

---

# Kuliah 8 dan 9

## Data Mining

### Association Analysis: Basic Concepts and Algorithms

Lecture Notes for Chapter 6  
Introduction to Data Mining  
by  
Tan, Steinbach, Kumar

## 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, Diaper, Juice, Eggs
3	Milk, Diaper, Juice, Coke
4	Bread, Milk, Diaper, Juice
5	Bread, Milk, Diaper, Coke

### Example of Association Rules

$\{\text{Diaper}\} \rightarrow \{\text{Juice}\},$   
 $\{\text{Milk, Bread}\} \rightarrow \{\text{Eggs, Coke}\},$   
 $\{\text{Juice, Bread}\} \rightarrow \{\text{Milk}\},$

Implication means co-occurrence,  
not causality!

## Definition: Frequent Itemset

### Itemset

- ▶ A collection of one or more items
  - ▶ Example: {Milk, Bread, Diaper}
- ▶ k-itemset
  - ▶ An itemset that contains k items

### Support count ( $\sigma$ )

- ▶ Frequency of occurrence of an itemset
- ▶ E.g.  $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

### Support

- ▶ 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

TID	Items
1	Bread, Milk
2	Bread, Diaper, Juice, Eggs
3	Milk, Diaper, Juice, Coke
4	Bread, Milk, Diaper, Juice
5	Bread, Milk, Diaper, Coke

## Definition: Association Rule

### Association Rule

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

TID	Items
1	Bread, Milk
2	Bread, Diaper, Juice, Eggs
3	Milk, Diaper, Juice, Coke
4	Bread, Milk, Diaper, Juice
5	Bread, Milk, Diaper, Coke

### 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

Example:

$\{\text{Milk, Diaper}\} \Rightarrow \text{Juice}$

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

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

## Association Rule Mining Task

- ▶ Given a set of transactions  $T$ , the goal of association rule mining is to find all rules having
    - ▶ support  $\geq \text{minsup}$  threshold
    - ▶ 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
- ⇒ **Computationally prohibitive!**

## Mining Association Rules

TID	Items
1	Bread, Milk
2	Bread, Diaper, Juice, Eggs
3	Milk, Diaper, Juice, Coke
4	Bread, Milk, Diaper, Juice
5	Bread, Milk, Diaper, Coke

### Example of Rules:

$\{\text{Milk}, \text{Diaper}\} \rightarrow \{\text{Juice}\}$  ( $s=0.4, c=0.67$ )  
 $\{\text{Milk}, \text{Juice}\} \rightarrow \{\text{Diaper}\}$  ( $s=0.4, c=1.0$ )  
 $\{\text{Diaper}, \text{Juice}\} \rightarrow \{\text{Milk}\}$  ( $s=0.4, c=0.67$ )  
 $\{\text{Juice}\} \rightarrow \{\text{Milk}, \text{Diaper}\}$  ( $s=0.4, c=0.67$ )  
 $\{\text{Diaper}\} \rightarrow \{\text{Milk}, \text{Juice}\}$  ( $s=0.4, c=0.5$ )  
 $\{\text{Milk}\} \rightarrow \{\text{Diaper}, \text{Juice}\}$  ( $s=0.4, c=0.5$ )

### Observations:

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

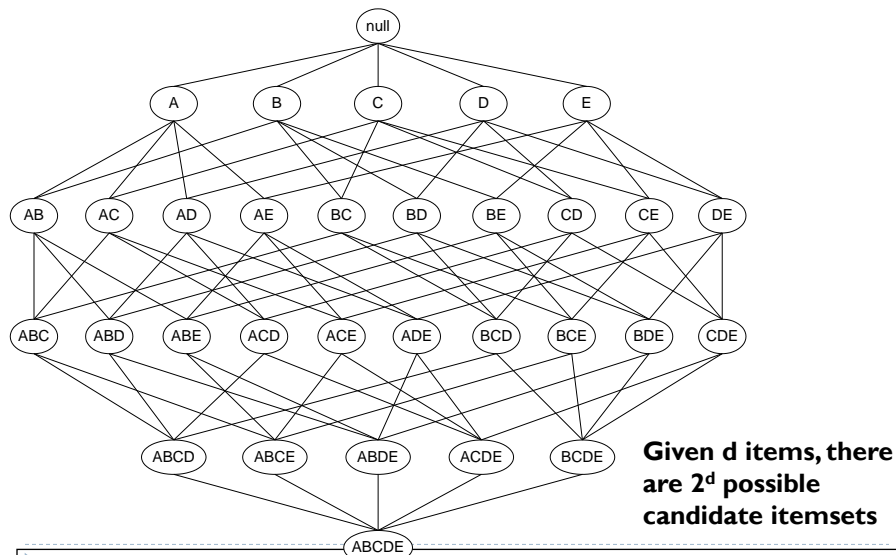
- Generate all itemsets whose support  $\geq$  minsup

