

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, Butter, Coffee, Eggs
3	Milk, Butter, Coffee, Coke
4	Bread, Milk, Butter, Coffee
5	Bread, Milk, Butter, Coke

Example of Association Rules

$\{\text{Butter}\} \rightarrow \{\text{Coffee}\},$
 $\{\text{Milk, Bread}\} \rightarrow \{\text{Eggs, Coke}\},$
 $\{\text{Coffee, Bread}\} \rightarrow \{\text{Milk}\},$

Implication means co-occurrence,
not causality!

Definition: Frequent Itemset

- **Itemset**
 - A collection of one or more items
 - Example: {Milk, Bread, Butter}
 - k-itemset
 - An itemset that contains k items
- **Support count (σ)**
 - Frequency of occurrence of an itemset
 - E.g. $\sigma(\{\text{Milk, Bread, Butter}\}) = 2$
- **Support**
 - Fraction of transactions that contain an itemset
 - E.g. $s(\{\text{Milk, Bread, Butter}\}) = 2/5$
- **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, Butter, Coffee, Eggs
3	Milk, Butter, Coffee, Coke
4	Bread, Milk, Butter, Coffee
5	Bread, Milk, Butter, Coke

Definition: Association Rule

- Association Rule

- An implication expression of the form $X \Rightarrow Y$, where X and Y are itemsets
- Example:
 $\{\text{Milk, Butter}\} \rightarrow \{\text{Coffee}\}$

- 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

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

Example:

$$\{\text{Milk, Butter}\} \Rightarrow \text{Coffee}$$

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

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

Association Rule Mining Task

- Given a set of transactions T , the goal of association rule mining is to find all rules having
 - support $\geq \textit{minsup}$ threshold
 - confidence $\geq \textit{minconf}$ 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

Example:

$\{\text{Milk, Butter}\} \Rightarrow \text{Coffee}$

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

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

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

Example of Rules:

$\{\text{Milk, Bread}\} \rightarrow \{\text{Coffee}\}$ ($s=1/5=0.2$, $c=1/3=0.33$)

$\{\text{Milk, Coffee}\} \rightarrow \{\text{Bread}\}$ ($s=1/5=0.2$, $c=1/2=0.5$)

$\{\text{Bread, Coffee}\} \rightarrow \{\text{Milk}\}$ ($s=1/5=0.2$, $c=1/2=0.5$)

$\{\text{Coffee}\} \rightarrow \{\text{Milk, Bread}\}$ ($s=1/5=0.2$, $c=1/3=0.33$)

$\{\text{Bread}\} \rightarrow \{\text{Milk, Coffee}\}$ ($s=1/5=0.2$, $c=1/4=0.25$)

$\{\text{Milk}\} \rightarrow \{\text{Bread, Coffee}\}$ ($s=1/5=0.2$, $c=1/4=0.25$)

Example of Rules:

$\{\text{Milk}, \text{Bread}\} \rightarrow \{\text{Coffee}\} \text{ (s=1/5=0.2, c=1/3=0.33)}$
 $\{\text{Milk}, \text{Coffee}\} \rightarrow \{\text{Bread}\} \text{ (s=1/5=0.2, c=1/2=0.5)}$
 $\{\text{Bread}, \text{Coffee}\} \rightarrow \{\text{Milk}\} \text{ (s=1/5=0.2, c=1/2=0.5)}$
 $\{\text{Coffee}\} \rightarrow \{\text{Milk}, \text{Bread}\} \text{ (s=1/5=0.2, c=1/3=0.33)}$
 $\{\text{Bread}\} \rightarrow \{\text{Milk}, \text{Coffee}\} \text{ (s=1/5=0.2, c=1/4=0.25)}$
 $\{\text{Milk}\} \rightarrow \{\text{Bread}, \text{Coffee}\} \text{ (s=1/5=0.2, c=1/4=0.25)}$

