# Association Rule Mining - II

## Dr. Faisal Kamiran

# 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)
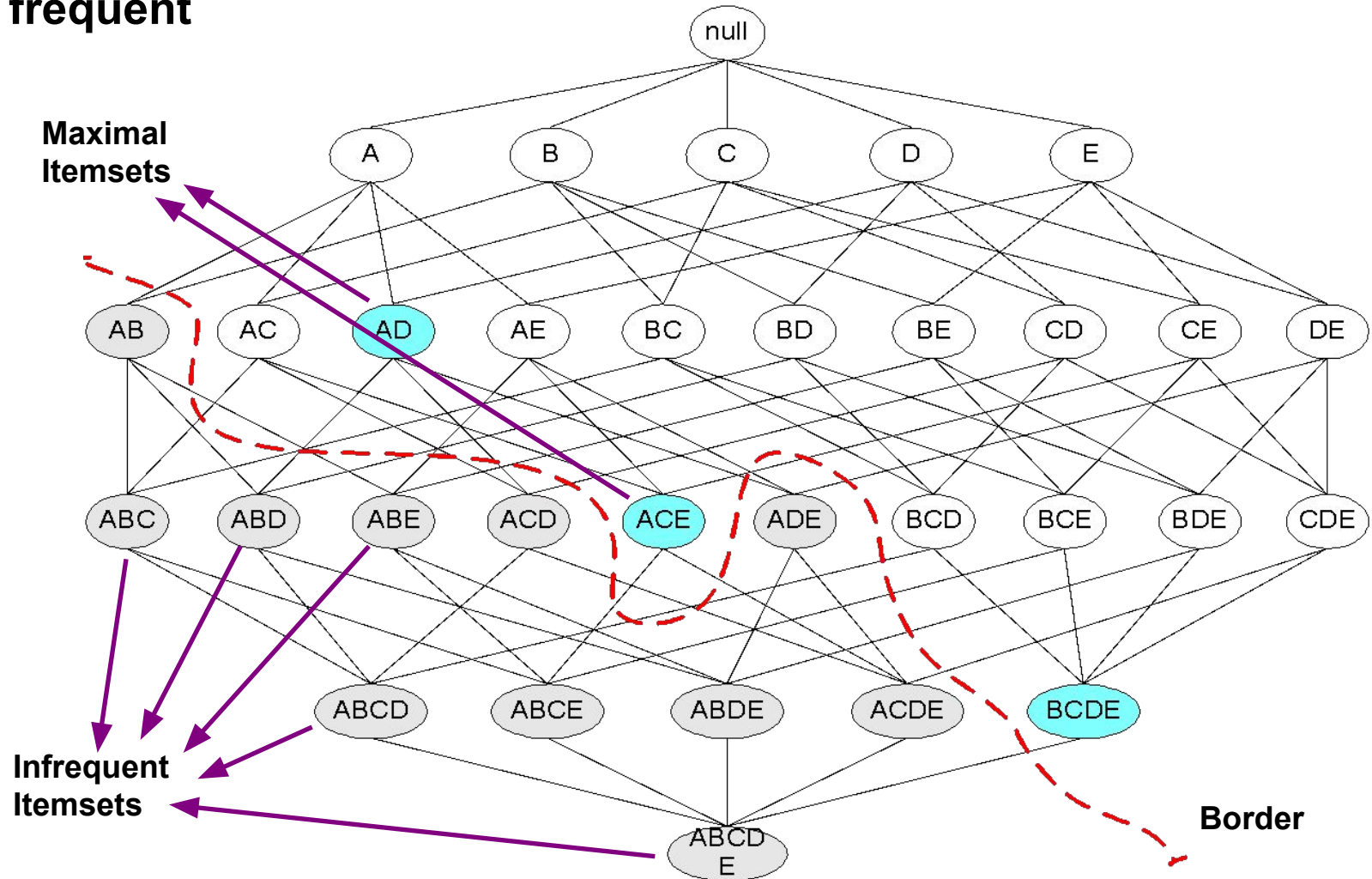
# Compact Representation

- Practically, too many frequent itemsets

- Need to identify a small representative set of itemsets to derive all itemsets

- Need to have compact representation of frequent itemsets:
  - **Maximal frequent itemsets**
  - **Closed frequent itemsets**

# Maximal Frequent Itemset

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

# Maximal Frequent Itemset

- Provide a valuable representation of dataset that can generate too many frequent itemsets

- We need an efficient algorithms to find maximal frequent itemsets

- Maximal frequent  itemsets do not contain support information about their subsets

