# Data Mining and Data Warehousing
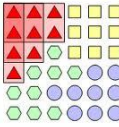
**Chapter 2**

Data warehousing

Instructor: Suresh Pokharel

ME in ICT (Asian Institute of Technology, Thailand)

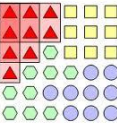BE in Computer ( NCIT, Pokhara University)

# Data Mining Tasks

1. **Classification:** learning a function that maps an item into one of a set of predefined classes

2. **Regression:** learning a function that maps an item to a real value

3. **Clustering:** identify a set of groups of similar items

4. **Dependencies and associations:** identify significant dependencies between data attributes

5. **Summarization:** find a compact description of the dataset or a subset of the dataset

**Data Mining and Data Warehousing** By: Suresh Pokharel

# Data Mining Methods

1. Decision Tree Classifiers:

   Used for modeling, classification

2. Association Rules:

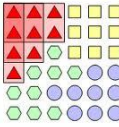   Used to find associations between sets of attributes

3. Sequential patterns:

   Used to find temporal associations in time series
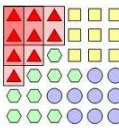
4. Hierarchical clustering:

   Used to group customers, web users, etc

# What Is Frequent Pattern Analysis?

- Frequent pattern: a pattern (a set of items, subsequences, substructures, etc.) 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?— 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, Web log (click stream) analysis, and DNA sequence analysis.
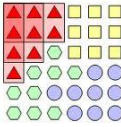
# Association Rule Mining

▪ 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},

**Data Mining and Data Warehousing** By: Suresh Pokharel

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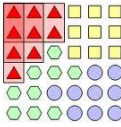    - E.g. $\sigma(\{$Milk, Bread,Diaper$\}) = 2$

- **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}

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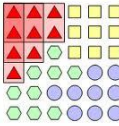    - ◆ Measures how often items in Y appear in transactions that contain X

| 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:

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

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

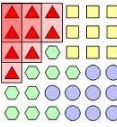$$c = \frac{\sigma(Milk, Diaper, Beer)}{\sigma(Milk, Diaper)} = \frac{2}{3} = 0.67$$

# 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

- 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!

**Data Mining and Data Warehousing** By: Suresh Pokharel

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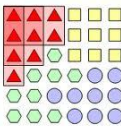{Milk} $\rightarrow$ {Diaper,Beer} (s=0.4, c=0.5)

