



# 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

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

# 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

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

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

## ■ Support count ( $\sigma$ )

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

$$\frac{\# \{ \}}{5}$$

$$= \text{minsup} = 40\% \quad \sigma_n = 2$$

$$s \geq \text{minsup} \\ 6 \geq 6m$$

## Support

Fraction of transactions that contain an itemset

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

## Frequent Itemset

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

# DEFINITION: ASSOCIATION RULE

## ● Association Rule

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

- Example:

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

X Y

## ● Rule Evaluation Metrics

- Support (s)  $= \frac{\#X \cup Y}{5} = \frac{2}{5}$ 
  - ◆ Fraction of transactions that contain both  $X$  and  $Y$

- Confidence (c)  $= \frac{\#X \cup Y}{\#X} = \frac{2}{3}$ 
  - ◆ 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:

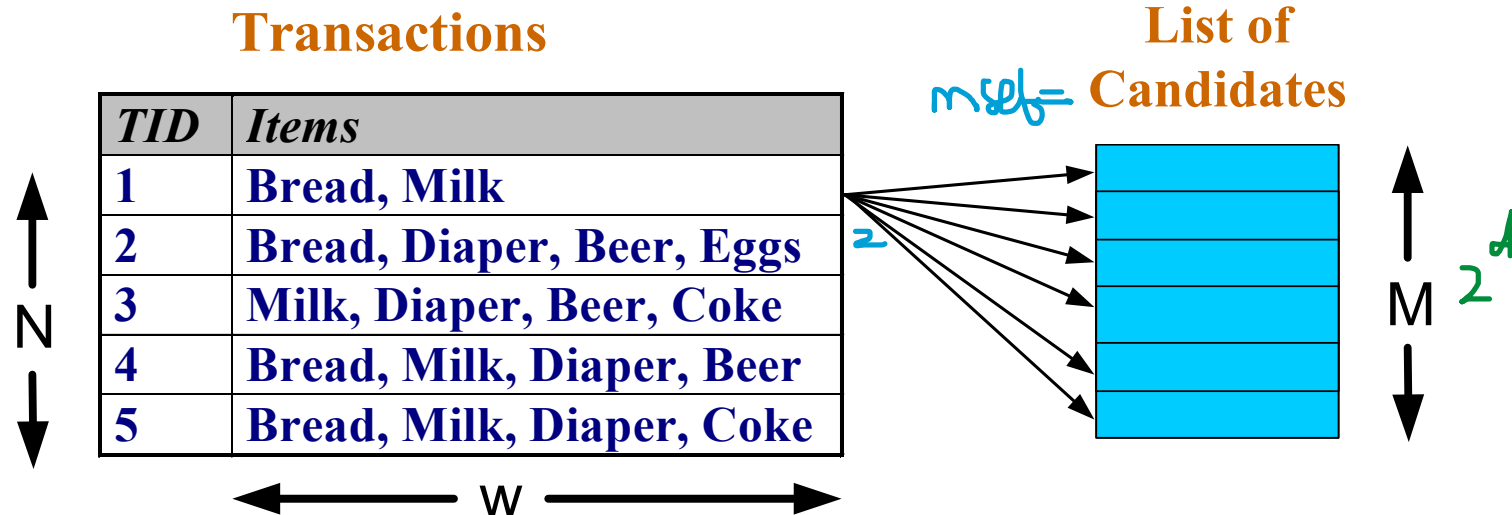
$\{\text{Milk, Diaper}\} \Rightarrow \{\text{Beer}\}$

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4 \quad \text{= \# of itemset / total \# transactions}$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67 \quad \text{= \# of itemset of X and Y / \# of transactions containing X}$$

