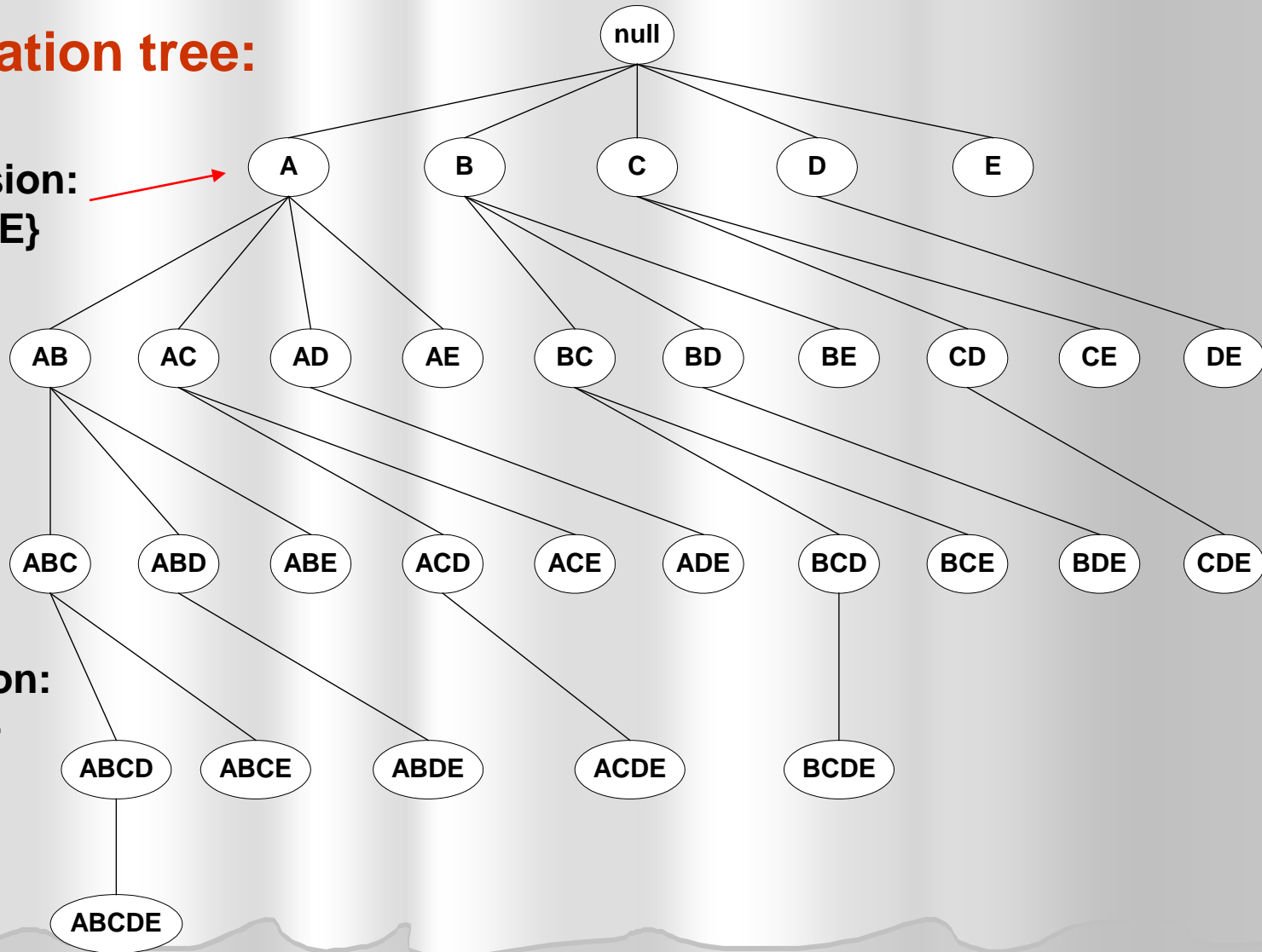# Association Rule Mining

Poonam Goyal
Computer Science
BITS, Pilani

# Tree Projection

**Set enumeration tree:**

**Possible Extension:**
E(A) = {B,C,D,E}

**Possible Extension:**
E(ABC) = {D,E}

# Tree Projection

- Items are listed in lexicographic order
- Each node P stores the following information:
  - Itemset for node P
  - List of possible lexicographic extensions of P: E(P)
  - Pointer to projected database of its ancestor node
  - Bitvector containing information about which transactions in the projected database contain the itemset

# Projected Database

**Original Database:**

| 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} |

**Projected Database for node A:**

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

For each transaction T, projected transaction at node A is T $\cap$ E(A)

