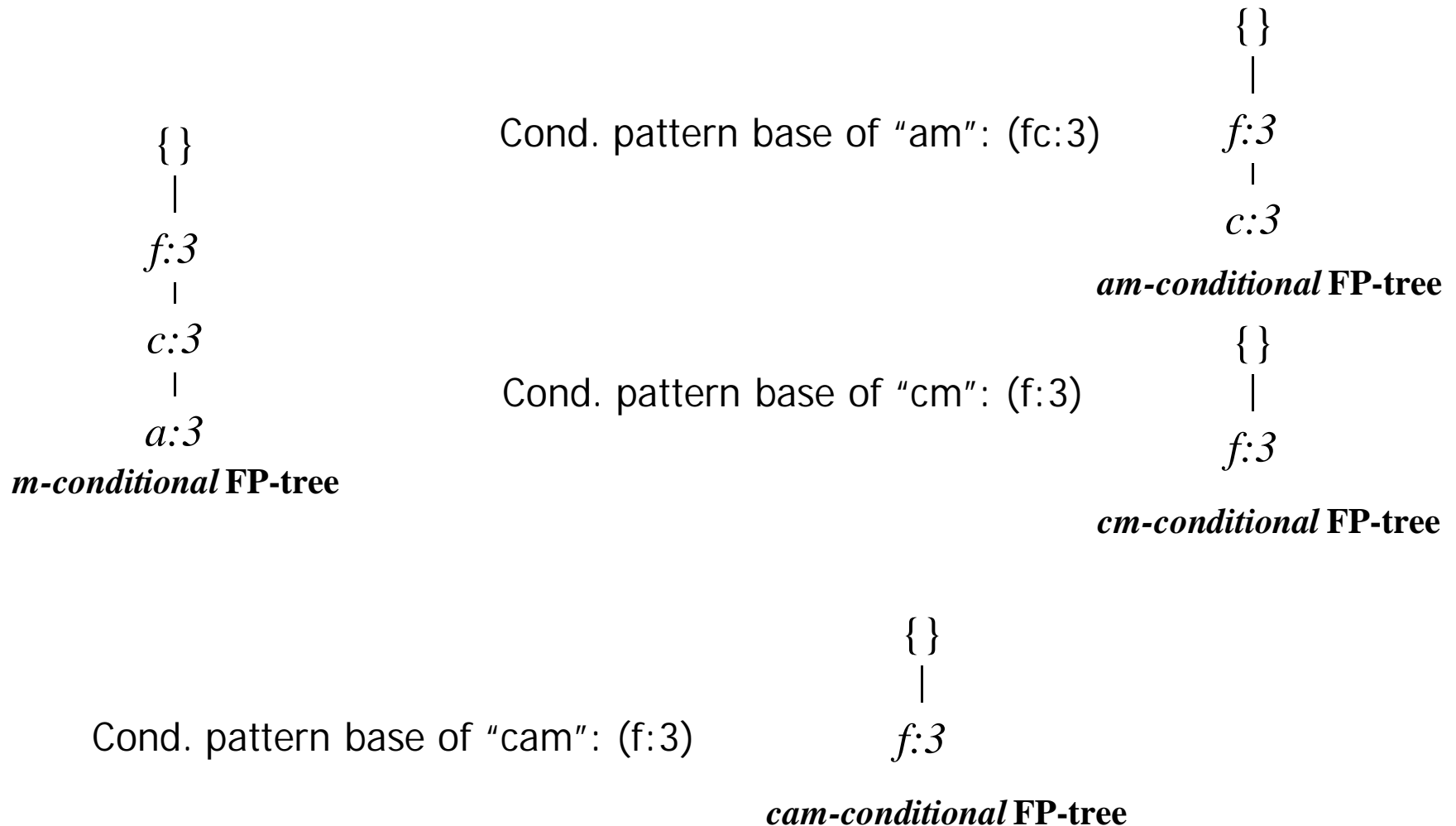
- For each pattern-base
 - Accumulate the count for each item in the base
 - Construct the FP-tree for the frequent items of the pattern base



Mining Frequent Patterns by Creating Conditional Pattern-Bases

Item	Conditional pattern-base	Conditional FP-tree
p	{(fcam:2), (cb:1)}	{(c:3)} p
m	{(fca:2), (fcab:1)}	{(f:3, c:3, a:3)} m
b	{(fca:1), (f:1), (c:1)}	Empty
a	{(fc:3)}	{(f:3, c:3)} a
c	{(f:3)}	{(f:3)} c
f	Empty	Empty

Step 3: Recursively mine the conditional FP-tree



Single FP-tree Path Generation

- Suppose an FP-tree T has a single path P
- The complete set of frequent pattern of T can be generated by enumeration of all the combinations of the sub-paths of P

{
|
f:3
|
c:3
|
a:3



**All frequent patterns
concerning *m***

m,
fm, cm, am,
fcm, fam, cam,
fcam

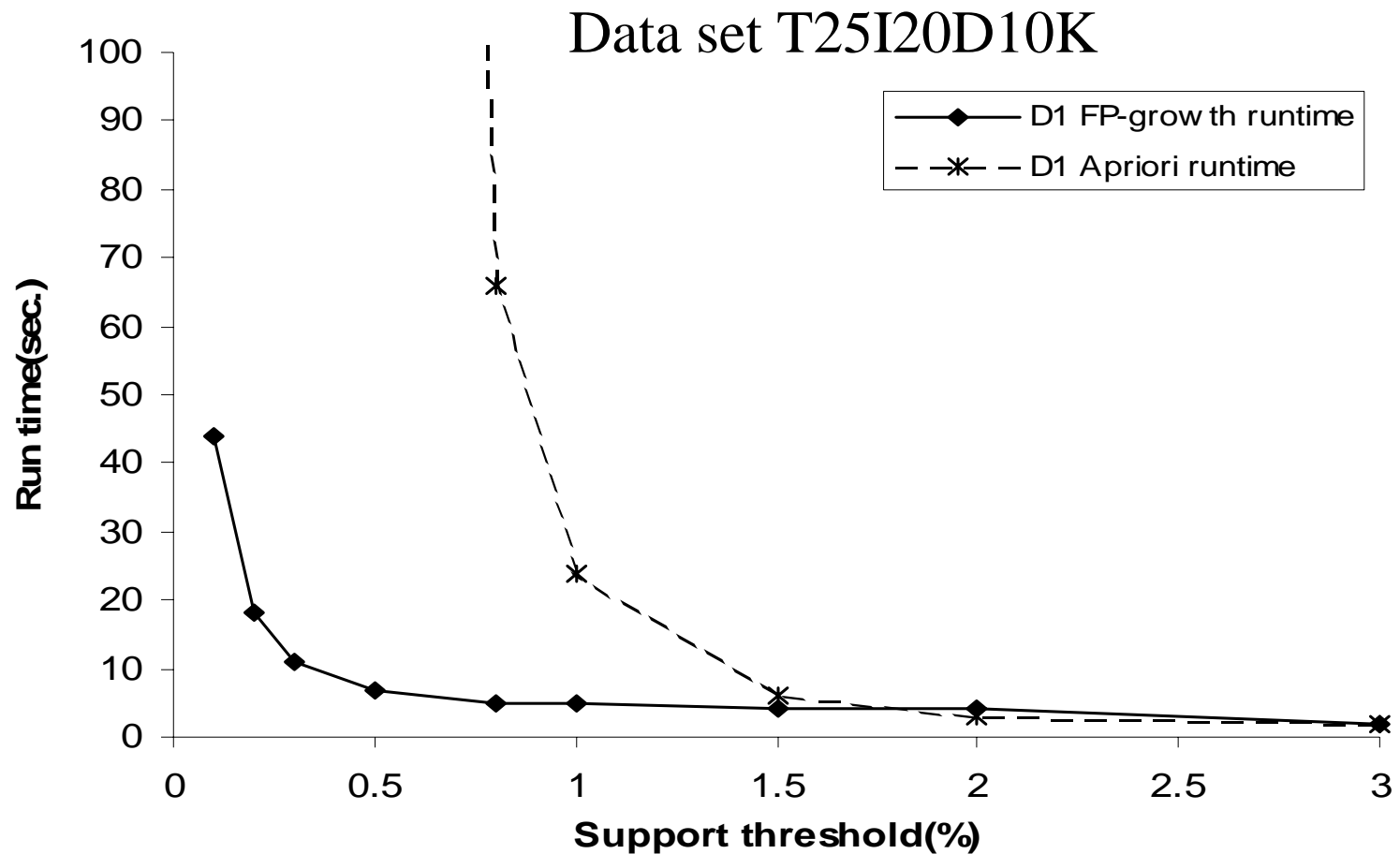
Principles of Frequent Pattern Growth

- Pattern growth property
 - Let α be a frequent itemset in DB, B be α 's conditional pattern base, and β be an itemset in B. Then $\alpha \cup \beta$ is a frequent itemset in DB iff β is frequent in B.
- "*abcdef*" is a frequent pattern, if and only if
 - "*abcde*" is a frequent pattern, and
 - "*f*" is frequent in the set of transactions containing "*abcde*"

Why Is Frequent Pattern Growth Fast?

