

$$\begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ 0 & & & \sigma_n \end{bmatrix}$$

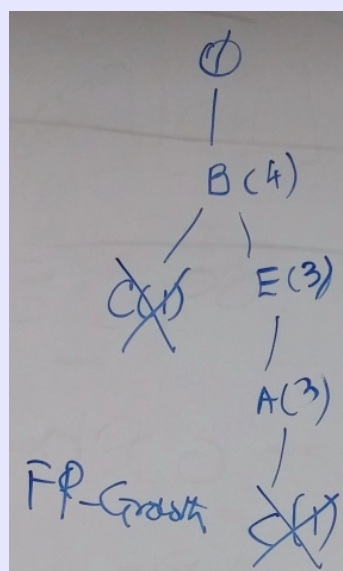
$$X = \sum_{i=1}^{\text{rank}(X)} \sigma_i u_i v_i^T = U \Sigma V^T$$

$\sigma_i$ :  $i^{\text{th}}$  singular value of  $X$   
 $u_i$ :  $i^{\text{th}}$  left singular value of  $X$  ( $i^{\text{th}}$  column of  $U$ )  
 $v_i^T$ :  $i^{\text{th}}$  right singular vector of  $X$  ( $i^{\text{th}}$  column of  $V^T$ )

Captures the patterns among attributes  
 Captures the patterns among the objects

CS 422: Data Mining  
 Vijay K. Gurbani, Ph.D.,  
 Illinois Institute of Technology

## Association Analysis (Rules)



CS 422  
 vgurbani@iit.edu



# Association Rule Mining

- Goal of Association Rule Mining: Given a set of transactions,  $T$ , find all rules having:
  - support  $\geq \textit{minsup}$
  - confidence  $\geq \textit{minconf}$
- How do we get there?
- Two steps:
  - Frequent itemset generation: find all items that satisfy *minsup* threshold (frequent itemsets). (Is computationally expensive!!)
  - Rule generation: extract all high-confidence rules from the frequent itemsets (strong rules).

# Frequent Itemset Generation: The Apriori Approach

- We looked at the brute-force approach to generate itemsets. Can we do better?
- Brute-force approach wastes computations because many of the candidates that it generates may not be frequent.
- Instead, note:

Let  $X$  and  $Y$  be two itemsets  $\in \mathcal{I}$  such that  $X \subseteq Y$ .

If so, then  $\text{sup}(X) \geq \text{sup}(Y)$

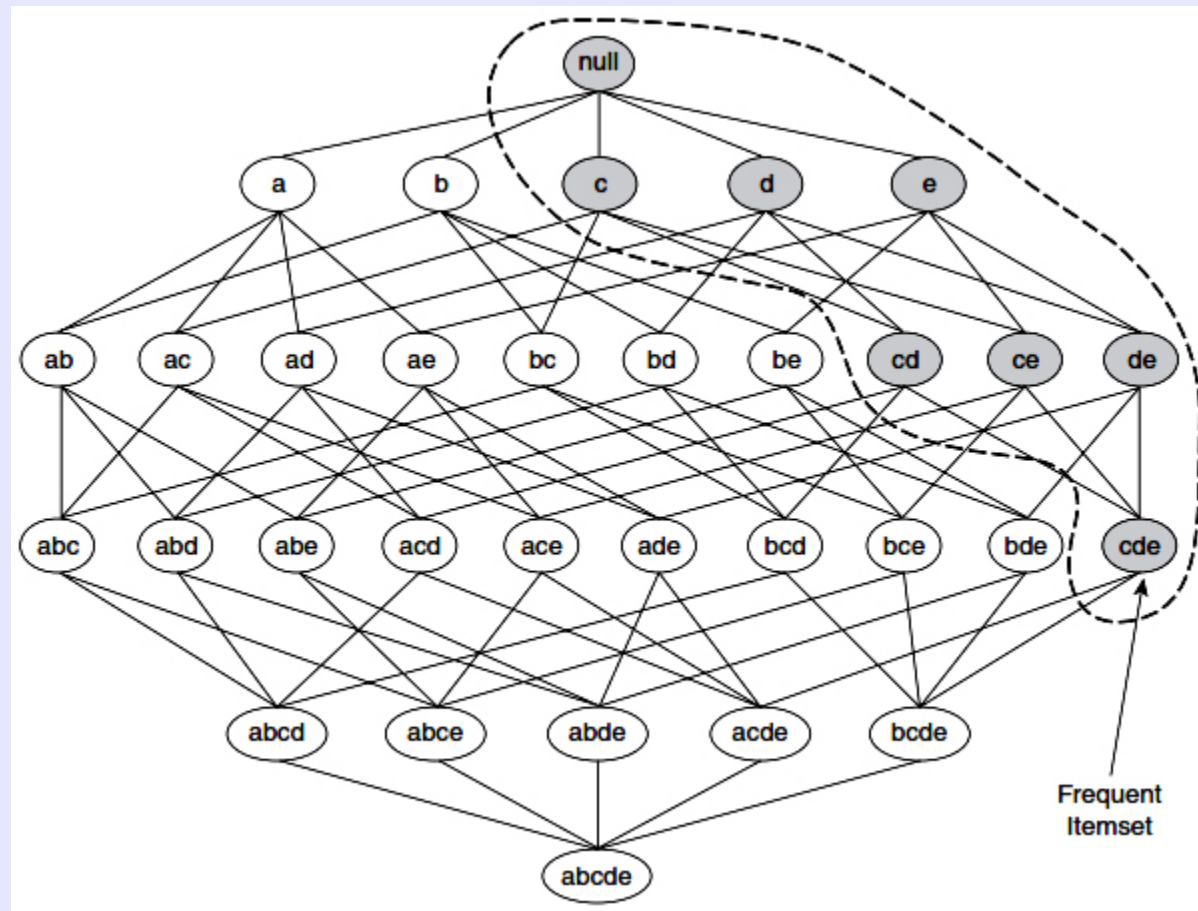
E.g.  $X=ABCD$ ,  $Y=ABCDE$ , then  $\text{sup}(ABCD) \geq \text{sup}(ABCDE)$ .

- This leads to ...

