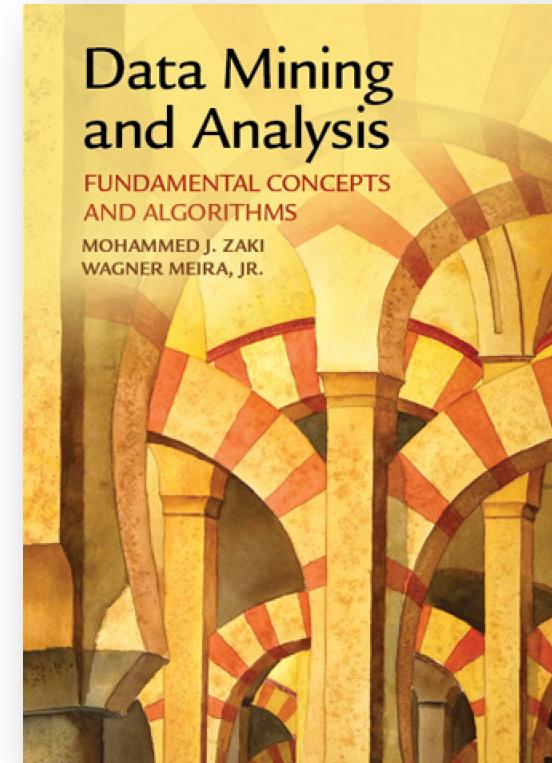


Association Rules

Data Mining and Text Mining



- “Data Mining and Analysis” by Zaki & Meira
 - Chapter 8
 - Section 9.1
 - Chapter 10 up to
Section 10.2.1 included
 - Section 12.1
- <http://www.dataminingbook.info>





Bread
Peanuts
Milk
Fruit
Jam

Bread
Jam
Soda
Chips
Milk
Fruit

Steak
Jam
Soda
Chips
Bread

Jam
Soda
Peanuts
Milk
Fruit

Jam
Soda
Chips
Milk
Bread

Fruit
Soda
Chips
Milk

Fruit
Soda
Peanuts
Milk

Fruit
Peanuts
Cheese
Yogurt

- Finding frequent patterns and associations among sets of items in transaction databases, relational databases, or other information repositories
- Applications
 - Basket data analysis
 - Cross-marketing
 - Catalog design
 - ...

- Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in a transaction
- Consider the following data set with 8 transactions
- Examples of itemsets
 - {Bread, Milk}
 - {Soda, Chips}
- Examples of association rules
 - {Bread} \Rightarrow {Milk}
 - {Soda} \Rightarrow {Chips}
 - {Bread} \Rightarrow {Jam}

TID	Items
1	Bread, Peanuts, Milk, Fruit, Jam
2	Bread, Jam, Soda, Chips, Milk, Fruit
3	Steak, Jam, Soda, Chips, Bread
4	Jam, Soda, Peanuts, Milk, Fruit
5	Jam, Soda, Chips, Milk, Bread
6	Fruit, Soda, Chips, Milk
7	Fruit, Soda, Peanuts, Milk
8	Fruit, Peanuts, Cheese, Yogurt

“Itemsets” are the fundamental pattern in association rule mining

- Implication of the form $X \Rightarrow Y$, where X and Y are itemsets
- Example, $\{\text{Bread}\} \Rightarrow \{\text{Milk}\}$
- Two rule evaluation metrics
 - Support
 - Confidence

evaluation metrics for association rules

$\{\text{Bread}\} \Rightarrow \{\text{Milk}\}$

Support

Fraction of transactions that contain both X and Y

$$s = \frac{\sigma(\{\text{Bread, Milk}\})}{\# \text{ of transactions}} = 0.38$$

Confidence

Measures how often items in Y appear in transactions that contain X

$$c = \frac{\sigma(\{\text{Bread, Milk}\})}{\sigma(\{\text{Bread}\})} = 0.75$$

- Given a set of transactions T , the goal of association rule mining is to find all rules having
 - $\text{support} \geq \text{minsup}$ threshold
 - $\text{confidence} \geq \text{minconf}$ 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
- Brute-force approach is computationally prohibitive!

- $\{\text{Bread, Jam}\} \Rightarrow \{\text{Milk}\}$ s=0.4 c=0.75
- $\{\text{Milk, Jam}\} \Rightarrow \{\text{Bread}\}$ s=0.4 c=0.75
- $\{\text{Bread}\} \Rightarrow \{\text{Milk, Jam}\}$ s=0.4 c=0.75
- $\{\text{Jam}\} \Rightarrow \{\text{Bread, Milk}\}$ s=0.4 c=0.6
- $\{\text{Milk}\} \Rightarrow \{\text{Bread, Jam}\}$ s=0.4 c=0.5
- They are all binary partitions of the itemset $\{\text{Milk, Bread, Jam}\}$
- Thus, rules originating from the same itemset have,
 - Same support but
 - Can have different confidence

We can decouple the support
and confidence requirements!

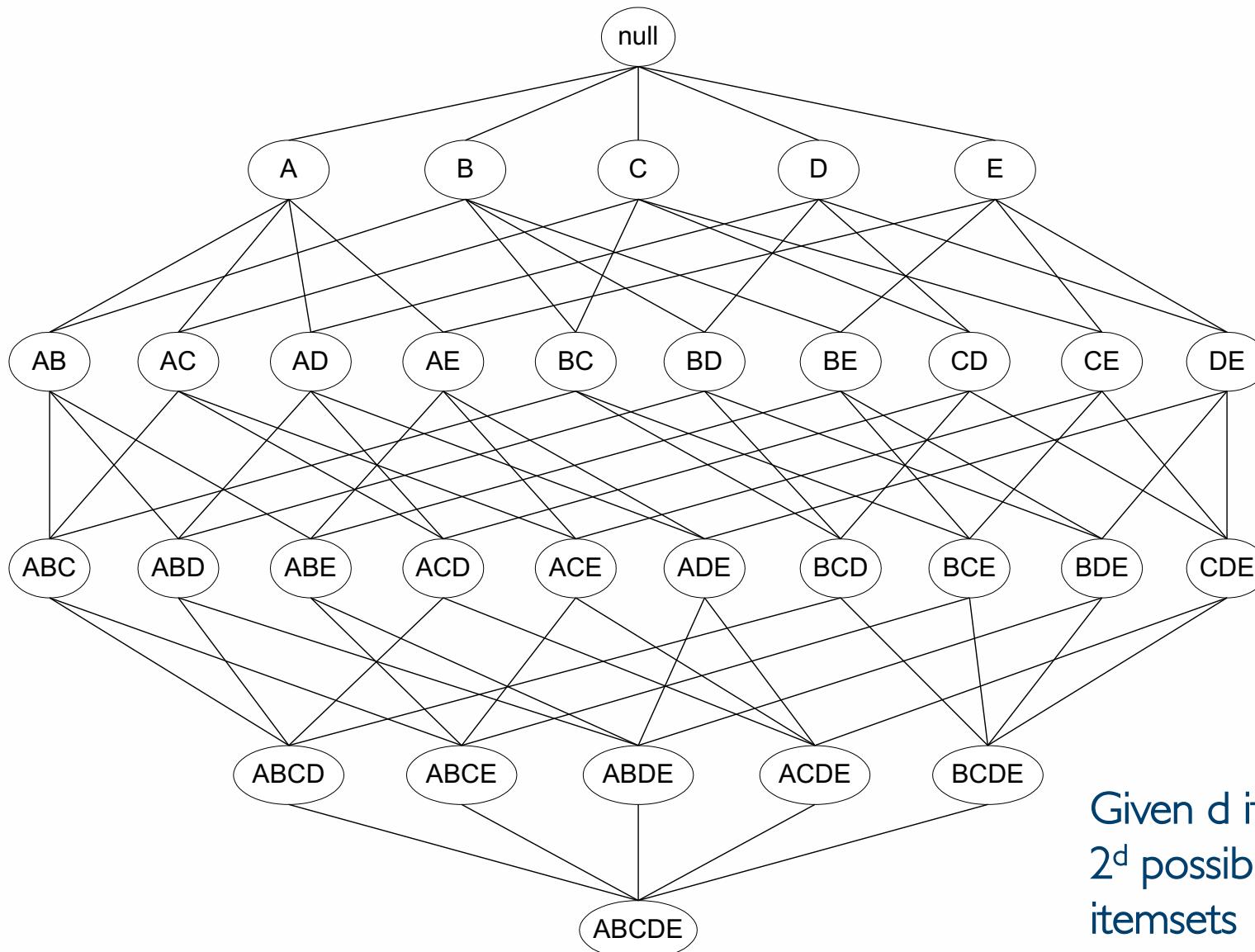
First, find the itemsets with
 $\text{support} \geq \text{minsup}$

Next, generate the rules with
 $\text{confidence} \geq \text{minconf}$

- **Frequent Itemset**
 - An itemset whose support is greater than or equal to the minsup threshold
- **Support**
 - Fraction of transactions that contain an itemset
 - $s(\{\text{Milk, Bread}\}) = 3/8$
 - $s(\{\text{Soda, Chips}\}) = 4/8$
- **Support count (σ)**
 - Raw count of occurrence of an itemset
 - For example, $\sigma(\{\text{Milk, Bread}\}) = 3$ $\sigma(\{\text{Soda, Chips}\}) = 4$

TID	Items
1	Bread, Peanuts, Milk, Fruit, Jam
2	Bread, Jam, Soda, Chips, Milk, Fruit
3	Steak, Jam, Soda, Chips, Bread
4	Jam, Soda, Peanuts, Milk, Fruit
5	Jam, Soda, Chips, Milk, Bread
6	Fruit, Soda, Chips, Milk
7	Fruit, Soda, Peanuts, Milk
8	Fruit, Peanuts, Cheese, Yogurt

- Frequent Itemset Generation
 - Generate all itemsets whose support $\geq \text{minsup}$
- Rule Generation
 - Generate high confidence rules from frequent itemset
 - Each rule is a binary partitioning of a frequent itemset
- Frequent itemset generation is computationally expensive



