# Data Mining
# Association Analysis: Basic Concepts and Algorithms

Lecture Notes for Chapter 6

Introduction to Data Mining

Fang Zhou

# What is Pattern Discovery?

- **What are patterns?**
  - **Patterns**: A set of items, subsequences, or substructures that occur frequently together (or strongly correlated) in a data set
  - Patterns represent **intrinsic** and **important properties** of datasets
- **Pattern discovery**: Uncovering patterns from massive data sets
- Motivation examples:
  - What products were often purchased together?
  - What are the subsequent purchases after buying an iPad?
  - What word sequences likely form phrases in this corpus?

# Why is it important?

- Finding inherent regularities in a data set

- Foundation for many essential data mining tasks

  – Association, correlation, and causality analysis

  – Mining sequential, structural (e.g., sub-graph) patterns

  – Pattern analysis in spatiotemporal, multimedia, time-series, and stream data

  – Classification: Discriminative pattern-based analysis

  – Cluster analysis: Pattern-based subspace clustering

- Broad applications

  – Market basket analysis, cross-marketing, sale campaign analysis, Web log analysis, biological sequence analysis

# Association Rule Discovery

**Supermarket shelf management – Market-basket model:**

- **Goal:** Identify items that are bought together by sufficiently many customers

- **Approach:** Process the sales data collected with barcode scanners to find dependencies among items

- **A classic rule:**
  - If someone buys diaper and milk, then he/she is likely to buy beer
  - Don't be surprised if you find six-packs next to diapers!

# The Market-Basket Model

- A large set of **items**
  - e.g., things sold in a supermarket

- A **large set** of **baskets**

- Each basket is a **small subset of items**
  - e.g., the things one customer buys on one day

- Want to discover **association rules**
  - People who bought {x,y,z} tend to buy {v,w}
    - ◆ Amazon!

**Input:**

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

**Output:**

**Rules Discovered:**
{Milk} --> {Coke}
{Diaper, Milk} --> {Beer}

# Applications – (1)

- **Items** = products; **Baskets** = sets of products someone bought in one trip to the store

- **Real market baskets:** Chain stores keep TBs of data about what customers buy together
  - Tells how typical customers navigate stores, lets them position tempting items
  - Suggests tie-in "tricks", e.g., run sale on diapers and raise the price of beer

  **Amazon's people who bought *X* also bought *Y***

# Applications – (2)

- **Baskets** = patients; **Items** = drugs & side-effects
  - Has been used to detect combinations
    of drugs that result in particular side-effects
  - **But requires extension:** Absence of an item
    needs to be observed as well as presence

# Association Rule Mining

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

**Market-Basket transactions**

| 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 |

**Example of Association Rules**

{Diaper} → {Beer},
{Milk, Bread} → {Eggs,Coke},
{Beer, Bread} → {Milk},

Implication means co-occurrence, not causality!

# Two key issues

- Discovering patterns from a large transaction data set can be computationally expensive.


- Some of the discovered patterns may be spurious
  - Evaluating the discovered pattern

# Outline

- ## First: Define
  - Frequent itemsets
  - Association rules:
    - Confidence, Support

- ## Then: Algorithms for finding frequent itemsets
  - A-Priori algorithm

# Definition: Frequent Itemset

- ## Itemset
  - A collection of one or more items
    - Example: {Milk, Bread, Diaper}
  - k-itemset
    - An itemset that contains k items

| 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 |

- ## Support count ($\sigma$)
  - Frequency of occurrence of an itemset
  - E.g.   $\sigma(\{Milk, Bread, Diaper\}) = 2$

$$\sigma(X) = \left|\{t_i | X \subseteq t_i, \ t_i \in T\}\right|,$$

# Definition: Frequent Itemset

- **Simplest question:** Find sets of items that appear together "frequently" in baskets

- **Support**
  - Fraction of transactions that contain an itemset
  - E.g. s({Milk, Bread, Diaper}) = 2/5

- **Frequent Itemset**
  - An itemset whose support is greater than or equal to a *minsup* **threshold**

| 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 |

# Example: Frequent Itemsets

- **Items** = {milk, coke, pepsi, beer, juice}
- **Support threshold** = 3 baskets

  $B_1$ = {m, c, b}    $B_2$ = {m, p, j}

  $B_3$ = {m, b}    $B_4$ = {c, j}

  $B_5$ = {m, p, b}    $B_6$ = {m, c, b, j}

  $B_7$ = {c, b, j}    $B_8$ = {b, c}

- **Frequent itemsets:** {m}, {c}, {b}, {j},

  **{m,b}, {b,c}, {c,j}.**

# Definition: Association Rule

- **Association Rule**

  – An implication expression of the form $X \rightarrow Y$, where X and Y are <span style="color:red">disjoint itemsets</span>

