

---

Università degli Studi di Milano  
Master Degree in Computer Science

# Information Management course

Teacher: Alberto Ceselli

Lecture 16: 25/11/2014

# Data Mining:

---

## Concepts and Techniques

(3<sup>rd</sup> ed.)

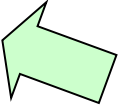
### — Chapter 6 —

Jiawei Han, Micheline Kamber, and Jian Pei  
University of Illinois at Urbana-Champaign &  
Simon Fraser University

©2011 Han, Kamber & Pei. All rights reserved.

# Scalable Frequent Itemset Mining Methods

---

- Apriori: A Candidate Generation-and-Test Approach
- Improving the Efficiency of Apriori
- FPGrowth: A Frequent Pattern-Growth Approach 
- ECLAT: Frequent Pattern Mining with Vertical Data Format
- Mining Close Frequent Patterns and Maxpatterns

# Pattern-Growth Approach: Mining Frequent Patterns Without Candidate Generation

---

- Bottlenecks of the Apriori approach
  - Breadth-first (i.e., level-wise) search
  - Candidate generation and test
    - Often generates a huge number of candidates
- The FPGrowth Approach (J. Han, J. Pei, and Y. Yin, SIGMOD' 00)
  - Depth-first search
  - Avoid explicit candidate generation
- Major philosophy: Grow long patterns from short ones using local frequent items only
  - Divide et impera
    - Compress the DB using a tree data structure
    - Divide the compressed DB in *conditional DBs* and proceed recursively

# Construct FP-tree from a Transaction Database, then **mine the tree**

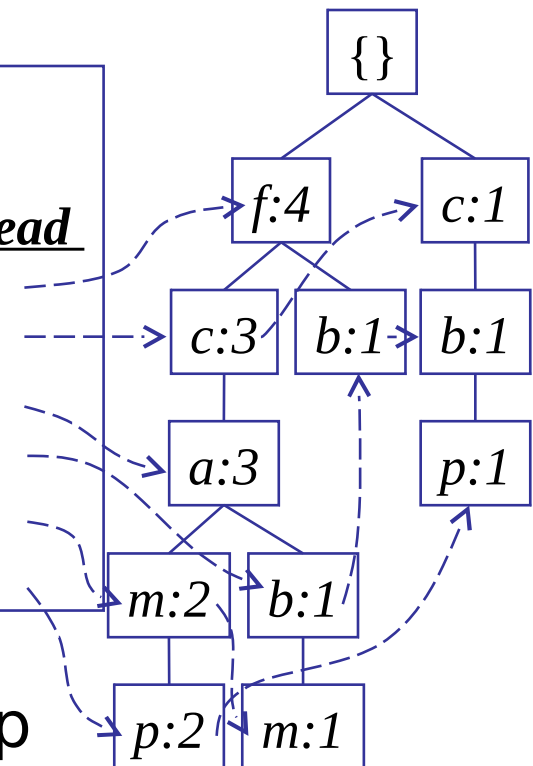
<u>TID</u>	<u>Items bought</u>	<u>(ordered) frequent items</u>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o, w}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

$min\_support = 3$

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

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

**F-list** = f-c-a-b-m-p



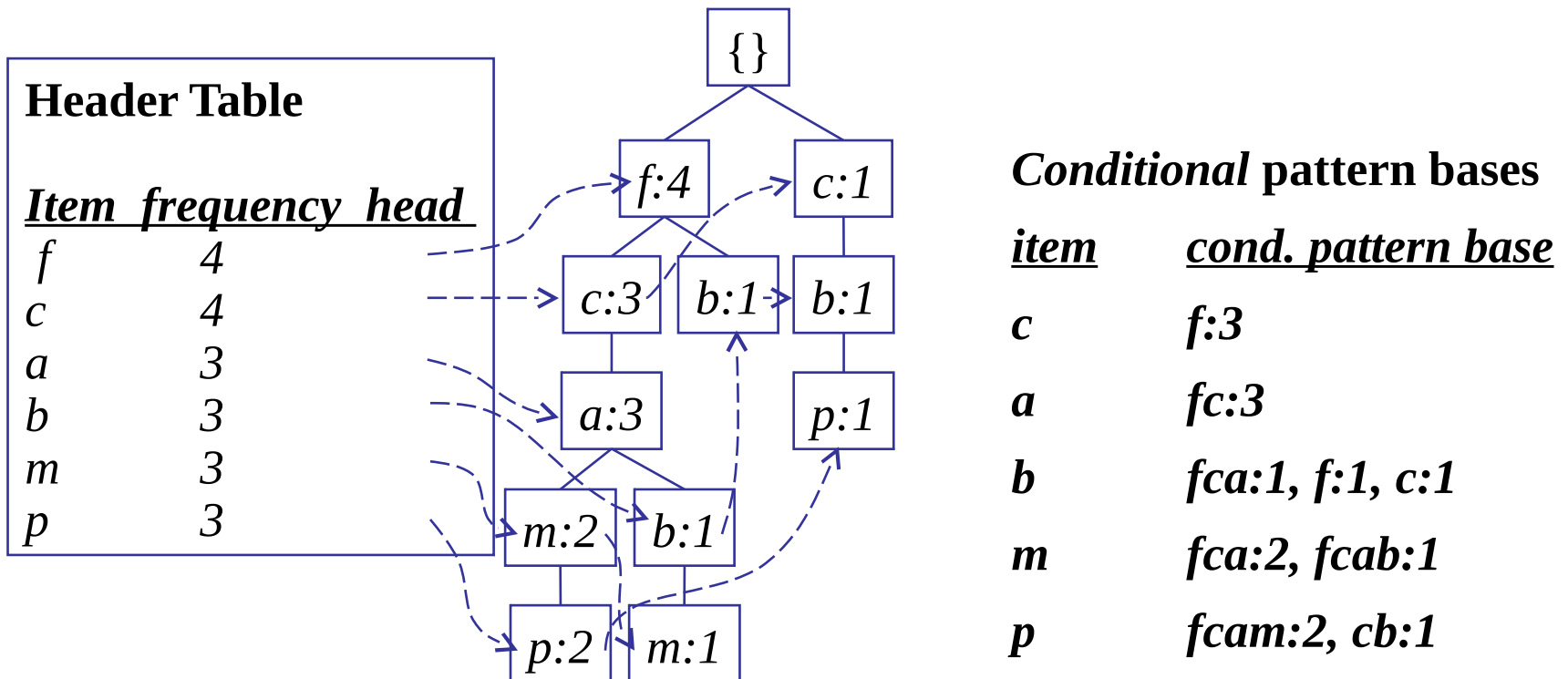
# Partition Patterns and Databases

---

- Start with each frequent length-1 pattern  
→ as an initial *suffix pattern*
- Build its conditional pattern base  
→ a sub-DB containing all its prefix paths in the FP-tree
- The conditional pattern base is a DB  
→ build its (conditional) FP-tree
- Apply FP growth recursively  
(stop with empty support)
- Append to the mining results the suffix pattern

# Find Patterns Having x From x-conditional Database

- Starting at the frequent item header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item  $x$
- Accumulate all of *transformed prefix paths* of item  $x$  to form  $x$ 's conditional pattern base



# Starting with least frequent items is better!

---

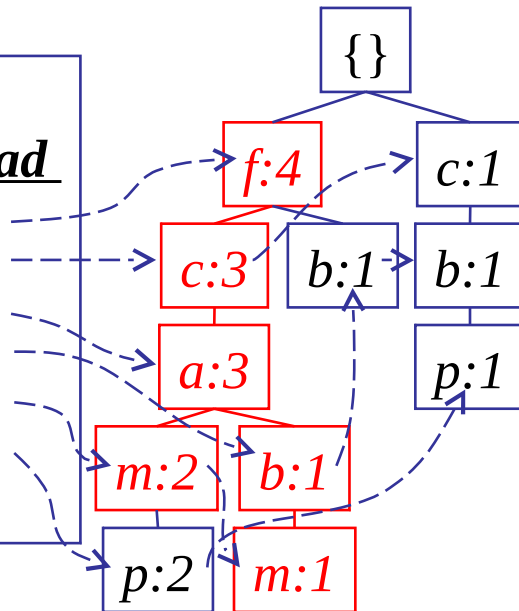
- Frequent patterns can be partitioned into subsets according to f-list (Completeness and non-redundancy)
  - F-list = f-c-a-b-m-p
  - Patterns containing p
  - Patterns having m but no p
  - ...
  - Patterns having c but no a nor b, m, p
  - Pattern f
- Least frequent items are leaves of the FP-tree (once processed can be removed ...)