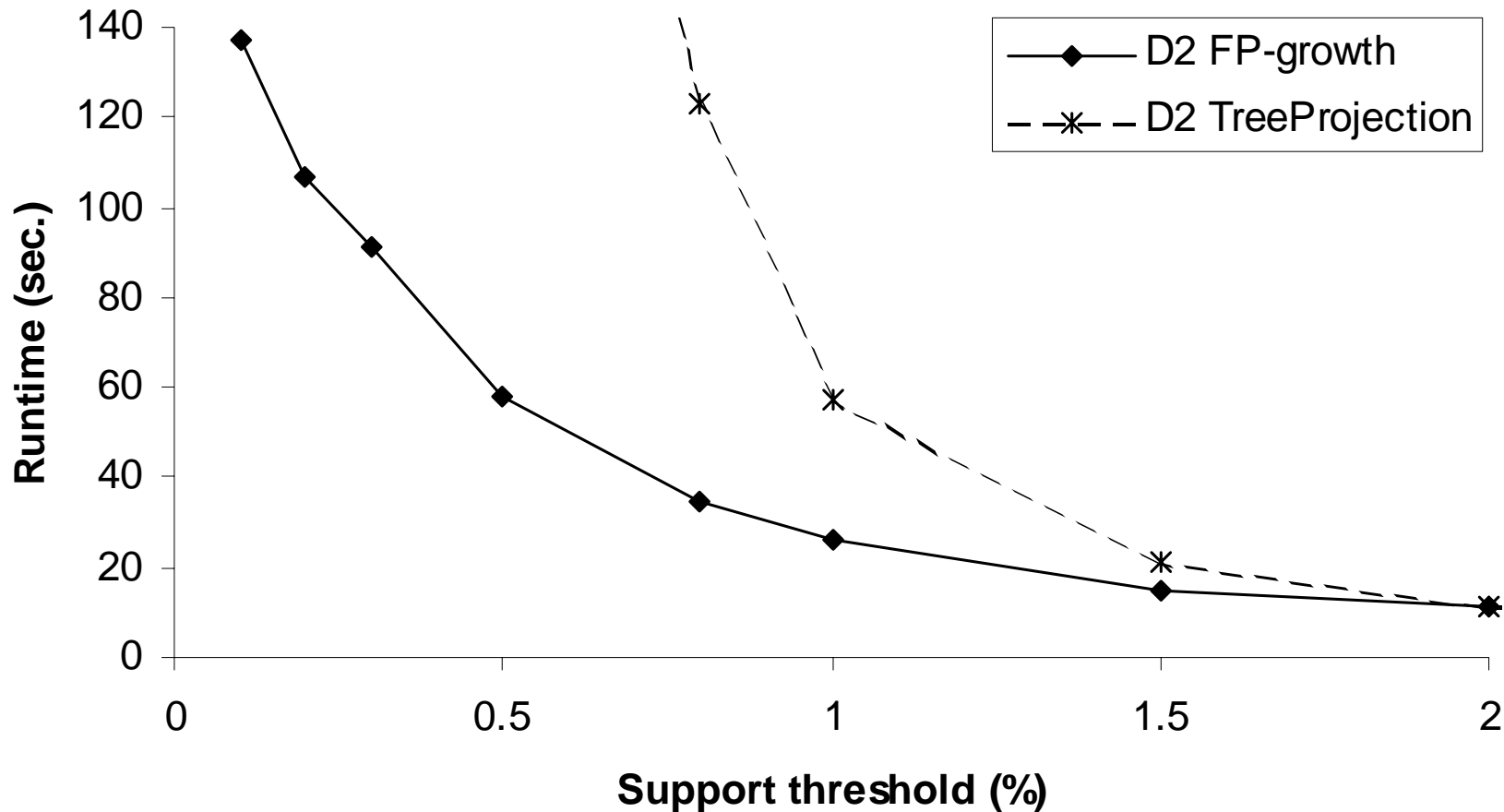
- Our performance study shows
 - FP-growth is an order of magnitude faster than Apriori, and is also faster than tree-projection
- Reasoning
 - No candidate generation, no candidate test
 - Use compact data structure
 - Eliminate repeated database scan
 - Basic operation is counting and FP-tree building

