**Data Mining**

**Session 6 – Main Theme**
**Mining Frequent Patterns,**
**Association, and Correlations**

**Dr. Jean-Claude Franchitti**

*New York University*
*Computer Science Department*
*Courant Institute of Mathematical Sciences*

*Adapted from course textbook resources*
*Data Mining Concepts and Techniques (2nd Edition)*
*Jiawei Han and Micheline Kamber*

---

## Agenda

1 **Session Overview**

2 **Mining Frequent Patterns, Association, and Correlations**

3 **Summary and Conclusion**

- Course description and syllabus:
  - » http://www.nyu.edu/classes/jcf/g22.3033-002/
  - » http://www.cs.nyu.edu/courses/spring10/G22.3033-002/index.html

- Textbooks:
  - » *Data Mining: Concepts and Techniques (2nd Edition)*
    Jiawei Han, Micheline Kamber
    Morgan Kaufmann
    ISBN-10: 1-55860-901-6, ISBN-13: 978-1-55860-901-3, (2006)
  - » *Microsoft SQL Server 2008 Analysis Services Step by Step*
    Scott Cameron
    Microsoft Press
    ISBN-10: 0-73562-620-0, ISBN-13: 978-0-73562-620-31 1st Edition (04/15/09)

- Basic concepts and a roadmap

- Scalable frequent itemset mining methods

- Mining various kinds of association rules

- From association to correlation analysis

- Constraint-based association mining

- Mining colossal patterns

- Summary

## Icons / Metaphors

Information

Common Realization

Knowledge/Competency Pattern

Governance

Alignment

Solution Approach

## Agenda

| 1 | Session Overview |
| 2 | Mining Frequent Patterns, Association, and Correlations |
| 3 | Summary and Conclusion |

## Mining Frequent Patterns, Association and Correlations – Sub-Topics

➡ ▪ Basic concepts and a road map

▪ Scalable frequent itemset mining methods

▪ Mining various kinds of association rules

▪ From association to correlation analysis

▪ Constraint-based association mining

▪ Mining colossal patterns

▪ 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.
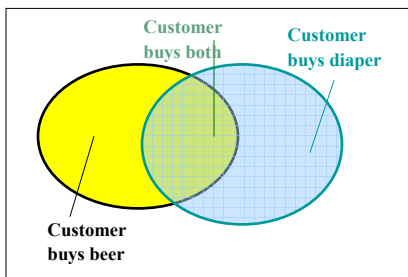
## Why Is Freq. Pattern Mining Important?

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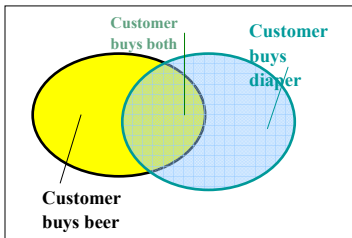
- itemset: A set of one or more items
- k-itemset $X = \{x_1, \ldots, x_k\}$
- *(absolute) support*, or, *support count* of X: Frequency or occurrence of an itemset X
- *(relative) support*, *s*, is the fraction of transactions that contains X (i.e., the probability that a transaction contains X)
- An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

## 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

- Find all the rules $X \rightarrow Y$ with minimum support and confidence
  - support, *s*, probability that a transaction contains $X \cup Y$
  - confidence, *c,* conditional probability that a transaction having X also contains *Y*
- *Let minsup = 50%, minconf = 50%*
- *Freq. Pat.:* Beer:3, Nuts:3, Diaper:4, Eggs:3, {Beer, Diaper}:3

- Association rules: (many more!)
  - *Beer → Diaper* (60%, 100%)
  - *Diaper → Beer* (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 $\binom{100}{1} + \binom{100}{2} + \ldots + \binom{100}{100} = 2^{100} - 1 = 1.27 \times 10^{30}$ sub-patterns!
- Solution: *Mine closed patterns and max-patterns instead*
- An itemset X is closed if X is *frequent* and there exists *no super-pattern* $Y \supset X$, *with the same support* as X (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 = {$<a_1, ..., a_{100}>$, $< a_1, ..., a_{50}>$}
  - Min_sup = 1.
- What is the set of closed itemset?
  - $<a_1, ..., a_{100}>$: 1
  - $< a_1, ..., a_{50}>$: 2
- What is the set of max-pattern?
  - $<a_1, ..., a_{100}>$: 1
- What is the set of all patterns?
  - !!

## 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 senstive 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?

- Basic concepts and a road map
➡ - Scalable frequent itemset mining methods

- Mining various kinds of association rules

- From association to correlation analysis

- Constraint-based association mining

- Mining colossal patterns

- Summary

---

**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)
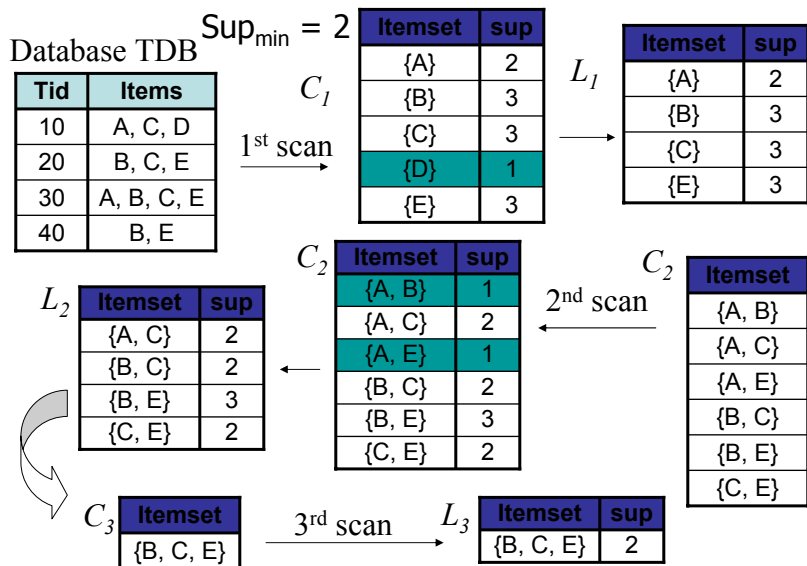  - Vertical data format approach (Charm—Zaki & Hsiao @SDM'02)

## Apriori: A Candidate Generation & Test Approach

- Apriori pruning principle: 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 | sup |
|---------|-----|
| {A, B} | 1 |
| {A, C} | 2 |
| {A, E} | 1 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

2nd scan

$C_2$

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

$L_2$

| 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:
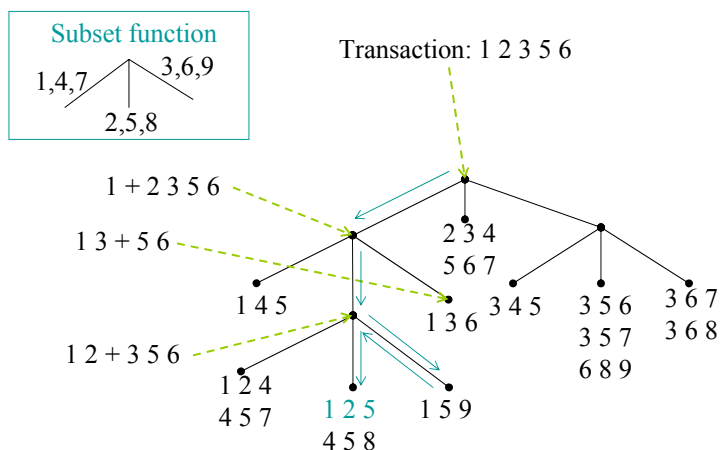    - acde is removed because ade is not in $L_3$
  - $C_4$ = {abcd}

## 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: Counting Supports of Candidates

Subset function

1,4,7     3,6,9

2,5,8

Transaction: 1 2 3 5 6

1 + 2 3 5 6

1 3 + 5 6

1 2 + 3 5 6

1 4 5

1 2 4
4 5 7

1 2 5
4 5 8

1 5 9

2 3 4
5 6 7

1 3 6

3 4 5

3 5 6
3 5 7
6 8 9

3 6 7
3 6 8

## Candidate Generation: An SQL Implementation

- SQL Implementation of candidate generation
  - Suppose the items in $L_{k-1}$ are listed in an order
  - Step 1: self-joining $L_{k-1}$
    - insert into $C_k$
    - select $p.item_1, p.item_2, \ldots, p.item_{k-1}, q.item_{k-1}$
    - from $L_{k-1}\ p, L_{k-1}\ q$
    - where $p.item_1 = q.item_1, \ldots, p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$
  - 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** $C_k$
- Use object-relational extensions like UDFs, BLOBs, and Table functions for efficient implementation [S. Sarawagi, S. Thomas, and R. Agrawal. Integrating association rule mining with relational database systems: Alternatives and implications. SIGMOD'98]