# Frequent Itemset Generation: The Apriori Approach

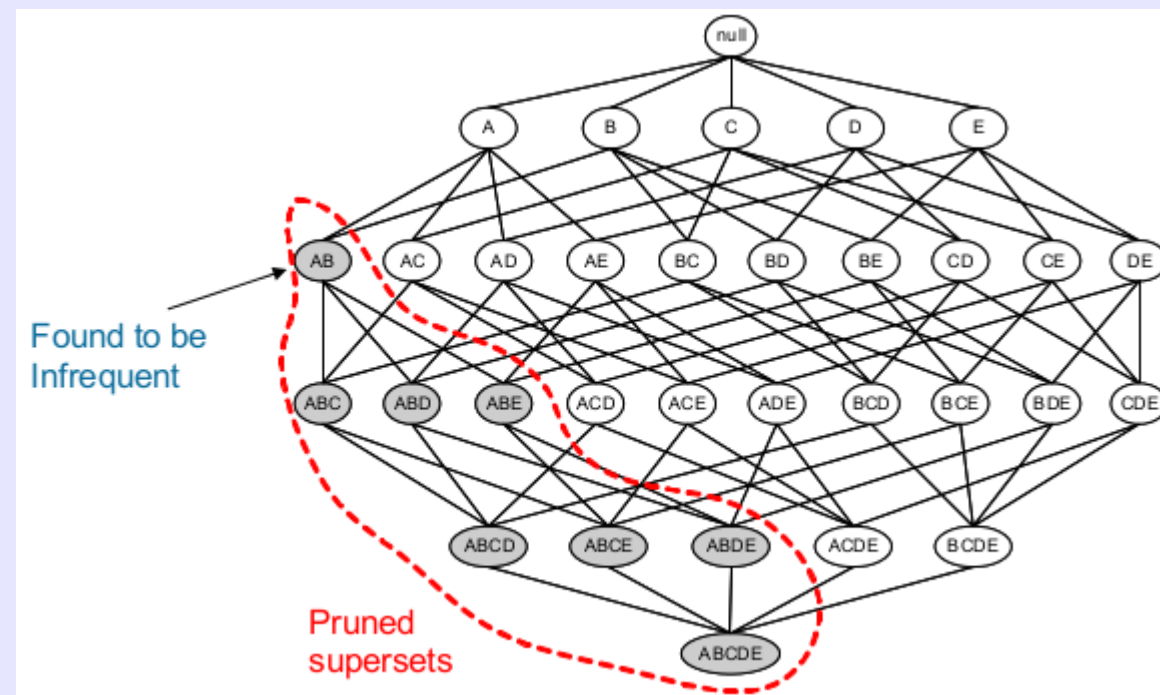
- The *Apriori* principle:

- If an itemset is frequent, then all of its subsets must be frequent as well.

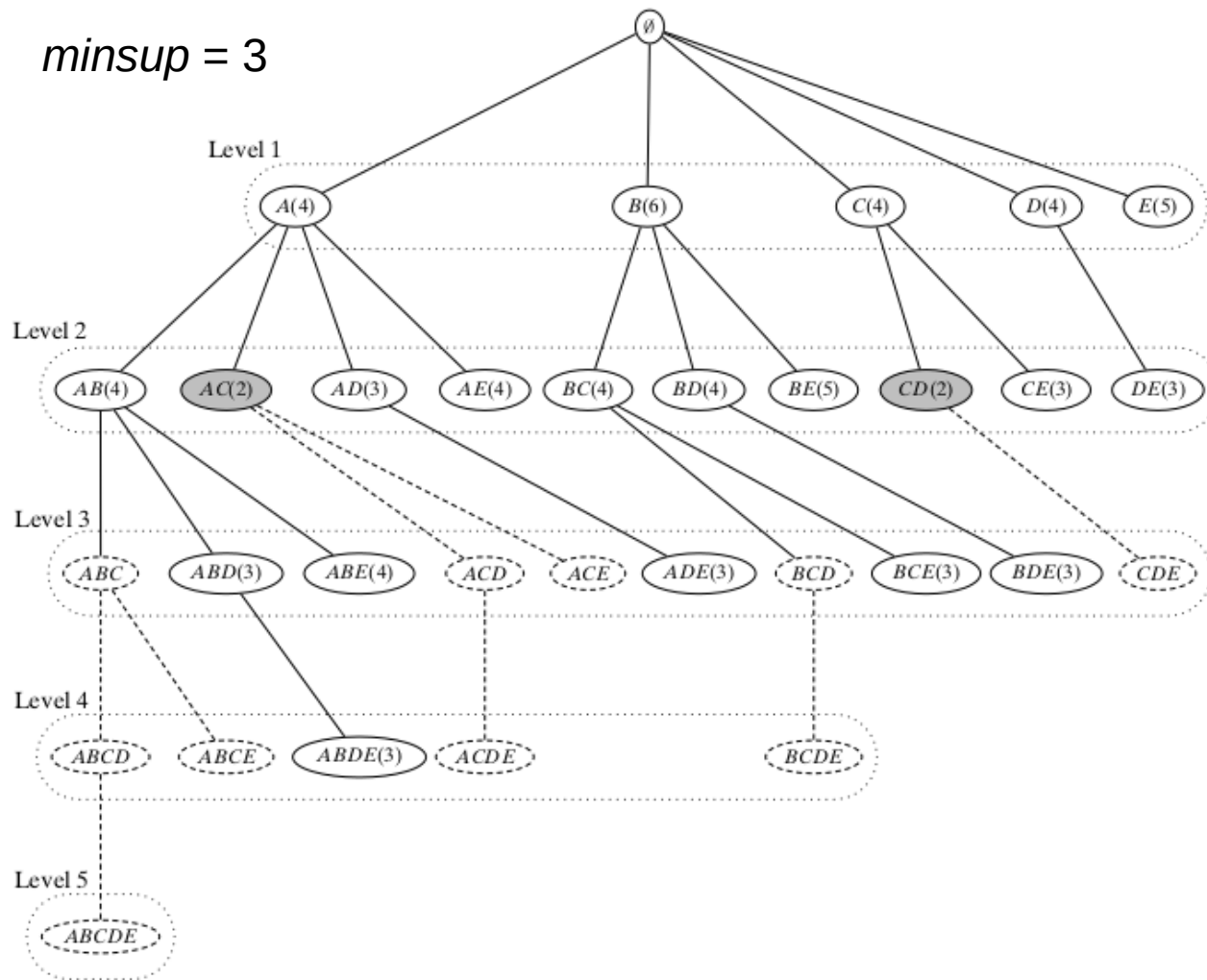


# Frequent Itemset Generation: The Apriori Approach

- The *Apriori* principle:
  - Conversely, if an itemset is infrequent, then all of its supersets must be infrequent as well.



# Frequent Itemset Generation: The Apriori Approach



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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with  $minsup = 3$

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

Figure 8.3. Apriori: prefix search tree and effect of pruning. Shaded nodes indicate infrequent itemsets, whereas dashed nodes and lines indicate all of the pruned nodes and branches. Solid lines indicate frequent itemsets.



