

# Data Mining and Data Warehousing

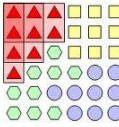
## Unit 7

### Association Rule

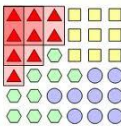
Instructor: Suresh Pokharel



# 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

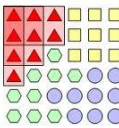
- Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction.

## Market-Basket transactions

<i><b>TID</b></i>	<i><b>Items</b></i>
<b>1</b>	<b>Bread, Milk</b>
<b>2</b>	<b>Bread, Diaper, Beer, Eggs</b>
<b>3</b>	<b>Milk, Diaper, Beer, Coke</b>
<b>4</b>	<b>Bread, Milk, Diaper, Beer</b>
<b>5</b>	<b>Bread, Milk, Diaper, Coke</b>

## Example of Association Rules

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



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

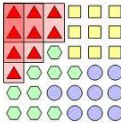
## ■ Support

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

## ■ Frequent Itemset

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

<i>TID</i>	<i>Items</i>
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:  
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

<i>TID</i>	<i>Items</i>
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  $Y$

- **Confidence (c)**

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

Example:

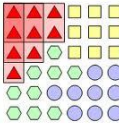
$\{\text{Milk, Diaper}\} \Rightarrow \text{Beer}$

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

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



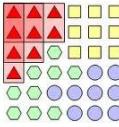
# Association Rule Mining Task



- Given a set of transactions  $T$ , the goal of association rule mining is to find all rules having
    - support  $\geq \textit{minsup}$  threshold
    - confidence  $\geq \textit{minconf}$  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



<i>TID</i>	<i>Items</i>
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:

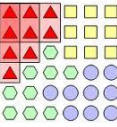
$\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\} (s=0.4, c=0.67)$   
 $\{\text{Milk, Beer}\} \rightarrow \{\text{Diaper}\} (s=0.4, c=1.0)$   
 $\{\text{Diaper, Beer}\} \rightarrow \{\text{Milk}\} (s=0.4, c=0.67)$   
 $\{\text{Beer}\} \rightarrow \{\text{Milk, Diaper}\} (s=0.4, c=0.67)$   
 $\{\text{Diaper}\}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

## Observations:

- All the above rules are binary partitions of the same itemset:  
 $\{\text{Milk, Diaper, Beer}\}$
- Rules originating from the same itemset have identical support but can have different confidence
- Thus, we may decouple the support and confidence requirements