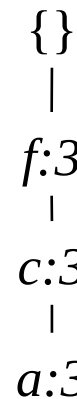
# From Conditional Pattern-bases to Conditional FP-trees

- 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

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



*m*-conditional pattern base:  
*fca:2, fcab:1*



All frequent patterns related to *m*

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

*m*-conditional FP-tree

# Recursion: Mining Each Conditional FP-tree

{}

|

*f*:3

|

*c*:3

|

*a*:3

*m*-conditional FP-tree

Cond. pattern base of “am”: (fc:3)

{}

|

*f*:3

|

*c*:3

*am*-conditional FP-tree

Cond. pattern base of “cm”: (f:3)

{}

|

*f*:3

*cm*-conditional FP-tree

Cond. pattern base of “cam”: (f:3)

{}

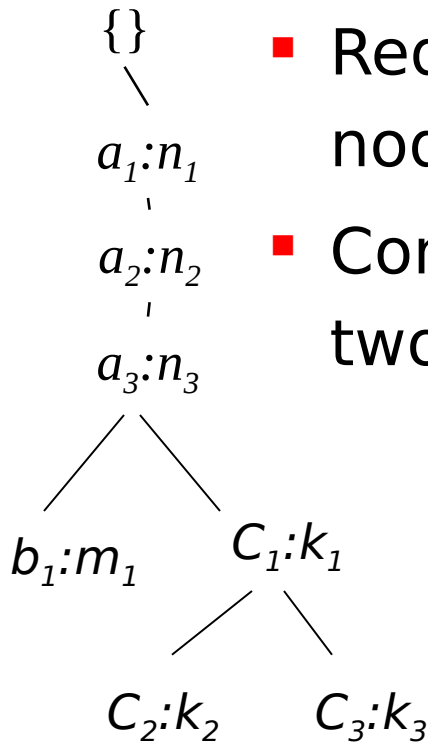
|

*f*:3

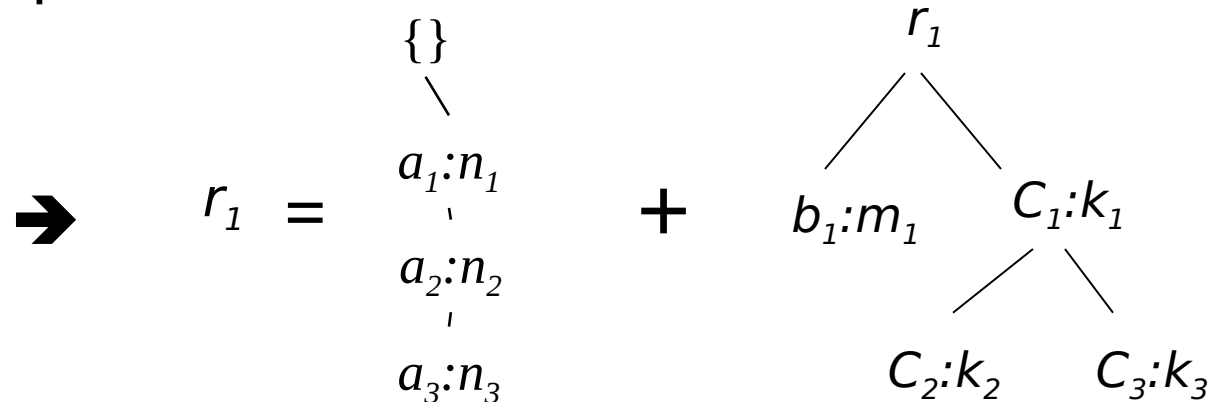
*cam*-conditional FP-tree

# A Special Case: Single Prefix Path in FP-tree

- Suppose a (conditional) FP-tree  $T$  has a shared single prefix-path  $P$
- Mining can be decomposed into two parts



- Reduction of the single prefix path into one node
- Concatenation of the mining results of the two parts



# Benefits of the FP-tree Structure

---

- Completeness
  - Preserve complete information for frequent pattern mining
  - Never break a long pattern of any transaction
- Compactness
  - Reduce irrelevant info—infrequent items are gone
  - Items in frequency descending order: the more frequently occurring, the more likely to be shared
  - Never be larger than the original database (not count node-links and the *count* field)

# The Frequent Pattern Growth Mining Method (summary)

---

- Idea: Frequent pattern growth
  - Recursively grow frequent patterns by pattern and database partition
- Method
  - For each frequent 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
  - Generate all the combinations of sub-paths from this single path: each of them is a frequent pattern

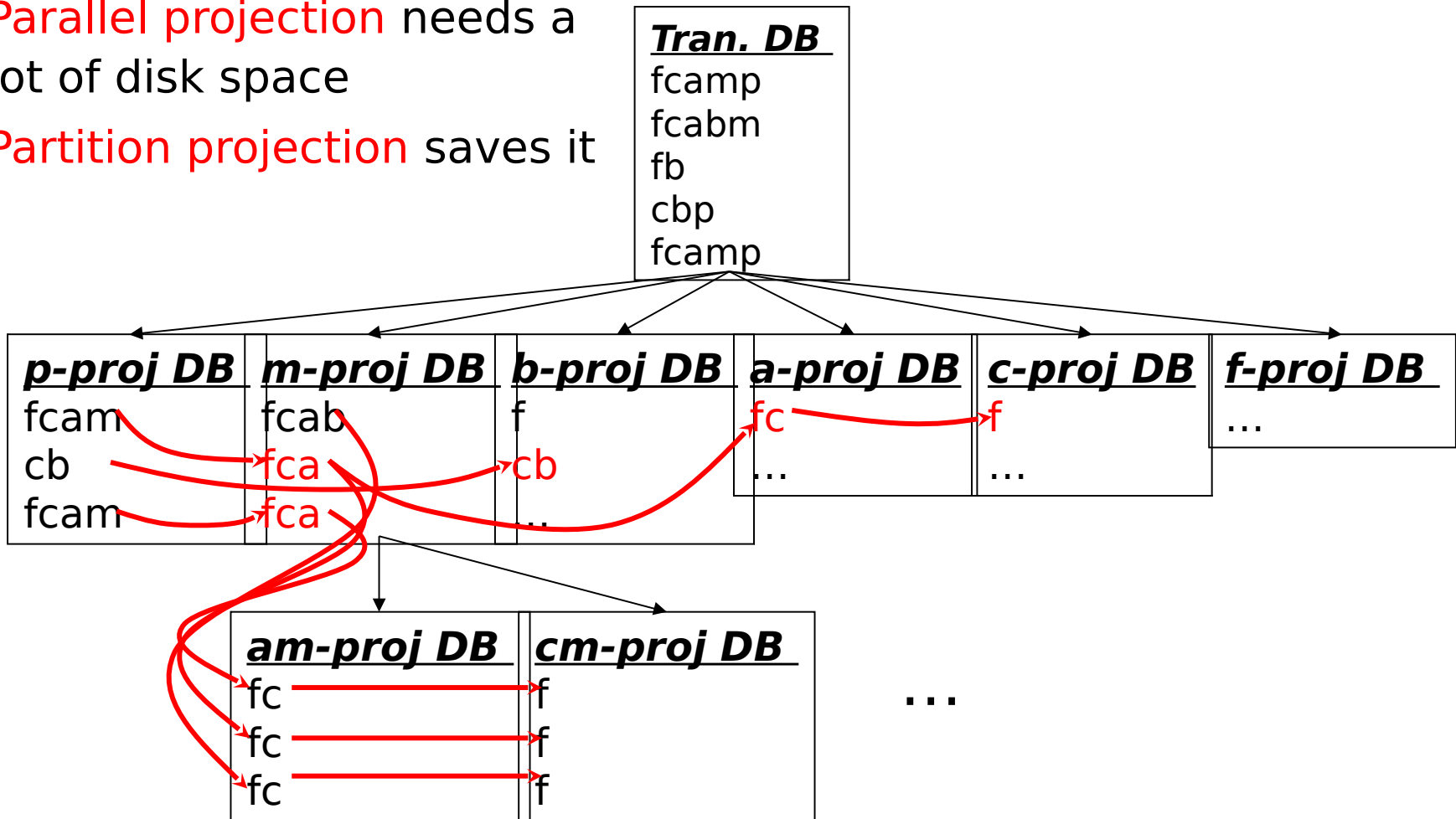
# Scaling FP-growth by Database Projection