# Frequent Itemset Generation: The Apriori Approach

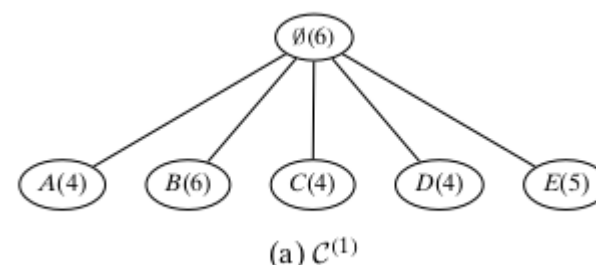
D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

(a) Binary database

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

(b) Transaction database

minsup = 3



## ALGORITHM 8.2. Algorithm APRIORI

**APRIORI (D, I, 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 ( $\mathcal{C}^{(k)}$ , 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 ( $\mathcal{C}^{(k)}$ , D):**

```

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

```

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

```

16 foreach leaf  $X_a \in \mathcal{C}^{(k)}$  do
17   foreach leaf  $X_b \in \text{SIBLING}(X_a)$ , such that  $b > a$  do
18      $X_{ab} \leftarrow X_a \cup X_b$ 
19     // prune candidate if there are any infrequent subsets
20     if  $X_j \in \mathcal{C}^{(k)}$ , for all  $X_j \subset X_{ab}$ , such that  $|X_j| = |X_{ab}| - 1$  then
21       Add  $X_{ab}$  as child of  $X_a$  with  $\text{sup}(X_{ab}) \leftarrow 0$ 
22   if no extensions from  $X_a$  then
23     remove  $X_a$ , and all ancestors of  $X_a$  with no extensions, from  $\mathcal{C}^{(k)}$ 
23 return  $\mathcal{C}^{(k)}$ 

```

# Frequent Itemset Generation: The Apriori Approach

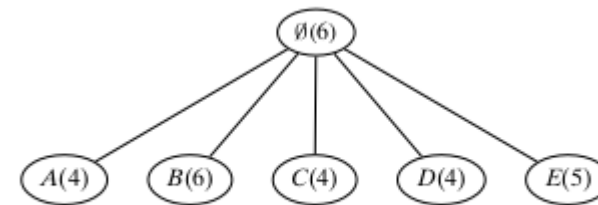
D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

(a) Binary database

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

(b) Transaction database

minsup = 3



(a)  $C^{(1)}$

We now extend the prefix tree from Level k to Level k+1: given two frequent k-itemsets ( $X_a$  and  $X_b$ ), with common k-1 length prefix (i.e., two siblings with common parent), we generate (k+1) length candidates

$$X_{ab} = X_a \cup X_b.$$

-  $X_{ab}$  retained only if it has no infrequent subset.

## ALGORITHM 8.2. Algorithm APRIORI

**APRIORI (D, I, 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 ( $\mathcal{C}^{(k)}, \mathbf{D}$ ):**

```

13 foreach  $(t, \mathbf{i}(t)) \in \mathbf{D}$  do
14   foreach  $k$ -subset  $X \subseteq \mathbf{i}(t)$  do
15     if  $X \in \mathcal{C}^{(k)}$  then  $\text{sup}(X) \leftarrow \text{sup}(X) + 1$ 
  
```

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

```

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



# Frequent Itemset Generation: The Apriori Approach

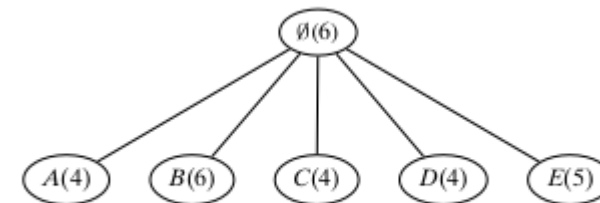
D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

(a) Binary database

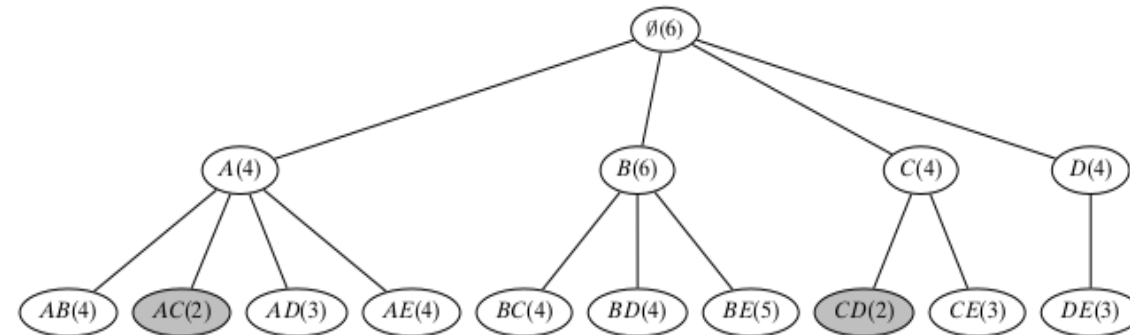
t	i(t)
1	ABDE
2	BCE
3	ABDE
4	ABCE
5	ABCDE
6	BCD

(b) Transaction database

minsup = 3



(a)  $\mathcal{C}^{(1)}$



(b)  $\mathcal{C}^{(2)}$

## ALGORITHM 8.2. Algorithm APRIORI

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

```

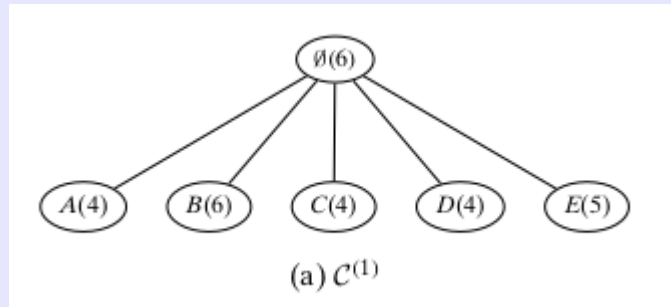
13 foreach  $(t, \mathbf{i}(t)) \in \mathbf{D}$  do
14   foreach  $k$ -subset  $X \subseteq \mathbf{i}(t)$  do
15     if  $X \in \mathcal{C}^{(k)}$  then  $\text{sup}(X) \leftarrow \text{sup}(X) + 1$ 
  
```

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