## 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
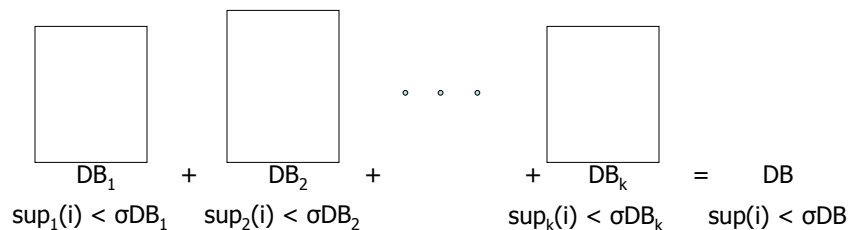
- TD-Close (Liu, et al, SDM'06)

## 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

## 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) < σDB$_1$ | | sup$_2$(i) < σDB$_2$ | | | | sup$_k$(i) < σDB$_k$ | | sup(i) < σDB |

## DHP: Reduce the Number of Candidates

- A *k*-itemset whose corresponding hashing bucket count is below the threshold cannot be frequent
  - Candidates: a, b, c, d, e
  - Hash entries: {ab, ad, ae} {bd, be, de} …
  - Frequent 1-itemset: a, b, d, e
  - ab is not a candidate 2-itemset if the sum of count of {ab, ad, ae} is below support threshold
- J. Park, M. Chen, and P. Yu. An effective hash-based algorithm for mining association rules. In *SIGMOD'95*
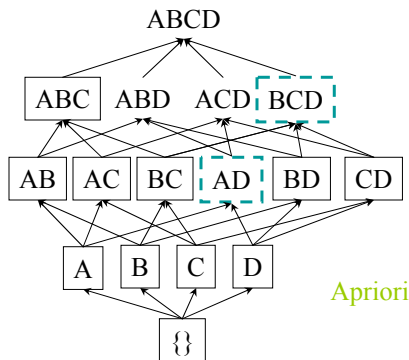
## Sampling for Frequent Patterns

- 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*
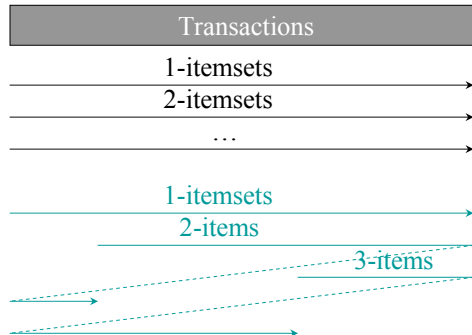
## DIC: Reduce Number of Scans

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. In
*SIGMOD'97*

- 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

---

## 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

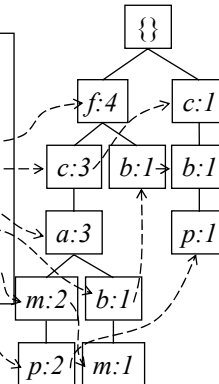| TID | Items bought | (ordered) frequent items |
|-----|--------------|--------------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o, w} | {f, b} |
| 400 | {b, c, k, s, p} | {c, b, p} |
| 500 | {a, f, c, e, l,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

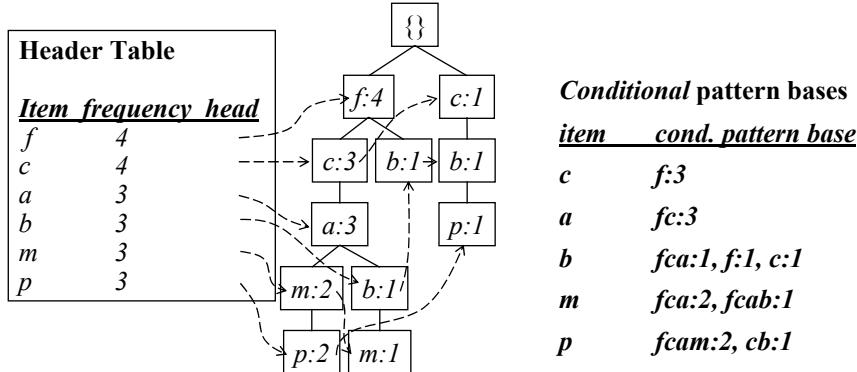| Item | frequency | head |
|------|-----------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |
| p | 3 | |

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



31

---

## 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
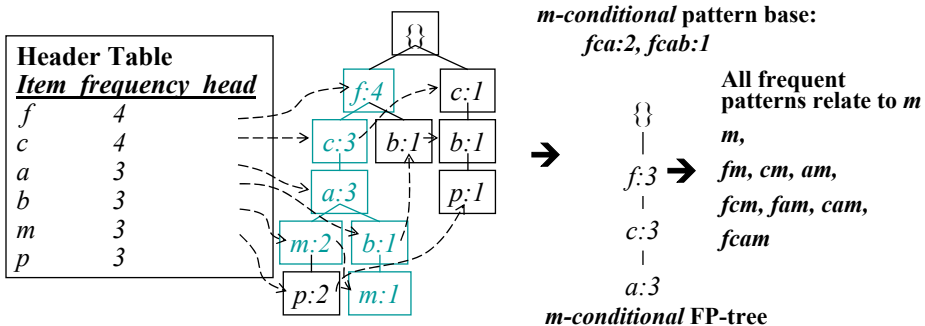- Completeness and non-redundancy

32

- 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 | |

*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 |

33

---

- 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 | |

*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*

34

Cond. pattern base of "am": (fc:3)

$$\{\} \\ | \\ f:3 \\ | \\ c:3$$

*am-conditional* **FP-tree**

$$\{\} \\ | \\ f:3 \\ | \\ c:3 \\ | \\ a:3$$

*m-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**

35

---

- 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



36

## 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
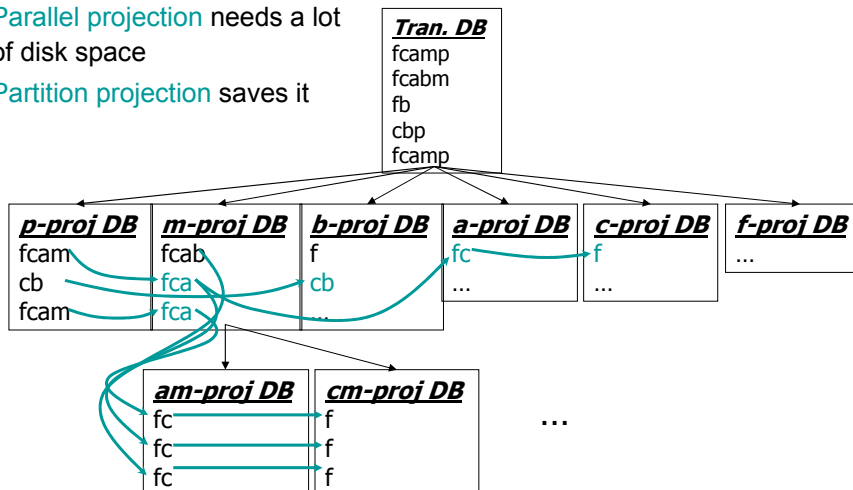
## Scaling FP-growth by Database Projection

- 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

...

## FP-Growth vs. Apriori: Scalability With the Support Threshold

Data set T25I20D10K



Legend:
- D1 FP-grow th runtime
- D1 Apriori runtime

Y-axis: Run time(sec.)
X-axis: Support threshold(%)

## FP-Growth vs. Tree-Projection: Scalability with the Support Threshold

Data set T25I20D100K



Legend:
- D2 FP-growth
- D2 TreeProjection

Y-axis: Runtime (sec.)
X-axis: Support threshold (%)

## 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)

43

## 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)

44

## MaxMiner: Mining Max-patterns

- 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,
- 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*

Potential max-patterns

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

45

## Mining Frequent Closed Patterns: CLOSET

- Flist: list of all frequent items in support ascending 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", DMKD'00.

