



# ASSOCIATION RULE MINING

BEIYU LIN



# 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

## Market 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

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

## REVIEW: SET AND SUBSET

- $\{a, b, c, d\} \Leftrightarrow$  a set (there are one or more than one items)
- Subset  $\Leftrightarrow$  possible combinations of the items in a set
  - Possible sets:
  - $\{a\}, \{b\}, \{c\}, \{d\}, \{a, b\}, \{a, c\}, \{a, d\}, \{b, c\}, \{b, d\}, \{c, d\}, \{a, b, c\}, \{a, b, d\}, \{a, c, d\}, \{b, c, d\}, \{a, b, c, d\}, \{\}$
  - What is the total number of the subset:

$$2^4 = 16$$

$$2^4 - \underbrace{1} = 15$$

Not Count  $\{\}$

# ASSOCIATION RULE MINING

- **Itemset (set / subset)**

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

<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

## Support

Fraction of transactions that contain an itemset

E.g.  $s(\{\text{Milk, Bread, Diaper}\}) = \underline{2/5}$

## Frequent Itemset

An itemset whose support is greater than or equal to a

minsup threshold

||

min-support

# DEFINITION: ASSOCIATION RULE

## ● Association Rule

- An implication expression of the form  $X \rightarrow Y$ , where  $X$  and  $Y$  are itemsets

- Example:

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

## ● 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:

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

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

$\frac{2}{5} = \frac{\text{\# of itemset}}{\text{total \# transactions}}$

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

$\frac{2}{3} = \frac{\text{\# of itemset of X and Y}}{\text{\# of transactions containing X}}$

# ASSOCIATION RULE MINING TASK

TID = 1, 2, 3, 4, 5

- 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 expensive / prohibitive!**

