

## Chapter 6:

### What Is Frequent Pattern Analysis?

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

April 26, 2011

Data Mining: Concepts and Techniques

1

### 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

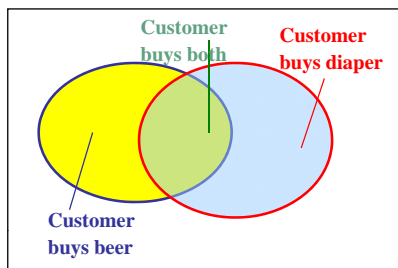
April 26, 2011

Data Mining: Concepts and Techniques

2

### 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



- **itemset**: A set of one or more items
- **k-itemset**  $X = \{x_1, \dots, x_k\}$
- **(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

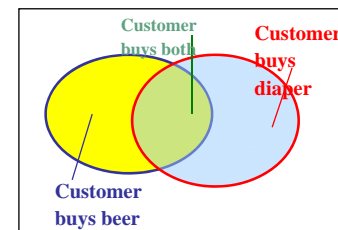
April 26, 2011

Data Mining: Concepts and Techniques

3

### Basic Concepts: Association Rules

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



- Find all the rules  $X \rightarrow Y$  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!)
    - $Beer \rightarrow Diaper$  (60%, 100%)
    - $Diaper \rightarrow Beer$  (60%, 75%)

April 26, 2011

Data Mining: Concepts and Techniques

4

## Closed Patterns and Max-Patterns

- A long pattern contains a combinatorial number of sub-patterns, e.g.,  $\{a_1, \dots, a_{100}\}$  contains  $\binom{100}{1} + \binom{100}{2} + \dots + \binom{100}{100} = 2^{100} - 1 = 1.27 \cdot 10^{30}$  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  $X$
- An itemset  $X$  is a **max-pattern** if  $X$  is frequent and there exists no frequent super-pattern  $Y \supset X$
- Closed pattern is a lossless compression of freq. patterns
  - Reducing the # of patterns and rules

## Closed Patterns and Max-Patterns

- Exercise.  $DB = \{ \langle a_1, \dots, a_{100} \rangle, \langle a_1, \dots, a_{50} \rangle \}$ 
  - $Min\_sup = 1.$
- What is the set of **closed itemset**?
  - $\langle a_1, \dots, a_{100} \rangle: 1$
  - $\langle a_1, \dots, a_{50} \rangle: 2$
- What is the set of **max-pattern**?
  - $\langle a_1, \dots, a_{100} \rangle: 1$
- 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?
  - number of frequent itemsets to be generated is sensitive to minsup threshold
  - If minsup is low, there is an exponential number of frequent itemsets
  - **The worst case:  $M^N$  where  $M$ : # distinct items, and  $N$ : max length of transactions**
- The worst case complexity 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?

## Scalable Frequent Itemset Mining Methods

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

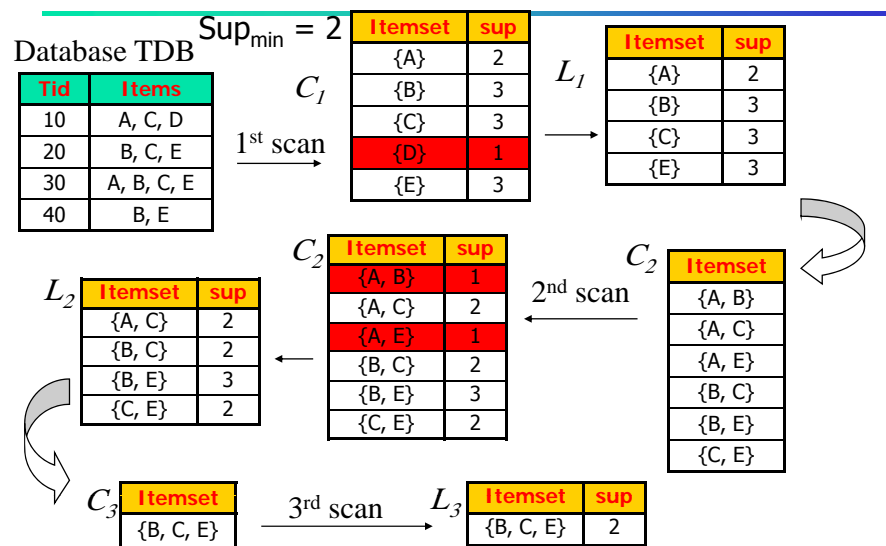
## The Downward Closure Property and Scalable Mining Methods

- 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
  - Freq. pattern growth
  - Vertical data format approach

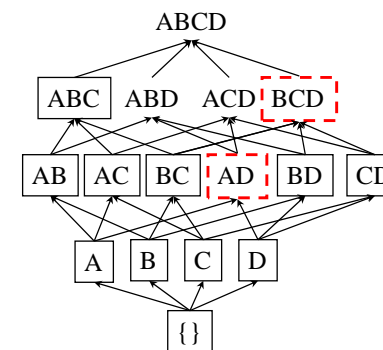
## Apriori: A Candidate Generation & Test Approach