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 |

46

## 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

## CHARM / 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 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

## Mining Frequent Patterns, Association and Correlations – Sub-Topics

- Basic concepts and a road map
- Scalable frequent itemset mining methods
➡ - Mining various kinds of association rules
- From association to correlation analysis
- Constraint-based association mining
- From association to correlation analysis
- Mining colossal patterns
- Summary

## Mining Various Kinds of Association Rules

- Mining multilevel association

- Mining multidimensional association

- Mining quantitative association

- Mining interesting correlation patterns

## Mining Multiple-Level Association Rules

- Items often form hierarchies
- Flexible support settings
  - Items at the lower level are expected to have lower support
- Exploration of *shared* multi-level mining (Agrawal & Srikant@VLB'95, Han & Fu@VLDB'95)

uniform support                                  reduced support

Level 1          Milk                   Level 1
min_sup = 5%     [support = 10%]        min_sup = 5%

Level 2     2% Milk        Skim Milk    Level 2
min_sup = 5%  [support = 6%]   [support = 4%]   min_sup = 3%

## 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

## Mining Multi-Dimensional Association

- Single-dimensional rules:
  - buys(X, "milk") $\Rightarrow$ buys(X, "bread")
- Multi-dimensional rules: $\geq$ 2 dimensions or predicates
  - Inter-dimension assoc. rules (*no repeated predicates*)
    - age(X,"19-25") $\wedge$ occupation(X,"student") $\Rightarrow$ buys(X, "coke")
  - hybrid-dimension assoc. rules (*repeated predicates*)
    - age(X,"19-25") $\wedge$ buys(X, "popcorn") $\Rightarrow$ buys(X, "coke")
- Categorical Attributes: finite number of possible values, no ordering among values—data cube approach
- Quantitative Attributes: Numeric, implicit ordering among values—discretization, clustering, and gradient approaches

## Mining Quantitative Associations

- Techniques can be categorized by how numerical attributes, such as age or salary are treated:

  1. Static discretization based on predefined concept hierarchies (data cube methods)
  2. Dynamic discretization based on data distribution (quantitative rules, e.g., Agrawal & Srikant@SIGMOD96)
  3. Clustering: Distance-based association (e.g., Yang & Miller@SIGMOD97)
     - One dimensional clustering then association
  4. Deviation: (such as Aumann and Lindell@KDD99)
     Sex = female => Wage: mean=$7/hr (overall mean = $9)

## 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



()

(age)  (income)  (buys)

(age, income)  (age,buys)  (income,buys)

(age,income,buys)

## Quantitative Association Rules

- Proposed by Lent, Swami and Widom ICDE'97
- Numeric attributes are *dynamically* discretized
  - » Such that the confidence or compactness of the rules mined is maximized
- 2-D quantitative association rules: $A_{quan1} \wedge A_{quan2} \Rightarrow A_{cat}$
- Cluster *adjacent* association rules to form general rules using a 2-D grid
- Example

age(X, "34-35") $\wedge$ income(X, "30-50K") $\Rightarrow$ buys(X, "high resolution TV")



---

## Mining Other Interesting Patterns

- Flexible support constraints (Wang, et al. @ VLDB'02)
  - Some items (e.g., diamond) may occur rarely but are valuable
  - Customized $sup_{min}$ specification and application
- Top-K closed frequent patterns (Han, et al. @ ICDM'02)
  - Hard to specify $sup_{min}$, but top-k with $length_{min}$ is more desirable
  - Dynamically raise $sup_{min}$ in FP-tree construction and mining, and select most promising path to mine

- Basic concepts and a road map

- Scalable frequent itemset mining methods

- Mining various kinds of association rules

➡ - From association to correlation analysis

- Constraint-based association mining

- Mining colossal patterns

- Summary

61

---

## 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

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

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

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

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

62

## 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 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 | $-1\ldots1$ | $\frac{P(A,B)-P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$ |
| $Q$ | Yule's Q | $-1\ldots1$ | $\frac{P(A,B)P(\overline{A},\overline{B})-P(A,\overline{B})P(\overline{A},B)}{P(A,B)P(\overline{A},\overline{B})+P(A,\overline{B})P(\overline{A},B)}$ |
| $Y$ | Yule's Y | $-1\ldots1$ | $\frac{\sqrt{P(A,B)P(\overline{A},\overline{B})}-\sqrt{P(A,\overline{B})P(\overline{A},B)}}{\sqrt{P(A,B)P(\overline{A},\overline{B})}+\sqrt{P(A,\overline{B})P(\overline{A},B)}}$ |
| $k$ | Cohen's | $-1\ldots1$ | $\frac{P(A,B)+P(\overline{A},\overline{B})-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A)P(B)-P(\overline{A})P(\overline{B})}$ |
| $PS$ | Piatetsky-Shapiro's | $-0.25\ldots0.25$ | $P(A,B)-P(A)P(B)$ |
| $F$ | Certainty factor | $-1\ldots1$ | $\max(\frac{P(B|A)-P(B)}{1-P(B)},\frac{P(A|B)-P(A)}{1-P(A)})$ |
| $AV$ | added value | $-0.5\ldots1$ | $\max(P(B|A)-P(B),P(A|B)-P(A))$ |
| $K$ | Kloegen's Q | $-0.33\ldots0.38$ | $\sqrt{P(A,B)}\max(P(B|A)-P(B),P(A|B)-P(A))$ |
| $g$ | Goodman-kruskal's | $0\ldots1$ | $\frac{\sum_j \max_k P(A_j,B_k)+\sum_k \max_j P(A_j,B_k)-\max_j P(A_j)-\max_k P(B_k)}{2-\max_j P(A_j)-\max_k P(B_k)}$ |
| $M$ | Mutual Information | $0\ldots1$ | $\frac{\sum_i\sum_j P(A_i,B_j)\log\frac{P(A_i,B_j)}{P(A_i)P(B_j)}}{\min(-\sum_i P(A_i)\log P(A_i),-\sum_j P(B_j)\log P(B_j))}$ |
| $J$ | J-Measure | $0\ldots1$ | $\max(P(A,B)\log(\frac{P(B|A)}{P(B)})+P(A\overline{B})\log(\frac{P(\overline{B}|A)}{P(\overline{B})}),$ $P(A,B)\log(\frac{P(A|B)}{P(A)})+P(\overline{A}B)\log(\frac{P(\overline{A}|B)}{P(A)}))$ |
| $G$ | Gini index | $0\ldots1$ | $\max(P(A)[P(B|A)^2+P(\overline{B}|A)^2]+P(\overline{A})[P(B|\overline{A})^2+P(\overline{B}|\overline{A})^2]-P(B)^2-P(\overline{B})^2,$ $P(B)[P(A|B)^2+P(\overline{A}|B)^2]+P(\overline{B})[P(A|\overline{B})^2+P(\overline{A}|\overline{B})^2]-P(A)^2-P(\overline{A})^2)$ |
| $s$ | support | $0\ldots1$ | $P(A,B)$ |
| $c$ | confidence | $0\ldots1$ | $\max(P(B|A),P(A|B))$ |
| $L$ | Laplace | $0\ldots1$ | $\max(\frac{NP(A,B)+1}{NP(A)+2},\frac{NP(A,B)+1}{NP(B)+2})$ |
| $IS$ | Cosine | $0\ldots1$ | $\frac{P(A,B)}{\sqrt{P(A)P(B)}}$ |
| $\gamma$ | coherence(Jaccard) | $0\ldots1$ | $\frac{P(A,B)}{P(A)+P(B)-P(A,B)}$ |
| $\alpha$ | all_confidence | $0\ldots1$ | $\frac{P(A,B)}{\max(P(A),P(B))}$ |
| $o$ | odds ratio | $0\ldots\infty$ | $\frac{P(A,B)P(\overline{A},\overline{B})}{P(A,\overline{B})P(\overline{A},B)}$ |
| $V$ | Conviction | $0.5\ldots\infty$ | $\max(\frac{P(A)P(\overline{B})}{P(A\overline{B})},\frac{P(B)P(\overline{A})}{P(B\overline{A})})$ |
| $\lambda$ | lift | $0\ldots\infty$ | $\frac{P(A,B)}{P(A)P(B)}$ |
| $S$ | Collective strength | $0\ldots\infty$ | $\frac{P(A,B)+P(\overline{AB})}{P(A)P(B)+P(\overline{A})P(\overline{B})}\times\frac{1-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A,B)-P(\overline{AB})}$ |
| $\chi^2$ | $\chi^2$ | $0\ldots\infty$ | $\sum_i\frac{(P(A_i)-E_i)^2}{E_i}$ |

## 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 | No | 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: $O(M)=0$ 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

- 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 | 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 |
| $Kulc(a,b)$ | $\frac{sup(ab)}{2}(\frac{1}{sup(a)}+\frac{1}{sup(b)})$ | $[0,1]$ | 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 | $mc$ | $\overline{m}c$ | $m\overline{c}$ | $\overline{mc}$ | $\chi^2$ | $Lift$ | $AllConf$ | $Coherence$ | $Cosine$ | $Kulc$ | $MaxConf$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $D_1$ | 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

- 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 | $mc$ | $\overline{m}c$ | $m\overline{c}$ | $\overline{mc}$ | 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 |