  – Example:

  $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

| 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 |

- **Rule Evaluation Metrics**

  – **Support (s)**

    ◆ Fraction of transactions that contain both X and Y

    $$s(X \longrightarrow Y) = \frac{\sigma(X \cup Y)}{N};$$

  – **Confidence (c)**

    ◆ Measures how often items in Y appear in transactions that contain X

    $$c(X \longrightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)};$$

Example:

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

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

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

# Why use support and confidence?

- Support is often used to eliminate uninteresting rules.

- Confidence measures the reliability of the inference made by a rule.
  - The higher the confidence, the more likely it is for Y to be present in transactions that contains X.

- The inference made by an association rule does not necessary imply causality. It suggests a strong co-occurrence relationship between items in X and Y.

# Association Rule Mining Task

- **Association Rules:**
  If-then rules about the contents of baskets

- $\{i_1, i_2, \ldots, i_k\} \to j$ means: "if a basket contains all of $i_1, \ldots, i_k$ then it is *likely* to contain $j$"

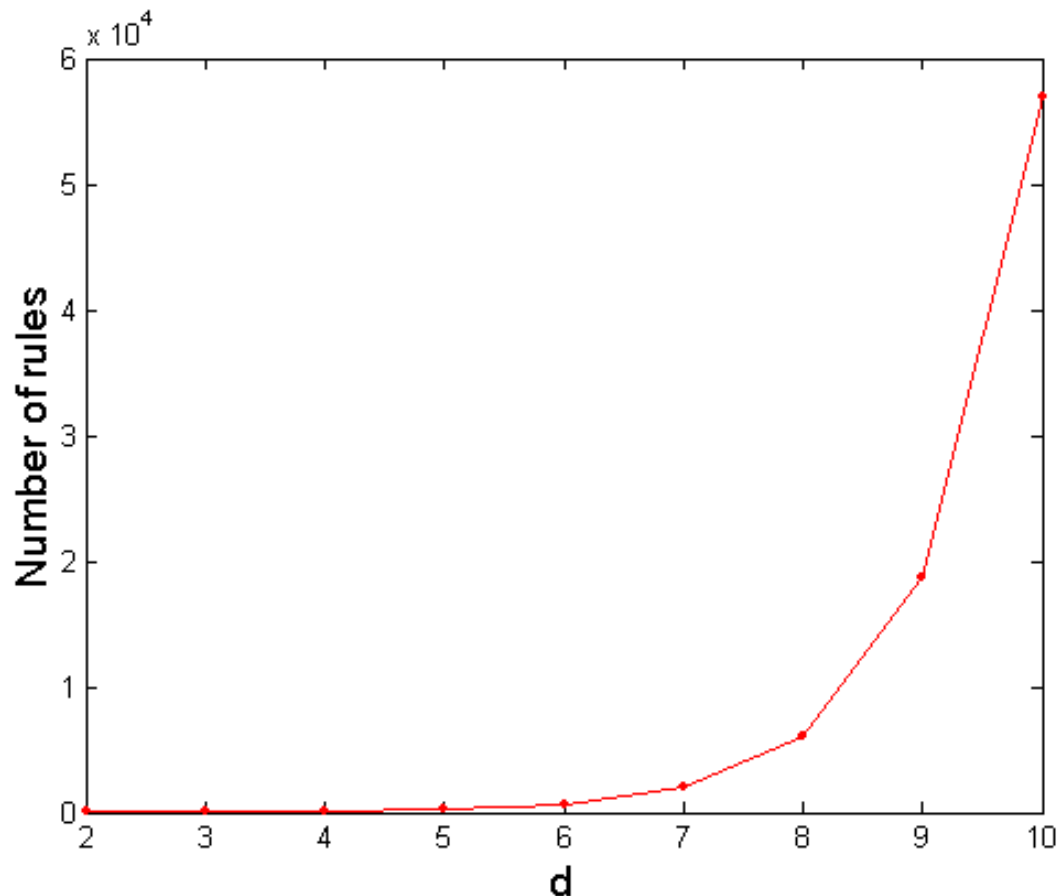- **Problem:** Given a set of transactions T, the goal of association rule mining is to find all rules having

  - support ≥ *minsup* threshold
  - confidence ≥ *minconf* threshold

# Association Rule Mining Task

- 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

  - $\Rightarrow$ Computationally prohibitive!

