# Big Data:
# From Theory to Practice

**Xing Wu**

xingwu@shu.edu.cn

**Shanghai University**

# Mining Frequent Patterns, Association and Correlations

- Basic concepts

- Efficient and scalable frequent itemset mining methods

- Mining various kinds of association rules

- From association mining to correlation analysis

- Constraint-based association mining
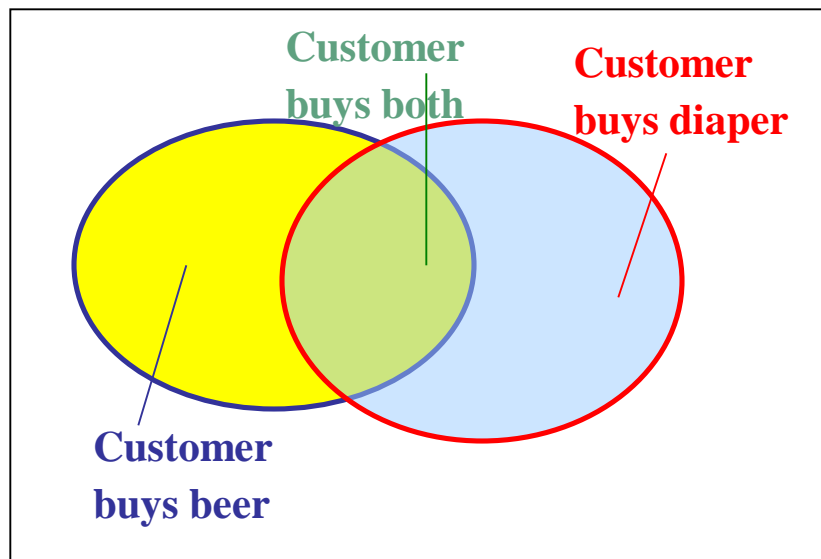
- Summary

# 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?

- Discloses an intrinsic and important property of data sets
- Forms the 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: associative classification
  - Cluster analysis: frequent pattern-based clustering
  - Broad applications

# Basic Concepts: Frequent Patterns and Association Rules

| Transaction-id | Items bought |
|:---:|:---:|
| 10 | A, B, D |
| 20 | A, C, D |
| 30 | A, D, E |
| 40 | B, E, F |
| 50 | B, C, D, E, F |



Customer buys both

Customer buys diaper

Customer buys beer

- Itemset X = $\{x_1, ..., x_k\}$
- 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* $sup_{min} = 50\%$, $conf_{min} = 50\%$
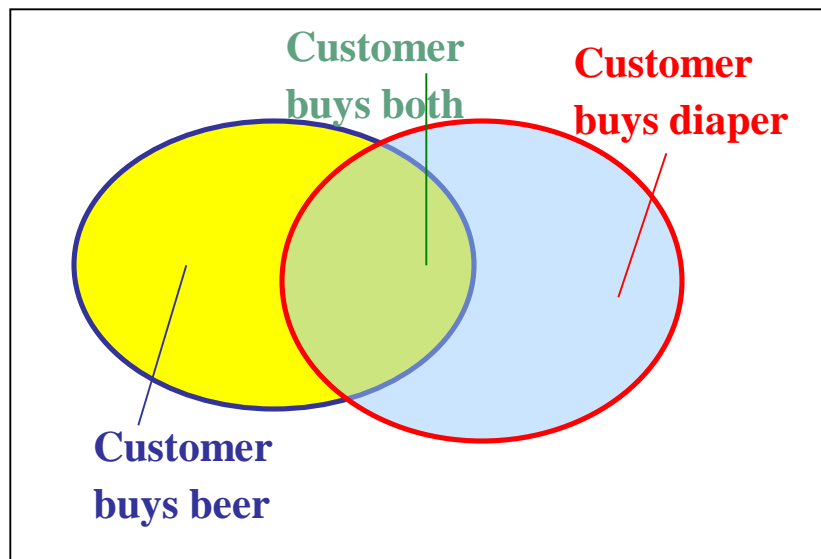**Freq. Pat**.: $\{A:3, B:3, D:4, E:3, AD:3\}$
Association rules:

    ***A → D* (?, ?)**
    ***D → A* (?, ?)**

# Basic Concepts: Frequent Patterns and Association Rules

| Transaction-id | Items bought |
|---|---|
| 10 | A, B, D |
| 20 | A, C, D |
| 30 | A, D, E |
| 40 | B, E, F |
| 50 | B, C, D, E, F |

**Customer buys both**

**Customer buys diaper**

**Customer buys beer**

- Itemset X = $\{x_1, ..., x_k\}$
- 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* $sup_{min}$ = 50%, $conf_{min}$ = 50%
**Freq. Pat.**: {A:3, B:3, D:4, E:3, AD:3}
Association rules:
  **A → D (60%, 100%)**
  **D → 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 $\binom{100}{1} + \binom{100}{2} + \ldots + \binom{100}{100} = 2^{100} - 1 = 1.27 \times 10^{30}$ 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 ⊃ 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 ⊃ 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 <span style="color:red">closed itemset</span>?

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

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

  - !!

