# Data Mining:
## Concepts and Techniques
### (3rd ed.)

— Chapter 6 —

Jianjun Cheng

School of Information Science & Engineering

Lanzhou University

# Chapter 6: Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

- **Basic Concepts**

- **Frequent Itemset Mining Methods**

- **Which Patterns Are Interesting?—Pattern Evaluation Methods**
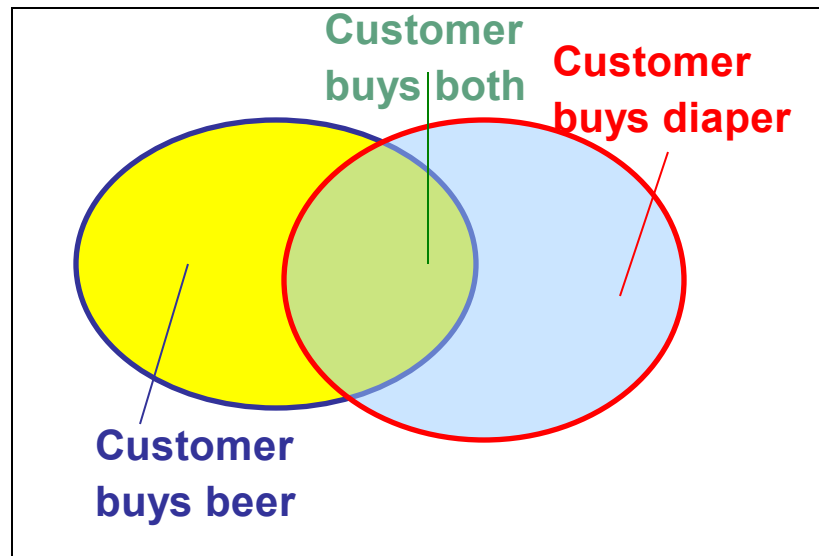
- **Summary**

# What Is Frequent Pattern Analysis?

- Frequent pattern: a pattern (a set of items, subsequences, substructures, etc.) that occurs frequently in a data set

- First proposed by Agrawal, Imielinski, and Swami [AIS93] in the context of frequent itemsets and association rule mining

- Motivation: Finding inherent regularities in data

  - What products were often purchased together?— Beer and diapers?!

  - What are the subsequent purchases after buying a PC?

  - What kinds of DNA are sensitive to this new drug?

  - Can we automatically classify web documents?

- Applications

  - Basket data analysis, cross-marketing, catalog design, sale campaign analysis, Web log (click stream) analysis, and DNA sequence analysis.

# Why Is Freq. Pattern Mining Important?

- Freq. pattern: An intrinsic and important property of datasets
- Foundation for many essential data mining tasks
  - Association, correlation, and causality analysis
  - Sequential, structural (e.g., sub-graph) patterns
  - Pattern analysis in spatiotemporal, multimedia, time-series, and stream data
  - Classification: discriminative, frequent pattern analysis
  - Cluster analysis: frequent pattern-based clustering
  - Data warehousing: iceberg cube and cube-gradient
  - Semantic data compression: fascicles
  - Broad applications
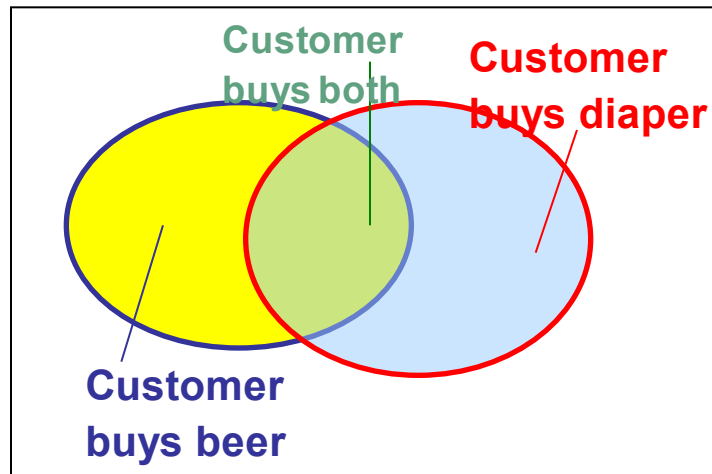
# Basic Concepts: Frequent Patterns

| 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 |



Customer buys both

Customer buys diaper

Customer buys beer

- **itemset**: A set of one or more items
- **k-itemset** $X = \{x_1, \ldots, x_k\}$
- **(absolute) support**, or, **support count** of X: Frequency or occurrence of an itemset X
- **(relative) support**, $s$, is the fraction of transactions that contains X (i.e., the probability that a transaction contains X)
- An itemset X is **frequent** if X's support is no less than a *minsup* threshold

5

# 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 |

**Customer buys both**

**Customer buys diaper**

**Customer buys beer**

- Itemset $X = \{x_1, \ldots, x_k\}$
- Find all the rules $X \rightarrow Y$ with minimum support and confidence
  - support, probability that a transaction contains $X \cup Y$
    $$\text{sup}(X \rightarrow Y) = P(X \cup Y)$$
  - confidence, conditional probability that a transaction having X also contains Y
    $$\text{conf}(X \rightarrow Y) = P(Y|X)$$

Let $\text{sup}_{min} = 50\%$, $\text{conf}_{min} = 50\%$

Freq. Pat.: {A:3, B:3, D:4, E:3, AD:3}

Association rules:

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

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

# Closed Patterns and Max-Patterns

- A long pattern contains a combinatorial number of sub-patterns, e.g., $\{a_1, \ldots, a_{100}\}$ contains $\binom{100}{1} + \binom{100}{2} + \ldots + \binom{100}{100} = 2^{100} - 1 = 1.27 * 10^{30}$ sub-patterns!

- Solution: *Mine closed patterns and max-patterns instead*

- An itemset X is closed in D if there exists no proper super-itemset Y such that Y has the same support count as X in D

- An itemset X is closed frequent itemsets in D if X is both closed and frequent (proposed by Pasquier, et al. @ ICDT'99)

- An itemset X is a max-pattern if X is frequent and there exists no frequent super-pattern $Y \supset X$ (proposed by Bayardo @ SIGMOD'98)

- Closed pattern is a lossless compression of freq. patterns
  - Reducing the # of patterns and rules

# Closed Patterns and Max-Patterns

- Exercise. DB = $\{<a_1, \ldots, a_{100}>, <a_1, \ldots, a_{50}>\}$
  - Min_sup = 1.
- What is the set of <span style="color:red">closed itemset</span>?

  - $<a_1, \ldots, a_{100}>$: 1
  - $<a_1, \ldots, a_{50}>$: 2
- What is the set of <span style="color:red">max-pattern?</span>

  - $<a_1, \ldots, a_{100}>$: 1
- What is the set of <span style="color:red">all patterns</span>?
  - !!

# Computational Complexity of Frequent Itemset Mining

- How many itemsets are potentially to be generated in the worst case?

  - The number of frequent itemsets to be generated is sensitive to the minsup threshold

  - When minsup is low, there exist potentially an exponential number of frequent itemsets

  - The worst case: $M^N$ where M: # distinct items, and N: max length of transactions

- The worst case complexty vs. the expected probability

  - Ex. Suppose Walmart has $10^4$ kinds of products

    - The chance to pick up one product $10^{-4}$

    - The chance to pick up a particular set of 10 products: $\sim 10^{-40}$

    - What is the chance this particular set of 10 products to be frequent $10^3$ times in $10^9$ transactions?

