# Data Mining

Association Analysis

Introduction to Data Mining, 2$^{nd}$ Edition
by
Tan, Steinbach, Karpatne, Kumar

# Introduction

- Business enterprises accumulate large quantities of data from their day-to-day operations

    - *huge amounts of customer purchase data* collected daily at the checkout counters of grocery stores.

**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 |

# Introduction

- Business enterprises accumulate large quantities of data from their day-to-day operations

  - *huge amounts of customer purchase data* collected daily at the checkout counters of grocery stores.

### 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 |

- Analyze the data to learn about the **purchasing behaviour** of the customers.

- Valuable information can be used to support variety of **business-related applications**:

  → *Marketing promotions*

  → *Inventory managment*

  → *CRM*

  → *Pricing*

# 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**

| 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!

# 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 |

# 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. $\sigma$ ({Milk, Bread,Diaper}) = 2

| 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 |

# 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. $\sigma$ ({Milk, Bread,Diaper}) = 2
- **Support**
  - Fraction of transactions that contain an itemset
  - E.g.   s({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 |

# 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. $\sigma$ ({Milk, Bread,Diaper}) = 2
- **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 |

# Definition: Association Rule

- **Association Rule**
  - An implication expression of the form X → Y, where X and Y are itemsets
  - Example:
    {Milk, Diaper} → {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 |

# 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

| 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 |

# Definition: Association Rule

- **Association Rule**
  - An implication expression of the form X→ Y, where X and Y are itemsets
  - Example:
    {Milk, Diaper}→{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

| 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:

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

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

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

# Association Rule Mining Task

- Given a set of transactions T, the goal of association rule mining is to find all rules having
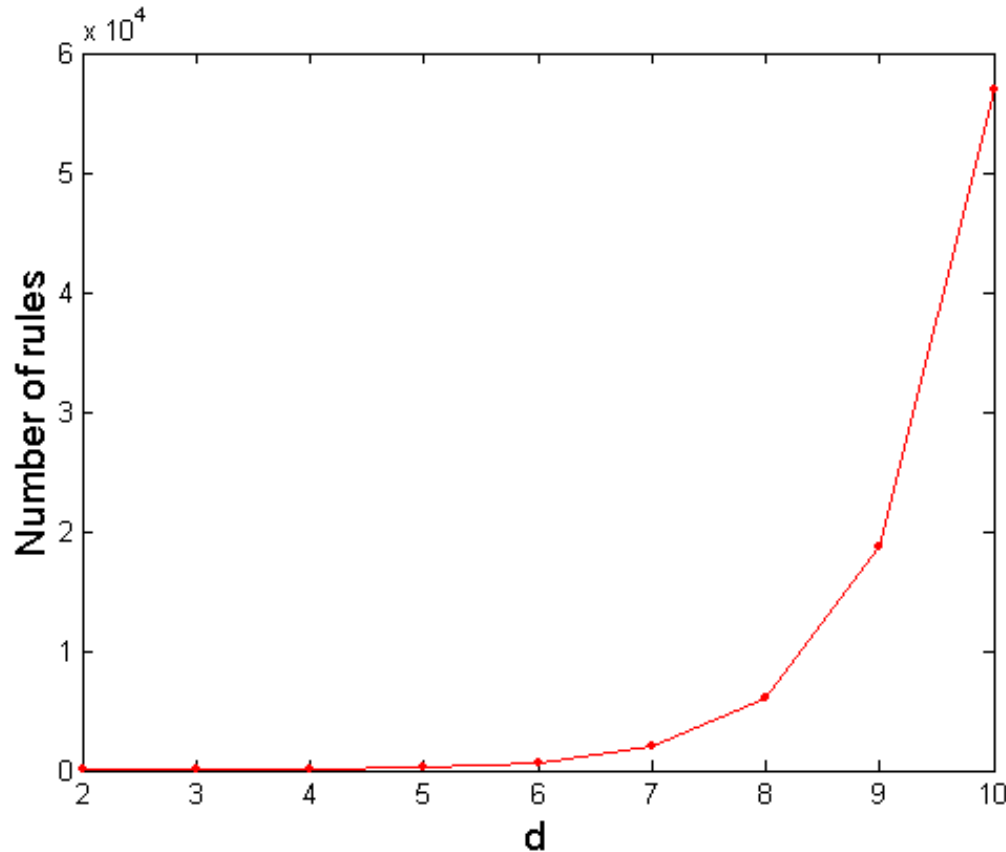  - support ≥ *minsup* threshold
  - confidence ≥ *minconf* threshold

# Association Rule Mining Task

- Given a set of transactions T, the goal of association rule mining is to find all rules having
  - support ≥ *minsup* threshold
  - confidence ≥ *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!

# 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} → {Beer} (s=0.4, c=0.67)
{Milk,Beer} → {Diaper} (s=0.4, c=1.0)
{Diaper,Beer} → {Milk} (s=0.4, c=0.67)
{Beer} → {Milk,Diaper} (s=0.4, c=0.67)
{Diaper} → {Milk,Beer} (s=0.4, c=0.5)
{Milk} → {Diaper,Beer} (s=0.4, c=0.5)

# Mining Association Rules

Observations:

✂ All the 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

Divide the problem in 2 sub tasks:

1. Frequent Itemset Generation
   - Find all the itemsets where support >= *minsup* threshold

2. Rule Generation
   - From above set find all itemsets with high confidence. These are strong rules

- Frequent Itemset Generation is still computationally expensive than Rule Generation