# Mining Frequent Patterns, Association and Correlations

- ■ Basic concepts

- ■ Efficient and scalable frequent itemset mining methods ⟵

- ■ Mining various kinds of association rules

- ■ From association mining to correlation analysis

- ■ Constraint-based association mining

- ■ Summary

# Scalable Methods for Mining Frequent Patterns

- 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}**
  - **Why?**

- Scalable mining methods: Three major approaches
  - Apriori (Agrawal & Srikant@VLDB'94)
  - Freq. pattern growth (FPgrowth—Han, Pei & Yin @SIGMOD'00)

# Scalable Methods for Mining Frequent Patterns

- 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)

# Apriori: A Candidate Generation-and-Test Approach

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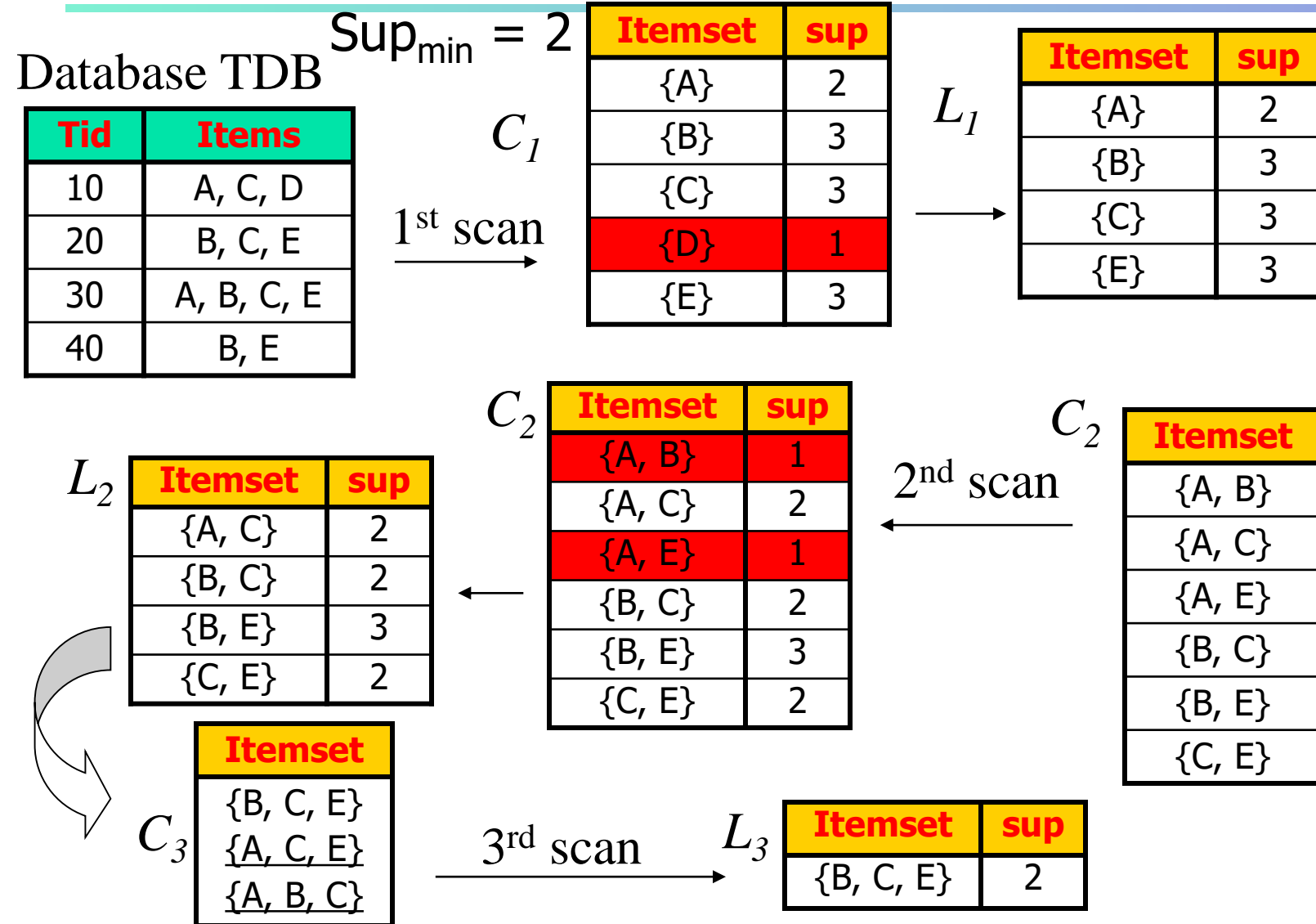
- Method:

   Can we use only smaller itemsets to generate larger ones rather than explore all larger ones?

# Apriori: A Candidate Generation-and-Test Approach

- Apriori pruning principle: If there is any itemset which is infrequent, its superset should not be generated/tested! **Why?** (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 |

$1^{st}$ 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 |

$2^{nd}$ 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} |
| {A, C, E} |
| {A, B, C} |

$3^{rd}$ 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** $\cup_k L_k$;

