



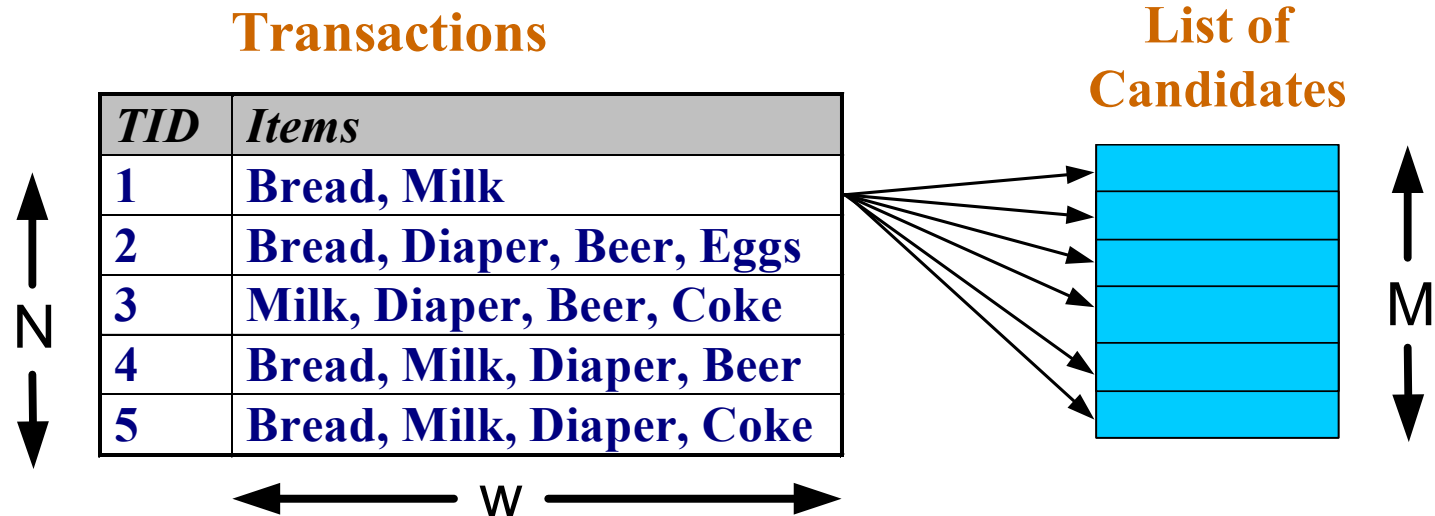
# ASSOCIATION RULE MINING

BEIYU LIN



# FREQUENT ITEMSET GENERATION

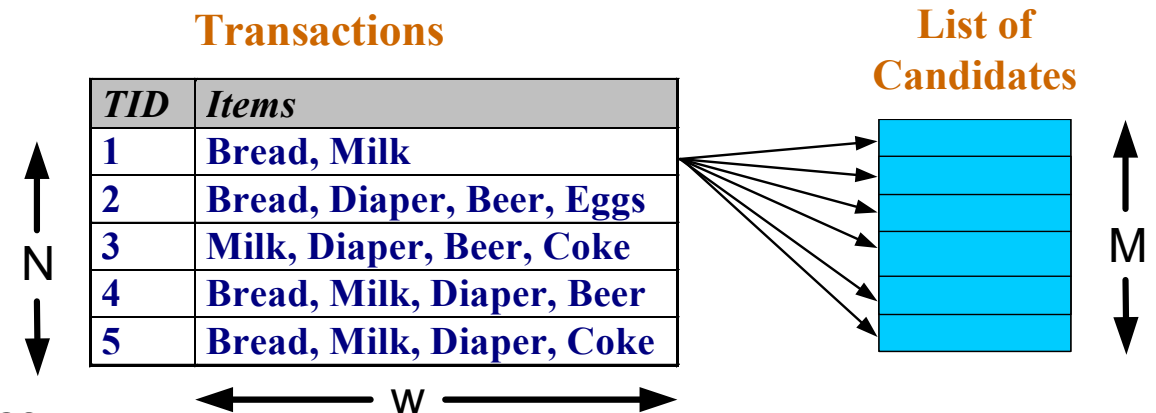
- Brute-force approach:
  - Each itemset in the lattice is a **candidate** frequent itemset
  - Count the support of each candidate by scanning the database