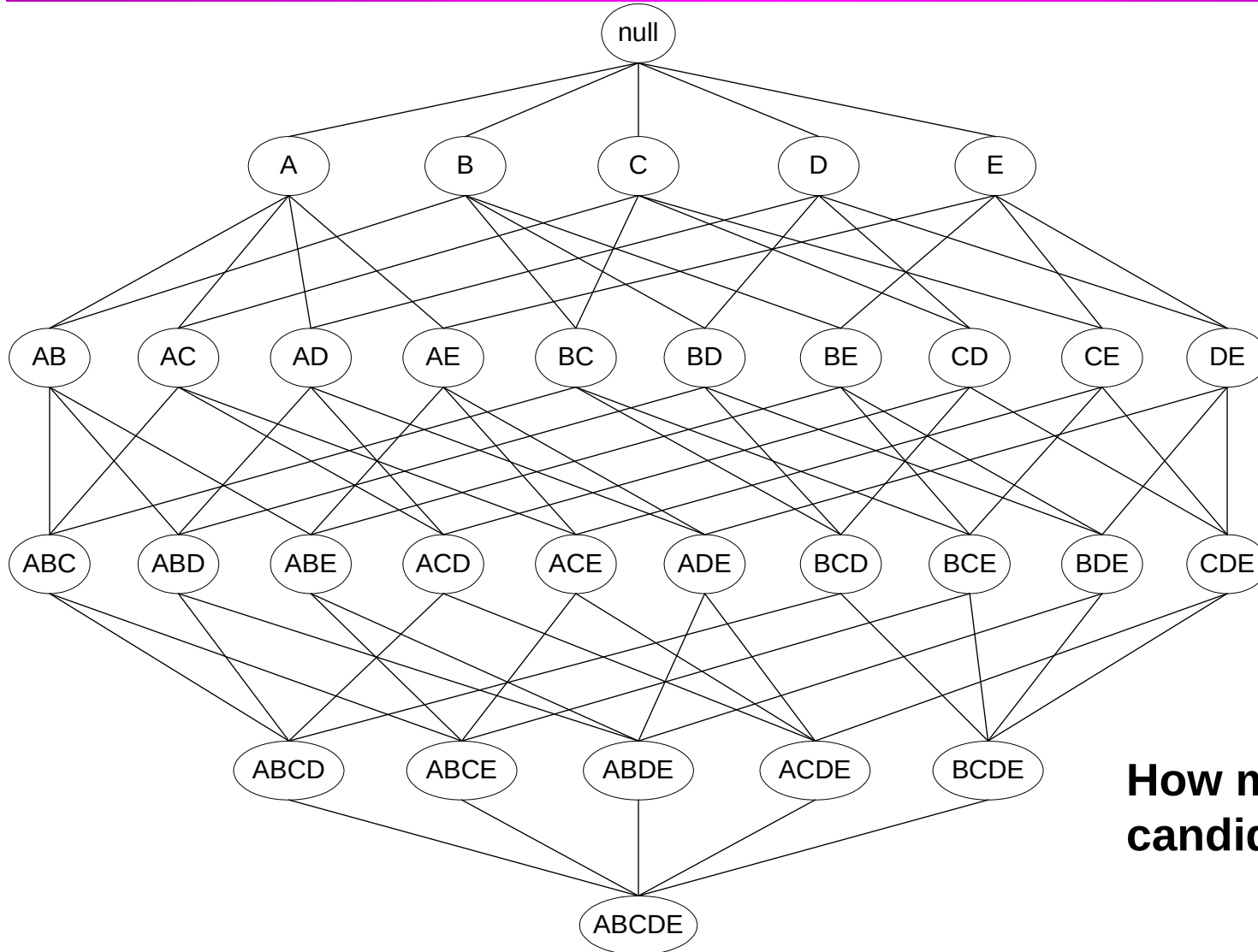# Frequent Itemset Generation

What are the list of possible itemsets?