## 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

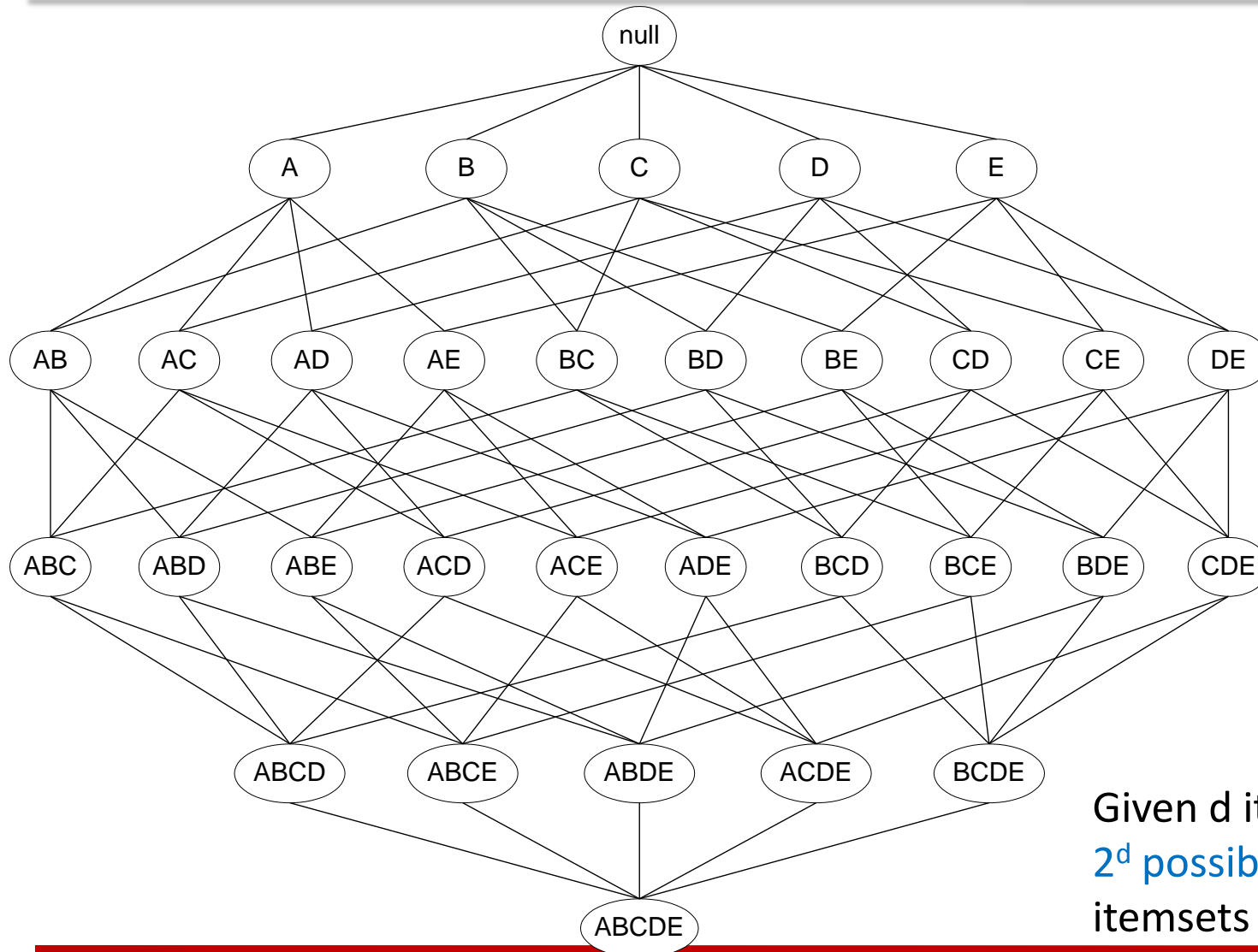**Data Mining and Data Warehousing**
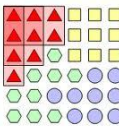
# Mining Association Rules

- Two-step approach:

  - **Frequent Itemset Generation**
    - Generate all itemsets whose support $\geq$ minsup

  - **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

# Frequent Itemset Generation

null

A    B    C    D    E

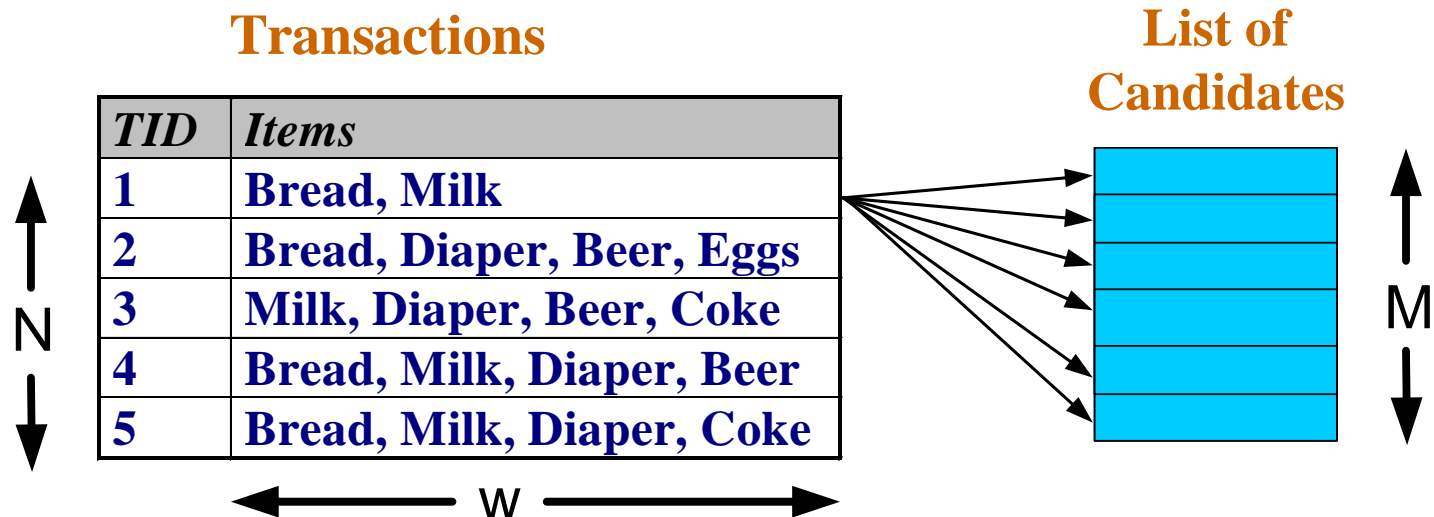AB   AC   AD   AE   BC   BD   BE   CD   CE   DE