# Chapter 6: Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

- Basic Concepts

☞ - Frequent Itemset Mining Methods

- Which Patterns Are Interesting?—Pattern Evaluation Methods

- Summary

# 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

# The Downward Closure Property and Scalable Mining Methods

- The downward closure property of frequent patterns
  - Any subset of a frequent itemset must be frequent
  - If **{beer, diaper, nuts}** is frequent, so is **{beer, diaper}**
  - i.e., every transaction having {beer, diaper, nuts} also contains {beer, diaper}
- Scalable mining methods: Three major approaches
  - Apriori (Agrawal & Srikant@VLDB'94)
  - Freq. pattern growth (FPgrowth—Han, Pei & Yin @SIGMOD'00)
  - Vertical data format approach (Charm—Zaki & Hsiao @SDM'02)

# Apriori: A Candidate Generation & Test Approach

- <u>Apriori pruning principle</u>: If there is any itemset which is infrequent, its superset should not be generated/tested! (Agrawal & Srikant @VLDB'94, Mannila, et al. @ KDD' 94)
- Method:

    - Initially, scan DB once to get frequent 1-itemset
    - Generate length (k+1) candidate itemsets from length k frequent itemsets
    - Test the candidates against DB
    - Terminate when no frequent or candidate set can be generated

# The Apriori Algorithm—An Example

$Sup_{min} = 2$

### Database TDB

| Tid | Items |
|---|---|
| 10 | A, C, D |
| 20 | B, C, E |
| 30 | A, B, C, E |
| 40 | B, E |

1st scan →

$C_1$

| Itemset | sup |
|---|---|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {D} | 1 |
| {E} | 3 |

$L_1$

| Itemset | sup |
|---|---|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {E} | 3 |

$C_2$

| Itemset |
|---|
| {A, B} |
| {A, C} |
| {A, E} |
| {B, C} |
| {B, E} |
| {C, E} |

2nd scan

$C_2$

| Itemset | sup |
|---|---|
| {A, B} | 1 |
| {A, C} | 2 |
| {A, E} | 1 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$L_2$

| Itemset | sup |
|---|---|
| {A, C} | 2 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$C_3$

| Itemset |
|---|
| {B, C, E} |

3rd scan →

$L_3$

| Itemset | sup |
|---|---|
| {B, C, E} | 2 |

# The Apriori Algorithm (Pseudo-Code)

$C_k$: Candidate itemset of size k
$L_k$: frequent itemset of size k

L$_1$ = {frequent items};
**for** ($k$ = 1; $L_k$ !=$\varnothing$; $k$++) **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** $\cup_k L_k$;

# Implementation of Apriori

- How to generate candidates?
  - Step 1: self-joining $L_k$
  - Step 2: pruning
- Example of Candidate-generation
  - $L_3 = \{abc, abd, acd, ace, bcd\}$
  - 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 $L_3$
  - $C_4 = \{abcd\}$

# Algorithm Apriori

**Algorithm: Apriori.** Find frequent itemsets using an iterative level-wise approach based on candidate generation.

**Input:**

- $D$, a database of transactions;

- $min\_sup$, the minimum support count threshold.

**Output:** $L$, frequent itemsets in $D$.

**Method:**

(1)     $L_1 = $ find_frequent_1-itemsets(D);
(2)     **for** $(k = 2; L_{k-1} \neq \phi; k{+}{+})$ {
(3)         $C_k = $ apriori_gen$(L_{k-1})$;
(4)         **for each** transaction $t \in D$ { // scan $D$ for counts
(5)             $C_t = $ subset$(C_k, t)$; // get the subsets of $t$ that are candidates
(6)             **for each** candidate $c \in C_t$
(7)                 c.count++;
(8)         }
(9)         $L_k = \{c \in C_k | c.count \geq min\_sup\}$
(10)    }
(11)    **return** $L = \cup_k L_k$;

# Algorithm Apriori

**procedure apriori_gen**($L_{k-1}$:frequent $(k-1)$-itemsets)
(1)      **for each** itemset $l_1 \in L_{k-1}$
(2)          **for each** itemset $l_2 \in L_{k-1}$
(3)              **if** $(l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2])$
                  $\wedge ... \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$ **then** {
(4)                  $c = l_1 \bowtie l_2$; // join step: generate candidates
(5)                  **if** has_infrequent_subset($c, L_{k-1}$) **then**
(6)                      **delete** $c$; // prune step: remove unfruitful candidate
(7)                  **else add** $c$ **to** $C_k$;
(8)              }
(9)      **return** $C_k$;

**procedure has_infrequent_subset**($c$: candidate $k$-itemset;
          $L_{k-1}$: frequent $(k-1)$-itemsets); // use prior knowledge
(1)      **for each** $(k-1)$-subset $s$ **of** $c$
(2)          **if** $s \notin L_{k-1}$ **then**
(3)              **return** TRUE;
(4)      **return** FALSE;

# Another Example

| TID | item-ID list |
|------|--------------|
| 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 |

# Another Example

C₁

Scan D for count of each candidate →

| Itemset | Sup. count |
|---------|------------|
| {I1} | 6 |
| {I2} | 7 |
| {I3} | 6 |
| {I4} | 2 |
| {I5} | 2 |

Compare candidate support count with minimum support count →

L₁

| Itemset | Sup. count |
|---------|------------|
| {I1} | 6 |
| {I2} | 7 |
| {I3} | 6 |
| {I4} | 2 |
| {I5} | 2 |

Generate C₂ candidates from L₁ →

C₂

| Itemset |
|---------|
| {I1, I2} |
| {I1, I3} |
| {I1, I4} |
| {I1, I5} |
| {I2, I3} |
| {I2, I4} |
| {I2, I5} |
| {I3, I4} |
| {I3, I5} |
| {I4, I5} |

Scan D for count of each candidate →

C₂

| Itemset | Sup. count |
|---------|------------|
| {I1, I2} | 4 |
| {I1, I3} | 4 |
| {I1, I4} | 1 |
| {I1, I5} | 2 |
| {I2, I3} | 4 |
| {I2, I4} | 2 |
| {I2, I5} | 2 |
| {I3, I4} | 0 |
| {I3, I5} | 1 |
| {I4, I5} | 0 |

Compare candidate support count with minimum support count →

L₂

| Itemset | Sup. count |
|---------|------------|
| {I1, I2} | 4 |
| {I1, I3} | 4 |
| {I1, I5} | 2 |
| {I2, I3} | 4 |
| {I2, I4} | 2 |
| {I2, I5} | 2 |

# Another Example



Generate $C_3$ candidates from $L_2$

**$C_3$**

| Itemset |
|---|
| {I1, I2, I3} |
| {I1, I2, I5} |

Scan $D$ for count of each candidate

**$C_3$**

| Itemset | Sup. count |
|---|---|
| {I1, I2, I3} | 2 |
| {I1, I2, I5} | 2 |

Compare candidate support count with minimum support count

**$L_3$**

| Itemset | Sup. count |
|---|---|
| {I1, I2, I3} | 2 |
| {I1, I2, I5} | 2 |

# How to generate the association rules?

■it is straightforward to generate strong association rules from the frequent itemsets

■strong association rules satisfy both minimum support and minimum confidence

$$\text{conf}(A \Rightarrow B) = P(B|A)$$

■method

- For each frequent itemset I, generate all nonempty subsets of I

- For every nonempty subset s of I, output the rule

$$s \Rightarrow (I - s)$$

if  support_count(I)/support_count(s) ≥min_conf, where min_conf is the minimum confidence threshold

# 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

# Further Improvement of the Apriori Method

- Major computational challenges
  - Multiple scans of transaction database
  - Huge number of candidates
  - Tedious workload of support counting for candidates
- Improving Apriori: general ideas
  - Reduce passes of transaction database scans
  - Shrink number of candidates
  - Facilitate support counting of candidates

# DHP: Reduce the Number of Candidates

- Used to reduce the size of the candidate k-itemsets, $C_k$ , for k > 1

- A *k*-itemset whose corresponding hashing bucket count is below the threshold cannot be frequent

$H_2$

Create hash table $H_2$
using hash function
$h(x, y) = ((order\ of\ x) \times 10 + (order\ of\ y))\ mod\ 7$