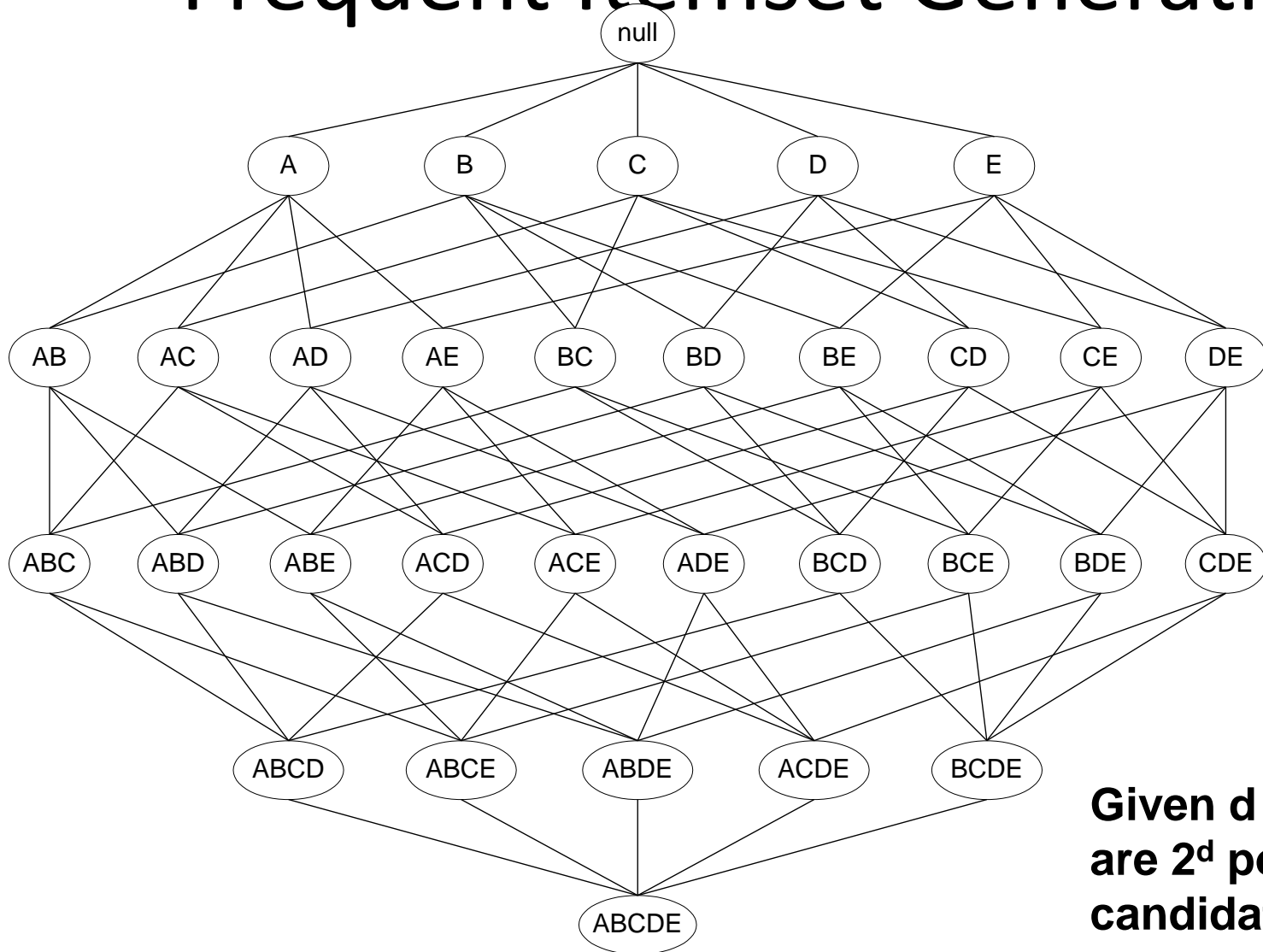
Observations:

- All the above rules are binary partitions of the same itemset:
 $\{\text{Milk}, \text{Bread}, \text{Coffee}\}$
- 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