ABC  ABD  ABE  ACD  ACE  ADE  BCD  BCE  BDE  CDE

ABCD  ABCE  ABDE  ACDE  BCDE

ABCDE

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

Data Mining and Data Warehousing
By: Suresh Pokharel

# 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
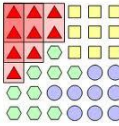
Match each transaction against every candidate

Complexity ~ O(NMw) => Expensive since M = $2^d$ !!!

# Reducing Number of Candidates

**Apriori principle**:

If an itemset is frequent, then all of its subsets must also be frequent

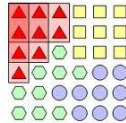Apriori principle holds due to the following property of the support measure:

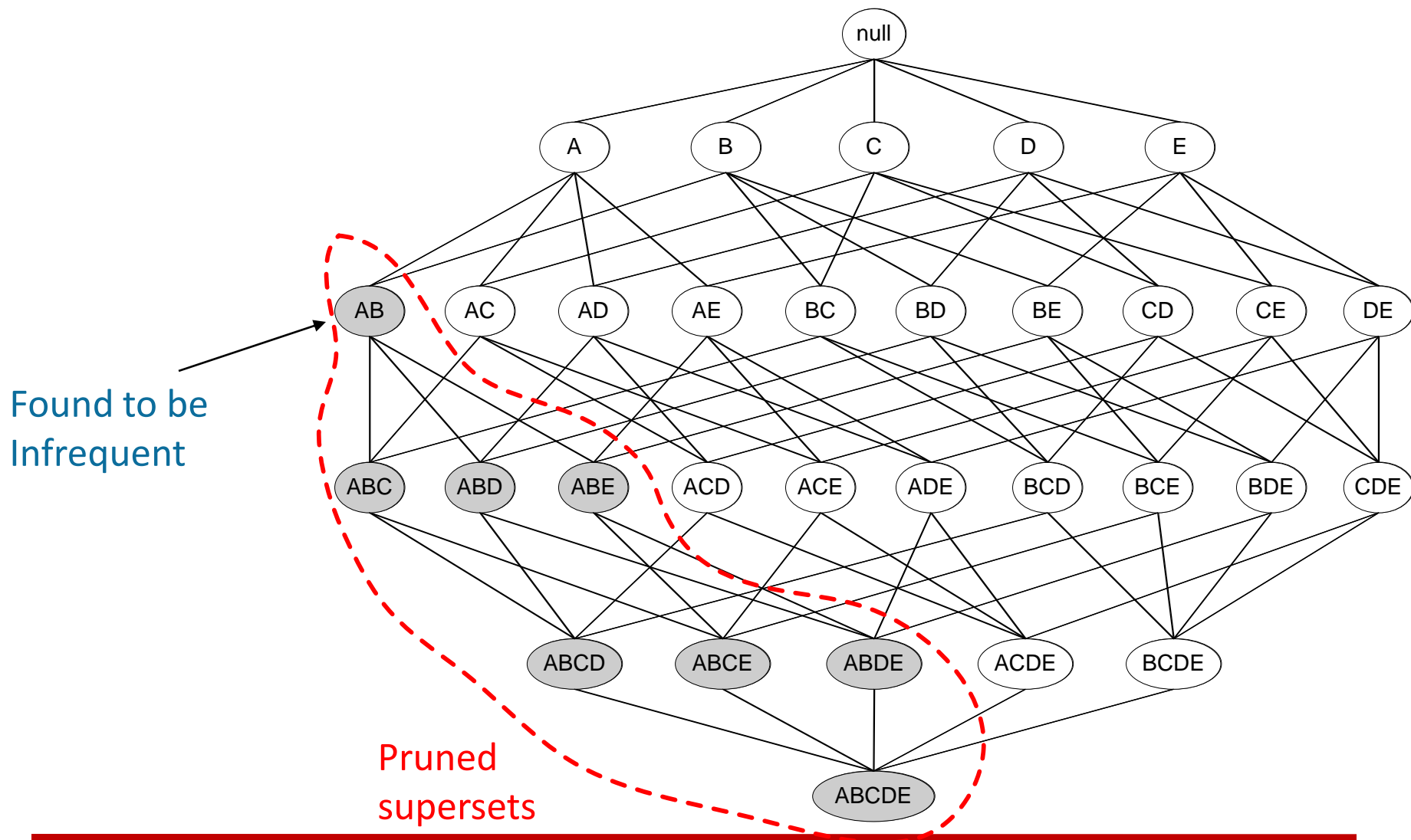$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