# Computational Complexity

● Given *d* unique items:

   – Total number of itemsets = $2^d$

   – Total number of possible association rules:



$$R = \sum_{k=1}^{d-1}\left[\binom{d}{k} \times \sum_{j=1}^{d-k}\binom{d-k}{j}\right]$$

$$= 3^d - 2^{d+1} + 1$$

If *d*=6, R = 602 rules

# Mining 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 |

## Example of Rules:

{Milk,Diaper} $\rightarrow$ {Beer} (s=0.4, c=0.67)
{Milk,Beer} $\rightarrow$ {Diaper} (s=0.4, c=1.0)
{Diaper,Beer} $\rightarrow$ {Milk} (s=0.4, c=0.67)
{Beer} $\rightarrow$ {Milk,Diaper} (s=0.4, c=0.67)
{Diaper} $\rightarrow$ {Milk,Beer} (s=0.4, c=0.5)
{Milk} $\rightarrow$ {Diaper,Beer} (s=0.4, c=0.5)

## Observations:

• All the above rules are binary partitions of the same itemset:
    {Milk, Diaper, Beer}

• Rules originating from the same itemset have identical support but can have different confidence

• Thus, we may decouple the support and confidence requirements

# Mining Association Rules

- **Problem: Find all association rules with support $\geq s$ and confidence $\geq c$**

- **Hard part: Finding the frequent itemsets!**
  - If $\{i_1, i_2, \ldots, i_k\} \rightarrow j$ has high support and confidence, then both $\{i_1, i_2, \ldots, i_k\}$ and $\{i_1, i_2, \ldots, i_k, j\}$ will be "frequent"

# Mining Association Rules

- Two-step approach:
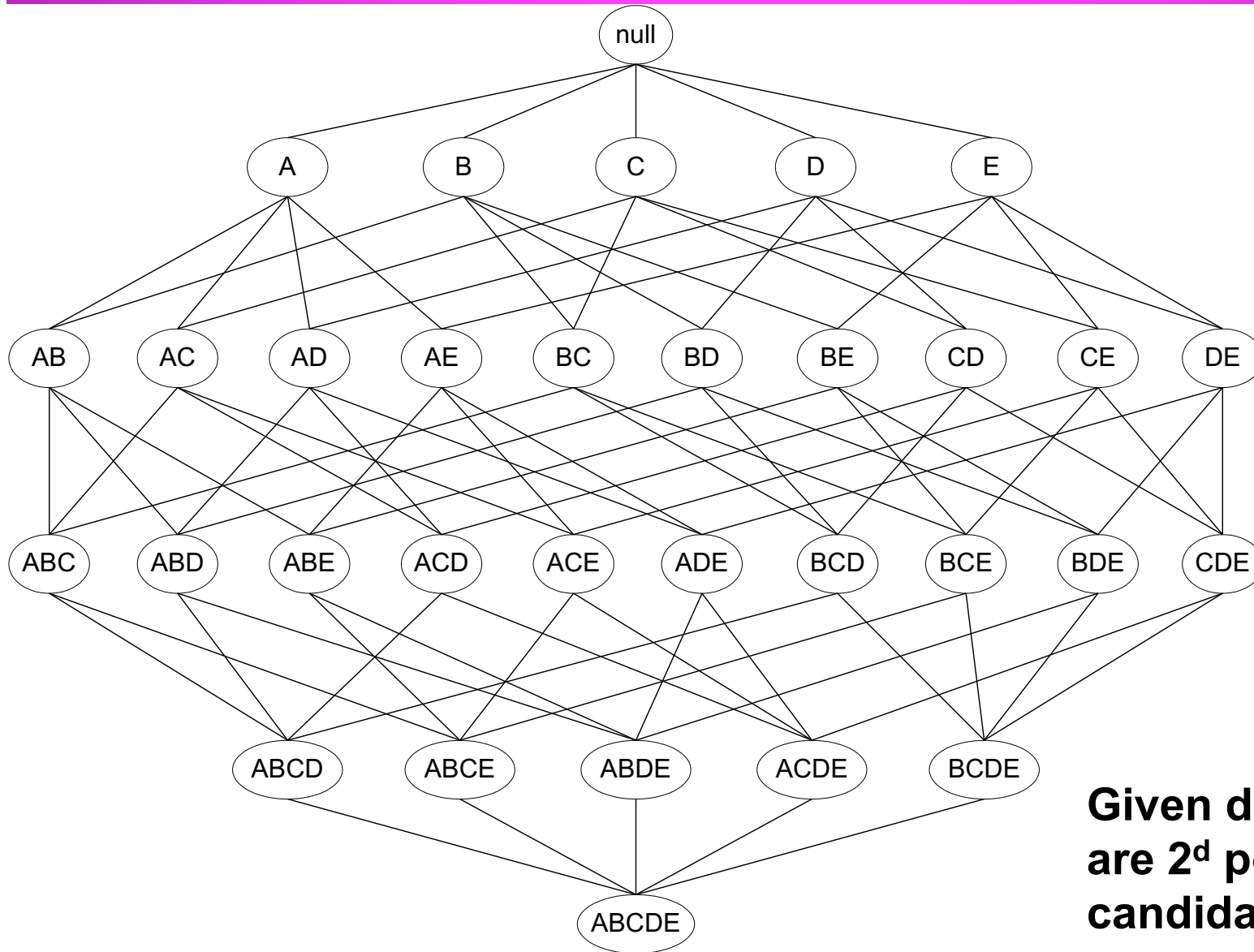  1. Frequent Itemset Generation
     - Generate all itemsets $I$ whose support $\geq$ minsup

