

Data Mining

Association Analysis: Basic Concepts and Algorithms

Lecture Notes for Chapter 6

Introduction to Data Mining

by

Tan, Steinbach, Kumar

Why Mining Association Rules?

► Objective:

- Finding interesting co-occurring items (or objects, events) in a given data set.

► Examples:

- Given a database of transactions, each transaction is a list of items (purchased by a customer in a visit), you may find:

computer → financial_management_software
[support=2%, confidence=60%]

- From a student database, you may find

major(x, "CS") ^ gpa(x, "A") → has_taken(x, "DB") [1%, 75%]

What Kind of Databases?

Transaction database TDB

TID	Items
100	f, a, c, d, g, i, m, p
200	a, b, c, f, l,m, o
300	b, f, h, j, o
400	b, c, k, s, p
500	a, f, c, e, l, p, m, n

- Itemset: a set of items
- A *transaction* is a tuple (tid, X)
 - Transaction ID tid
 - Itemset X
- A *transaction database* is a set of transactions
 - In many cases, a transaction database can be treated as a set of itemsets (ignore TIDs)
- Association rule from TDB (relates two itemsets):
 - $\{a, c, m\} \rightarrow \{l\}$ [support=40%, confidence=66.7%]

Association Rule Mining

- Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction

Market-Basket transactions

Example of Association Rules

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

$\{\text{Diaper}\} \rightarrow \{\text{Beer}\}$,
 $\{\text{Milk}, \text{Bread}\} \rightarrow \{\text{Eggs}, \text{Coke}\}$,
 $\{\text{Beer}, \text{Bread}\} \rightarrow \{\text{Milk}\}$,

Implication means co-occurrence,
not causality!

Why Mining Association Rules? (Cont'd)

Popular application: *Basket data analysis*

- place items frequently bought together close to each other to increase sales of these items.
- ? → *a particular product* (What the store should do to boost sales of the particular product)
- *Home Electronics* → ? (What other products should the store stock up?)

Definition: Association Rule

- **Association Rule**

- An implication expression of the form $X \rightarrow Y$, where X and Y are itemsets
- Example:
 $\{Milk, Diaper\} \rightarrow \{Beer\}$

- **Rule Evaluation Metrics**

- Support (s)
 - ◆ Fraction of transactions that contain both X and Y
- Confidence (c)
 - ◆ Measures how often items in Y appear in transactions that contain X

Example:

$$\{Milk, Diaper\} \Rightarrow Beer$$

$$s = \frac{\sigma(Milk, Diaper, Beer)}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(Milk, Diaper, Beer)}{\sigma(Milk, Diaper)} = \frac{2}{3} = 0.67$$

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Mining Association Rules

► Problem statement

Given a *minimum support* (min_sup), also called *support threshold*, and a *minimum confidence* (min_conf), also called *confidence threshold*, find all association rules that satisfy both min_sup and min_conf from a data set D .

Definition: Frequent Itemset

- **Itemset**
 - A collection of one or more items
 - ◆ Example: {Milk, Bread, Diaper}
 - k-itemset
 - ◆ An itemset that contains k items
- **Support count (σ)**
 - Frequency of occurrence of an itemset
 - E.g. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$
- **Support**
 - Fraction of transactions that contain an itemset
 - E.g. $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Association Rule Mining Task

- Given a set of transactions T , the goal of association rule mining is to find all rules having
 - support $\geq \text{minsup}$ threshold
 - confidence $\geq \text{minconf}$ threshold
 - Brute-force approach:
 - List all possible association rules
 - Compute the support and confidence for each rule
 - Prune rules that fail the minsup and minconf thresholds
- ⇒ Computationally prohibitive!