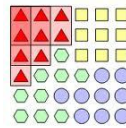Support of an itemset never exceeds the support of its subsets

This is known as the anti-monotone property of support

Anti-monotone: if a set can't pass a test, all of its superset will fail the same test as well

**Data Mining and Data Warehousing**         By: Suresh Pokharel

# Illustrating Apriori Principle



Found to be Infrequent

Pruned supersets

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

Items (1-itemsets)

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

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

**Minimum Support = 3**

Triplets (3-itemsets)

| Itemset | Count |
|---------|-------|
| **{Bread,Milk,Diaper}** | **3** |

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

Q: Total number of possible frequent itemsets ???

# 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 subsets of length k that are infrequent
  - Count the support of each candidate by scanning the DB
  - Eliminate(prune) candidates that are infrequent, leaving only those that are frequent

**Data Mining and Data Warehousing** By: Suresh Pokharel

Database D

| TID | Items |
|-----|-------|
| 100 | 1 3 4 |
| 200 | 2 3 5 |
| 300 | 1 2 3 5 |
| 400 | 2 5 |

Scan D →

$C_1$

| itemset | sup. |
|---------|------|
| {1} | 2 |
| {2} | 3 |
| {3} | 3 |
| {4} | 1 |
| {5} | 3 |

$L_1$

| itemset | sup. |
|---------|------|
| {1} | 2 |
| {2} | 3 |
| {3} | 3 |
| {5} | 3 |

$C_2$

| itemset |
|---------|
| {1 2} |
| {1 3} |
| {1 5} |
| {2 3} |
| {2 5} |
| {3 5} |

$C_2$

| itemset | sup |
|---------|-----|
| {1 2} | 1 |
| {1 3} | 2 |
| {1 5} | 1 |
| {2 3} | 2 |
| {2 5} | 3 |
| {3 5} | 2 |

Scan D →

$L_2$

| itemset | sup |
|---------|-----|
| {1 3} | 2 |
| {2 3} | 2 |
| {2 5} | 3 |
| {3 5} | 2 |