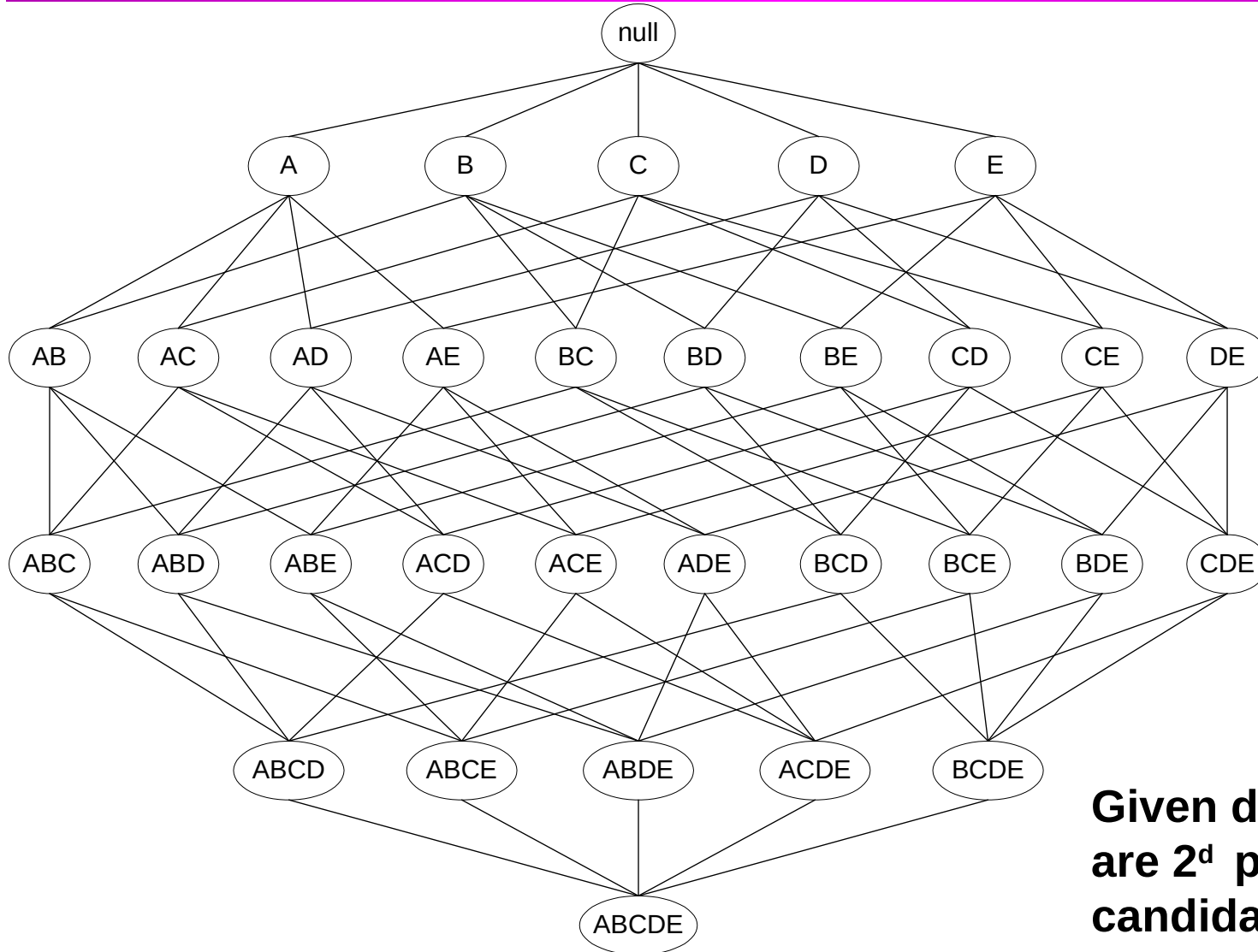For items, I = { a, b, c, d, e}

- Itemsets could be **a, b, c, d, e**

- **ab, ac, ad, ae**

- **bc, bd, be**

- **cd, ce**

- **de**

- **abc, abd, abe, acd, ace, and all such combinations**

# Frequent Itemset Generation



**How many possible candidate itemsets?**

**Introduction to Data Mining, 2nd Edition**

# 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**

**List of Candidates**

| 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

M

  - Match (compare) each transaction against every candidate
  - Complexity ~ O(NMw) => Expensive since M = $2^d$-1 !!!

# 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

# The Apriori Principle

- Apriori principle:
  - If an itemset is frequent, then all of its subsets must also be frequent

- For example:

  - {Milk, Bread, Diapers} is a frequent itemset
    - means Milk, Bread, Diapers <u>support is high</u>
    - thus {Milk, Bread, Diapers} occur in many transaction

# The Apriori Principle

- **Apriori principle**:

  If {Milk, Bread, Diapers} occurs in many transactions

  So does {Milk} , {Bread}, {Diapers}

  and

  {Milk, Bread}, {Milk, Diapers}, {Bread, Diapers}

# The 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.

# The Apriori Principle

- Apriori principle:
  - If an itemset is infrequent, then all of its supersets must also be infrequent

- For example:

  - if {Coke} , {Jam} is infrequent itemset
    - means {Coke} , {Jam} <u>support is low</u>
    - thus {Coke, Jam} also don't occur in many transaction

# The Apriori Principle



**Figure 6.4.** An illustration of support-based pruning. If $\{a, b\}$ is infrequent, then all supersets of $\{a, b\}$ are infrequent.

# The Apriori Principle

- **Support-Based Pruning**.

  - Trimming the <u>exponential search space</u> based on the <u>support measure</u>.

  - Based on **Anti-Monotone Property**

    - Support of an itemset never exceeds the support of it's subsets:

    - **I : set of items,**

    - **J = $2^I$ (power set) // all possible combinations**

    - **S (support) is anti-monotone (or downward closed)**

    $$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

    - Any anti-monotone measured can be used in mining

# The Apriori Principle

- **Monotonicity Property**

  - **I : set of items,**

  - **J = $2^I$ (power set) // all possible combinations**

  - **f is monotone (or upward closed)**

    - **$\forall X, Y \in J : (X \subseteq Y) \rightarrow f(X) \leq f(Y)$**

# Apriori Algorithm

- Uses Apriori Principle

- Support Based Pruning

- Control exponential growth of candidate itemset

# Apriori Algorithm

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

Candidate
1-Itemsets

| Item | Count |
|------|-------|
| Beer | 3 |
| Bread | 4 |
| Cola | 2 |
| Diapers | 4 |
| Milk | 4 |
| Eggs | 1 |

Minimum support count = 3

Candidate
2-Itemsets

| Itemset | Count |
|---------|-------|
| {Beer, Bread} | 2 |
| {Beer, Diapers} | 3 |
| {Beer, Milk} | 2 |
| {Bread, Diapers} | 3 |
| {Bread, Milk} | 3 |
| {Diapers, Milk} | 3 |

Itemsets removed
because of low
support

Candidate
3-Itemsets

| Itemset | Count |
|---------|-------|
| {Bread, Diapers, Milk} | 3 |

**Figure 6.5.** Illustration of frequent itemset generation using the *Apriori* algorithm.

# Apriori Algorithm

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

Candidate
1-Itemsets

| Item | Count |
|------|-------|
| Beer | 3 |
| Bread | 4 |
| Cola | 2 |
| Diapers | 4 |
| Milk | 4 |
| Eggs | 1 |

Minimum support count = 3

Candidate
2-Itemsets

| Itemset | Count |
|---------|-------|
| {Beer, Bread} | 2 |
| {Beer, Diapers} | 3 |
| {Beer, Milk} | 2 |
| {Bread, Diapers} | 3 |
| {Bread, Milk} | 3 |
| {Diapers, Milk} | 3 |