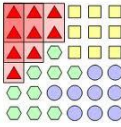


# Mining Association Rules

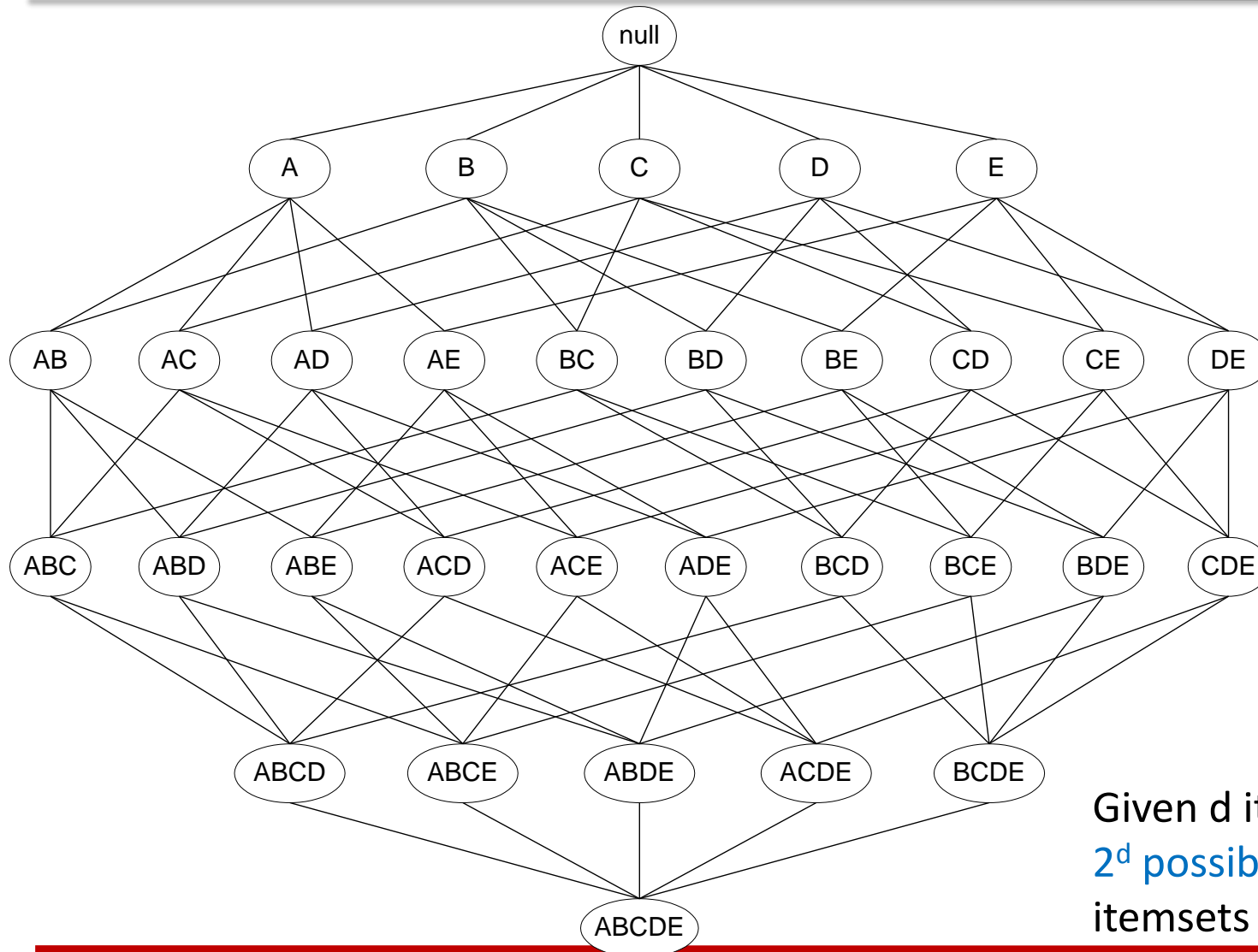


- Two-step approach:
  - Frequent Itemset Generation
    - Generate all itemsets whose support  $\geq$  minsup
  - Rule Generation
    - Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset
- ❖ Frequent itemset generation is still computationally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

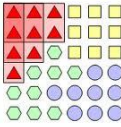




# Frequent Itemset Generation



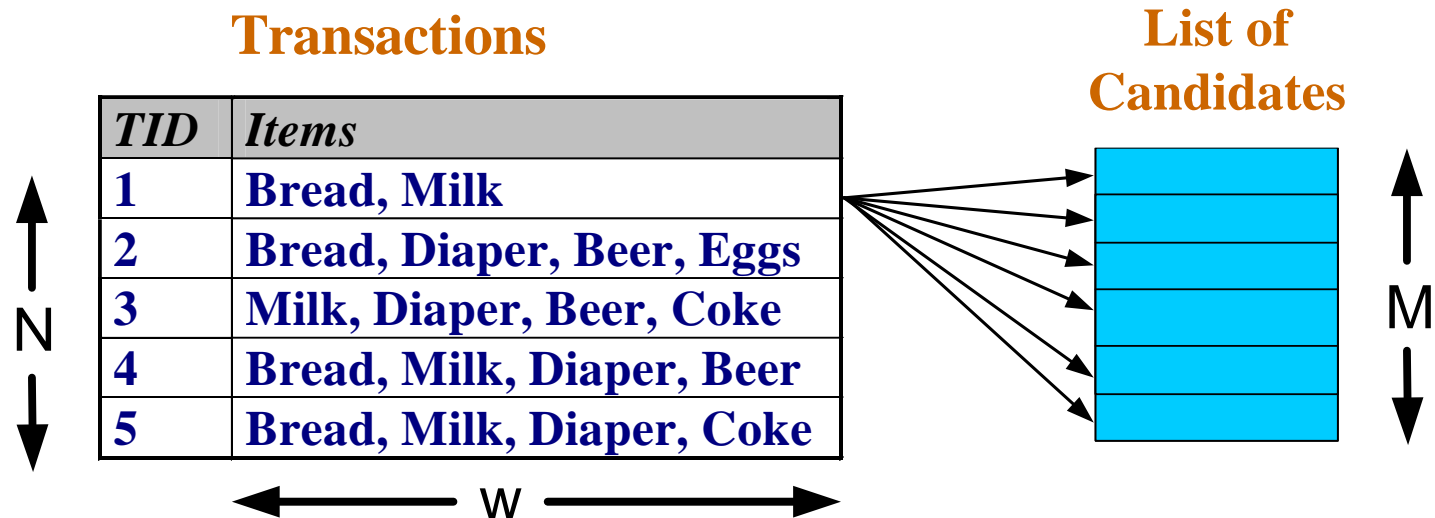
Given  $d$  items, there are  $2^d$  possible candidate itemsets