# Closed Itemset

- Minimal representation without losing the support information

- 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} |

| 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 |

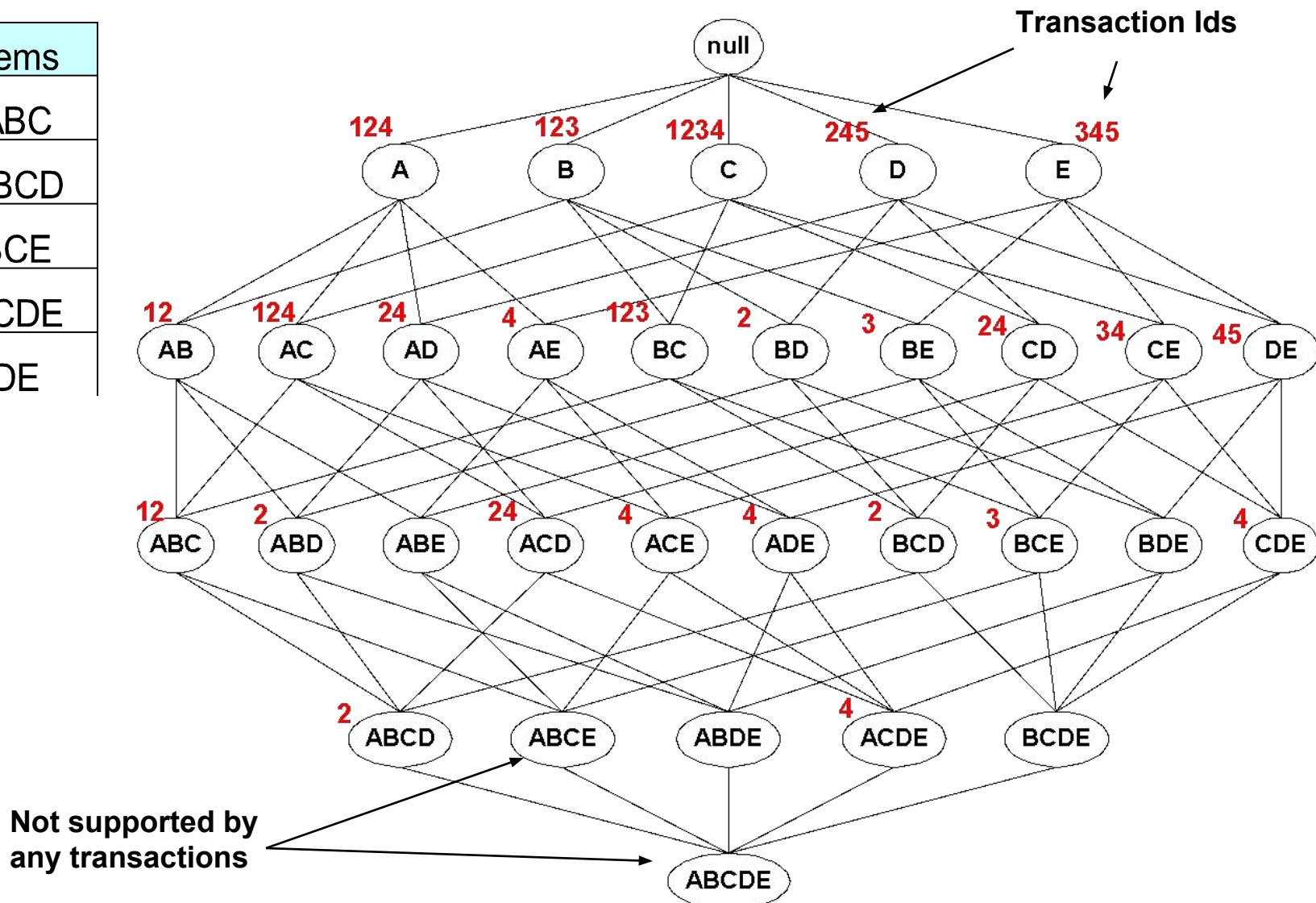| Itemset | Support |
|---------|---------|
| {A,B,C} | 2 |
| {A,B,D} | 3 |
| {A,C,D} | 2 |
| {B,C,D} | 3 |
| {A,B,C,D} | 2 |

# Closed Itemset

- An closed itemset is closed frequent itemset if its support is greater than the min support threshold