#### 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

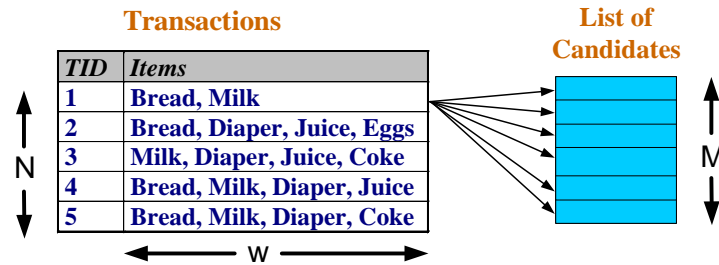
## Frequent Itemset Generation



## 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

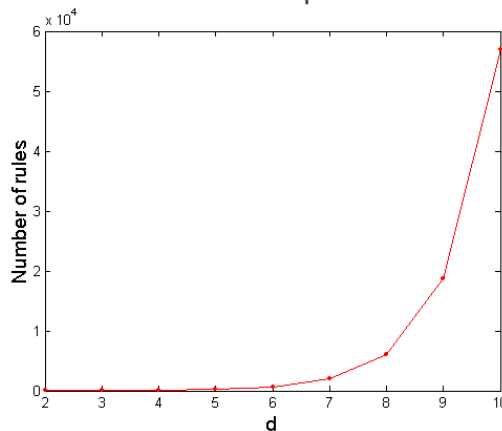


- Match each transaction against every candidate
- Complexity  $\sim O(NMw) \Rightarrow$  **Expensive since  $M = 2^d$  !!!**

## 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
  - ▶ Use vertical-partitioning of the data to apply the 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

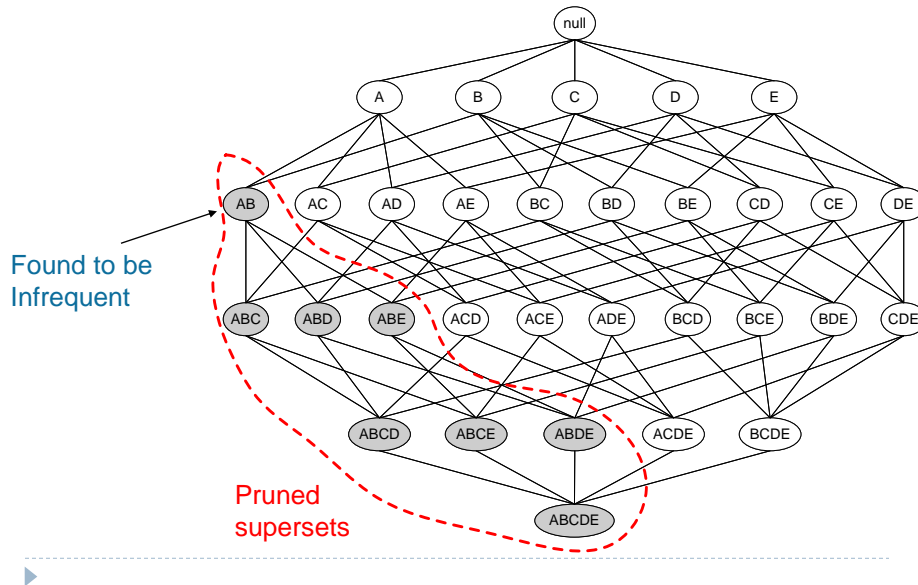
---

- ▶ **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:

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

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

## Illustrating Apriori Principle



## Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Juice	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread, Milk}	3
{Bread, Juice}	2
{Bread, Diaper}	3
{Milk, Juice}	2
{Milk, Diaper}	3
{Juice, Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)



Triplets (3-itemsets)