```

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

# Frequent Itemset Generation: The Apriori Approach

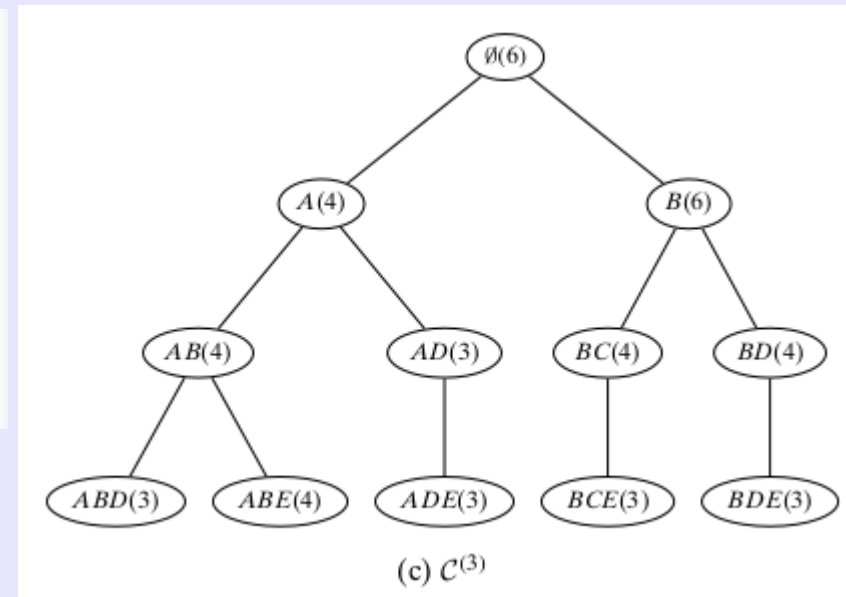
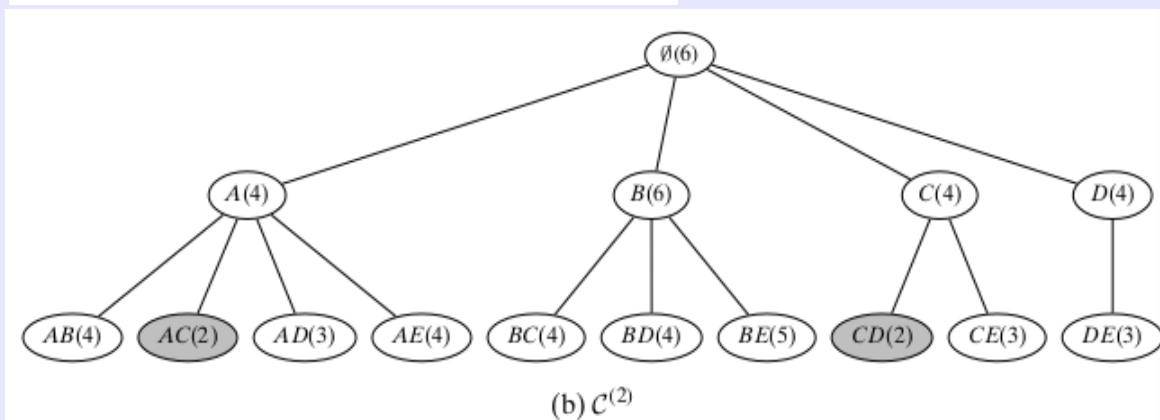


D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

(a) Binary database

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

(b) Transaction database



# Frequent Itemset Generation: The Apriori Approach

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

(a) Binary database

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

(b) Transaction database

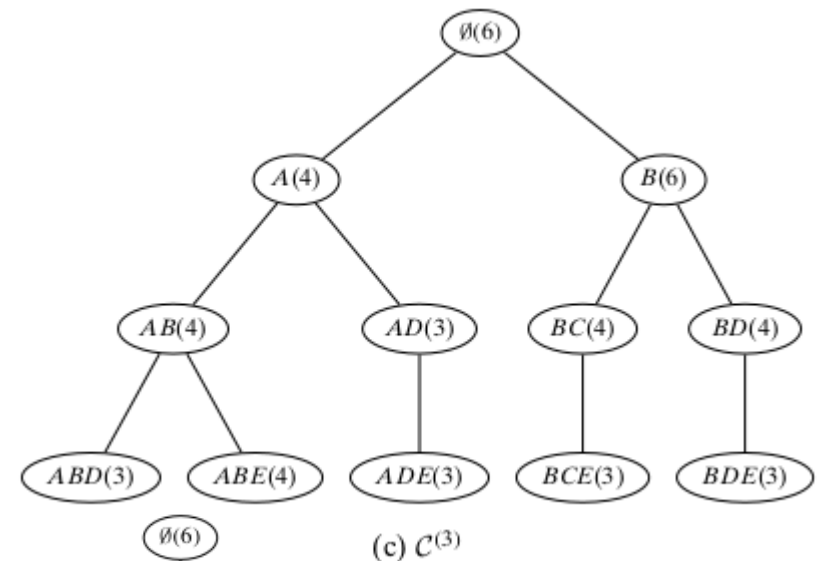
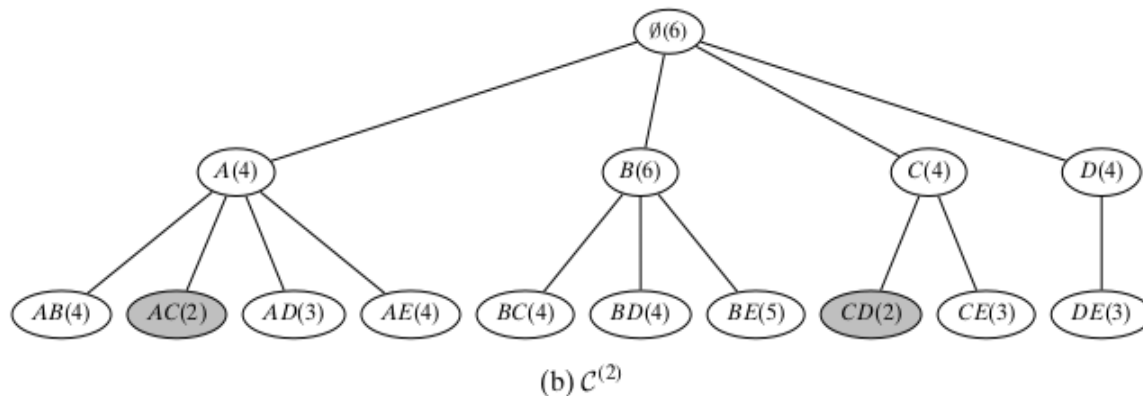
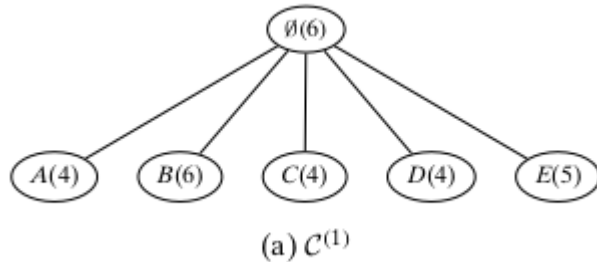
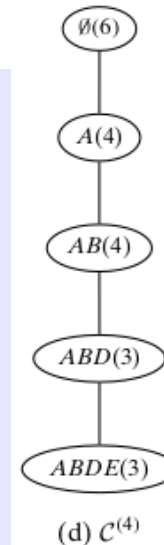


Table 8.1. Frequent itemsets with  $minsup = 3$

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



# Frequent Itemset Generation: The Apriori Approach

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

(a) Binary database

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

(b) Transaction database

## ALGORITHM 8.2. Algorithm APRIORI

```

APRIORI (D,  $\mathcal{I}$ , 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 ( $\mathcal{C}^{(k)}$ , 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 ( $\mathcal{C}^{(k)}$ , D):
13 foreach  $(t, \mathbf{i}(t)) \in \mathbf{D}$  do
14   foreach  $k$ -subset  $X \subseteq \mathbf{i}(t)$  do
15     if  $X \in \mathcal{C}^{(k)}$  then  $\text{sup}(X) \leftarrow \text{sup}(X) + 1$ 

EXTENDPREFIXTREE ( $\mathcal{C}^{(k)}$ ):
16 foreach leaf  $X_a \in \mathcal{C}^{(k)}$  do
17   foreach leaf  $X_b \in \text{SIBLING}(X_a)$ , such that  $b > a$  do
18      $X_{ab} \leftarrow X_a \cup X_b$ 
19     // prune candidate if there are any infrequent subsets
20     if  $X_j \in \mathcal{C}^{(k)}$ , for all  $X_j \subset X_{ab}$ , such that  $|X_j| = |X_{ab}| - 1$  then
21       Add  $X_{ab}$  as child of  $X_a$  with  $\text{sup}(X_{ab}) \leftarrow 0$ 
22   if no extensions from  $X_a$  then
23     remove  $X_a$ , and all ancestors of  $X_a$  with no extensions, from  $\mathcal{C}^{(k)}$ 
23 return  $\mathcal{C}^{(k)}$ 

```

Worst case complexity:

Complexity of Apriori:  $\mathcal{O}(|\mathcal{I}| * D * 2^{|\mathcal{I}|})$  as all itemsets may be frequent. In practice, much lower due to pruning.