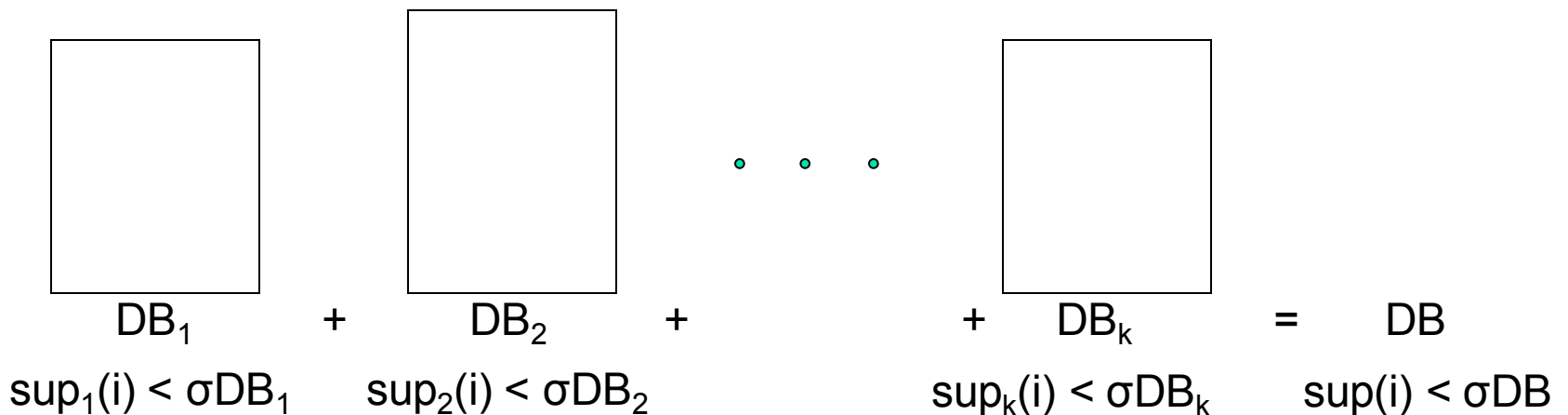| bucket address | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| bucket count | 2 | 2 | 4 | 2 | 2 | 4 | 4 |
| bucket contents | {I1, I4}<br>{I3, I5} | {I1, I5}<br>{I1, I5} | {I2, I3}<br>{I2, I3}<br>{I2, I3}<br>{I2, I3} | {I2, I4}<br>{I2, I4} | {I2, I5}<br>{I2, I5} | {I1, I2}<br>{I1, I2}<br>{I1, I2}<br>{I1, I2} | {I1, I3}<br>{I1, I3}<br>{I1, I3}<br>{I1, I3} |

- J. Park, M. Chen, and P. Yu. An effective hash-based algorithm for mining association rules. *SIGMOD'95*

# Transaction reduction

- reducing the number of transactions scanned in future iterations

  - A transaction that does not contain any frequent k-itemsets cannot contain any frequent (k + 1)-itemsets.

  - such a transaction can be marked or removed from further consideration because subsequent database scans for j-itemsets, where j > k, will not need to consider such a transaction
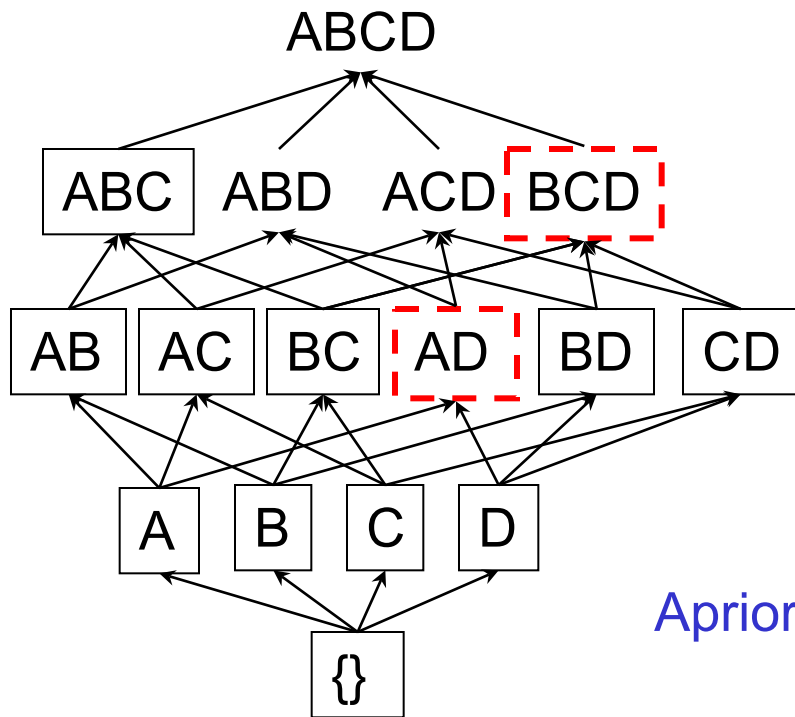
# Partition: Scan Database Only Twice

- Any itemset that is potentially frequent in DB must be frequent in at least one of the partitions of DB
  - Scan 1: partition database and find local frequent patterns
  - Scan 2: consolidate global frequent patterns
- A. Savasere, E. Omiecinski and S. Navathe, *VLDB'95*

|  |  |  |  |  |  |
|---|---|---|---|---|---|
| $DB_1$ | + | $DB_2$ | + | ... + $DB_k$ | = DB |
| $sup_1(i) < \sigma DB_1$ | | $sup_2(i) < \sigma DB_2$ | | $sup_k(i) < \sigma DB_k$ | $sup(i) < \sigma DB$ |

# Sampling for Frequent Patterns

- Select a sample of original database, mine frequent patterns within sample using Apriori

- Scan database once to verify frequent itemsets found in sample, only *borders* of closure of frequent patterns are checked

  - Example: check *abcd* instead of *ab, ac, ..., etc.*

- Scan database again to find missed frequent patterns

- H. Toivonen. Sampling large databases for association rules. In *VLDB'96*
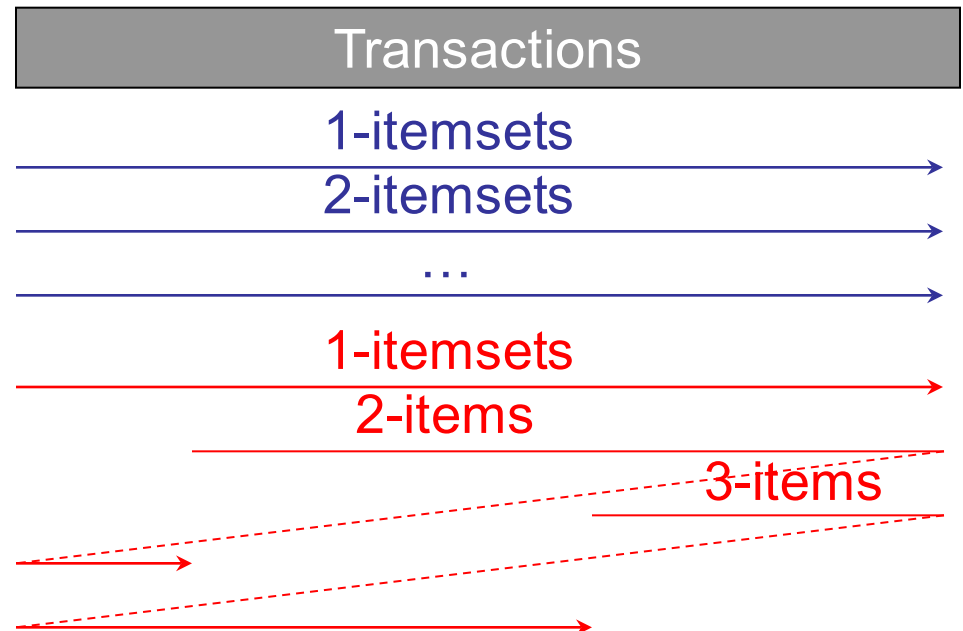
# DIC: Reduce Number of Scans

ABCD

ABC  ABD  ACD  BCD

AB  AC  BC  AD  BD  CD

A  B  C  D

{}

Itemset lattice

S. Brin R. Motwani, J. Ullman, and S. Tsur. Dynamic itemset counting and implication rules for market basket data. *SIGMOD'97*

- the database is partitioned into blocks marked by start points
- new candidate itemsets can be added at any start point
- Once both A and D are determined frequent, the counting of AD begins
- Once all length-2 subsets of BCD are determined frequent, the counting of BCD begins