Itemset	Count
{Bread, Milk, Diaper}	3



Minimum Support = 3

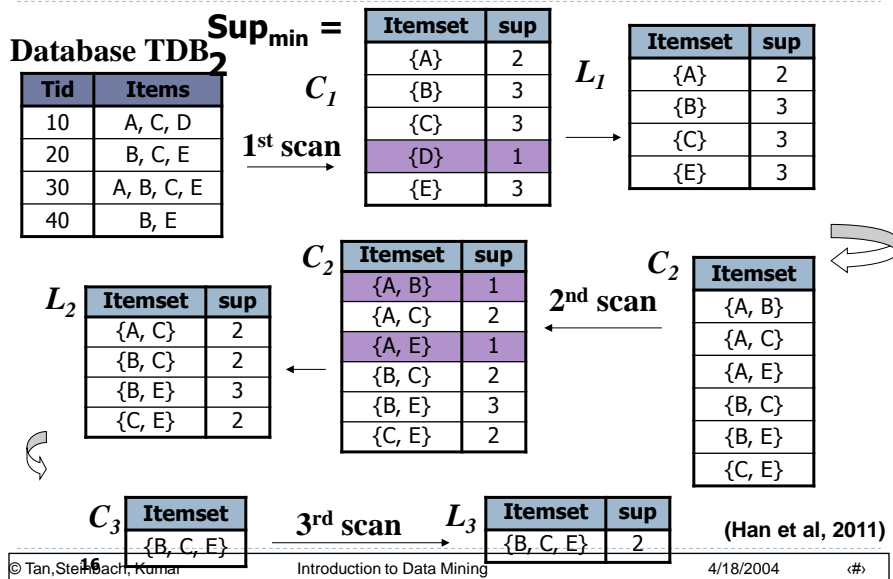
If every subset is considered,  
 ${}^6C_1 + {}^6C_2 + {}^6C_3 = 6 + 15 + 20 = 41$   
 With support-based pruning,  
 $6 + 6 + 1 = 13$

# Apriori Algorithm

## Method:

- ▶ Let  $k=1$
- ▶ Generate frequent itemsets of length 1
- ▶ 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 DB
  - ▶ Eliminate candidates that are infrequent, leaving only those that are frequent

## The Apriori Algorithm—An Example

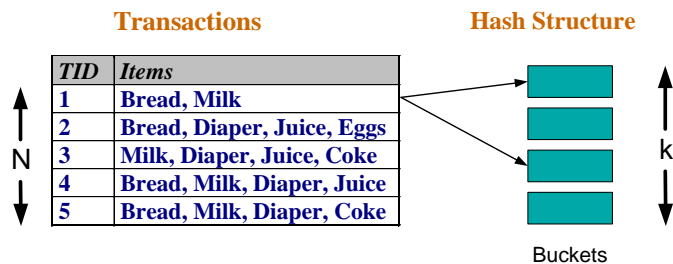




## Reducing Number of Comparisons

### ▶ Candidate counting:

- ▶ Scan the database of transactions to determine the support of each candidate itemset
- ▶ To reduce the number of comparisons, store the candidates in a hash structure
  - ▶ Instead of matching each transaction against every candidate, match it against candidates contained in the hashed buckets



## Factors Affecting Complexity

- ▶ **Choice of minimum support threshold**
  - ▶ lowering support threshold results in more frequent itemsets
  - ▶ this may increase number of candidates and max length of frequent itemsets
- ▶ **Dimensionality (number of items) of the data set**
  - ▶ more space is needed to store support count of each item
  - ▶ if number of frequent items also increases, both computation and I/O costs may also increase
- ▶ **Size of database**
  - ▶ since Apriori makes multiple passes, run time of algorithm may increase with number of transactions
- ▶ **Average transaction width**
  - ▶ transaction width increases with denser data sets
  - ▶ This may increase max length of frequent itemsets and traversals of hash tree (number of subsets in a transaction increases with its width)

## Compact Representation of Frequent Itemsets

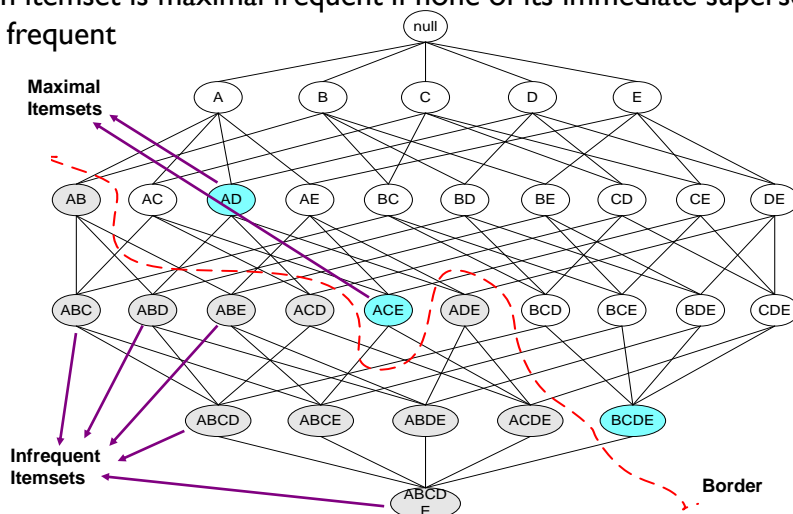
- Some itemsets are redundant because they have identical support as their supersets

TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1

- Number of frequent itemsets  $= 3 \times \sum_{k=1}^{10} \binom{10}{k}$
- Need a compact representation