I/O costs are much lower, to the tune of  $\mathcal{O}(|\mathcal{I}|)$  database scans as opposed to  $\mathcal{O}(2^{|\mathcal{I}|})$  scans for brute-force. In practice, the algorithm only requires  $l$  database scans, where  $l$  is the length of the longest frequent itemset.

# Association Rule Mining

- Rule mining: Preliminaries and foundation

## Association Rules

An *association rule* is an expression  $X \xrightarrow{s,c} Y$ , where  $X$  and  $Y$  are itemsets and they are disjoint, that is,  $X, Y \subseteq \mathcal{I}$ , and  $X \cap Y = \emptyset$ . Let the itemset  $X \cup Y$  be denoted as  $XY$ . The *support* of the rule is the number of transactions in which both  $X$  and  $Y$  co-occur as subsets:

$$s = \text{sup}(X \longrightarrow Y) = |\mathbf{t}(XY)| = \text{sup}(XY)$$

The *relative support* of the rule is defined as the fraction of transactions where  $X$  and  $Y$  co-occur, and it provides an estimate of the joint probability of  $X$  and  $Y$ :

$$\text{rsup}(X \longrightarrow Y) = \frac{\text{sup}(XY)}{|\mathbf{D}|} = P(X \wedge Y)$$

The *confidence* of a rule is the conditional probability that a transaction contains  $Y$  given that it contains  $X$ :

$$c = \text{conf}(X \longrightarrow Y) = P(Y|X) = \frac{P(X \wedge Y)}{P(X)} = \frac{\text{sup}(XY)}{\text{sup}(X)}$$

A rule is *frequent* if the itemset  $XY$  is frequent, that is,  $\text{sup}(XY) \geq \text{minsup}$  and a rule is *strong* if  $\text{conf} \geq \text{minconf}$ , where *minconf* is a user-specified minimum confidence threshold.

# Association Rule Mining

- Rule mining: Preliminaries and foundation

## Association Rules

An *association rule* is an expression  $X \xrightarrow{s,c} Y$ , where  $X$  and  $Y$  are itemsets and they are disjoint, that is,  $X, Y \subseteq \mathcal{I}$ , and  $X \cap Y = \emptyset$ . Let the itemset  $X \cup Y$  be denoted as  $XY$ . The *support* of the rule is the number of transactions in which both  $X$  and  $Y$  co-occur as subsets:

$$s = \text{sup}(X \longrightarrow Y) = |\mathbf{t}(XY)| = \text{sup}(XY)$$

The *relative support* of the rule is defined as the fraction of transactions where  $X$  and  $Y$  co-occur, and it provides an estimate of the joint probability of  $X$  and  $Y$ :

$$\text{rsup}(X \longrightarrow Y) = \frac{\text{sup}(XY)}{|\mathbf{D}|} = P(X \wedge Y)$$

The *confidence* of a rule is the conditional probability that a transaction contains  $Y$  given that it contains  $X$ :

$$c = \text{conf}(X \longrightarrow Y) = P(Y|X) = \frac{P(X \wedge Y)}{P(X)} = \frac{\text{sup}(XY)}{\text{sup}(X)}$$

A rule is **frequent** if the itemset  $XY$  is frequent, that is,  $\text{sup}(XY) \geq \text{minsup}$  and a rule is **strong** if  $\text{conf} \geq \text{minconf}$ , where *minconf* is a user-specified minimum confidence threshold.

Table 8.1. Frequent itemsets with *minsup* = 3

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

Example:  $BC \rightarrow E$ .

$$\text{sup}(BC \rightarrow E) = \text{sup}(BCE) = 3$$

$$\begin{aligned} \text{conf}(BC \rightarrow E) &= \frac{\text{sup}(BCE)}{\text{sup}(BC)} \\ &= \frac{3}{4} = 0.75 \end{aligned}$$



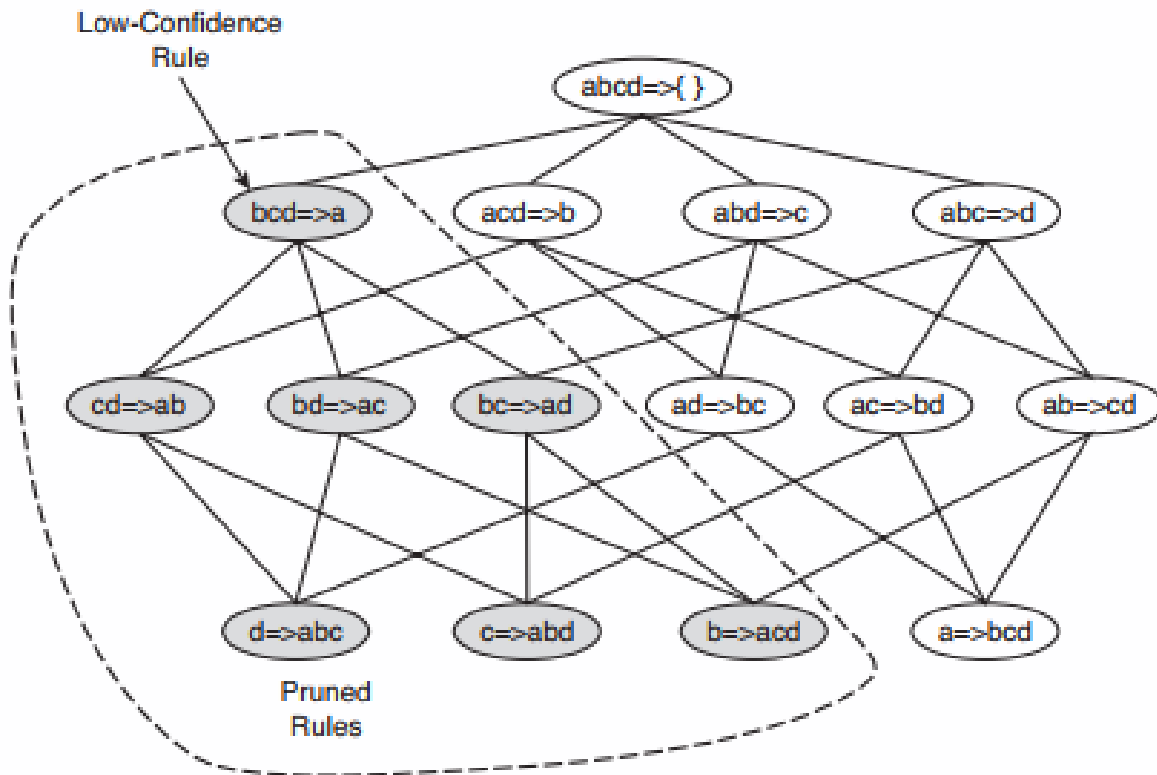
# Generating association rules

- Rules are pruned using confidence.