Transactions

Apriori

1-itemsets

2-itemsets

…

DIC

1-itemsets

2-items

3-items

# 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
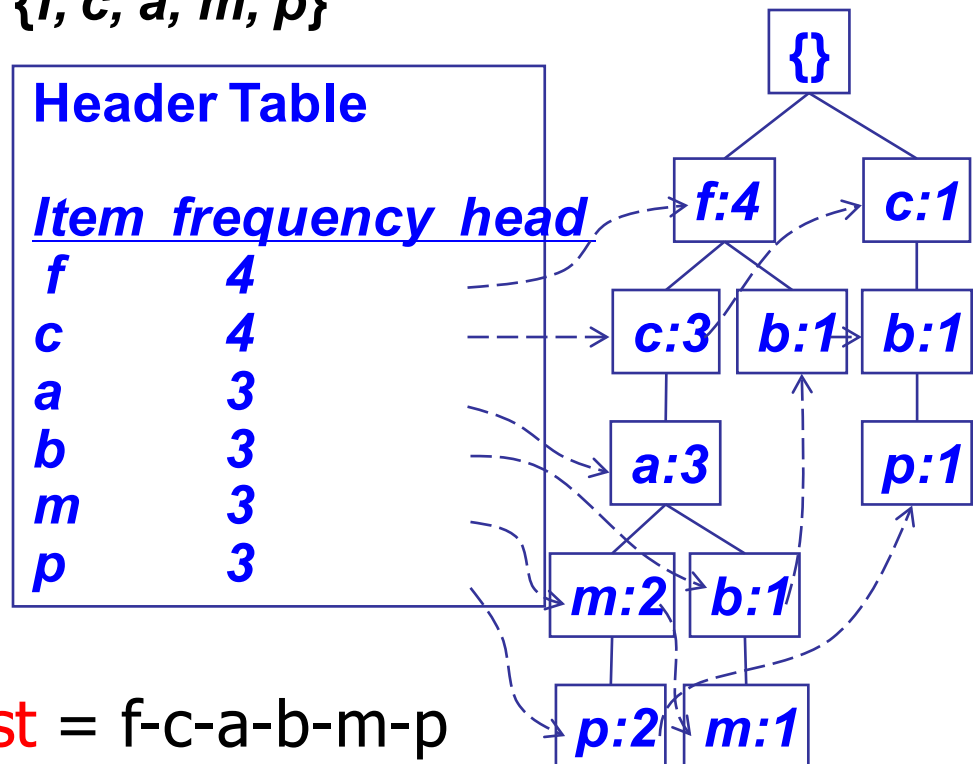
- Bottlenecks of the Apriori approach
  - Breadth-first (i.e., level-wise) search
  - Candidate generation and test
    - Often generates a huge number of candidates
- The FPGrowth Approach (J. Han, J. Pei, and Y. Yin, SIGMOD' 00)
  - Depth-first search
  - Avoid explicit candidate generation
- Major philosophy: Grow long patterns from short ones using local frequent items only
  - "abc" is a frequent pattern
  - Get all transactions having "abc", i.e., project DB on abc: DB|abc
  - "d" is a local frequent item in DB|abc → abcd is a frequent pattern

# Construct FP-tree from a Transaction Database

| TID | Items bought | (ordered) frequent items |
|-----|-------------|--------------------------|
| 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**

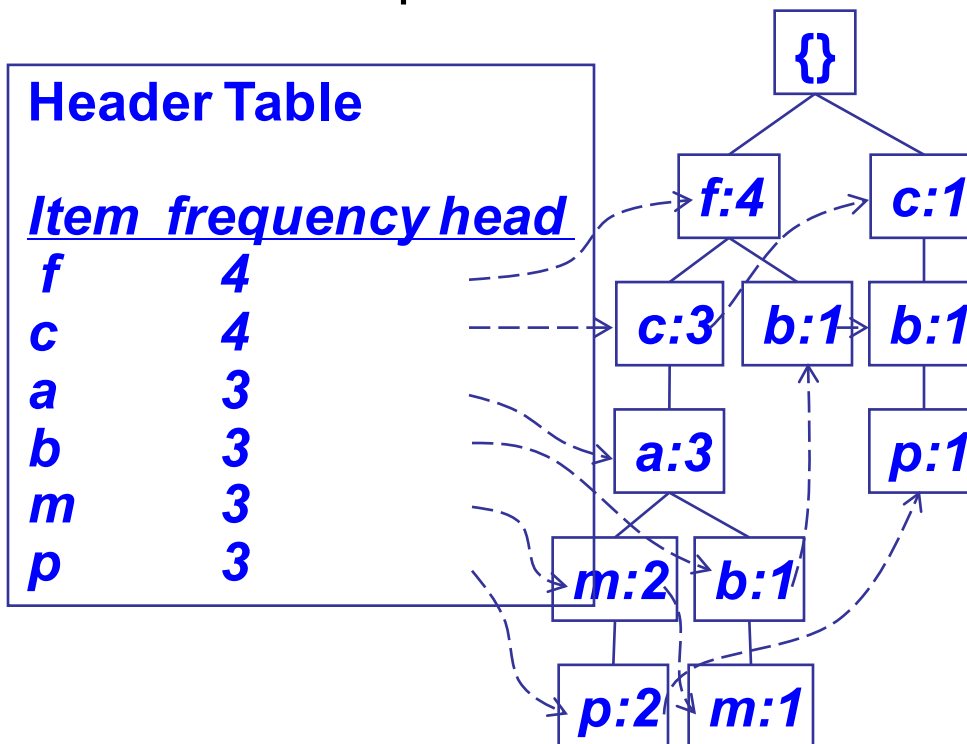| Item | frequency | head |
|------|-----------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |
| p | 3 | |

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

# Partition Patterns and Databases

- Frequent patterns can be partitioned into subsets according to f-list
  - 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
- Completeness and non-redundancy

# Find Patterns Having P From P-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 *p*
- Accumulate all of *transformed prefix paths* of item *p* to form *p's* conditional pattern base

**Header Table**

| Item | frequency | head |
|------|-----------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |
| p | 3 | |

{}

f:4    c:1

c:3  b:1  b:1

a:3        p:1

m:2  b:1

p:2  m:1

**Conditional pattern bases**

| item | cond. pattern base |
|------|--------------------|
| c | f:3 |
| a | fc:3 |
| b | fca:1, f:1, c:1 |
| m | fca:2, fcab:1 |
| p | fcam:2, cb:1 |

# From Conditional Pattern-bases to Conditional FP-trees

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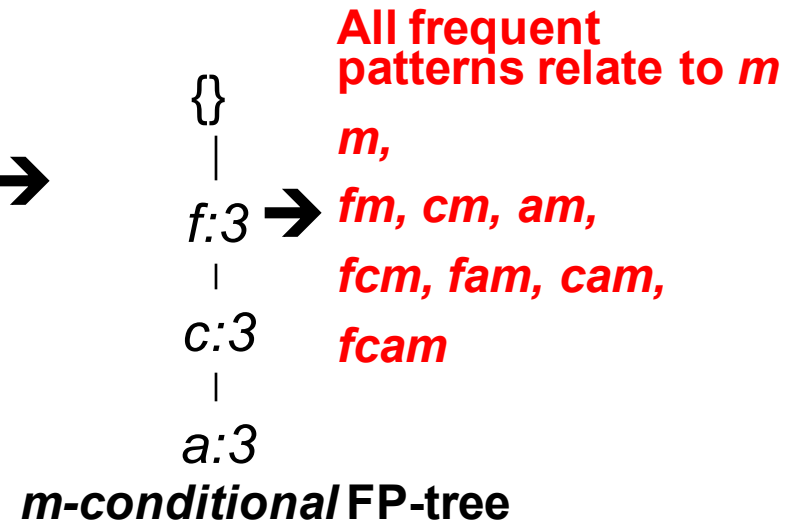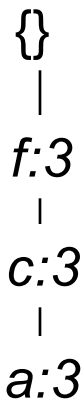  - Accumulate the count for each item in the base
  - Construct the FP-tree for the frequent items of the pattern base

**Header Table**

| Item | frequency | head |
|------|-----------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |
| p | 3 | |

```
              {}
         /         \
      f:4           c:1
      /  \            \
    c:3   b:1         b:1
     |                  \
    a:3                 p:1
    /  \
  m:2   b:1
   |
  p:2   m:1
```