- Basic concepts and a road map

- Scalable frequent itemset mining methods

- Mining various kinds of association rules

- From association to correlation analysis

➡ - Constraint-based association mining

- Mining colossal patterns

- Summary

## Constraint-based (Query-Directed) Mining

- Finding all the patterns in a database autonomously? — unrealistic!
  - The patterns could be too many but not focused!
- Data mining should be an interactive process
  - User directs what to be mined using a data mining query language (or a graphical user interface)
- Constraint-based mining
  - User flexibility: provides constraints on what to be mined
  - System optimization: explores such constraints for efficient mining — constraint-based mining: constraint-pushing, similar to push selection first in DB query processing
  - Note: still find all the answers satisfying constraints, not finding some answers in "heuristic search"

## Constraints in Data Mining

- Knowledge type constraint:
  - classification, association, etc.
- Data constraint — using SQL-like queries
  - find product pairs sold together in stores in Chicago in Dec.'02
- Dimension/level constraint
  - in relevance to region, price, brand, customer category
- Rule (or pattern) constraint
  - small sales (price < $10) triggers big sales (sum > $200)
- Interestingness constraint
  - strong rules: min_support $\geq$ 3%, min_confidence $\geq$ 60%

## Constraint-Based Frequent Pattern Mining

- Classification of constraints based on their constraint-pushing capabilities
  - Anti-monotonic: If constraint c is violated, its further mining can be terminated
  - Monotonic: If c is satisfied, no need to check c again
  - Data anti-monotonic: If a transaction t does not satisfy c, t can be pruned from its further mining
  - Succinct: c must be satisfied, so one can start with the data sets satisfying c
  - Convertible: c is not monotonic nor anti-monotonic, but it can be converted into it if items in the transaction can be properly ordered

## Anti-Monotonicity in Constraint Pushing

TDB (min_sup=2)

| TID | Transaction |
|-----|-------------|
| 10 | a, b, c, d, f |
| 20 | b, c, d, f, g, h |
| 30 | a, c, d, e, f |
| 40 | c, e, f, g |

| Item | Profit |
|------|--------|
| a | 40 |
| b | 0 |
| c | -20 |
| d | 10 |
| e | -30 |
| f | 30 |
| g | 20 |
| h | -10 |

- A constraint C is *antimonotone* if the super pattern satisfies C, all of its sub-patterns do so too
- In other words, *anti-monotonicity:* If an itemset S **violates** the constraint, so does any of its superset
- Ex. 1. *sum(S.price)* $\leq v$ is anti-monotone
- Ex. 2. range(S.profit) $\leq 15$ is anti-monotone
  - Itemset *ab* violates C
  - So does every superset of *ab*
- Ex. 3. *sum(S.Price)* $\geq v$ is not anti-monotone
- Ex. 4. *support count* is anti-monotone: core property used in Apriori

## Monotonicity for Constraint Pushing

TDB (min_sup=2)

| TID | Transaction |
|-----|-------------|
| 10 | a, b, c, d, f |
| 20 | b, c, d, f, g, h |
| 30 | a, c, d, e, f |
| 40 | c, e, f, g |

| Item | Profit |
|------|--------|
| a | 40 |
| b | 0 |
| c | -20 |
| d | 10 |
| e | -30 |
| f | 30 |
| g | 20 |
| h | -10 |

- A constraint C is *monotone* if the pattern satisfies C, we do not need to check C in subsequent mining
- Alternatively, monotonicity: *If an itemset S **satisfies** the constraint, so does any of its superset*
- Ex. 1. *sum(S.Price) $\geq v$* is monotone
- Ex. 2. *min(S.Price) $\leq v$* is monotone
- Ex. 3. C: range(S.profit) $\geq 15$
  - Itemset *ab* satisfies C
  - So does every superset of *ab*

73

## Data Antimonotonicity: Pruning Data Space

TDB (min_sup=2)

| TID | Transaction |
|-----|-------------|
| 10 | a, b, c, d, f, h |
| 20 | b, c, d, f, g, h |
| 30 | b, c, d, f, g |
| 40 | c, e, f, g |

| Item | Profit |
|------|--------|
| a | 40 |
| b | 0 |
| c | -20 |
| d | -15 |
| e | -30 |
| f | -10 |
| g | 20 |
| h | -5 |

- A constraint c is *data antimonotone* if for a pattern p cannot satisfy a transaction t under c, p's superset cannot satisfy t under c either
- The key for data antimonotone is *recursive data reduction*
- Ex. 1. *sum(S.Price) $\geq v$* is data antimonotone
- Ex. 2. *min(S.Price) $\leq v$* is data antimonotone
- Ex. 3. C: *range(S.profit) $\geq 25$* is data antimonotone
  - Itemset {b, c}'s projected DB:
    - T10': {d, f, h}, T20': {d, f, g, h}, T30': {d, f, g}
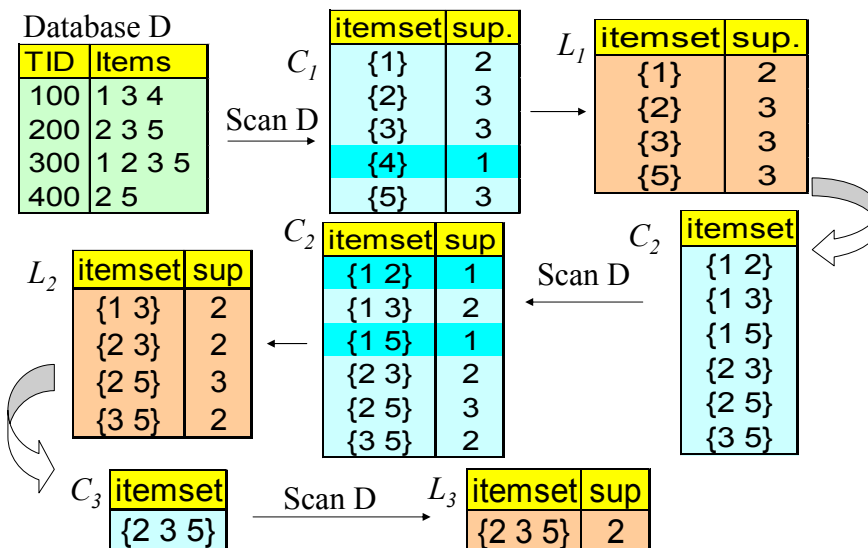  - since C cannot satisfy T10', T10' can be pruned

74

- Succinctness:
  - Given $A_1$, the set of items satisfying a succinctness constraint $C$, then any set $S$ satisfying $C$ is based on $A_1$, i.e., $S$ contains a subset belonging to $A_1$
  - Idea: Without looking at the transaction database, whether an itemset $S$ satisfies constraint C can be determined based on the selection of items
  - $min(S.Price) \leq v$ is succinct
  - $sum(S.Price) \geq v$ is not succinct
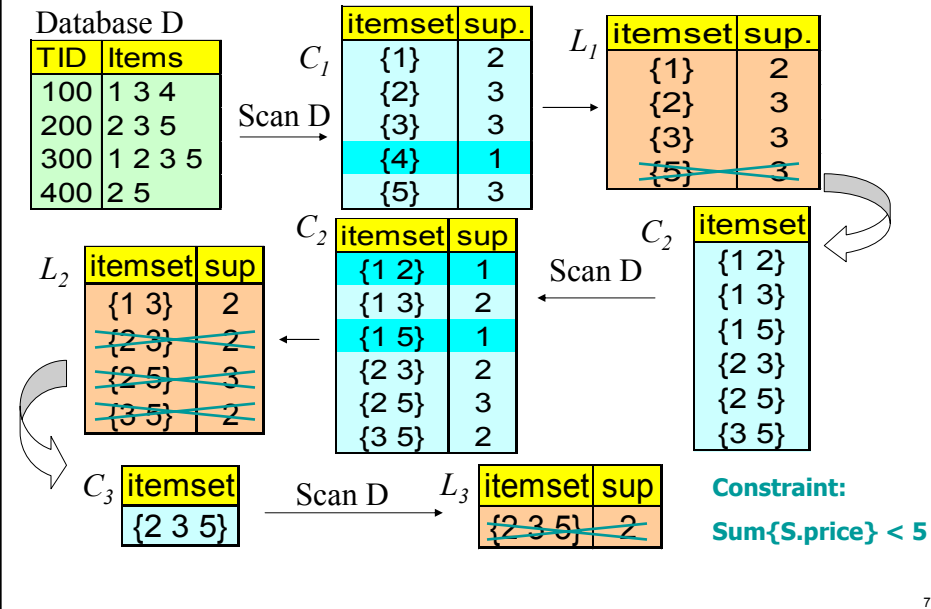- Optimization: If $C$ is succinct, $C$ is pre-counting pushable
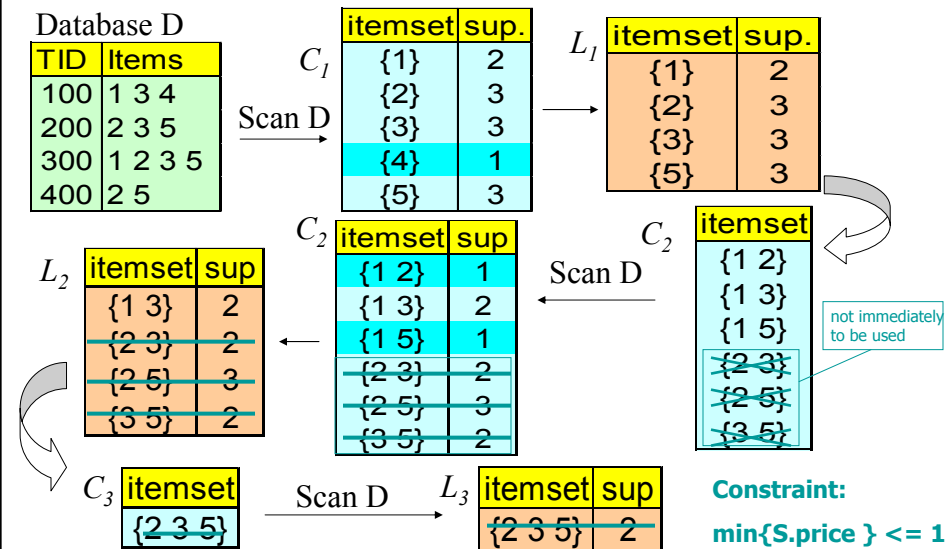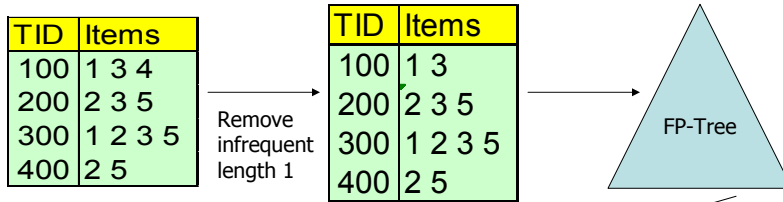
75