## Maximal Frequent Itemset

An itemset is maximal frequent if none of its immediate supersets is frequent



## Maximal Frequent Itemset

- ▶ {A, D}, {A, C, E}, dan {B, C, D, E} are considered to be maximal frequent itemsets because their immediate supersets are infrequent.
- ▶ For example, {A,D} is maximal frequent since all of its immediate supersets, {A,B,D}, {A,C,D}, and {A,D,E} are infrequent.
- ▶ In contrast, an itemset such as {A,C} is non-maximal because one of its immediate supersets, {A,C,E}, is frequent
- ▶ Maximal frequent itemsets effectively provide a compact representation of frequent itemsets.
- ▶ For example, the frequent itemsets can be divided into two groups:
  - ▶ Frequent itemsets that begin with item a and that may contain items c, d, or e. This group includes itemsets such as {A}, {A, C}, {A, D}, {A, E}, and {A, C, E}
  - ▶ Frequent itemsets that begin with items b, c, d, or e. This group includes itemsets such as {B}, {B, C}, {C, D}, {B, C, D, E}, etc.

## Closed Itemset

- ▶ An itemset is closed if none of its immediate supersets has the same support as the itemset

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,B,C,D}
4	{A,B,D}
5	{A,B,C,D}

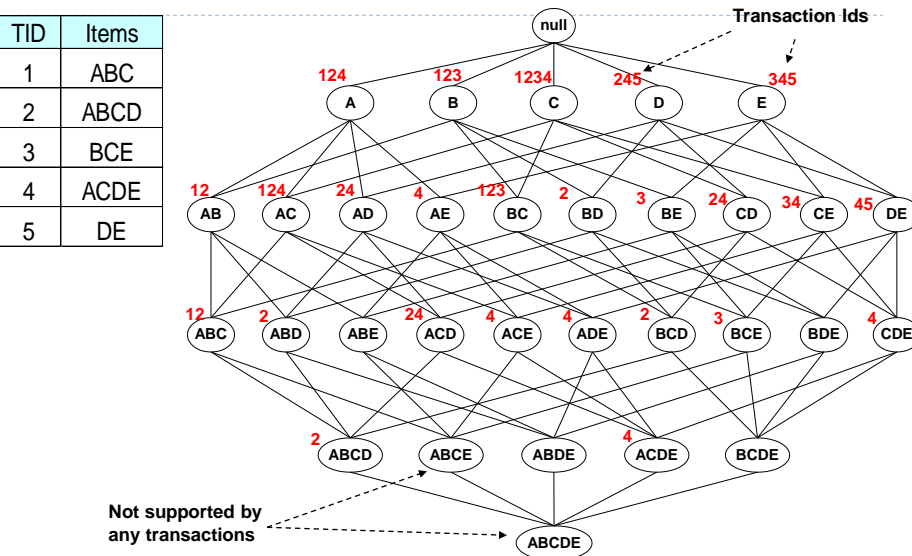
**Closed itemset**

Itemset	Support
{A,B,C}	2
{A,B,D}	3
{A,C,D}	2
{B,C,D}	3
{A,B,C,D}	2

Itemset	Support
{A}	4
{B}	5
{C}	3
{D}	4
{A,B}	4
{A,C}	2
{A,D}	3
{B,C}	3
{B,D}	4
{C,D}	3