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

➔

```
  {}
   |
  f:3   ➔
   |
  c:3
   |
  a:3
```

**m-conditional** FP-tree

**All frequent patterns relate to m**

**m,**

**fm, cm, am,**

**fcm, fam, cam,**

**fcam**

# 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)
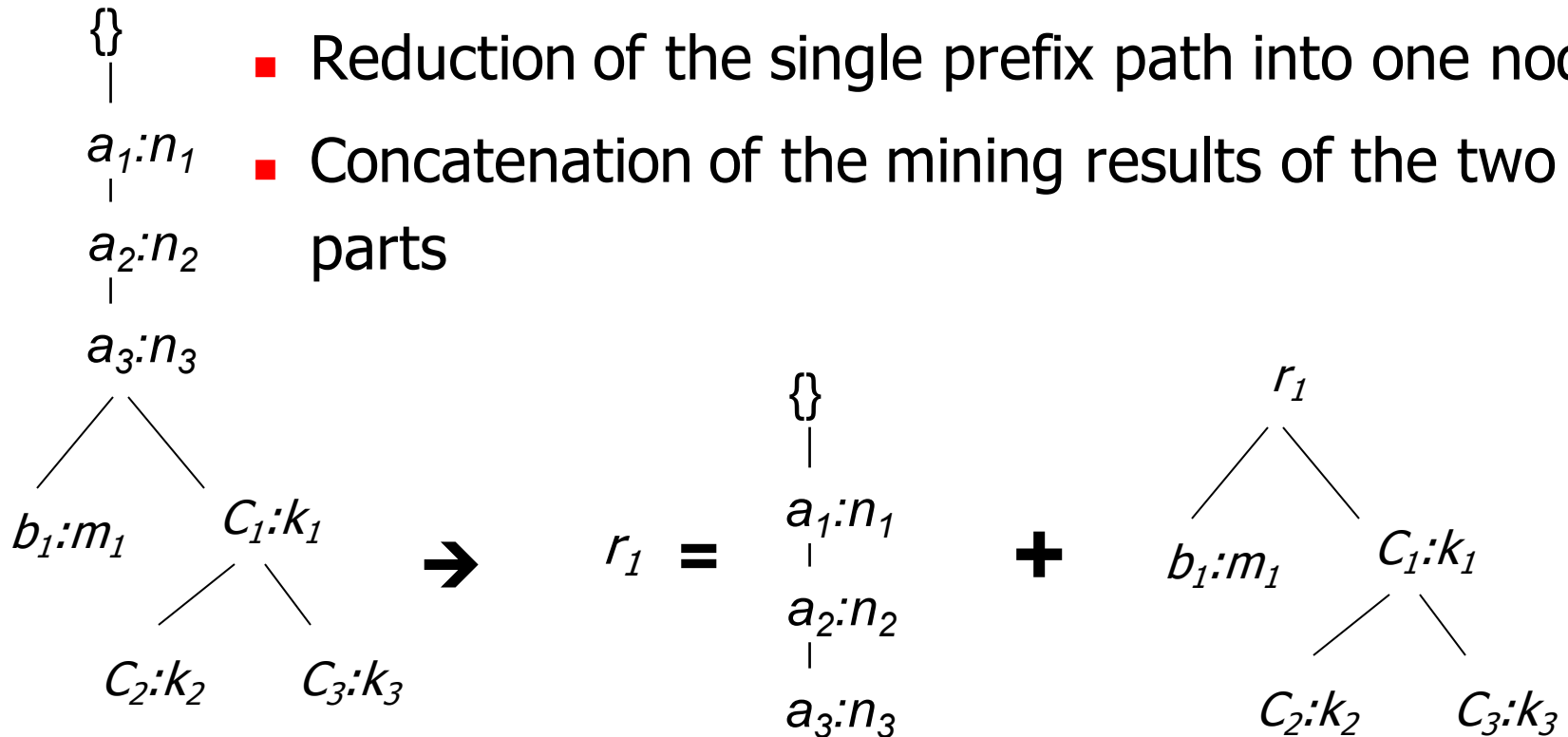
{}
|
f:3

**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

{}
|
$a_1{:}n_1$
|
$a_2{:}n_2$
|
$a_3{:}n_3$

$b_1{:}m_1$          $C_1{:}k_1$

$C_2{:}k_2$          $C_3{:}k_3$

➔

$r_1$ =

{}
|
$a_1{:}n_1$
|
$a_2{:}n_2$
|
$a_3{:}n_3$

**+**

$r_1$

$b_1{:}m_1$          $C_1{:}k_1$

$C_2{:}k_2$          $C_3{:}k_3$

# An Example

| TID | item-ID list |
|---|---|
| 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 |

# An Example



| Item | Conditional Pattern Base | Conditional FP-tree | Frequent Patterns Generated |
|------|--------------------------|---------------------|-----------------------------|
| I5 | {{I2, I1: 1}, {I2, I1, I3: 1}} | ⟨I2: 2, I1: 2⟩ | {I2, I5: 2}, {I1, I5: 2}, {I2, I1, I5: 2} |
| I4 | {{I2, I1: 1}, {I2: 1}} | ⟨I2: 2⟩ | {I2, I4: 2} |
| I3 | {{I2, I1: 2}, {I2: 2}, {I1: 2}} | ⟨I2: 4, I1: 2⟩, ⟨I1: 2⟩ | {I2, I3: 4}, {I1, I3: 4}, {I2, I1, I3: 2} |
| I1 | {{I2: 4}} | ⟨I2: 4⟩ | {I2, I1: 4} |

# 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

- 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—single path will generate all the combinations of its sub-paths, each of which is a frequent pattern

# The Frequent Pattern Growth Mining Method

**Algorithm: FP_growth.** Mine frequent itemsets using an FP-tree by pattern fragment growth.

**Input:**

- $D$, a transaction database;
- $min\_sup$, the minimum support count threshold.

**Output:** The complete set of frequent patterns.

**Method:**

1. The FP-tree is constructed in the following steps:

   (a) Scan the transaction database $D$ once. Collect $F$, the set of frequent items, and their support counts. Sort $F$ in support count descending order as $L$, the *list* of frequent items.

   (b) Create the root of an FP-tree, and label it as "null." For each transaction *Trans* in $D$ do the following.
   Select and sort the frequent items in *Trans* according to the order of $L$. Let the sorted frequent item list in *Trans* be $[p|P]$, where $p$ is the first element and $P$ is the remaining list. Call **insert_tree**($[p|P]$, $T$), which is performed as follows. If $T$ has a child $N$ such that $N.item\text{-}name = p.item\text{-}name$, then increment $N$'s count by 1; else create a new node $N$, and let its count be 1, its parent link be linked to $T$, and its node-link to the nodes with the same *item-name* via the node-link structure. If $P$ is nonempty, call **insert_tree**($P$, $N$) recursively.