Itemsets removed
because of low
support

Candidate
3-Itemsets

| Itemset | Count |
|---------|-------|
| {Bread, Diapers, Milk} | 3 |

If every subset is considered,
$$^6C_1 + {}^6C_2 + {}^6C_3$$
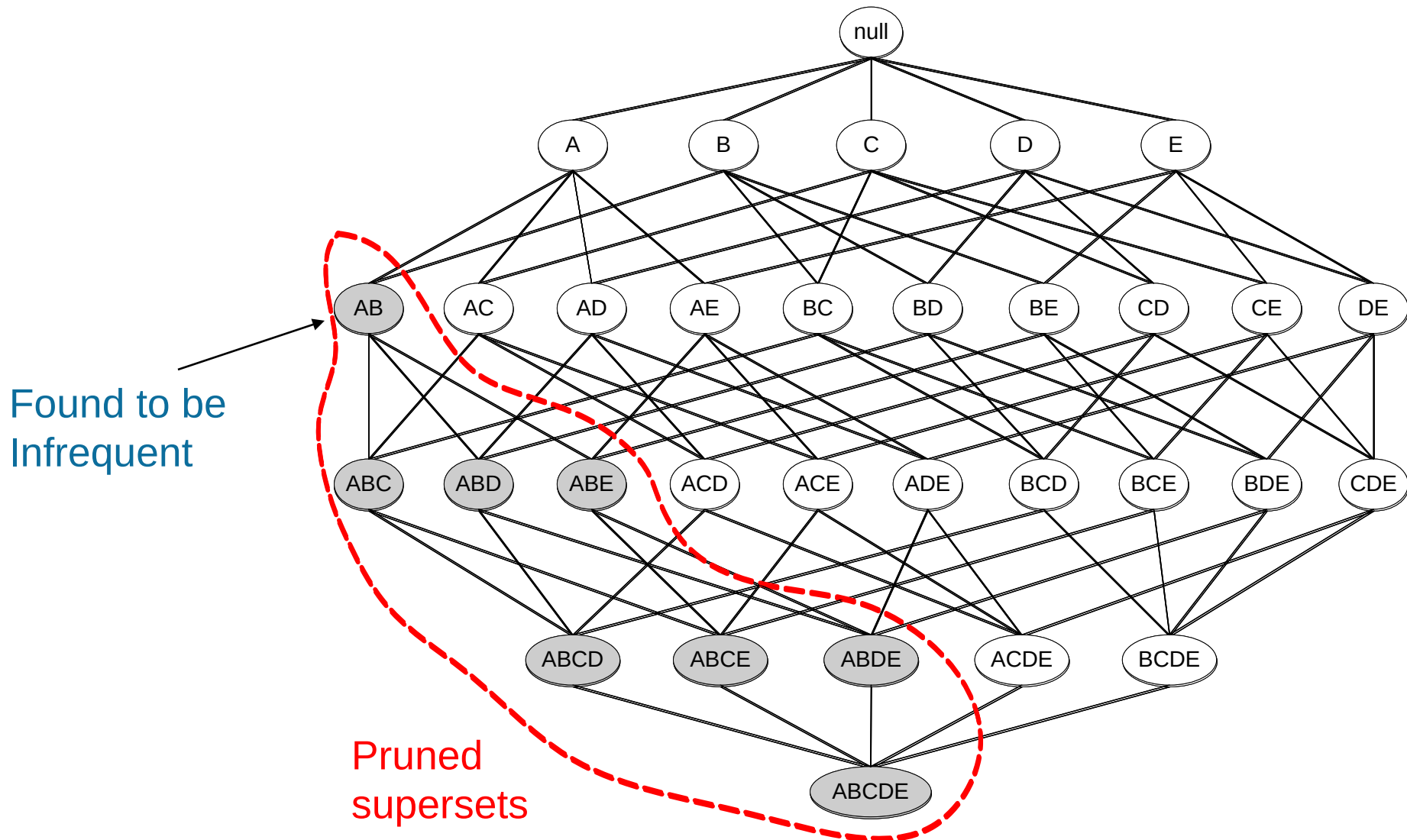$$6 + 15 + 20 = 41$$
With support-based pruning,
$$6 + 6 + 1 = 13$$

...mset generation using the *Apriori* algorithm.

...ining, 2nd Edition                                    32

# Illustrating Apriori Principle



Found to be Infrequent

Pruned supersets

# Illustrating Apriori Principle

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

Items (1-itemsets)

| Item | Count |
|------|-------|
| Bread | 4 |
| Coke | 2 |
| Milk | 4 |
| Beer | 3 |
| Diaper | 4 |
| Eggs | 1 |

Minimum Support = 3

If every subset is considered,
$$^6C_1 + {}^6C_2 + {}^6C_3$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
$$6 + 6 + 4 = 16$$

# Illustrating Apriori Principle

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

Items (1-itemsets)

| Item | Count |
|------|-------|
| Bread | 4 |
| Coke | 2 |
| Milk | 4 |
| Beer | 3 |
| Diaper | 4 |
| Eggs | 1 |

Minimum Support = 3

If every subset is considered,
$$^6C_1 + {}^6C_2 + {}^6C_3$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
$$6 + 6 + 4 = 16$$

# Illustrating Apriori Principle

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

Items (1-itemsets)

| Item | Count |
|------|-------|
| Bread | 4 |
| Coke | 2 |
| Milk | 4 |
| Beer | 3 |
| Diaper | 4 |
| Eggs | 1 |

| Itemset |
|---------|
| {Bread,Milk} |
| {Bread, Beer } |
| {Bread,Diaper} |
| {Beer, Milk} |
| {Diaper, Milk} |
| {Beer,Diaper} |

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
$$^6C_1 + {}^6C_2 + {}^6C_3$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
$$6 + 6 + 4 = 16$$

# Illustrating Apriori Principle

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

Items (1-itemsets)

| Item | Count |
|------|-------|
| Bread | 4 |
| Coke | 2 |
| Milk | 4 |
| Beer | 3 |
| Diaper | 4 |
| Eggs | 1 |

Pairs (2-itemsets)

| Itemset | Count |
|---------|-------|
| {Bread,Milk} | 3 |
| {Beer, Bread} | 2 |
| {Bread,Diaper} | 3 |
| {Beer,Milk} | 2 |
| {Diaper,Milk} | 3 |
| {Beer,Diaper} | 3 |

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
$$^6C_1 + {}^6C_2 + {}^6C_3$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
$$6 + 6 + 4 = 16$$

# Illustrating Apriori Principle

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

Items (1-itemsets)

| Item | Count |
|---|---|
| Bread | 4 |
| Coke | 2 |
| Milk | 4 |
| Beer | 3 |
| Diaper | 4 |
| Eggs | 1 |

Pairs (2-itemsets)

| Itemset | Count |
|---|---|
| {Bread,Milk} | 3 |
| {Bread,Beer} | 2 |
| {Bread,Diaper} | 3 |
| {Milk,Beer} | 2 |
| {Milk,Diaper} | 3 |
| {Beer,Diaper} | 3 |

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
$^6C_1 + {}^6C_2 + {}^6C_3$
$6 + 15 + 20 = 41$
With support-based pruning,
$6 + 6 + 4 = 16$

Triplets (3-itemsets)

| Itemset |
|---|
| { Beer, Diaper, Milk} |
| { Beer,Bread,Diaper} |
| {Bread,Diaper,Milk} |
| { Beer, Bread, Milk} |