- Useful for removing some of the redundant association rules

- The rule X → Y is redundant if another rule X` → Y` exists with where $X \in X`$ $and$ $Y \in Y`$ and the support and confidence for both rules are identical.
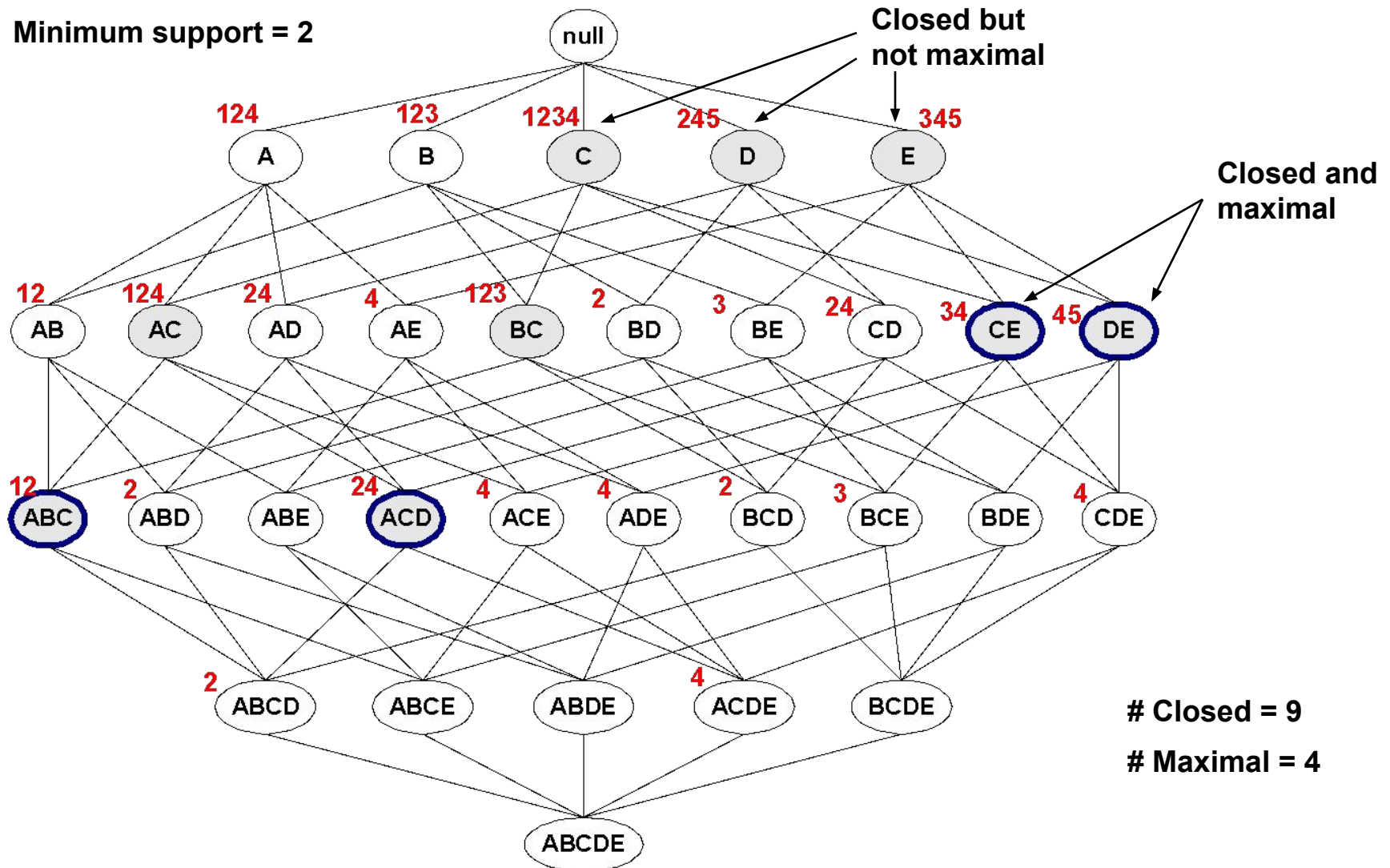
# Maximal vs Closed Itemsets



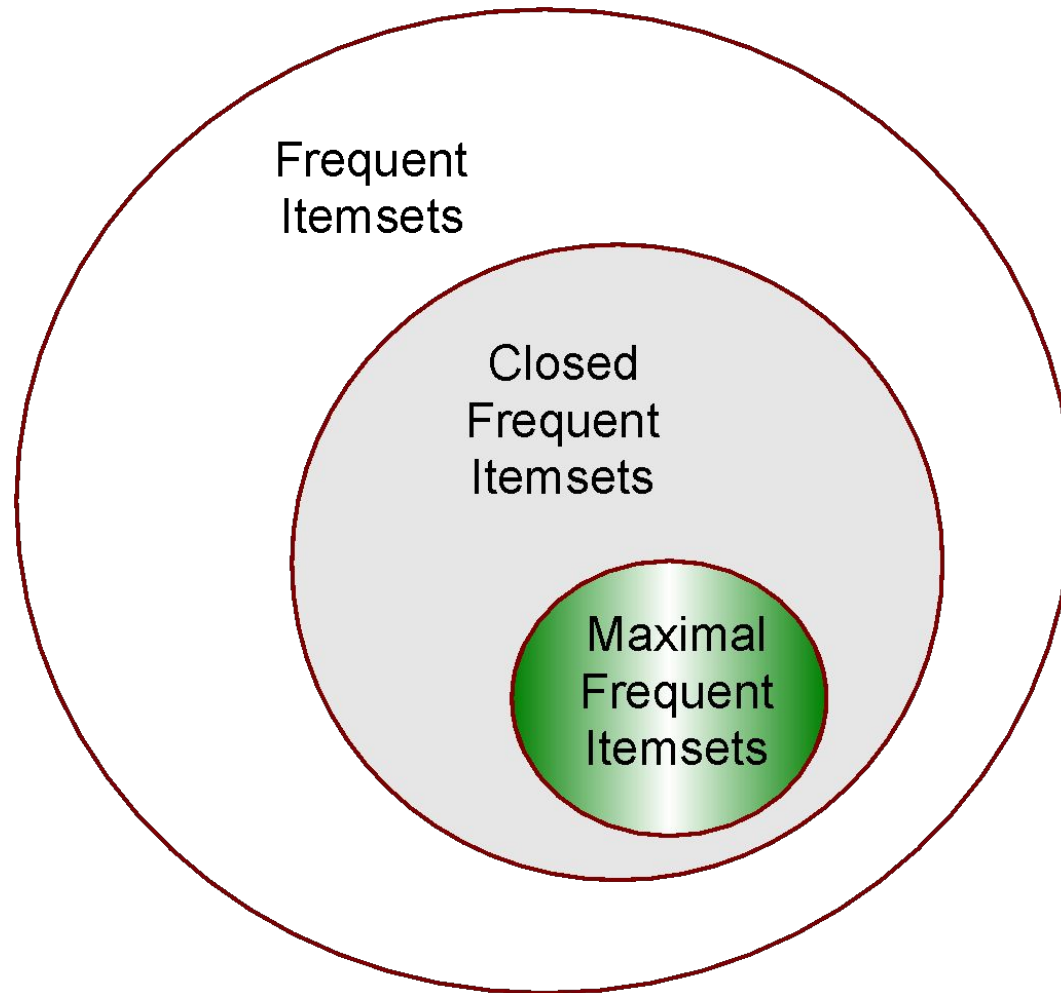| TID | Items |
|-----|-------|
| 1 | ABC |
| 2 | ABCD |
| 3 | BCE |
| 4 | ACDE |
| 5 | DE |

Transaction Ids

Not supported by any transactions

# Maximal vs Closed Frequent Itemsets



Minimum support = 2

Closed but not maximal

Closed and maximal

# Closed = 9

# Maximal = 4

# Maximal vs Closed Itemsets



Frequent Itemsets

Closed Frequent Itemsets

Maximal Frequent Itemsets

# Alternative Methods for Frequent Itemset Generation

- Performance of Apriori algorithm degrades over dense data because of increased width of transactions.