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

# FREQUENT ITEMSET GENERATION STRATEGIES

- 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



Given a transaction {B, M, D, C}, find all possible subset with size 3 from this transaction.

# REDUCING NUMBER OF CANDIDATES

TID	Items
1	<del>B</del> read, <del>M</del> ilk
2	<del>B</del> read, <del>D</del> iaper, <del>B</del> eer, <del>E</del> ggs
3	Milk, Diaper, Beer, Coke
4	<del>B</del> read, Milk, Diaper, Beer
5	<del>B</del> read, Milk, Diaper, Coke

## 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

$s(\{\text{Milk, Bread, Diaper}\}) = 2/5 = \#$   
itemsets / total # of transaction

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

**Support count:** # of the itemsets that show in the transaction

$$\begin{matrix} \{A\} & \{AB\} & X & \subseteq & Y \\ X & Y & & & \end{matrix}$$

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

$$P(AB|A) = \frac{P(AB)}{P(A)}$$

$$s(X) = \frac{\sigma(X)}{\#} = \frac{\sigma(\{A\})}{\#}$$

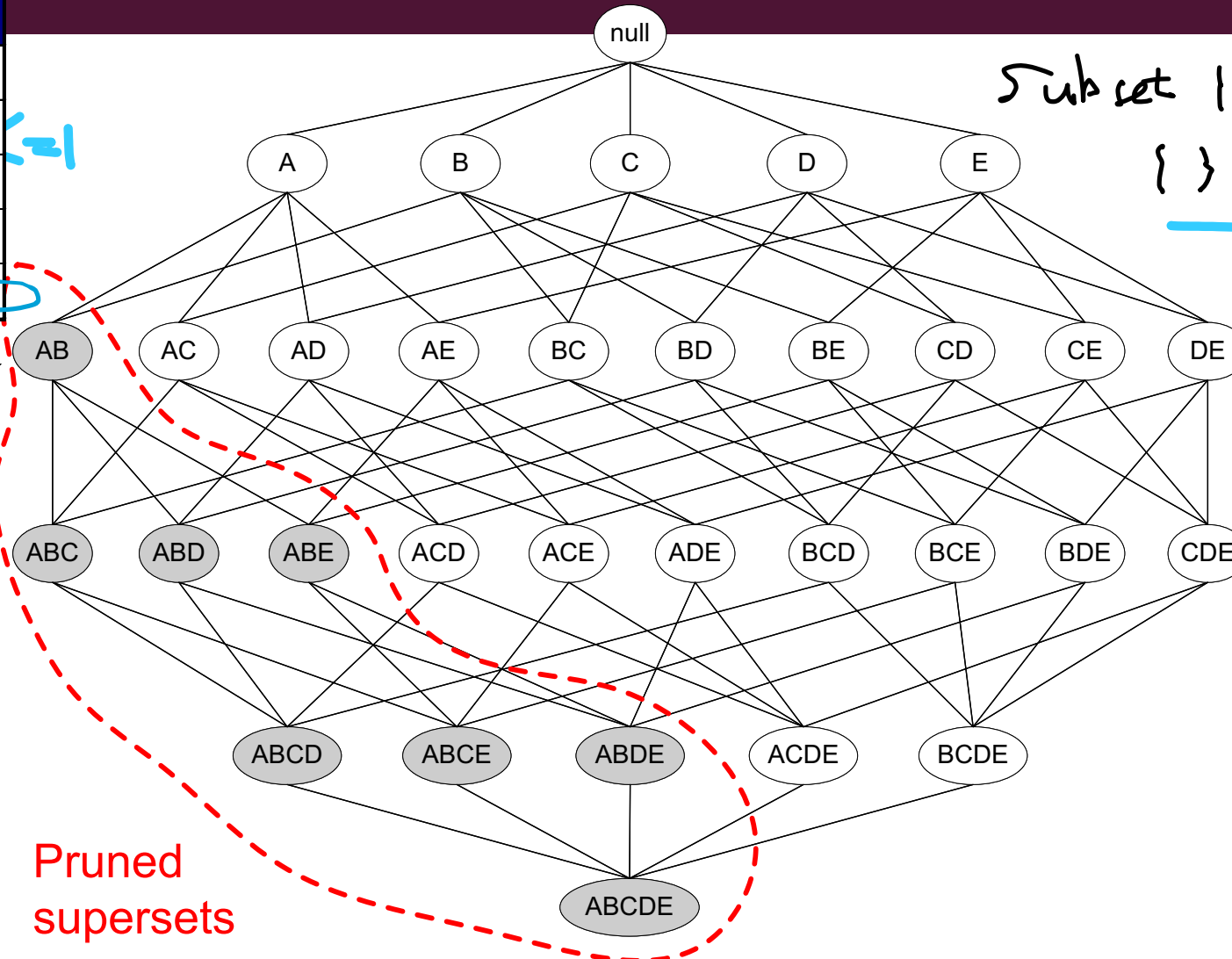
$$s(Y) = \frac{\sigma(Y)}{\#} = \frac{\sigma(\{AB\})}{\#}$$