# FREQUENT ITEMSET GENERATION

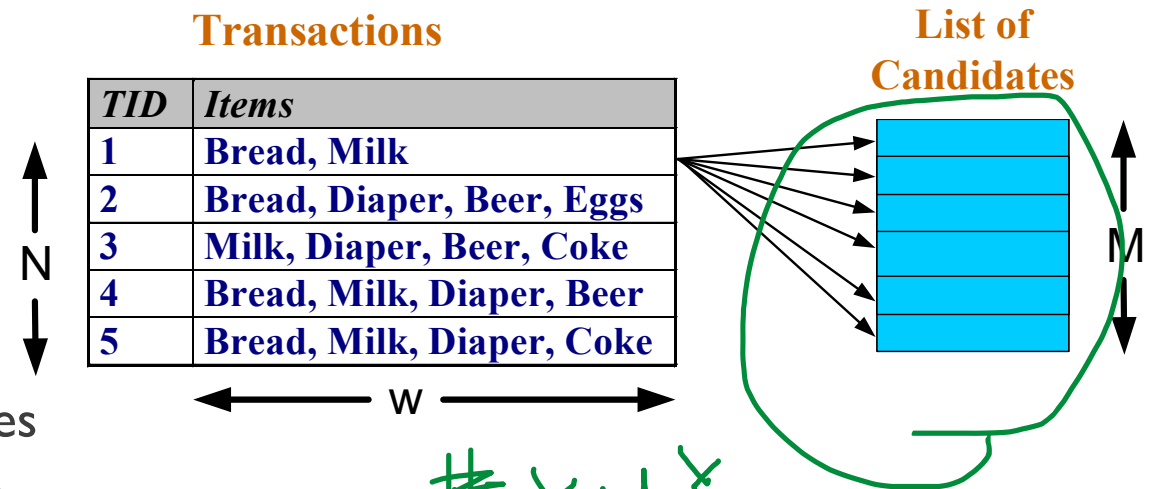
- Brute-force approach:
  - Each itemset in the lattice is a **candidate** frequent itemset
  - Count the support of each candidate by scanning the database



- Match each transaction against every candidate
- Complexity  $\sim O(NMw) \Rightarrow$  **Expensive since  $M = 2^d$  !!!**

# FREQUENT ITEMSET GENERATION STRATEGIES

- OL N M M
  - 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



$$C = \frac{\#XUY}{\#X}$$

$\updownarrow$

F, 6

# REDUCING NUMBER OF CANDIDATES

$A \rightarrow B \Leftarrow \text{NOT } B \Rightarrow \text{NOT } A$

## Apriori principle:

- If an itemset is frequent, then all of its subsets must also be frequent  
If the subset of an itemset is not frequent, then the itemset is not 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