## 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 →

$C_1$

| 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 |

Scan D ←

$C_2$

| 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 |

76

## Naïve Algorithm: Apriori + Constraint



Database D

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

Scan D →

$C_1$

| 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 |
|---------|
| {1 2} |
| {1 3} |
| {1 5} |
| {2 3} |
| {2 5} |
| {3 5} |

Scan D ←

$C_2$

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

$L_2$

| 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 |

**Constraint:**
**Sum{S.price} < 5**

---

## The Constrained Apriori Algorithm: Push a Succinct Constraint Deep



Database D

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

Scan D →

$C_1$

| 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 |
|---------|
| {1 2} |
| {1 3} |
| {1 5} |
| {2 3} |
| {2 5} |
| {3 5} |

not immediately to be used

Scan D ←

$C_2$

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

$L_2$

| 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 |

**Constraint:**
**min{S.price } <= 1**

## The Constrained FP-Growth Algorithm: Push a Succinct Constraint Deep

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

Remove infrequent length 1 →

| TID | Items |
|-----|-------|
| 100 | 1 3 |
| 200 | 2 3 5 |
| 300 | 1 2 3 5 |
| 400 | 2 5 |

→ FP-Tree

1-Projected DB

| TID | Items |
|-----|-------|
| 100 | 3 4 |
| 300 | 2 3 5 |

No Need to project on 2, 3, or 5

**Constraint:**

**min{S.price } <= 1**

## The Constrained FP-Growth Algorithm:
### Push a Data Antimonotonic Constraint Deep

Remove from data

| TID | Items |
|-----|-------|
| 100 | 1 3 4 |
| ~~200~~ | ~~2 3 5~~ |
| 300 | 1 2 3 5 |
| ~~400~~ | ~~2 5~~ |

→

| TID | Items |
|-----|-------|
| 100 | 1 3 |
| 300 | 1 3 |

→ FP-Tree

Single branch, we are done

**Constraint:**

**min{S.price } <= 1**

## The Constrained FP-Growth Algorithm:
### Push a Data Antimonotonic Constraint Deep

| TID | Transaction |
|-----|-------------|
| 10 | a, b, c, d, f, h |
| 20 | b, c, d, f, g, h |
| 30 | b, c, d, f, g |
| 40 | a, c, e, f, g |

| TID | Transaction |
|-----|-------------|
| 10 | a, b, c, d, f, h |
| 20 | b, c, d, f, g, |
| 30 | b, c, d, f, g |
| 40 | a, c, e, f, g |

FP-Tree

Recursive Data Pruning

### B-Projected DB

| TID | Transaction |
|-----|-------------|
| 10 | a, c, d, f, h |
| 20 | c, d, f, g, h |
| 30 | c, d, f, g |

B FP-Tree

Single branch:

bcdfg: 2

| Item | Profit |
|------|--------|
| a | 40 |
| b | 0 |
| c | -20 |
| d | -15 |
| e | -30 |
| f | -10 |
| g | 20 |
| h | -5 |

**Constraint:**
**range{S.price } > 25**
**min_sup >= 2**

---

## Converting "Tough" Constraints

- Convert tough constraints into anti-monotone or monotone by properly ordering items
- Examine C: avg($S$.profit) ≥ 25
  - Order items in value-descending order
    - <$a, f, g, d, b, h, c, e$>
  - If an itemset $afb$ violates C
    - So does $afbh, afb*$
    - It becomes anti-monotone!

TDB (min_sup=2)

| TID | Transaction |
|-----|-------------|
| 10 | a, b, c, d, f |
| 20 | b, c, d, f, g, h |
| 30 | a, c, d, e, f |
| 40 | c, e, f, g |

| Item | Profit |
|------|--------|
| a | 40 |
| b | 0 |
| c | -20 |
| d | 10 |
| e | -30 |
| f | 30 |
| g | 20 |
| h | -10 |

## Strongly Convertible Constraints

- avg(X) ≥ 25 is convertible anti-monotone w.r.t. item value descending order R: *<a, f, g, d, b, h, c, e>*
  - If an itemset *af* violates a constraint C, so does every itemset with *af* as prefix, such as *afd*
- avg(X) ≥ 25 is convertible monotone w.r.t. item value ascending order $R^{-1}$: *<e, c, h, b, d, g, f, a>*
  - If an itemset *d* satisfies a constraint *C*, so does itemsets *df* and *dfa*, which having *d* as a prefix
- Thus, avg(X) ≥ 25 is strongly convertible

| Item | Profit |
|------|--------|
| a | 40 |
| b | 0 |
| c | -20 |
| d | 10 |
| e | -30 |
| f | 30 |
| g | 20 |
| h | -10 |

83

## Can Apriori Handle Convertible Constraints?

- A convertible, not monotone nor anti-monotone nor succinct constraint cannot be pushed deep into the an Apriori mining algorithm
  - Within the level wise framework, no direct pruning based on the constraint can be made
  - Itemset df violates constraint C: avg(X) >= 25
  - Since adf satisfies C, Apriori needs df to assemble adf, df cannot be pruned
- But it can be pushed into frequent-pattern growth framework!

| Item | Value |
|------|-------|
| a | 40 |
| b | 0 |
| c | -20 |
| d | 10 |
| e | -30 |
| f | 30 |
| g | 20 |
| h | -10 |

84

## Mining With Convertible Constraints

- C: avg(X) >= 25, min_sup=2
- List items in every transaction in value descending order R: <a, f, g, d, b, h, c, e>
  - C is convertible anti-monotone w.r.t. R
- Scan TDB once
  - remove infrequent items
    - Item h is dropped
  - Itemsets a and f are good, …
- Projection-based mining
  - Imposing an appropriate order on item projection
  - Many tough constraints can be converted into (anti)-monotone

| Item | Value |
|------|-------|
| a | 40 |
| f | 30 |
| g | 20 |
| d | 10 |
| b | 0 |
| h | -10 |
| c | -20 |
| e | -30 |

TDB (min_sup=2)

| TID | Transaction |
|-----|-------------|
| 10 | a, f, d, b, c |
| 20 | f, g, d, b, c |
| 30 | a, f, d, c, e |
| 40 | f, g, h, c, e |

85

## Handling Multiple Constraints

- Different constraints may require different or even conflicting item-ordering
- If there exists an order $R$ s.t. both $C_1$ and $C_2$ are convertible w.r.t. $R$, then there is no conflict between the two convertible constraints
- If there exists conflict on order of items
  - Try to satisfy one constraint first
  - Then using the order for the other constraint to mine frequent itemsets in the corresponding projected database