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

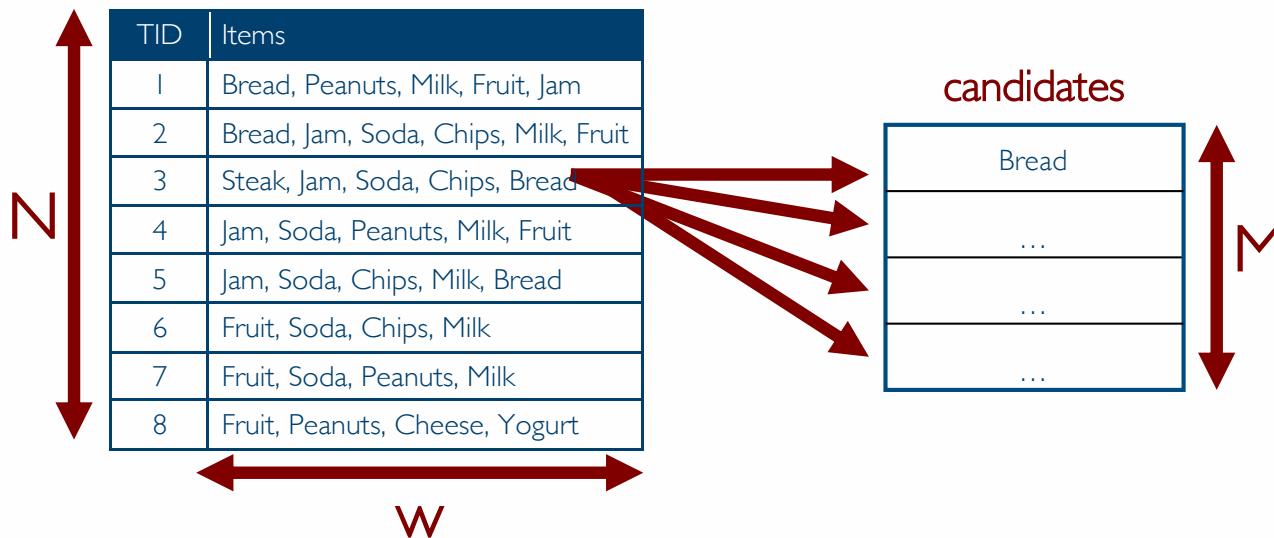
BRUTEFORCE ($\mathbf{D}, \mathcal{I}, \text{minsup}$):

```
1  $\mathcal{F} \leftarrow \emptyset$  // set of frequent itemsets
2 foreach  $X \subseteq \mathcal{I}$  do
3    $\text{sup}(X) \leftarrow \text{COMPUTESUPPORT}(X, \mathbf{D})$ 
4   if  $\text{sup}(X) \geq \text{minsup}$  then
5      $\mathcal{F} \leftarrow \mathcal{F} \cup \{(X, \text{sup}(X))\}$ 
6 return  $\mathcal{F}$ 
```

COMPUTESUPPORT (X, \mathbf{D}):

```
1  $\text{sup}(X) \leftarrow 0$ 
2 foreach  $\langle t, \mathbf{i}(t) \rangle \in \mathbf{D}$  do
3   if  $X \subseteq \mathbf{i}(t)$  then
4      $\text{sup}(X) \leftarrow \text{sup}(X) + 1$ 
5 return  $\text{sup}(X)$ 
```

- Each itemset in the lattice is a candidate frequent itemset
- Count the support of each candidate by scanning the database



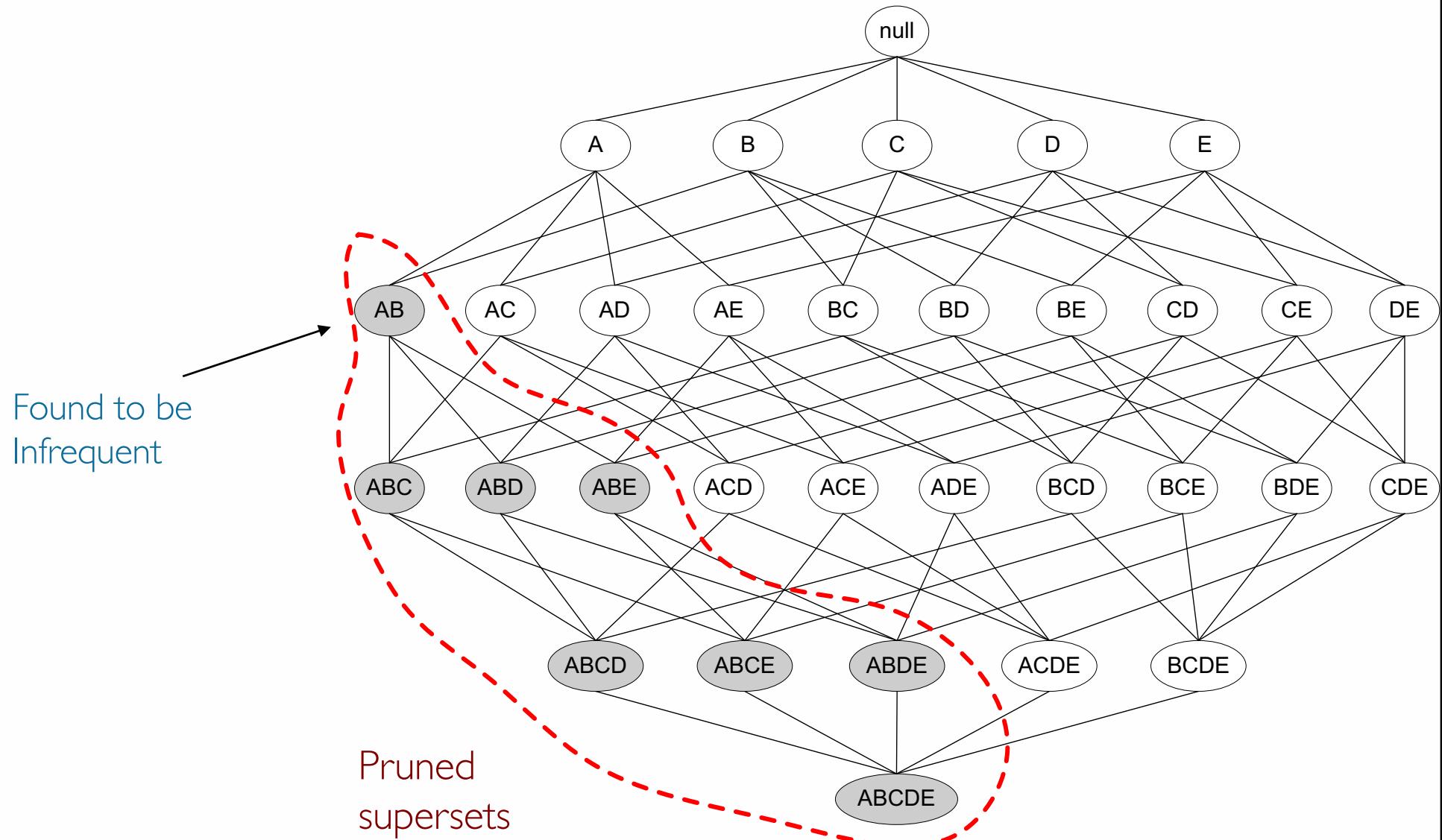
- Match each transaction against every candidate
- Complexity $\sim O(NMw)$ => Expensive since $M = 2^d$

- Reduce the number of candidates (M)
 - Complete search has $M=2^d$
 - Use pruning techniques to reduce M
- Reduce the number of transactions (N)
 - Reduce size of N as the size of itemset increases
- 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 the Number of Candidates: The Apriori Principle

19

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



How does the Apriori principle work?

21

Minimum Support = 4/8

Item	Support
Bread	4/8
Peanuts	4/8
Milk	6/8
Fruit	6/8
Jam	5/8
Soda	6/8
Chips	4/8
Steak	1/8
Cheese	1/8
Yogurt	1/8

1-itemsets

2-Itemset	Support
Bread, Jam	4/8
Peanuts, Fruit	4/8
Milk, Fruit	5/8
Milk, Jam	4/8
Milk, Soda	5/8
Fruit, Soda	4/8
Jam, Soda	4/8
Soda, Chips	4/8

2-itemsets

3-Itemset	Support
Milk, Fruit, Soda	4/8

3-itemsets

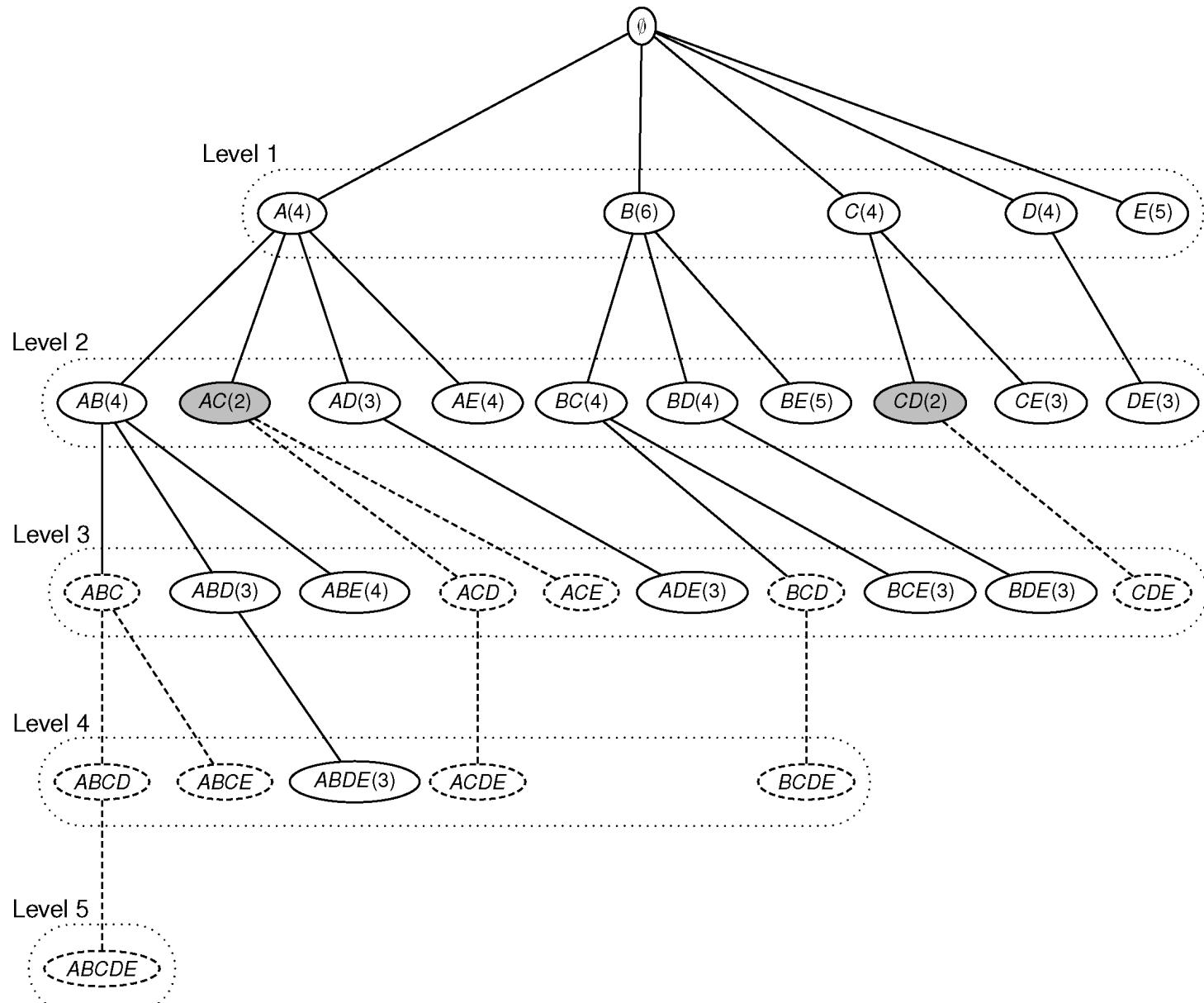
- Given the following database and a min support of 3, generate all the frequent itemsets