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

Subset: {milk, bread}

$$s(\{\text{milk, bread}\}) = \# \{\text{milk, bread}\} / 5 = \frac{3}{5}$$

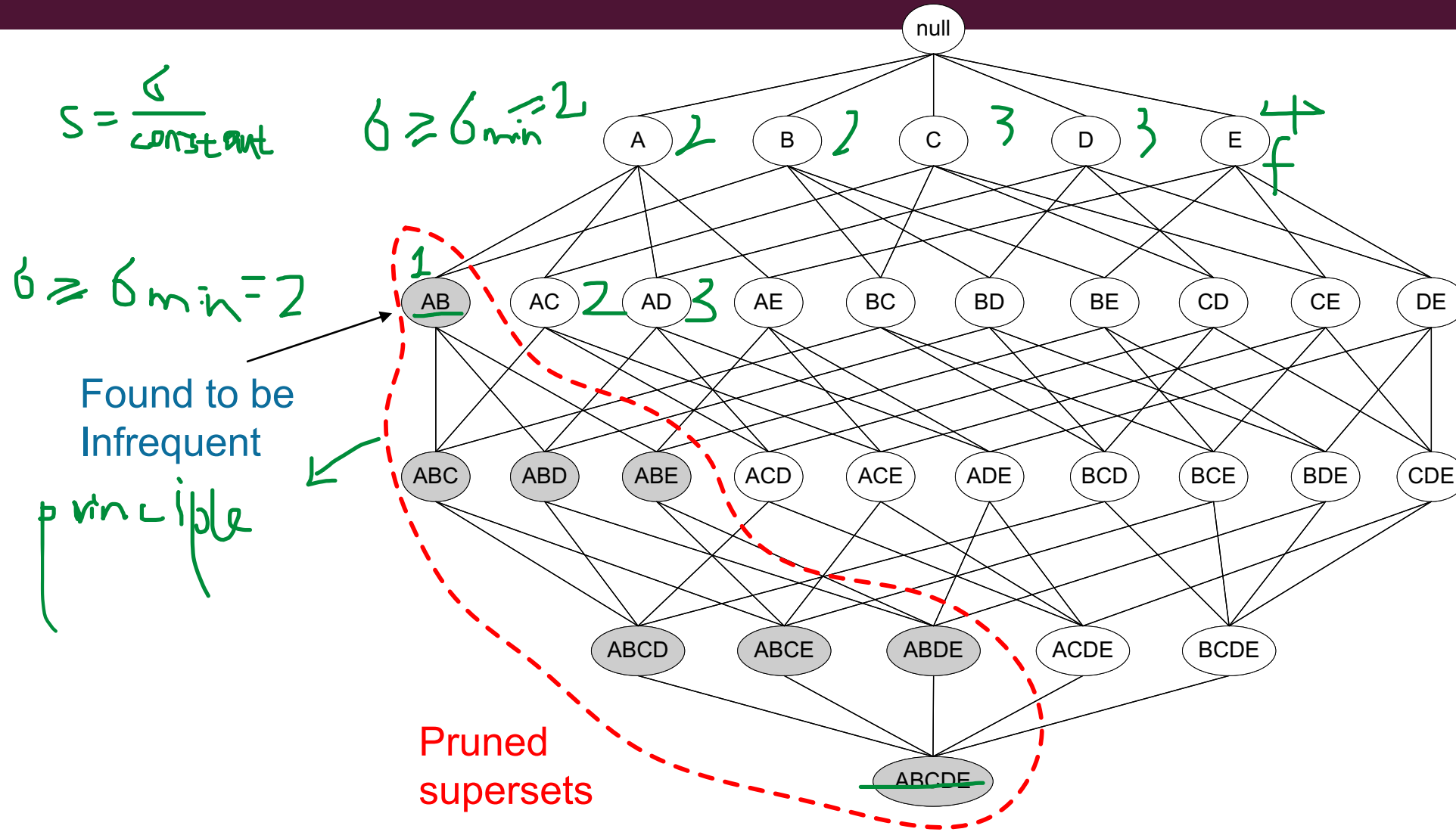
$$\frac{3}{5} > \frac{2}{5}$$

$$s(X) \geq s(Y) \\ X \subseteq Y$$

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

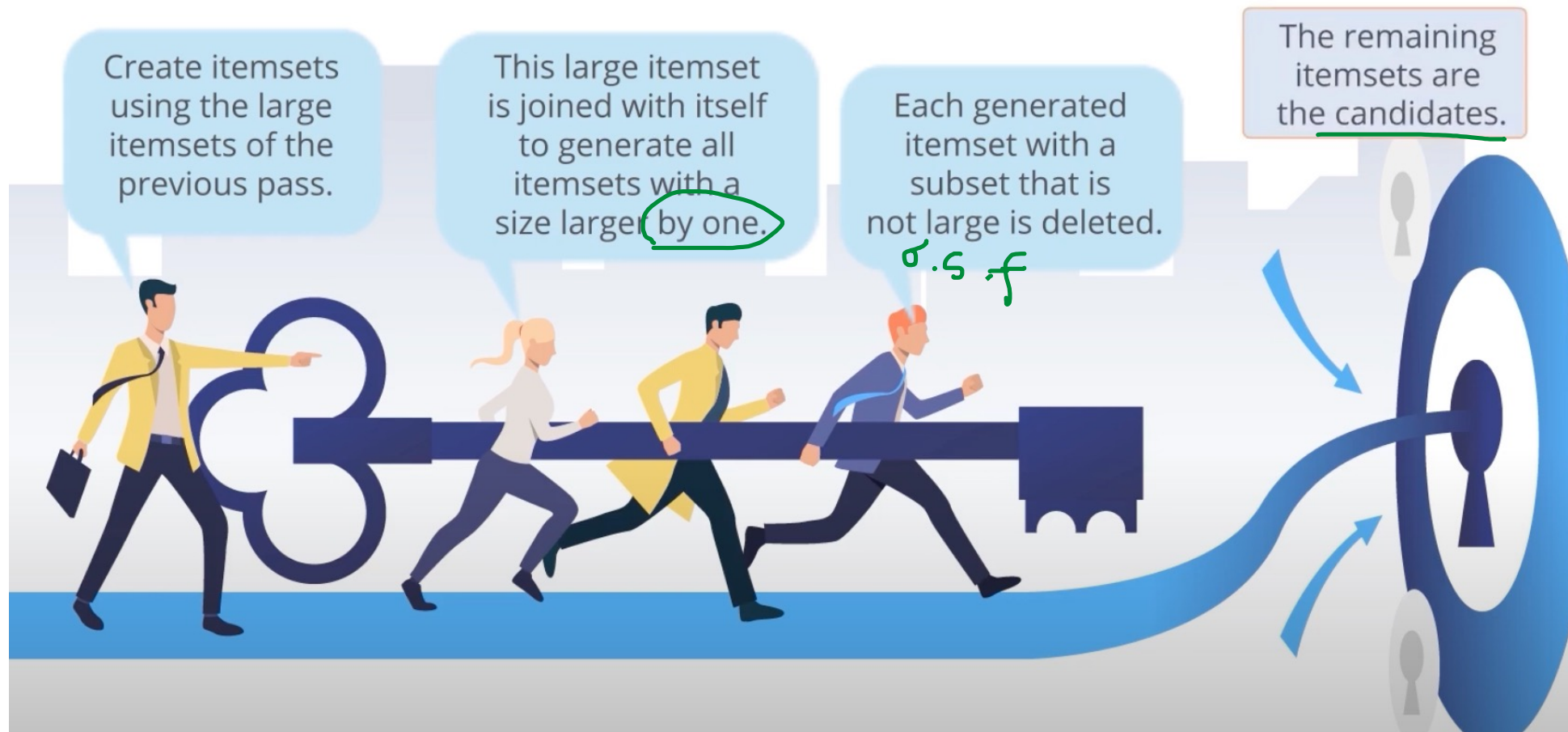
- Support of an itemset never exceeds the support of its subsets
- This is known as the **anti-monotone** property of support