Mining Association Rules

- Two-step approach:
 1. Frequent Itemset Generation
 - Generate all itemsets whose support $\geq \text{minsup}$
 2. Rule Generation
 - Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset
- Frequent itemset generation is still computationally expensive

Support and Confidence: Example

$\text{support}(X \rightarrow Y) = P(X \cup Y)$

$\text{confidence}(X \rightarrow Y) = P(Y|X)$

Transaction ID	Items Bought
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F

Relative frequency is used to estimate the Probability.

- $A \rightarrow C$
- $C \rightarrow A$
- $A \& C \rightarrow B$
- $A \& B \rightarrow E$

(10)⁸

Support and Confidence: Example

$$\text{support}(A \rightarrow B) = P(A \cup B)$$

$$\text{confidence}(A \rightarrow B) = P(B|A)$$

Transaction ID	Items Bought
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F

- $A \rightarrow C$ (50%, 66.7%)
- $C \rightarrow A$ (50%, 100%)
- $A \& C \rightarrow B$ (25%, 50%)
- $A \& B \rightarrow E$ (0%, 0%)

(iv)

How to Mine Association Rules

- A two-step process:

- Find all frequent itemsets ---- the key step
- Generate strong association rules from frequent itemsets.
- Example: given min_sup=50% and min_conf=50%

Transaction ID	Items Bought
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F



Frequent Itemset	Support
{A}	75%
{B}	50%
{C}	50%
{A, C}	50%

- Generate strong rules:

- $\{A\} \rightarrow \{C\}$ [support=50%, confidence=66.6%]
- $\{C\} \rightarrow \{A\}$ [support=50%, confidence=100%]