- Need to overcome this limitation and to improve the efficiency

- Conceptually the search for frequent itemset can be viewed as a traversal on the itemset lattice

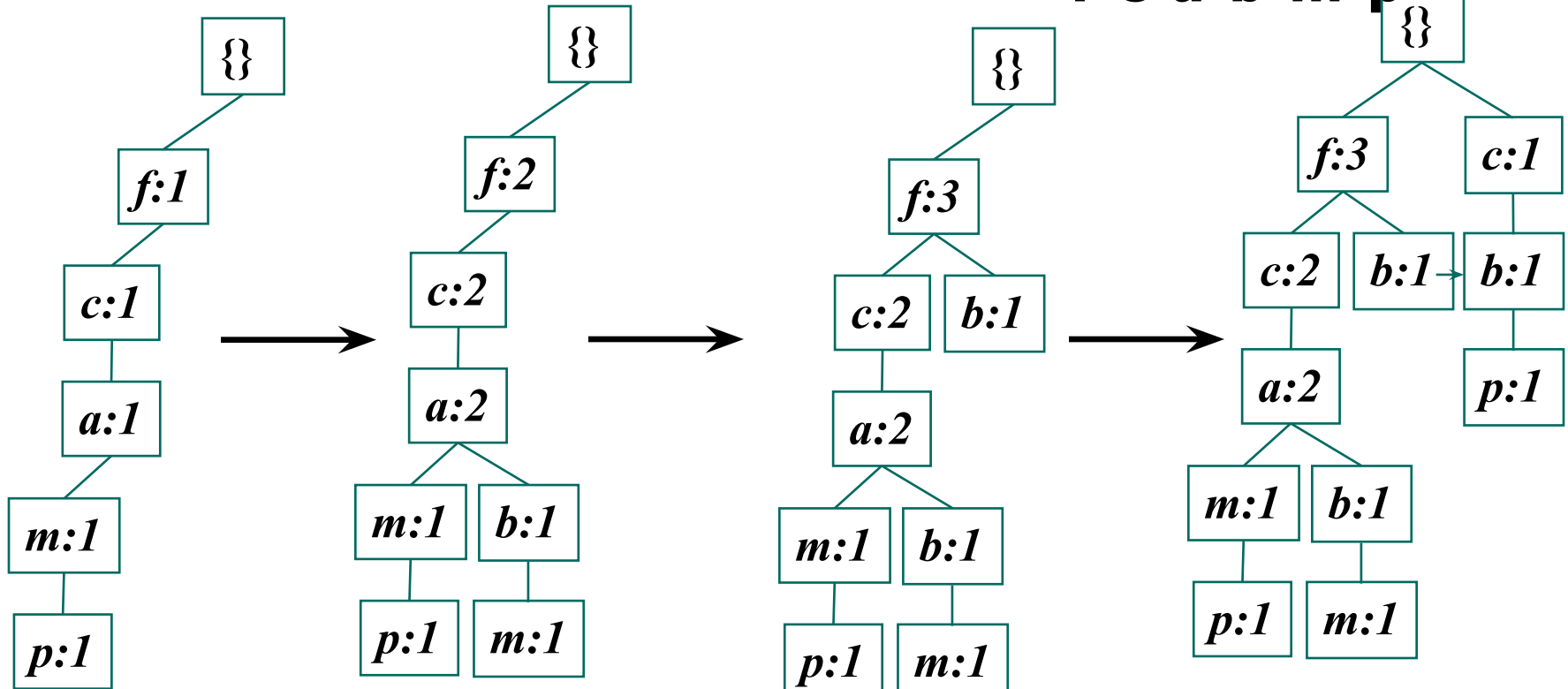- Need better search and traversal strategies

# 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

- Steps:
  - FP tree construction
  - Conditional pattern base construction
  - Conditional FP-trees
  - Frequent pattern generation

# FP-Tree Construction

| TID | Items bought | (ordered) frequent items |
|-----|--------------|--------------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o, w} | {f, b} |
| 400 | {b, c, k, s, p} | {c, b, p} |
| 500 | {a, f, c, e, l, p, m, n} | {f, c, a, m, p} |