# The Frequent Pattern Growth Mining Method

**2.** The FP-tree is mined by calling **FP_growth**(*FP_tree, null*), which is implemented as follows.

**procedure FP_growth**(*Tree, α*)
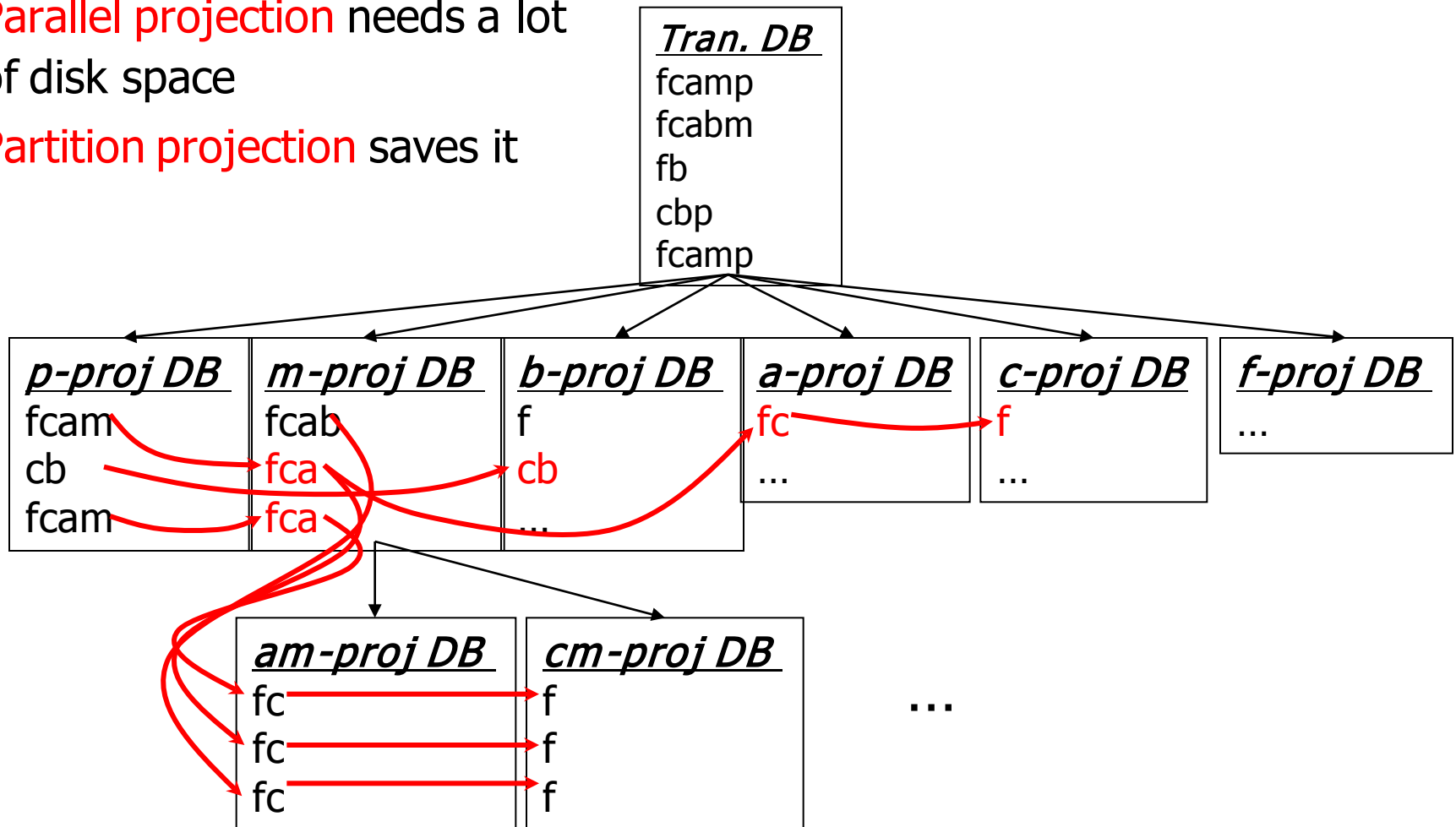
(1)    **if** *Tree* contains a single path *P* **then**

(2)        **for each** combination (denoted as *β*) of the nodes in the path *P*

(3)            generate pattern $\beta \cup \alpha$ with *support_count = minimum support count of nodes in β*;

(4)    **else for each** $a_i$ in the header of *Tree* {

(5)        generate pattern $\beta = a_i \cup \alpha$ with *support_count* = $a_i$.*support_count*;

(6)        construct *β*'s conditional pattern base and then *β*'s conditional FP_tree $Tree_\beta$;

(7)        **if** $Tree_\beta \neq \emptyset$ **then**

(8)            call **FP_growth**($Tree_\beta, \beta$); }
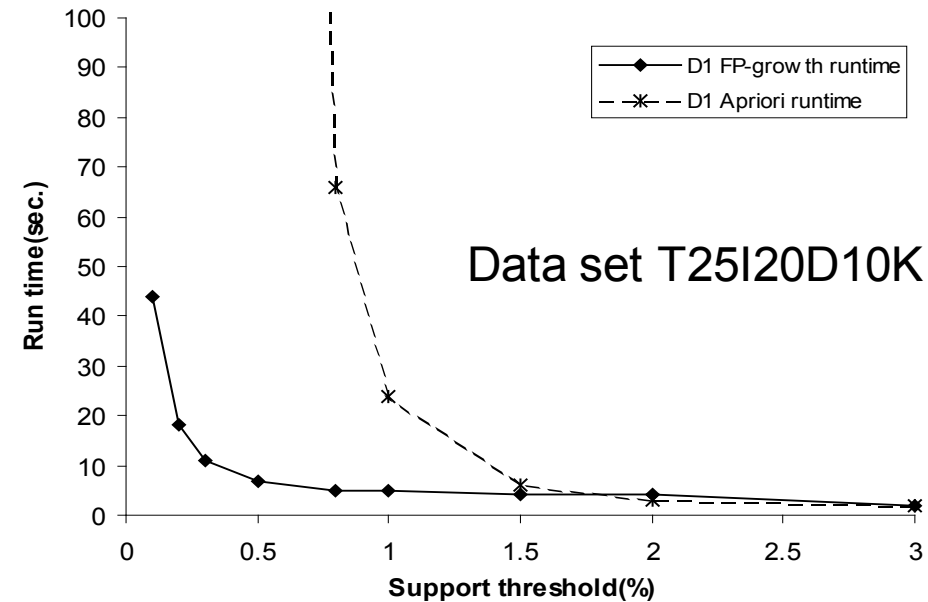
# Scaling FP-growth by Database Projection

- 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

**Tran. DB**
fcamp
fcabm
fb
cbp
fcamp

**p-proj DB**
fcam
cb
fcam

**m-proj DB**
fcab
fca
fca

**b-proj DB**
f
cb
...

**a-proj DB**
fc
...

**c-proj DB**
f
...

**f-proj DB**
...

**am-proj DB**
fc
fc
fc

**cm-proj DB**
f
f
f

...

# Performance of FPGrowth in Large Datasets

Data set T25I20D10K

Data set T25I20D100K

FP-Growth vs. Apriori

FP-Growth vs. Tree-Projection

# 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 FPMax (Grahne & Zhu, 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

# ECLAT: Mining by Exploring Vertical Data Format

- Vertical format: $t(AB) = \{T_{11}, T_{25}, \ldots\}$

    - 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

    - $t(X) = \{T_1, T_2, T_3\}$, $t(XY) = \{T_1, T_3\}$

    - $Diffset(XY, X) = \{T_2\}$

- Eclat  (Zaki et al. @KDD'97)

- Mining Closed patterns using vertical format:  CHARM (Zaki & Hsiao@SDM'02)

# Mining frequent itemses using vertical data format

- ECLAT (Equivalence CLAss Transformation)

    Mining is performed on this data set by intersecting the TID sets of every pair of frequent single items

| TID | Items |
|-----|-------|
| 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 |

| itemset | TID-set |
|---------|---------|
| 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 |

# Mining frequent itemses using vertical data format

- if min_support_count=2, then all the 1-itemset are frequent

| itemset | TID-set |
|---------|---------|
| 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 |

| itemset | TID-set |
|---------|---------|
| 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 Maxpatterns

# Mining Frequent Closed Patterns: CLOSET

- Flist: list of all frequent items in support ascending order
  - Flist: d-a-f-e-c

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 |

- 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", DMKD'00.

# CLOSET+: Mining Closed Itemsets by Pattern-Growth

- Itemset merging: if Y appears in every occurrence of X, then Y is merged with X

- Sub-itemset pruning: if Y ⊃ X, and sup(X) = sup(Y), X and all of X's descendants in the set enumeration tree can be pruned

- Hybrid tree projection
    - Bottom-up physical tree-projection
    - Top-down pseudo tree-projection