$C_3$

| itemset |
|---------|
| {2 3 5} |

Scan D →

$L_3$

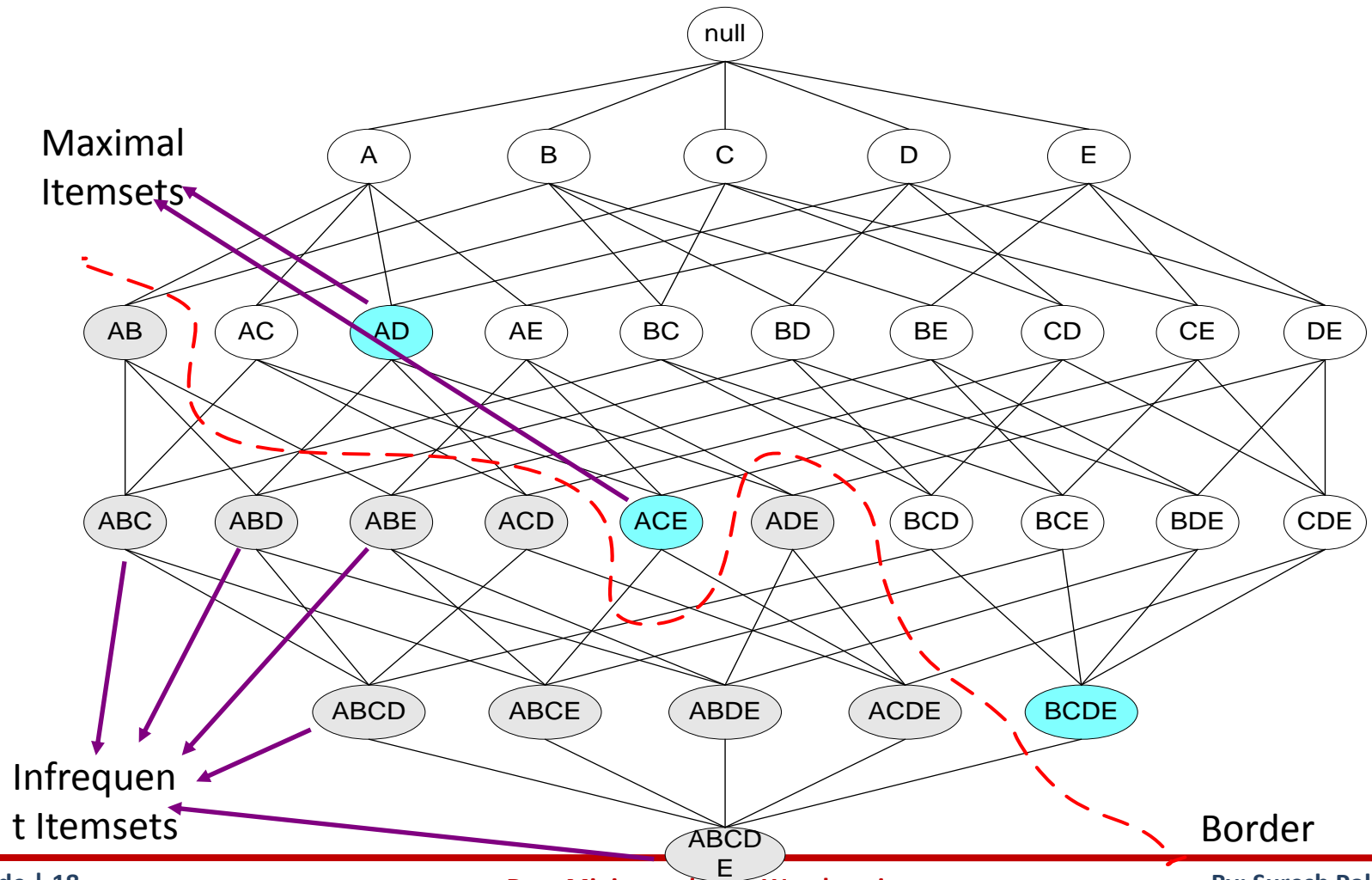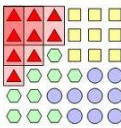| itemset | sup |
|---------|-----|
| {2 3 5} | 2 |

Why {1 2 3}, {1 2 5}, {1 3 5} are not listed in C3???