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



## الطريقة الفهلوية



Itemset lattice

## The Apriori Algorithm (Pseudo-Code)

$C_k$ : Candidate itemset of size  $k$

$L_k$ : frequent itemset of size  $k$

$L_1 = \{\text{frequent items}\};$

**for** ( $k = 1; L_k \neq \emptyset; k++$ ) **do begin**

$C_{k+1}$  = candidates generated from  $L_k$ ;

**for each** transaction  $t$  in database **do**

increment the count of all candidates in  $C_{k+1}$  that are contained in  $t$

$L_{k+1}$  = candidates in  $C_{k+1}$  with min\_support

**end**

**return**  $\cup_k L_k$ ;

## 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

## Candidate Generation: An SQL Implementation

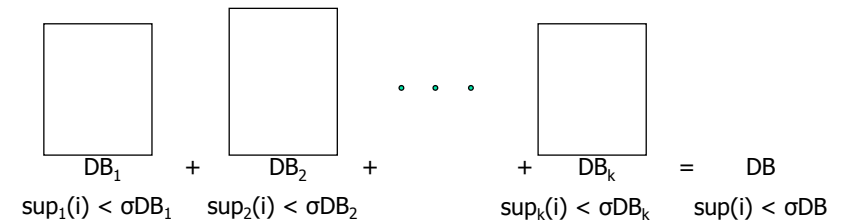
- Suppose the items in  $L_{k-1}$  are listed in an order
- Step 1: self-joining  $L_{k-1}$ 
  - insert into  $C_k$
  - select  $p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$
  - from  $L_{k-1} p, L_{k-1} q$
  - where  $p.item_1 = q.item_1, \dots, p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$
- Step 2: pruning
  - forall *itemsets*  $c$  in  $C_k$  do
  - forall  $(k-1)$ -subsets  $s$  of  $c$  do
  - if ( $s$  is not in  $L_{k-1}$ ) then delete  $c$  from  $C_k$
- Use object-relational extensions like UDFs, BLOBs, and Table functions for efficient implementation

## 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



## 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

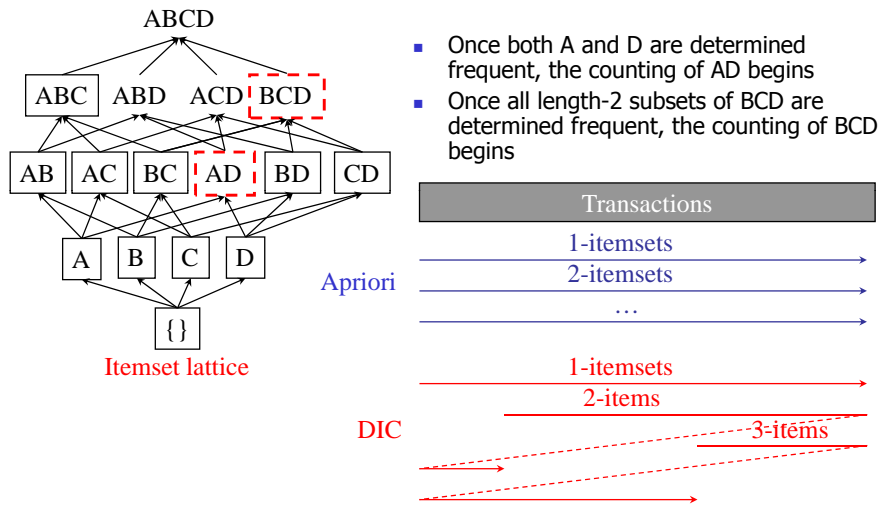
count	itemsets
35	{ab, ad, ae}
88	{bd, be, de}
⋮	⋮
⋮	⋮
102	{yz, qs, wt}

Hash Table

## 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 pattern

## DIC: Reduce Number of Scans



- 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

April 26, 2011

Data Mining: Concepts and Techniques

21

## 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**
  - 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

April 26, 2011

Data Mining: Concepts and Techniques

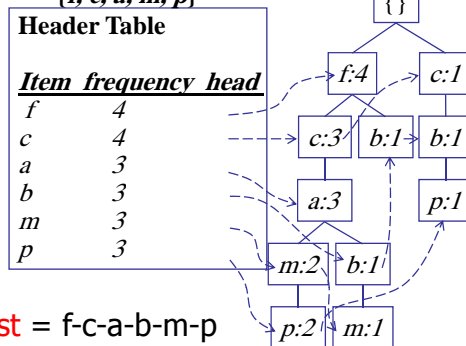
22

## 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

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



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

April 26, 2011

Data Mining: Concepts and Techniques

23

## 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

April 26, 2011

Data Mining: Concepts and Techniques

24