D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

Binary Database

t	i(t)
1	ABDE
2	BCE
3	ABDE
4	ABCE
5	ABCDE
6	BCD

Transaction Database



- Let $k=1$
- Generate frequent itemsets of length l
- Repeat until no new frequent itemsets are identified
 - Generate length $(k+1)$ candidate itemsets from length k frequent itemsets
 - Prune candidate itemsets containing subsets of length k that are infrequent
 - Count the support of each candidate by scanning the database
 - Eliminate candidates that are infrequent, leaving only those that are frequent

APRIORI ($\mathbf{D}, \mathcal{I}, \text{minsup}$):

```
1  $\mathcal{F} \leftarrow \emptyset$ 
2  $\mathcal{C}^{(1)} \leftarrow \{\emptyset\}$  // Initial prefix tree with single items
3 foreach  $i \in \mathcal{I}$  do Add  $i$  as child of  $\emptyset$  in  $\mathcal{C}^{(1)}$  with  $\text{sup}(i) \leftarrow 0$ 
4  $k \leftarrow 1$  //  $k$  denotes the level
5 while  $\mathcal{C}^{(k)} \neq \emptyset$  do
6   COMPUTESUPPORT ( $\mathcal{C}^{(k)}, \mathbf{D}$ )
7   foreach leaf  $X \in \mathcal{C}^{(k)}$  do
8     if  $\text{sup}(X) \geq \text{minsup}$  then  $\mathcal{F} \leftarrow \mathcal{F} \cup \{(X, \text{sup}(X))\}$ 
9     else remove  $X$  from  $\mathcal{C}^{(k)}$ 
10   $\mathcal{C}^{(k+1)} \leftarrow \text{EXTENDPREFIXTREE } (\mathcal{C}^{(k)})$ 
11   $k \leftarrow k + 1$ 
12 return  $\mathcal{F}^{(k)}$ 
```

COMPUTESUPPORT ($\mathcal{C}^{(k)}$, \mathbf{D}):

```
1 foreach  $\langle t, \mathbf{i}(t) \rangle \in \mathbf{D}$  do
2   foreach  $k$ -subset  $X \subseteq \mathbf{i}(t)$  do
3     if  $X \in \mathcal{C}^{(k)}$  then  $sup(X) \leftarrow sup(X) + 1$ 
```

EXTENDPREFIXTREE ($\mathcal{C}^{(k)}$):

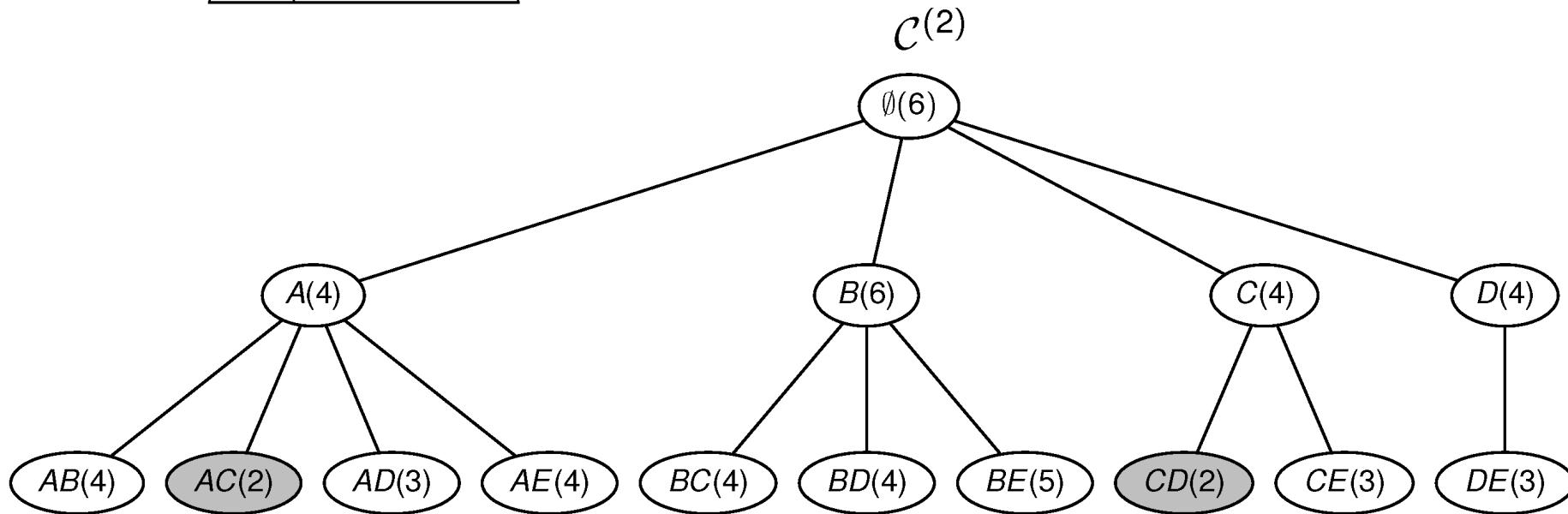
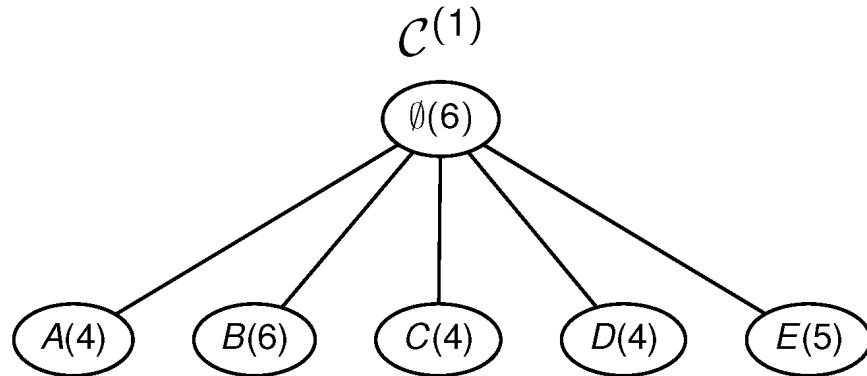
```
1 foreach leaf  $X_a \in \mathcal{C}^{(k)}$  do
2   foreach leaf  $X_b \in \text{SIBLING}(X_a)$ , such that  $b > a$  do
3      $X_{ab} \leftarrow X_a \cup X_b$ 
        // prune candidate if there are any infrequent
        subsets
4     if  $X_j \in \mathcal{C}^{(k)}$ , for all  $X_j \subset X_{ab}$ , such that  $|X_j| = |X_{ab}| - 1$  then
5       Add  $X_{ab}$  as child of  $X_a$  with  $sup(X_{ab}) \leftarrow 0$ 
6   if no extensions from  $X_a$  then
7     remove  $X_a$ , and all ancestors of  $X_a$  with no extensions, from  $\mathcal{C}^{(k)}$ 
8 return  $\mathcal{C}^{(k)}$ 
```

Apriori Algorithm Example (min sup=3):

First two levels of the search

27

D	
t	$\mathbf{i}(t)$
1	$ABDE$
2	BCE
3	$ABDE$
4	$ABCE$
5	$ABCDE$
6	BCD



Eclat

- Leverages the tidsets directly for support computation.
- The support of a candidate itemset can be computed by intersecting the tidsets of suitably chosen subsets.
- Given $t(X)$ and $t(Y)$ for any two frequent itemsets X and Y ,
then $t(XY) = t(X) \wedge t(Y)$
- And $\text{sup}(XY) = |t(XY)|$

D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

Binary Database

<i>t</i>	<i>i(t)</i>
1	<i>ABDE</i>
2	<i>BCE</i>
3	<i>ABDE</i>
4	<i>ABCE</i>
5	<i>ABCDE</i>
6	<i>BCD</i>

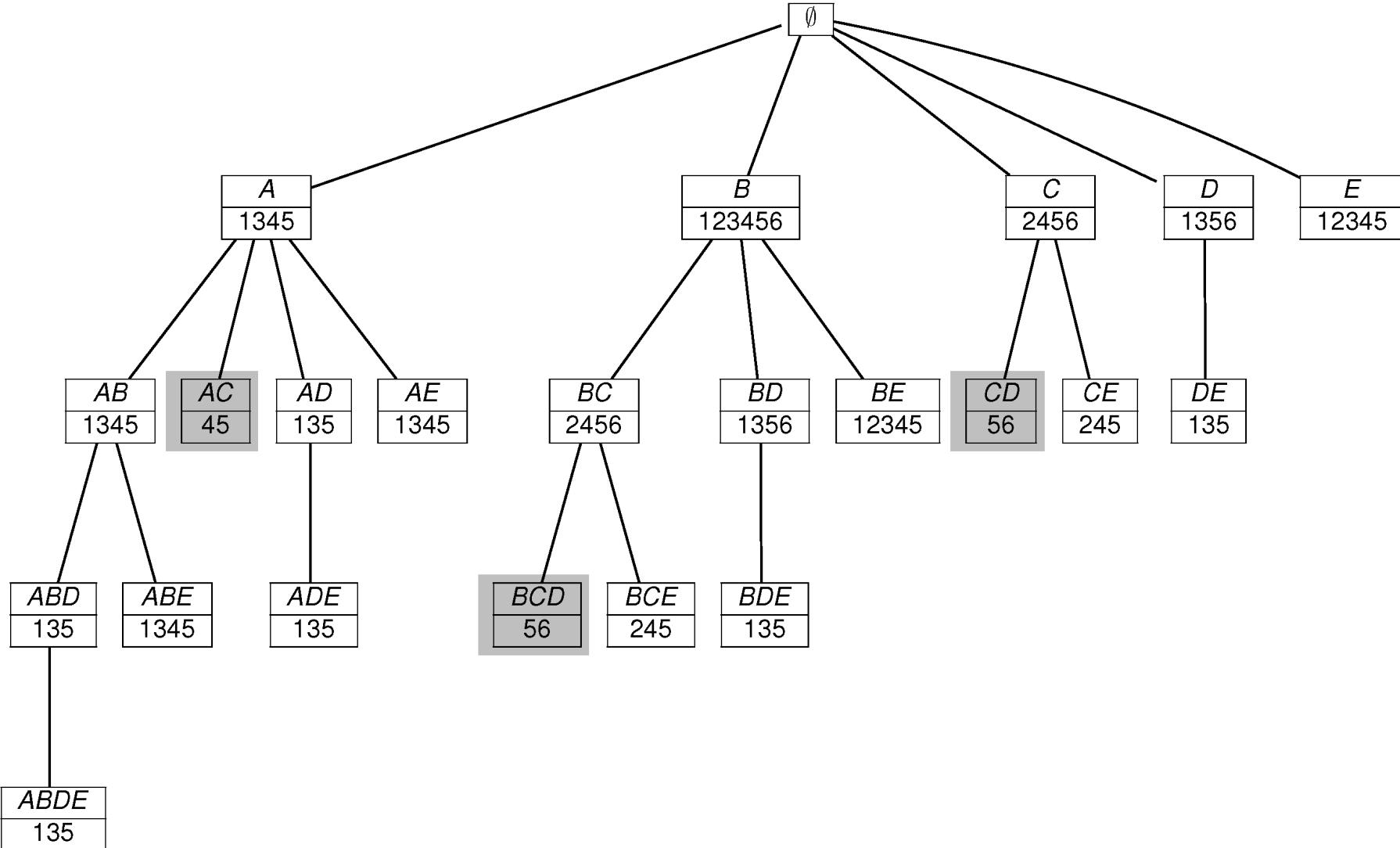
Transaction Database

<i>t(x)</i>					
	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
1	1	1	2	1	1
3	2	2	4	3	2
4	3	3	5	5	3
5	4	4	6	6	4
		5			5
		6			