# 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

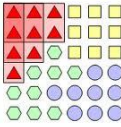


Match each transaction against every candidate

Complexity  $\sim O(NMw) \Rightarrow$  **Expensive since  $M = 2^d$  !!!**



# Reducing Number of Candidates



## Apriori principle:

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

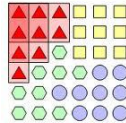
Apriori principle holds due to the following property of the support measure:

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

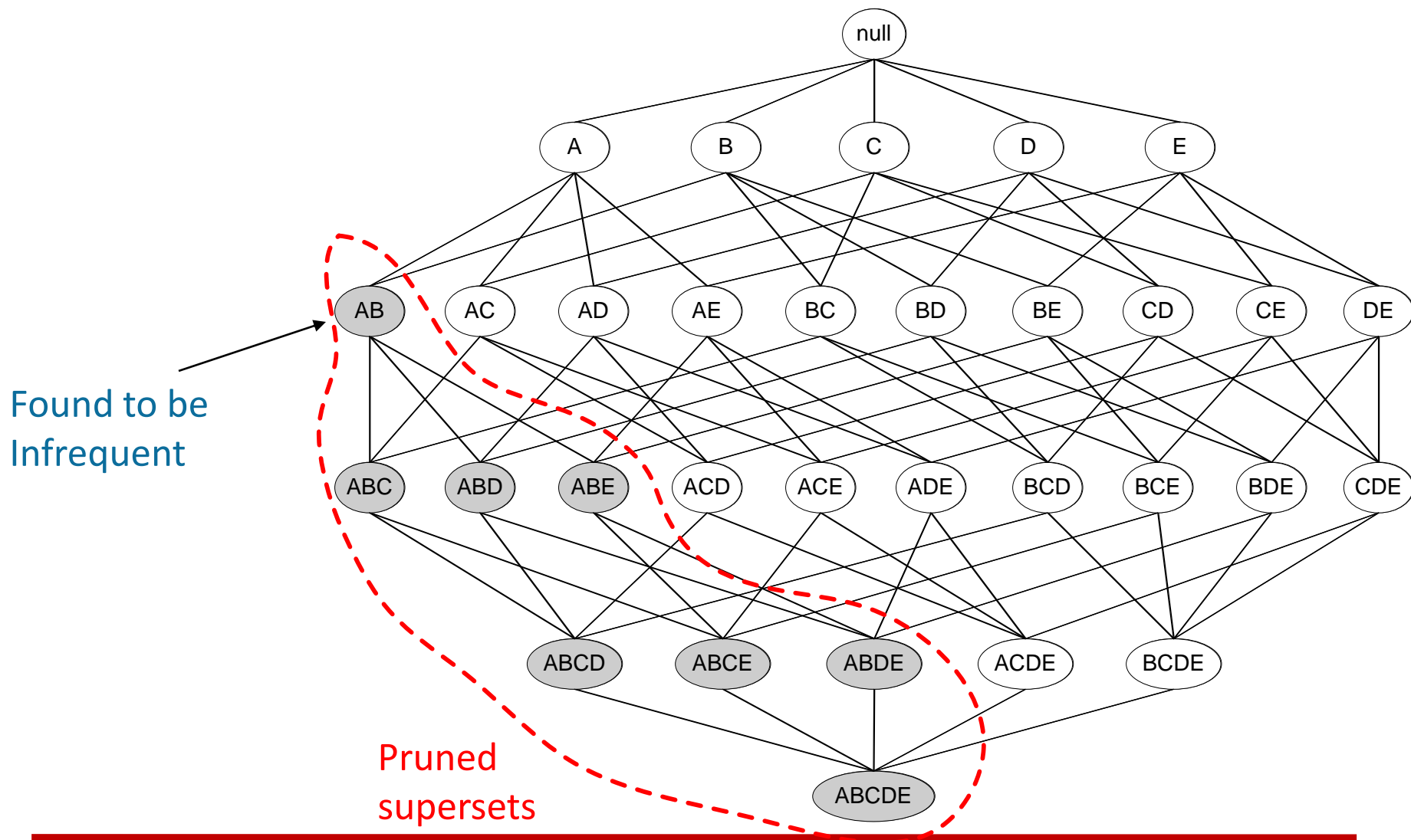
Support of an itemset never exceeds the support of its subsets