# Important Details of Apriori

- How to generate candidates?
  - Step 1: self-joining $L_k$
  - Step 2: pruning
- How to count supports of candidates?
- 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:
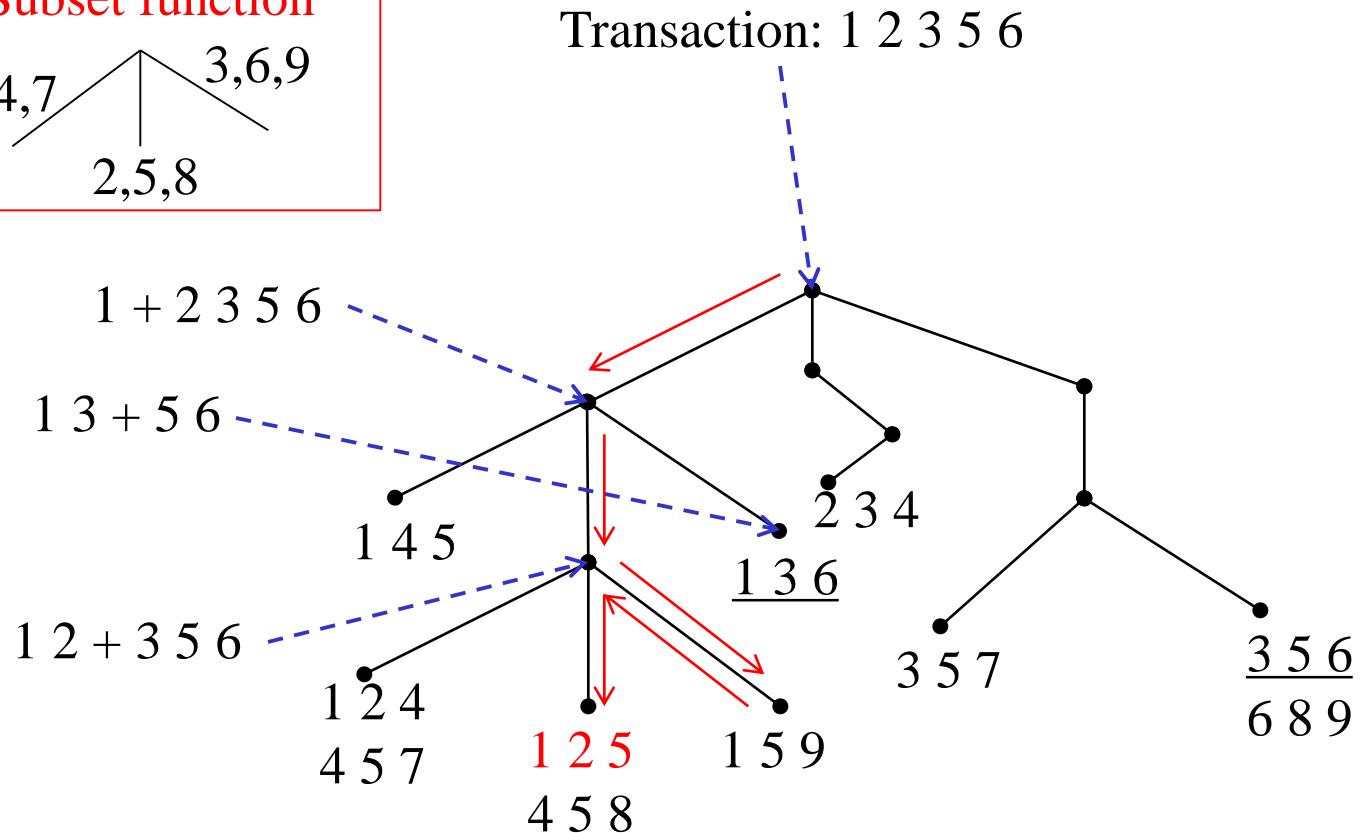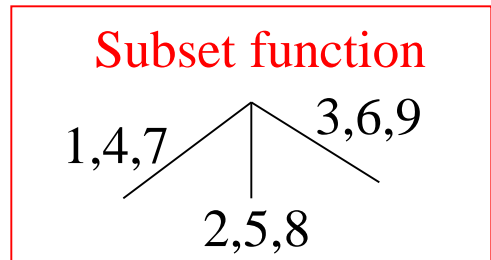    - *acde* is removed because *ade* is not in $L_3$
  - $C_4$={*abcd*}

# How to Generate Candidates?

- Suppose the items in $L_{k-1}$ are listed in an order
- Step 1: self-joining $L_{k-1}$

  insert into $C_k$

  select $p.item_1, p.item_2, ..., p.item_{k-1}, q.item_{k-1}$

  from $L_{k-1}\ p, L_{k-1}\ q$

  where $p.item_1 = q.item_1, ..., p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$

- Step 2: pruning

  forall *itemsets c in $C_k$* do

      forall *(k-1)-subsets s of c* do

          **if** *(s is not in $L_{k-1}$)* **then delete** *c* **from** $C_k$

# How to Count Supports of Candidates?

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

# Example: Counting Supports of Candidates

Subset function

1,4,7    2,5,8    3,6,9

Transaction: 1 2 3 5 6

1 + 2 3 5 6

1 3 + 5 6

1 4 5

1 2 + 3 5 6

1 2 4
4 5 7

1 2 5
4 5 8

1 5 9

2 3 4

1 3 6

3 5 7

3 5 6
6 8 9

**3-item candidates**
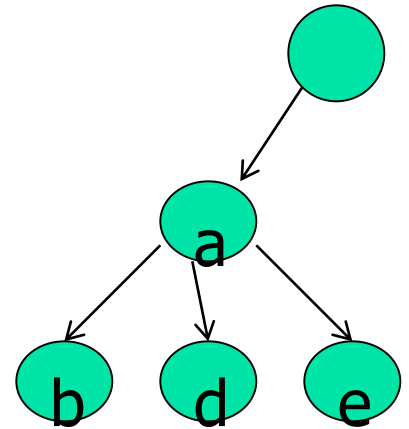
# Challenges of Frequent Pattern Mining

- 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 if the support is defined in terms of percent of transaction

  - Scan 1: partition database and find local frequent patterns

  - Scan 2: consolidate global frequent patterns

- A. Savasere, E. Omiecinski, and S. Navathe. An efficient algorithm for mining association in large databases. In *VLDB'95*
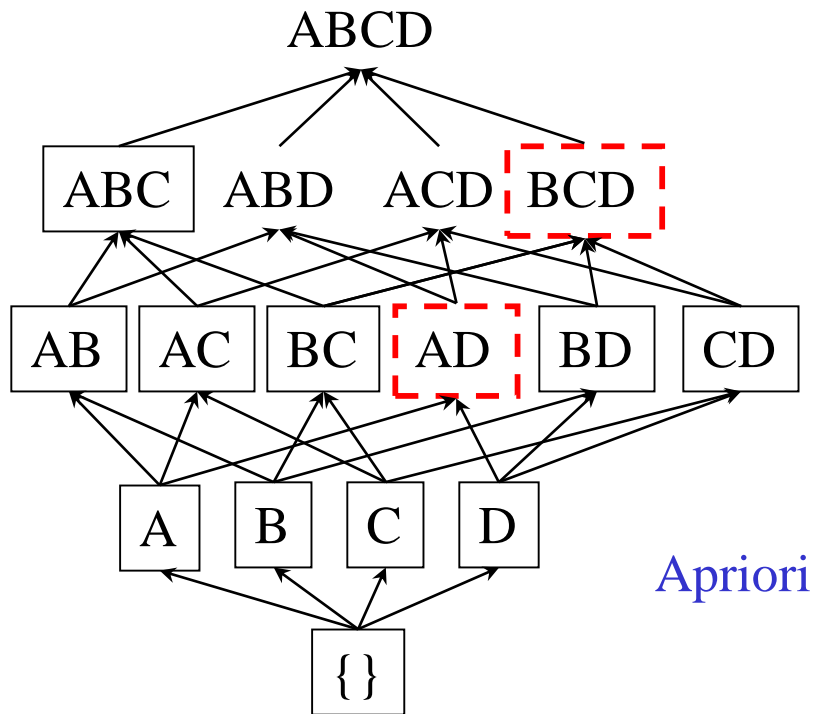
# DHP: Reduce the Number of Candidates