Frequent Itemset Generation



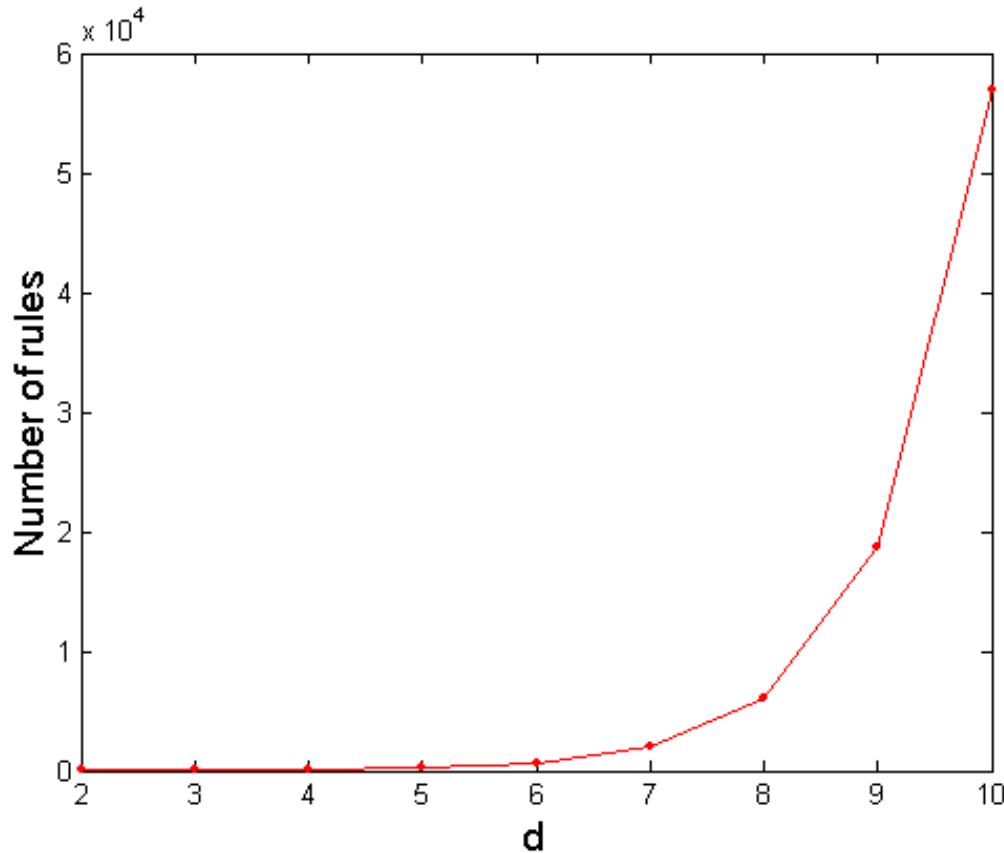
**Given d items, there
are 2^d possible
candidate itemsets**

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
 - N The number of transactions
 - M the list of candidates
 - W – The number of items (the width of the transaction)
 - Match each transaction against every candidate
 - Complexity $\sim O(NMw) \Rightarrow$ **Expensive since $M = 2^d$!!!**

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

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 vertical-dataset 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