---

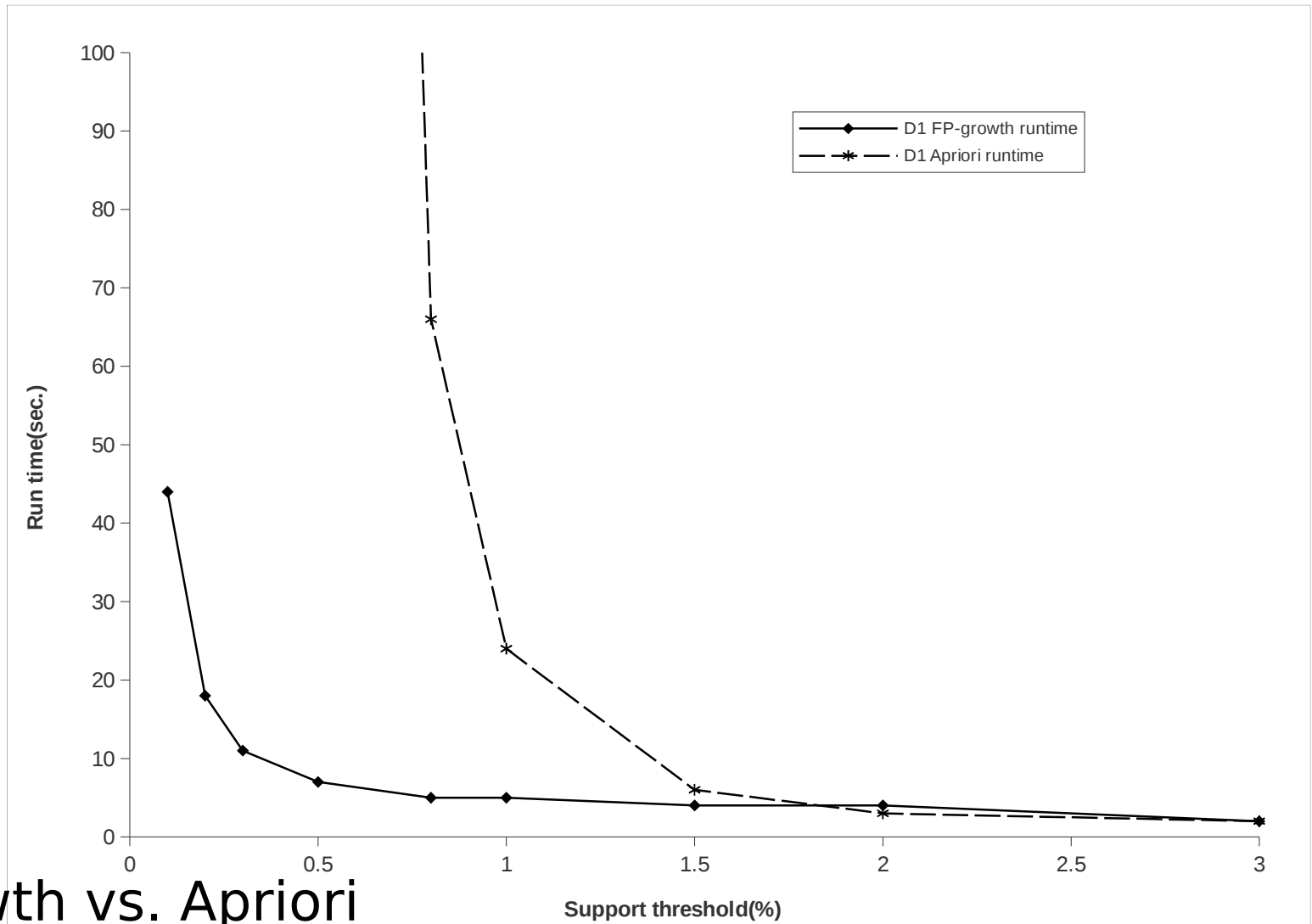
- What about if FP-tree cannot fit in memory?
  - DB projection
- First partition a database into a set of projected DBs
- Then construct and mine FP-tree for each projected DB
- **Parallel projection** vs. **partition projection** techniques
  - Parallel projection
    - Project the DB in parallel for each frequent item
    - Parallel projection is space costly
    - All the partitions can be processed in parallel
  - Partition projection
    - Partition the DB based on the ordered frequent items
    - Passing the unprocessed parts to the subsequent partitions

# Partition-Based Projection

- **Parallel projection** needs a lot of disk space
- **Partition projection** saves it



# Performance of FPGrowth in Large Datasets



FP-Growth vs. Apriori



# Advantages of the Pattern Growth Approach

---

- Divide-and-conquer:
  - Decompose both the mining task and DB according to the frequent patterns obtained so far
  - Lead to focused search of smaller databases
- Other factors
  - No candidate generation, no candidate test
  - Compressed database: FP-tree structure
  - No repeated scan of entire database
  - Basic ops: counting local freq items and building sub FP-tree, no pattern search and matching
- A good open-source implementation and refinement of FPGrowth
  - FPGrowth+ (Grahne and J. Zhu, FIMI'03)

# Further Improvements of Mining Methods

---

- AFOPT (Liu, et al. @ KDD'03)
  - A “push-right” method for mining condensed frequent pattern (CFP) tree
- Carpenter (Pan, et al. @ KDD'03)
  - Mine data sets with small rows but numerous columns
  - Construct a row-enumeration tree for efficient mining
- FPgrowth+ (Grahne and Zhu, FIMI'03)
  - Efficiently Using Prefix-Trees in Mining Frequent Itemsets, Proc. ICDM'03 Int. Workshop on Frequent Itemset Mining Implementations (FIMI'03), Melbourne, FL, Nov. 2003
- TD-Close (Liu, et al, SDM'06)

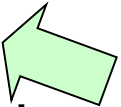
# Extension of Pattern Growth Mining Methodology