- 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*
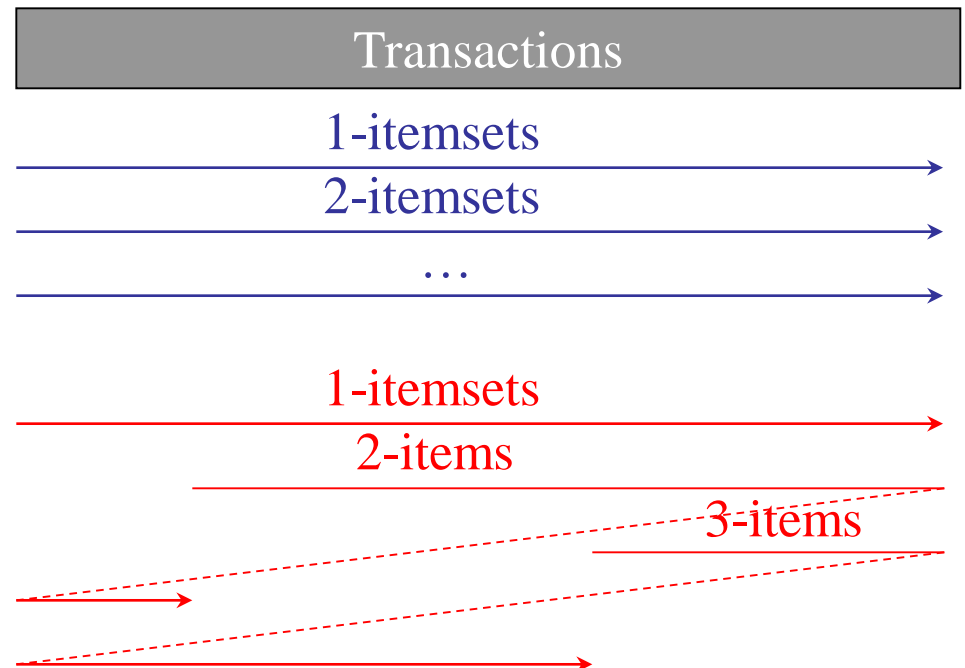
# DIC: Reduce Number of Scans



Itemset lattice

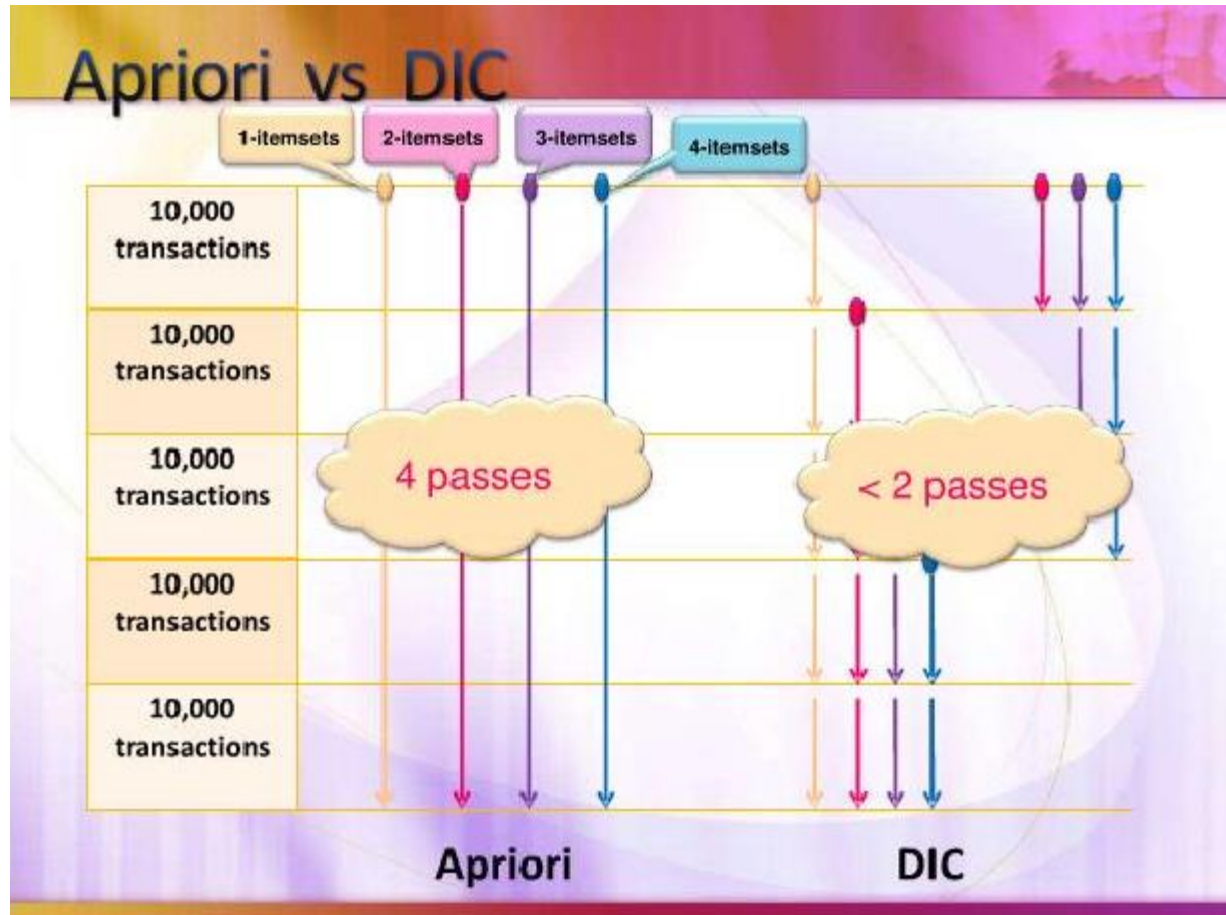S. Brin R. Motwani, J. Ullman, and S. Tsur. Dynamic itemset counting and implication rules for market basket data. In *SIGMOD'97*