- 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

# MaxMiner: Mining Max-Patterns

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

- 1st scan: find frequent items
  - A, B, C, D, E
- 2nd scan: find support for
  - AB, AC, AD, AE, ABCDE
  - BC, BD, BE, BCDE
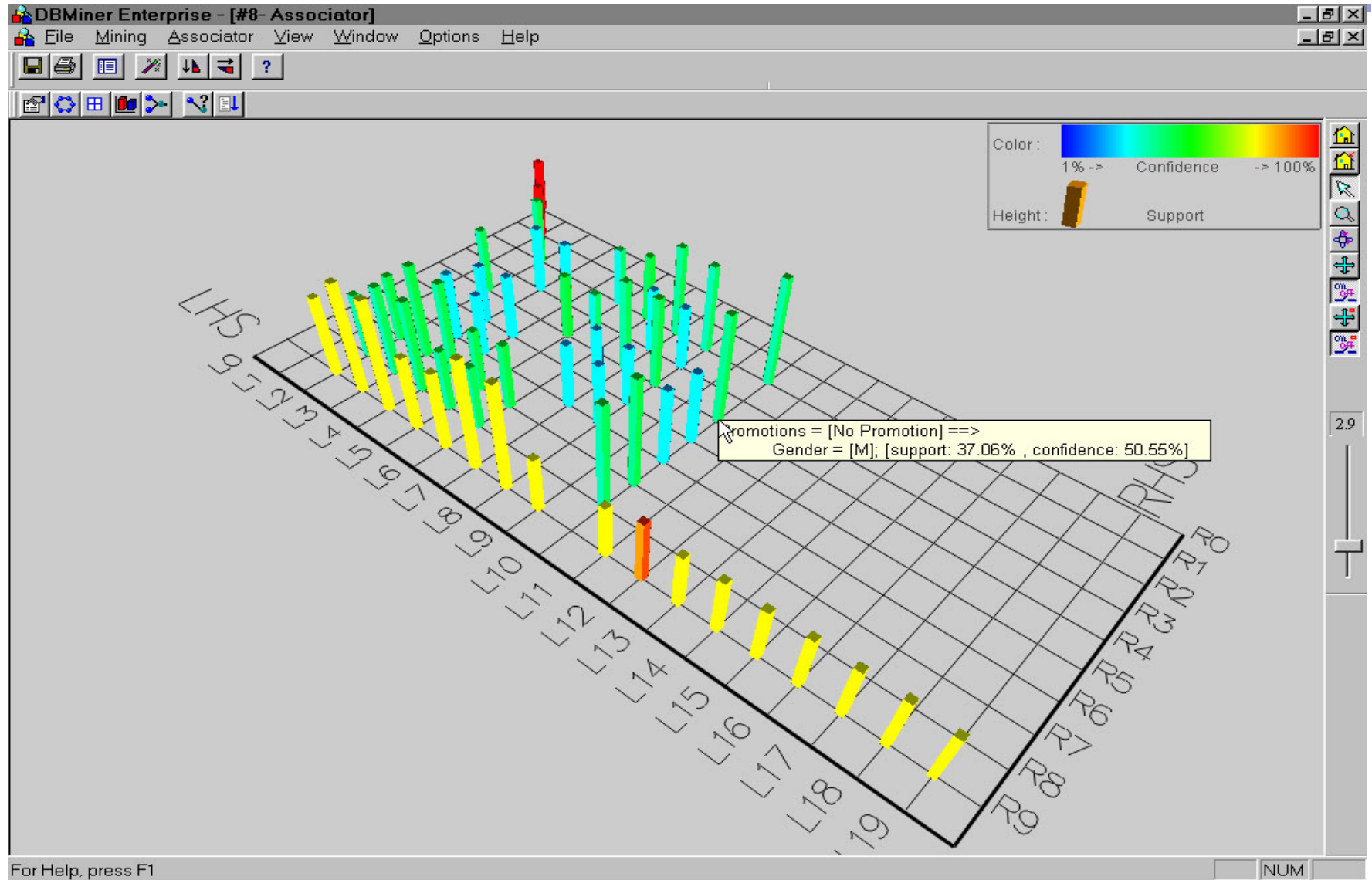  - CD, CE, CDE, DE

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*

# CHARM: Mining by Exploring Vertical Data Format

- Vertical format: $t(AB) = \{T_{11}, T_{25}, \ldots\}$
  - tid-list: list of trans.-ids containing an itemset
- Deriving closed 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
  - $t(X) = \{T_1, T_2, T_3\}$, $t(XY) = \{T_1, T_3\}$
  - Diffset $(XY, X) = \{T_2\}$
- Eclat/MaxEclat (Zaki et al. @KDD'97), VIPER(P. Shenoy et al.@SIGMOD'00), CHARM (Zaki & Hsiao@SDM'02)

# Visualization of Association Rules: Plane Graph

# Visualization of Association Rules: Rule Graph

# Visualization of Association Rules (SGI/MineSet 3.0)

# Chapter 6: Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

- Basic Concepts

- Frequent Itemset Mining Methods

☞ - Which Patterns Are Interesting?—Pattern

Evaluation Methods

- Summary

# 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

- Measure of dependent/correlated events: lift

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

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

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

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

# Are *lift* and χ² Good Measures of Correlation?

- *"Buy walnuts ⇒ buy milk* [1%, 80%]"  is misleading if 85% of customers buy milk

- Support and confidence are not good to indicate correlations

- Over 20 interestingness measures have been proposed  (see Tan, Kumar, Sritastava @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 | $\frac{P(A,B)P(\overline{A},\overline{B})-P(A,\overline{B})P(\overline{A},B)}{P(A,B)P(\overline{A},\overline{B})+P(A,\overline{B})P(\overline{A},B)}$ |
| $Y$ | Yule's Y | -1 ... 1 | $\frac{\sqrt{P(A,B)P(\overline{A},\overline{B})}-\sqrt{P(A,\overline{B})P(\overline{A},B)}}{\sqrt{P(A,B)P(\overline{A},\overline{B})}+\sqrt{P(A,\overline{B})P(\overline{A},B)}}$ |
| $k$ | Cohen's | -1 ... 1 | $\frac{P(A,B)+P(\overline{A},\overline{B})-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A)P(B)-P(\overline{A})P(\overline{B})}$ |
| $PS$ | 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 | $\max(P(B|A)-P(B), P(A|B)-P(A))$ |
| $K$ | Klosgen'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 | $\frac{\Sigma_j \max_k P(A_j,B_k)+\Sigma_k \max_j P(A_j,B_k)-\max_j P(A_j)-\max_k P(B_k)}{2-\max_j P(A_j)-max_k P(B_k)}$ |
| $M$ | Mutual Information | 0 ... 1 | $\frac{\Sigma_i\Sigma_j P(A_i,B_j)\log\frac{P(A_i,B_j)}{P(A_i)P(B_j)}}{\min(-\Sigma_i P(A_i)\log P(A_i)\log P(A_i), -\Sigma_i P(B_i)\log P(B_i)\log P(B_i))}$ |
| $J$ | J-Measure | 0 ... 1 | $\max(P(A,B)\log(\frac{P(B|A)}{P(B)}) + P(A\overline{B})\log(\frac{P(\overline{B}|A)}{P(\overline{B})}))$ $P(A,B)\log(\frac{P(A|B)}{P(A)}) + P(\overline{A}B)\log(\frac{P(\overline{A}|B)}{P(\overline{A})})$ |
| $G$ | Gini index | 0 ... 1 | $\max(P(A)[P(B|A)^2 + P(\overline{B}|A)^2] + P(\overline{A})[P(B|\overline{A})^2 + P(\overline{B}|\overline{A})^2] - P(B)^2 - P(\overline{B})^2,$ $P(B)[P(A|B)^2 + P(\overline{A}|B)^2] + P(\overline{B})[P(A|\overline{B})^2 + P(\overline{A}|\overline{B})^2] - P(A)^2 - P(\overline{A})^2)$ |
| $s$ | 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 ... ∞ | $\frac{P(A,B)P(\overline{A},\overline{B})}{P(\overline{A},B)P(A,\overline{B})}$ |
| $V$ | Conviction | 0.5 ... ∞ | $\max(\frac{P(A)P(\overline{B})}{P(A\overline{B})}, \frac{P(B)P(\overline{A})}{P(B\overline{A})})$ |
| $\lambda$ | lift | 0 ... ∞ | $\frac{P(A,B)}{P(A)P(B)}$ |
| $S$ | Collective strength | 0 ... ∞ | $\frac{P(A,B)+P(\overline{AB})}{P(A)P(B)+P(\overline{A})P(\overline{B})} \times \frac{1-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A,B)-P(\overline{AB})}$ |
| $\chi^2$ | $\chi^2$ | 0 ... ∞ | $\Sigma_i\frac{(P(A_i)-E_i)^2}{E_i}$ |