# Maximal Frequent Itemset

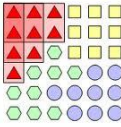An itemset is maximal frequent if none of its immediate supersets is frequent



**Data Mining and Data Warehousing**   **By: Suresh Pokharel**

# Closed Itemset

An itemset is closed if none of its immediate supersets has the same support as the itemset

| TID | Items |
|-----|-------|
| 1 | {A,B} |
| 2 | {B,C,D} |
| 3 | {A,B,C,D} |
| 4 | {A,B,D} |
| 5 | {A,B,C,D} |

| Itemset | Support |
|---------|---------|
| {A} | 4 |
| {B} | 5 |
| {C} | 3 |
| {D} | 4 |
| {A,B} | 4 |
| {A,C} | 2 |
| {A,D} | 3 |
| {B,C} | 3 |
| {B,D} | 4 |
| {C,D} | 3 |

| Itemset | Support |
|---------|---------|
| {A,B,C} | 2 |
| {A,B,D} | 3 |
| {A,C,D} | 2 |
| {B,C,D} | 3 |
| {A,B,C,D} | 2 |

**Data Mining and Data Warehousing**

By: Suresh Pokharel

# Maximal vs Closed Itemsets

| TID | Items |
|-----|-------|
| 1 | ABC |
| 2 | ABCD |
| 3 | BCE |
| 4 | ACDE |
| 5 | DE |

Transaction Ids

null

124 A    123 B    1234 C    245 D    345 E

12 AB   124 AC   24 AD   4 AE   123 BC   2 BD   3 BE   24 CD   34 CE   45 DE

12 ABC   2 ABD   ABE   24 ACD   4 ACE   4 ADE   2 BCD   3 BCE   BDE   4 CDE

2 ABCD   ABCE   ABDE   4 ACDE   BCDE

ABCDE

Not supported by any transactions

**Data Mining and Data Warehousing**

**By: Suresh Pokharel**

# Maximal vs Closed Frequent Itemsets

Minimum support = 2

Closed but not maximal

Closed and maximal

**null**

| | | | | |
|---|---|---|---|---|
| **124** | **123** | **1234** | **245** | **345** |
| A | B | C | D | E |

| **12** | **124** | **24** | **4** | **123** | **2** | **3** | **24** | **34** | **45** |
|---|---|---|---|---|---|---|---|---|---|
| AB | AC | AD | AE | BC | BD | BE | CD | CE | DE |

| **12** | **2** | | **24** | **4** | **4** | **2** | **3** | | **4** |
|---|---|---|---|---|---|---|---|---|---|
| ABC | ABD | ABE | ACD | ACE | ADE | BCD | BCE | BDE | CDE |

| **2** | | | **4** | |
|---|---|---|---|---|
| ABCD | ABCE | ABDE | ACDE | BCDE |

ABCDE

# Closed = 9

# Maximal = 4

| TID | Items |
|-----|-------|
| 1 | ABC |
| 2 | ABCD |
| 3 | BCE |
| 4 | ACDE |
| 5 | DE |

**Data Mining and Data Warehousing**

By: Suresh Pokharel

# Maximal vs Closed Itemsets

Frequent
Itemsets

Closed
Frequent
Itemsets

Maximal
Frequent
Itemsets

**Data Mining and Data Warehousing**
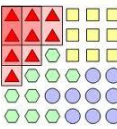
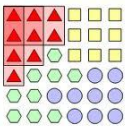**By: Suresh Pokharel**