- 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

Provided by Kiran

# Bottleneck of Frequent-pattern Mining

- Multiple database scans are <span style="color:red">costly</span>

- Mining long patterns needs many passes of scanning and generates lots of candidates

  - To find frequent itemset $i_1 i_2 \ldots i_{100}$

    - \# of scans: <span style="color:red">100</span>

    - \# of Candidates: $\binom{100}{1} + \binom{100}{2} + \ldots + \binom{100}{100} = 2^{100} - 1 = $ <span style="color:red">$1.27 * 10^{30}$</span> !

- Bottleneck: candidate-generation-and-test

- Can we avoid candidate generation?

# Mining Frequent Patterns Without Candidate Generation

- Grow long patterns from short ones using local frequent items

  - "abc" is a frequent pattern

  - Get all transactions having "abc": DB|abc

  - "d" is a local frequent item (in term of count of occurrences) in DB|abc → abcd is a frequent pattern

# Construct FP-tree from a Transaction Database

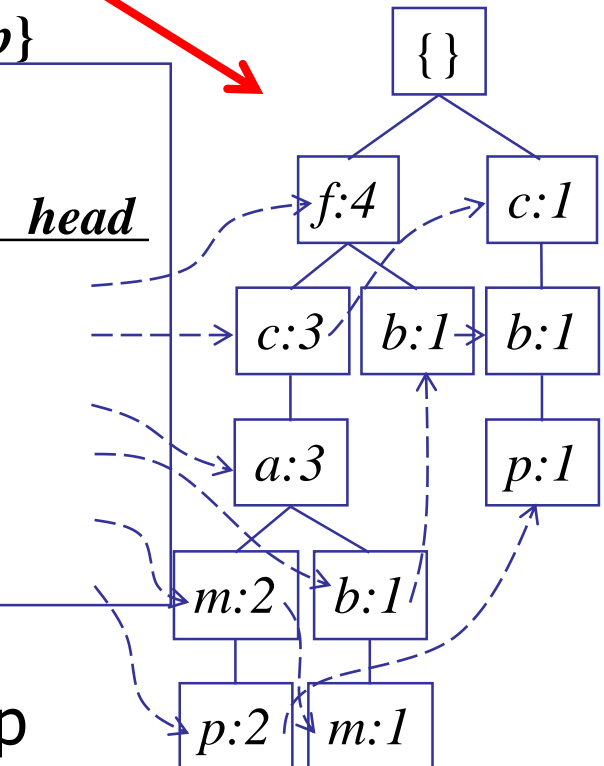| TID | Items bought | (ordered) frequent items |
|-----|--------------|--------------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o, w} | {f, b} |
| 400 | {b, c, k, s, p} | {c, b, p} |
| 500 | {a, f, c, e, l,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

*min_support = 3*

**Prefix Tree**

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:4    c:1

c:3   b:1   b:1

a:3         p:1

m:2   b:1

p:2   m:1

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

# 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)
  - For Connect-4 DB, compression ratio could be over 100

# Partition Patterns and Databases

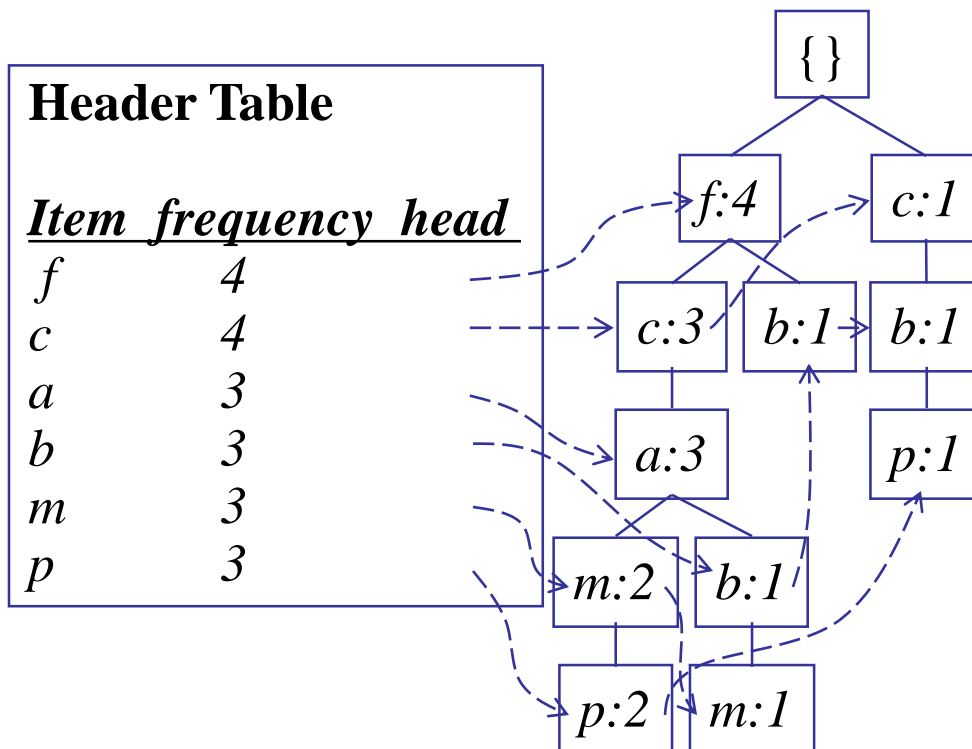- Frequent patterns can be partitioned into subsets according to f-list **Peeling of Onion**