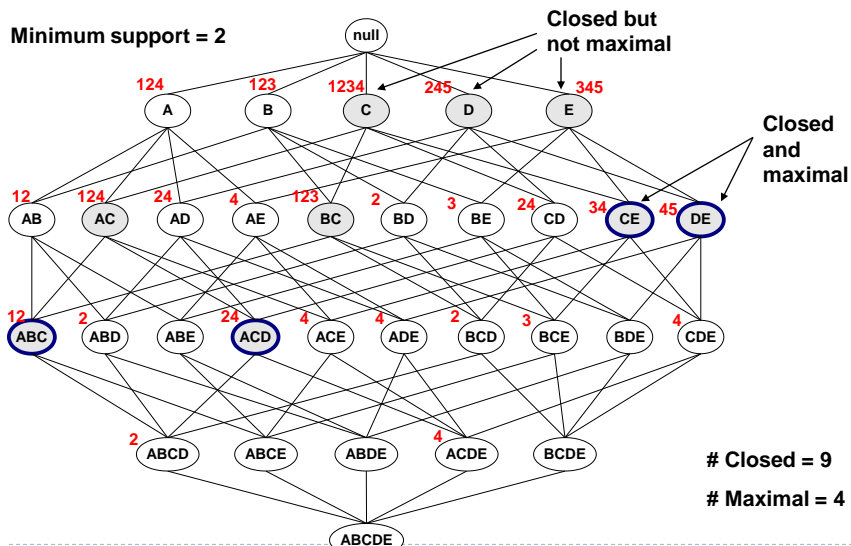
## Maximal vs Closed Itemsets

TID	Items
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE

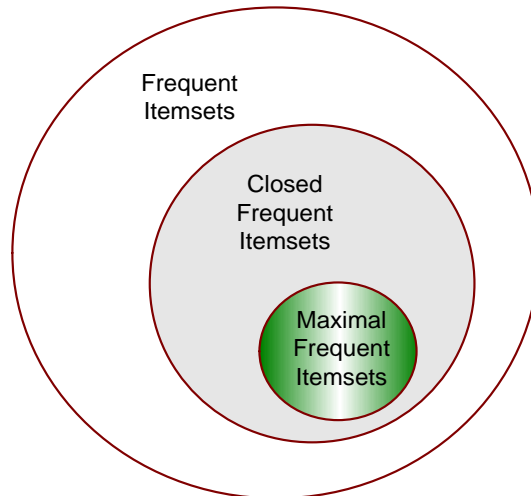


## Maximal vs Closed Frequent Itemsets

Minimum support = 2



## Maximal vs Closed Itemsets



## Alternative Methods for Frequent Itemset Generation

- Representation of Database
  - horizontal vs vertical data layout

Horizontal  
Data Layout

TID	Items
1	A,B,E
2	B,C,D
3	C,E
4	A,C,D
5	A,B,C,D
6	A,E
7	A,B
8	A,B,C
9	A,C,D
10	B

Vertical Data Layout

A	B	C	D	E
1	1	2	2	1
4	2	3	4	3
5	5	4	5	6
6	7	8	9	
7	8	9		
8	10			
9				

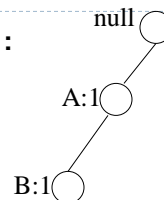
## FP-growth Algorithm

- ▶ Use a compressed representation of the database using an **FP-tree**
- ▶ Once an FP-tree has been constructed, it uses a recursive divide-and-conquer approach to mine the frequent itemsets

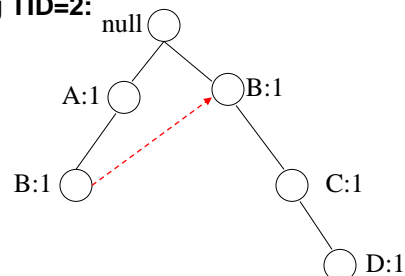
## FP-tree construction

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

After reading TID=1:



After reading TID=2:



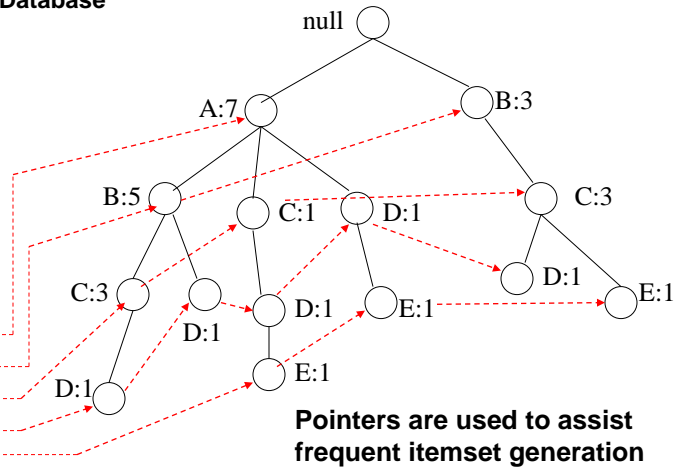
## FP-Tree Construction

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

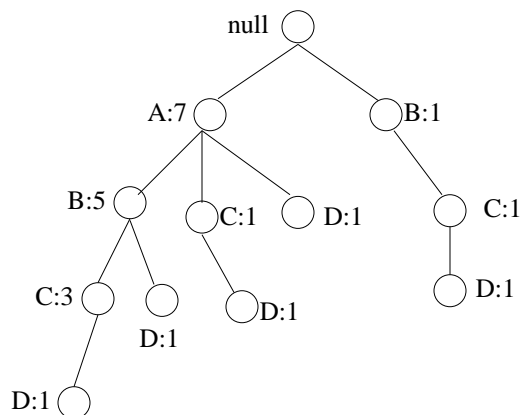
**Transaction Database**

**Header table**

Item	Pointer
A	
B	
C	
D	
E	



## FP-growth



Conditional Pattern base for D:

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

Recursively apply FP-growth on P

Frequent Itemsets found (with sup > 1):

AD, BD, CD, ACD, BCD

# ECLAT

- For each item, store a list of transaction ids (tids)

Horizontal  
Data Layout

TID	Items
1	A,B,E
2	B,C,D
3	C,E
4	A,C,D
5	A,B,C,D
6	A,E
7	A,B
8	A,B,C
9	A,C,D
10	B

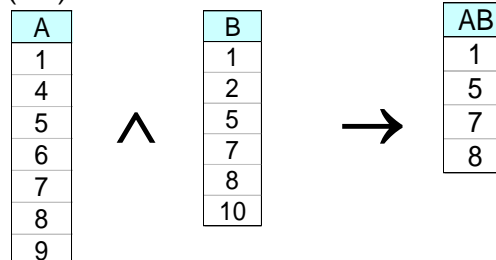
Vertical Data Layout

A	B	C	D	E
1	1	2	2	1
4	2	3	4	3
5	5	4	5	6
6	7	8	9	
7	8	9		
8	10			
9				

↓  
**TID-list**

# ECLAT

- Determine support of any k-itemset by intersecting tid-lists of two of its (k-1) subsets.