Vertical Database

Example

30



```

// Initial Call:  $\mathcal{F} \leftarrow \emptyset, P \leftarrow \{\langle i, \mathbf{t}(i) \rangle \mid i \in \mathcal{I}, |\mathbf{t}(i)| \geq \text{minsup}\}$ 
ECLAT ( $P, \text{minsup}, \mathcal{F}$ ):
1 foreach  $\langle X_a, \mathbf{t}(X_a) \rangle \in P$  do
2    $\mathcal{F} \leftarrow \mathcal{F} \cup \{(X_a, \text{sup}(X_a))\}$ 
3    $P_a \leftarrow \emptyset$ 
4   foreach  $\langle X_b, \mathbf{t}(X_b) \rangle \in P$ , with  $X_b > X_a$  do
5      $X_{ab} = X_a \cup X_b$ 
6      $\mathbf{t}(X_{ab}) = \mathbf{t}(X_a) \cap \mathbf{t}(X_b)$ 
7     if  $\text{sup}(X_{ab}) \geq \text{minsup}$  then
8        $P_a \leftarrow P_a \cup \{\langle X_{ab}, \mathbf{t}(X_{ab}) \rangle\}$ 
9   if  $P_a \neq \emptyset$  then ECLAT ( $P_a, \text{minsup}, \mathcal{F}$ )

```

Frequent Patterns Mining Without Candidate Generation

- The core of the Apriori algorithm
 - Use frequent $(k-1)$ -itemsets to generate candidate frequent k -itemsets
 - Use database scan and pattern matching to collect counts for the candidate itemsets
- It may need to generate huge candidate sets
 - 10^4 frequent 1-itemset will generate 10^7 candidate 2-itemsets
 - To discover a frequent pattern of size 100, e.g., $\{a_1, a_2, \dots, a_{100}\}$, needs to generate $2^{100} \sim 10^{30}$ candidates.
 - Multiple scans of database, it needs $(n+1)$ scans, n is the length of the longest pattern

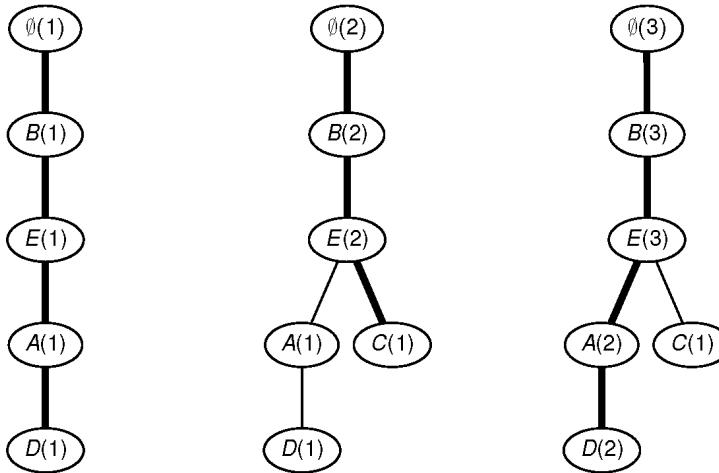
Mining Frequent Patterns Without Candidate Generation

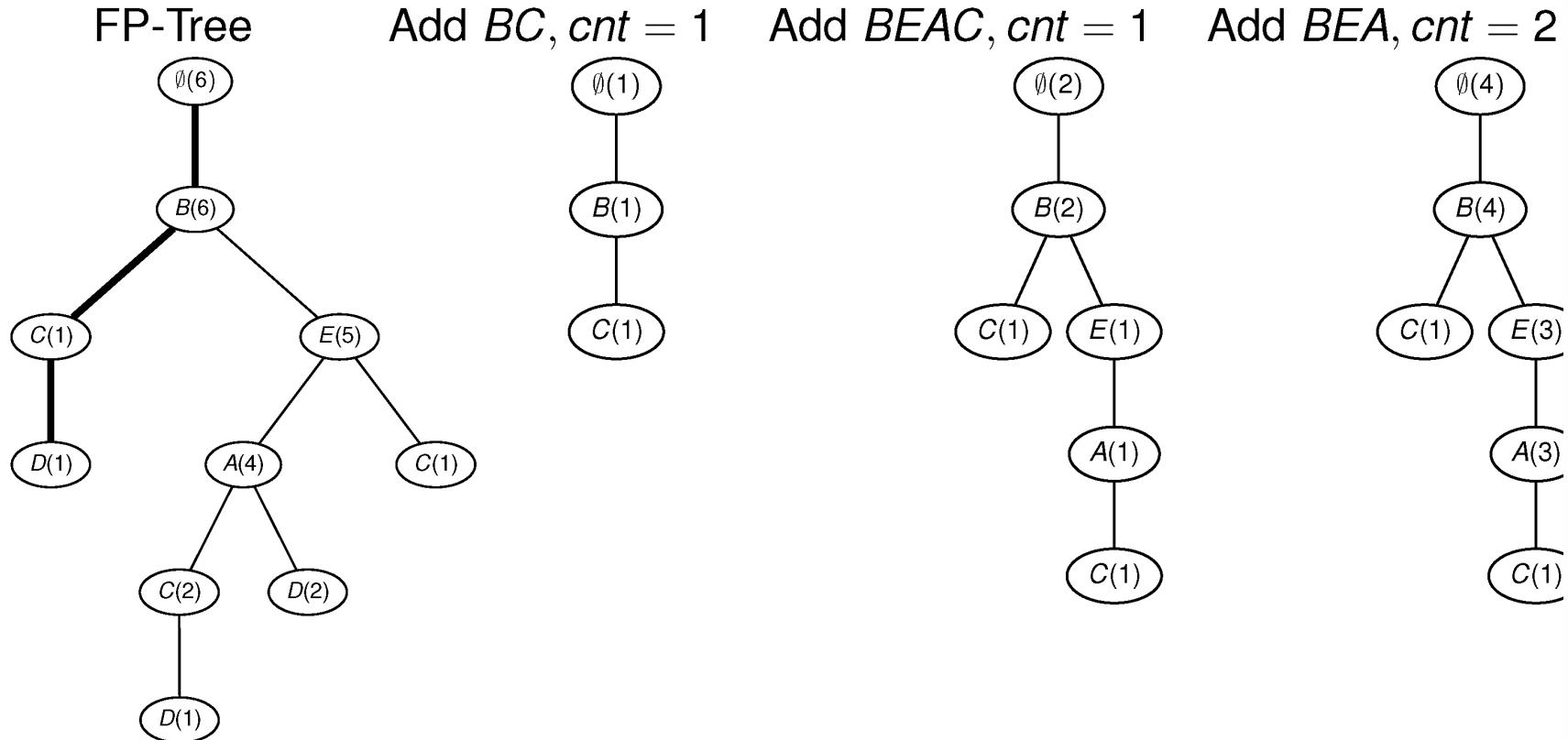
34

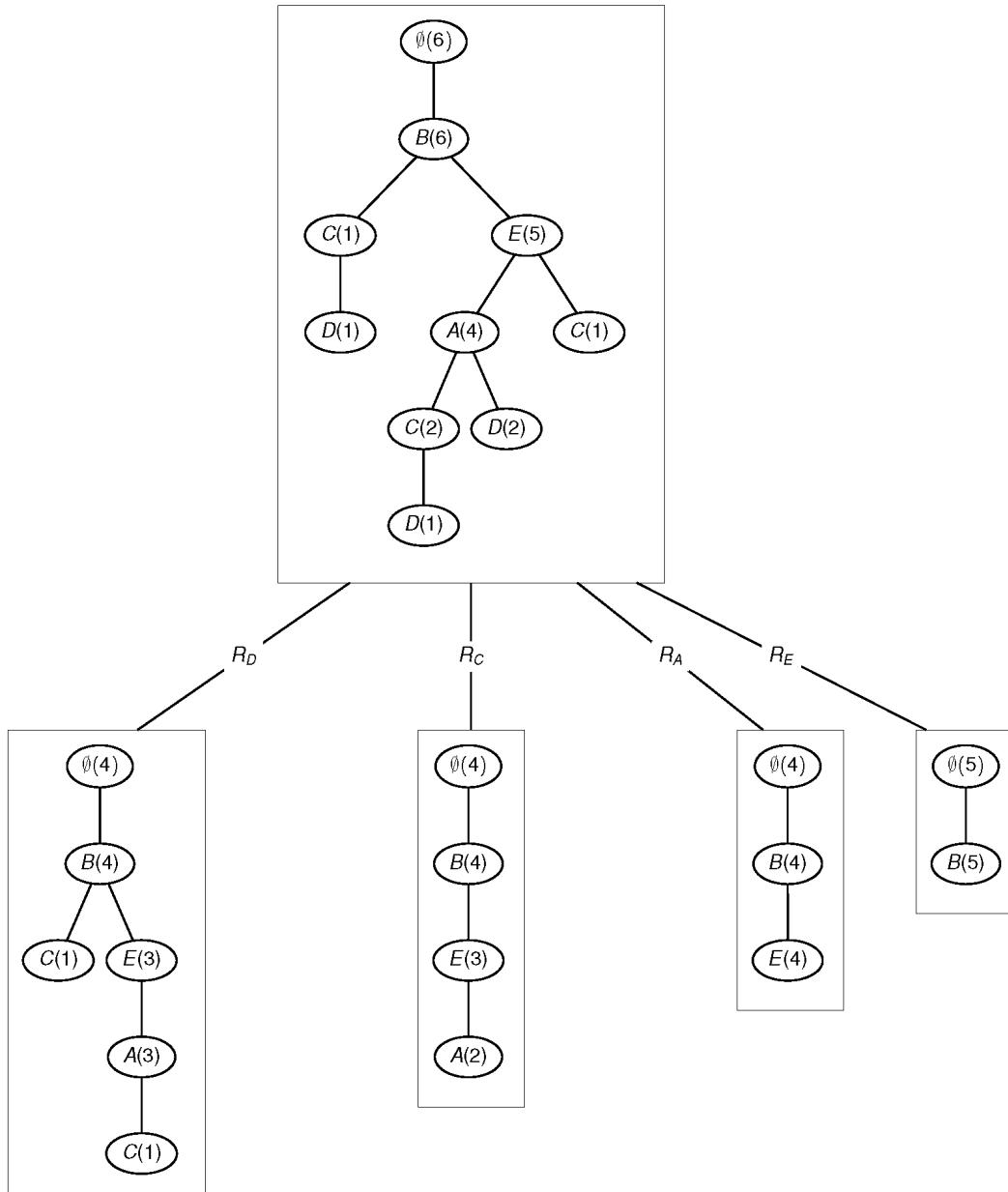
- Compress a large database into a compact, Frequent-Pattern tree (FP-tree) structure
 - Highly condensed, but complete for frequent pattern mining
 - Avoid costly database scans
- Use an efficient, FP-tree-based frequent pattern mining method
- A divide-and-conquer methodology: decompose mining tasks into smaller ones
- Avoid candidate generation: sub-database test only