$min\_support = 3$

**F-list =**
**f-c-a-b-m-p**

# FP-Tree Construction

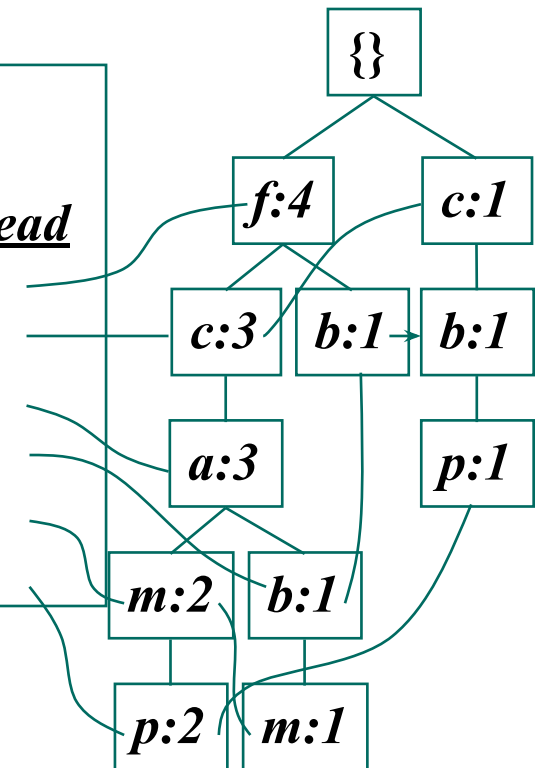| TID | Items bought | (ordered) frequent items |
|-----|-------------|--------------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o, w} | {f, b} |
| 400 | {b, c, k, s, p} | {c, b, p} |
| 500 | {a, f, c, e, l, p, m, n} | {f, c, a, m, p} |

*min_support = 3*

1. **Scan DB once, find frequent 1-itemset (single item pattern)**

2. **Sort frequent items in frequency descending order, f-list**

3. **Scan DB again, construct FP-tree**

**Header Table**

| Item | frequency | head |
|------|-----------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |
| p | 3 | |

**F-list = f-c-a-b-m-p**



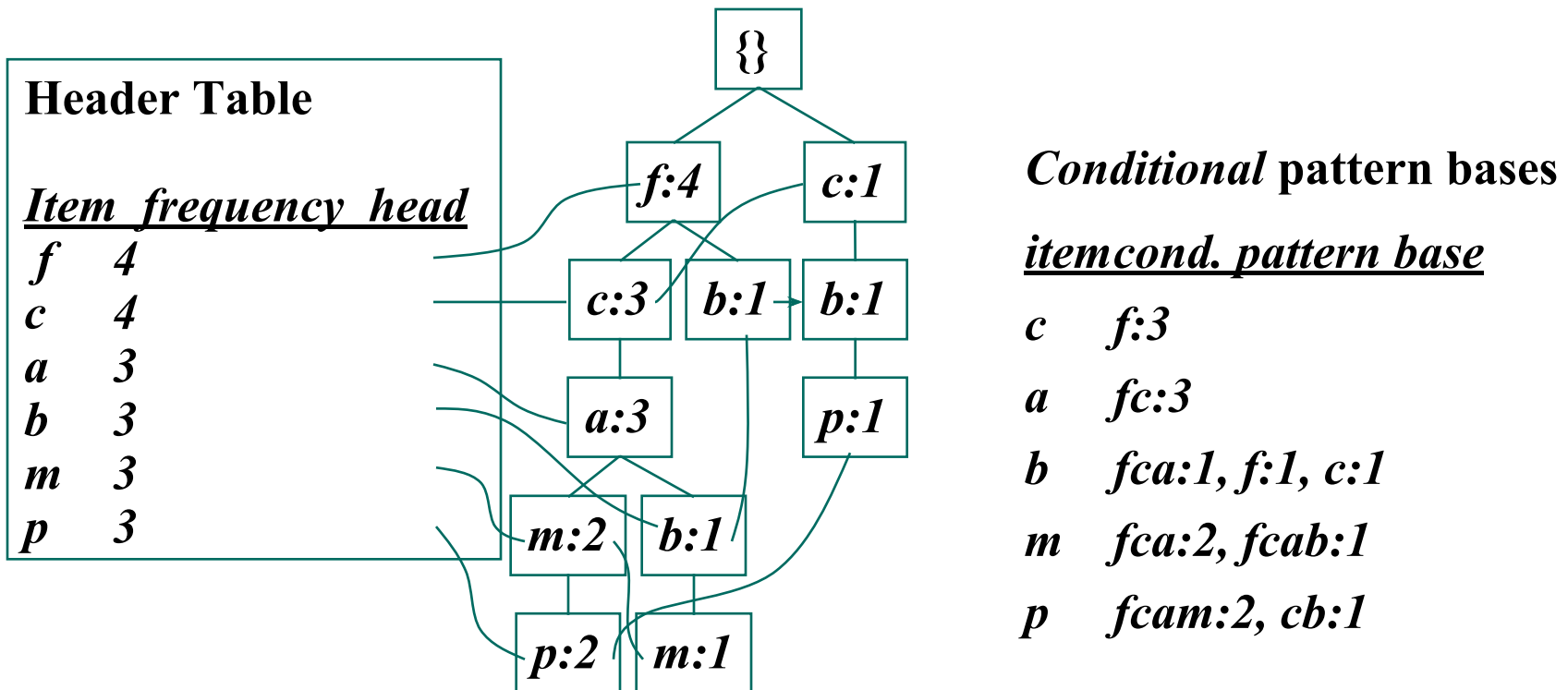© Tan,Steinbach, Kumar    Introduction to Data Mining    4/18/2004