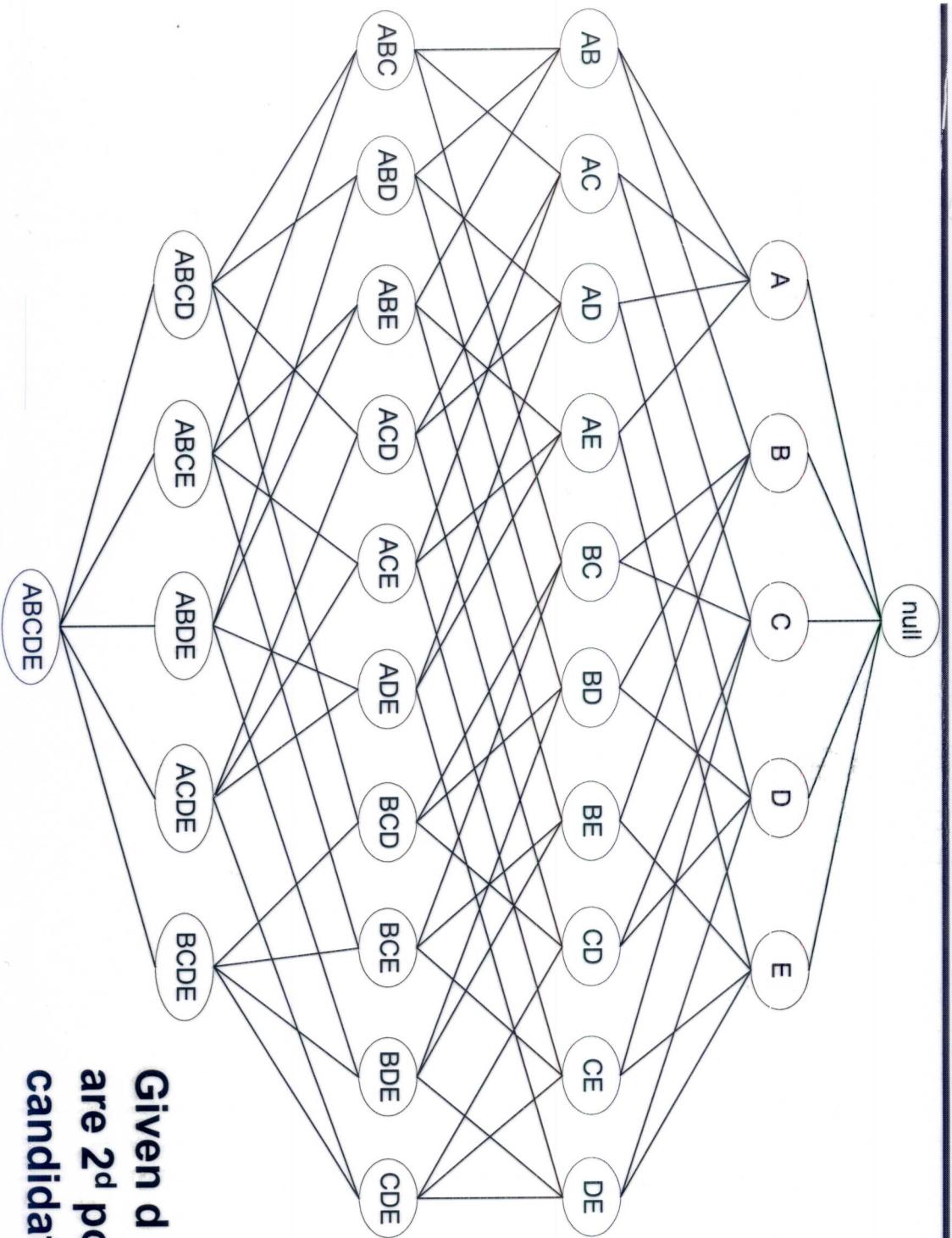
Finding Frequent Itemsets

- Objective: find the itemsets that satisfy the minimum support count.
- Algorithms
 - Apriori
 - FP-Growth
 - H-Mine
 - Partition
 - CLOSET
 - CHARM
 - etc.

Frequent Itemset Generation Strategies

- Reduce the number of candidates (M)
 - Complete search: $M=2^d$
 - Use pruning techniques to reduce M
- Reduce the number of transactions (N)
 - Reduce size of N as the size of itemset increases
 - Used by DHP and vertical-based mining algorithms
- Reduce the number of comparisons (NM)
 - Use efficient data structures to store the candidates or transactions
 - No need to match every candidate against every transaction