## Confidence

$$C(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$$

# ILLUSTRATING APRIORI PRINCIPLE

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

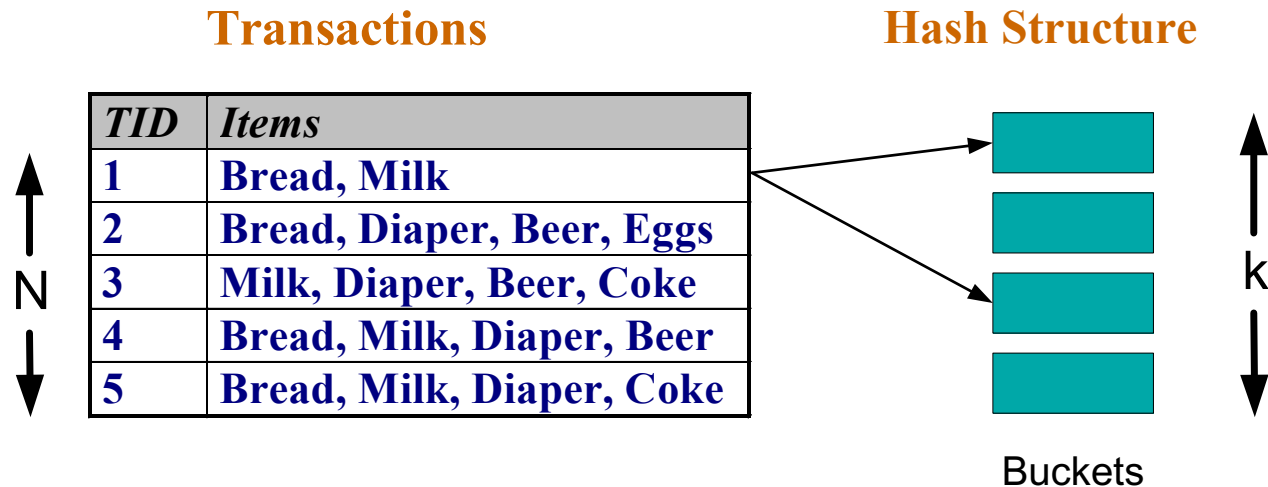


Subset  $\{ABC\}$   $2^3$   
 $\{ \}, \{A\}, \dots$

supsets  
 $\{ABC\}$   
 $\{ABCD\}$   
 $\{ABCDE\}$

# SUPPORT COUNTING OF CANDIDATE ITEMSETS

- To reduce number of comparisons, store the candidate itemsets in a hash structure / hash function
- Instead of matching each transaction against every candidate, match it against candidates contained in the hashed buckets