# Partition Patterns and Databases

- Frequent patterns can be partitioned into subsets according to f-list
  - F-list = f-c-a-b-m-p
  - Patterns containing p
  - Patterns having m but no p
  - …
  - Patterns having c but no a nor b, m, p
  - Pattern f
- Completeness and non-redundancy
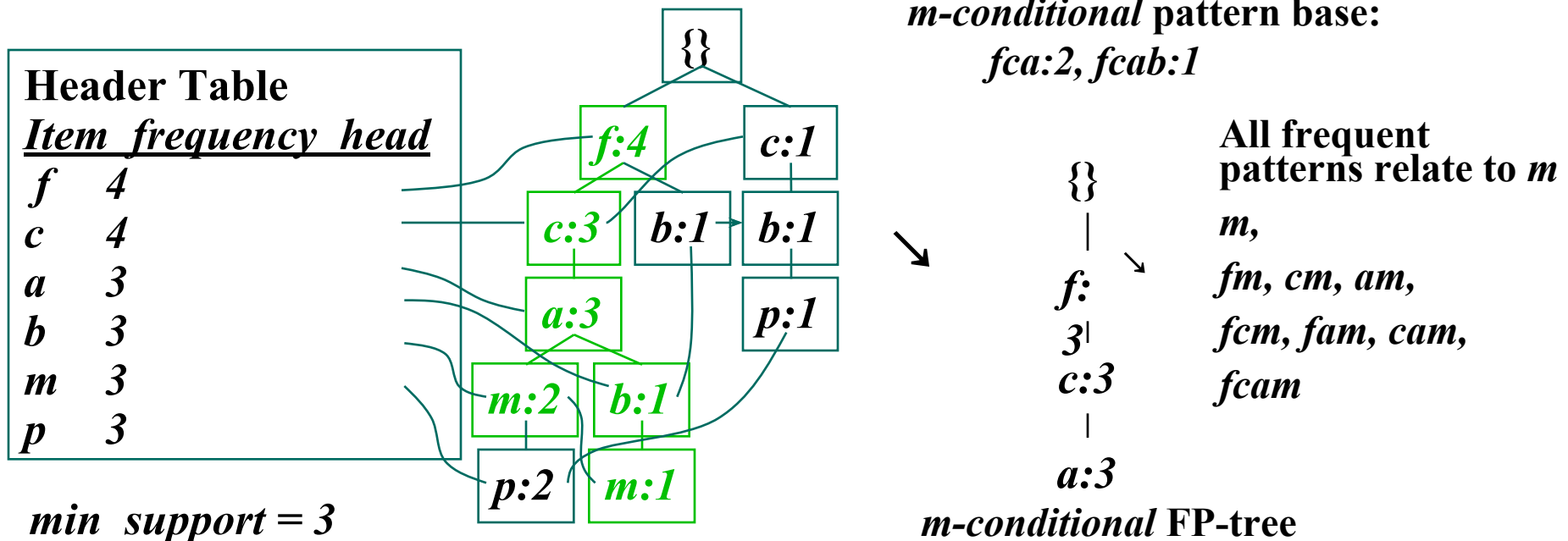
# Conditional Pattern Base Construction

- **Conditional Pattern Base** is a sub-database which consists of the set of prefix paths in the FP tree co-occurring with the suffix pattern or item. For example *prefix paths (fcam, cb)* of item *p* to form *p*'s conditional pattern base.

**Header Table**

*Item  frequency  head*

| Item | frequency |
|------|-----------|
| f | 4 |
| c | 4 |
| a | 3 |
| b | 3 |
| m | 3 |
| p | 3 |

FP tree structure:

{}
- f:4
  - c:3
    - a:3
      - m:2
        - p:2
      - b:1
        - m:1
- c:1
  - b:1
    - p:1

b:1 (under f:4)

*Conditional* **pattern bases**

*item  cond. pattern base*

| item | cond. pattern base |
|------|--------------------|
| c | f:3 |
| a | fc:3 |
| b | fca:1, f:1, c:1 |
| m | fca:2, fcab:1 |
| p | fcam:2, cb:1 |

# From Conditional Pattern-bases to Conditional FP-trees

- For each pattern-base
  - Accumulate the count for each item in the base
  - Construct the FP-tree for the frequent items of the pattern base
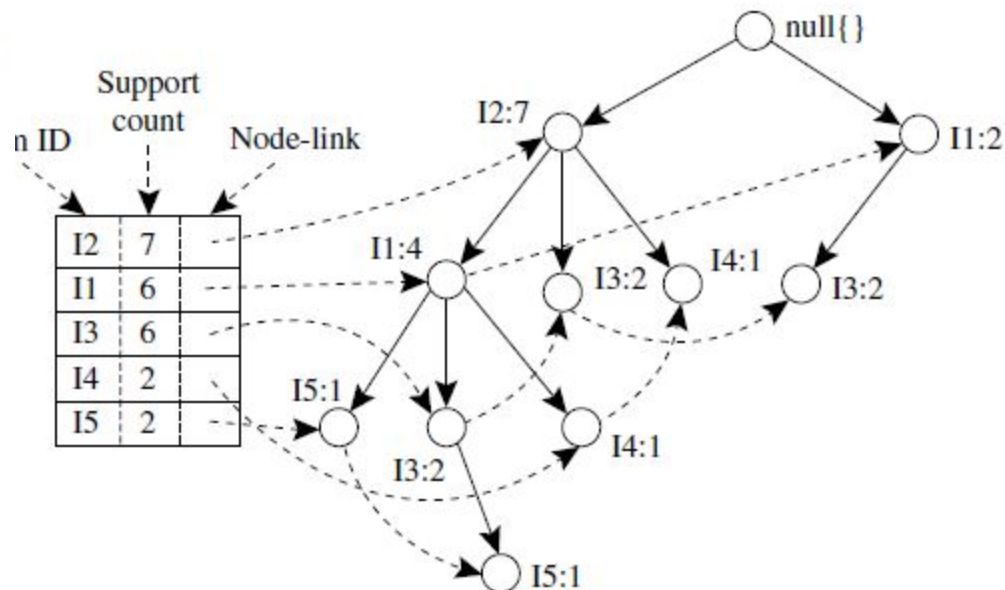


