



# **Session 54**

## **Association Rules**



# Table of Content

## What will We Learn Today?

1. Basic concept
2. Frequent itemset
3. Apriori algorithm





# Basic Concept





# 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>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Vegetables, Fruits, Eggs
3	Milk, Vegetables, Fruits, Coke
4	Bread, Milk, Vegetables, Fruits
5	Bread, Milk, Vegetables, Coke

## Example of Association Rules

$\{\text{Vegetables}\} \rightarrow \{\text{Fruits}\},$   
 $\{\text{Milk, Bread}\} \rightarrow \{\text{Eggs, Coke}\},$   
 $\{\text{Fruits, Bread}\} \rightarrow \{\text{Milk}\},$

Implication means co-occurrence, not causality!







# Frequent Itemset





# Definition: Frequent Itemset

## Itemset

- A collection of one or more items
  - ⑩ Example: {Milk, Bread, Vegetables}
- k-itemset
  - ⑩ An itemset that contains k items

## Support count ( $\sigma$ )

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

## Support

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

## 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, Vegetables, Fruits, Eggs
3	Milk, Vegetables, Fruits, Coke
4	Bread, Milk, Vegetables, Fruits,
5	Bread, Milk, Vegetables, Coke



# Definition: Association Rule

## ● Association Rule

- An implication expression of the form  $X \rightarrow Y$ , where  $X$  and  $Y$  are itemsets
- $X$  = antecedent,  $Y$  = consequence
- Antecedent (if) and a consequent (then).
- An antecedent is an item found within the data.
- A consequent is an item found in combination with the antecedent.
- Example:  
 $\{\text{Milk, Vegetables}\} \rightarrow \{\text{Fruits}\}$

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

Example:

$\{\text{Milk, Vegetables}\} \Rightarrow \text{Fruits}$







# Definition: Association Rule

## 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
  - ◆ Indicates the number of times the if-then statements are found true.
- Lift
  - ◆ How many times an if-then statement is expected to be found true

$$\text{support}(A \rightarrow C) = \text{support}(A \cup C), \quad \text{range: } [0, 1]$$

$$\text{confidence}(A \rightarrow C) = \frac{\text{support}(A \rightarrow C)}{\text{support}(A)}, \quad \text{range: } [0, 1]$$

$$\text{lift}(A \rightarrow C) = \frac{\text{confidence}(A \rightarrow C)}{\text{support}(C)}, \quad \text{range: } [0, \infty]$$





# Definition: Association Rule

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Vegetables, Fruits, Eggs
3	Milk, Vegetables, Fruits, Coke
4	Bread, Milk, Vegetables, Fruits
5	Bread, Milk, Vegetables, Coke

**Example:**

$\{ \text{Milk , Vegetables} \} \Rightarrow \text{Fruits}$

$\{ X \} \Rightarrow Y$

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

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

$$l = \frac{\sigma(\text{Milk, Vegetables , Fruits})}{\sigma(\text{Milk , Vegetables}) * \sigma(\text{Fruits})} = \frac{2}{(3 * 3)} = 0.22$$



# Association Rule Mining Task

- Given a set of transactions  $T$ , the goal of association rule mining is to find all rules having
  - support  $\geq$  *minsup* threshold
  - confidence  $\geq$  *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, Vegetables, Fruits, Eggs
3	Milk, Vegetables, Fruits, Coke
4	Bread, Milk, Vegetables, Fruits
5	Bread, Milk, Vegetables, Coke

## Example of Rules:

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

## Observations:

- All the above rules are binary partitions of the same itemset:  
 $\{\text{Milk, Vegetables, Fruits}\}$
- 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:
  1. **Frequent Itemset Generation**
    - Generate all itemsets 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
- Frequent itemset generation is still computationally expensive