---

- Mining closed frequent itemsets and max-patterns
  - CLOSET (DMKD'00), FPclose, and FPMMax (Fimi'03)
- Mining sequential patterns
  - PrefixSpan (ICDE'01), CloSpan (SDM'03), BIDE (ICDE'04)
- Mining graph patterns
  - gSpan (ICDM'02), CloseGraph (KDD'03)
- Constraint-based mining of frequent patterns
  - Convertible constraints (ICDE'01), gPrune (PAKDD'03)
- Computing iceberg data cubes with complex measures
  - H-tree, H-cubing, and Star-cubing (SIGMOD'01, VLDB'03)
- Pattern-growth-based Clustering
  - MaPle (Pei, et al., ICDM'03)
- Pattern-Growth-Based Classification
  - Mining frequent and discriminative patterns (Cheng, et al, ICDE'07)

# Scalable Frequent Itemset Mining Methods

---

- Apriori: A Candidate Generation-and-Test Approach
- Improving the Efficiency of Apriori
- FPGrowth: A Frequent Pattern-Growth Approach
- ECLAT: Frequent Pattern Mining with Vertical Data Format
- Mining Close Frequent Patterns and Maxpatterns 

# Equivalence CLASS Transformation: Mining by Exploring Vertical Data Format

---

- Vertical format:  $t(ab) = \{T_{11}, T_{25}, \dots\}$ 
  - tid-list: list of trans.-ids containing an itemset
- Deriving frequent patterns based on vertical intersections
  - $t(X) = t(Y)$ : X and Y always happen together
  - $t(X) \subset t(Y)$ : transaction having X always has Y
- Using **diffset** to accelerate mining
  - Only keep track of differences of tids
  - $t(X) = \{T_1, T_2, T_3\}$ ,  $t(XY) = \{T_1, T_3\}$
  - $\text{Diffset}(XY, X) = \{T_2\}$
- ECLAT (Zaki et al. @KDD'97)
- Mining Closed patterns using vertical format: CHARM (Zaki & Hsiao@SDM'02)

# Equivalence CLAss Transformation: Mining by Exploring Vertical Data Format

- DB:

T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

min\_supp = 2

I1 T100, T400, T500, T700, T800, T900

I2 T100, T200, T300, T400, T600, T800, T900

I3 T300, T500, T600, T700, T800, T900

I4 T200, T400

I5 T100, T800

I1, I2 T100, T400, T800, T900

I1, I3 T500, T700, T800, T900

I1, I4 T400

I1, I5 T100, T800

I2, I3 T300, T600, T800, T900

I2, I4 T200, T400

I2, I5 T100, T800

I3, I5 T800

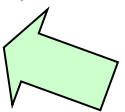
I1, I2, I3 T800, T900

I1, I2, I5 T100, T800

# Scalable Frequent Itemset Mining Methods

---

- Apriori: A Candidate Generation-and-Test Approach
- Improving the Efficiency of Apriori
- FPGrowth: A Frequent Pattern-Growth Approach
- ECLAT: Frequent Pattern Mining with Vertical Data Format
- Mining Close Frequent Patterns and Maximal patterns



# Closed Patterns and Max-Patterns

---

- *closed* and *maximal* patterns
  - An itemset  $X$  is a **closed p.** if  $X$  is *frequent* and there exists *no super-pattern*  $Y \supset X$ , with the *same support* as  $X$
  - An itemset  $X$  is a **maximal p.** if  $X$  is *frequent* and there exists *no super-pattern*  $Y \supset X$ , which is *also frequent*



# Mining Frequent Closed Patterns: CLOSET

- Flist: list of all frequent items in support asc. order
  - Flist: d-a-f-e-c
- Divide search space
  - Patterns having d
  - Patterns having d but no a, etc.
- Find frequent closed pattern recursively
  - Every transaction having d also has *cfa* → *cfad* is a frequent closed pattern
- J. Pei, J. Han & R. Mao. "CLOSET: An Efficient Algorithm for Mining Frequent Closed Itemsets" (2000)

Min\_sup=2

TID	Items
10	a, c, d, e, f
20	a, b, e
30	c, e, f
40	a, c, d, f
50	c, e, f

# CLOSET+: Mining Closed Itemsets by Pattern-Growth

---

- Itemset merging: if every occurrence of  $X$  contains  $Y$ , but not every proper subset of  $Y$ , then  $X \cup Y$  is a frequent closed itemset (no need to search for itemsets containing  $X$  but not  $Y$ )
- Sub-itemset pruning: if  $Y \supset X$ , and  $\text{sup}(X) = \text{sup}(Y)$ ,  $X$  and all of  $X$ 's descendants in the set enumeration tree can be pruned
- Item skipping: if a local frequent item has the same support in several header tables at different levels, one can prune it from the header table at higher levels
- Efficient subset checking
- Hybrid tree projection
  - bottomup physical / topdown pseudo tree-projection