86

## What Constraints Are Convertible?

| Constraint | Convertible anti-monotone | Convertible monotone | Strongly convertible |
|---|---|---|---|
| avg(S) ≤ , ≥ v | Yes | Yes | Yes |
| median(S) ≤ , ≥ v | Yes | Yes | Yes |
| sum(S) ≤ v (items could be of any value, v ≥ 0) | Yes | No | No |
| sum(S) ≤ v (items could be of any value, v ≤ 0) | No | Yes | No |
| sum(S) ≥ v (items could be of any value, v ≥ 0) | No | Yes | No |
| sum(S) ≥ v (items could be of any value, v ≤ 0) | Yes | No | No |
| …… | | | |

## Constraint-Based Mining — A General Picture

| Constraint | Antimonotone | Monotone | Succinct |
|---|---|---|---|
| v ∈ S | no | yes | yes |
| S ⊇ V | no | yes | yes |
| S ⊆ V | yes | no | yes |
| min(S) ≤ v | no | yes | yes |
| min(S) ≥ v | yes | no | yes |
| max(S) ≤ v | yes | no | yes |
| max(S) ≥ v | no | yes | yes |
| count(S) ≤ v | yes | no | weakly |
| count(S) ≥ v | no | yes | weakly |
| sum(S) ≤ v ( a ∈ S, a ≥ 0 ) | yes | no | no |
| sum(S) ≥ v ( a ∈ S, a ≥ 0 ) | no | yes | no |
| range(S) ≤ v | yes | no | no |
| range(S) ≥ v | no | yes | no |
| avg(S) θ v, θ ∈ { =, ≤, ≥ } | convertible | convertible | no |
| support(S) ≥ ξ | yes | no | no |
| support(S) ≤ ξ | no | yes | no |

## A Classification of Constraints

Antimonotone

Monotone

Strongly convertible

Succinct

Convertible anti-monotone

Convertible monotone

Inconvertible

89

## Mining Frequent Patterns, Association and Correlations – Sub-Topics

- Basic concepts and a road map

- Scalable frequent itemset mining methods

- Mining various kinds of association rules

- From association to correlation analysis

- Constraint-based association mining

➡ - Mining colossal patterns

- Summary

90

## Why Mining Colossal Frequent Patterns?

- F. Zhu, X. Yan, J. Han, P. S. Yu, and H. Cheng, "Mining Colossal Frequent Patterns by Core Pattern Fusion", ICDE'07.
- We have many algorithms, but can we mine large (i.e., colossal) patterns? – such as just size around 50 to 100? Unfortunately, not!
- Why not? – the curse of "downward closure" of frequent patterns
  - The "downward closure" property
    - Any sub-pattern of a frequent pattern is frequent.
  - Example. If $(a_1, a_2, \ldots, a_{100})$ is frequent, then $a_1, a_2, \ldots, a_{100}, (a_1, a_2),$ $(a_1, a_3), \ldots, (a_1, a_{100}), (a_1, a_2, a_3), \ldots$ are all frequent! There are about $2^{100}$ such frequent itemsets!
  - No matter using breadth-first search (e.g., Apriori) or depth-first search (FPgrowth), we have to examine so many patterns
- Thus the downward closure property leads to explosion!

## Colossal Patterns: A Motivating Example

**Let's make a set of 40 transactions**
  **T1 = 1 2 3 4 ….. 39 40**
  **T2 = 1 2 3 4 ….. 39 40**
  :       .
  :         .
  :           .
  :             .
  **T40=1 2 3 4 ….. 39 40**

**Then delete the items on the diagonal**

  **T1 = 2 3 4 ….. 39 40**
  **T2 = 1 3 4 ….. 39 40**
  :    .
  :      .
  :        .
  :          .
  **T40=1 2 3 4 …… 39**

Closed/maximal patterns may partially alleviate the problem but not really solve it: We often need to mine scattered large patterns!

Let the minimum support threshold σ= 20

There are $\binom{40}{20}$ frequent patterns of size 20

Each is closed and maximal

# patterns = $\binom{n}{n/2} \approx \sqrt{2/\pi} \dfrac{2^n}{\sqrt{n}}$

The size of the answer set is exponential to n

## Colossal Pattern Set: Small but Interesting

- It is often the case that only a small number of patterns are colossal, i.e., of large size

- Colossal patterns are usually attached with greater importance than those of small pattern sizes

## Mining Colossal Patterns: Motivation and Philosophy

- Motivation: Many real-world tasks need mining colossal patterns
  - Micro-array analysis in bioinformatics (when support is low)
  - Biological sequence patterns
  - Biological/sociological/information graph pattern mining
- *No hope for completeness*
  - If the mining of mid-sized patterns is explosive in size, there is no hope to find colossal patterns efficiently by insisting "complete set" mining philosophy
- *Jumping out of the swamp of the mid-sized results*
  - What we may develop is a philosophy that may jump out of the swamp of mid-sized results that are explosive in size and jump to reach colossal patterns
- *Striving for mining almost complete colossal patterns*
  - The key is to develop a mechanism that may quickly reach colossal patterns and discover most of them

## Alas, A Show of Colossal Pattern Mining!

$T_1$ = 2 3 4 ..... 39 40
$T_2$ = 1 3 4 ..... 39 40
⋮      .
⋮      .
⋮       .
⋮        .
$T_{40}$=1 2 3 4 ...... 39
$T_{41}$= 41 42 43 ..... 79
$T_{42}$= 41 42 43 ..... 79
⋮      .
⋮      .
$T_{60}$= 41 42 43 ... 79

Let the min-support threshold σ= 20

Then there are $\binom{40}{20}$ closed/maximal frequent patterns of size 20

However, there is only one with size greater than 20, (*i.e.,* colossal):

α= {41,42,…,79} of size 39

The existing fastest mining algorithms (*e.g.,* FPClose, LCM) fail to complete running

The algorithm outputs this colossal pattern in seconds

## Methodology of Pattern-Fusion Strategy

- Pattern-Fusion traverses the tree in a bounded-breadth way
  - Always pushes down a frontier of a bounded-size candidate pool
  - Only a fixed number of patterns in the current candidate pool will be used as the starting nodes to go down in the pattern tree – thus avoids the exponential search space
- Pattern-Fusion identifies "shortcuts" whenever possible
  - Pattern growth is not performed by single-item addition but by leaps and bounded: agglomeration of multiple patterns in the pool
  - These shortcuts will direct the search down the tree much more rapidly towards the colossal patterns

## Observation: Colossal Patterns and Core Patterns

A colossal pattern α

| | |
|---|---|
| ▭ ▭ | α |
| ▭ ▭ | α₁ |
| ▭ ▭ | α₂ |
| ⋮ | |
| ▭ | αₖ |

Transaction Database D

D

Dα_k

Dα1
Dα2

**Subpatterns $\alpha_1$ to $\alpha_k$ cluster tightly around the colossal pattern α by sharing a similar support. We call such subpatterns *core patterns* of α**

---

## Robustness of Colossal Patterns

- Core Patterns

  Intuitively, for a frequent pattern α, a subpattern β is a τ-core
  pattern of α if β shares a similar support set with α, i.e.,

  $$\frac{|D_\alpha|}{|D_\beta|} \geq \tau \qquad 0 < \tau \leq 1$$

  where τ is called the core ratio

- Robustness of Colossal Patterns

  A colossal pattern is robust in the sense that it tends to have much
  more core patterns than small patterns

## Example: Core Patterns

- A colossal pattern has far more core patterns than a small-sized pattern
- A colossal pattern has far more core descendants of a smaller size c
- A random draw from a complete set of pattern of size c would more likely to pick a core descendant of a colossal pattern
- A colossal pattern can be generated by merging a set of core patterns

| Transaction (# of Ts) | Core Patterns ($\tau$ = 0.5) |
|---|---|
| (abe) (100) | (abe), (ab), (be), (ae), (e) |
| (bcf) (100) | (bcf), (bc), (bf) |
| (acf) (100) | (acf), (ac), (af) |
| (abcef) (100) | (ab), (ac), (af), (ae), (bc), (bf), (be) (ce), (fe), (e), (abc), (abf), (abe), (ace), (acf), (afe), (bcf), (bce), (bfe), (cfe), (abcf), (abce), (bcfe), (acfe), (abfe), (abcef) |