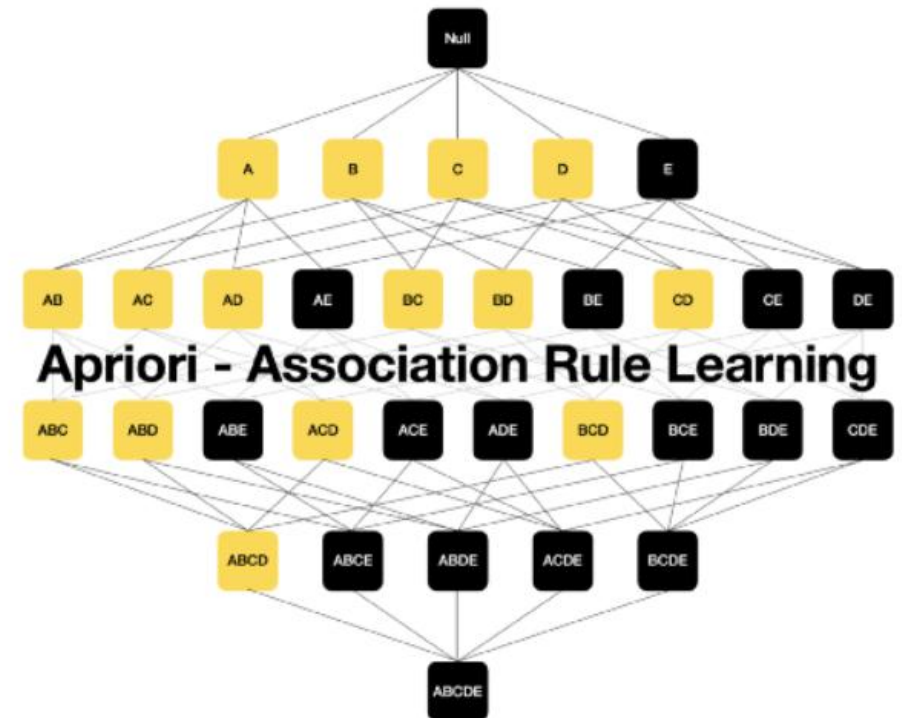
# Apriori





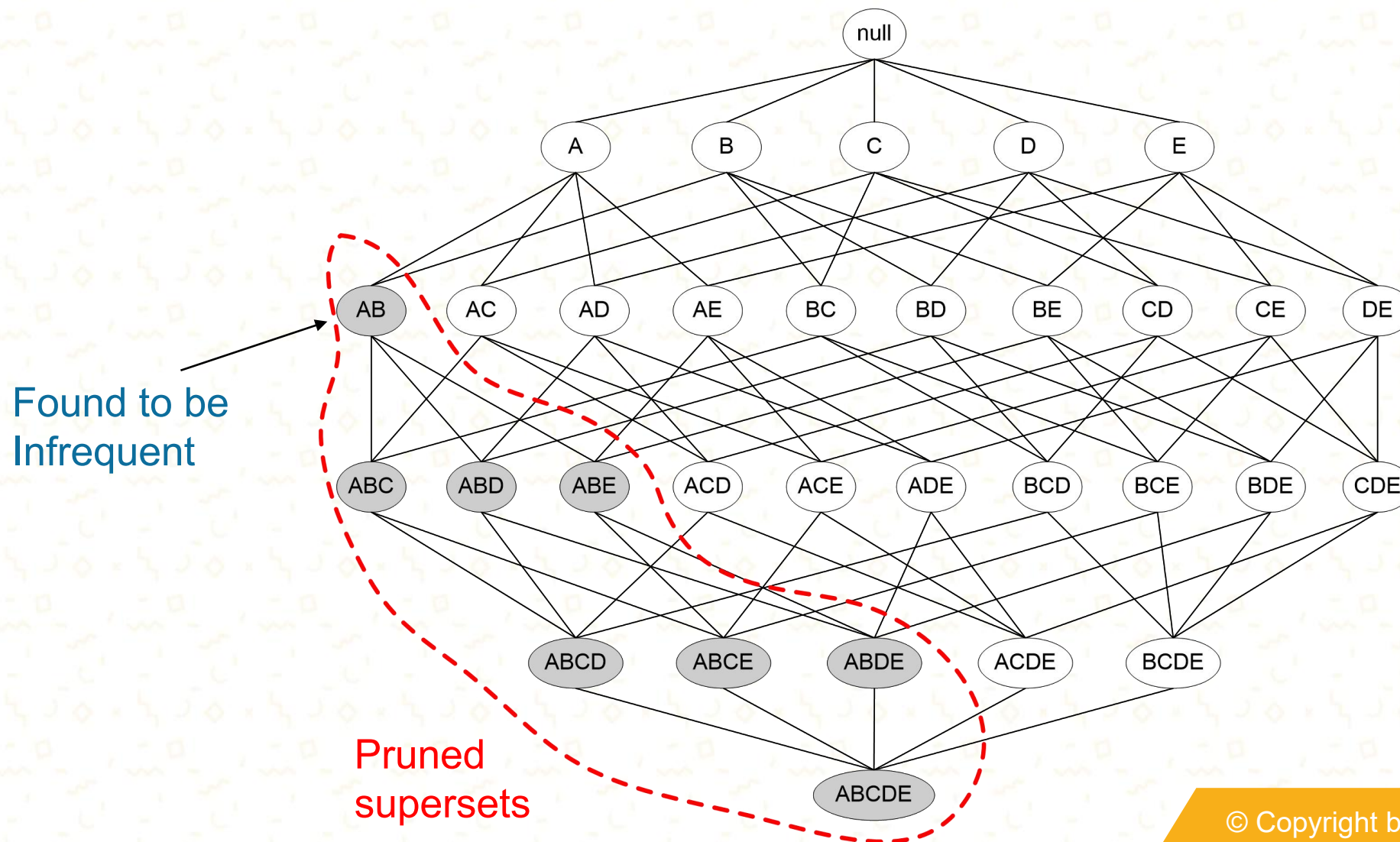
# Reducing Number of Candidates

- Apriori uses a "bottom up" approach, where frequent subsets are extended one item at a time (a step known as candidate generation), and groups of candidates are tested against the data.
- Apriori principle:
  - If an itemset is frequent, then all of its subsets must also be frequent
  - Support of an itemset never exceeds the support of its subsets