- Leave the generate-and-test paradigm of Apriori
- Data sets are encoded using a compact structure, the FP-tree
- Frequent itemsets are extracted directly from the FP-tree
- Major Steps to mine FP-tree
 - Construct the frequent pattern tree
 - For each frequent item i compute the projected FP-tree
 - Recursively mine conditional FP-trees and grow frequent patterns obtained so far
 - If the conditional FP-tree contains a single path, simply enumerate all the patterns

Transactions
BEAD
BEC
BEAD
BEAC
BEACD
BCD







- **Completeness**
 - Preserve complete information for frequent pattern mining
 - Never break a long pattern of any transaction
- **Compactness**
 - Reduce irrelevant info—infrequent items are gone
 - Items in frequency descending order: the more frequently occurring, the more likely to be shared
 - Never be larger than the original database (not count node-links and the count field)

- **Divide-and-conquer:**
 - decompose both the mining task and DB according to the frequent patterns obtained so far
 - leads to focused search of smaller databases
- **Other factors**
 - No candidate generation, no candidate test
 - Compressed database: FP-tree structure
 - No repeated scan of entire database
 - Basic ops—counting local freq items and building sub FP-tree, no pattern search and matching

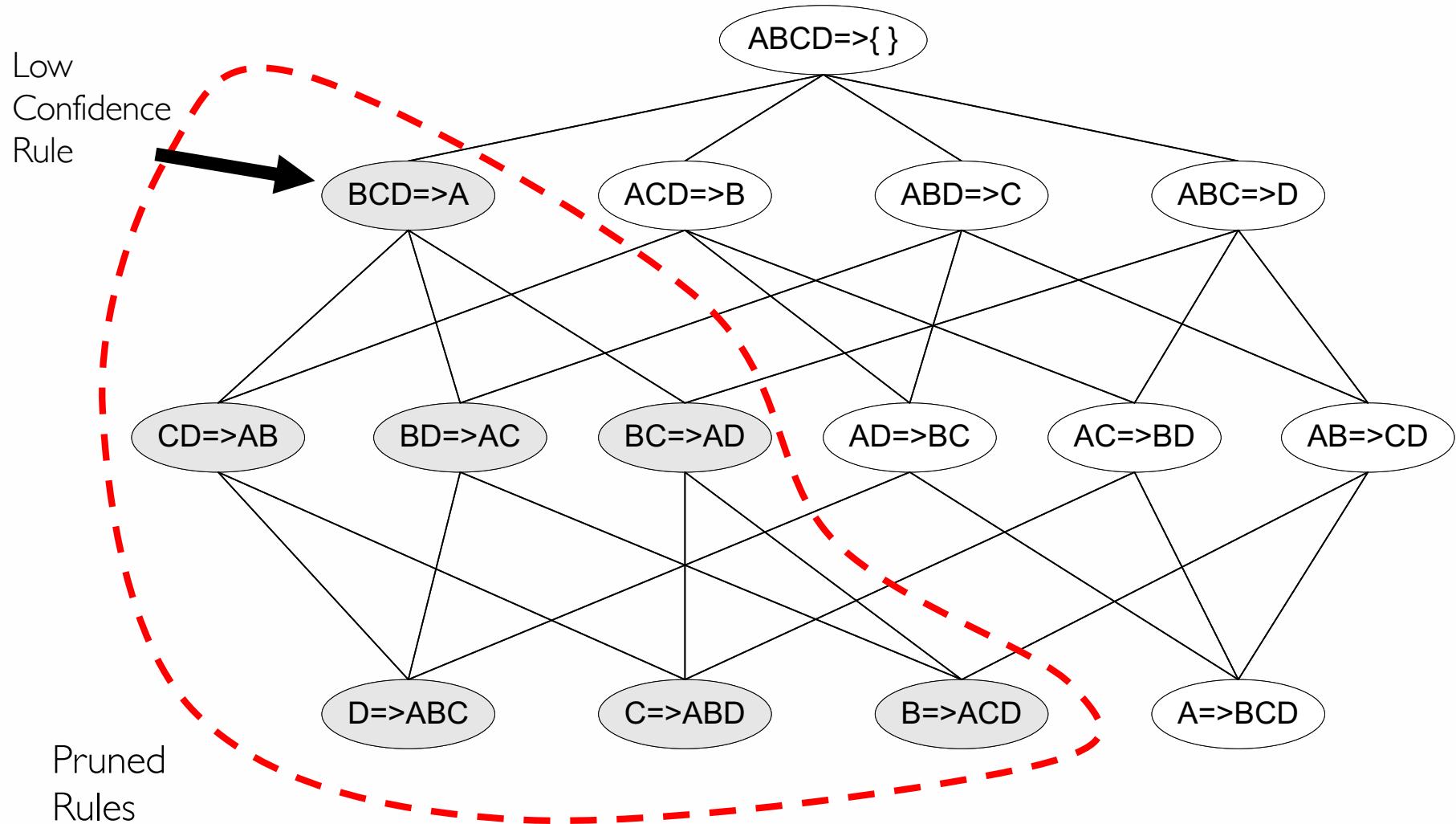
Rule Generation

- 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:
 $ABC \Rightarrow D$, $ABD \Rightarrow C$, $ACD \Rightarrow B$, $BCD \Rightarrow A$, $A \Rightarrow BCD$,
 $B \Rightarrow ACD$, $C \Rightarrow ABD$, $D \Rightarrow ABC$, $AB \Rightarrow CD$, $AC \Rightarrow BD$,
 $AD \Rightarrow BC$, $BC \Rightarrow AD$, $BD \Rightarrow AC$, $CD \Rightarrow AB$
- If $|L| = k$, then there are $2^k - 2$ candidate association rules
(ignoring $L \rightarrow \emptyset$ and $\emptyset \rightarrow L$)

How to efficiently generate rules from frequent itemsets?

44

- Confidence does not have an anti-monotone property
- $c(ABC \Rightarrow D)$ can be larger or smaller than $c(AB \Rightarrow D)$
- However, confidence of rules generated from the same itemset has an anti-monotone property
- $L = \{A, B, C, D\}$: $c(ABC \Rightarrow D) \geq c(AB \Rightarrow CD) \geq c(A \Rightarrow BCD)$
- Confidence is anti-monotone with respect to the number of items on the right-hand side of the rule



```
ASSOCIATIONRULES ( $\mathcal{F}$ ,  $minconf$ ):  
1 foreach  $Z \in \mathcal{F}$ , such that  $|Z| \geq 2$  do  
2    $\mathcal{A} \leftarrow \{X \mid X \subset Z, X \neq \emptyset\}$   
3   while  $\mathcal{A} \neq \emptyset$  do  
4      $X \leftarrow$  maximal element in  $\mathcal{A}$   
5      $\mathcal{A} \leftarrow \mathcal{A} \setminus X$  // remove  $X$  from  $\mathcal{A}$   
6      $c \leftarrow sup(Z)/sup(X)$   
7     if  $c \geq minconf$  then  
8       print  $X \rightarrow Y$ ,  $sup(Z)$ ,  $c$   
9     else  
10       $\mathcal{A} \leftarrow \mathcal{A} \setminus \{W \mid W \subset X\}$   
           // remove all subsets of  $X$  from  $\mathcal{A}$ 
```

- Suppose that we computed the following frequent itemsets and we have a minconf of 0.9

sup	itemsets
6	B
5	E, BE
4	$A, C, D, AB, AE, BC, BD, ABE$
3	$AD, CE, DE, ABD, ADE, BCE, BDE, ABDE$

- We start from the largest itemset (ABDE) and consider all the possible subsets $\mathcal{A} = \{ABDE(3), ABD(3), ABE(4), ADE(3), BDE(3), AB(4), AD(3), AE(4), BD(4), BE(5), DE(3), A(4), B(6), D(4), E(5)\}$

- The first subset is $X = ABD$, and the confidence of $ABD \Rightarrow E$ is $3/3 = 1.0$, so we output the rule.
- The next subset is $X = ABE$, but the corresponding rule $ABE \Rightarrow D$ is not strong since $c(ABE \Rightarrow D) = 3/4 = 0.75$
- We can thus remove from \mathcal{A} all subsets of ABE
- The updated set of antecedents is therefore
 $\mathcal{A} = \{ADE(3), BDE(3), AD(3), BD(4), DE(3), D(4)\}$
- When the algorithm ends, it will output the rules,
 $ABD \Rightarrow E$ conf=1.0; $ADE \Rightarrow B$ conf=1.0; $BDE \Rightarrow A$, conf=1.0;
 $AD \Rightarrow BE$ conf=1.0; $DE \Rightarrow AB$ conf=1.0

Rule Assessment Measures

- If minsup is set too high, we could miss itemsets involving interesting rare items (e.g., expensive products)
- If minsup is set too low, it is computationally expensive and the number of itemsets is very large
- A single minimum support threshold may not be effective

- Lift is the ratio of the observed joint probability of X and Y to the expected joint probability if they were statistically independent,

$$\text{lift}(X \Rightarrow Y) = \text{sup}(X \cup Y) / \text{sup}(X)\text{sup}(Y) = \text{conf}(X \Rightarrow Y) / \text{sup}(Y)$$

- Lift is a measure of the deviation from stochastic independence (if it is 1 then X and Y are independent)
- Lift also measures the surprise of the rule. A lift close to 1 means that the support of a rule is expected considering the supports of its components.
- We typically look for values of lift that are much larger (i.e., above expectation) or much smaller (i.e., below expectation) than one.