Frequent Itemset Generation

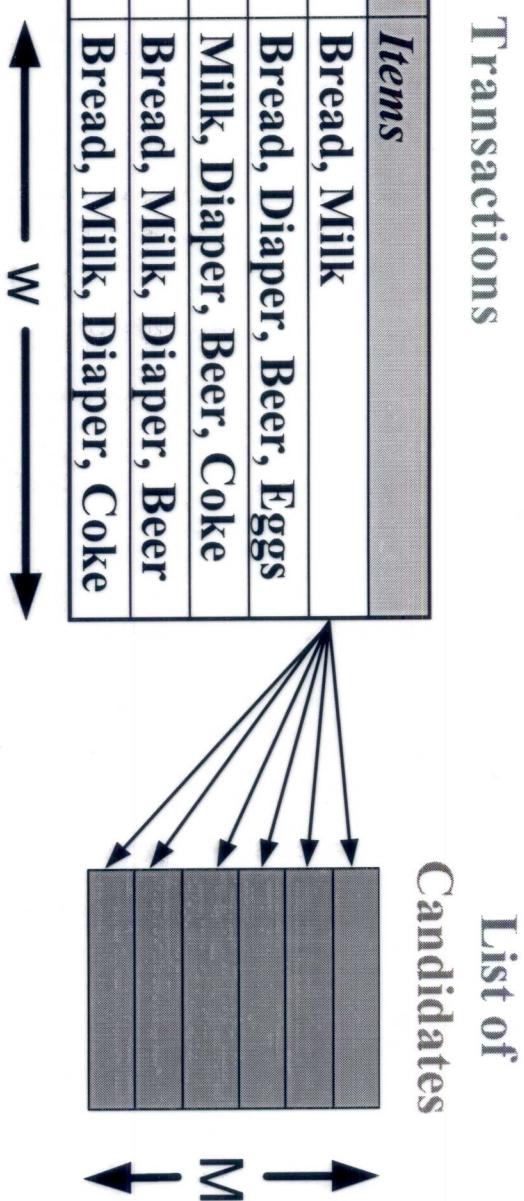


Given d items, there
are 2^d possible
candidate itemsets

Frequent Itemset Generation

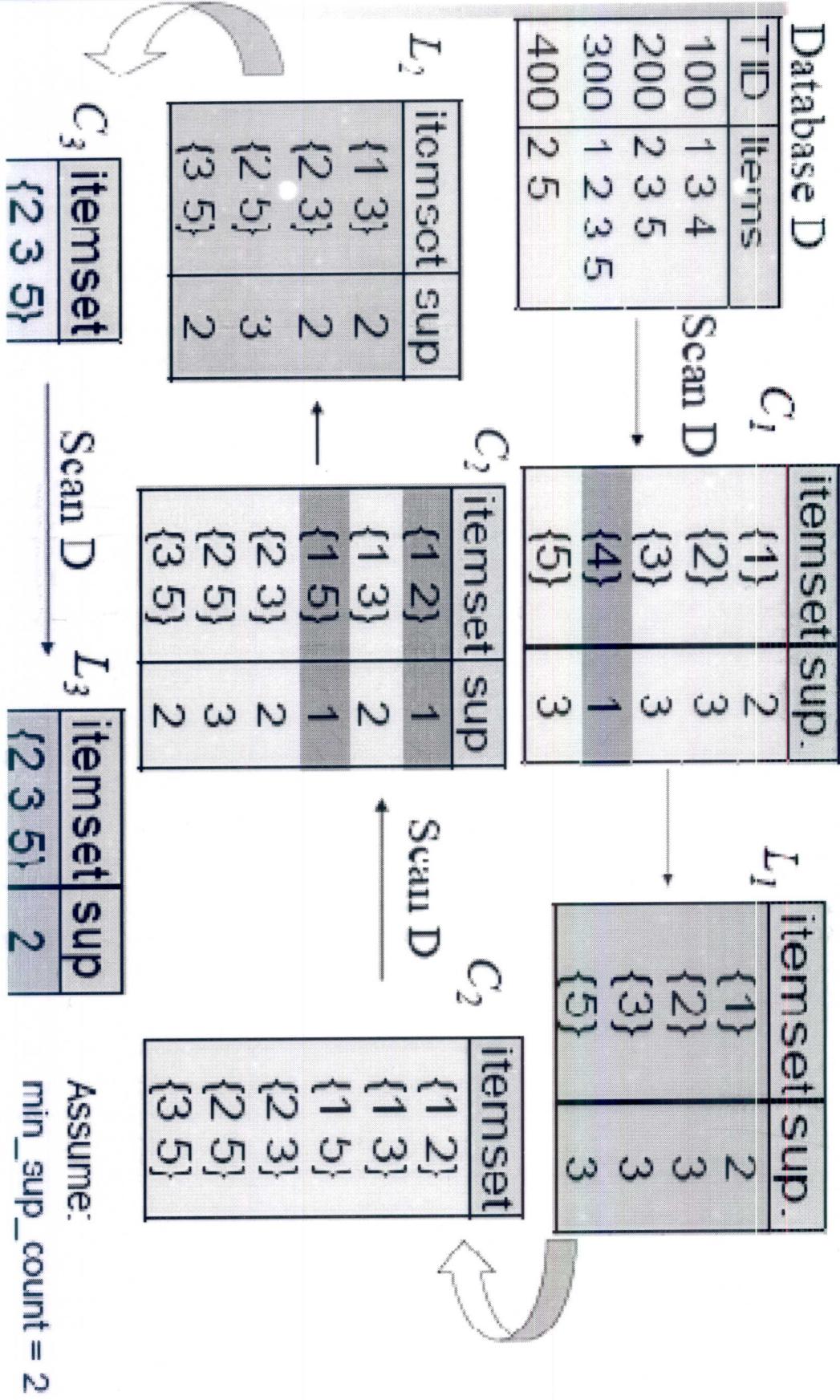
- Brute-force approach:

- Each itemset in the lattice is a candidate frequent itemset
- Count the support of each candidate by scanning the database



- Match each transaction against every candidate
- Complexity $\sim O(NMw) \Rightarrow$ Expensive since $M = 2^d$!!!