**Header Table**
*Item frequency head*

| | |
|---|---|
| *f* | *4* |
| *c* | *4* |
| *a* | *3* |
| *b* | *3* |
| *m* | *3* |
| *p* | *3* |

*min_support = 3*

*m-conditional* **pattern base:**
*fca:2, fcab:1*

**All frequent patterns relate to** *m*

*m,*

*fm, cm, am,*

*fcm, fam, cam,*

*fcam*

*m-conditional* **FP-tree**

{}
|
*f:*
*3*|
*c:3*
|
*a:3*

# FP Growth: Example

## Data:

## FP Tree:

| TID | List of item_IDs |
|-----|------------------|
| T100 | I1, I2, I5 |
| T200 | I2, I4 |
| T300 | I2, I3 |
| T400 | I1, I2, I4 |
| T500 | I1, I3 |
| T600 | I2, I3 |
| T700 | I1, I3 |
| T800 | I1, I2, I3, I5 |
| T900 | I1, I2, I3 |



### Mining the FP-Tree

| Item | Conditional Pattern Base | Conditional FP-tree | Frequent Patterns Generated |
|------|--------------------------|---------------------|------------------------------|
| I5 | {{I2, I1: 1}, {I2, I1, I3: 1}} | ⟨I2: 2, I1: 2⟩ | {I2, I5: 2}, {I1, I5: 2}, {I2, I1, I5: 2} |
| I4 | | | |
| I3 | | | |
| I1 | | | |

# Benefits of the FP-tree Structure

- 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)

# 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.

| A |
|---|
| 1 |
| 4 |
| 5 |
| 6 |
| 7 |
| 8 |
| 9 |

∧

| B |
|---|
| 1 |
| 2 |
| 5 |
| 7 |
| 8 |
| 10 |

→

| AB |
|---|
| 1 |
| 5 |
| 7 |
| 8 |

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

# Rule Generation

- Given a frequent itemset L, find all non-empty subsets f ⊂ L such that f → L – f satisfies the minimum confidence requirement
  - If {A,B,C,D} is a frequent itemset, candidate rules:

    ABC →D,   ABD →C,   ACD →B,   BCD →A,
    A →BCD,   B →ACD,   C →ABD,   D →ABC
    AB →CD,   AC → BD,  AD → BC,  BC →AD,
    BD →AC,   CD →AB,

- If |L| = k, then there are $2^k – 2$ candidate association rules (ignoring L → ∅ and ∅ → L)
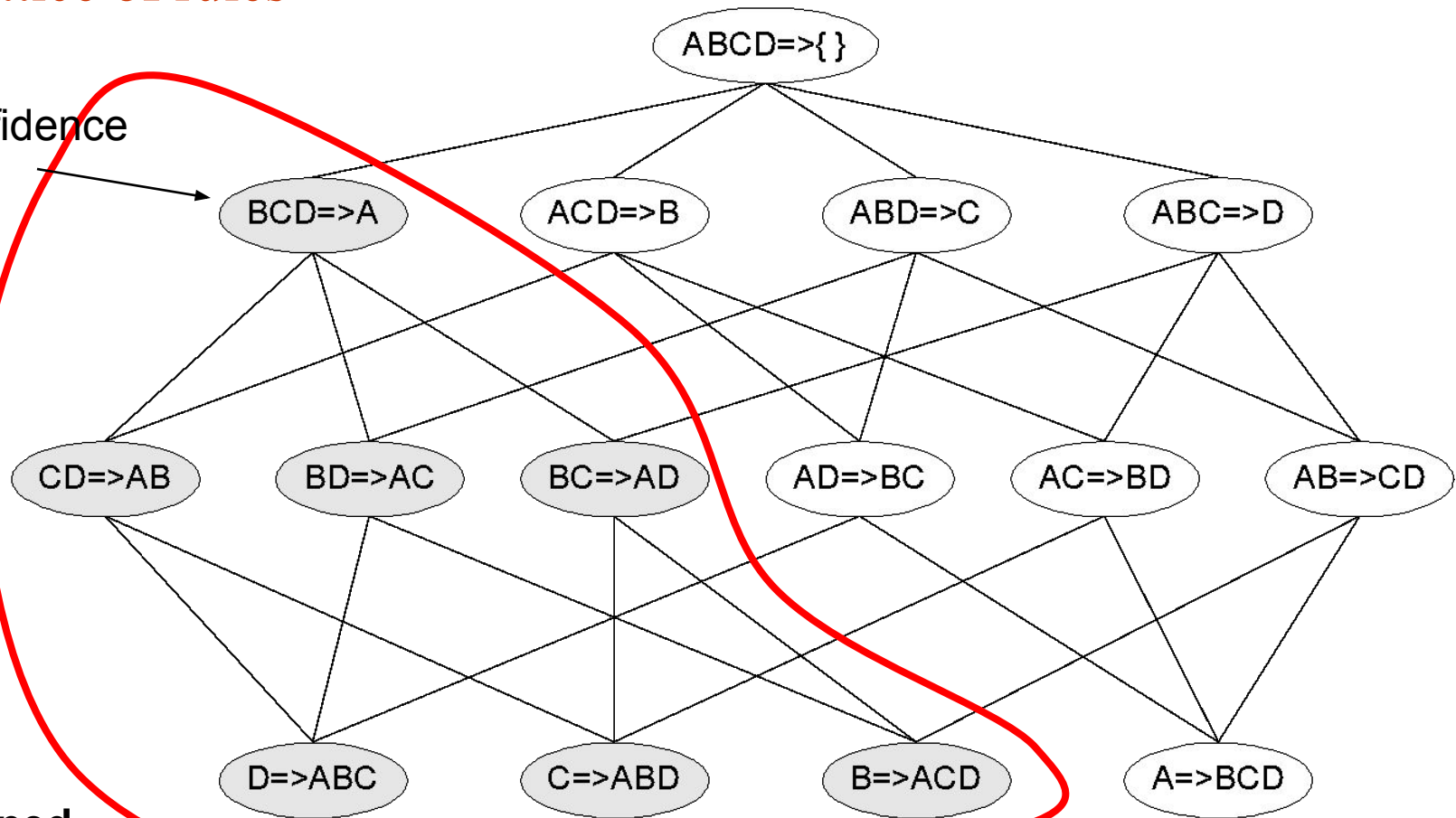
# Rule Generation