## Robustness of Colossal Patterns

- (d,τ)-robustness: A pattern α is *(d, τ)-robust* if *d* is the maximum number of items that can be removed from α for the resulting pattern to remain a τ-core pattern of α
- For a (d,τ)-robust pattern α, it has $\Omega(2^d)$ core patterns
  - » A colossal patterns tend to have a large number of core patterns
- Pattern distance: For patterns α and β, the pattern distance of α and β is defined to be

$$Dist(\alpha, \beta) = 1 - \frac{|D_\alpha \cap D_\beta|}{|D_\alpha \cup D_\beta|}$$

- If two patterns α and β are both core patterns of a same pattern, they would be bounded by a "ball" of a radius specified by their core ratio τ

$$Dist(\alpha, \beta) \leq 1 - \frac{1}{2/\tau - 1} = r(\tau)$$

- Once we identify one core pattern, we will be able to find all the other core patterns by a bounding ball of radius r(τ)

## Colossal Patterns Correspond to Dense Balls

- Due to their robustness, colossal patterns correspond to dense balls
  - $\Omega(2^d)$ in population
- A random draw in the pattern space will hit somewhere in the ball with high probability

## Idea of Pattern-Fusion Algorithm

- Generate a complete set of frequent patterns up to a small size
- Randomly pick a pattern β, and β has a high probability to be a core-descendant of some colossal pattern α
- Identify all α's descendants in this complete set, and merge all of them — This would generate a much larger core-descendant of α
- In the same fashion, we select K patterns. This set of larger core-descendants will be the candidate pool for the next iteration

## Pattern-Fusion: The Algorithm

- Initialization (Initial pool): Use an existing algorithm to mine all frequent patterns up to a small size, e.g., 3
- Iteration (Iterative Pattern Fusion):
  - At each iteration, k seed patterns are randomly picked from the current pattern pool
  - For each seed pattern thus picked, we find all the patterns within a bounding ball centered at the seed pattern
  - All these patterns found are fused together to generate a set of super-patterns. All the super-patterns thus generated form a new pool for the next iteration
- Termination: when the current pool contains no more than K patterns at the beginning of an iteration

## Why Is Pattern-Fusion Efficient?

- A bounded-breadth pattern tree traversal
  - It avoids explosion in mining mid-sized ones
  - Randomness comes to help to stay on the right path
- Ability to identify "short-cuts" and take "leaps"
  - fuse small patterns together in one step to generate new patterns of significant sizes
  - Efficiency

## Pattern-Fusion Leads to Good Approximation

- Gearing toward colossal patterns
  - The larger the pattern, the greater the chance it will be generated
- Catching outliers
  - The more distinct the pattern, the greater the chance it will be generated

## Experimental Setting

- Synthetic data set
  - $Diag_n$ an n x (n-1) table where $i^{th}$ row has integers from 1 to n except i. Each row is taken as an itemset. min_support is n/2.
- Real data set
  - Replace: A program trace data set collected from the "replace" program, widely used in software engineering research
  - ALL: A popular gene expression data set, a clinical data on ALL-AML leukemia (www.broad.mit.edu/tools/data.html).
    - Each item is a column, representing the activitiy level of gene/protein in the same
    - Frequent pattern would reveal important correlation between gene expression patterns and disease outcomes

## Experiment Results on Diag$_n$

- LCM run time increases exponentially with pattern size n
- Pattern-Fusion finishes efficiently
- The approximation error of Pattern-Fusion (with min-sup 20) in comparison with the complete set) is rather close to uniform sampling (which randomly picks K patterns from the complete answer set)
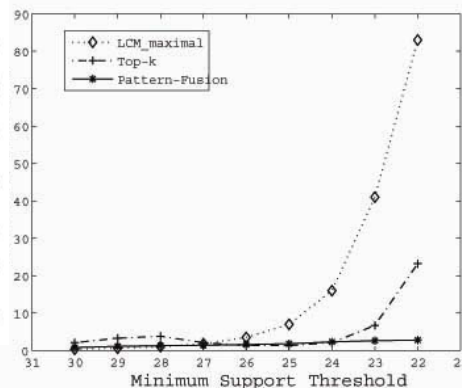


107

## Experimental Results on ALL

- ALL: A popular gene expression data set with 38 transactions, each with 866 columns
  - There are 1736 items in total
  - The table shows a high frequency threshold of 30

| Pattern Size | 110 | 107 | 102 | 91 | 86 | 84 | 83 |
|---|---|---|---|---|---|---|---|
| The complete set | 1 | 1 | 1 | 1 | 1 | 2 | 6 |
| Pattern-Fusion | 1 | 1 | 1 | 1 | 1 | 1 | 4 |
| Pattern Size | 82 | 77 | 76 | 75 | 74 | 73 | 71 |
| The complete set | 1 | 2 | 1 | 1 | 1 | 2 | 1 |
| Pattern-Fusion | 0 | 2 | 0 | 1 | 1 | 1 | 1 |



108

- REPLACE
  - A program trace data set, recording 4395 calls and transitions
  - The data set contains 4395 transactions with 57 items in total
  - With support threshold of 0.03, the largest patterns are of size 44
  - They are all discovered by Pattern-Fusion with different settings of K and τ, when started with an initial pool of 20948 patterns of size <=3
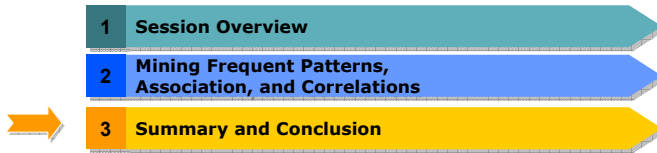
109

- Approximation error when compared with the complete mining result
- Example. Out of the total 98 patterns of size >=42, when K=100, Pattern-Fusion returns 80 of them
- A good approximation to the colossal patterns in the sense that any pattern in the complete set is on average at most 0.17 items away from one of these 80 patterns



110

## Agenda

| | |
|---|---|
| 1 | **Session Overview** |
| 2 | **Mining Frequent Patterns, Association, and Correlations** |
| 3 | **Summary and Conclusion** |

## Frequent-Pattern Mining: Summary

- Frequent pattern mining—an important task in data mining

- Scalable frequent pattern mining methods
    - Apriori (Candidate generation & test)
    - Projection-based (FPgrowth, CLOSET+, ...)
    - Vertical format approach (CHARM, ...)

- Mining a variety of rules and interesting patterns

- Constraint-based mining

- Mining sequential and structured patterns

- Extensions and applications

- Mining fault-tolerant frequent, sequential and structured patterns
  - Patterns allows limited faults (insertion, deletion, mutation)
- Mining truly interesting patterns
  - Surprising, novel, concise, …
- Application exploration
  - E.g., DNA sequence analysis and bio-pattern classification
  - "Invisible" data mining

---

**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:02.
- 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.
- G. Liu, H. Lu, W. Lou, J. X. Yu. On Computing, Storing and Querying Frequent Patterns. KDD'03.
- G. Grahne and J. Zhu, Efficiently Using Prefix-Trees in Mining Frequent Itemsets, Proc. ICDM'03 Int. Workshop on Frequent Itemset Mining Implementations (FIMI'03), Melbourne, FL, Nov. 2003

## 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.
- Zaki and 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 Multi-Level and Quantitative Rules

- R. Srikant and R. Agrawal. Mining generalized association rules. VLDB'95.
- J. Han and Y. Fu. Discovery of multiple-level association rules from large databases. VLDB'95.
- R. Srikant and R. Agrawal. Mining quantitative association rules in large relational tables. SIGMOD'96.
- T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama. Data mining using two-dimensional optimized association rules: Scheme, algorithms, and visualization. SIGMOD'96.
- K. Yoda, T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama. Computing optimized rectilinear regions for association rules. KDD'97.
- R.J. Miller and Y. Yang.  Association rules over interval data. SIGMOD'97.
- Y. Aumann and Y. Lindell. A Statistical Theory for Quantitative Association Rules KDD'99.

## Ref: Mining Correlations and Interesting Rules

- M. Klemettinen, H. Mannila, P. Ronkainen, H. Toivonen, and A. I. Verkamo.  Finding interesting rules from large sets of discovered association rules.  CIKM'94.
- S. Brin, R. Motwani, and C. Silverstein.  Beyond market basket: Generalizing association rules to correlations.  SIGMOD'97.
- 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, "Association Mining in Large Databases: A Re-Examination of Its Measures", PKDD'07

## Ref: Mining Other Kinds of Rules

- R. Meo, G. Psaila, and S. Ceri.  A new SQL-like operator for mining association rules. VLDB'96.
- B. Lent, A. Swami, and J. Widom.  Clustering association rules. ICDE'97.
- A. Savasere, E. Omiecinski, and S. Navathe.  Mining for strong negative associations in a large database of customer transactions. ICDE'98.
- D. Tsur, J. D. Ullman, S. Abitboul, C. Clifton, R. Motwani, and S. Nestorov.  Query flocks: A generalization of association-rule mining. SIGMOD'98.
- F. Korn, A. Labrinidis, Y. Kotidis, and C. Faloutsos.  Ratio rules: A new paradigm for fast, quantifiable data mining. VLDB'98.
- F. Zhu, X. Yan, J. Han, P. S. Yu, and H. Cheng, "Mining Colossal Frequent Patterns by Core Pattern Fusion", ICDE'07.