# Illustrating Apriori Principle

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

Items (1-itemsets)

| Item | Count |
|------|-------|
| Bread | 4 |
| Coke | 2 |
| Milk | 4 |
| Beer | 3 |
| Diaper | 4 |
| Eggs | 1 |

Pairs (2-itemsets)

| Itemset | Count |
|---------|-------|
| {Bread,Milk} | 3 |
| {Bread,Beer} | 2 |
| {Bread,Diaper} | 3 |
| {Milk,Beer} | 2 |
| {Milk,Diaper} | 3 |
| {Beer,Diaper} | 3 |

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

Triplets (3-itemsets)

| Itemset | Count |
|---------|-------|
| { Beer, Diaper, Milk} | 2 |
| { Beer,Bread, Diaper} | 2 |
| {Bread, Diaper, Milk} | 2 |
| {Beer, Bread, Milk} | 1 |

If every subset is considered,
$$^6C_1 + {}^6C_2 + {}^6C_3$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
$$6 + 6 + 4 = 16$$

# Illustrating Apriori Principle

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

Items (1-itemsets)

| Item | Count |
|---|---|
| Bread | 4 |
| Coke | 2 |
| Milk | 4 |
| Beer | 3 |
| Diaper | 4 |
| Eggs | 1 |

Pairs (2-itemsets)

| Itemset | Count |
|---|---|
| {Bread,Milk} | 3 |
| {Bread,Beer} | 2 |
| {Bread,Diaper} | 3 |
| {Milk,Beer} | 2 |
| {Milk,Diaper} | 3 |
| {Beer,Diaper} | 3 |

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

Triplets (3-itemsets)

| Itemset | Count |
|---|---|
| { Beer, Diaper, Milk} | 2 |
| { Beer,Bread, Diaper} | 2 |
| {Bread, Diaper, Milk} | 2 |
| {Beer, Bread, Milk} | 1 |

If every subset is considered,
$$^6C_1 + {}^6C_2 + {}^6C_3$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
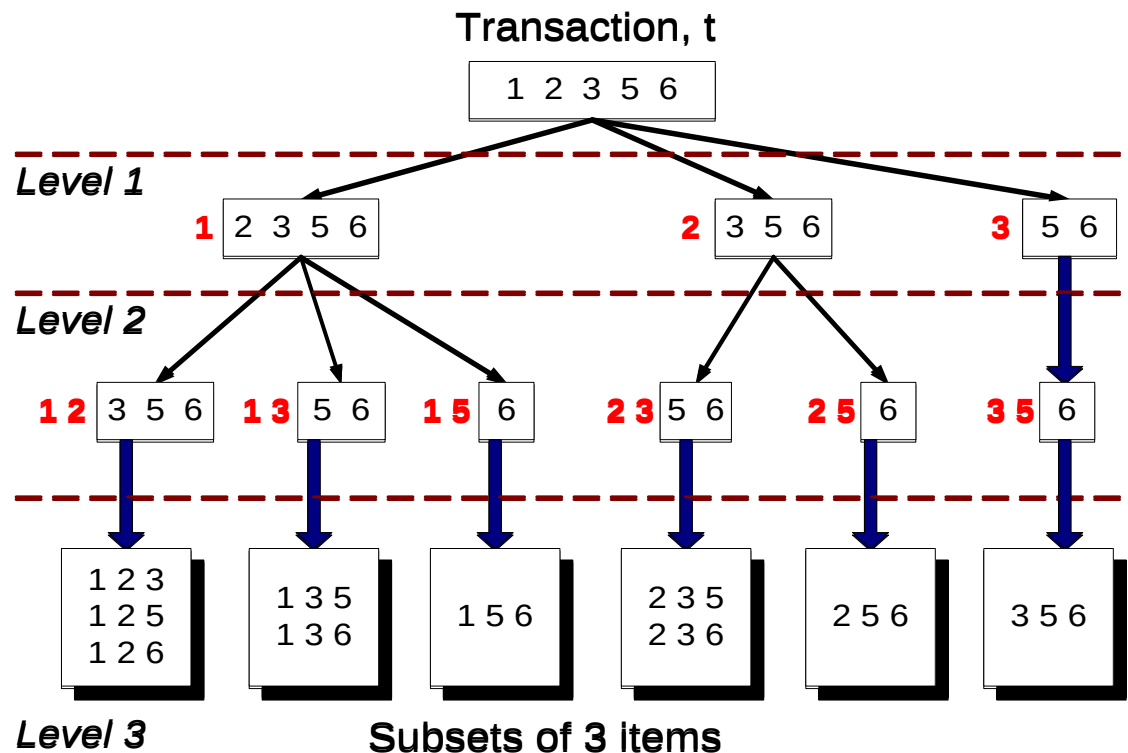$$6 + 6 + 4 = 16$$
$$6 + 6 + 1 = 13$$

# Apriori Algorithm

- Method:

  - Let k=1

  - Generate frequent itemsets of length 1

  - Repeat until no new frequent itemsets are identified

    - Generate length(k+1) candidate itemsets from length k frequent itemsets

    - Prune candidate itemsets containing subset of length k that are infrequent

    - Count the support of each candidate by scanning the DB

    - Eliminate candidates that are infrequent, leaving only those that are frequent

# Support Counting: An Example

Suppose you have 15 candidate itemsets of length 3:

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

How many of these itemsets are supported by transaction (1,2,3,5,6)?