- 3 traversal approaches:
  - top-down, bottom-up and hybrid
- Advantage: very fast support counting
- Disadvantage: intermediate tid-lists may become too large for memory



## The Apriori Algorithm

### ► Pseudo-code:

$C_k$ : Candidate itemset of size  $k$

$L_k$ : frequent itemset of size  $k$

$L_1 = \{\text{frequent items}\};$

**for** ( $k = 1; L_k \neq \emptyset; k++$ ) **do begin**

$C_{k+1}$  = candidates generated from  $L_k$ ;

**for each** transaction  $t$  in database **do**

increment the count of all candidates in  $C_{k+1}$

that are contained in  $t$

$L_{k+1}$  = candidates in  $C_{k+1}$  with min\_support

**end**

**return**  $\cup_k L_k$ ;

## Important Details of Apriori

### ► How to generate candidates?

► Step 1: self-joining  $L_k$

► Step 2: pruning

### ► How to count supports of candidates?

### ► Example of Candidate-generation

►  $L_3 = \{abc, abd, acd, ace, bcd\}$

► Self-joining:  $L_3 * L_3$

►  $abcd$  from  $abc$  and  $abd$

►  $acde$  from  $acd$  and  $ace$

► Pruning:

►  $acde$  is removed because  $ade$  is not in  $L_3$

►  $C_4 = \{abcd\}$

## 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  $\emptyset \rightarrow L$ )

## Confidence based pruning

- ▶ **Theorema 6.2.** If a rule  $X \rightarrow Y$  does not satisfy the confidence threshold, the any rule  $X' \rightarrow Y - X'$ , where  $X'$  is a subset of  $X$ , must not satisfy the confidence threshold as well.

## Apriori Algorithm

---

---

**Algorithm 6.2** Rule generation of the *Apriori* algorithm.

---

```
1: for each frequent  $k$ -itemset  $f_k$ ,  $k \geq 2$  do
2:    $H_1 = \{i \mid i \in f_k\}$       {1-item consequents of the rule.}
3:   call ap-genrules( $f_k, H_1$ .)
4: end for
```

---

## Apriori Algorithm

---

---

**Algorithm 6.3** Procedure ap-genrules( $f_k, H_m$ ).

---

```
1:  $k = |f_k|$     {size of frequent itemset.}
2:  $m = |H_m|$     {size of rule consequent.}
3: if  $k > m + 1$  then
4:    $H_{m+1} = \text{apriori-gen}(H_m)$ .
5:   for each  $h_{m+1} \in H_{m+1}$  do
6:      $conf = \sigma(f_k) / \sigma(f_k - h_{m+1})$ .
7:     if  $conf \geq minconf$  then
8:       output the rule  $(f_k - h_{m+1}) \rightarrow h_{m+1}$ .
9:     else
10:      delete  $h_{m+1}$  from  $H_{m+1}$ .
11:    end if
12:  end for
13:  call ap-genrules( $f_k, H_{m+1}$ .)
14: end if
```

---

## Rule Generation

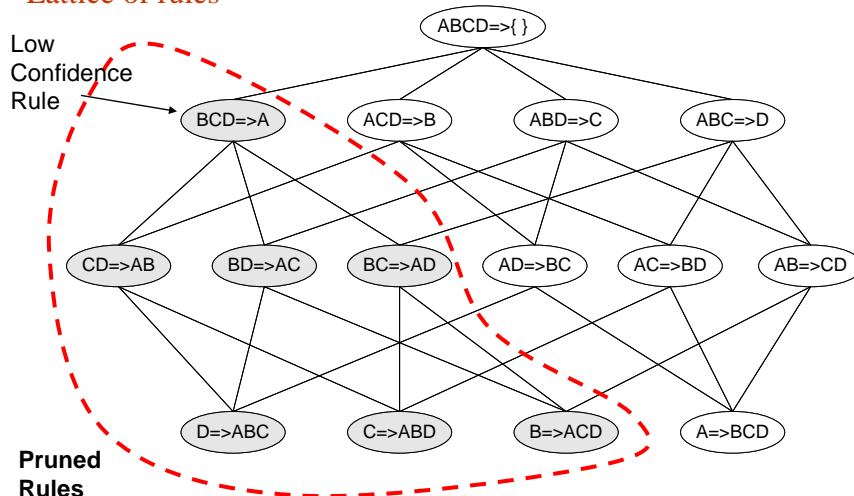
- ▶ How to efficiently generate rules from frequent itemsets?
  - ▶ In general, confidence does not have an anti-monotone property  
 $c(ABC \rightarrow D)$  can be larger or smaller than  $c(AB \rightarrow D)$
  - ▶ But confidence of rules generated from the same itemset has an anti-monotone property
  - ▶ e.g.,  $L = \{A, B, C, D\}$ :