    - 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, no others

- Completeness and non-redundency**?**

| F only No others | | | All with b | All with m | All with P |
|---|---|---|---|---|---|

■ Starting at the frequent item header table in the FP-tree

**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

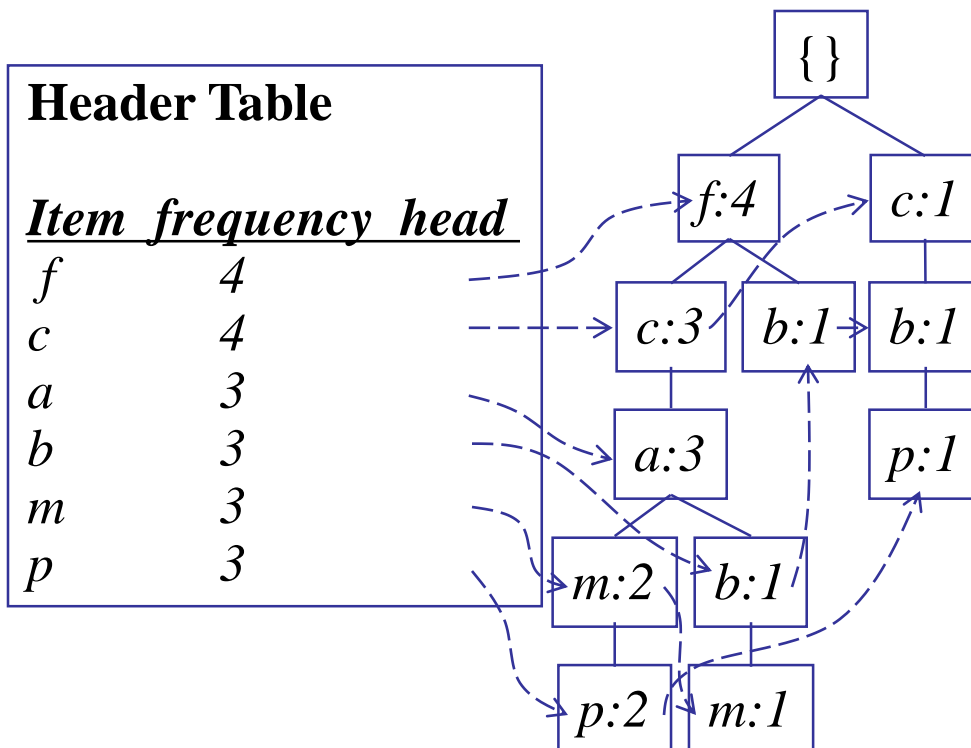p:2  m:1

*Output Frequent Items:*

*f, c, a, b, m, p*

*Use each of them as a condition to partition data:*

*Collect all prefixes end at each node*

# Generate Frequent Item Sets Using Conditional Database Recursively – Step 1

- Starting at the frequent item header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item *x*
- Accumulate all of *prefix paths* of item *x* to form *x*'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* |
|--------|----------------------|
| *f* | *{}* |
| *c* | *f:3* |
| *a* | *fc:3* |
| *b* | *fca:1, f:1, c:1* |
| *m* | *fca:2, fcab:1* |
| *p* | *fcam:2, cb:1* |

Recursion

# Construct FP Tree for Each Conditional Database

*Conditional* **pattern bases**

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

Output frequent
1-item set

Empty, no item, not tree, stop

Header table: F    3    {}    fc {}

Output: cf

f:3

Header Table:  f    3    {}
               c    3    fa    {}

Output: af, ac

f:3

c:3    ca    f:3

Header Table: f    3

Output: acf    fca    {}

# Construct FP Tree for Each Conditional Database

*Conditional* **pattern bases**

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

Header Table: f    2
                      c    2
                      a    1
None of them is frequent, stop!

*Conditional* **pattern bases**

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

Header Table: f  3
              c  3
              a  3

Output:
mf, mc, ma

{}
↓
f:3 → fm: {}
↓
c:3 → cm: f:3 → Header Table: {}
                 f  3
                 Output: fcm
↓
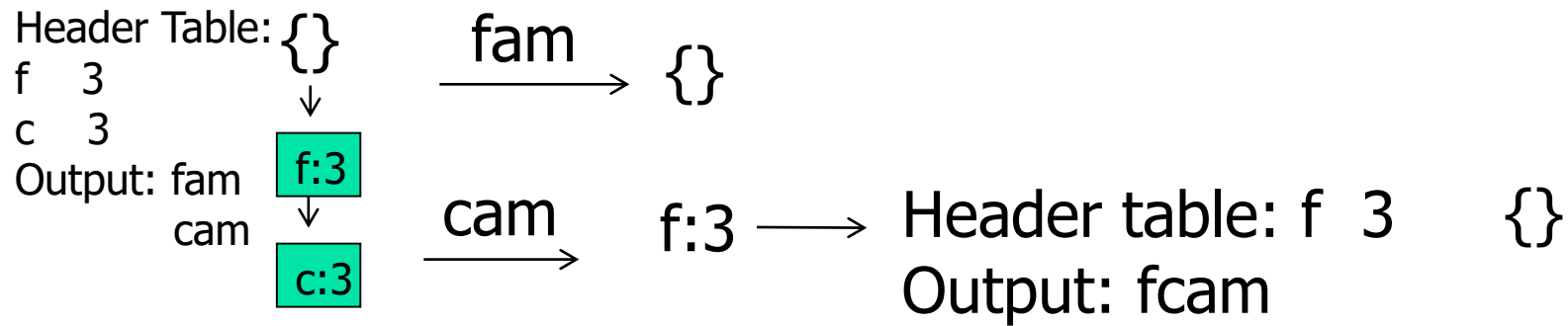a:3 → am: fc:3    Header Table: {}
                  f  3
                  c  3
                  Output: fam
                         cam