FP-growth vs. Apriori: Scalability With the Support Threshold



FP-growth vs. Tree-Projection: Scalability with Support Threshold

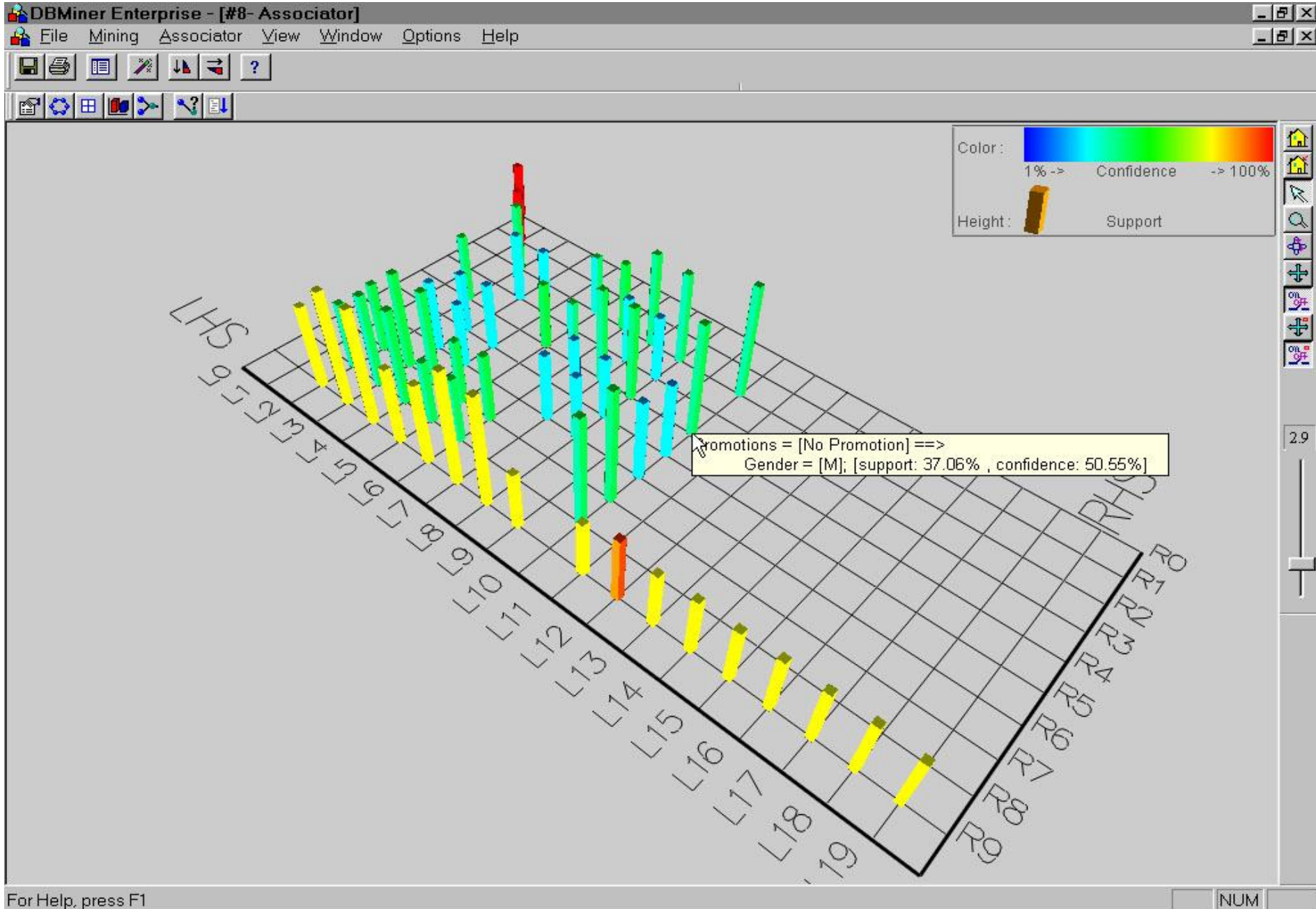
Data set T25I20D100K



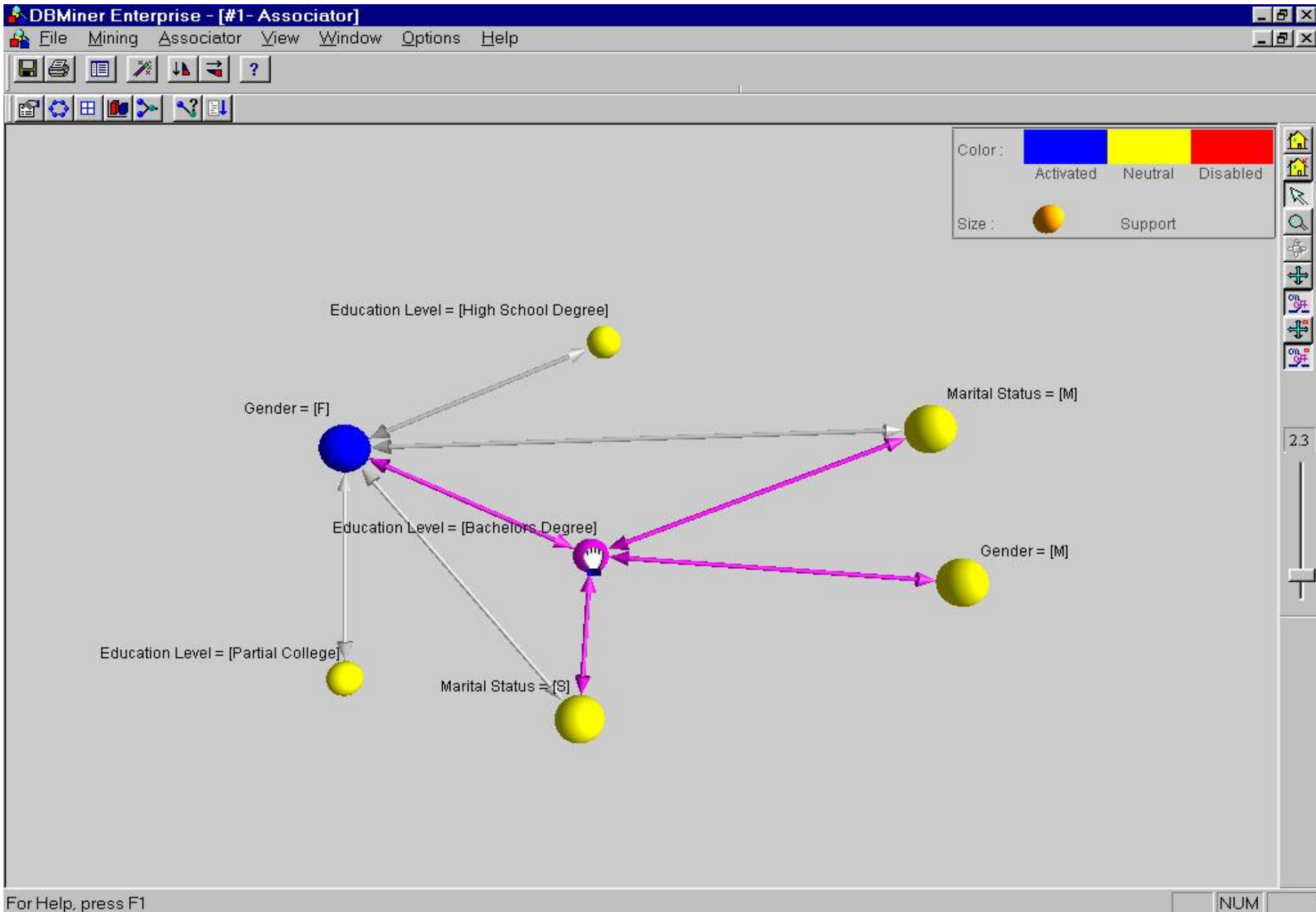
Presentation of Association Rules (Table Form)

	Body	Implies	Head	Supp (%)	Conf (%)	F	G	H	I
1	cost(x) = '0.00~1000.00'	==>	revenue(x) = '0.00~500.00'	28.45	40.4				
2	cost(x) = '0.00~1000.00'	==>	revenue(x) = '500.00~1000.00'	20.46	29.05				
3	cost(x) = '0.00~1000.00'	==>	order_qty(x) = '0.00~100.00'	59.17	84.04				
4	cost(x) = '0.00~1000.00'	==>	revenue(x) = '1000.00~1500.00'	10.45	14.84				
5	cost(x) = '0.00~1000.00'	==>	region(x) = 'United States'	22.56	32.04				
6	cost(x) = '1000.00~2000.00'	==>	order_qty(x) = '0.00~100.00'	12.91	69.34				
7	order_qty(x) = '0.00~100.00'	==>	revenue(x) = '0.00~500.00'	28.45	34.54				
8	order_qty(x) = '0.00~100.00'	==>	cost(x) = '1000.00~2000.00'	12.91	15.67				
9	order_qty(x) = '0.00~100.00'	==>	region(x) = 'United States'	25.9	31.45				
10	order_qty(x) = '0.00~100.00'	==>	cost(x) = '0.00~1000.00'	59.17	71.86				
11	order_qty(x) = '0.00~100.00'	==>	product_line(x) = 'Tents'	13.52	16.42				
12	order_qty(x) = '0.00~100.00'	==>	revenue(x) = '500.00~1000.00'	19.67	23.88				
13	product_line(x) = 'Tents'	==>	order_qty(x) = '0.00~100.00'	13.52	98.72				
14	region(x) = 'United States'	==>	order_qty(x) = '0.00~100.00'	25.9	81.94				
15	region(x) = 'United States'	==>	cost(x) = '0.00~1000.00'	22.56	71.39				
16	revenue(x) = '0.00~500.00'	==>	cost(x) = '0.00~1000.00'	28.45	100				
17	revenue(x) = '0.00~500.00'	==>	order_qty(x) = '0.00~100.00'	28.45	100				
18	revenue(x) = '1000.00~1500.00'	==>	cost(x) = '0.00~1000.00'	10.45	96.75				
19	revenue(x) = '500.00~1000.00'	==>	cost(x) = '0.00~1000.00'	20.46	100				
20	revenue(x) = '500.00~1000.00'	==>	order_qty(x) = '0.00~100.00'	19.67	96.14				
21									
22									
23	cost(x) = '0.00~1000.00'	==>	revenue(x) = '0.00~500.00' AND order_qty(x) = '0.00~100.00'	28.45	40.4				
24	cost(x) = '0.00~1000.00'	==>	revenue(x) = '0.00~500.00' AND order_qty(x) = '0.00~100.00'	28.45	40.4				
25	cost(x) = '0.00~1000.00'	==>	revenue(x) = '500.00~1000.00' AND order_qty(x) = '0.00~100.00'	19.67	27.93				
26	cost(x) = '0.00~1000.00'	==>	revenue(x) = '500.00~1000.00' AND order_qty(x) = '0.00~100.00'	19.67	27.93				
27	cost(x) = '0.00~1000.00' AND order_qty(x) = '0.00~100.00'	==>	revenue(x) = '500.00~1000.00'	19.67	33.23				

Visualization of Association Rule Using Plane Graph



Visualization of Association Rule Using Rule Graph



Iceberg Queries

- **Iceberg query**: Compute aggregates over one or a set of attributes only for those whose aggregate values is above certain threshold
- Example:

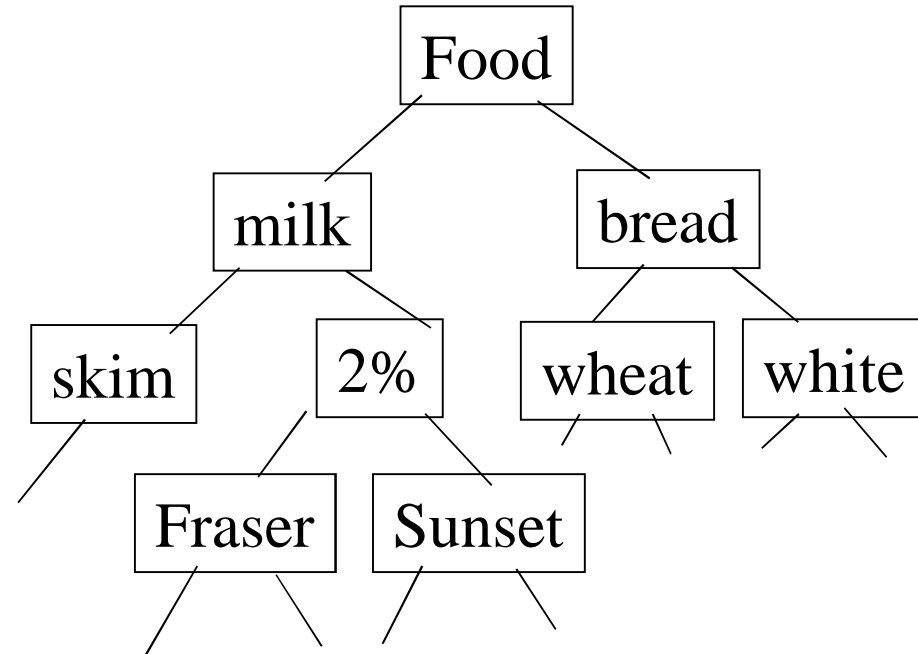
```
select P.custID, P.itemID, sum(P.qty)
from purchase P
group by P.custID, P.itemID
having sum(P.qty) >= 10
```
- **Compute** iceberg queries efficiently **by Apriori**:
 - First compute lower dimensions
 - Then compute higher dimensions only when **all** the lower ones are above the threshold

Mining Association Rules in Large Databases

- Association rule mining
- Mining single-dimensional Boolean association rules from transactional databases
- Mining multilevel association rules from transactional databases
- Mining multidimensional association rules from transactional databases and data warehouse
- From association mining to correlation analysis
- Constraint-based association mining
- Summary

Multiple-Level Association Rules

- Items often form hierarchy.
- Items at the lower level are expected to have lower support.
- Rules regarding itemsets at appropriate levels could be quite useful.
- Transaction database can be encoded based on dimensions and levels
- We can explore shared multi-level mining



TID	Items
T1	{111, 121, 211, 221}
T2	{111, 211, 222, 323}
T3	{112, 122, 221, 411}
T4	{111, 121}
T5	{111, 122, 211, 221, 413}

Mining Multi-Level Associations

- A top_down, progressive deepening approach:
 - First find high-level strong rules:
milk → bread [20%, 60%].
 - Then find their lower-level “weaker” rules:
2% milk → wheat bread [6%, 50%].
- Variations at mining multiple-level association rules.
 - Level-crossed association rules:
2% milk → *Wonder* wheat bread
 - Association rules with multiple, alternative hierarchies:
2% milk → *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.
 - – Lower level items do not occur as frequently. 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
 - There are 4 search strategies:
 - Level-by-level independent
 - Level-cross filtering by k-itemset
 - Level-cross filtering by single item
 - Controlled level-cross filtering by single item

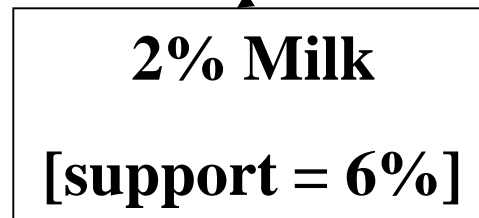
Uniform Support

Multi-level mining with uniform support

Level 1
min_sup = 5%



Level 2
min_sup = 5%



[Back](#)

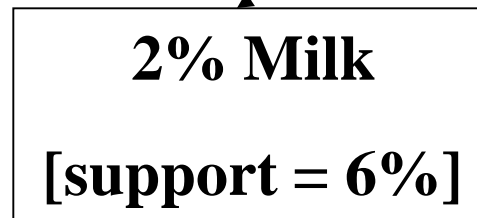
Reduced Support

Multi-level mining with reduced support

Level 1
min_sup = 5%



Level 2
min_sup = 3%



[Back](#)

Multi-level Association: Redundancy Filtering

- Some rules may be redundant due to “ancestor” relationships between items.
- Example
 - milk \Rightarrow wheat bread [support = 8%, confidence = 70%]
 - 2% milk \Rightarrow wheat bread [support = 2%, confidence = 72%]
- We say the first rule is an ancestor of the second rule.
- A rule is redundant if its support is close to the “expected” value, based on the rule’s ancestor.

Multi-Level Mining: Progressive Deepening

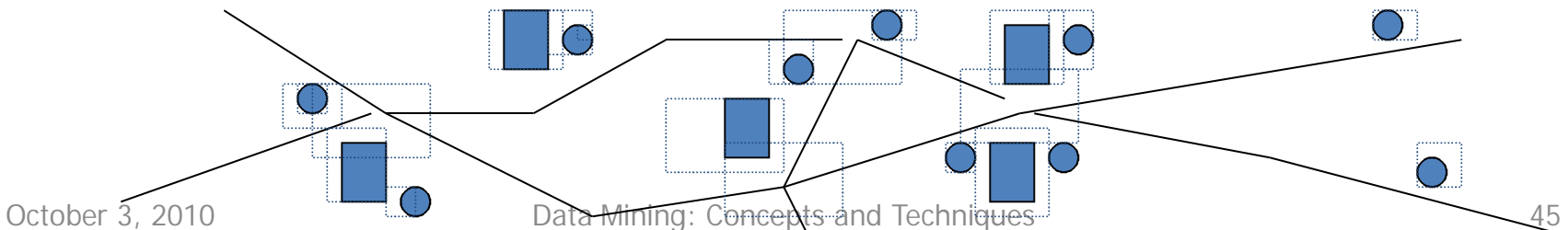
- A top-down, progressive deepening approach:
 - First mine high-level frequent items:
milk (15%), bread (10%)
 - Then mine their lower-level “weaker” frequent itemsets:
2% milk (5%), wheat bread (4%)
- Different $min_support$ threshold across multi-levels lead to different algorithms:
 - If adopting the same $min_support$ across multi-levels then toss t if any of t 's ancestors is infrequent.
 - If adopting reduced $min_support$ at lower levels then examine only those descendents whose ancestor's support is frequent/non-negligible.

Progressive Refinement of Data Mining Quality

- Why progressive refinement?
 - Mining operator can be expensive or cheap, fine or rough
 - Trade speed with quality: step-by-step refinement.
- Superset coverage property:
 - Preserve all the positive answers—allow a positive false test but not a false negative test.
- Two- or multi-step mining:
 - First apply rough/cheap operator (superset coverage)
 - Then apply expensive algorithm on a substantially reduced candidate set (Koperski & Han, [SSD'95](#)).