Transaction, t

1 2 3 5 6

Level 1

Level 2

Subsets of 3 items

Level 3

**Introduction to Data Mining, 2nd Edition**
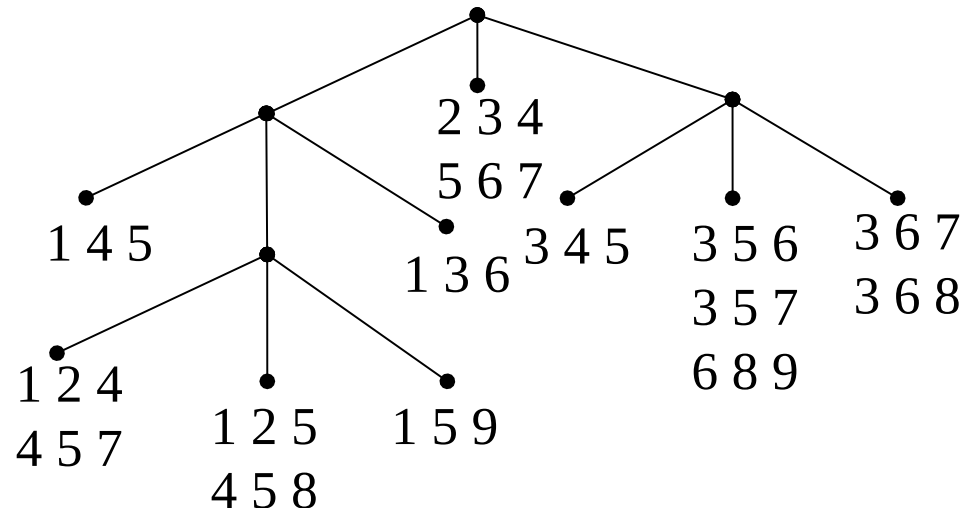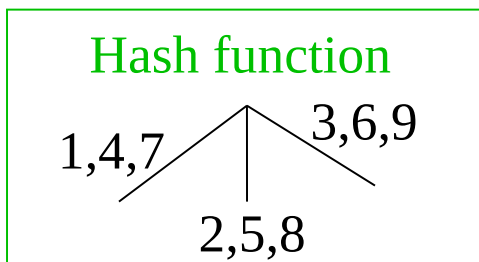
# Support Counting Using a Hash Tree

**Suppose you have 15 candidate itemsets of length 3:**

**{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4}, {5 6 7}, {3 4 5}, {3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}**
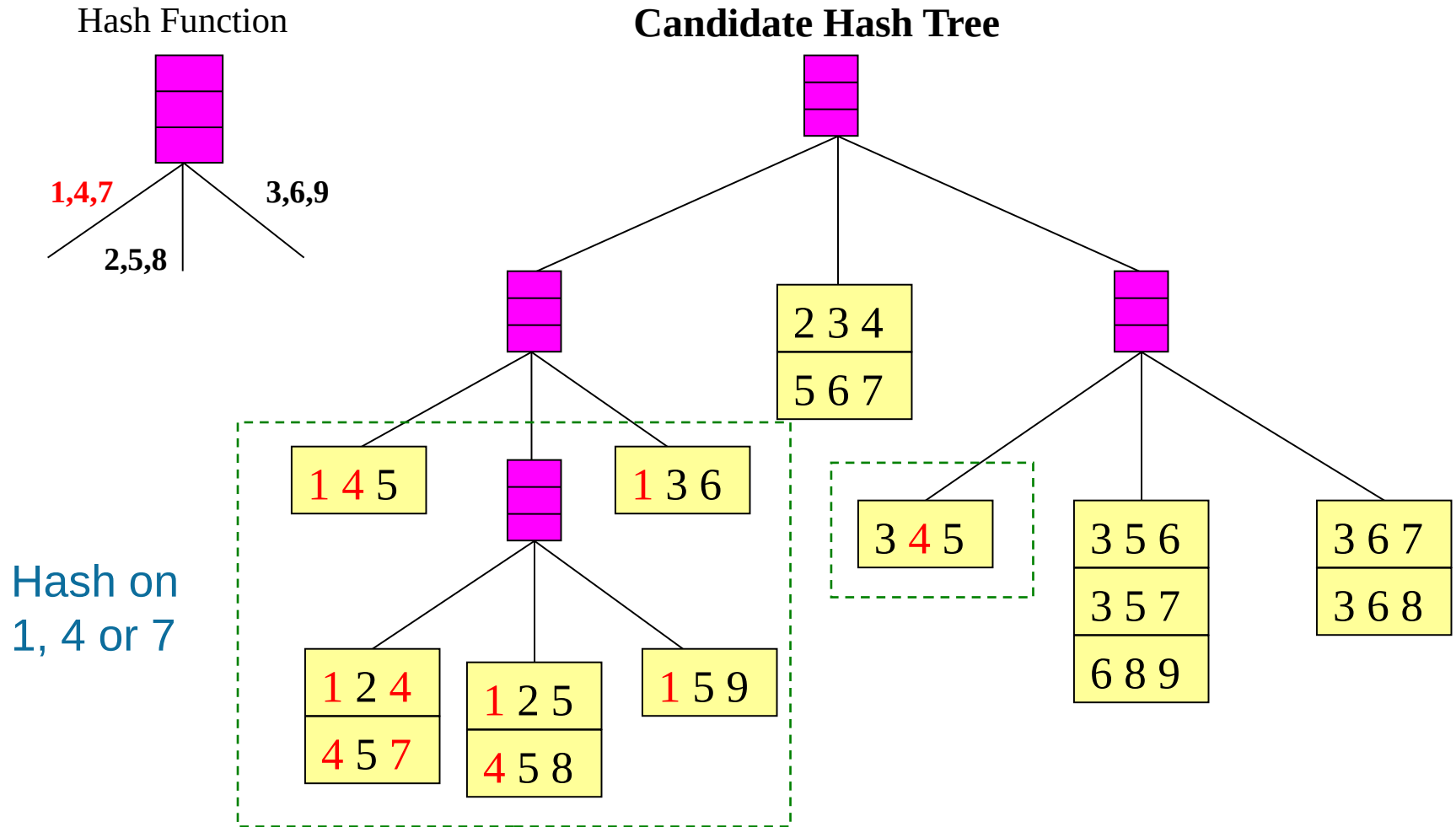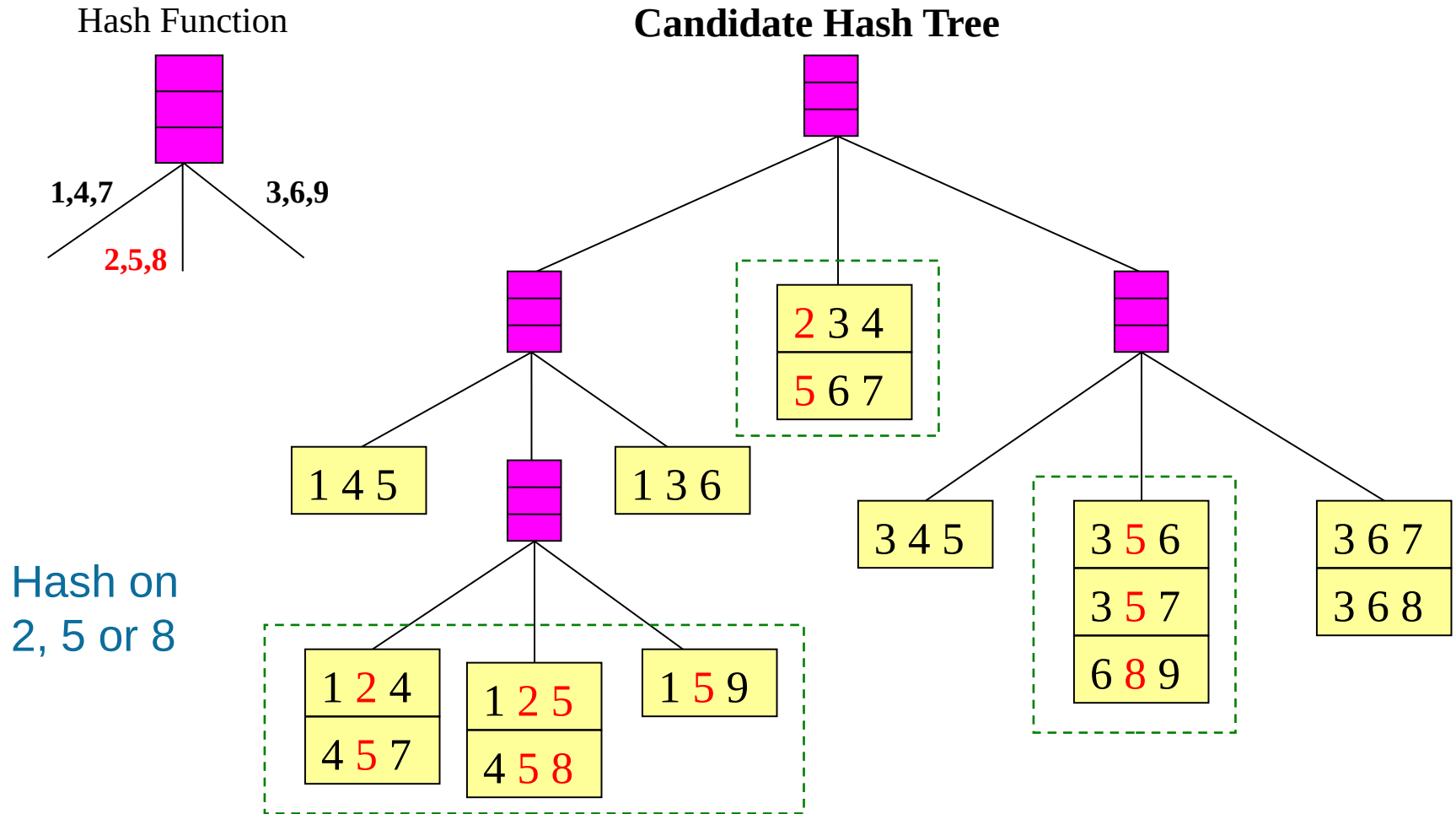
**You need:**