f:3
↓
c:3

# Construct FP Tree for Each Conditional Database

Header Table:
f    3
c    3
Output: fam
        cam

{}
 ↓
f:3
 ↓
c:3

—— fam ——→  {}

—— cam ——→  f:3 ——→  Header table: f   3     {}
                      Output: fcam

*Conditional* **pattern bases**

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

Header Table: c   3
Output: cp

{}

$\xrightarrow{\text{cp}}$

{}

c

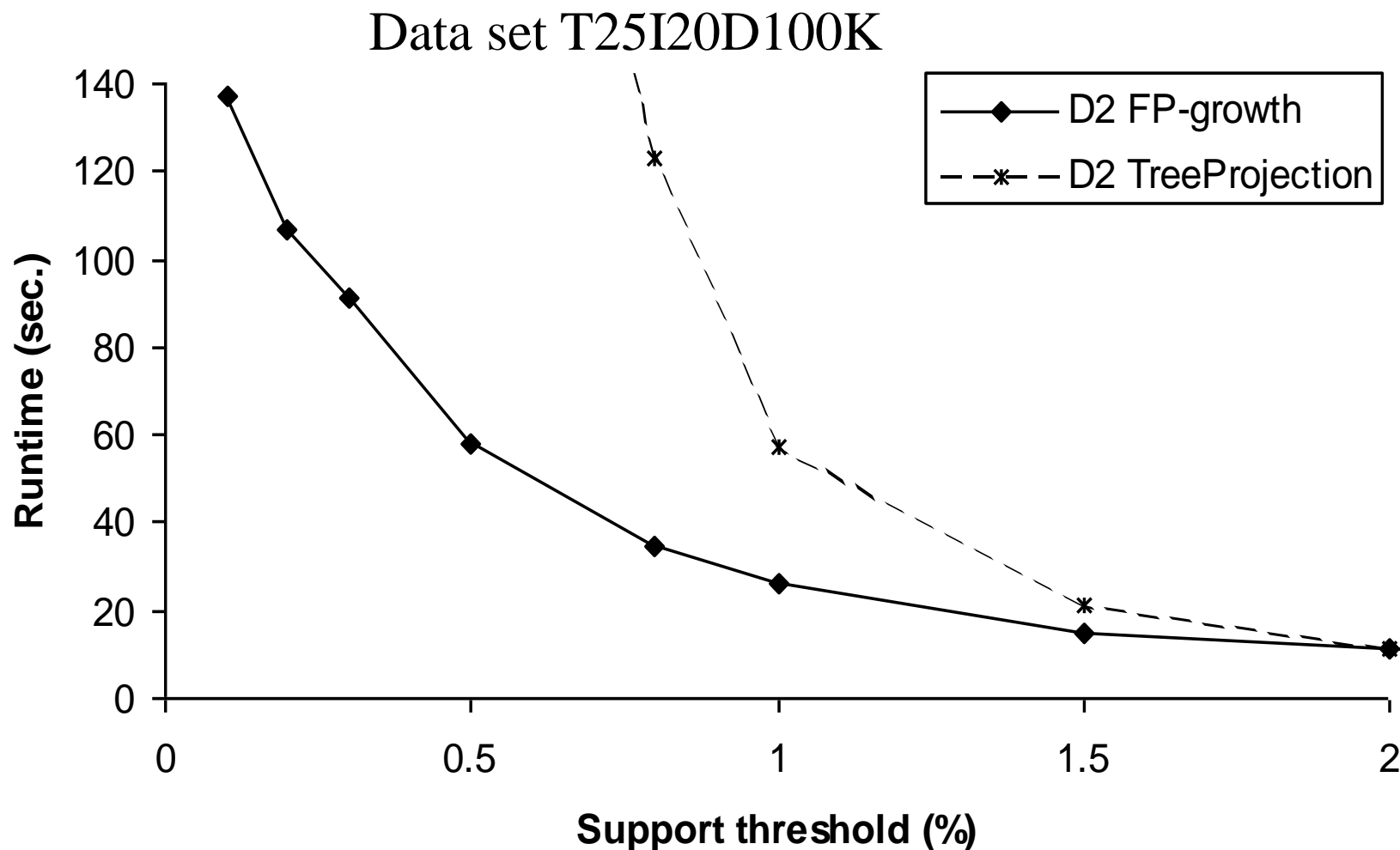# Mining Frequent Patterns With FP-trees

- 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
  - Output frequent patterns found at the current step
  - Repeat the process on each newly created conditional FP-tree
  - Until the resulting FP-tree is empty

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



Data set T25I20D10K

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



Data set T25I20D100K

# Why Is FP-Growth the Winner?

- **Divide-and-conquer:**
    - decompose both the mining task and DB according to the frequent patterns obtained so far
    - leads 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

# Visualization of Association Rules: Plane Graph

# Visualization of Association Rules
# (SGI/MineSet 3.0)

# Visualization of Association Rules: Rule Graph

# Mining Frequent Patterns, Association and Correlations

- Basic concepts and a road map

- Efficient and scalable frequent itemset mining methods

- Mining various kinds of association rules

- From association mining to correlation analysis

- Constraint-based association mining

- Summary

# Mining Various Kinds of Association Rules

- Mining multilevel association

- Miming 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**
**min_sup = 5%**