Progressive Refinement Mining of Spatial Association Rules

- Hierarchy of spatial relationship:
 - “g_close_to”: near_by, touch, intersect, contain, etc.
 - First search for rough relationship and then refine it.
- Two-step mining of spatial association:
 - Step 1: rough spatial computation (as a filter)
 - Using MBR or R-tree for rough estimation.
 - Step2: Detailed spatial algorithm (as refinement)
 - Apply only to those objects which have passed the rough spatial association test (no less than *min_support*)



Mining Association Rules in Large Databases

- Association rule mining
- Mining single-dimensional Boolean association rules from transactional databases
- Mining multilevel association rules from transactional databases
- Mining multidimensional association rules from transactional databases and data warehouse
- From association mining to correlation analysis
- Constraint-based association mining
- Summary

Multi-Dimensional Association: Concepts

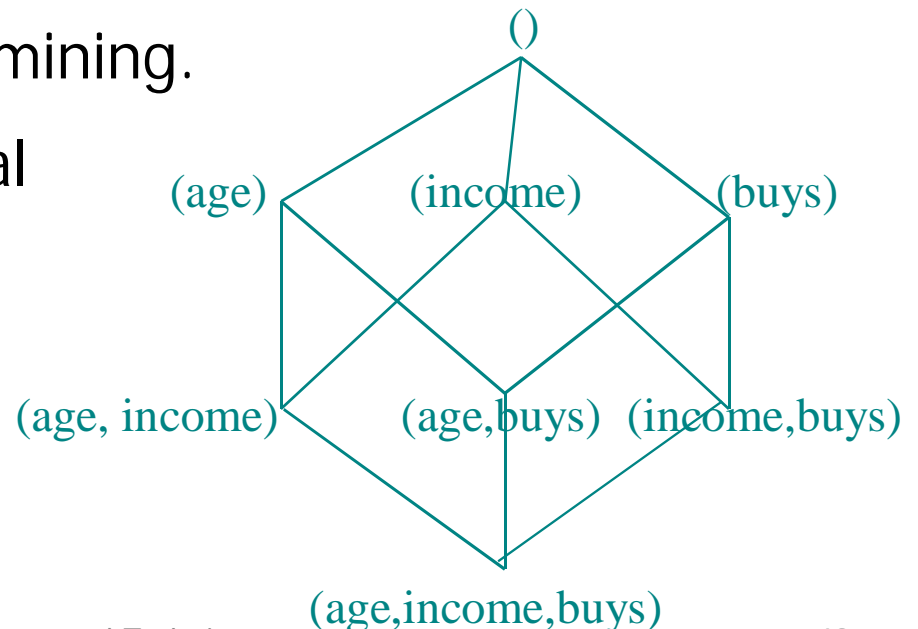
- Single-dimensional rules:
 $\text{buys}(X, \text{"milk"}) \Rightarrow \text{buys}(X, \text{"bread"})$
- Multi-dimensional rules: ○ 2 dimensions or predicates
 - Inter-dimension association rules (*no repeated predicates*)
 $\text{age}(X, \text{"19-25"}) \wedge \text{occupation}(X, \text{"student"}) \Rightarrow \text{buys}(X, \text{"coke"})$
 - hybrid-dimension association rules (*repeated predicates*)
 $\text{age}(X, \text{"19-25"}) \wedge \text{buys}(X, \text{"popcorn"}) \Rightarrow \text{buys}(X, \text{"coke"})$
- Categorical Attributes
 - finite number of possible values, no ordering among values
- Quantitative Attributes
 - numeric, implicit ordering among values

Techniques for Mining MD Associations

- Search for frequent k -predicate set:
 - Example: {age, occupation, buys} is a 3-predicate set.
 - Techniques can be categorized by how age are treated.
- 1. Using static discretization of quantitative attributes
 - Quantitative attributes are statically discretized by using predefined concept hierarchies.
- 2. Quantitative association rules
 - Quantitative attributes are dynamically discretized into “bins” based on the distribution of the data.
- 3. Distance-based association rules
 - This is a dynamic discretization process that considers the distance between data points.

Static Discretization of Quantitative Attributes

- Discretized prior to mining using concept hierarchy.
- Numeric values are replaced by ranges.
- In relational database, finding all frequent k -predicate sets will require k or $k+1$ table scans.
- Data cube is well suited for mining.
- The cells of an n -dimensional cuboid correspond to the predicate sets.
- Mining from data cubes can be much faster.



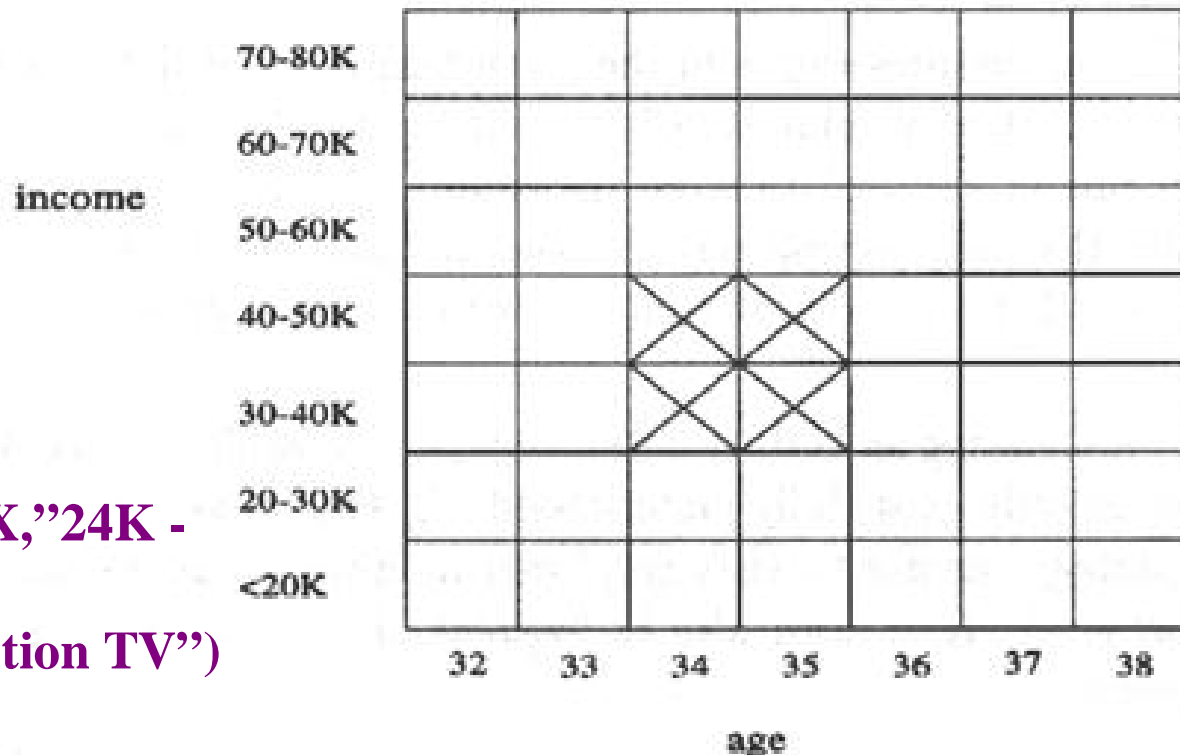
Quantitative Association Rules

- Numeric attributes are *dynamically* discretized
 - Such that the confidence or compactness of the rules mined is maximized.
- 2-D quantitative association rules: $A_{\text{quan1}} \wedge A_{\text{quan2}} \Rightarrow A_{\text{cat}}$
- Cluster "adjacent" association rules to form general rules using a 2-D grid.

- Example:

$\text{age}(X, "30-34") \wedge \text{income}(X, "24K - 48K")$

$\Rightarrow \text{buys}(X, "high\ resolution\ TV")$



ARCS (Association Rule Clustering System)