Reducing Number of Candidates

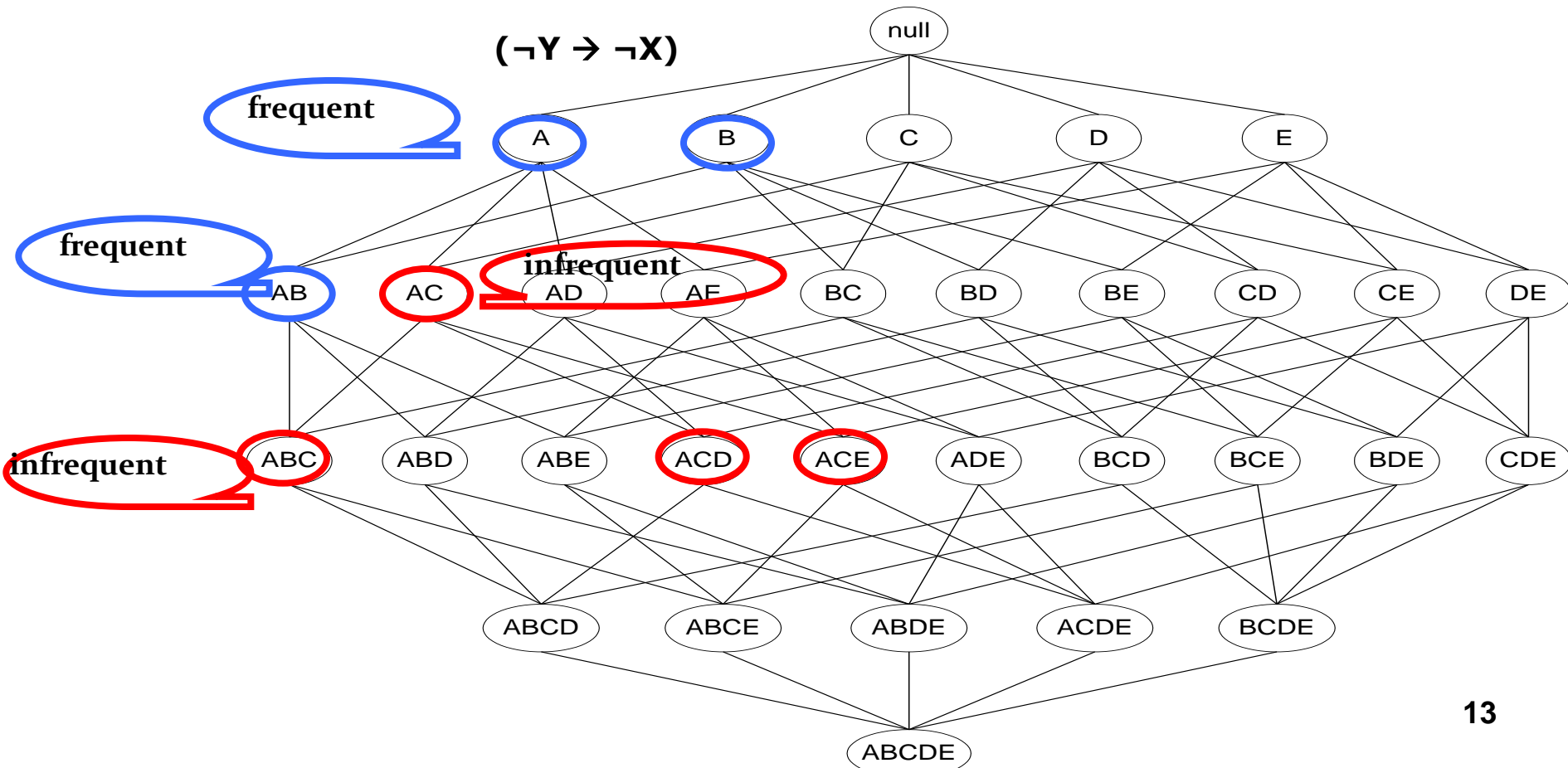
- **Apriori principle:**
 - If an itemset is frequent, then all of its subsets must also be frequent
- Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