- 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., $L = \{A,B,C,D\}$:

    $$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$
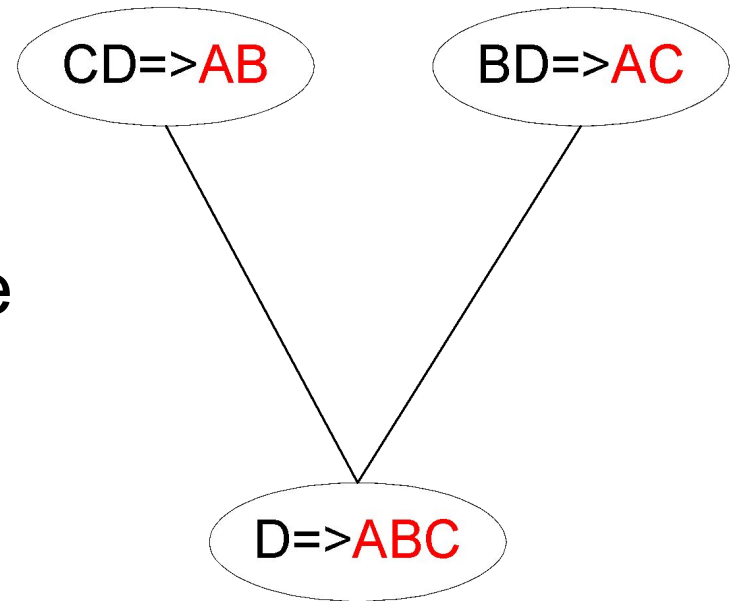
# Rule Generation for Apriori Algorithm

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

- join(CD=>AB,BD=>AC)
  would produce the candidate
  rule D => ABC

- Prune rule D=>ABC if its
  subset CD=>AB does not have
  high confidence

```
    ( CD=>AB )        ( BD=>AC )
          \             /
           \           /
            \         /
             \       /
              \     /
             ( D=>ABC )
```

# Pattern Evaluation

- Association rule algorithms tend to produce too many rules
    - many of them are uninteresting or redundant
    - Redundant if {A,B,C} → {D} and {A,B} → {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

# 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 | $\overline{Y}$ |  |
|---|---|---|---|
| X | $f_{11}$ | $f_{10}$ | $f_{1+}$ |
| $\overline{X}$ | $f_{01}$ | $f_{00}$ | $f_{o+}$ |
|  | $f_{+1}$ | $f_{+0}$ | $|T|$ |

$f_{11}$: support of X and Y
$f_{10}$: support of X and $\overline{Y}$
$f_{01}$: support of $\overline{X}$ and Y
$f_{00}$: support of $\overline{X}$ and $\overline{Y}$

Used to define various measures

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

# Drawback of Confidence

| | Coffee | $\overline{\text{Coffee}}$ | |
|---|---|---|---|
| Tea | 15 | 5 | 20 |
| $\overline{\text{Tea}}$ | 75 | 5 | 80 |
| | 90 | 10 | 100 |

Association Rule: Tea → Coffee

- Confidence= P(Coffee|Tea) = 0.75

- but P(Coffee) = 0.9

- Although confidence is high, rule is misleading

- P(Coffee|$\overline{\text{Tea}}$) = 0.9375

# 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)

  - $P(S \wedge B) = 420/1000 = 0.42$
  - $P(S) \times P(B) = 0.6 \times 0.7 = 0.42$

  - $P(S \wedge B) = P(S) \times P(B) =>$ Statistical independence
  - $P(S \wedge B) > P(S) \times P(B) =>$ Positively correlated
  - $P(S \wedge B) < P(S) \times P(B) =>$ Negatively correlated

# Statistical-based Measures

- Measures that take into account statistical dependence

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

# Example: Lift/Interest

|  | Coffee | $\overline{\text{Coffee}}$ |  |
|---|---|---|---|
| Tea | 15 | 5 | 20 |
| $\overline{\text{Tea}}$ | 75 | 5 | 80 |
|  | 90 | 10 | 100 |

Association Rule: Tea → Coffee

Confidence= P(Coffee|Tea) = 0.75

but P(Coffee) = 0.9

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