# MaxMiner: Mining Max-Patterns

- 1<sup>st</sup> scan: find frequent items

- A, B, C, D, E

- 2<sup>nd</sup> scan: find support for

- AB, AC, AD, AE, ABCDE

- BC, BD, BE, BCDE

- CD, CE, CDE

- DE

Tid	Items
10	A, B, C, D, E
20	B, C, D, E,
30	A, C, D, F

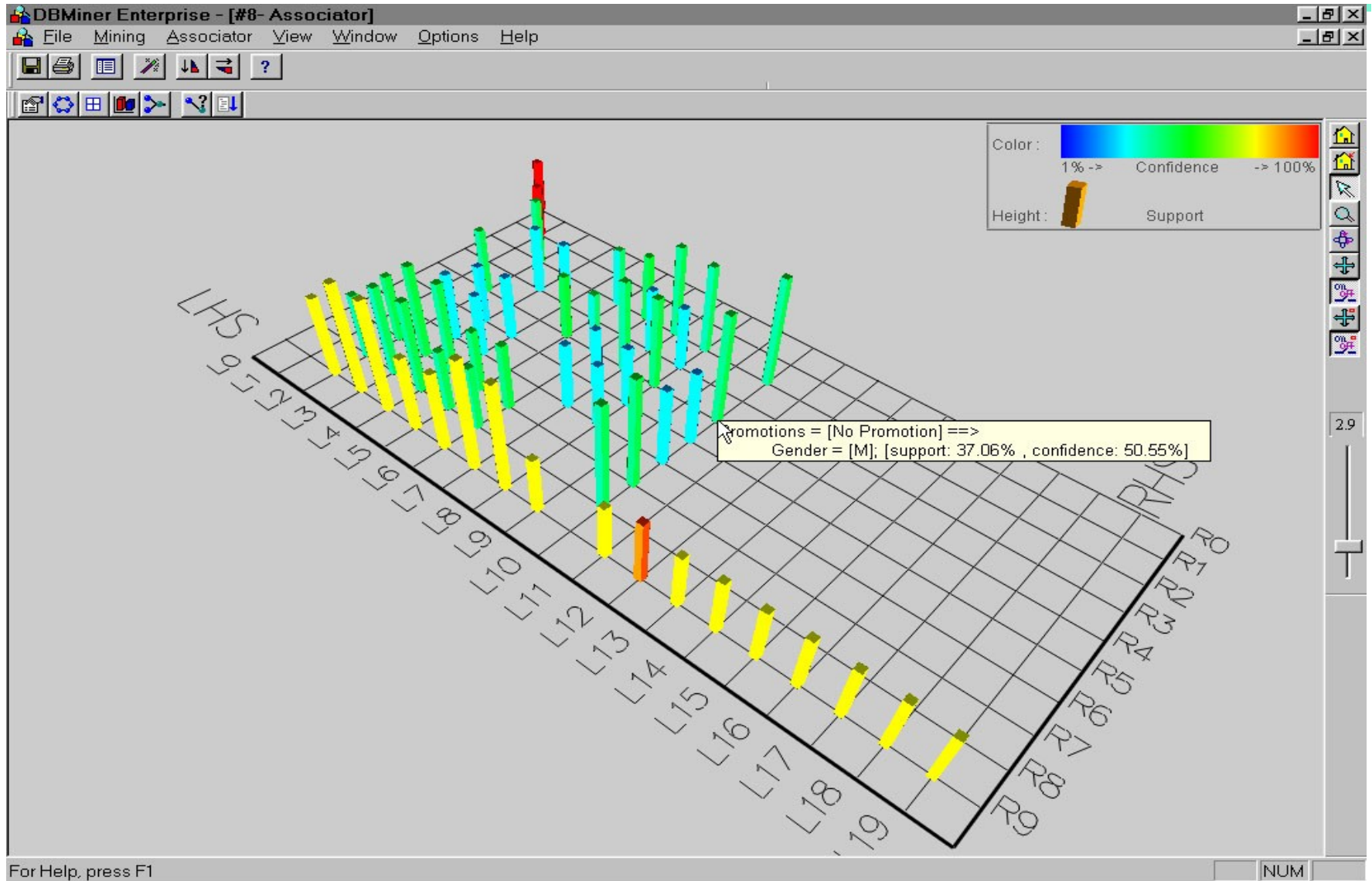
Potential  
max-  
patterns



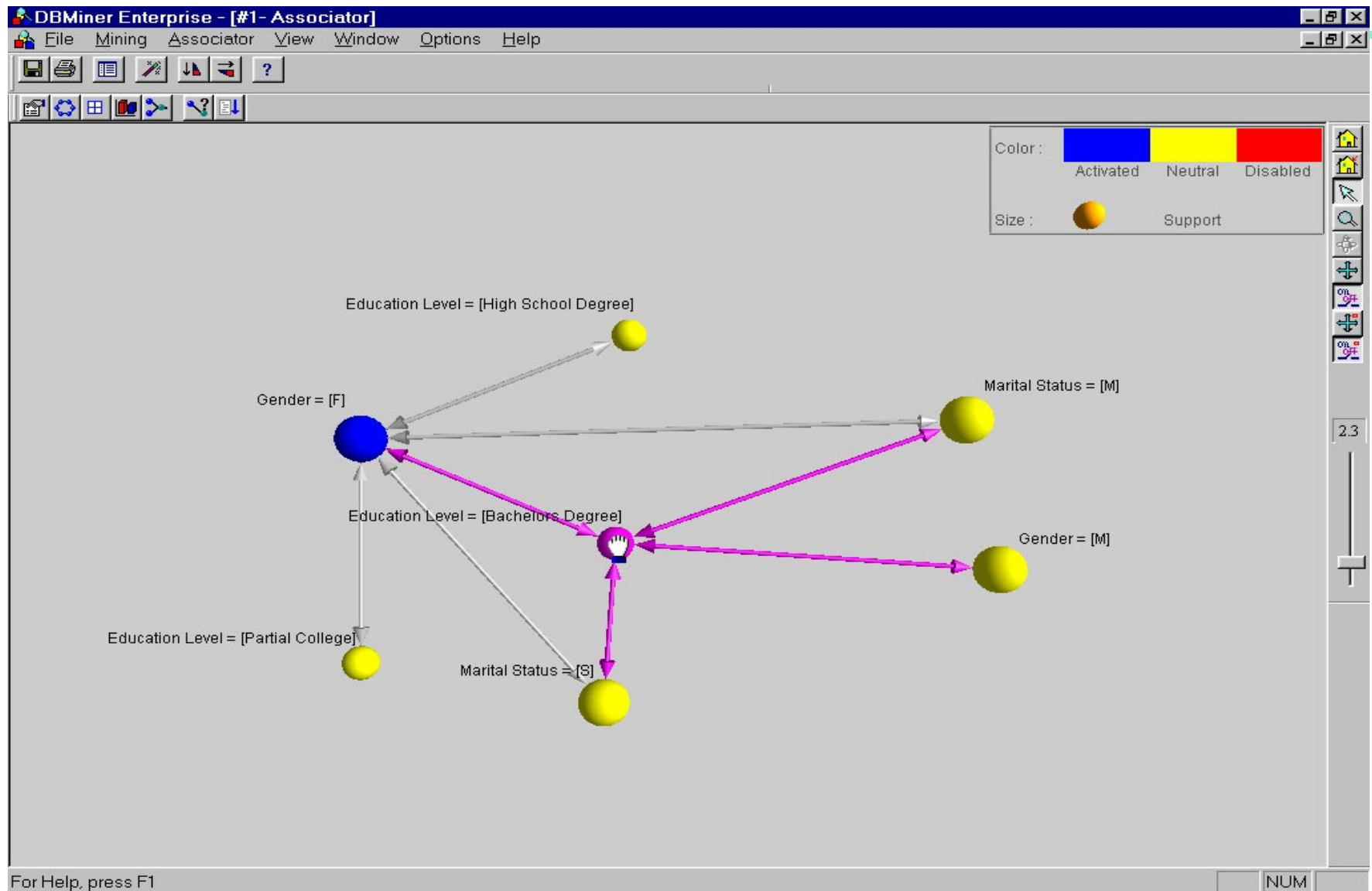
- Since BCDE is a max-pattern, no need to check BCD, BDE, CDE in later scan

- R. Bayardo. Efficiently mining long patterns from databases. *SIGMOD'98*

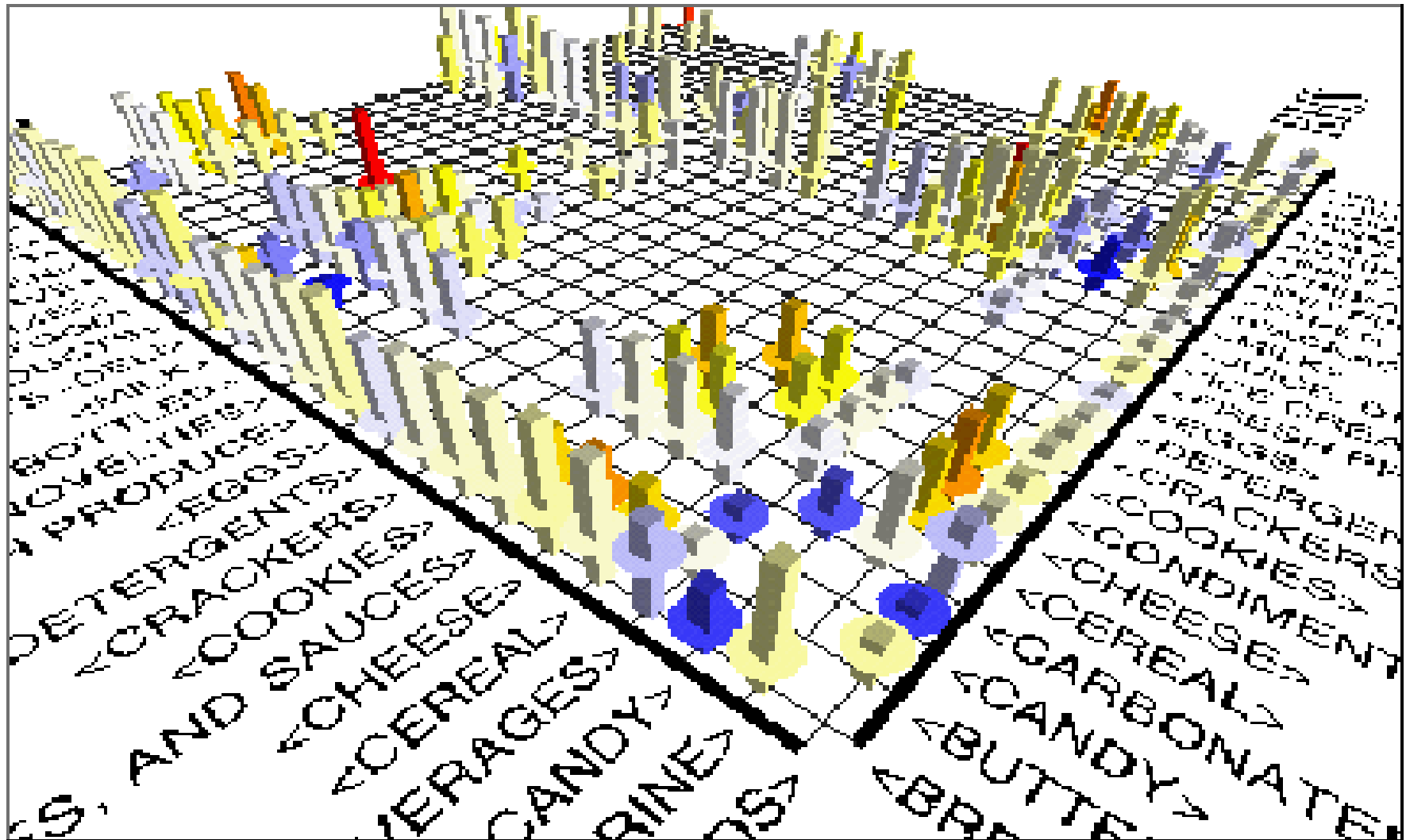
# Visualization of Association Rules: Plane Graph