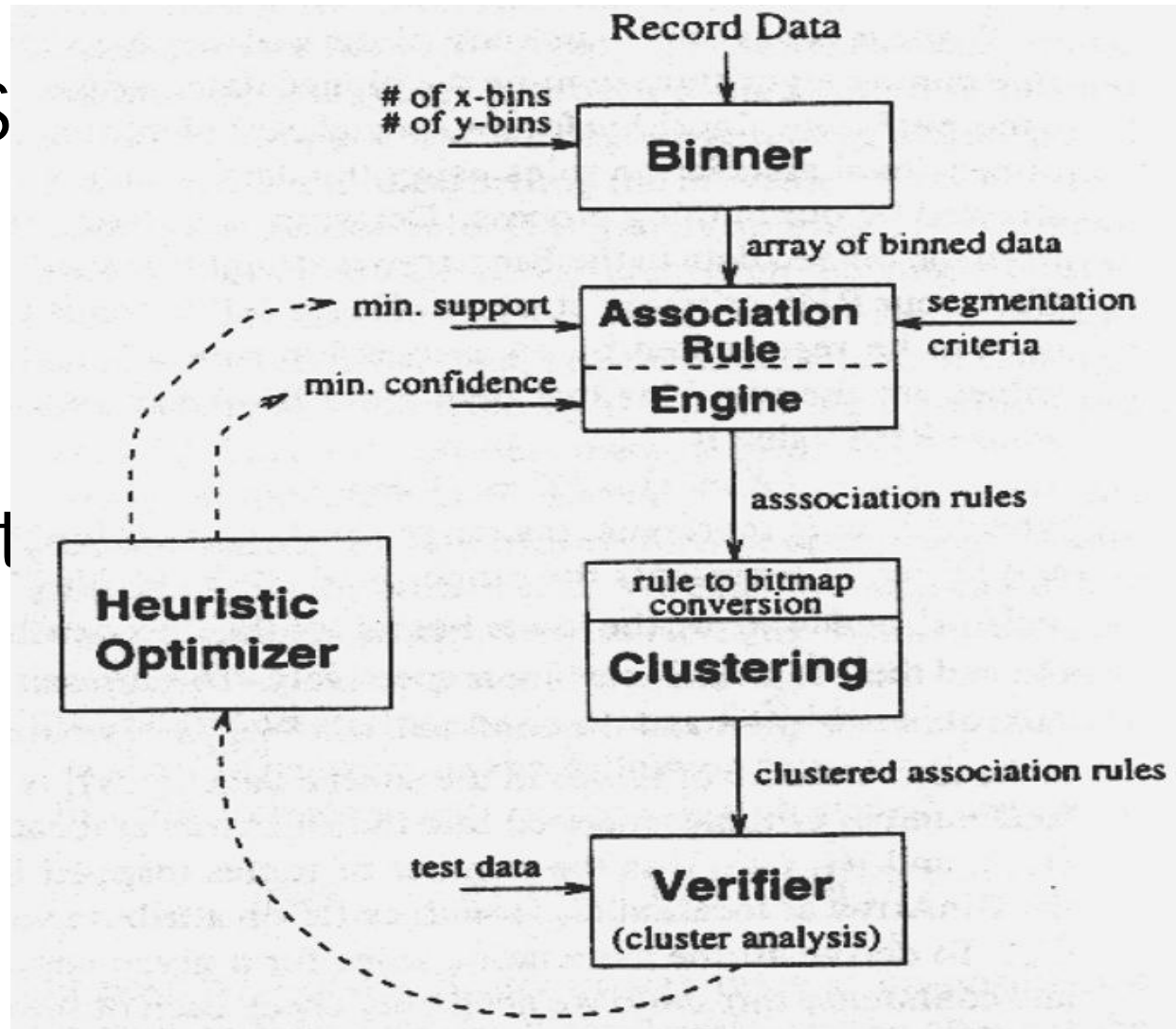
How does ARCS

1. Binning

2. Find frequent predicateset

3. Clustering

4. Optimize



Limitations of ARCS

- Only quantitative attributes on LHS of rules.
- Only 2 attributes on LHS. (2D limitation)
- An alternative to ARCS
 - Non-grid-based
 - equi-depth binning
 - clustering based on a measure of *partial completeness*.
 - "*Mining Quantitative Association Rules in Large Relational Tables*" by R. Srikant and R. Agrawal.

Mining Distance-based Association Rules

- Binning methods do not capture the semantics of interval data

Price(\$)	Equi-width (width \$10)	Equi-depth (depth 2)	Distance- based
7	[0,10]	[7,20]	[7,7]
20	[11,20]	[22,50]	[20,22]
22	[21,30]	[51,53]	[50,53]
50	[31,40]		
51	[41,50]		
53	[51,60]		

- Distance based partitioning, more meaningful discretization considering:
 - density/number of points in an interval
 - “closeness” of points in an interval

Clusters and Distance Measurements

- $S[X]$ is a set of N tuples t_1, t_2, \dots, t_N , projected on the attribute set X
- The diameter of $S[X]$:

$$d(S[X]) = \frac{\sum_{i=1}^N \sum_{j=1}^N dist_x(t_i[X], t_j[X])}{N(N-1)}$$

- $dist_x$: distance metric, e.g. Euclidean distance or Manhattan

Clusters and Distance Measurements(Cont.)

- The diameter, d , assesses the density of a cluster C_X , where

$$d(C_X) \leq d_0^X$$

$$|C_X| \geq s_0$$

- Finding clusters and distance-based rules
 - the density threshold, d_0 , replaces the notion of support
 - modified version of the BIRCH clustering algorithm

Mining Association Rules in Large Databases

- Association rule mining
- Mining single-dimensional Boolean association rules from transactional databases
- Mining multilevel association rules from transactional databases
- Mining multidimensional association rules from transactional databases and data warehouse
- From association mining to correlation analysis
- Constraint-based association mining
- Summary

Interestingness Measurements

- Objective measures
 - Two popular measurements:
 - ★ *support*; and
 - 🕒 *confidence*
- Subjective measures (Silberschatz & Tuzhilin, KDD95)
 - A rule (pattern) is interesting if
 - ★ it is *unexpected* (surprising to the user); and/or
 - 🕒 *actionable* (the user can do something with it)

Criticism to Support and Confidence

- Example 1: (Aggarwal & Yu, PODS98)
 - Among 5000 students
 - 3000 play basketball
 - 3750 eat cereal
 - 2000 both play basket ball and eat cereal
 - *play basketball* \Rightarrow *eat cereal* [40%, 66.7%] is misleading because the overall percentage of students eating cereal is 75% which is higher than 66.7%.
 - *play basketball* \Rightarrow *not eat cereal* [20%, 33.3%] is far more accurate, although with lower support and confidence

	basketball	not basketball	sum(row)
cereal	2000	1750	3750
not cereal	1000	250	1250
sum(col.)	3000	2000	5000

Criticism to Support and Confidence (Cont.)

- Example 2:
 - X and Y: positively correlated,
 - X and Z, negatively related
 - support and confidence of $X \Rightarrow Z$ dominates

X	1	1	1	1	0	0	0	0
Y	1	1	0	0	0	0	0	0
Z	0	1	1	1	1	1	1	1

- We need a measure of dependent or correlated events

$$corr_{A,B} = \frac{P(A \cup B)}{P(A)P(B)}$$

Rule	Support	Confidence
$X \Rightarrow Y$	25%	50%
$X \Rightarrow Z$	37.50%	75%

- $P(B|A)/P(B)$ is also called the **lift** of rule $A \Rightarrow B$

Other Interestingness Measures: Interest

- Interest (correlation, lift)

$$\frac{P(A \wedge B)}{P(A)P(B)}$$
 - taking both $P(A)$ and $P(B)$ in consideration
 - $P(A \wedge B) = P(B) * P(A)$, if A and B are independent events
 - A and B negatively correlated, if the value is less than 1; otherwise A and B positively correlated

X	1	1	1	1	0	0	0	0
Y	1	1	0	0	0	0	0	0
Z	0	1	1	1	1	1	1	1

Itemset	Support	Interest
X,Y	25%	2
X,Z	37.50%	0.9
Y,Z	12.50%	0.57

Mining Association Rules in Large Databases

- Association rule mining
- Mining single-dimensional Boolean association rules from transactional databases
- Mining multilevel association rules from transactional databases
- Mining multidimensional association rules from transactional databases and data warehouse
- From association mining to correlation analysis
- **Constraint-based association mining**
- Summary

Constraint-Based Mining

- Interactive, exploratory mining giga-bytes of data?
 - Could it be real? — Making good use of constraints!
- What kinds of constraints can be used in mining?
 - **Knowledge type constraint**: classification, association, etc.
 - **Data constraint**: SQL-like queries
 - Find product pairs sold together in *Vancouver* in *Dec.'98*.
 - **Dimension/level constraints**:
 - in relevance to *region, price, brand, customer category*.
 - **Rule constraints**
 - small sales ($\text{price} < \$10$) triggers big sales ($\text{sum} > \200).
 - **Interestingness constraints**:
 - strong rules ($\text{min_support} \geq 3\%$, $\text{min_confidence} \geq 60\%$).