  2. Rule Generation
     - Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

     - For every subset $A$ of $I$, generate a rule $A \rightarrow I \setminus A$
       - Since $I$ is frequent, $A$ is also frequent

- Frequent itemset generation is still computationally expensive

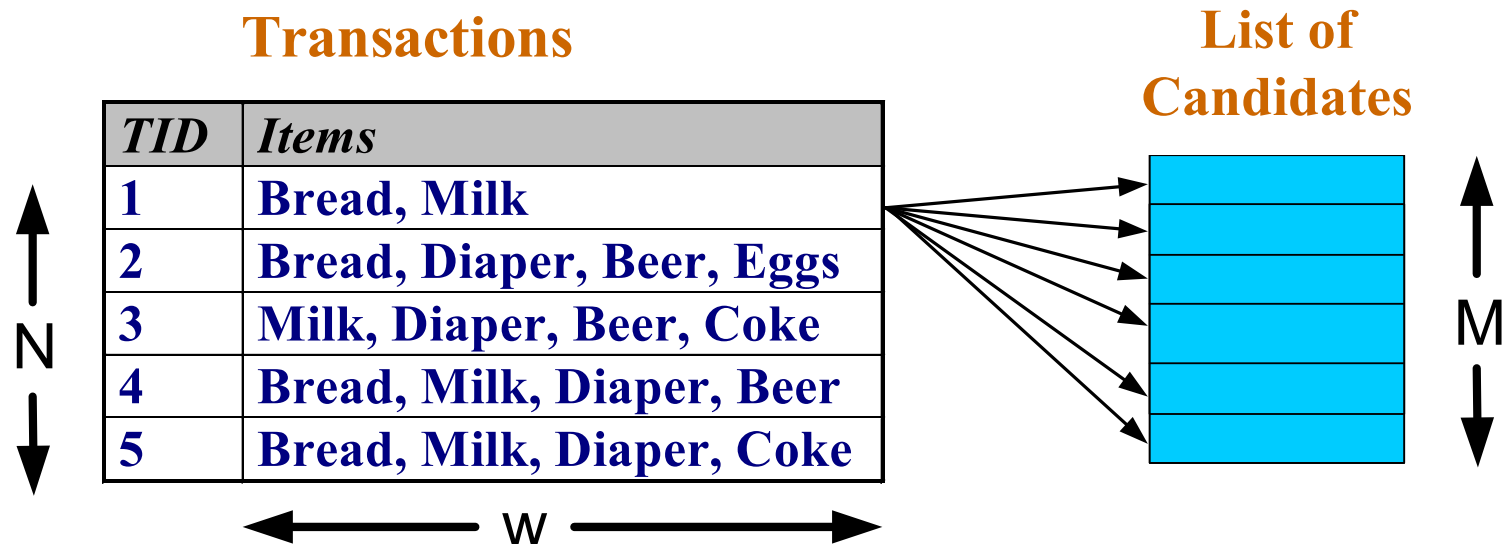# Frequent Itemset Generation



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

**Transactions**

| 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 |

N

w

**List of Candidates**

M

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

# 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 comparisons (NM)
  - Use efficient data structures to store the candidates or transactions
  - No need to match every candidate against every transaction

# Quiz

- Let $c_1$, $c_2$, and $c_3$ be the confidence values of the rules $\{p\} \rightarrow \{q\}$, $\{p\} \rightarrow \{q,r\}$, and $\{p,r\} \rightarrow \{q\}$, respectively.