# Frequent Pattern Tree
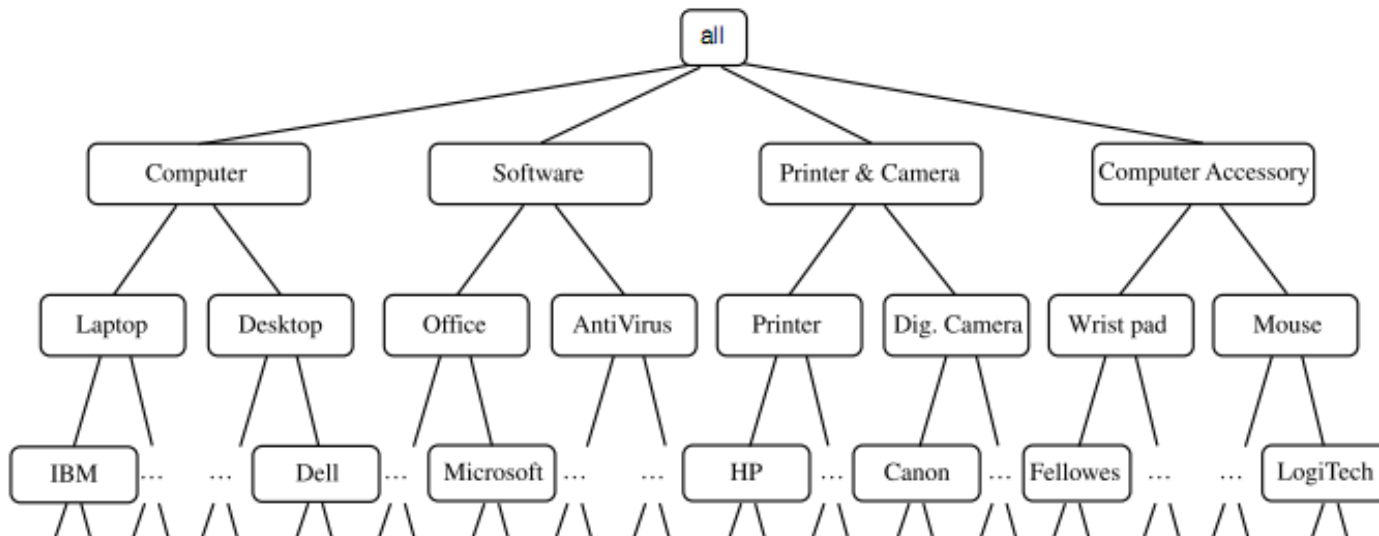
# Generating Association Rule (Example)

□ **Given a frequent itemset L**

- ■ Find all non-empty subsets F in L, such that the association rule $F \Rightarrow \{L-F\}$ satisfies the minimum confidence
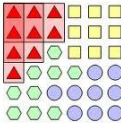
- ■ Create the rule $F \Rightarrow \{L-F\}$

□ **If L={A,B,C}**

- ■ The candidate itemsets are: $AB \Rightarrow C$, $AC \Rightarrow B$, $BC \Rightarrow A$, $A \Rightarrow BC$, $B \Rightarrow AC$, $C \Rightarrow AB$

- ■ In general, there are $2^K - 2$ candidate solutions, where k is the length of the itemset L