✁ **Hash function**

✁ **Max leaf size: max number of itemsets stored in a leaf node (if number of candidate itemsets exceeds max leaf size, split the node)**

# Support Counting Using a Hash Tree

Hash Function

**Candidate Hash Tree**

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

Hash on
1, 4 or 7

2 3 4
5 6 7

1 4 5    1 3 6

1 2 4
4 5 7

1 2 5
4 5 8

1 5 9

3 4 5

3 5 6
3 5 7
6 8 9

3 6 7
3 6 8

# Support Counting Using a Hash Tree



Hash Function

Candidate Hash Tree

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

Hash on
2, 5 or 8

2 3 4
5 6 7

1 4 5    1 3 6

1 2 4    1 2 5    1 5 9
4 5 7    4 5 8

3 4 5    3 5 6    3 6 7
         3 5 7    3 6 8
         6 8 9

# Support Counting Using a Hash Tree

Hash Function

Candidate Hash Tree

1,4,7      3,6,9

2,5,8

Hash on
3, 6 or 9

2 3 4
5 6 7

1 4 5

1 3 6

3 4 5

3 5 6
3 5 7
6 8 9

3 6 7
3 6 8

1 2 4
4 5 7

1 2 5
4 5 8

1 5 9

# Support Counting Using a Hash Tree

# Support Counting Using a Hash Tree

# Support Counting Using a Hash Tree

1 2 3 5 6  transaction

Hash Function

1,4,7          3,6,9

2,5,8

1 + 2 3 5 6

2 + 3 5 6

1 2 + 3 5 6

3 + 5 6

1 3 + 5 6

1 5 + 6

2 3 4
5 6 7

1 4 5

1 3 6

3 4 5

3 5 6
3 5 7
6 8 9

3 6 7
3 6 8

1 2 4
4 5 7

1 2 5
4 5 8

1 5 9

Match transaction against 11 out of 15 candidates

# Rule Generation

- Given a frequent itemset L, find all non-empty subsets f ⊂ L such that f → L – f satisfies the minimum confidence requirement

    – If {A,B,C,D} is a frequent itemset, candidate rules:

    ABC → D,         ABD → C,         ACD → B,         BCD → A,
    A → BCD, B → ACD, C → ABD,        D → ABC
    AB → CD, AC → BD,        AD → BC,        BC → AD,
    BD → AC,        CD → AB,

- If |L| = k, then there are $2^k – 2$ candidate association rules (ignoring L → ∅ and ∅ → L)

# Rule Generation

- In general, confidence does not have an anti-monotone property

    c(ABC → D) can be larger or smaller than c(AB → D)

- But confidence of rules generated from the same itemset has an anti-monotone property