# 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 $L \rightarrow \varnothing$ and $\varnothing \rightarrow L$)
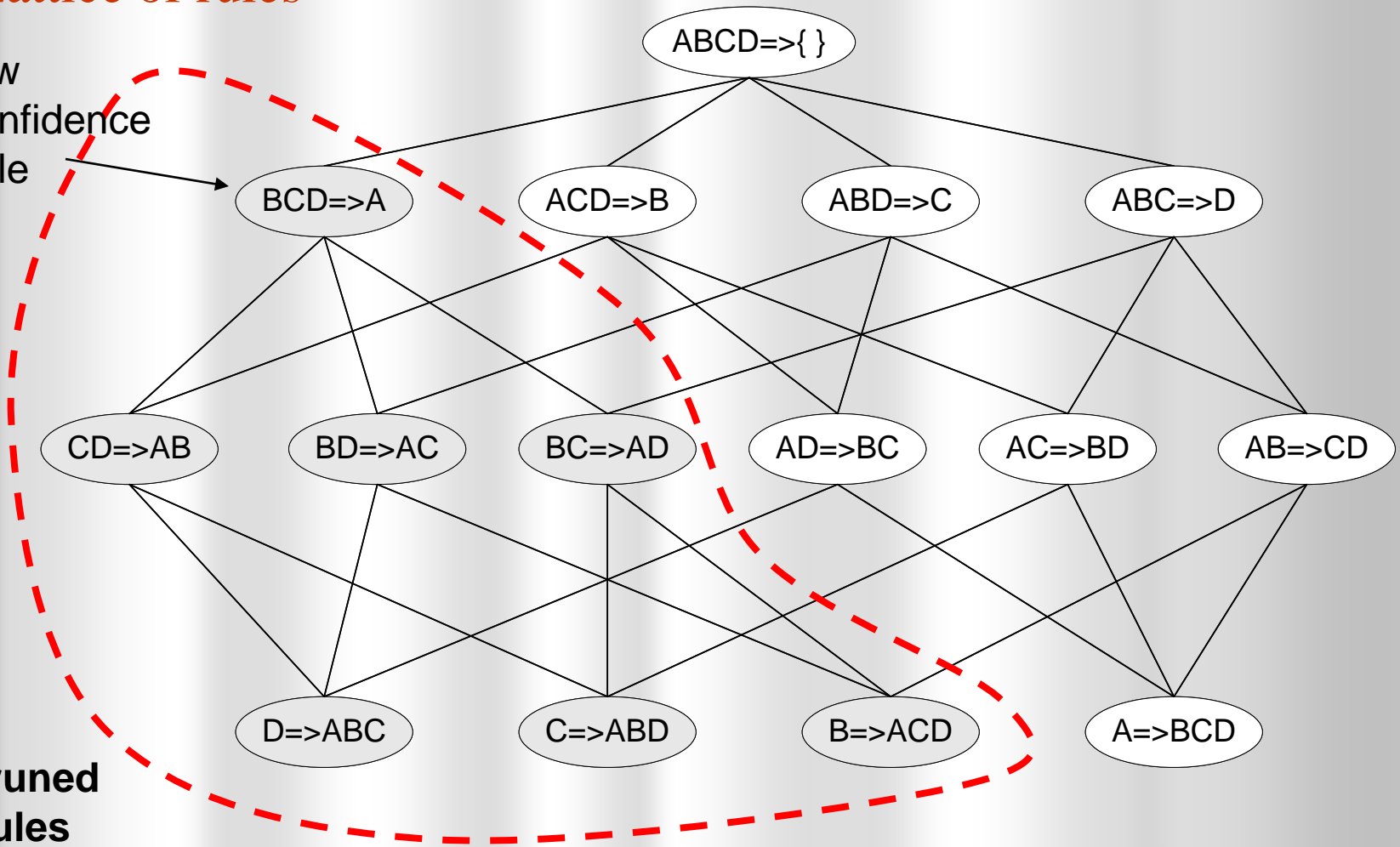
# Rule Generation

- 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}:

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

    - Confidence is anti-monotone w.r.t. number of items on the RHS of the rule
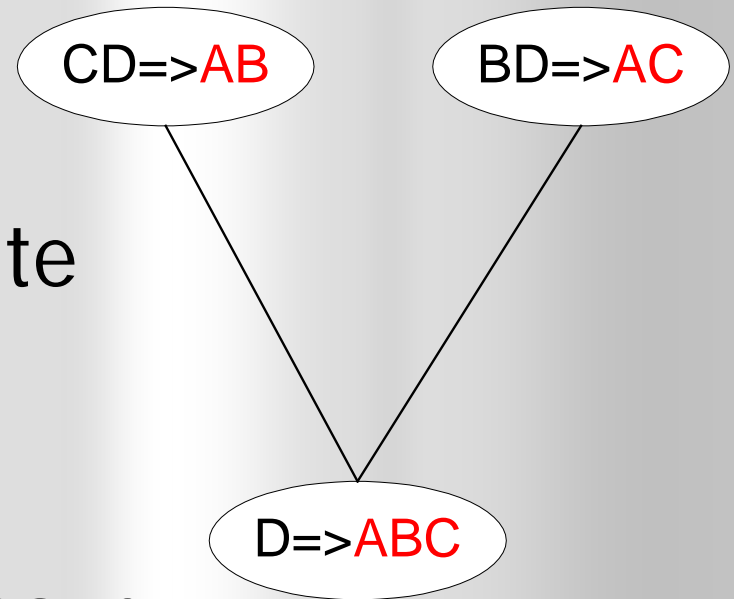
# Rule Generation for Apriori Algorithm

Lattice of rules

Low
Confidence
Rule

**Pruned
Rules**

ABCD=>{ }

BCD=>A   ACD=>B   ABD=>C   ABC=>D

CD=>AB   BD=>AC   BC=>AD   AD=>BC   AC=>BD   AB=>CD

D=>ABC   C=>ABD   B=>ACD   A=>BCD

# 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 AD=>BC does not have high confidence

```
  CD=>AB          BD=>AC


            D=>ABC
```
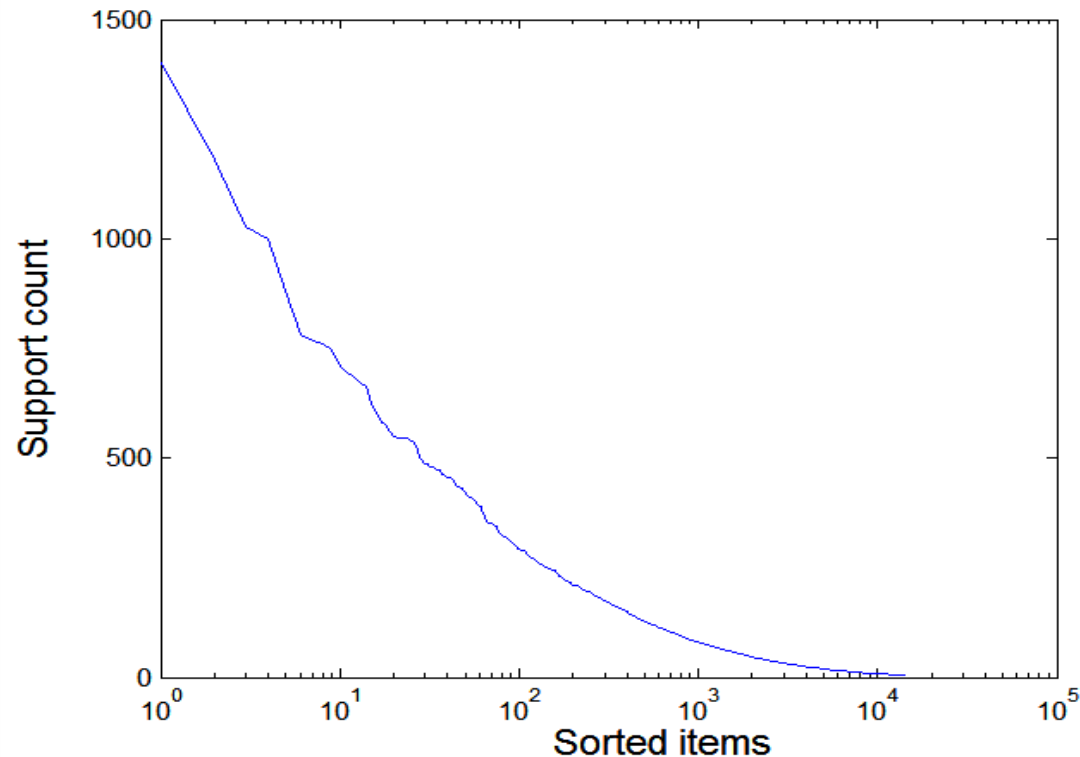
# Effect of Support Distribution

- Many real data sets have skewed support distribution

**Support distribution of a retail data set**

# 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(Milk)=5%,          MS(Coke) = 3%,
              MS(Broccoli)=0.1%,           MS(Salmon)=0.5%

  MS({Milk, Broccoli}) = min (MS(Milk), MS(Broccoli))
                                  = 0.1%

  Challenge: Support is no longer anti-monotone

  - Suppose:  Support(Milk, Coke) = 1.5% and
                      Support(Milk, Coke, Broccoli) = 0.5%