# ILLUSTRATING APRIORI PRINCIPLE





# APRIORI PRINCIPLE



# APRIORI PRINCIPLE



Uses frequent itemsets to generate association rules



Support value of frequent itemsets is greater than the threshold value

$$\text{count}(X \Rightarrow Y) = \frac{\sigma}{n}$$

The algorithm reduces the number of candidates being considered by only exploring the itemsets whose support count is greater than the minimum support count.

# A PRIORI PRINCIPLE EXAMPLE

TID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5
500	1 3 5

IF

Itemset	Support
{1}	3 $\geq 2$ ✓
{2}	3 $\geq 2$ ✓
{3}	4
{4}	1
{5}	4

$\delta \geq 6_{min} = 2 \quad \text{and} \quad 5_{min} = 5 = \%$

Itemset	Support
{1}	3
{2}	3
{3}	4
{5}	4

TID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5
500	1 3 5



# APRIORI PRINCIPLE EXAMPLE

The length of the itemset is extended with 1 ( $k = k+1$ ).

F11

Itemset	Support
{1}	3
{2}	3
{3}	4
{5}	4

CI2

Itemset	Support
<del>{1,2}</del>	<del>1</del>
{1,3}	3
{1,5}	2
{2,3}	1
{2,5}	3
{3,5}	3

FI2

Itemset	Support
{1,3}	3
{1,5}	2
{2,3}	2
{2,5}	3
{3,5}	3

$$6 \geq 6_{\min} = 2$$

# APRIORI PRINCIPLE EXAMPLE

The length of the itemset is extended with 1 ( $k = k+1$ ).

F12

Itemset	Support
{1,3}	3
{1,5}	2
{2,3}	2
{2,5}	3
{3,5}	3

{1,2,3,5}

C13

IF

Itemset	Support
{1,2,3}	1
{1,2,5}	1
{1,3,5}	2
{2,3,5}	2



Itemset	Support
{1,3,5}	2
{2,3,5}	2

$$5 \geq 6_{min} = 2$$

# APRIORI PRINCIPLE EXAMPLE

Divide your itemset to check if there are any other subsets whose support you haven't calculated yet.

C13	
Itemset	Support
<u>{1,2,3}</u>	1
{1,2,5}	1
F {1,3,5}	2
F {2,3,5}	2

Itemset	In F12?
IF {1,2}, {1,3}, {2,3}, {1}, {2}, {3}	IF
IF {1,2}, {1,5}, {2,5}, {1}, {2}, {5}	IF
{1,3}, {1,5}, {3,5}, {1}, {3}, {5}	F