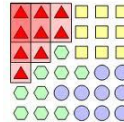
This is known as the **anti-monotone** property of support

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



# Illustrating Apriori Principle





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

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



Triplets (3-itemsets)

Itemset	Count
{Bread,Milk,Diaper}	3

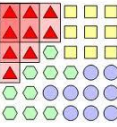


Q: Total number of possible frequent itemsets ???



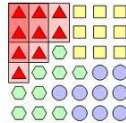
# Apriori Algorithm

---



## 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



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

$C_2$

$L_2$

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

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

Scan D

$C_2$

itemset
{1 2}
{1 3}
{1 5}
{2 3}
{2 5}
{3 5}

itemset
{1,2, 3}
{1,3, 5}
{2,3, 5}

$C_3$

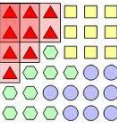
itemset
{2 3 5}

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???

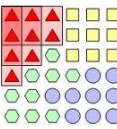


# Frequent Pattern Tree

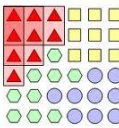




# Generating Association Rule (Example)

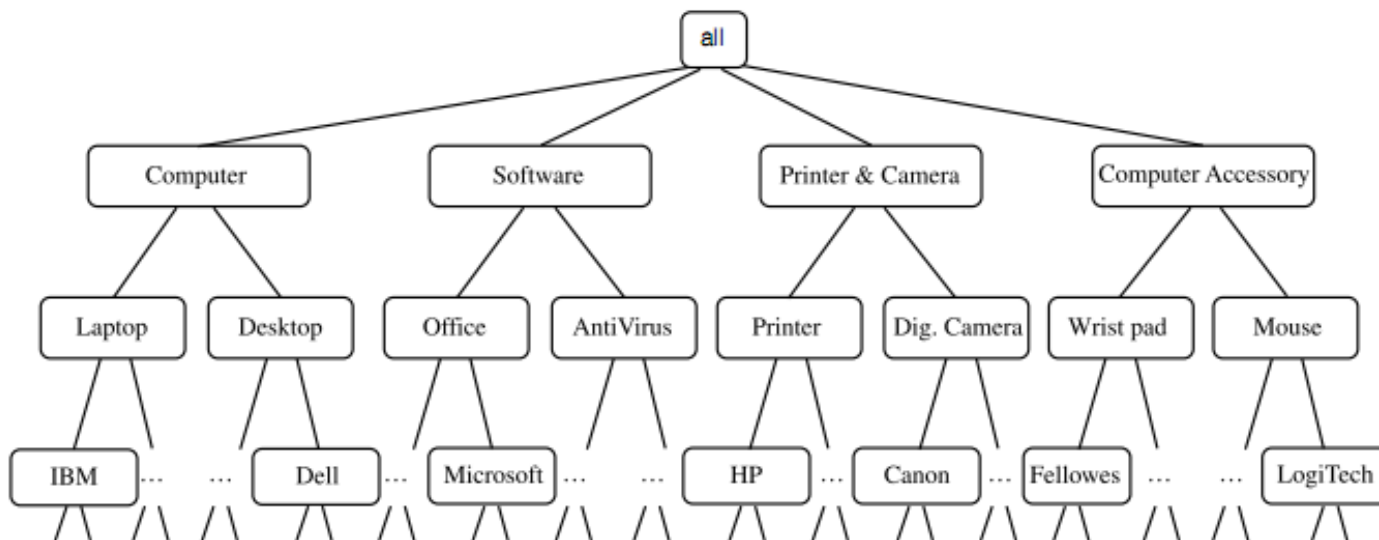


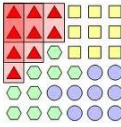
- **Given a frequent itemset L**
  - Find all non-empty subsets F in L, such that the association rule  $F \Rightarrow \{L-F\}$  satisfies the minimum confidence
  - Create the rule  $F \Rightarrow \{L-F\}$
  