    – E.g., Suppose {A,B,C,D} is a frequent 4-itemset:

    $$c(ABC \to D) \geq c(AB \to CD) \geq c(A \to BCD)$$

    – Confidence is anti-monotone w.r.t. number of items on the RHS of the rule

# Rule Generation for Apriori Algorithm

Lattice of rules



Low Confidence Rule

**Pruned Rules**

# Association Analysis: Basic Concepts and Algorithms

Algorithms and Complexity

# Factors Affecting Complexity of Apriori

- Choice of minimum support threshold

- Dimensionality (number of items) of the data set

- Size of database

- Average transaction width
  -

# Factors Affecting Complexity of Apriori

- Choice of minimum support threshold
  - lowering support threshold results in more frequent itemsets
  - this may increase number of candidates and max length of frequent itemsets
- Dimensionality (number of items) of the data set
  –

- Size of database
  –

- Average transaction width
  –

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

# Impact of Support Based Pruning

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

Items (1-itemsets)

| Item | Count |
|------|-------|
| Bread | 4 |
| Coke | 2 |
| Milk | 4 |
| Beer | 3 |
| Diaper | 4 |
| Eggs | 1 |

Minimum Support = 3

If every subset is considered,

$$^6C_1 + {}^6C_2 + {}^6C_3$$
$$6 + 15 + 20 = 41$$

With support-based pruning,

$$6 + 6 + 4 = 16$$

Minimum Support = 2

If every subset is considered,
$$^6C_1 + {}^6C_2 + {}^6C_3 + {}^6C_4$$
$$6 + 15 + 20 + 15 = 56$$

# Factors Affecting Complexity of Apriori

- Choice of minimum support threshold
    - lowering support threshold results in more frequent itemsets
    - this may increase number of candidates and max length of frequent itemsets
- Dimensionality (number of items) of the data set
    - More space is needed to store support count of  itemsets
    - if number of frequent itemsets also increases, both computation and I/O costs may also increase
- Size of database

- Average transaction width
    -

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

# Factors Affecting Complexity of Apriori

- Choice of minimum support threshold
  - lowering support threshold results in more frequent itemsets
  - this may increase number of candidates and max length of frequent itemsets
- Dimensionality (number of items) of the data set
  - More space is needed to store support count of itemsets
  - if number of frequent itemsets also increases, both computation and I/O costs may also increase
- Size of database
  - run time of algorithm increases with number of transactions
- Average transaction width

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

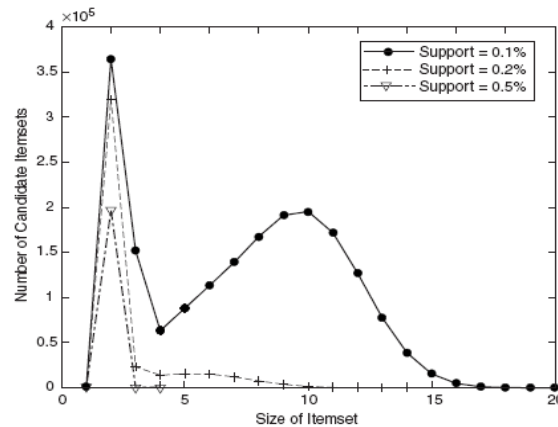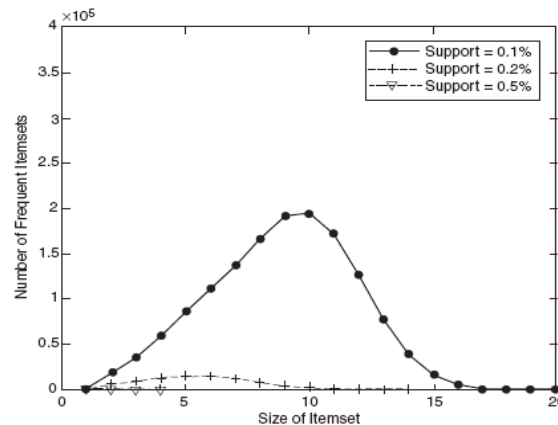# Factors Affecting Complexity of Apriori

- Choice of minimum support threshold
  - lowering support threshold results in more frequent itemsets
  - this may increase number of candidates and max length of frequent itemsets
- Dimensionality (number of items) of the data set
  - More space is needed to store support count of itemsets
  - if number of frequent itemsets also increases, both computation and I/O costs may also increase
- Size of database
  - run time of algorithm increases with number of transactions
- Average transaction width
  - transaction width increases the max length of frequent itemsets
  - number of subsets in a transaction increases with its width, increasing computation time for support counting

# Factors Affecting Complexity of Apriori



(a) Number of candidate itemsets.

(b) Number of frequent itemsets.

**Figure 6.13.** Effect of support threshold on the number of candidate and frequent itemsets.
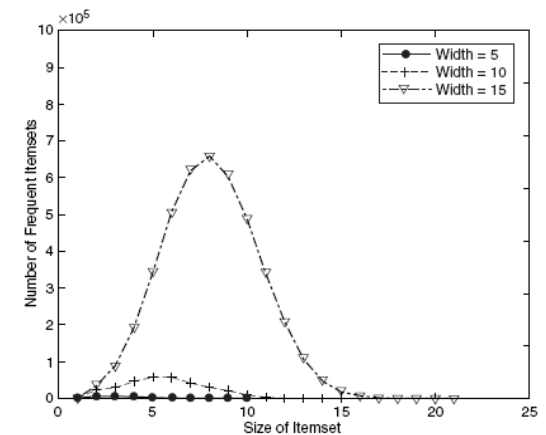


(a) Number of candidate itemsets.

(b) Number of Frequent Itemsets.

**Figure 6.14.** Effect of average transaction width on the number of candidate and frequent itemsets.

# Compact Representation of Frequent Itemsets

- Some frequent itemsets are redundant because their supersets are also frequent

Consider the following data set.  Assume support threshold =5

| TID | A1 | A2 | A3 | A4 | A5 | A6 | A7 | A8 | A9 | A10 | B1 | B2 | B3 | B4 | B5 | B6 | B7 | B8 | B9 | B10 | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 |
|-----|----|----|----|----|----|----|----|----|----|-----|----|----|----|----|----|----|----|----|----|-----|----|----|----|----|----|----|----|----|----|-----|
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

Number of frequent itemsets $= 3 \times \sum_{k=1}^{10} \binom{10}{k}$

- Need a compact representation

# Maximal Frequent Itemset

**An itemset is maximal frequent if it is frequent and none of its immediate supersets is frequent**