Itemset	Support
{1,3}	3
{1,5}	2
{2,3}	2
{2,5}	3
{3,5}	3

# APRIORI PRINCIPLE EXAMPLE

TID	Items
100	1 3 4
200	2 3 5
300	<u>1 2 3 5</u>
400	2 5
500	1 3 5

$k=3$

Itemset	Support
<u>{1,3,5}</u>	2
<u>{2,3,5}</u>	2

now

1 2 3 4 5

prune



# APRIORI PRINCIPLE EXAMPLE

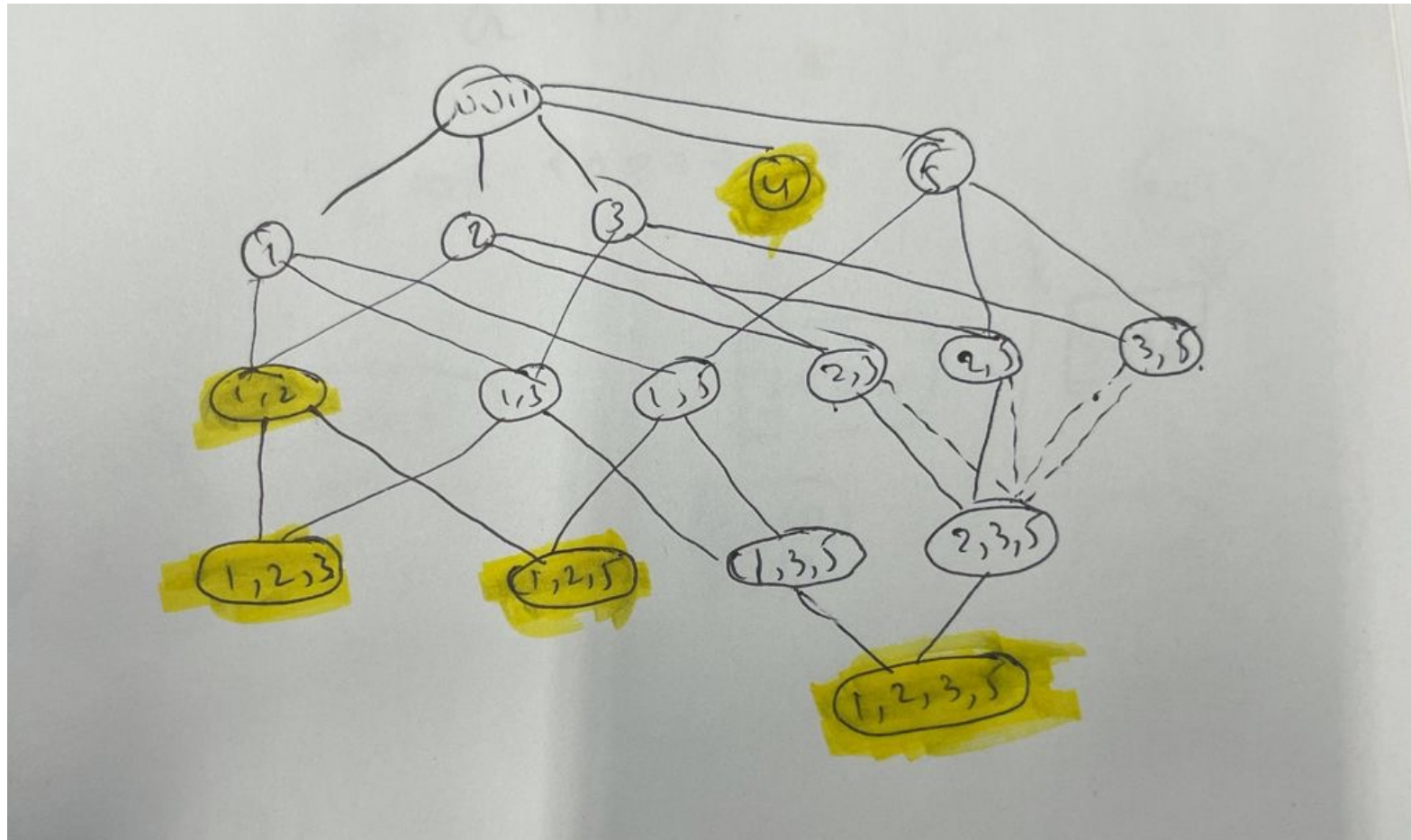
The length of the itemset is extended with 1 ( $k = k+1$ ).

Itemset	Support
{1,3,5}	2
{2,3,5}	2



C4

Itemset	Support
{1,2,3,5}	1



# APRIORI ALGORITHM

- $F_k$ : frequent k-itemsets
- $L_k$ : candidate k-itemsets
- Algorithm
  - Let  $k=1$
  - Generate  $F_1 = \{\text{frequent 1-itemsets}\}$
  - Repeat until  $F_k$  is empty
    - **Candidate Generation:** Generate  $L_{k+1}$  from  $F_k$
    - **Candidate Pruning:** Prune candidate itemsets in  $L_{k+1}$  containing subsets of length  $k$  that are infrequent
    - **Support Counting:** Count the support of each candidate in  $L_{k+1}$  by scanning the DB
    - **Candidate Elimination:** Eliminate candidates in  $L_{k+1}$  that are infrequent, leaving only those that are frequent  $\Rightarrow F_{k+1}$