- **If  $L=\{A,B,C\}$** 
  - The candidate itemsets are:  $AB \Rightarrow C$ ,  $AC \Rightarrow B$ ,  $BC \Rightarrow A$ ,  $A \Rightarrow BC$ ,  $B \Rightarrow AC$ ,  $C \Rightarrow AB$
  - In general, there are  $2^k - 2$  candidate solutions, where k is the length of the itemset L



# Recap : A Concept Hierarchy

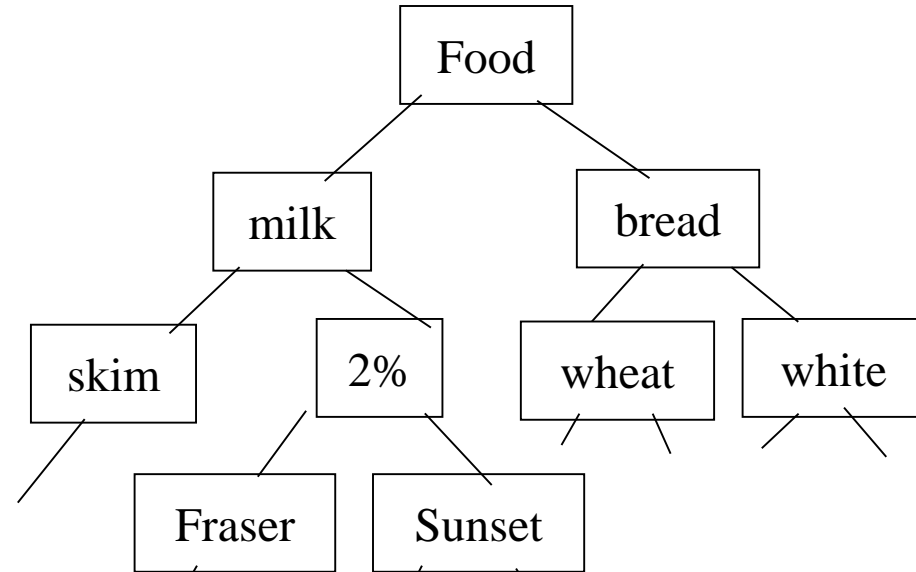
<i>TID</i>	<i>Items Purchased</i>
T100	IBM-ThinkPad-T40/2373, HP-Photosmart-7660
T200	Microsoft-Office-Professional-2003, Microsoft-Plus!-Digital-Media
T300	Logitech-MX700-Cordless-Mouse, Fellowes-Wrist-Rest
T400	Dell-Dimension-XPS, Canon-PowerShot-S400
T500	IBM-ThinkPad-R40/P4M, Symantec-Norton-Antivirus-2003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



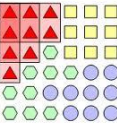


# Multiple-Level Association Rules

- 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



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}

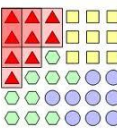


# Mining Multi-Level Associations

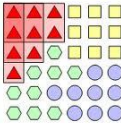
- A top\_down, progressive deepening approach:
  - First find high-level strong rules:  
milk  $\rightarrow$  bread [20%, 60%].
  - Then find their lower-level “weaker” rules:  
2% milk  $\rightarrow$  wheat bread [6%, 50%].
- Variations at mining multiple-level association rules.
  - Association rules with multiple, alternative hierarchies:  
2% milk  $\rightarrow$  *Wonder* bread



# Multi-level Association: Uniform Support vs. Reduced Support



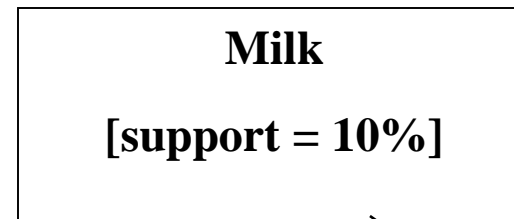
- Uniform Support: the same minimum support for all levels
  - + One minimum support threshold. No need to examine itemsets containing any item whose ancestors do not have minimum support.
  - – Lower level items do not occur as frequently. If support threshold
    - too high  $\Rightarrow$  miss low level associations
    - too low  $\Rightarrow$  generate too many high level associations
- Reduced Support: reduced minimum support at lower levels