Recall:  $c = \text{conf}(X \rightarrow Y) = \frac{\text{sup}(XY)}{\text{sup}(X)}$

- Unlike support, confidence does not exhibit the monotone property ... but ...
- If we consider rules generated from the same frequent itemset  $Y$ , then...
- Theorem:** If a rule  $X \rightarrow Y \setminus X$  does not satisfy the confidence threshold, then any rule  $X' \rightarrow Y \setminus X'$ , where  $X' \subset X$ , must not satisfy the confidence threshold as well.

# Generating association rules



- Theorem:** Assume rules are generated from the same frequent itemset. Then, if a rule  $X \rightarrow Y \setminus X$  does not satisfy the confidence threshold, then any rule  $X' \rightarrow Y \setminus X'$ , where  $X' \subset X$ , must not satisfy the confidence threshold as well.

Example:  $Y = \{a, b, c, d\}$ ,  $X = \{b, c, d\}$ ;  
 $X \rightarrow Y \setminus X \Rightarrow$

$\{b, c, d\} \rightarrow \{a, b, c, d\} \setminus \{b, c, d\} \Rightarrow$

$\{b, c, d\} \rightarrow \{a\}$

If  $\{b, c, d\} \rightarrow \{a\}$  does not satisfy the confidence threshold, then any  $X' \subset X$  must not satisfy the confidence threshold as well.

# Generating association rules

- Empirical example for the theorem:

Transaction list: {A,B,C}, {A,B,D}, {A,C}. Below,  $\sigma$  implies support count, and  $Y - X$  is the same as  $Y \setminus X$ .

$$\sigma(\{A\}) = 3$$

$$\sigma(\{A, B\}) = 2$$

$$\sigma(\{A, B, C\}) = 1$$

$$\sigma(\{A, C\}) = 2$$

$$\sigma(\{B\}) = 2$$

To begin with,  $Y = \{A, B, C\}$ ,  $X = \{\}$ , so we have the rule  $X \rightarrow Y$ , or  $\{\} \rightarrow \{A, B, C\}$ .

Let the antecedent  $X = \{A, B\}$ , then  $Y - X = \{C\}$ .

The confidence of the rule  $X \rightarrow Y - X$ , or  $C(X \rightarrow Y - X) = \sigma(X \cup (Y - X)) / \sigma(X) = \sigma(\{A, B, C\}) / \sigma(\{A, B\}) = 1/2 = 0.5$ .

Let  $X'$  be a subset of  $X$ , ( $X' \subset X$ ). Now,  $C(X' \rightarrow Y - X') \leq C(X \rightarrow Y - X)$ , or  $\leq 0.5$  (Theorem 6.2, Tan).

Let  $X' = \{A\}$ , then  $Y - X' = \{B, C\}$ .

The confidence of the rule  $X' \rightarrow Y - X' = C(X' \rightarrow Y - X') = \sigma(X' \cup (Y - X')) / \sigma(X') = \sigma(\{A, B, C\}) / \sigma(\{A\}) = 1/3 = 0.3$ .

If you let  $X' = \{B\}$  and  $Y - X' = \{B, C\}$ , you should also get that Theorem 6.2 holds.

# Generating association rules

- Consider the frequent itemset ABDE(3).

- Assume  $minconf = 0.9$ .

- Recall:  $c = conf(X \rightarrow Y) = \frac{sup(XY)}{sup(X)}$

- Step 1: Let  $Z = \{ABDE\}$ , and enumerate all proper subsets  $X \subset Z$ :  
 $A = \{ABD(3), ABE(4), ADE(3), BDE(3), AB(4), AD(3), AE(4), BD(4), BE(5), DE(3), A(4), B(6), D(4), E(5)\}$

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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with  $minsup = 3$

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

# Generating association rules

- $A = \{ABD(3), ABE(4), ADE(3), BDE(3), AB(4), AD(3), AE(4), BD(4), BE(5), DE(3), A(4), B(6), D(4), E(5)\}$

- Step 2: Take the first subset,  $X = \{ABD\}$ , which becomes the antecedent. The consequent is the missing element(s),  $Y = \{E\}$ .

$$\text{conf}(X \rightarrow Y) = \text{sup}(ABDE) / \text{sup}(ABD) = 3/3 = 1.0$$

- Calculated confidence  $\geq \text{minconf}$ . (Recall,  $\text{minconf}=0.9$ )

Output rule:  $\{ABD\} \rightarrow \{E\}$ ,  $\text{conf} = 1.0$

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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with  $\text{minsup} = 3$

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

$$\text{Recall: } c = \text{conf}(X \rightarrow Y) = \frac{\text{sup}(XY)}{\text{sup}(X)}$$

# Generating association rules

- $\mathbf{A} = \{ABE(4), ADE(3), BDE(3), AB(4), AD(3), AE(4), BD(4), BE(5), DE(3), A(4), B(6), D(4), E(5)\}$
- Step 2: Take the first subset,  $X = \{ABE\}$ , which becomes the antecedent. The consequent is the missing element(s),  $Y = \{D\}$ .

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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with  $minsup = 3$

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

$$\text{Recall: } c = \text{conf}(X \rightarrow Y) = \frac{\text{sup}(XY)}{\text{sup}(X)}$$

$$\text{conf}(X \rightarrow Y) = \text{sup}(ABED)/\text{sup}(ABE) = 3/4 = 0.75$$

- Calculated confidence  $< minconf$ .
  - REJECT rule  $\{ABE\} \rightarrow \{D\}$ ,  $\text{conf} = 0.75$ .
- Because we rejected  $ABE \rightarrow D$ , we can remove from  $\mathbf{A}$  all subsets of ABE (the antecedent).



# Generating association rules

- $A = \{ABE(4), ADE(3), BDE(3), AB(4), AD(3), AE(4), BD(4), BE(5), DE(3), A(4), B(6), D(4), E(5)\}$
- Because we rejected  $ABE \rightarrow D$ , we can remove from  $A$  all subsets of  $ABE$  (the antecedent).

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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with  $minsup = 3$

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

$$\text{Recall: } c = \text{conf}(X \rightarrow Y) = \frac{\text{sup}(XY)}{\text{sup}(X)}$$

~~$A = \{ADE(3), BDE(3), AB(4), AD(3), AE(4), BD(4), BE(5), DE(3), A(4), B(6), D(4), E(5)\}$~~

$A = \{ADE(3), BDE(3), AD(3), BD(4), DE(3), D(4)\}$

# Generating association rules

- $A = \{ADE(3), BDE(3), AD(3), BD(4), DE(3), D(4)\}$
- Step 2: Take the first subset,  $X = \{ADE\}$ , which becomes the antecedent. The consequent is the missing element(s),  $Y = \{B\}$ .