## Ref: Constraint-Based Pattern Mining

- R. Srikant, Q. Vu, and R. Agrawal. Mining association rules with item constraints. KDD'97
- R. Ng, L.V.S. Lakshmanan, J. Han & A. Pang. Exploratory mining and pruning optimizations of constrained association rules. SIGMOD'98
- G. Grahne, L. Lakshmanan, and X. Wang. Efficient mining of constrained correlated sets. ICDE'00
- J. Pei, J. Han, and L. V. S. Lakshmanan. Mining Frequent Itemsets with Convertible Constraints. ICDE'01
- J. Pei, J. Han, and W. Wang, Mining Sequential Patterns with Constraints in Large Databases, CIKM'02
- F. Bonchi, F. Giannotti, A. Mazzanti, and D. Pedreschi. ExAnte: Anticipated Data Reduction in Constrained Pattern Mining, PKDD'03
- F. Zhu, X. Yan, J. Han, and P. S. Yu, "gPrune: A Constraint Pushing Framework for Graph Pattern Mining", PAKDD'07

## Ref: Mining Sequential and Structured Patterns

- R. Srikant and R. Agrawal. Mining sequential patterns: Generalizations and performance improvements. EDBT'96.
- H. Mannila, H Toivonen, and A. I. Verkamo. Discovery of frequent episodes in event sequences. DAMI:97.
- M. Zaki. SPADE: An Efficient Algorithm for Mining Frequent Sequences. Machine Learning:01.
- J. Pei, J. Han, H. Pinto, Q. Chen, U. Dayal, and M.-C. Hsu. PrefixSpan: Mining Sequential Patterns Efficiently by Prefix-Projected Pattern Growth. ICDE'01.
- M. Kuramochi and G. Karypis. Frequent Subgraph Discovery. ICDM'01.
- X. Yan, J. Han, and R. Afshar. CloSpan: Mining Closed Sequential Patterns in Large Datasets. SDM'03.
- X. Yan and J. Han. CloseGraph: Mining Closed Frequent Graph Patterns. KDD'03.

## Ref: Mining Spatial, Multimedia, and Web Data

- K. Koperski and J. Han, Discovery of Spatial Association Rules in Geographic Information Databases,  SSD'95.
- O. R. Zaiane, M. Xin, J. Han, Discovering Web Access Patterns and Trends by Applying OLAP and Data Mining Technology on Web Logs. ADL'98.
- O. R. Zaiane, J. Han, and H. Zhu, Mining Recurrent Items in Multimedia with Progressive Resolution Refinement. ICDE'00.
- D. Gunopulos and I. Tsoukatos.  Efficient Mining of Spatiotemporal Patterns.   SSTD'01.

## Ref: Mining Frequent Patterns in Time-Series Data

- B. Ozden, S. Ramaswamy, and A. Silberschatz. Cyclic association rules. ICDE'98.
- J. Han, G. Dong and Y. Yin, Efficient Mining of Partial Periodic Patterns in Time Series Database, ICDE'99.
- H. Lu, L. Feng, and J. Han.  Beyond Intra-Transaction Association Analysis: Mining Multi-Dimensional Inter-Transaction Association Rules. TOIS:00.
- B.-K. Yi, N. Sidiropoulos, T. Johnson, H. V. Jagadish, C. Faloutsos, and A. Biliris. Online Data Mining for Co-Evolving Time Sequences. ICDE'00.
- W. Wang, J. Yang, R. Muntz. TAR: Temporal Association Rules on Evolving Numerical Attributes. ICDE'01.
- J. Yang, W. Wang, P. S. Yu. Mining Asynchronous Periodic Patterns in Time Series Data. TKDE'03.

## Ref: Iceberg Cube and Cube Computation

- S. Agarwal, R. Agrawal, P. M. Deshpande, A. Gupta, J. F. Naughton, R. Ramakrishnan, and S. Sarawagi. On the computation of multidimensional aggregates. VLDB'96.
- Y. Zhao, P. M. Deshpande, and J. F. Naughton. An array-based algorithm for simultaneous multidi-mensional aggregates. SIGMOD'97.
- J. Gray, et al. Data cube: A relational aggregation operator generalizing group-by, cross-tab and sub-totals. DAMI: 97.
- M. Fang, N. Shivakumar, H. Garcia-Molina, R. Motwani, and J. D. Ullman. Computing iceberg queries efficiently. VLDB'98.
- S. Sarawagi, R. Agrawal, and N. Megiddo. Discovery-driven exploration of OLAP data cubes. EDBT'98.
- K. Beyer and R. Ramakrishnan. Bottom-up computation of sparse and iceberg cubes. SIGMOD'99.

## Ref: Iceberg Cube and Cube Exploration

- J. Han, J. Pei, G. Dong, and K. Wang, Computing Iceberg Data Cubes with Complex Measures. SIGMOD' 01.
- W. Wang, H. Lu, J. Feng, and J. X. Yu. Condensed Cube: An Effective Approach to Reducing Data Cube Size. ICDE'02.
- G. Dong, J. Han, J. Lam, J. Pei, and K. Wang. Mining Multi-Dimensional Constrained Gradients in Data Cubes. VLDB'01.
- T. Imielinski, L. Khachiyan, and A. Abdulghani. Cubegrades: Generalizing association rules. DAMI:02.
- L. V. S. Lakshmanan, J. Pei, and J. Han. Quotient Cube: How to Summarize the Semantics of a Data Cube. VLDB'02.
- D. Xin, J. Han, X. Li, B. W. Wah. Star-Cubing: Computing Iceberg Cubes by Top-Down and Bottom-Up Integration. VLDB'03.

## Ref: FP for Classification and Clustering

- G. Dong and J. Li. Efficient mining of emerging patterns: Discovering trends and differences. KDD'99.
- B. Liu, W. Hsu, Y. Ma. Integrating Classification and Association Rule Mining. KDD'98.
- W. Li, J. Han, and J. Pei. CMAR: Accurate and Efficient Classification Based on Multiple Class-Association Rules. ICDM'01.
- H. Wang, W. Wang, J. Yang, and P.S. Yu. Clustering by pattern similarity in large data sets. SIGMOD′02.
- J. Yang and W. Wang. CLUSEQ: efficient and effective sequence clustering. ICDE'03.
- X. Yin and J. Han. CPAR: Classification based on Predictive Association Rules. SDM'03.
- H. Cheng, X. Yan, J. Han, and C.-W. Hsu, Discriminative Frequent Pattern Analysis for Effective Classification", ICDE'07.

## Ref: Stream and Privacy-Preserving FP Mining

- A. Evfimievski, R. Srikant, R. Agrawal, J. Gehrke. Privacy Preserving Mining of Association Rules. KDD′02.
- J. Vaidya and C. Clifton. Privacy Preserving Association Rule Mining in Vertically Partitioned Data. KDD′02.
- G. Manku and R. Motwani. Approximate Frequency Counts over Data Streams. VLDB′02.
- Y. Chen, G. Dong, J. Han, B. W. Wah, and J. Wang. Multi-Dimensional Regression Analysis of Time-Series Data Streams. VLDB'02.
- C. Giannella, J. Han, J. Pei, X. Yan and P. S. Yu. Mining Frequent Patterns in Data Streams at Multiple Time Granularities, Next Generation Data Mining:03.
- A. Evfimievski, J. Gehrke, and R. Srikant. Limiting Privacy Breaches in Privacy Preserving Data Mining. PODS'03.

## Ref: Other Freq. Pattern Mining Applications

- Y. Huhtala, J. Kärkkäinen, P. Porkka, H. Toivonen. Efficient Discovery of Functional and Approximate Dependencies Using Partitions. ICDE'98.

- H. V. Jagadish, J. Madar, and R. Ng. Semantic Compression and Pattern Extraction with Fascicles. VLDB'99.

- T. Dasu, T. Johnson, S. Muthukrishnan, and V. Shkapenyuk. Mining Database Structure; or How to Build a Data Quality Browser. SIGMOD'02.

- K. Wang, S. Zhou, J. Han.  Profit Mining: From Patterns to Actions. EDBT'02.

## 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)

## Assignments & Readings

- Readings
  - » Chapter 5
- Individual Project #1
  - » Ongoing

## Next Session: Classification and Prediction