

Data Mining

Mining Frequent Patterns

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

What Is Frequent Pattern Analysis?

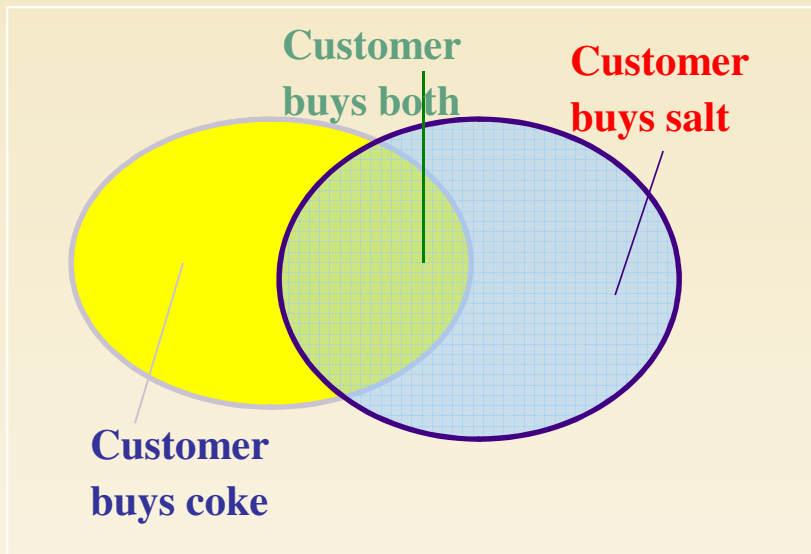
- ❑ **Frequent pattern**: a pattern (a set of items) that occurs frequently in a data set
- ❑ First proposed by Agrawal, Imielinski, and Swami [AIS93] in the context of frequent itemsets and association rule mining
- ❑ Motivation: **Finding inherent regularities in data**
 - What products were often purchased together?— Sabzi and Masala?!
 - What are the subsequent purchases after buying a PC?
 - What kinds of DNA are sensitive to this new drug?
- ❑ 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 and correlation
 - Sequential and structural (e.g., sub-graph) patterns
 - Pattern analysis in multimedia, time-series, and stream data

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



- Itemset $X = \{x_1, \dots, x_k\}$
- 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:

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

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

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

Scalable Methods for Mining Frequent Patterns

- ❑ The **downward closure property** of frequent patterns
 - Any subset of a frequent itemset must be frequent
 - If **{coke, salt, nuts}** is frequent, so is **{coke, salt}**
 - i.e., every transaction having {coke, salt, nuts} also contains {coke, salt}
- ❑ Scalable mining methods: Two 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! (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

- 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
that are contained in

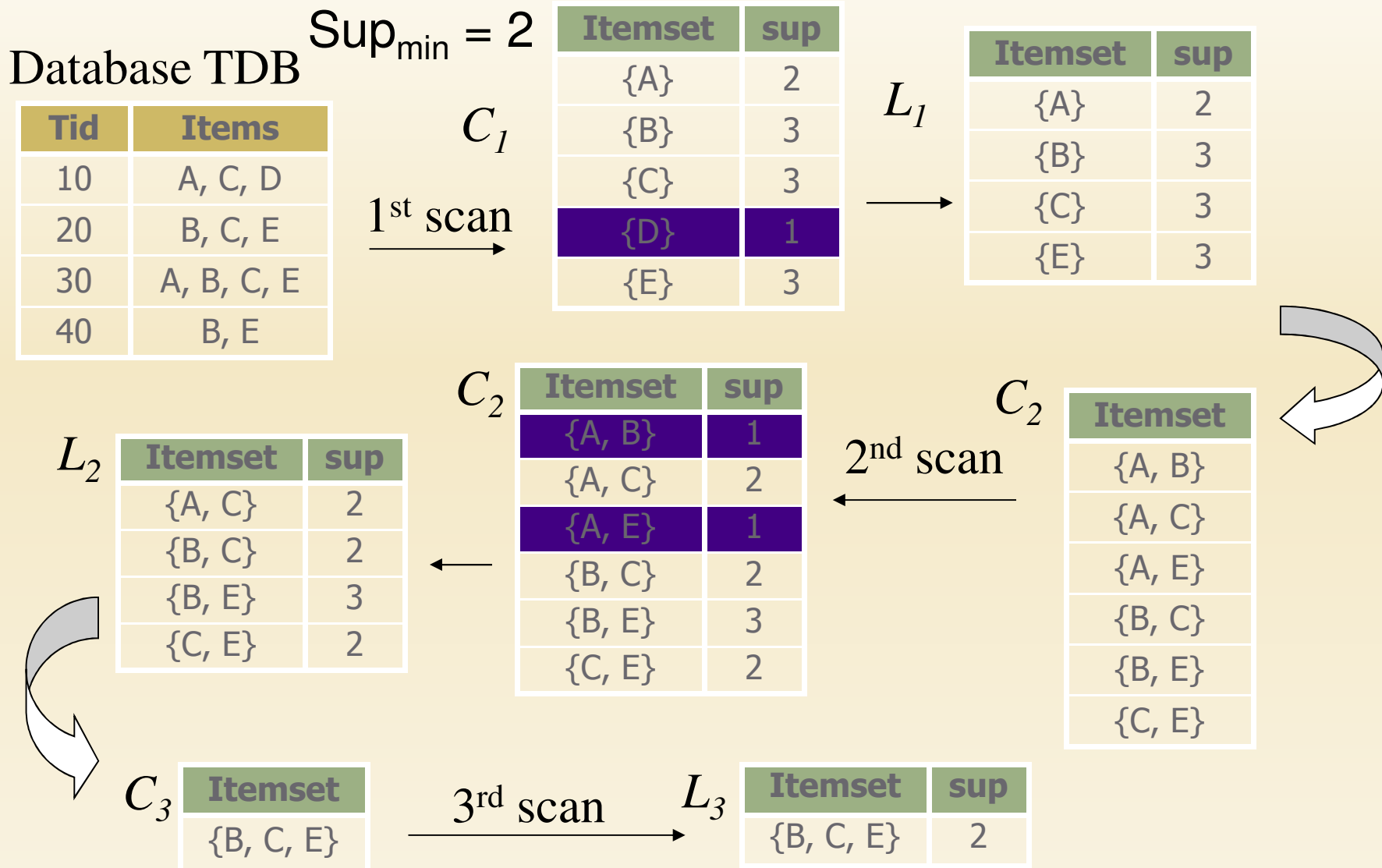
C_{k+1}
 t

L_{k+1} = candidates in C_{k+1} with min_support

end

return $\cup_k L_k$;

The Apriori Algorithm—An Example



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

Bottleneck of Frequent-pattern Mining

- ❑ Multiple database scans are costly
- ❑ Mining long patterns needs many passes of scanning and generates lots of candidates
 - To find frequent itemset $i_1 i_2 \dots i_{100}$
 - # of scans: 100
 - # of Candidates: $\binom{100}{1} + \binom{100}{2} + \dots + \binom{100}{100}$
 $= 2^{100} - 1 = 1.27 * 10^{30} !$
- ❑ Bottleneck: candidate-generation-and-test
- ❑ Can we avoid candidate generation?

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. An efficient algorithm for mining association in large databases. In *VLDB'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*

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 DB|abc → abcd is a frequent pattern

Closed Patterns and Max-Patterns

- A pattern $\{a_1, a_2 \dots, a_n\}$ contains $2^n - 1$ sub-patterns
 - If $n=100$, 1.27×10^{30} sub-patterns!
- A manager may be only interested in patterns involving some items being managed (not all)
- 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 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?
 - !!