$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$

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

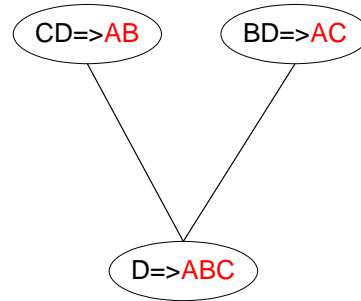
## Rule Generation for Apriori Algorithm

### Lattice of rules



## Rule Generation for Apriori Algorithm

- ▶ Candidate rule is generated by merging two rules that share the same prefix in the rule consequent
- ▶  $\text{join}(\text{CD} \Rightarrow \text{AB}, \text{BD} \Rightarrow \text{AC})$  would produce the candidate rule  $\text{D} \Rightarrow \text{ABC}$
- ▶ Prune rule  $\text{D} \Rightarrow \text{ABC}$  if its subset  $\text{AD} \Rightarrow \text{BC}$  does not have high confidence



## Effect of Support Distribution

- ▶ How to set the appropriate *minsup* threshold?
  - ▶ 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
- ▶ Using a single minimum support threshold may not be effective

## Multiple Minimum Support

- ▶ How to apply multiple minimum supports?
  - ▶  $MS(i)$ : minimum support for item  $i$
  - ▶ e.g.:  $MS(\text{Milk})=5\%$ ,  $MS(\text{Coke}) = 3\%$ ,  
 $MS(\text{Broccoli})=0.1\%$ ,  $MS(\text{Salmon})=0.5\%$
  - ▶  $MS(\{\text{Milk}, \text{Broccoli}\}) = \min (MS(\text{Milk}), MS(\text{Broccoli}))$   
 $= 0.1\%$
- ▶ Challenge: Support is no longer anti-monotone
  - ▶ Suppose:  $\text{Support}(\text{Milk}, \text{Coke}) = 1.5\%$  and  
 $\text{Support}(\text{Milk}, \text{Coke}, \text{Broccoli}) = 0.5\%$
  - ▶  $\{\text{Milk}, \text{Coke}\}$  is infrequent but  $\{\text{Milk}, \text{Coke}, \text{Broccoli}\}$  is frequent

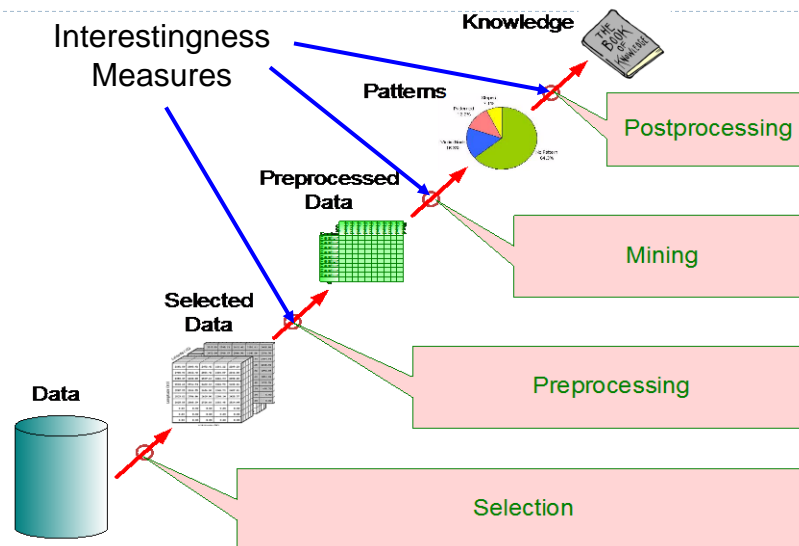
## Multiple Minimum Support (Liu 1999)

- ▶ Order the items according to their minimum support (in ascending order)
  - ▶ e.g.:  $MS(\text{Milk})=5\%$ ,  $MS(\text{Coke}) = 3\%$ ,  
 $MS(\text{Broccoli})=0.1\%$ ,  $MS(\text{Salmon})=0.5\%$
  - ▶ Ordering: Broccoli, Salmon, Coke, Milk
- ▶ Need to modify Apriori such that:
  - ▶  $L_1$ : set of frequent items
  - ▶  $F_1$ : set of items whose support is  $\geq MS(I)$   
where  $MS(I)$  is  $\min_i (MS(i))$
  - ▶  $C_2$ : candidate itemsets of size 2 is generated from  $F_1$   
instead of  $L_1$