**Milk**
**[support = 10%]**

**Level 1**
**min_sup = 5%**

**Level 2**
**min_sup = 5%**

**2% Milk**
**[support = 6%]**

**Skim Milk**
**[support = 4%]**

**Level 2**
**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

- Quantitative Attributes: numeric, implicit ordering among values—discretization

# 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

2. Dynamic discretization based on data distribution (Agrawal & Srikant@SIGMOD96)

# Quantitative Association Rules

- Proposed by Lent, Swami and Widom ICDE'97
- Numeric attributes are *dynamically* discretized
  - Such that the confidence of the rules mined is maximized
- 2-D quantitative association rules: $A_{quan1} \wedge A_{quan2} \Rightarrow A_{cat}$
- Example

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

# Mining Frequent Patterns, Association and Correlations

- Basic concepts and a road map

- Efficient and scalable frequent itemset mining methods

- Mining various kinds of association rules

- From association mining to correlation analysis

- Constraint-based association mining

- 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

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

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

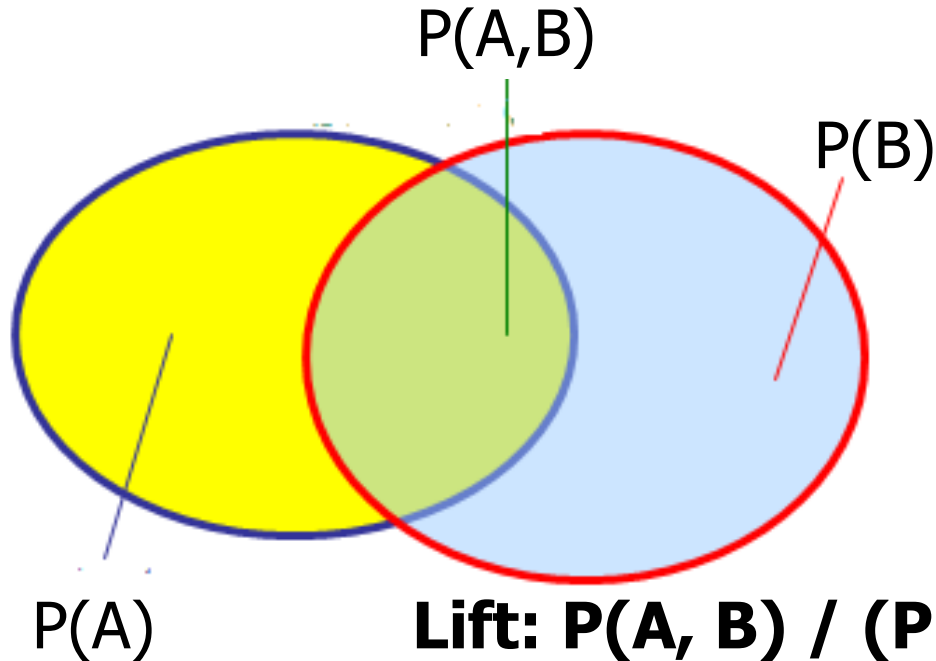$$lift(B,C) = \frac{2000/5000}{3000/5000*3750/5000} = 0.89$$

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

# Which Measures Should Be Used?

- **lift** *and* χ² are not good measures for correlations in large transactional DBs

- **all-conf** or **coherence** could be good measures (Omiecinski@TKDE'03)

- Both **all-conf** and **coherence** have the downward closure property

- Efficient algorithms can be derived for mining (Lee et al. @ICDM'03sub)

| 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…∞ | $\Sigma_i\frac{(P(A_i)-E_i)^2}{E_i}$ |

# Difference Between Confidence, Lift, All-Confidence and Coherence

P(A,B)

P(B)

P(A)

**Lift: P(A, B) / (P(A) * P(B))**

**Confidence:  P(A,B) / P(A)**

**All-Conf: P(A, B) / max(P(A), P(B))**

**Coherence: P(A,B) / (P(A)+P(B)-P(A,B))**

# 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 represent correlations
- So many interestingness measures?  (Tan, Kumar, Sritastava @KDD'02)

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

$$all\_conf = \frac{\sup(X)}{\max\_item\_\sup(X)}$$

$$coh = \frac{\sup(X)}{|universe(X)|}$$

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

| DB | m, c | ~m, c | m~c   | ~m~c     | lift | all-conf | coh  | χ2   |
|----|------|-------|-------|----------|------|----------|------|------|
| A1 | 1000 | 100   | 100   | 10,000   | 9.26 | 0.91     | 0.83 | 9055 |
| A2 | 100  | 1000  | 1000  | 100,000  | 8.44 | 0.09     | 0.05 | 670  |
| A3 | 1000 | 100   | 10000 | 100,000  | 9.18 | 0.09     | 0.09 | 8172 |
| A4 | 1000 | 1000  | 1000  | 1000     | 1    | 0.5      | 0.33 | 0    |

# Mining Frequent Patterns, Association and Correlations

- Basic concepts and a road map

- Efficient and scalable frequent itemset mining methods

- Mining various kinds of association rules

- From association mining to correlation analysis

- Constraint-based association mining

- 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

# Constraints in Data Mining

- Data constraint
  - 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%

# 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 |
|---------|
| {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 |

# 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} |

$C_2$

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

← Scan D

$L_2$

| 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 an Anti-monotone 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 | 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 →
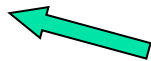
$L_3$

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

**Constraint:**

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

# Mining Frequent Patterns, Association and Correlations

- Basic concepts and a road map

- Efficient and scalable frequent itemset mining methods

- Mining various kinds of association rules

- From association mining to correlation analysis

- Constraint-based association mining

- Summary

# 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)

- Mining a variety of rules and interesting patterns

- Constraint-based mining

- Mining sequential and structured patterns

- Mining truly interesting patterns

  - Surprising, novel, concise, …