**Data Mining and Data Warehousing**
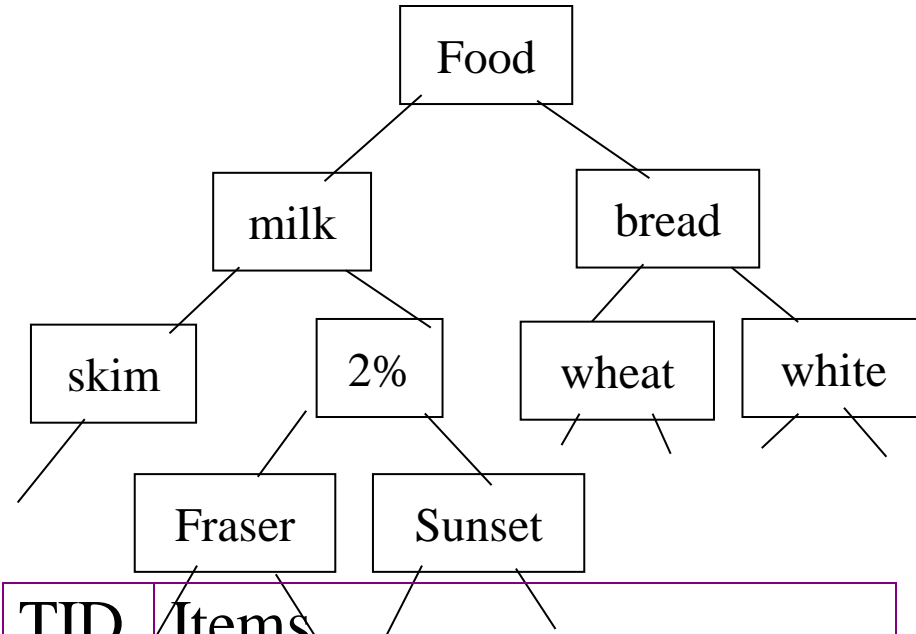
**By: Suresh Pokharel**

# Recap : A Concept Hierarchy

| TID | Items Purchased |
|-----|-----------------|
| T100 | IBM-ThinkPad-T40/2373, HP-Photosmart-7660 |
| T200 | Microsoft-Office-Professional-2003, Microsoft-Plus!-Digital-Media |
| T300 | Logitech-MX700-Cordless-Mouse, Fellowes-Wrist-Rest |
| T400 | Dell-Dimension-XPS, Canon-PowerShot-S400 |
| T500 | IBM-ThinkPad-R40/P4M, Symantec-Norton-Antivirus-2003 |
| … | … |

Data Mining and Data Warehousing
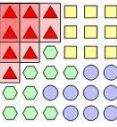
lecture 2

By : Suresh Pokharel

# Multiple-Level Association Rules

- Items often form hierarchy.

- Items at the lower level are expected to have lower support.

- Rules regarding itemsets at appropriate levels could be quite useful.
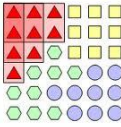
- We can explore shared multi-level mining

```
                    Food
                   /    \
                milk      bread
               /    \     /    \
            skim    2%  wheat  white
                   /  \
               Fraser  Sunset
```

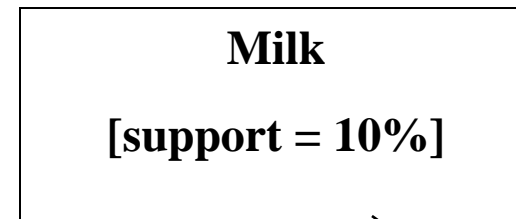| TID | Items |
|-----|-------|
| T1 | {111, 121, 211, 221} |
| T2 | {111, 211, 222, 323} |
| T3 | {112, 122, 221, 411} |
| T4 | {111, 121} |
| T5 | {111, 122, 211, 221, 413} |

# Mining Multi-Level Associations

- A top_down, progressive deepening approach:
  - First find high-level strong rules:

    milk $\rightarrow$ bread [20%, 60%].
  - Then find their lower-level "weaker" rules:

    2% milk $\rightarrow$ wheat bread [6%, 50%].

- Variations at mining multiple-level association rules.

  - Association rules with multiple, alternative hierarchies:
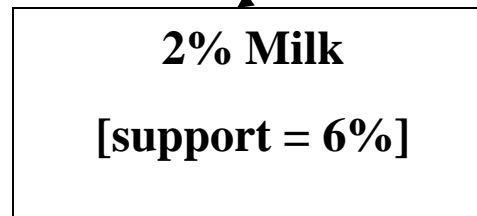
    2% *milk* $\rightarrow$ *Wonder bread*

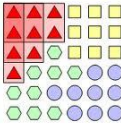**Data Mining and Data Warehousing**     **By: Suresh Pokharel**

# Multi-level mining with uniform support

**Level 1**
**min_sup = 5%**

**Milk**

**[support = 10%]**

**Level 2**
**min_sup = 5%**

**2% Milk**

**[support = 6%]**

**Skim Milk**

**[support = 4%]**

# Multi-level mining with reduced support

**Level 1**
**min_sup = 5%**

**Level 2**
**min_sup = 3%**

```
                    ┌─────────────────────┐
                    │        Milk         │
                    │  [support = 10%]    │
                    └─────────────────────┘
                       ↙              ↘
    ┌─────────────────────┐    ┌─────────────────────┐
    │      2% Milk        │    │     Skim Milk       │
    │  [support = 6%]     │    │  [support = 4%]     │
    └─────────────────────┘    └─────────────────────┘
```

# Interestingness Measurements

- Objective measures

  Two popular measurements:

  ☆ *support;*  and

  🕐 *confidence*

- Subjective  measures

  A rule (pattern) is interesting if

  ☆ it is *unexpected* (surprising to the user); and/or

  🕐 *actionable* (the user can do something with it)

**Data Mining and Data Warehousing**

**Data Mining and Data Warehousing**