k=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

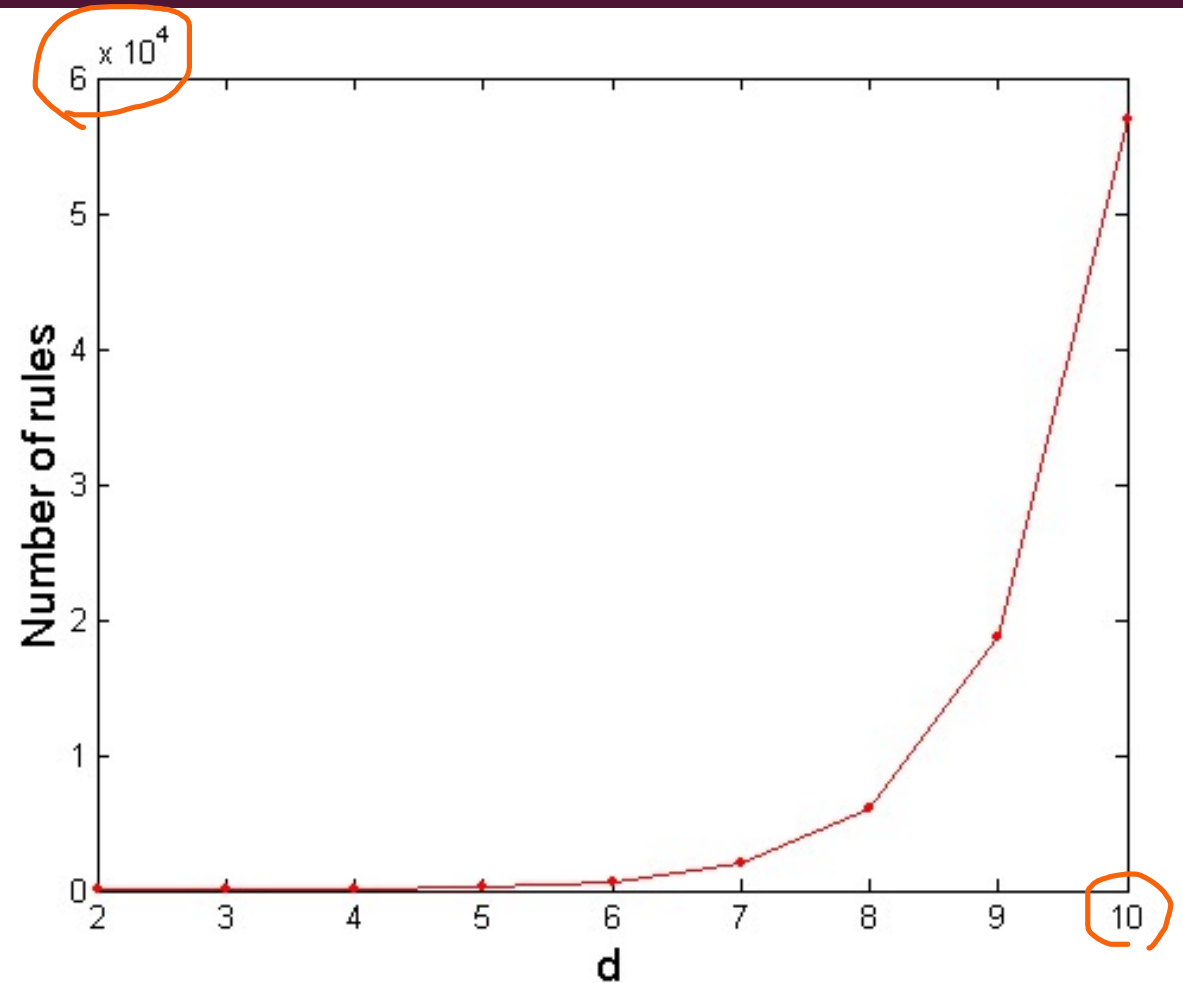
n=4 {a,b,c,d} → # of subset:  $2^n$

# 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)

$\{a, b, c\}$   
 $= \{b, c, a\}$   
 $= \{c, a, b\}$   
 $= \{c, b, a\}$

Observations:  $X = \{B\}$   
 $Y = \{M, D\}$

$$S(X \rightarrow Y) = \frac{\# \{X \cup Y\}}{\# \text{trans}} = \frac{\# \{B, M, D\} = 2}{5} = \frac{2}{5} = 0.4$$

- 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

$$C(X \rightarrow Y) = \frac{\# \{X \cup Y\}}{\# X} = \frac{2}{3} = 0.67$$



# MINING ASSOCIATION RULES

- **Two-step approach:**

- 1. **Frequent Itemset Generation**

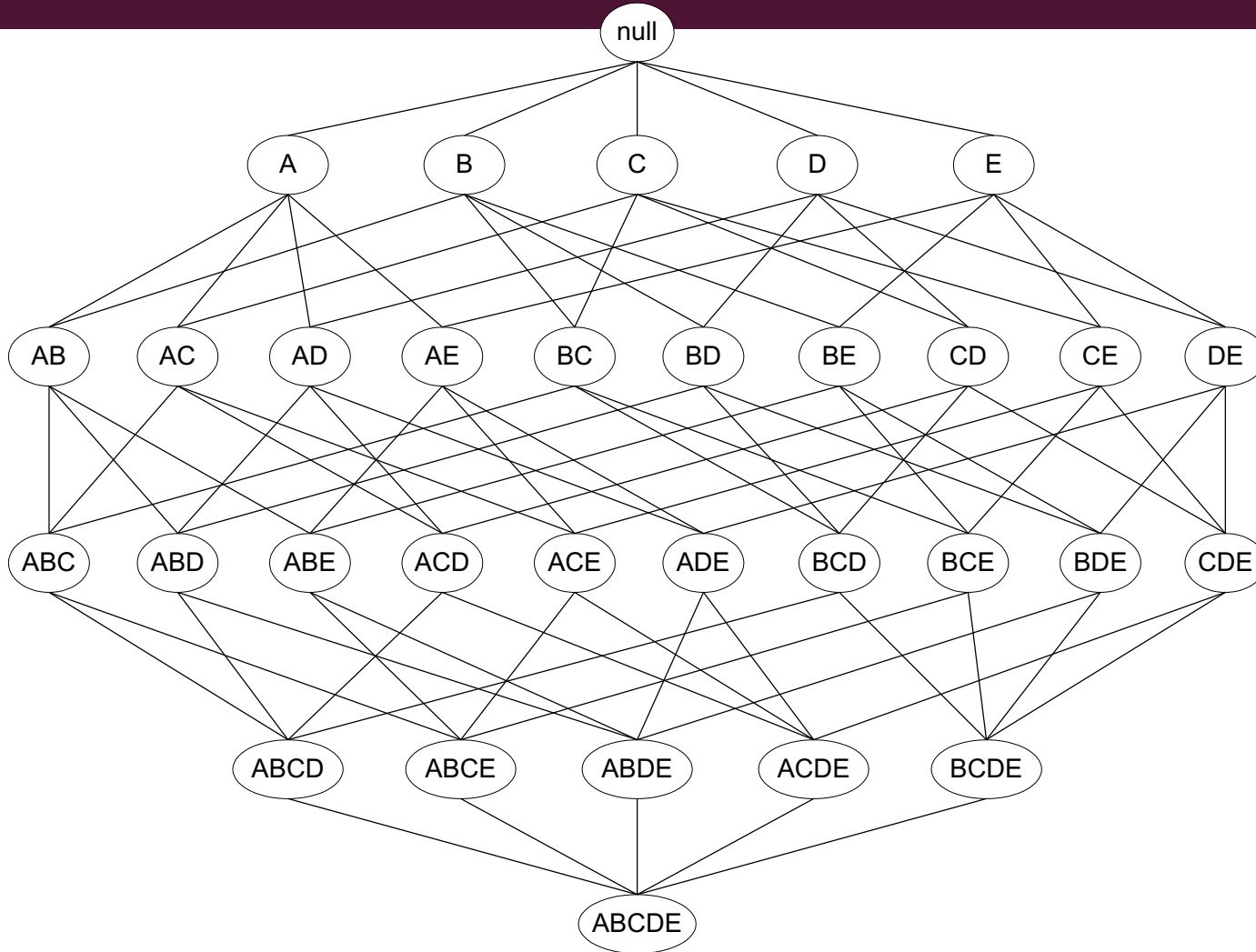
- Generate all itemsets whose support  $\geq$  minsup