  If we assume that $c_1$, $c_2$, and $c_3$ have different values, what are the possible relationships that may exist among $c_1$, $c_2$, and $c_3$? Which rule has the lowest confidence?

# Reducing Number of Candidates

● Apriori principle:

- If an itemset is frequent, then all of its subsets must also be frequent
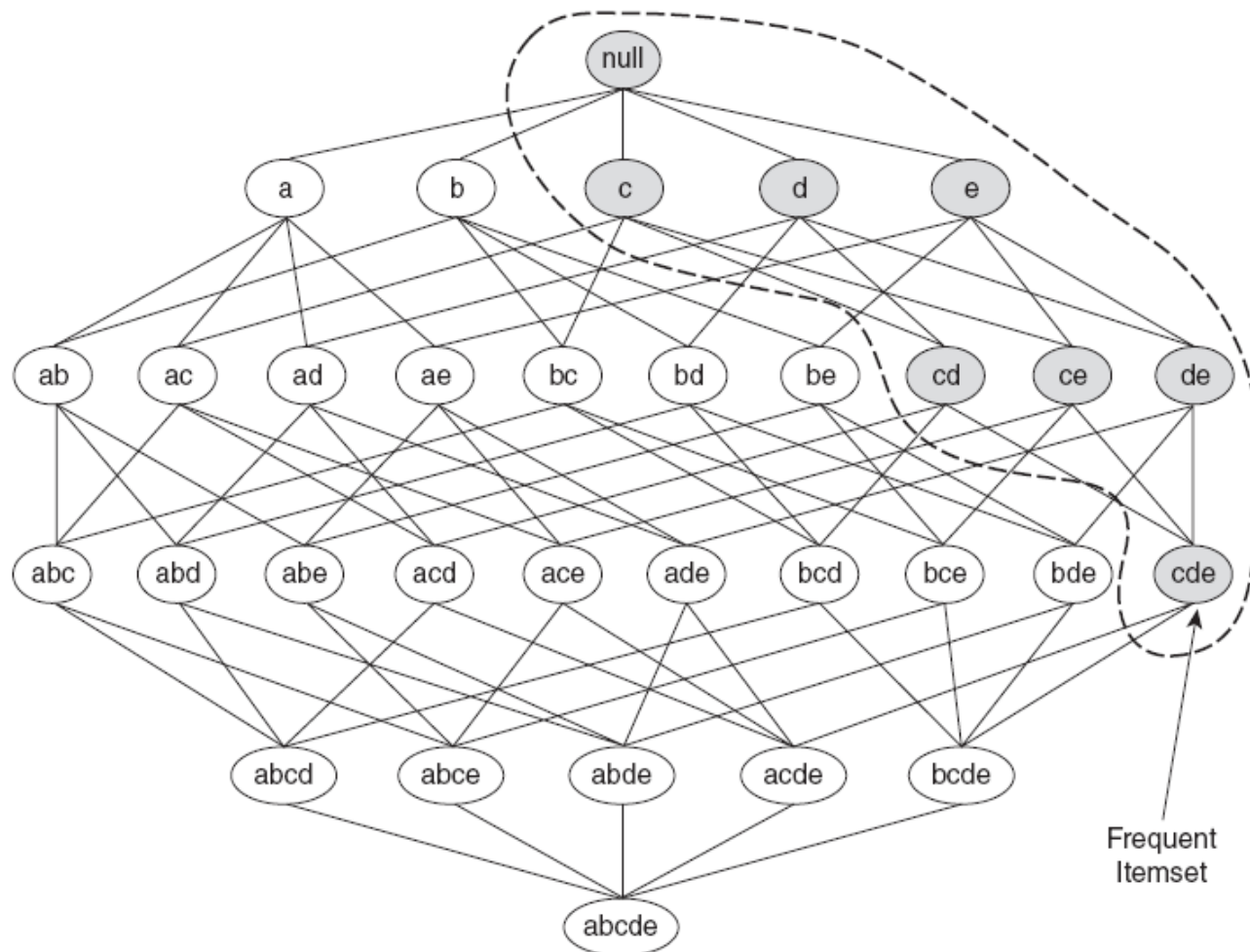
# Illustrating Apriori Principle



**Figure 6.3.** An illustration of the *Apriori* principle. If $\{c, d, e\}$ is frequent, then all subsets of this itemset are frequent.