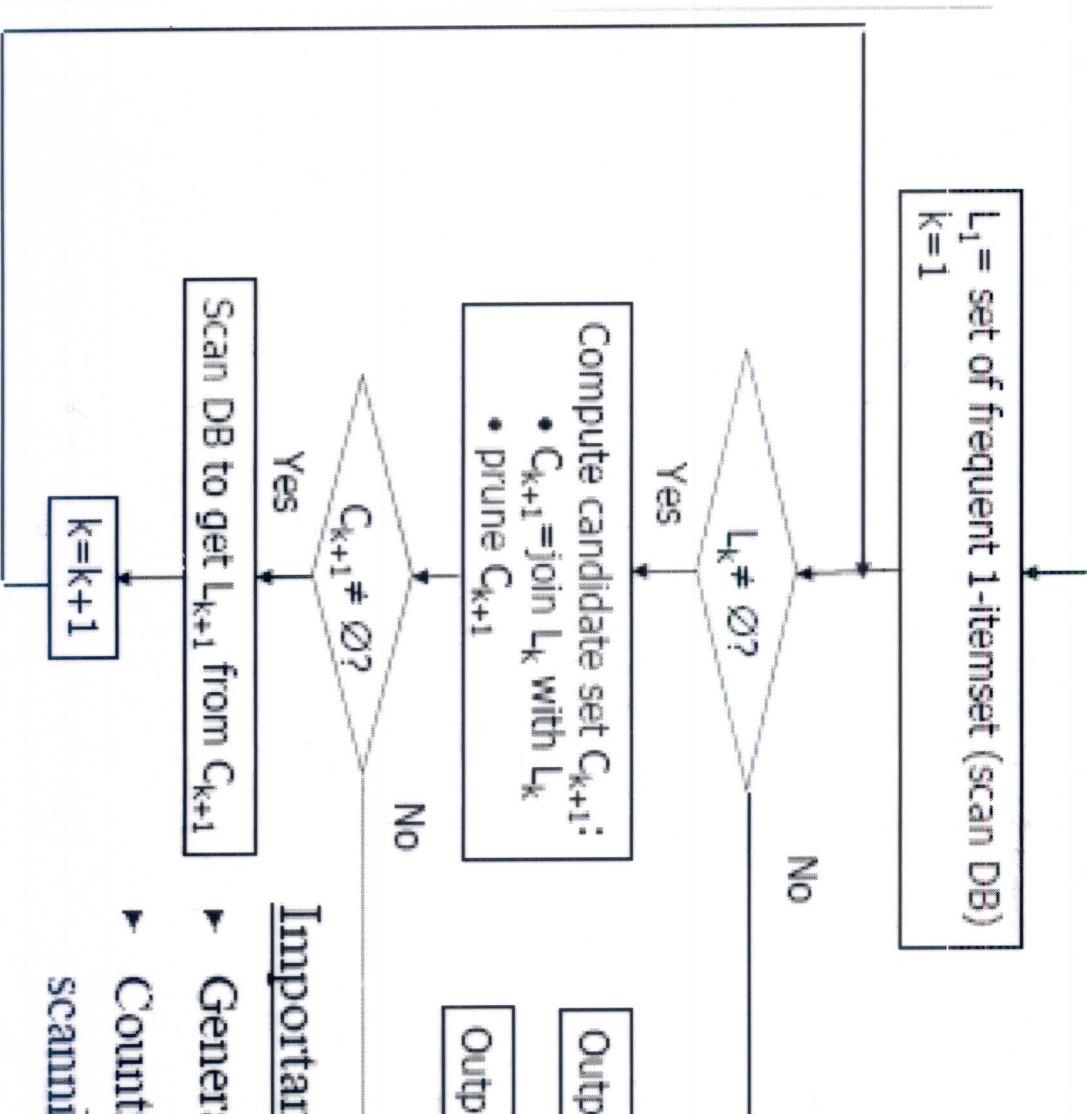
The Apriori Algorithm — Example



Apriori

- The Apriori property (an anti-monotone property), also called *downward closure property*:
 - Any nonempty subset of a frequent itemset must be frequent
- i.e., if $\{A,B,C\}$ is a frequent itemset, $\{A,B\}$, $\{A,C\}$, $\{B,C\}$, $\{A\}$, $\{B\}$ and $\{C\}$ must also be frequent.
- This is because a transaction containing $\{A, B, C\}$ also contains $\{A,B\}$, $\{B,C\}$, ...
- No superset of any infrequent itemset should be checked.
- Many item combinations can be pruned from the search space.
- Apriori-based mining (level-wise iterations)
 - Generate length $(k+1)$ candidate itemsets from frequent k -itemsets
 - Test the candidates against DB

Apriori Algorithm (Flow Chart)



Important steps:

- Generating candidates
- Counting supports of candidates by scanning DB

The Apriori Algorithm

- Based on the Apriori property, use *iterative level-wise approach* and *candidate generation-and-test*

Pseudo-code:

C_k : a set of candidate itemsets of size k
 L_k : the set of frequent itemsets of size k

```
 $L_1 = \{\text{frequent items}\}$ ;           (Scan database to find all frequent 1-itemsets)
for ( $k = 1$ ;  $L_k \neq \emptyset$ ;  $k++$ ) do begin
     $C_{k+1} = \underline{\text{candidates generated from } L_k}$ ;      Scan database to calculate support for each itemset in Ck+1