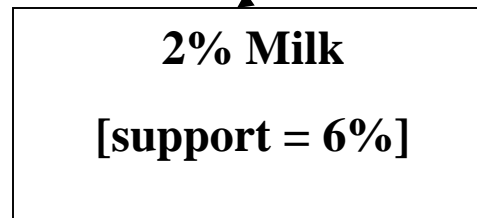
# Uniform Support

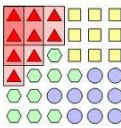
## Multi-level mining with uniform support

**Level 1**  
**min\_sup = 5%**



**Level 2**  
**min\_sup = 5%**



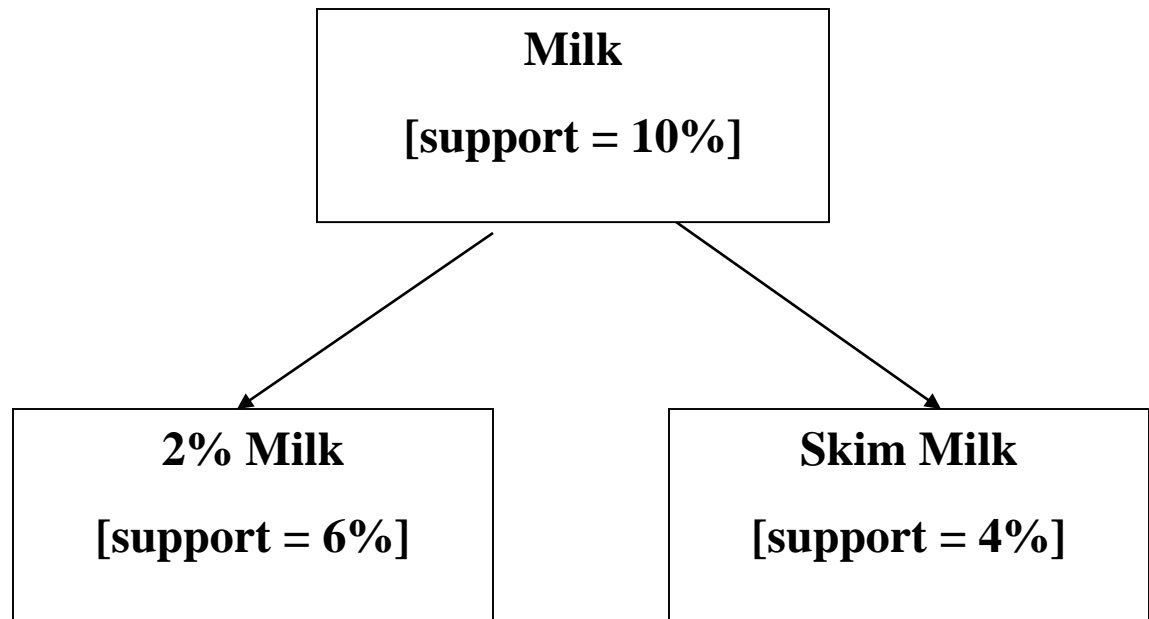


# Reduced Support

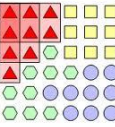
## Multi-level mining with reduced support

**Level 1**  
**min\_sup = 5%**

**Level 2**  
**min\_sup = 3%**

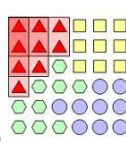


# Multi-level Association: Redundancy Filtering



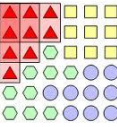
- Some rules may be redundant due to “ancestor” relationships between items.
- Example
  - milk  $\Rightarrow$  wheat bread [support = 8%, confidence = 70%]
  - 2% milk  $\Rightarrow$  wheat bread [support = 2%, confidence = 72%]
- We say the first rule is an ancestor of the second rule.





# Multi-Dimensional Association: Concepts

- Single-dimensional rules:  
 $\text{buys}(X, \text{"milk"}) \Rightarrow \text{buys}(X, \text{"bread"})$
- Multi-dimensional rules: ○ 2 dimensions or predicates
  - Inter-dimension association rules (*no repeated predicates*)  
 $\text{age}(X, \text{"19-25"}) \wedge \text{occupation}(X, \text{"student"}) \Rightarrow \text{buys}(X, \text{"coke"})$
  - hybrid-dimension association rules (*repeated predicates*)  
 $\text{age}(X, \text{"19-25"}) \wedge \text{buys}(X, \text{"popcorn"}) \Rightarrow \text{buys}(X, \text{"coke"})$



# 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)