## Pattern Evaluation

- ▶ Association rule algorithms tend to produce too many rules
  - ▶ many of them are uninteresting or redundant
  - ▶ Redundant if  $\{A,B,C\} \rightarrow \{D\}$  and  $\{A,B\} \rightarrow \{D\}$  have same support & confidence
- ▶ Interestingness measures can be used to prune/rank the derived patterns
- ▶ In the original formulation of association rules, support & confidence are the only measures used

## Application of Interestingness Measure



## Computing Interestingness Measure

- ▶ Given a rule  $X \rightarrow Y$ , information needed to compute rule interestingness can be obtained from a contingency table

Contingency table for  $X \rightarrow Y$

	Y	$\bar{Y}$	
X	$f_{11}$	$f_{10}$	$f_{1+}$
$\bar{X}$	$f_{01}$	$f_{00}$	$f_{0+}$
	$f_{+1}$	$f_{+0}$	$ T $

$f_{11}$ : support of X and Y

$f_{10}$ : support of  $\bar{X}$  and  $\bar{Y}$

$f_{01}$ : support of  $\bar{X}$  and Y

$f_{00}$ : support of X and  $\bar{Y}$

Used to define various measures

- ◆ support, confidence, lift, Gini, J-measure, etc.

## Statistical Independence

- ▶ Population of 1000 students
  - ▶ 600 students know how to swim (S)
  - ▶ 700 students know how to bike (B)
  - ▶ 420 students know how to swim and bike (S,B)
- ▶  $P(S \wedge B) = 420/1000 = 0.42$
- ▶  $P(S) \times P(B) = 0.6 \times 0.7 = 0.42$
- ▶  $P(S \wedge B) = P(S) \times P(B) \Rightarrow$  Statistical independence
- ▶  $P(S \wedge B) > P(S) \times P(B) \Rightarrow$  Positively correlated
- ▶  $P(S \wedge B) < P(S) \times P(B) \Rightarrow$  Negatively correlated



## Statistical-based Measures

- Measures that take into account statistical dependence

$$Lift = \frac{P(Y | X)}{P(Y)}$$

$$Interest = \frac{P(X, Y)}{P(X)P(Y)}$$

$$PS = P(X, Y) - P(X)P(Y)$$

$$\phi - coefficient = \frac{P(X, Y) - P(X)P(Y)}{\sqrt{P(X)[1 - P(X)]P(Y)[1 - P(Y)]}}$$

## Example: Lift/Interest

	Coffee	<u>Coffee</u>	
Tea	15	5	20
<u>Tea</u>	75	5	80
	90	10	100

Association Rule: Tea → Coffee

Confidence =  $P(\text{Coffee} | \text{Tea}) = 0.75$

but  $P(\text{Coffee}) = 0.9$

⇒ Lift =  $0.75/0.9 = 0.8333$  (< 1, therefore is negatively associated)

There are lots of measures proposed in the literature

Some measures are good for certain applications, but not for others

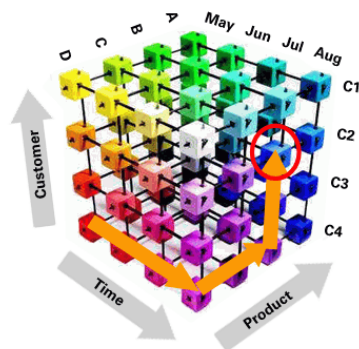
What criteria should we use to determine whether a measure is good or bad?

What about Apriori-style support based pruning? How does it affect these measures?

#	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 $\lambda$	$\frac{\sum_j \max_k P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
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 ( $L$ )	$\max \left( \frac{NP(A,B)+1}{NP(A)+2}, \frac{NP(A,B)+1}{NP(B)+2} \right)$
13	Conviction ( $V$ )	$\max \left( \frac{P(A)P(\bar{B})}{P(A\bar{B})}, \frac{P(B)P(\bar{A})}{P(\bar{A}B)} \right)$
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 ( $S$ )	$\frac{P(A,B) + P(\bar{A}\bar{B})}{P(A,B)} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A,B) - P(\bar{A}\bar{B})}$
20	Jaccard ( $\zeta$ )	$\frac{P(A,B)}{P(A) + P(B) - P(A,B)}$
21	Klosgen ( $K$ )	$\sqrt{P(\bar{A}, \bar{B})} \max(P(B A) - P(B), P(A B) - P(A))$

## Referensi

Tan P, Michael S., & Vipin K. 2006. *Introduction to Data mining*. Pearson Education, Inc.



Sumber: [www.oracle.com](http://www.oracle.com)

Materi selanjutnya:

Data Warehouse dan OLAP