    for each transaction  $t$  in database do
        increment the count of all candidates in  $C_{k+1}$ 
        that are contained in  $t$ 
     $L_{k+1} = \text{candidates in } C_{k+1} \text{ satisfying min\_support}$ 
end
return  $\bigcup_k L_k$ .
```

Level-wise Generation process: $L_k \rightarrow C_{k+1} \rightarrow L_{k+1}$

How to Generate Candidates?

(i.e., How to Generate C_{k+1} from L_k)

- Given L_k = the set of frequent k -itemsets
- List the items in each itemset of L_k in an order
- Given L_k , generate C_{k+1} in two steps:
 - Join Step: Join L_k with L_k by joining two k -itemsets in L_k . Two k -itemsets are joinable if their first $(k-1)$ items are the same [and the last item in the first itemset is smaller than the last item in the second itemset] (the condition for joining two members of L_k).
Now, $C_4 = \{\{1 2 3 4\}, \{1 3 4 5\}\}$
 - Prune Step: Delete all candidates in C_{k+1} that have a non-frequent subset by checking all length- k subsets of a candidate
 - Now, $C_4 = \{\{1 2 3 4\}\}$

$$L_3 =$$

{1 2 3}
{1 2 4}
{1 3 4}
{1 3 5}
{2 3 4}

Example of Candidate-generation

- ▶ $L_4 = \{abcd, abcg, abdg, abef, abeh, acdg, bcdg\}$
- ▶ Self-joining: $L_4 * L_4$
 - ▶ $abcdg$ from $abcd$ and $abcg$
 - ▶ $abefh$ from $abef$ and $aceh$
- ▶ Pruning:
 - ▶ $abefh$ is removed because $abfh$ or $aefh$ or $befh$ is not in L_4
 - ▶ $C_5 = \{abcdg\}$