Summarizing Itemsets

- A frequent itemset X is called maximal if it has no frequent supersets
- The set of all maximal frequent itemsets, given as

$$M = \{X \mid X \in F \text{ and } \nexists Y \supset X, \text{ such that } Y \in F\}$$

- M is a condensed representation of the set of all frequent itemset F , because we can determine whether any itemset is frequent or not using M
- If there is a maximal itemset Z such that $X \subseteq Z$, then X must be frequent, otherwise Z cannot be frequent
- However, M alone cannot be used to determine $\text{sup}(X)$, we can only use to have a lower-bound, that is, $\text{sup}(X) \geq \text{sup}(Z)$ if $X \subseteq Z \in M$.

- An itemset X is closed if all supersets of X have strictly less support, that is,
$$\text{sup}(X) > \text{sup}(Y), \text{ for all } Y \supset X$$
- The set of all closed frequent itemsets C is a condensed representation, as we can determine whether an itemset X is frequent, as well as the exact support of X using C alone

- A frequent itemset X is a minimal generator if it has no subsets with the same support

$$G = \{ X \mid X \in F \text{ and } \nexists Y \subset X, \text{ such that } \text{sup}(X) = \text{sup}(Y) \}$$

- Thus, all subsets of X have strictly higher support, that is,

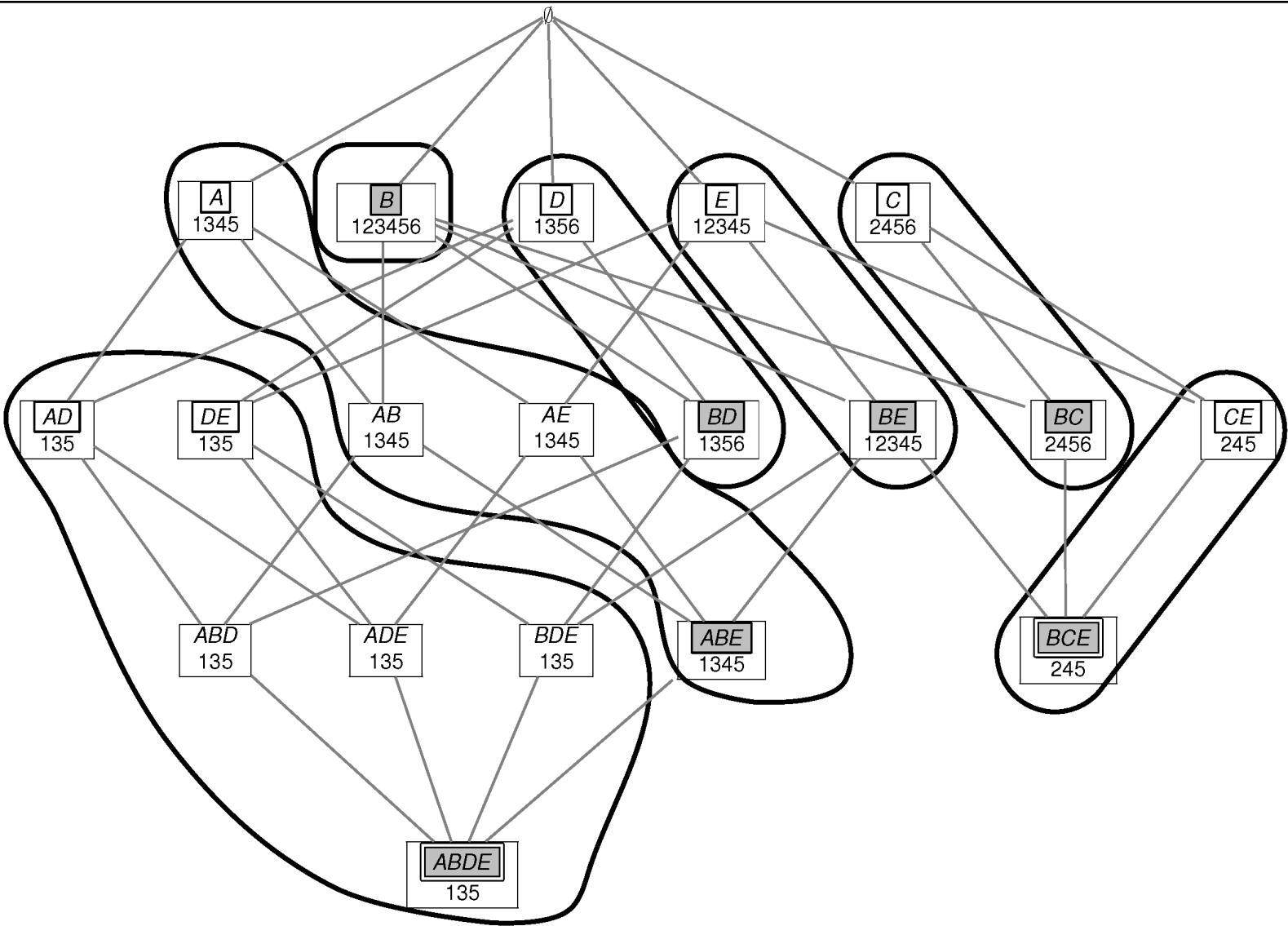
$$\text{sup}(X) < \text{sup}(Y)$$

Transaction database

Tid	Itemset
1	<i>ABDE</i>
2	<i>BCE</i>
3	<i>ABDE</i>
4	<i>ABCE</i>
5	<i>ABCDE</i>
6	<i>BCD</i>

Frequent itemsets ($\text{minsup} = 3$)

sup	Itemsets
6	<i>B</i>
5	<i>E, BE</i>
4	<i>A, C, D, AB, AE, BC, BD, ABE</i>
3	<i>AD, CE, DE, ABD, ADE, BCE, BDE, ABDE</i>



Shaded boxes are closed itemsets

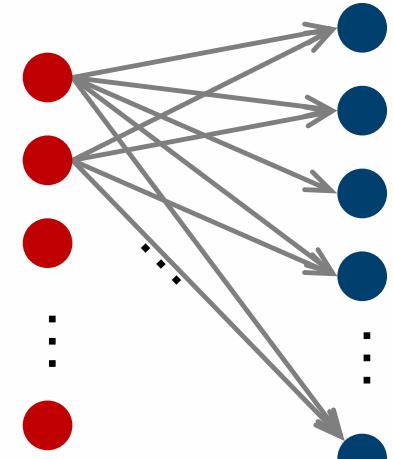
Simple boxes are generators

Shaded boxes with double lines are maximal itemsets

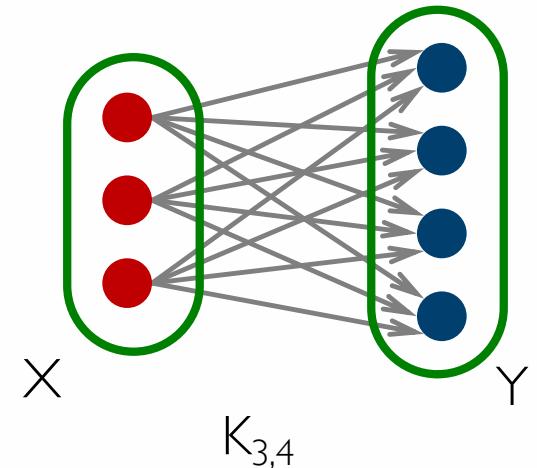
Searching for Small Communities (Trawling)

Searching for small communities in the Web graph (Trawling)

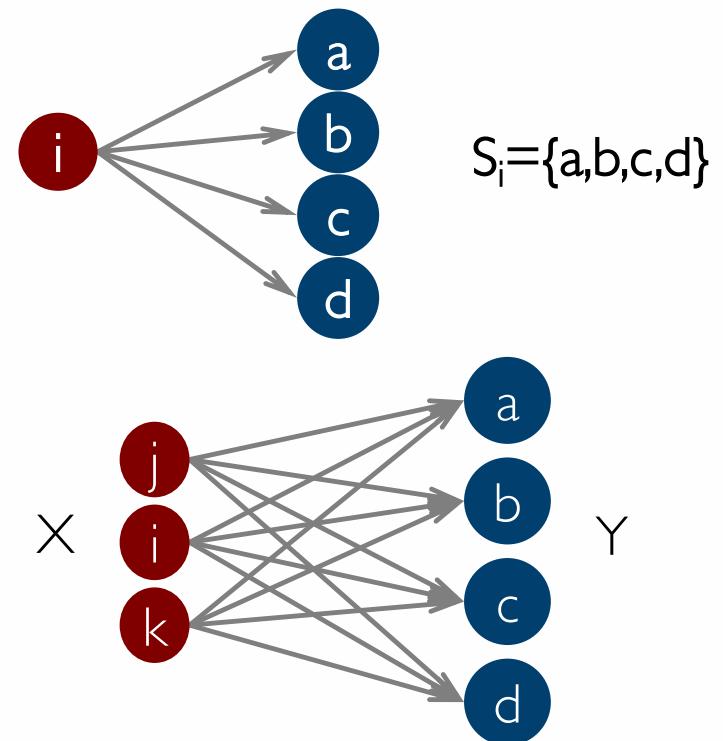
- **Trawling**
 - What is the signature of a community in a Web graph?
 - The underlying intuition, that small communities involve many people talking about the same things
 - Use this to define “topics”: what the same people on the left talk about on the right?
- **More formally**
 - Enumerate complete bipartite subgraphs $K_{s,t}$
 - $K_{s,t}$ has s nodes on the “left” and t nodes on the “right”
 - The left nodes link to the same node of on the right, forming a fully connected bipartite graph



Dense 2-layer

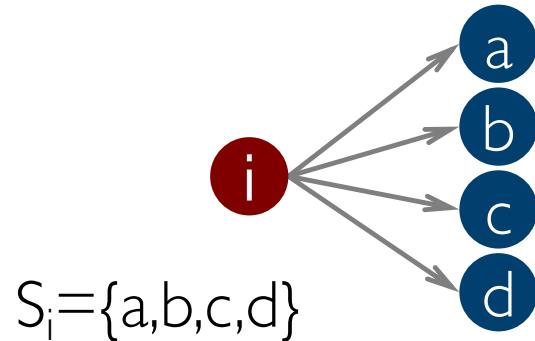


- Searching for such complete bipartite graphs can be viewed as a frequent itemset mining problem
- View each node i as a set S_i of nodes i points to
- $K_{s,t} = \text{a set } Y \text{ of size } t \text{ that occurs in } s \text{ sets } S_i$
- Looking for $K_{s,t}$ is equivalent to setting the frequency threshold to s and look at layer t (i.e., all frequent sets of size t)

