

Association Rule Mining

- **Itemset (set / subset)** . {Milk, Bread, Diaper}
 - A collection of one or more items
 - Example: {Milk, Bread, Diaper}
 - k-itemset
 - An itemset that contains k items
- **Support count (σ)**
 - Frequency of occurrence of an itemset
 - E.g. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

<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 $S = \frac{s}{\#} = \frac{2}{5}$

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

$there = minsup \quad minsup \text{ count} = 2$
 $\sigma(\text{itemset}) \geq 2$

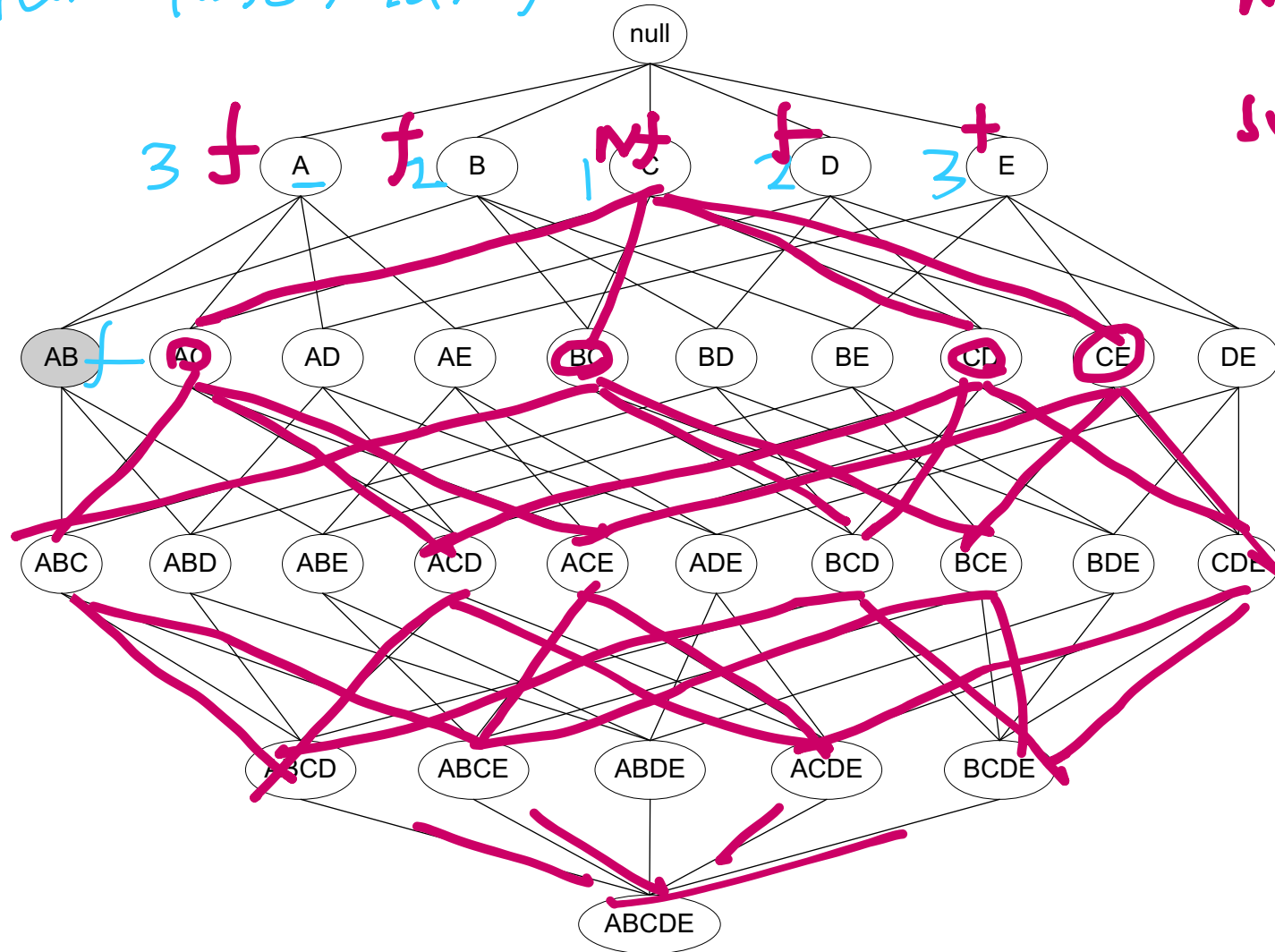
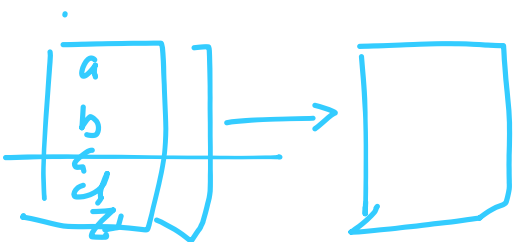
ILLUSTRATING APRIORI PRINCIPLE

 \Rightarrow unique items $\{a, b, c, d, e\}$

min G=2

$\{A\}, \{B\}, \{D\}, \{E\}$
 $\{AB\},$

min sup count $\rightarrow 2$



nt:

supers M

$\{a, b\}$
 \uparrow
 $\{a\} \{b\}$

Rule Generation

$$AR: X \rightarrow Y$$

$k=1$
 $k=2$
 $k=3$

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

- Given a frequent itemset L , find all non-empty subsets $f \subset L$ such that $f \rightarrow L - f$ satisfies the minimum confidence requirement

- If $\{A, B, C, D\}$ is a frequent itemset, candidate rules:

$ABC \rightarrow D,$	$ABD \rightarrow C,$	$ACD \rightarrow B,$	$BCD \rightarrow A,$	} $k=3$
$A \rightarrow BCD,$	$B \rightarrow ACD,$	$C \rightarrow ABD,$	$D \rightarrow ABC$	
$AB \rightarrow CD,$	$AC \rightarrow BD,$	$AD \rightarrow BC,$	$BC \rightarrow AD,$	} $k=2$
$BD \rightarrow AC,$	$CD \rightarrow AB,$			

- If $|L| = k$, then there are $2^k - 2$ candidate association rules (ignoring $L \rightarrow \emptyset$ and $\emptyset \rightarrow L$)

Rule Generation

$$c(ABC \rightarrow D) = \frac{\sigma(ABCD)}{\sigma(ABC)} \neq c(AB \rightarrow D) = \frac{\sigma(ABD)}{\sigma(AB)}$$

- In general, confidence does not have an anti-monotone property
 $c(ABC \rightarrow D)$ can be larger or smaller than $c(AB \rightarrow D)$

- But confidence of rules generated from the same itemset has an anti-monotone property

- E.g., Suppose $\{A, B, C, D\}$ is a frequent 4-itemset:

$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD) \quad \Leftarrow \frac{\sigma(ABCD)}{\sigma(A)}$$

- Confidence is anti-monotone w.r.t. number of items on the RHS of the rule

Rule Generation for Apriori Algorithm

$X \rightarrow Y$
 $k=3$

Lattice of rules