# 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 | $-1\cdots0\cdots1$ | Yes | Yes | Yes | Yes | No | Yes | Yes | No |
| $\lambda$ | Goodman-Kruskal's | $0\cdots1$ | Yes | No | No | Yes | No | No* | Yes | No |
| $\alpha$ | odds ratio | $0\cdots1\cdots\infty$ | Yes* | Yes | Yes | Yes | Yes | Yes* | Yes | No |
| $Q$ | Yule's $Q$ | $-1\cdots0\cdots1$ | Yes | Yes | Yes | Yes | Yes | Yes | Yes | No |
| $Y$ | Yule's $Y$ | $-1\cdots0\cdots1$ | Yes | Yes | Yes | Yes | Yes | Yes | Yes | No |
| $\kappa$ | Cohen's | $-1\cdots0\cdots1$ | Yes | Yes | Yes | Yes | No | No | Yes | No |
| $M$ | Mutual Information | $0\cdots1$ | Yes | Yes | Yes | No** | No | No* | Yes | No |
| $J$ | J-Measure | $0\cdots1$ | Yes | No | No | No** | No | No | No | No |
| $G$ | Gini index | $0\cdots1$ | Yes | No | No | No** | No | No* | Yes | No |
| $s$ | Support | $0\cdots1$ | No | Yes | No | Yes | No | No | No | No |
| $c$ | Confidence | $0\cdots1$ | No | Yes | No | No** | No | No | No | Yes |
| $L$ | Laplace | $0\cdots1$ | No | Yes | No | No** | No | No | No | No |
| $V$ | Conviction | $0.5\cdots1\cdots\infty$ | No | Yes | No | No** | No | No | Yes | No |
| $I$ | Interest | $0\cdots1\cdots\infty$ | Yes* | Yes | Yes | Yes | No | No | No | No |
| $IS$ | Cosine | $0\cdots\sqrt{P(A,B)}\cdots1$ | No | Yes | Yes | Yes | No | No | No | Yes |
| $PS$ | Piatetsky-Shapiro's | $-0.25\cdots0\cdots0.25$ | Yes | Yes | Yes | Yes | No | Yes | Yes | No |
| $F$ | Certainty factor | $-1\cdots0\cdots1$ | Yes | Yes | Yes | No** | No | No | Yes | No |
| $AV$ | Added value | $-0.5\cdots0\cdots1$ | Yes | Yes | Yes | No** | No | No | No | No |
| $S$ | Collective strength | $0\cdots1\cdots\infty$ | No | Yes | Yes | Yes | No | Yes* | Yes | No |
| $\zeta$ | Jaccard | $0\cdots1$ | No | Yes | Yes | Yes | No | No | No | Yes |
| $K$ | Klosgen's | $(\frac{2}{\sqrt{3}}-1)^{1/2}[2-\sqrt{3}-\frac{1}{\sqrt{3}}]\cdots0\cdots\frac{2}{3\sqrt{3}}$ | Yes | Yes | Yes | No** | No | No | No | No |

where: P1: $O(\mathbf{M})=0$ if $det(\mathbf{M})=0$, *i.e.*, whenever $A$ and $B$ are statistically independent.
P2: $O(\mathbf{M_2})>O(\mathbf{M_1})$ if $\mathbf{M_2}=\mathbf{M_1}+[k\ -k;\ -k\ k]$.
P3: $O(\mathbf{M_2})<O(\mathbf{M_1})$ if $\mathbf{M_2}=\mathbf{M_1}+[0\ k;\ 0\ -k]$ or $\mathbf{M_2}=\mathbf{M_1}+[0\ 0;\ k\ -k]$.
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

- Null-(transaction) invariance is crucial for correlation analysis
- Lift and $\chi^2$ are not null-invariant
- 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 | $\Sigma$ |

| 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 |
| $Kulc(a, b)$ | $\frac{sup(ab)}{2}(\frac{1}{sup(a)} + \frac{1}{sup(b)})$ | $[0, 1]$ | Yes |
| $MaxConf(a,b)$ | $max\{\frac{sup(ab)}{sup(a)}, \frac{sup(ab)}{sup(b)}\}$ | $[0, 1]$ | Yes |

**Table 3.** Interestingness measure definitions.

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

Kulczynski measure (1927)

Null-invariant

| Data set | $mc$ | $\overline{m}c$ | $m\overline{c}$ | $\overline{m}\overline{c}$ | $\chi^2$ | 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

# Analysis of DBLP Coauthor Relationships

Recent DB conferences, removing balanced associations, low sup, etc.

| ID | Author $a$ | Author $b$ | $sup(ab)$ | $sup(a)$ | $sup(b)$ | $Coherence$ | $Cosine$ | $Kulc$ |
|----|-----------|-----------|-----------|----------|----------|-------------|----------|--------|
| 1 | Hans-Peter Kriegel | Martin Ester | 28 | 146 | 54 | 0.163 (2) | 0.315 (7) | 0.355 (9) |
| 2 | Michael Carey | Miron Livny | 26 | 104 | 58 | 0.191 (1) | 0.335 (4) | 0.349 (10) |
| 3 | Hans-Peter Kriegel | Joerg Sander | 24 | 146 | 36 | 0.152 (3) | 0.331 (5) | 0.416 (8) |
| 4 | Christos Faloutsos | Spiros Papadimitriou | 20 | 162 | 26 | 0.119 (7) | 0.308 (10) | 0.446 (7) |
| 5 | Hans-Peter Kriegel | Martin Pfeifle | 18 | 146 | 18 | 0.123 (6) | 0.351 (2) | 0.562 (2) |
| 6 | Hector Garcia-Molina | Wilburt Labio | 16 | 144 | 18 | 0.110 (9) | 0.314 (8) | 0.500 (4) |
| 7 | Divyakant Agrawal | Wang Hsiung | 16 | 120 | 16 | 0.133 (5) | 0.365 (1) | 0.567 (1) |
| 8 | Elke Rundensteiner | Murali Mani | 16 | 104 | 20 | 0.148 (4) | 0.351 (3) | 0.477 (6) |
| 9 | Divyakant Agrawal | Oliver Po | 12 | 120 | 12 | 0.100 (10) | 0.316 (6) | 0.550 (3) |
| 10 | Gerhard Weikum | Martin Theobald | 12 | 106 | 14 | 0.111 (8) | 0.312 (9) | 0.485 (5) |

Table 5. Experiment on DBLP data set.

Advisor-advisee relation: Kulc: high, coherence: low, cosine: middle

- Tianyi Wu, Yuguo Chen and Jiawei Han, "Association Mining in Large Databases: A Re-Examination of Its Measures", Proc. 2007 Int. Conf. Principles and Practice of Knowledge Discovery in Databases (PKDD'07), Sept. 2007

# Which Null-Invariant Measure Is Better?

- IR (Imbalance Ratio): measure the imbalance of two itemsets A and B in rule implications

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

- Kulczynski and Imbalance Ratio (IR) together present a clear picture for all the three datasets $D_4$ through $D_6$

  - $D_4$ is balanced & neutral

  - $D_5$ is imbalanced & neutral

  - $D_6$ is very imbalanced & neutral

| Data | $mc$ | $\overline{m}c$ | $m\overline{c}$ | $\overline{m}\overline{c}$ | all_conf. | max_conf. | Kulc. | cosine | IR |
|---|---|---|---|---|---|---|---|---|---|
| $D_1$ | 10,000 | 1,000 | 1,000 | 100,000 | 0.91 | 0.91 | 0.91 | 0.91 | 0.0 |
| $D_2$ | 10,000 | 1,000 | 1,000 | 100 | 0.91 | 0.91 | 0.91 | 0.91 | 0.0 |
| $D_3$ | 100 | 1,000 | 1,000 | 100,000 | 0.09 | 0.09 | 0.09 | 0.09 | 0.0 |
| $D_4$ | 1,000 | 1,000 | 1,000 | 100,000 | 0.5 | 0.5 | 0.5 | 0.5 | 0.0 |
| $D_5$ | 1,000 | 100 | 10,000 | 100,000 | 0.09 | 0.91 | 0.5 | 0.29 | 0.89 |
| $D_6$ | 1,000 | 10 | 100,000 | 100,000 | 0.01 | 0.99 | 0.5 | 0.10 | 0.99 |

# Chapter 6: 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 association rule mining with relational database systems: Alternatives and implications. SIGMOD'98

# 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 for Mining Frequent Closed Itemsets.  KDD'03

# 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