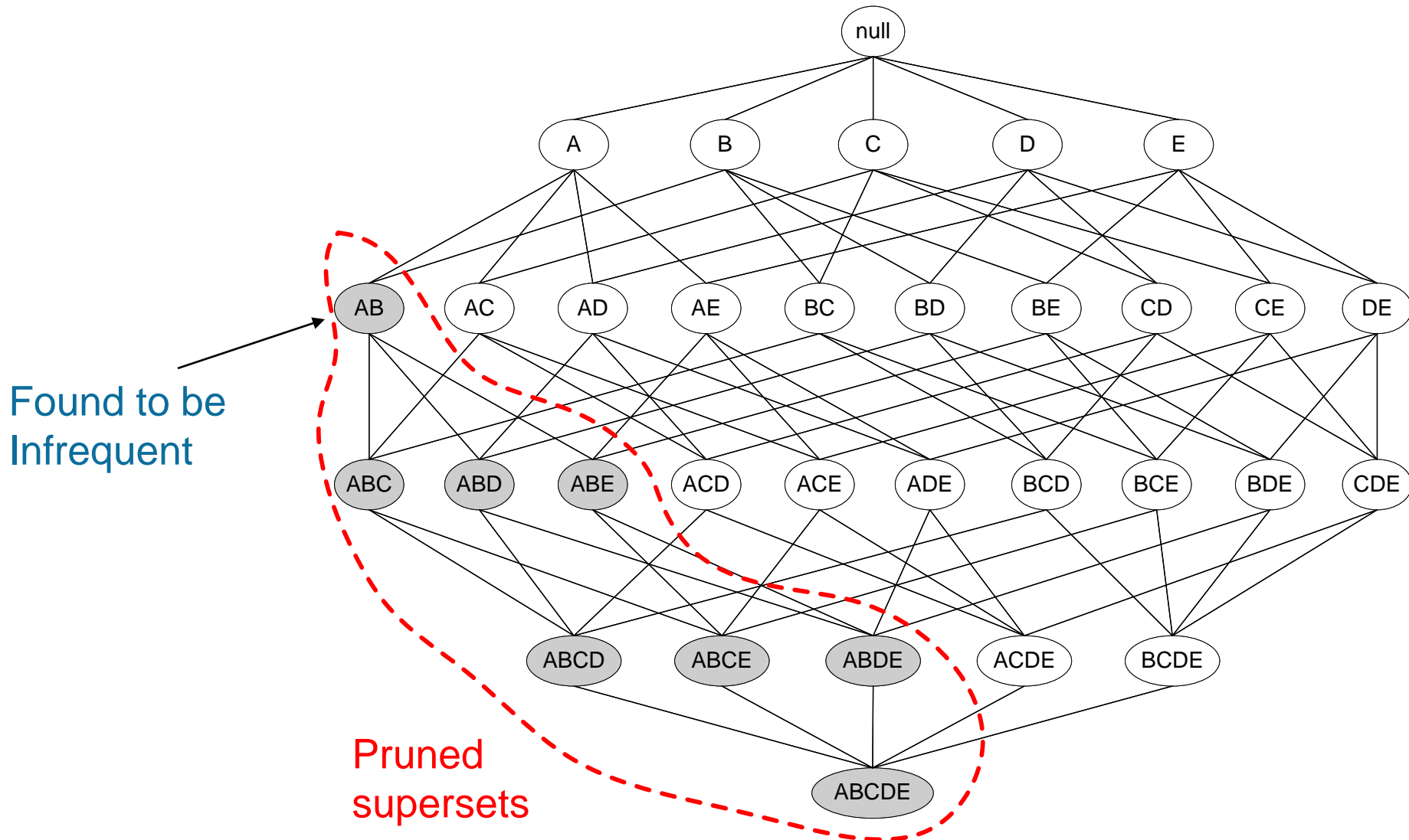
- For all X, Y in the dataset, support $(X) \geq \text{Support}(Y)$
 - Support of an itemset never exceeds the support of its subsets
 - This is known as the **anti-monotone** property of support

Apriori Principle

- If an itemset is frequent, then all of its subsets must also be frequent
 - If an itemset is infrequent, then all of its supersets must be infrequent too
- $(X \rightarrow Y)$
- $(\neg Y \rightarrow \neg X)$



Illustrating Apriori Principle



Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Coffee	3
Cookies	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Coffee}	2
{Bread,Cookies}	3
{Milk,Coffee}	2
{Milk, Cookies }	3
{Coffee, Cookies }	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3 = 41$
 With support-based pruning,
 $6 + 6 + 1 = 13$



Triplets (3-itemsets)

Itemset	Count
{Bread,Milk,Cookie s}	3



$$n C r = n! / (r!)(n-r)!$$