Rule Constraints in Association Mining

- Two kind of rule 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).
 - $\text{sum}(\text{LHS}) < 100 \wedge \text{min}(\text{LHS}) > 20 \wedge \text{count}(\text{LHS}) > 3 \wedge \text{sum}(\text{RHS}) > 1000$
- 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).
 - $\text{sum}(\text{LHS}) < \text{min}(\text{RHS}) \wedge \text{max}(\text{RHS}) < 5 * \text{sum}(\text{LHS})$

Constrain-Based Association Query

- Database: (1) trans (TID, Itemset), (2) itemInfo (Item, Type, Price)
- A constrained asso. query (CAQ) is in the form of $\{(S_1, S_2) / C\}$,
 - where C is a set of constraints on S_1, S_2 including frequency constraint
- A classification of (single-variable) constraints:
 - Class constraint: $S \subset A$. *e.g.* $S \subset Item$
 - Domain constraint:
 - $S \theta v, \theta \in \{=, \neq, <, \leq, >, \geq\}$. *e.g.* $S.Price < 100$
 - $v \theta S, \theta$ is \in or \notin . *e.g.* $snacks \notin S.Type$
 - $V \theta S$, or $S \theta V, \theta \in \{\subseteq, \subset, \not\subseteq, =, \neq\}$
 - *e.g.* $\{snacks, sodas\} \subseteq S.Type$
 - Aggregation constraint: $agg(S) \theta v$, where agg is in $\{min, max, sum, count, avg\}$, and $\theta \in \{=, \neq, <, \leq, >, \geq\}$.
 - *e.g.* $count(S_1.Type) = 1, avg(S_2.Price) < 100$

Constrained Association Query Optimization Problem

- Given a CAQ = $\{ (S_1, S_2) / C \}$, the algorithm should be :
 - **sound**: It only finds frequent sets that satisfy the given constraints C
 - **complete**: All frequent sets satisfy the given constraints C are found
- A naïve solution:
 - Apply Apriori for finding all frequent sets, and **then** to test them for constraint satisfaction one by one.
- Our approach:
 - Comprehensive analysis of the properties of constraints and try to **push them as deeply as possible inside** the frequent set computation.

Anti-monotone and Monotone Constraints

- A constraint C_a is **anti-monotone** iff. for any pattern S not satisfying C_a , none of the super-patterns of S can satisfy C_a
- A constraint C_m is **monotone** iff. for any pattern S satisfying C_m , every super-pattern of S also satisfies it

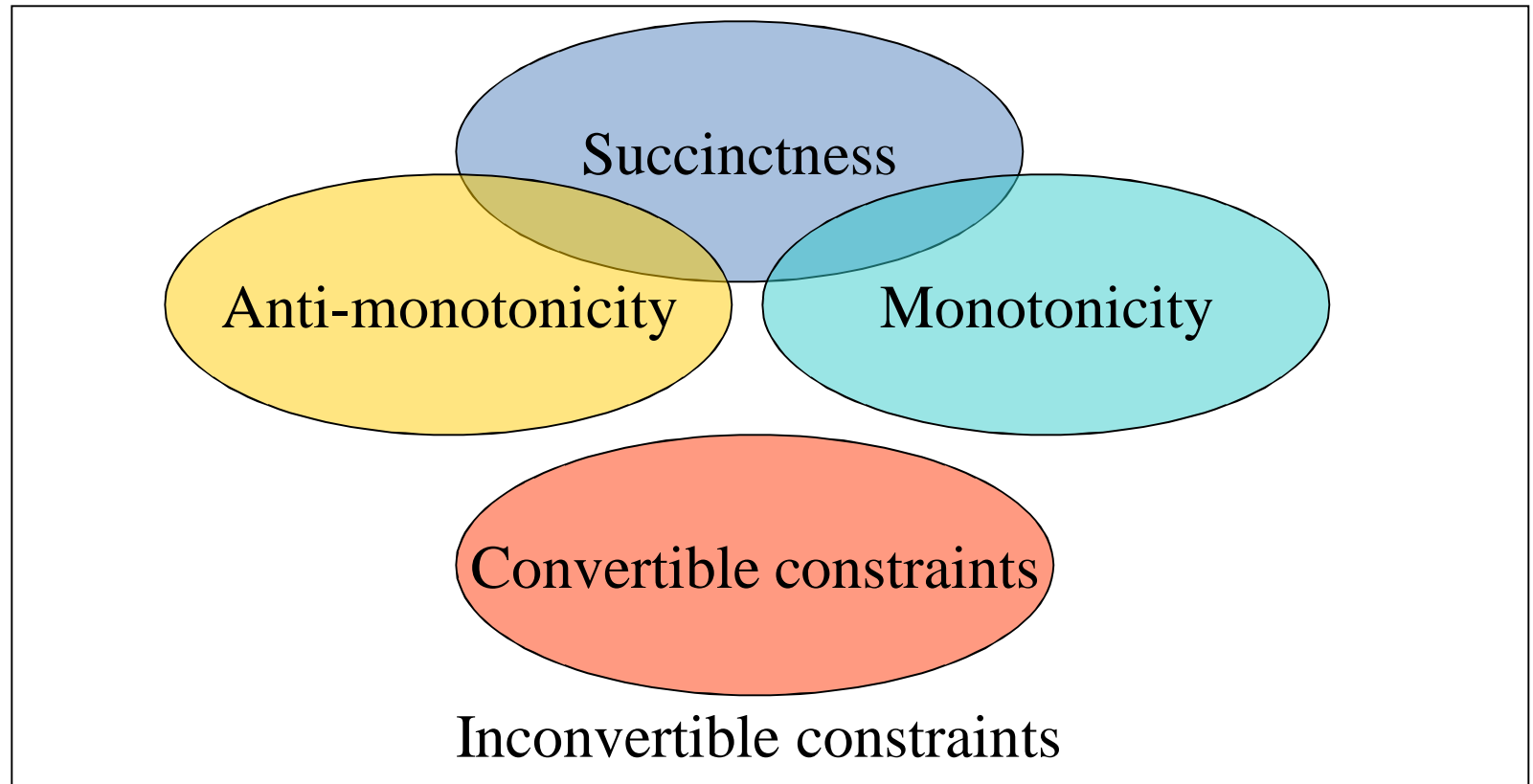
Succinct Constraint

- A subset of item I_s is a **succinct set**, if it can be expressed as $\sigma_p(I)$ for some selection predicate p , where σ is a selection operator
- $SP \subseteq 2^I$ is a succinct **power set**, if there is a fixed number of succinct set $I_1, \dots, I_k \subseteq I$, s.t. SP can be expressed in terms of the strict power sets of I_1, \dots, I_k using union and minus
- A constraint C_s is **succinct** provided $SAT_{C_s}(I)$ is a succinct power set

Convertible Constraint

- Suppose all items in patterns are listed in a total order R
- A constraint C is **convertible anti-monotone** iff a pattern S satisfying the constraint implies that each suffix of S w.r.t. R also satisfies C
- A constraint C is **convertible monotone** iff a pattern S satisfying the constraint implies that each pattern of which S is a suffix w.r.t. R also satisfies C

Relationships Among Categories of Constraints



Property of Constraints: Anti-Monotone

- Anti-monotonicity: *If a set S violates the constraint, any superset of S violates the constraint.*
- Examples:
 - $sum(S.Price) \leq v$ is **anti-monotone**
 - $sum(S.Price) \geq v$ is **not anti-monotone**
 - $sum(S.Price) = v$ is **partly anti-monotone**
- Application:
 - Push " $sum(S.price) \leq 1000$ " deeply into iterative frequent set computation.