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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with  $minsup = 3$

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

Recall:  $c = conf(X \rightarrow Y) = \frac{sup(XY)}{sup(X)}$

$$conf(X \rightarrow Y) = sup(ADEB)/sup(ADE) = 3/3 = 1.0$$

- Calculated confidence  $\geq minconf$ .
  - Output rule  $\{ADE\} \rightarrow \{B\}$ ,  $conf = 1.0$ .

# Generating association rules

- $A = \{BDE(3), AD(3), BD(4), DE(3), D(4)\}$
- Step 2: Take the first subset,  $X = \{BDE\}$ , which becomes the antecedent. The consequent is the missing element(s),  $Y = \{A\}$ .

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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with  $minsup = 3$

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

Recall:  $c = conf(X \rightarrow Y) = \frac{sup(XY)}{sup(X)}$

$$conf(X \rightarrow Y) = sup(BDEA)/sup(BDE) = 3/3 = 1.0$$

- Calculated confidence  $\geq minconf$ .
  - Output rule  $\{BDE\} \rightarrow \{A\}$ ,  $conf = 1.0$ .

# Generating association rules

- $A = \{AD(3), BD(4), DE(3), D(4)\}$
- Step 2: Take the first subset,  $X = \{AD\}$ , which becomes the antecedent. The consequent is the missing element(s),  $Y = \{BE\}$ .

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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with  $minsup = 3$

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

Recall:  $c = conf(X \rightarrow Y) = \frac{sup(XY)}{sup(X)}$

$$conf(X \rightarrow Y) = sup(ADBE)/sup(AD) = 3/3 = 1.0$$

- Calculated confidence  $\geq minconf$ .
  - Output rule  $\{AD\} \rightarrow \{BE\}$ ,  $conf = 1.0$ .

# Generating association rules

- $A = \{BD(4), DE(3), D(4)\}$
- Step 2: Take the first subset,  $X = \{BD\}$ , which becomes the antecedent. The consequent is the missing element(s),  $Y = \{AE\}$ .

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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with  $minsup = 3$

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

$$\text{Recall: } c = \text{conf}(X \rightarrow Y) = \frac{\text{sup}(XY)}{\text{sup}(X)}$$

$$\text{conf}(X \rightarrow Y) = \text{sup}(BDAE)/\text{sup}(BD) = 3/4 = 0.75$$

- Calculated confidence  $< minconf$ .
  - REJECT rule  $\{BD\} \rightarrow \{AE\}$ ,  $\text{conf} = 0.75$ .
- Because we rejected  $BD \rightarrow AE$ , we can remove from  $A$  all subsets of  $BD$  (the antecedent).

# Generating association rules

- $A = \{DE(3)\}$
- Step 2: Take the remaining subset,  $X = \{DE\}$ , which becomes the antecedent. The consequent is the missing element(s),  $Y = \{AB\}$ .

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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with  $minsup = 3$

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

Recall:  $c = conf(X \rightarrow Y) = \frac{sup(XY)}{sup(X)}$

$$conf(X \rightarrow Y) = sup(DEAB)/sup(DE) = 3/3 = 1.0$$

- Calculated confidence  $\geq minconf$ .
  - Output rule  $\{DE\} \rightarrow \{AB\}$ ,  $conf = 1.0$ .



# Generating association rules

- We perform Steps 1-2 for **all the frequent itemsets** meeting a *minsup* to derive all of the rules!

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

(a) Binary database

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

(b) Transaction database

Table 8.1. Frequent itemsets with *minsup* = 3

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

$$\text{Recall: } c = \text{conf}(X \rightarrow Y) = \frac{\text{sup}(XY)}{\text{sup}(X)}$$

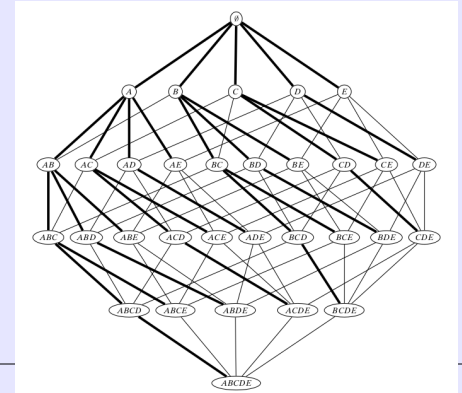
- Question: What about the null rule? That is:  $\{\text{null}\} \rightarrow \{\text{ABDE}\}$ ?

# How many itemsets and rules?

How many itemsets?

$k$  items generate up to  $2^k$  frequent itemsets.

Number of itemsets for  $k$  items =  
 $\binom{k}{1} + \binom{k}{2} + \dots + \binom{k}{k} = 2^k - 1$ ,  
 excluding the null set; with the null set,  $2^k$



How many rules?

For a 3 item dataset {a,b,c}, the number of candidate rules will be:

$d$  = No. of items

For  $d = 3$ ,  $R = 12$

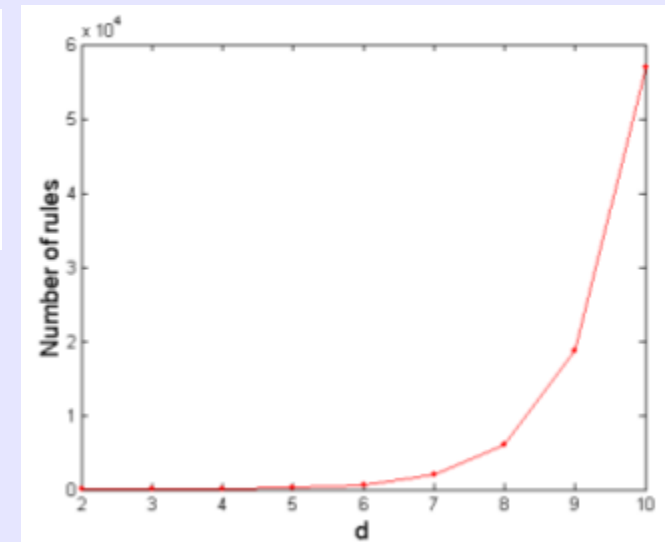
For  $d = 6$ ,  $R = 602$

$$R = \sum_{k=1}^{d-1} \left[ \binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right]$$

$$= 3^d - 2^{d+1} + 1$$

The actual rules are shown below:

```
{ }      => { C }
{ }      => { B }
{ }      => { A }
{ C }    => { B }
{ B }    => { C }
{ C }    => { A }
{ A }    => { C }
{ B }    => { A }
{ A }    => { B }
{ B, C } => { A }
{ A, C } => { B }
{ A, B } => { C }
```



# Limitations of support and confidence

- What value should we pick for support (*minsup*)?
  - If *minsup* too high (0.20), we may miss interesting low-support items. Such items may correspond to expensive products (jewelry) that are seldom purchased by customers, but whose patterns are interesting to mine for the retailer.
  - If *minsup* is too low (0.001), we get information overload: too many frequent itemsets and too many spurious rules.
    - While {milk} → {diapers} is interesting, {milk} → {eggs} is not.