- 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**

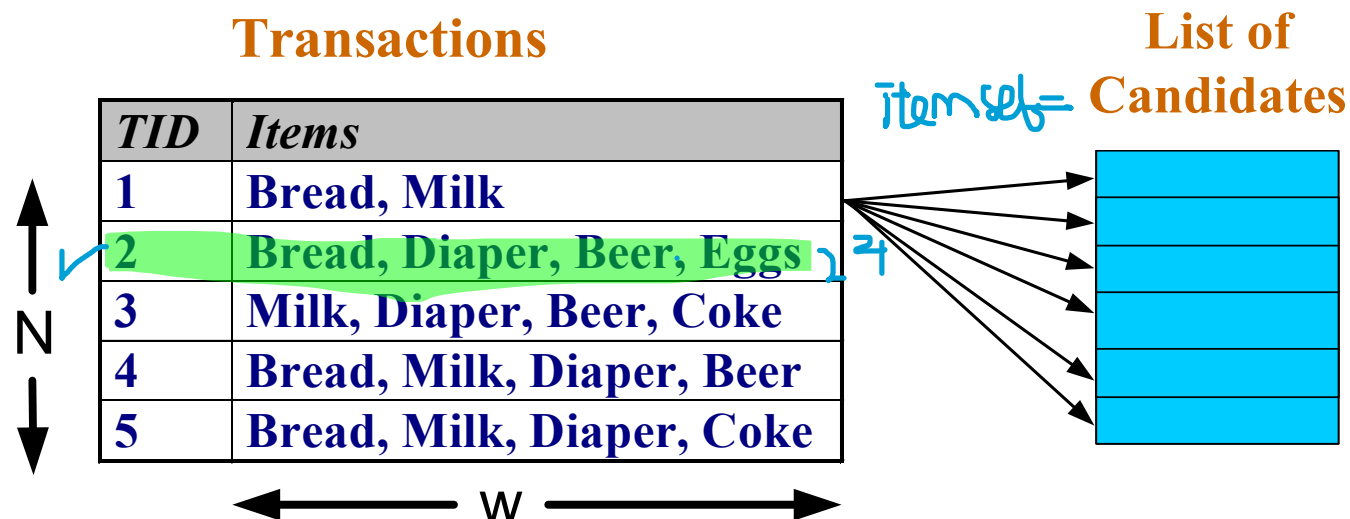
# FREQUENT ITEMSET GENERATION



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

# FREQUENT ITEMSET GENERATION

- Brute-force approach: expensive
  - 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  $\sim O(NMw)$  Expensive since  $M = 2^d$  !!!

$$N = 5$$

$$W = 4 \text{ (# items)}$$

$$M = 2^4$$

$$O(N \cdot W \cdot M):$$

computation

$$O(N \cdot W \cdot 2^d)$$

$O(N \cdot W \cdot 2^d)$

# 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 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