Characterization of Anti-Monotonicity Constraints

$S \theta v, \theta \in \{=, \leq, \geq\}$	yes
$v \in S$	no
$S \supseteq V$	no
$S \subseteq V$	yes
$S = V$	partly
$\min(S) \leq v$	no
$\min(S) \geq v$	yes
$\min(S) = v$	partly
$\max(S) \leq v$	yes
$\max(S) \geq v$	no
$\max(S) = v$	partly
$\text{count}(S) \leq v$	yes
$\text{count}(S) \geq v$	no
$\text{count}(S) = v$	partly
$\text{sum}(S) \leq v$	yes
$\text{sum}(S) \geq v$	no
$\text{sum}(S) = v$	partly
$\text{avg}(S) \theta v, \theta \in \{=, \leq, \geq\}$	convertible
(frequent constraint)	(yes)

Example of Convertible Constraints: $\text{Avg}(S) \geq V$

- Let R be the value descending order over the set of items
 - E.g. $I = \{9, 8, 6, 4, 3, 1\}$
- $\text{Avg}(S) \geq v$ is convertible monotone w.r.t. R
 - If S is a suffix of S_1 , $\text{avg}(S_1) \geq \text{avg}(S)$
 - $\{8, 4, 3\}$ is a suffix of $\{9, 8, 4, 3\}$
 - $\text{avg}(\{9, 8, 4, 3\}) = 6 \geq \text{avg}(\{8, 4, 3\}) = 5$
 - If S satisfies $\text{avg}(S) \geq v$, so does S_1
 - $\{8, 4, 3\}$ satisfies constraint $\text{avg}(S) \geq 4$, so does $\{9, 8, 4, 3\}$

Property of Constraints: Succinctness

- Succinctness:
 - For any set S_1 and S_2 satisfying C , $S_1 \cup S_2$ satisfies C
 - Given A_1 is the sets of size 1 satisfying C , then any set S satisfying C are based on A_1 , i.e., it contains a subset belongs to A_1 ,
- Example :
 - $sum(S.Price) \geq v$ is not succinct
 - $min(S.Price) \leq v$ is succinct
- Optimization:
 - If C is succinct, then C is pre-counting prunable. The satisfaction of the constraint alone is not affected by the iterative support counting.

Characterization of Constraints by Succinctness

$S \theta v, \theta \in \{=, \leq, \geq\}$	Yes
$v \in S$	yes
$S \supseteq V$	yes
$S \subseteq V$	yes
$S = V$	yes
$\min(S) \leq v$	yes
$\min(S) \geq v$	yes
$\min(S) = v$	yes
$\max(S) \leq v$	yes
$\max(S) \geq v$	yes
$\max(S) = v$	yes
$\text{count}(S) \leq v$	weakly
$\text{count}(S) \geq v$	weakly
$\text{count}(S) = v$	weakly
$\text{sum}(S) \leq v$	no
$\text{sum}(S) \geq v$	no
$\text{sum}(S) = v$	no
$\text{avg}(S) \theta v, \theta \in \{=, \leq, \geq\}$	no
(frequent constraint)	(no)

Mining Association Rules in Large Databases

- Association rule mining
- Mining single-dimensional Boolean association rules from transactional databases
- Mining multilevel association rules from transactional databases
- Mining multidimensional association rules from transactional databases and data warehouse
- From association mining to correlation analysis
- Constraint-based association mining
- [Summary](#)

Why Is the Big Pie Still There?

- More on constraint-based mining of associations
 - Boolean vs. quantitative associations
 - Association on discrete vs. continuous data
 - From association to correlation and causal structure analysis.
 - Association does not necessarily imply correlation or causal relationships
 - From intra-transaction association to inter-transaction associations
 - E.g., break the barriers of transactions (Lu, et al. TOIS'99).
 - From association analysis to classification and clustering analysis
 - E.g, clustering association rules

Mining Association Rules in Large Databases

- [Association rule mining](#)
- Mining single-dimensional Boolean association rules from transactional databases
- Mining multilevel association rules from transactional databases
- Mining multidimensional association rules from transactional databases and data warehouse
- From association mining to correlation analysis
- Constraint-based association mining
- Summary

Summary

- Association rule mining
 - probably the most significant contribution from the database community in 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.

References

- R. Agarwal, C. Aggarwal, and V. V. V. Prasad. A tree projection algorithm for generation of frequent itemsets. In *Journal of Parallel and Distributed Computing (Special Issue on High Performance Data Mining)*, 2000.
- 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.
- S. Brin, R. Motwani, J. D. Ullman, and S. Tsur. Dynamic itemset counting and implication rules for market basket analysis. *SIGMOD'97*, 255-264, Tucson, Arizona, May 1997.
- K. Beyer and R. Ramakrishnan. Bottom-up computation of sparse and iceberg cubes. *SIGMOD'99*, 359-370, Philadelphia, PA, June 1999.
- 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.
- M. Fang, N. Shivakumar, H. Garcia-Molina, R. Motwani, and J. D. Ullman. Computing iceberg queries efficiently. *VLDB'98*, 299-310, New York, NY, Aug. 1998.

References (2)

- G. Grahne, L. Lakshmanan, and X. Wang. Efficient mining of constrained correlated sets. ICDE'00, 512-521, San Diego, CA, Feb. 2000.
- Y. Fu and J. Han. Meta-rule-guided mining of association rules in relational databases. KDOOD'95, 39-46, Singapore, Dec. 1995.
- 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, G. Dong, and Y. Yin. Efficient mining of partial periodic patterns in time series database. ICDE'99, Sydney, Australia.
- J. Han and Y. Fu. Discovery of multiple-level association rules from large databases. VLDB'95, 420-431, Zurich, Switzerland.
- J. Han, J. Pei, and Y. Yin. Mining frequent patterns without candidate generation. SIGMOD'00, 1-12, Dallas, TX, May 2000.
- T. Imielinski and H. Mannila. A database perspective on knowledge discovery. Communications of ACM, 39:58-64, 1996.
- 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.

References (3)

- F. Korn, A. Labrinidis, Y. Kotidis, and C. Faloutsos. Ratio rules: A new paradigm for fast, quantifiable data mining. VLDB'98, 582-593, New York, NY.
- B. Lent, A. Swami, and J. Widom. Clustering association rules. ICDE'97, 220-231, Birmingham, England.
- H. Lu, J. Han, and L. Feng. Stock movement and n-dimensional inter-transaction association rules. SIGMOD Workshop on Research Issues on Data Mining and Knowledge Discovery (DMKD'98), 12:1-12:7, Seattle, Washington.
- H. Mannila, H. Toivonen, and A. I. Verkamo. Efficient algorithms for discovering association rules. KDD'94, 181-192, Seattle, WA, July 1994.
- H. Mannila, H Toivonen, and A. I. Verkamo. Discovery of frequent episodes in event sequences. Data Mining and Knowledge Discovery, 1:259-289, 1997.
- R. Meo, G. Psaila, and S. Ceri. A new SQL-like operator for mining association rules. VLDB'96, 122-133, Bombay, India.
- R.J. Miller and Y. Yang. Association rules over interval data. SIGMOD'97, 452-461, Tucson, Arizona.
- 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.
- N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Discovering frequent closed itemsets for association rules. ICDT'99, 398-416, Jerusalem, Israel, Jan. 1999.

References (4)

- 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, May 1995.
- J. Pei, J. Han, and R. Mao. CLOSET: An Efficient Algorithm for Mining Frequent Closed Itemsets. DMKD'00, Dallas, TX, 11-20, May 2000.
- J. Pei and J. Han. Can We Push More Constraints into Frequent Pattern Mining? KDD'00. Boston, MA. Aug. 2000.
- G. Piatetsky-Shapiro. Discovery, analysis, and presentation of strong rules. In G. Piatetsky-Shapiro and W. J. Frawley, editors, Knowledge Discovery in Databases, 229-238. AAAI/MIT Press, 1991.
- 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.
- A. Savasere, E. Omiecinski, and S. Navathe. An efficient algorithm for mining association rules in large databases. VLDB'95, 432-443, Zurich, Switzerland.
- A. Savasere, E. Omiecinski, and S. Navathe. Mining for strong negative associations in a large database of customer transactions. ICDE'98, 494-502, Orlando, FL, Feb. 1998.

References (5)

- 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, Sept. 1995.
- 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.
- H. Toivonen. Sampling large databases for association rules. VLDB'96, 134-145, Bombay, India, Sept. 1996.
- 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.
- K. Yoda, T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama. Computing optimized rectilinear regions for association rules. KDD'97, 96-103, Newport Beach, CA, Aug. 1997.
- M. J. Zaki, S. Parthasarathy, M. Ogihara, and W. Li. Parallel algorithm for discovery of association rules. Data Mining and Knowledge Discovery, 1:343-374, 1997.
- M. Zaki. Generating Non-Redundant Association Rules. KDD'00. Boston, MA. Aug. 2000.
- O. R. Zaiane, J. Han, and H. Zhu. Mining Recurrent Items in Multimedia with Progressive Resolution Refinement. ICDE'00, 461-470, San Diego, CA, Feb. 2000.