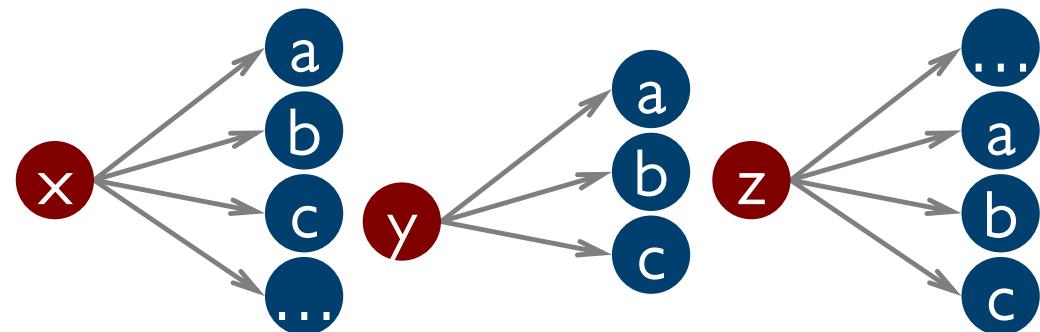


$s = \text{minimum support } (|X|=s)$
 $t = \text{itemset size } (|Y|=t)$

View each node i as a set S_i of nodes i points to

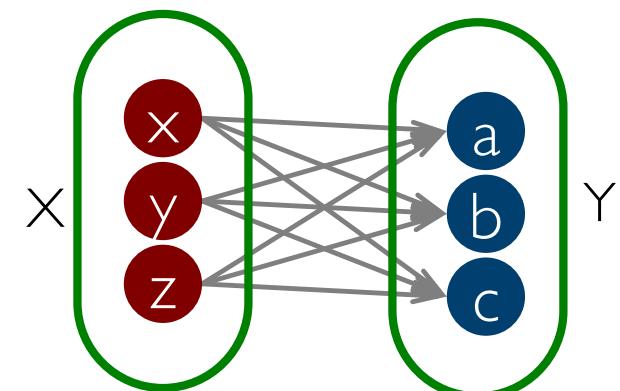


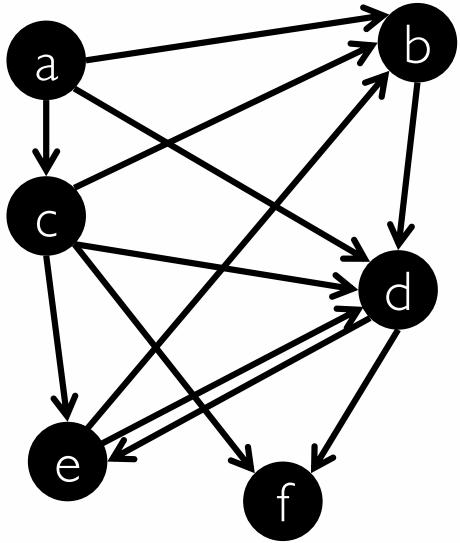
Suppose we find a frequent itemset $Y = \{a, b, c\}$ of supp s ; then, there are s nodes that link to all of $\{a, b, c\}$:



Find frequent itemsets with support s and itemset size t is equivalent to find $K_{s,t}$

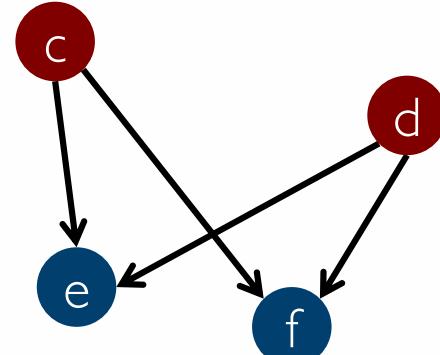
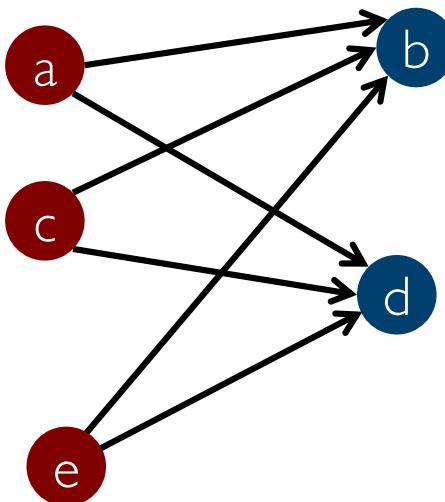
$K_{s,t}$ = a set Y of size t that occurs in s sets S_i





- Support threshold $s=2$
 - $\{b,d\}$: support 3
 - $\{e,f\}$: support 2
- And we just found 2 bipartite subgraphs:

- Itemsets
 - $a = \{b,c,d\}$
 - $b = \{d\}$
 - $c = \{b,d,e,f\}$
 - $d = \{e,f\}$
 - $e = \{b,d\}$
 - $f = \{\}$



Mining Frequent Sequences

- A sequence is an ordered list of symbols $s = s_1, s_2, \dots, s_k$ where s_i or $s[i]$ denote the symbol at position i
- Notation $s[i:j]$ denotes $s_{i+1}, s_{i+2}, \dots, s_j$
- Let $s = s_1, s_2, \dots, s_n$ and $r = r_1, r_2, \dots, r_m$ we say that r is a subsequence of s denoted $r \subseteq s$ if there is a function ϕ from $[1, m]$ to $[1, n]$ such that $r[i] = s[\phi(i)]$ and for any position i, j in r ,

$$i < j \Rightarrow \phi(i) < \phi(j)$$

- r is a consecutive subsequence when

$$r_1, r_2, \dots, r_m = s_j, s_{j+1}, \dots, s_{j+m-1}$$

- Given $s = \text{ACTGAAACG}$
- $r_1 = \text{CGAAG}$ is a subsequence of s
- $r_2 = \text{CTGA}$ is a consecutive subsequence of s
- $r_3 = \text{ACT}$ is a prefix of s
- $r_4 = \text{AACG}$ is a suffix of s

- Given a database D containing N sequences, the support of a sequence r in D is defined as the total number of sequences in D that contains r

$$sup(r) = |\{s_i \in D | r \subseteq s_i\}|$$

- The relative support of r is the percentage of sequences that contain r,

$$rsup(r) = sup(r)/N$$

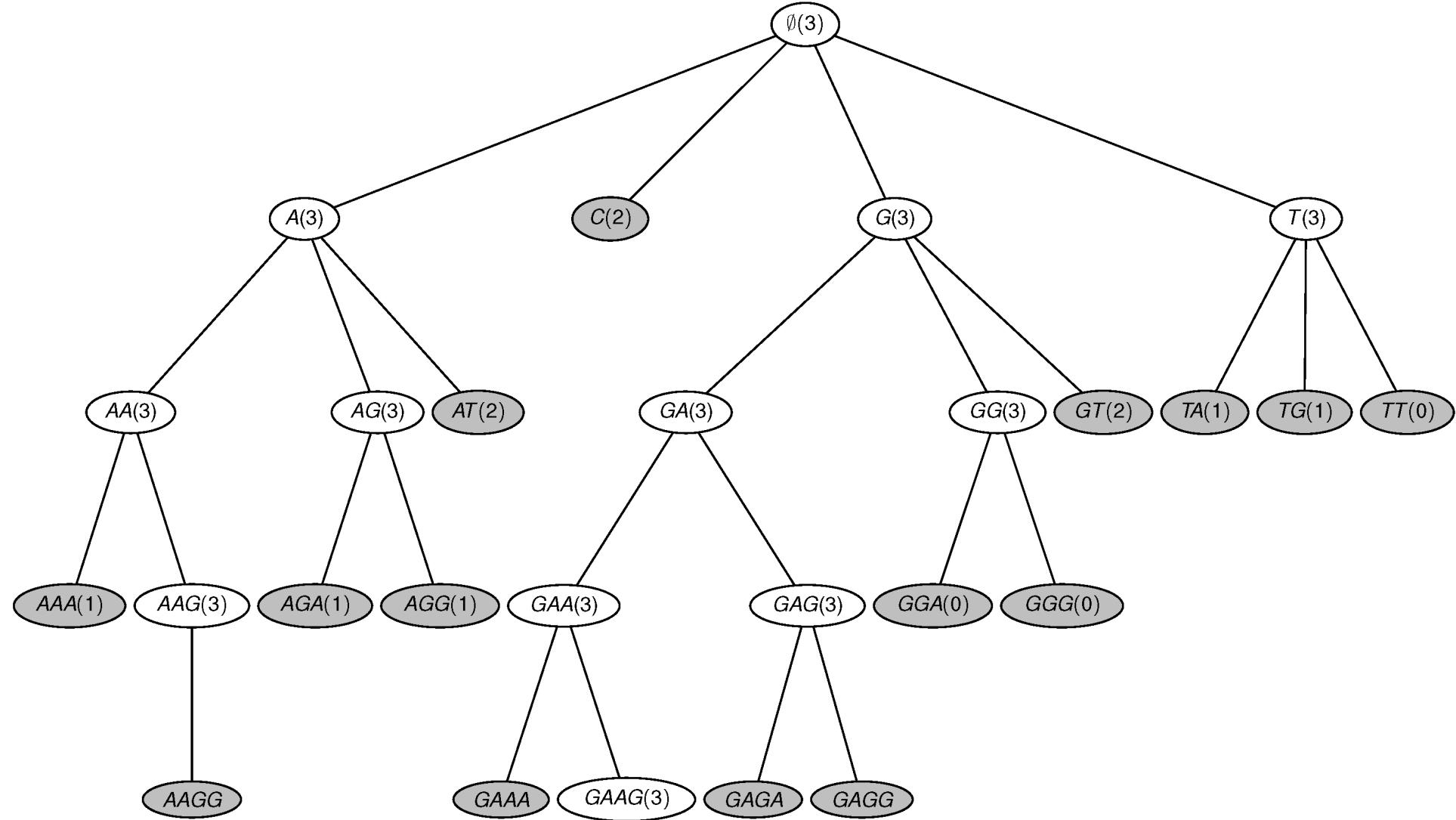
- Given a minsup threshold we that r is frequent if $sup(r) \geq \text{minsup}$
- Note that, in sequence mining, we need to consider all possible permutations of the items not just the combinations

- Given the following sequence database

Id	Sequence
s_1	CAGAAGT
s_2	TGACAG
s_3	GAAGT

- With a minsup of 3, the set of frequent subsequences is
 - A(3), G(3), T(3)
 - AA(3), AG(3), GA(3), GG(3)
 - AAG(3), GAA(3), GAG(3)
 - GAAG(3)

- Searches the sequence prefix tree using a level-wise (breadth-first search).
- Given the set of frequent sequences at level k , the algorithm generates the candidate for level $k+1$ and compute the support of each candidate and prune the not frequent ones.
- For each sequence s_i in D , we check if a candidate r is a subsequence of s , if it is the support of r is incremented. Once the frequent sequences at level k are computed the $k+1$ candidates are generated.
- For each leaf r_a , the sequence is extended with the last symbol of any other leaf r_b that shares the same prefix (it has the same parent) so $r_{ab} = r_a + r_b[k]$. If r_{ab} is infrequent, we prune it.