# Visualization of Association Rules: Rule Graph



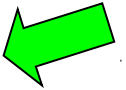
# Visualization of Association Rules (SGI/MineSet 3.0)



# Chapter 5: Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

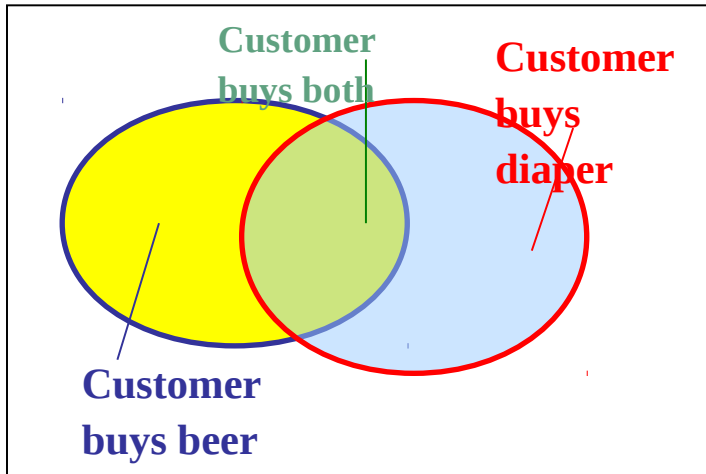
---

- Basic Concepts
- Frequent Itemset Mining Methods
- Which Patterns Are Interesting?—Pattern Evaluation Methods
- Summary



# Basic Concepts: Association Rules

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk



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

Let  $minsup = 50\%$ ,  $minconf = 50\%$   
 Freq. Pat.: Beer:3, Nuts:3, Diaper:4,  
 Eggs:3, {Beer, Diaper}:3

Association rules: (many more!)
 

- $Beer \rightarrow Diaper$  (60%, 100%)
- $Diaper \rightarrow Beer$  (60%, 75%)



# Interestingness Measure: Correlations (Lift)

- *play basketball*  $\Rightarrow$  *eat cereal* [40%, 66.7%] is misleading
  - The overall % of students eating cereal is 75% > 66.7%.
- *play basketball*  $\Rightarrow$  *not eat cereal* [20%, 33.3%] is more accurate, although with lower support and confidence
- Idea: put rules as “A  $\rightarrow$  B [support, confidence, correlation]”
- Measure of dependent/correlated events: **lift**

$$\text{lift} = \frac{P(A \cup B)}{P(A)P(B)} = \frac{P(B|A)}{P(B)}$$

$$\text{lift}(B, C) = \frac{2000/5000}{3000/5000 * 3750/5000} = 0.89$$

$$\text{lift}(B, \neg C) = \frac{1000/5000}{3000/5000 * 1250/5000} = 1.33$$

	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

- **N.B. A and B are independent if  $P(A \cup B) = P(A) * P(B)$**

# Interestingness Measure: Correlations (Chi-Square)

---

- Chi-square test:

$$\sum (\text{observed} - \text{expected})^2 / \text{expected}$$

and compare to tables

# Are *lift* and $\chi^2$ Good Measures of Correlation?

- “Buy walnuts  $\Rightarrow$  buy milk [1%, 80%]” is misleading if 85% of customers buy milk
- Support and confidence are not enough to indicate correlations
- Over 20 interestingness measures have been proposed (see Tan et al @KDD’02)
- Which are good ones?

symbol	measure	range	formula
$\phi$	$\phi$ -coefficient	-1 ... 1	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
$Q$	Yule’s Q	-1 ... 1	$\frac{P(A,B)P(\bar{A},\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A},\bar{B}) + P(A,\bar{B})P(\bar{A},B)}$
$Y$	Yule’s Y	-1 ... 1	$\frac{\sqrt{P(A,B)P(\bar{A},\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A},\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}}$
$k$	Cohen’s	-1 ... 1	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
$PS$	Piatetsky-Shapiro’s	-0.25 ... 0.25	$P(A, B) - P(A)P(B)$
$F$	Certainty factor	-1 ... 1	$\max(\frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)})$
$AV$	added value	-0.5 ... 1	$\max(P(B A) - P(B), P(A B) - P(A))$
$K$	Klogsen’s Q	-0.33 ... 0.38	$\sqrt{P(A, B)} \max(P(B A) - P(B), P(A B) - P(A))$
$g$	Goodman-kruskal’s	0 ... 1	$\frac{\sum_j \max_k P(A_j, B_k) + \sum_k \max_j P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
$M$	Mutual Information	0 ... 1	$\frac{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}$
$J$	J-Measure	0 ... 1	$\min(-\sum_i P(A_i) \log P(A_i) \log P(A_i), -\sum_i P(B_i) \log P(B_i) \log P(B_i))$ $\max(P(A, B) \log(\frac{P(B A)}{P(B)}) + P(\bar{A}\bar{B}) \log(\frac{P(\bar{B} \bar{A})}{P(\bar{B})}))$
$G$	Gini index	0 ... 1	$P(A, B) \log(\frac{P(A B)}{P(A)}) + P(\bar{A}\bar{B}) \log(\frac{P(\bar{A} \bar{B})}{P(\bar{A})})$ $\max(P(A)[P(B A)^2 + P(\bar{B} A)^2] + P(\bar{A})[P(B \bar{A})^2 + P(\bar{B} \bar{A})^2] - P(B)^2 - P(\bar{B})^2,$ $P(B)[P(A B)^2 + P(\bar{A} B)^2] + P(\bar{B})[P(A \bar{B})^2 + P(\bar{A} \bar{B})^2] - P(A)^2 - P(\bar{A})^2)$
$s$	support	0 ... 1	$P(A, B)$
$c$	confidence	0 ... 1	$\max(P(B A), P(A B))$
$L$	Laplace	0 ... 1	$\max(\frac{NP(A,B)+1}{NP(A)+2}, \frac{NP(A,B)+1}{NP(B)+2})$
$IS$	Cosine	0 ... 1	$\frac{P(A,B)}{\sqrt{P(A)P(B)}}$
$\gamma$	coherence(Jaccard)	0 ... 1	$\frac{P(A,B)}{P(A)+P(B)-P(A,B)}$
$\alpha$	all.confidence	0 ... 1	$\frac{P(A,B)}{\max(P(A), P(B))}$
$o$	odds ratio	0 ... $\infty$	$\frac{P(A,B)P(\bar{A},\bar{B})}{P(\bar{A},B)P(A,\bar{B})}$
$V$	Conviction	0.5 ... $\infty$	$\max(\frac{P(A)P(\bar{B})}{P(A\bar{B})}, \frac{P(B)P(\bar{A})}{P(\bar{B}\bar{A})})$
$\lambda$	lift	0 ... $\infty$	$\frac{P(A,B)}{P(A)P(B)}$
$S$	Collective strength	0 ... $\infty$	$\frac{P(A,B) + P(\bar{A}\bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A,B) - P(\bar{A}\bar{B})}$
$\chi^2$	$\chi^2$	0 ... $\infty$	$\sum_i \frac{(P(A_i) - E_i)^2}{E_i}$