# Limitations of support and confidence

- Drawback of confidence is more subtle.
- Some high confidence rules can be misleading.
- Example:

	<i>Coffee</i>	$\overline{Coffee}$	
<i>Tea</i>	150	50	200
$\overline{Tea}$	650	150	800
	800	200	1000

- Evaluate  $\{Tea\} \rightarrow \{Coffee\}$ .
  - $\text{sup}(\{Tea, Coffee\}) = 150$ ;  $\text{rsup}(\{Tea, Coffee\}) = 150/1000 = 0.15$  (15%).
  - $\text{conf}(\{Tea\} \rightarrow \{Coffee\}) = \text{sup}(\{Tea, Coffee\})/\text{sup}(\{Tea\}) = \text{?????}$

# Limitations of support and confidence

- Drawback of confidence is more subtle.
- Some high confidence rules can be misleading.
- Example:

	<i>Coffee</i>	$\overline{Coffee}$	
<i>Tea</i>	150	50	200
$\overline{Tea}$	650	150	800
	800	200	1000

- Evaluate  $\{Tea\} \rightarrow \{Coffee\}$ .
  - $\text{sup}(\{Tea, Coffee\}) = 150$ ;  $\text{rsup}(\{Tea, Coffee\}) = 150/1000 = 0.15$  (15%).
  - $\text{conf}(\{Tea\} \rightarrow \{Coffee\}) = \text{sup}(\{Tea, Coffee\})/\text{sup}(\{Tea\}) = 150/200 = 0.75$ .
  - The rule  $\{Tea\} \rightarrow \{Coffee\}$  appears to be robust. But is it really?

# Limitations of support and confidence

- Drawback of confidence is more subtle. Example:

	<i>Coffee</i>	$\overline{Coffee}$	
<i>Tea</i>	150	50	200
$\overline{Tea}$	650	150	800
	800	200	1000

rsup = 0.15  
conf = 0.75

- Evaluate  $\{Tea\} \rightarrow \{Coffee\}$ .
  - Fraction of people who like coffee, regardless of if they like tea = 0.80.
  - Fraction of people who like tea and also like coffee = 0.75.
  - Thus knowing that a person likes tea decreases her probability of liking coffee!
    - The rule  $\{Tea\} \rightarrow \{Coffee\}$  is, therefore, misleading.



# Limitations of support and confidence

- Previous example shows that high-confidence rules can sometimes be misleading.
  - Confidence measure ignores the support of the itemset appearing in the rule consequent.

$$\text{Recall: } c = \text{conf}(X \rightarrow Y) = \frac{\text{sup}(XY)}{\text{sup}(X)}$$
  - A new metric called **lift** accounts for the consequent and is defined as:  $\text{lift}(X \rightarrow Y) = \frac{\text{conf}(X \rightarrow Y)}{r\text{sup}(Y)}$
  - By this measure:  
 $\text{lift}(\{\text{Tea}\} \rightarrow \{\text{Coffee}\}) = 0.75/0.80 = 0.94$
  - Value of lift close to 1 implies that the support of the rule is expected; we look for values  $> 1$  (above expectation) and  $< 1$  (below expectation).
  - The rule  $\{\text{Tea}\} \rightarrow \{\text{Coffee}\}$  had a high confidence, but this was an aberration until we saw the lift associated with that rule. Lift normalizes the confidence of the rule using the support of the consequent (coffee). If the support of the consequent is high, i.e., the consequent appears many times in the *tidset*, then the lift will be low.

# Evaluating association rules

- Many to choose from!
- For now, focus on support, confidence, and lift.

#	Measure	Formula
1	$\phi$ -coefficient	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
2	Goodman-Kruskal's ( $\lambda$ )	$\frac{\sum_j \max_k P(A_j, B_k) + \sum_k \max_j P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
3	Odds ratio ( $\alpha$ )	$\frac{P(A,B)P(\bar{A},\bar{B})}{P(A,\bar{B})P(\bar{A},B)}$
4	Yule's $Q$	$\frac{P(A,B)P(\bar{A}\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A}\bar{B}) + P(A,\bar{B})P(\bar{A},B)} = \frac{\alpha - 1}{\alpha + 1}$
5	Yule's $Y$	$\frac{\sqrt{P(A,B)P(\bar{A}\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A}\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}} = \frac{\sqrt{\alpha} - 1}{\sqrt{\alpha} + 1}$
6	Kappa ( $\kappa$ )	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
7	Mutual Information ( $M$ )	$\frac{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\min(-\sum_i P(A_i) \log P(A_i), -\sum_j P(B_j) \log P(B_j))}$
8	J-Measure ( $J$ )	$\max \left( P(A, B) \log \left( \frac{P(B A)}{P(B)} \right) + P(\bar{A}\bar{B}) \log \left( \frac{P(\bar{B} \bar{A})}{P(\bar{B})} \right), \right. \\ \left. P(A, B) \log \left( \frac{P(A B)}{P(A)} \right) + P(\bar{A}\bar{B}) \log \left( \frac{P(\bar{A} \bar{B})}{P(\bar{A})} \right) \right)$
9	Gini index ( $G$ )	$\max \left( P(A)[P(B A)^2 + P(\bar{B} A)^2] + P(\bar{A})[P(B \bar{A})^2 + P(\bar{B} \bar{A})^2] \right. \\ \left. - P(B)^2 - P(\bar{B})^2, \right. \\ \left. P(B)[P(A B)^2 + P(\bar{A} B)^2] + P(\bar{B})[P(A \bar{B})^2 + P(\bar{A} \bar{B})^2] \right. \\ \left. - P(A)^2 - P(\bar{A})^2 \right)$
10	Support ( $s$ )	$P(A, B)$
11	Confidence ( $c$ )	$\max(P(B A), P(A B))$
12	Laplace ( $L$ )	$\max \left( \frac{NP(A,B)+1}{NP(A)+3}, \frac{NP(A,B)+1}{NP(B)+3} \right)$
13	Conviction ( $V$ )	$\max \left( \frac{P(A)P(\bar{B})}{P(A,B)}, \frac{P(B)P(\bar{A})}{P(A,B)} \right)$
14	Interest ( $I$ )	$\frac{P(A,B)}{P(A)P(B)}$
15	cosine ( $IS$ )	$\frac{P(A,B)}{\sqrt{P(A)P(B)}}$
16	Piatetsky-Shapiro's ( $PS$ )	$P(A, B) - P(A)P(B)$
17	Certainty factor ( $F$ )	$\max \left( \frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)} \right)$
18	Added Value ( $AV$ )	$\max(P(B A) - P(B), P(A B) - P(A))$
19	Collective strength ( $S$ )	$\frac{P(A,B) + P(\bar{A}\bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A,B) - P(\bar{A}\bar{B})}$
20	Jaccard ( $\zeta$ )	$\frac{P(A,B)}{P(A) + P(B) - P(A,B)}$
21	Kloggen ( $K$ )	$\sqrt{P(A,B) \max(P(B A) - P(B), P(A B) - P(A))}$