GSP (\mathbf{D} , Σ , $minsup$):

```
1  $\mathcal{F} \leftarrow \emptyset$ 
2  $\mathcal{C}^{(1)} \leftarrow \{\emptyset\}$  // Initial prefix tree with single symbols
3 foreach  $s \in \Sigma$  do Add  $s$  as child of  $\emptyset$  in  $\mathcal{C}^{(1)}$  with  $sup(s) \leftarrow 0$ 
4  $k \leftarrow 1$  //  $k$  denotes the level
5 while  $\mathcal{C}^{(k)} \neq \emptyset$  do
6   COMPUTESUPPORT ( $\mathcal{C}^{(k)}$ ,  $\mathbf{D}$ )
7   foreach leaf  $\mathbf{s} \in \mathcal{C}^{(k)}$  do
8     if  $sup(\mathbf{r}) \geq minsup$  then  $\mathcal{F} \leftarrow \mathcal{F} \cup \{(\mathbf{r}, sup(\mathbf{r}))\}$ 
9     else remove  $\mathbf{s}$  from  $\mathcal{C}^{(k)}$ 
10   $\mathcal{C}^{(k+1)} \leftarrow \text{EXTENDPREFIXTREE } (\mathcal{C}^{(k)})$ 
11   $k \leftarrow k + 1$ 
12 return  $\mathcal{F}^{(k)}$ 
```

```

COMPUTESUPPORT ( $\mathcal{C}^{(k)}$ ,  $\mathbf{D}$ ):
1 foreach  $s_i \in \mathbf{D}$  do
2   foreach  $r \in \mathcal{C}^{(k)}$  do
3     if  $r \subseteq s_i$  then  $sup(r) \leftarrow sup(r) + 1$ 

EXTENDPREFIXTREE ( $\mathcal{C}^{(k)}$ ):
1 foreach leaf  $r_a \in \mathcal{C}^{(k)}$  do
2   foreach leaf  $r_b \in \text{CHILDREN}(\text{PARENT}(r_a))$  do
3      $r_{ab} \leftarrow r_a + r_b[k]$  // extend  $r_a$  with last item of  $r_b$ 
      // prune if there are any infrequent
      subsequences
4     if  $r_c \in \mathcal{C}^{(k)}$ , for all  $r_c \subset r_{ab}$ , such that  $|r_c| = |r_{ab}| - 1$  then
5       Add  $r_{ab}$  as child of  $r_a$  with  $sup(r_{ab}) \leftarrow 0$ 
6     if no extensions from  $r_a$  then
7       remove  $r_a$ , and all ancestors of  $r_a$  with no extensions, from  $\mathcal{C}^{(k)}$ 
8 return  $\mathcal{C}^{(k)}$ 

```

Sequential Pattern Mining

- Association rules do not consider the order of transactions, however, in many applications ordering is significant
- In market basket analysis, it is interesting to know whether people buy some items in sequence
- Example
 - Buying bed first and then bed sheets later
 - Navigational patterns of users in a Web site from sequences of page visits of users

- **Web sequence**
 - < {Homepage} {Electronics} {Digital Cameras} {Canon Digital Camera} {Shopping Cart} {Order Confirmation} {Return to Shopping} >
- **Sequence of initiating events causing the accident at 3-mile Island:**
 - <{clogged resin} {outlet valve closure} {loss of feedwater} {condenser polisher outlet valve shut} {booster pumps trip} {main waterpump trips} {main turbine trips} {reactor pressure increases}>
- **Sequence of books checked out at a library:**
 - <{Fellowship of the Ring} {The Two Towers} {Return of the King}>

- Let $I = \{i_1, i_2, \dots, i_m\}$ be a set of items.
- A sequence is an ordered list of itemsets.
- We denote a sequence s by $\langle a_1, a_2, \dots, a_r \rangle$, where a_i is an itemset, also called an element of s .
- An element (or an itemset) of a sequence is denoted by $\{x_1, x_2, \dots, x_k\}$, where $x_j \in I$ is an item.
- We assume without loss of generality that items in an element of a sequence are in lexicographic order

- Sequence databases consist of sequences of ordered elements or events (with or without time)
- Sequential Pattern Mining is the mining of frequently occurring ordered events or subsequences as patterns
- Sequences consists of nominal (categorical) data. When sequence are based on numerical data, then we apply time series analysis

- The size of a sequence is the number of elements (or itemsets) in the sequence
- The length of a sequence is the number of items in the sequence (a sequence of length k is called k -sequence)
- Given two sequences, $s_1 = \langle a_1, a_2, \dots, a_r \rangle$ and $s_2 = \langle b_1, b_2, \dots, b_v \rangle$
- s_1 is a subsequence of s_2 if there exist integers

$$1 \leq j_1 < j_2 < \dots < j_r \leq v$$

such that

$$a_1 \subseteq b_{j_1} \wedge a_2 \subseteq b_{j_2} \wedge \dots \wedge a_r \subseteq b_{j_v}$$

- Let $I = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$.
- The sequence $\langle\{3\}\{4, 5\}\{8\}\rangle$ is contained in (or is a subsequence of) $\langle\{6\} \{3, 7\}\{9\}\{4, 5, 8\}\{3, 8\}\rangle$
- In fact, $\{3\} \subseteq \{3, 7\}$, $\{4, 5\} \subseteq \{4, 5, 8\}$, and $\{8\} \subseteq \{3, 8\}$
- However, $\langle\{3\}\{8\}\rangle$ is not contained in $\langle\{3, 8\}\rangle$ or vice versa
- The size of the sequence $\langle\{3\}\{4, 5\}\{8\}\rangle$ is 3, and the length of the sequence is 4.

- The input is a set S of input data sequences (or sequence database)
- The problem of mining sequential patterns is to find all the sequences that have a user-specified minimum support
- Each such sequence is called a frequent sequence, or a sequential pattern
- The support for a sequence is the fraction of total data sequences in S that contains this sequence

- **Itemset**
 - Non-empty set of items
 - Each itemset is mapped to an integer
- **Sequence**
 - Ordered list of itemsets
- **Support for a Sequence**
 - Fraction of total customers that support a sequence.
- **Maximal Sequence**
 - A sequence that is not contained in any other sequence
- **Large Sequence**
 - Sequence that meets minisup

- Given a database $D \{s_1, \dots, s_N\}$ of N sequences and a sequence r , we define the support count of r as

$$\text{sup}(r) = | \{s_i \mid r \text{ is a subsequence of } s_i\} |$$

- And the support (or relative support),

$$\text{rsup}(r) = \text{sup}(r)/N$$

Customer ID	Transaction Time	Items Bought
1	June 25 '93	30
1	June 30 '93	90
2	June 10 '93	10,20
2	June 15 '93	30
2	June 20 '93	40,60,70
3	June 25 '93	30,50,70
4	June 25 '93	30
4	June 30 '93	40,70
4	July 25 '93	90
5	June 12 '93	90

Customer ID	Customer Sequence
1	< (30) (90) >
2	< (10 20) (30) (40 60 70) >
3	< (30) (50) (70) >
4	< (30) (40 70) (90) >
5	< (90) >

Maximal seq with support > 40%

< (30) (90) >

< (30) (40 70) >

Note: Use Minisup of 40%, no less than two customers must support the sequence
 < (10 20) (30) > Does not have enough support (Only by Customer #2)
 < (30) >, < (70) >, < (30) (40) > ... are not maximal.

- Given the following sequence database and a minsup of 2

SID	Sequence
10	<(a)(<u>ab</u> c)(a <u>c</u>)(d)(cf)>
20	<(ad)(c)(bc)(ae)>
30	<(ef)(<u>ab</u>)(df)(<u>c</u>)(b)>
40	<(e)(g)(af)(c)(b)(c)>

- <(a)(bc)(d)(c)> is a subsequence of <(a)(abc)(ac)(d)(cf)>
- <(ab)(c)> is a frequent sequence or a sequential pattern

- Given the sequence database

SID	Sequence
10	<(bd)(c)(b)(ac)>
20	<(bf)(ce)(b)(fg)>
30	<(ah)(bf)(a)(b)(f)>
40	<(be)(ce)(d)>
50	<(a)(bd)(b)(c)(b)(ade)>

- Given a minsup of 2, <(bd)cb> is a sequential pattern (or frequent sequence)

Apriori principle still applies

If a sequence S is not frequent,
then every super-sequence of S is not
frequent

In the previous example,

Mining Sequential Patterns

(Outline of the basic principle)

87

- Level-by-level do
 - Generate candidate sequences
 - Scan database to collect support counts
 - Use Apriori property to prune candidates
- Only generate candidates satisfying Apriori property
- Limitations
 - It can generate a massive amount of candidates (1000 frequent 1-sequences can generate around 1.5 million candidates)
 - It requires many scans of the database
- Accordingly, as for association rules, there are several prefix-based extensions (e.g. PrefixSpan) that deal with all these limitations

Seq. ID	Sequence
10	<(bd)(c)(b)(ac)>
20	<(bf)(ce)(b)(fg)>
30	<(ah)(bf)(a)(b)(f)>
40	<(be)(ce)(d)>
50	<(a)(bd)(b)(c)(b)(ade)>

1st scan - 8 candidates
6 sequential patterns (length 1)

<a> <c> <d> <e> <f> **<g>** <h>

2nd scan - 51 candidates
19 sequential patterns (length 2)
10 candidates were not in the database

<aa> <ab> ... **<af>** <ba> <bb> ... <ff>
<(ab)> ... <(ef)>

3rd scan - 46 candidates
19 sequential patterns (length 3)
20 candidates were not in the database

<abb> <aab> <aba> <baa> **<bab>** ...

Association Rules for Classification

- Association rule mining assumes that the data consist of a set of transactions. Thus, the typical tabular representation of data used in classification must be mapped into such a format.
- Association rule mining is then applied to the new dataset and the search is focused on association rules in which the tail identifies a class label

$$X \Rightarrow c_i \text{ (where } c_i \text{ is a class label)}$$

- The association rules are pruned using the pessimistic error-based method used in C4.5
- Finally, rules are sorted to build the final classifier.

The Weather Dataset

91

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

outlook=overcast ==> play=yes

humidity=normal windy=FALSE ==> play=yes

outlook=rainy windy=FALSE ==> play=yes

outlook=sunny humidity=high ==> play=no

outlook=rainy windy=TRUE ==> play=no

(default class is the majority class)

Run the notebooks for this lecture

Download and run SPMF library

<http://www.philippe-fournier-viger.com/spmf/index.php?link=download.php>