Generate Association Rules from Frequent Itemsets

- ▶ Naïve algorithm:

```
for each frequent itemset  $l$  do  
    for each nonempty proper subset  $s$  of  $l$  do  
        if ( $\text{support}(l)/\text{support}(s) \geq \text{min\_conf}$ )  
            output the rule  $s \rightarrow l-s$ , with  
                support =  $\text{support}(l)$  and  
                confidence =  $\text{support}(l)/\text{support}(s)$ 
```

Generate Association Rules from Frequent Itemsets

- ▶ Example:

- ▶ Given a frequent itemset: $I = \{I_1, I_2, I_5\}$
- ▶ nonempty subsets of I are $\{I_1, I_2\}, \{I_1, I_5\}, \{I_2, I_5\}, \{I_1\}, \{I_2\}, \{I_5\}$

- ▶ resulting association rules:

- $I_1 \wedge I_2 \rightarrow I_5, \text{ confidence} = 50\%$
- $I_1 \wedge I_5 \rightarrow I_2, \text{ confidence} = 100\% \checkmark$
- $I_2 \wedge I_5 \rightarrow I_1, \text{ confidence} = 100\% \checkmark$
- $I_1 \rightarrow I_2 \wedge I_5, \text{ confidence} = 33\%$
- $I_2 \rightarrow I_1 \wedge I_5, \text{ confidence} = 29\%$
- $I_5 \rightarrow I_1 \wedge I_2, \text{ confidence} = 100\% \checkmark$

If minimum confidence threshold is 70%, only 3 rules are output.

Is Apriori Fast Enough? — Performance Bottlenecks

- ▶ The core of the Apriori algorithm:
 - ▶ Use frequent k -itemsets to generate candidate $(k+1)$ -itemsets
 - ▶ Use database scan and pattern matching to collect counts for the candidate itemsets to generate frequent $(k+1)$ -itemsets from $k+1$ candidate set
- ▶ The bottleneck of *Apriori*: candidate generation and test
- ▶ Huge candidate sets:
 - ▶ 10^4 frequent 1-itemsets will generate more than 10^7 candidate 2-itemsets
 - ▶ To discover a frequent pattern of size 100, e.g., $\{a_1, a_2, \dots, a_{100}\}$, one needs to generate $2^{100} \approx 10^{30}$ candidates.
- ▶ Multiple scans of database:
 - ▶ Needs n or $n+1$ scans, n is the length of the longest frequent pattern

Q3: For the following TDB, find all association rules with 50% support and 75% confidence.

TID	Items
100	Bread, cheese, Eggs, Juice
200	Bread, cheese, Juice
300	Bread, Milk, Yogurt
400	Bread, Juice, Milk
500	cheese, Juice, Milk

Ans:

Min-support count = 3

C _i	Item	Sup	L _i	Item	Sup
	Bread	4	⇒	Bread	4
	cheese	3		cheese	3
	Juice	4		Juice	4
	Milk	3		Milk	3
*	Eggs	1			
*	Yogurt	1			

(27)

all players.

\therefore All have minimum 75% confidence.

$$x_{SL} = \frac{3}{4}$$

Juice \leftrightarrow Cheese

(Cheese, Juice) \rightarrow Cheese \leftrightarrow Juice

$x_{SL} = \frac{3}{4} = 75\%$

Juice \leftrightarrow Bread

(Bread, Juice) \rightarrow Bread \leftrightarrow Juice

$x_{SL} = \frac{3}{4} = 75\%$

Conf

Generalization :-

Rule

No C3

	x	2	(Juice, Juice)
x	x	1	(Cheese, Juice)
x		3	(Cheese, Juice)
	x	2	(Bread, Juice)
x		3	(Bread, Juice)
	x	2	(Bread, Cheese)

greatest sup.

: z

sup. count

greatest

: z

Class Association Rules

- T --- Transaction data set with n transactions
- I --- set of all items in T
- Y --- set of all class labels
- A CAR is

$X \rightarrow y$, where $X \subseteq I$, and $y \in Y$.