# Null-invariant measures

---

- Bad property for a measure:
  - given A and B, and a rule  $A \rightarrow B$
  - Is the correlation measure affected by the transactions containing neither A nor B (null-transactions)?
- If it's not  $\rightarrow$  null-invariant property
- E.g. lift and chi-square are not null-invariant

# Null-invariant measures (range [0,1])

---

- All confidence:

$$\begin{aligned} \text{all\_conf}(A, B) &= \\ &= \text{sup}(A \cup B) / \max\{ \text{sup}(A), \text{sup}(B) \} = \\ &= \min\{ P(A|B), P(B|A) \} \end{aligned}$$

- Max confidence:

$$\text{max\_conf}(A, B) = \max\{ P(A|B), P(B|A) \}$$

- Kulczynski

$$\text{Kulc}(A, B) = \frac{1}{2} * (P(A|B) + P(B|A))$$

- Cosine

$$\begin{aligned} \text{cos}(A, B) &= P(A \cup B) / \text{sqrt}( P(A) * P(B) ) = \\ &= \text{sup}(A \cup B) / \text{sqrt}( \text{sup}(A) * \text{sup}(B) ) = \\ &= \text{sqrt}( P(A|B) * P(B|A) ) \end{aligned}$$

# Null-Invariant Measures

Table 6: Properties of interestingness measures. Note that none of the measures satisfies all the properties.

Symbol	Measure	Range	P1	P2	P3	O1	O2	O3	O3'	O4
$\phi$	$\phi$ -coefficient	$-1 \dots 0 \dots 1$	Yes	Yes	Yes	Yes	No	Yes	Yes	No
$\lambda$	Goodman-Kruskal's	$0 \dots 1$	Yes	No	No	Yes	No	No*	Yes	No
$\alpha$	odds ratio	$0 \dots 1 \dots \infty$	Yes*	Yes	Yes	Yes	Yes	Yes*	Yes	No
$Q$	Yule's $Q$	$-1 \dots 0 \dots 1$	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No
$Y$	Yule's $Y$	$-1 \dots 0 \dots 1$	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No
$\kappa$	Cohen's	$-1 \dots 0 \dots 1$	Yes	Yes	Yes	Yes	No	No	Yes	No
$M$	Mutual Information	$0 \dots 1$	Yes	Yes	Yes	No**	No	No*	Yes	No
$J$	J-Measure	$0 \dots 1$	Yes	No	No	No**	No	No	No	No
$G$	Gini index	$0 \dots 1$	Yes	No	No	No**	No	No*	Yes	No
$s$	Support	$0 \dots 1$	No	Yes	No	Yes	No	No	No	No
$c$	Confidence	$0 \dots 1$	No	Yes	No	No**	No	No	No	Yes
$L$	Laplace	$0 \dots 1$	No	Yes	No	No**	No	No	No	No
$V$	Conviction	$0.5 \dots 1 \dots \infty$	No	Yes	No	No**	No	No	Yes	No
$I$	Interest	$0 \dots 1 \dots \infty$	Yes*	Yes	Yes	Yes	No	No	No	No
$IS$	Cosine	$0 \dots \sqrt{P(A, B)} \dots 1$	No	Yes	Yes	Yes	No	No	No	Yes
$PS$	Piatetsky-Shapiro's	$-0.25 \dots 0 \dots 0.25$	Yes	Yes	Yes	Yes	No	Yes	Yes	No
$F$	Certainty factor	$-1 \dots 0 \dots 1$	Yes	Yes	Yes	No**	No	No	Yes	No
$AV$	Added value	$-0.5 \dots 0 \dots 1$	Yes	Yes	Yes	No**	No	No	No	No
$S$	Collective strength	$0 \dots 1 \dots \infty$	No	Yes	Yes	Yes	No	Yes*	Yes	No
$\zeta$	Jaccard	$0 \dots 1$	No	Yes	Yes	Yes	No	No	No	Yes
$K$	Klosgen's	$(\frac{2}{\sqrt{3}} - 1)^{1/2} [2 - \sqrt{3} - \frac{1}{\sqrt{3}}] \dots 0 \dots \frac{2}{3\sqrt{3}}$	Yes	Yes	Yes	No**	No	No	No	No

where: P1:  $O(M) = 0$  if  $\det(M) = 0$ , i.e., whenever  $A$  and  $B$  are statistically independent.

P2:  $O(M_2) > O(M_1)$  if  $M_2 = M_1 + [k \ -k; \ -k \ k]$ .

P3:  $O(M_2) < O(M_1)$  if  $M_2 = M_1 + [0 \ k; \ 0 \ -k]$  or  $M_2 = M_1 + [0 \ 0; \ k \ -k]$ .

O1: Property 1: Symmetry under variable permutation.

O2: Property 2: Row and Column scaling invariance.

O3: Property 3: Antisymmetry under row or column permutation.

O3': Property 4: Inversion invariance.

O4: Property 5: Null invariance.

Yes\*: Yes if measure is normalized.

No\*: Symmetry under row or column permutation.

No\*\*: No unless the measure is symmetrized by taking  $\max(M(A, B), M(B, A))$ .

# Comparison of Interestingness Measures (milk → coffee)

- Null-(transaction) invariance is crucial for correlation analysis
- 5 null-invariant measures

	Milk	No Milk	Sum (row)
Coffee	m, c	~m, c	c
No Coffee	m, ~c	~m, ~c	~c
Sum(col.)	m	~m	Σ

Measure	Definition	Range	Null-Invariant
$\chi^2(a, b)$	$\sum_{i,j=0,1} \frac{(e(a_i, b_j) - o(a_i, b_j))^2}{e(a_i, b_j)}$	$[0, \infty]$	No
$Lift(a, b)$	$\frac{P(ab)}{P(a)P(b)}$	$[0, \infty]$	No
$AllConf(a, b)$	$\frac{sup(ab)}{\max\{sup(a), sup(b)\}}$	$[0, 1]$	Yes
$Coherence(a, b)$	$\frac{sup(ab)}{sup(a) + sup(b) - sup(ab)}$	$[0, 1]$	Yes
$Cosine(a, b)$	$\frac{sup(ab)}{\sqrt{sup(a)sup(b)}}$	$[0, 1]$	Yes
$Kulc(a, b)$	$\frac{sup(ab)}{2} \left( \frac{1}{sup(a)} + \frac{1}{sup(b)} \right)$	$[0, 1]$	Yes
$MaxConf(a, b)$	$\max\left\{ \frac{sup(ab)}{sup(a)}, \frac{sup(ab)}{sup(b)} \right\}$	$[0, 1]$	Yes

Table 3. Interestingness measure definitions.