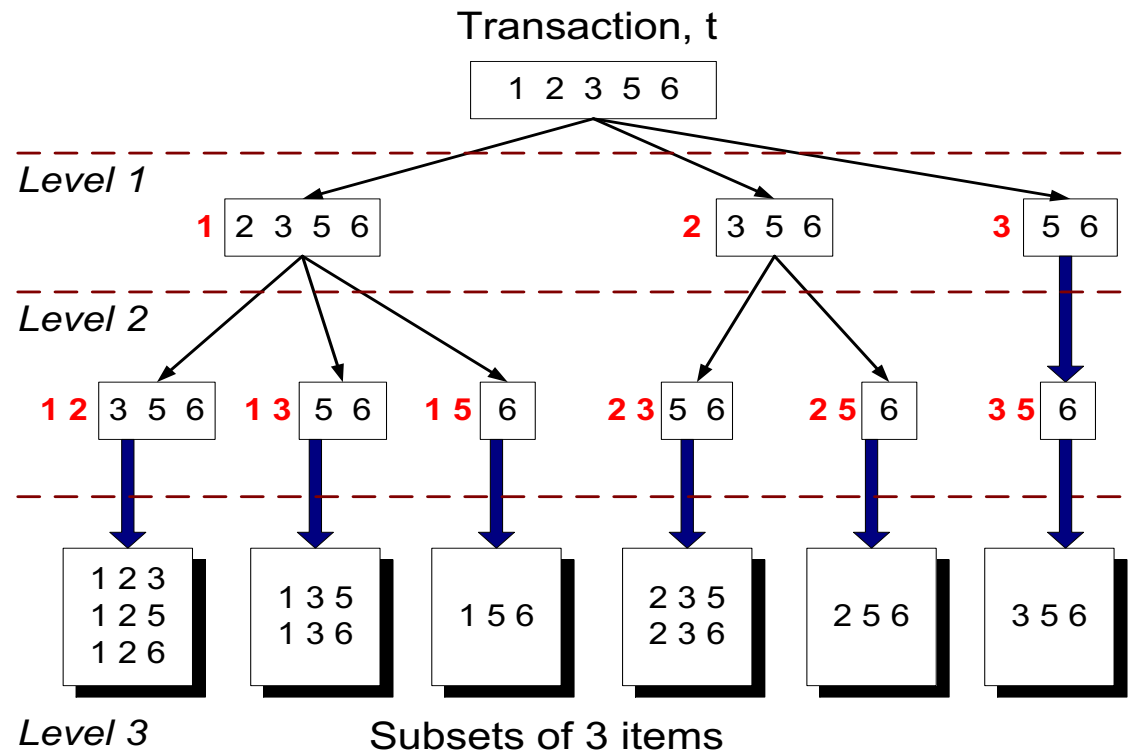


# SUPPORT COUNTING: AN EXAMPLE

Suppose you have 15 candidate itemsets of length 3:

{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4}, {5 6 7}, {3 4 5}, {3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}

How many of these itemsets are supported by transaction (1,2,3,5,6)?



# SUPPORT COUNTING:AN EXAMPLE

Suppose you have 15 candidate itemsets of length 3:

{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4}, {5 6 7}, {3 4 5}, {3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}

## How to find all the subsets with k items from an itemset?

Given an itemset  $\{1,2,3,5,6\} \Leftrightarrow$  one transaction, we want to reduce NM

How many of these itemsets are supported by transaction (1,2,3,5,6)?

Q) Given Transaction (1, 2, 3, 5, 6)

We need list all the possible itemsets with k = 3 from this transaction.

2