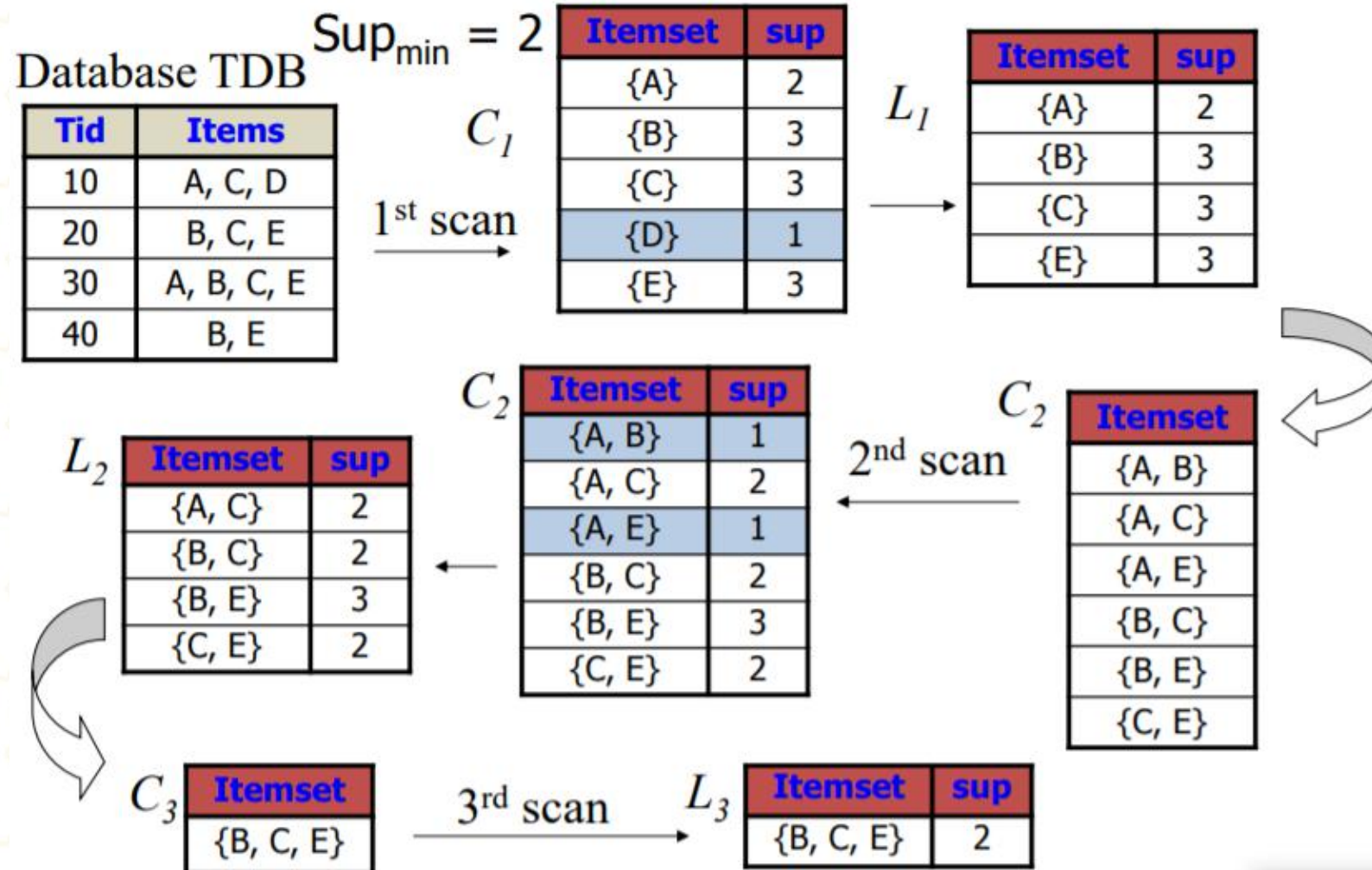
# Illustrating Apriori Principle







# Illustrating Apriori Principle







# 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 candidates that are infrequent, leaving only those that are frequent





**Lets Coding!**



Thank  
YOU