Null-transactions  
w.r.t. m and c

Null-invariant

Data set	$mc$	$\bar{m}\bar{c}$	$m\bar{c}$	$\bar{m}c$	$\chi^2$	$Lift$	$AllConf$	$Coherence$	$Cosine$	$Kulc$	$MaxConf$
$D_1$	10,000	1,000	1,000	100,000	90557	9.26	0.91	0.83	0.91	0.91	0.91
$D_2$	10,000	1,000	1,000	100	0	1	0.91	0.83	0.91	0.91	0.91
$D_3$	100	1,000	1,000	100,000	670	8.44	0.09	0.05	0.09	0.09	0.09
$D_4$	1,000	1,000	1,000	100,000	24740	25.75	0.5	0.33	0.5	0.5	0.5
$D_5$	1,000	100	10,000	100,000	8173	9.18	0.09	0.09	0.29	0.5	0.91
$D_6$	1,000	10	100,000	100,000	965	1.97	0.01	0.01	0.10	0.5	0.99

Table 2. Example data sets.

Subtle: They disagree

# Which Null-Invariant Measure Is Better?

- IR (Imbalance Ratio): measure the imbalance of two itemsets A and B in rule implications (IR = 0 if the two directional implications between A and B are the same)

$$IR(A, B) = \frac{|sup(A) - sup(B)|}{sup(A) + sup(B) - sup(A \cup B)}$$

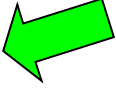
- Kulczynski and Imbalance Ratio (IR) together present a clear picture for all the three datasets D<sub>4</sub> through D<sub>6</sub>
  - D<sub>4</sub> is balanced & neutral
  - D<sub>5</sub> is imbalanced & neutral
  - D<sub>6</sub> is very imbalanced & neutral

<i>Data</i>	<i>mc</i>	$\overline{mc}$	$m\overline{c}$	$\overline{m\overline{c}}$	<i>all_conf.</i>	<i>max_conf.</i>	<i>Kulc.</i>	<i>cosine</i>	IR
D <sub>1</sub>	10,000	1,000	1,000	100,000	0.91	0.91	0.91	0.91	0.0
D <sub>2</sub>	10,000	1,000	1,000	100	0.91	0.91	0.91	0.91	0.0
D <sub>3</sub>	100	1,000	1,000	100,000	0.09	0.09	0.09	0.09	0.0
D <sub>4</sub>	1,000	1,000	1,000	100,000	0.5	0.5	0.5	0.5	0.0
D <sub>5</sub>	1,000	100	10,000	100,000	0.09	0.91	0.5	0.29	0.89
D <sub>6</sub>	1,000	10	100,000	100,000	0.01	0.99	0.5	0.10	0.99



# Chapter 5: Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

---

- Basic Concepts
- Frequent Itemset Mining Methods
- Which Patterns Are Interesting?—Pattern Evaluation Methods
- Summary 

# Summary

---

- Basic concepts: association rules, support-confident framework, closed and max-patterns
- Scalable frequent pattern mining methods
  - Apriori (Candidate generation & test)
  - Projection-based (FPgrowth, CLOSET+, ...)
  - Vertical format approach (ECLAT, CHARM, ...)
- Which patterns are interesting?
  - Pattern evaluation methods

# Ref: Basic Concepts of Frequent Pattern Mining

---

- (**Association Rules**) R. Agrawal, T. Imielinski, and A. Swami. Mining association rules between sets of items in large databases. SIGMOD'93
- (**Max-pattern**) R. J. Bayardo. Efficiently mining long patterns from databases. SIGMOD'98
- (**Closed-pattern**) N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Discovering frequent closed itemsets for association rules. ICDT'99
- (**Sequential pattern**) R. Agrawal and R. Srikant. Mining sequential patterns. ICDE'95

# Ref: Apriori and Its Improvements

---

- R. Agrawal and R. Srikant. Fast algorithms for mining association rules. VLDB'94
- H. Mannila, H. Toivonen, and A. I. Verkamo. Efficient algorithms for discovering association rules. KDD'94
- A. Savasere, E. Omiecinski, and S. Navathe. An efficient algorithm for mining association rules in large databases. VLDB'95
- J. S. Park, M. S. Chen, and P. S. Yu. An effective hash-based algorithm for mining association rules. SIGMOD'95
- H. Toivonen. Sampling large databases for association rules. VLDB'96
- S. Brin, R. Motwani, J. D. Ullman, and S. Tsur. Dynamic itemset counting and implication rules for market basket analysis. SIGMOD'97
- S. Sarawagi, S. Thomas, and R. Agrawal. Integrating

# Ref: Depth-First, Projection-Based FP Mining

---

- R. Agarwal, C. Aggarwal, and V. V. V. Prasad. A tree projection algorithm for generation of frequent itemsets. *J. Parallel and Distributed Computing*, 2002.
- G. Grahne and J. Zhu, Efficiently Using Prefix-Trees in Mining Frequent Itemsets, *Proc. FIMI'03*
- B. Goethals and M. Zaki. An introduction to workshop on frequent itemset mining implementations. *Proc. ICDM'03 Int. Workshop on Frequent Itemset Mining Implementations (FIMI'03)*, Melbourne, FL, Nov. 2003
- J. Han, J. Pei, and Y. Yin. Mining frequent patterns without candidate generation. *SIGMOD' 00*
- J. Liu, Y. Pan, K. Wang, and J. Han. Mining Frequent Item Sets by Opportunistic Projection. *KDD'02*
- J. Han, J. Wang, Y. Lu, and P. Tzvetkov. Mining Top-K Frequent Closed Patterns without Minimum Support. *ICDM'02*
- J. Wang, J. Han, and J. Pei. CLOSET+: Searching for the Best Strategies 47

# Ref: Vertical Format and Row Enumeration Methods

---

- M. J. Zaki, S. Parthasarathy, M. Ogihara, and W. Li. Parallel algorithm for discovery of association rules. DAMI:97.
- M. J. Zaki and C. J. Hsiao. CHARM: An Efficient Algorithm for Closed Itemset Mining, SDM'02.
- C. Bucila, J. Gehrke, D. Kifer, and W. White. DualMiner: A Dual-Pruning Algorithm for Itemsets with Constraints. KDD'02.
- F. Pan, G. Cong, A. K. H. Tung, J. Yang, and M. Zaki , CARPENTER: Finding Closed Patterns in Long Biological Datasets. KDD'03.
- H. Liu, J. Han, D. Xin, and Z. Shao, Mining Interesting Patterns from Very High Dimensional Data: A Top-Down Row Enumeration Approach, SDM'06.

# Ref: Mining Correlations and Interesting Rules

---

- S. Brin, R. Motwani, and C. Silverstein. Beyond market basket: Generalizing association rules to correlations. SIGMOD'97.
- M. Klemettinen, H. Mannila, P. Ronkainen, H. Toivonen, and A. I. Verkamo. Finding interesting rules from large sets of discovered association rules. CIKM'94.
- R. J. Hilderman and H. J. Hamilton. *Knowledge Discovery and Measures of Interest*. Kluwer Academic, 2001.
- C. Silverstein, S. Brin, R. Motwani, and J. Ullman. Scalable techniques for mining causal structures. VLDB'98.
- P.-N. Tan, V. Kumar, and J. Srivastava. Selecting the Right Interestingness Measure for Association Patterns. KDD'02.
- E. Omiecinski. Alternative Interest Measures for Mining Associations. TKDE'03.
- T. Wu, Y. Chen, and J. Han, "Re-Examination of Interestingness Measures in Pattern Mining: A Unified Framework", *Data Mining and Knowledge Discovery*, 21(3):371-397, 2010