  - {Milk,Coke} is infrequent but {Milk,Coke,Broccoli} is frequent

# Multiple Minimum Support (Liu 1999)

- Order the items according to their minimum support (in ascending order)
  - e.g.:  MS(Milk)=5%,      MS(Coke) = 3%,
           MS(Broccoli)=0.1%,    MS(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(1)
          where MS(1) is $\min_i(\, MS(i)\, )$
  - $C_2$ : candidate itemsets of size 2 is generated from $F_1$ instead of $L_1$
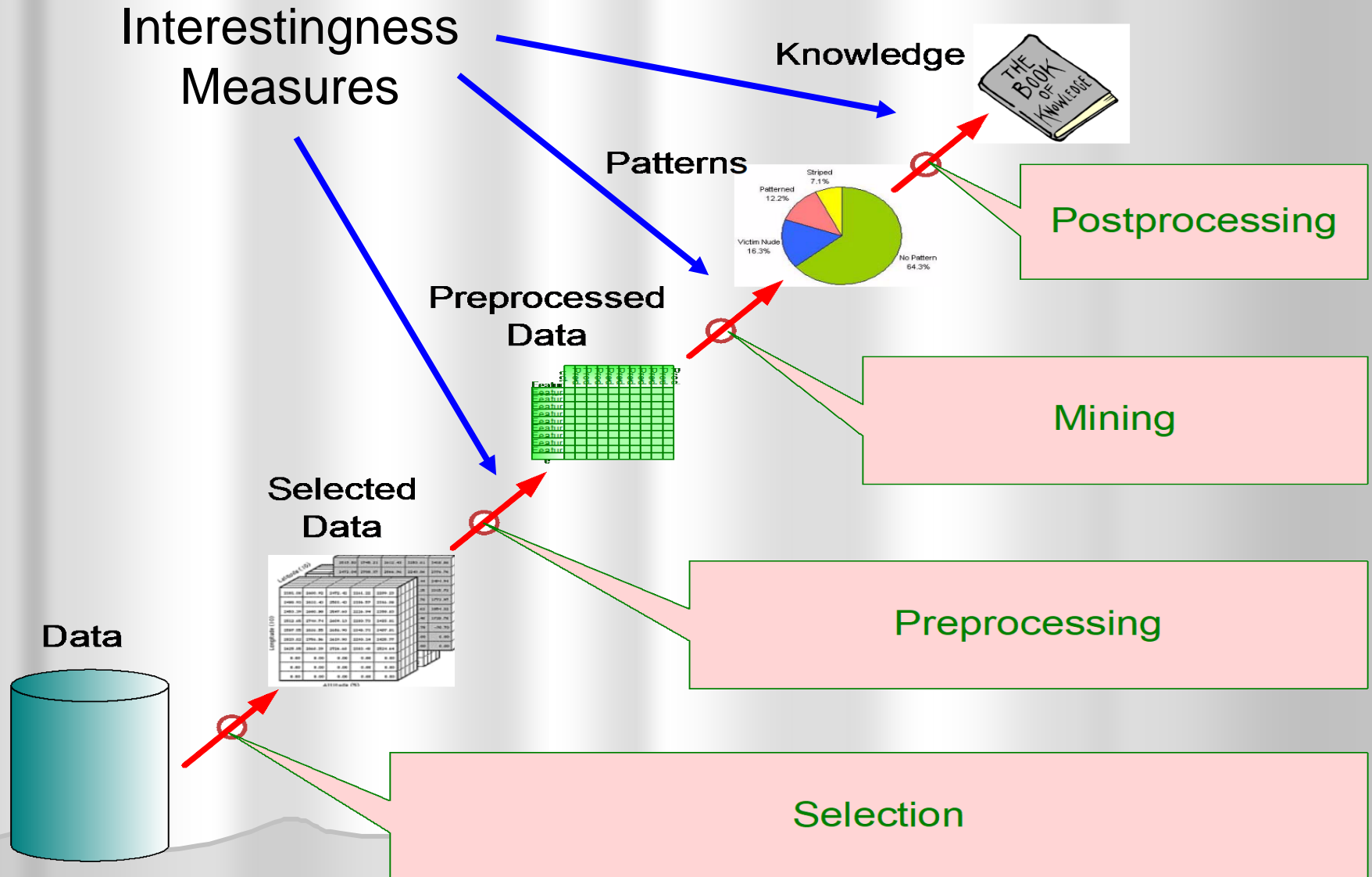
# Multiple Minimum Support (Liu 1999)

Modifications to Apriori:

- In traditional Apriori,
  - A candidate (k+1)-itemset is generated by merging two frequent itemsets of size k
  - The candidate is pruned if it contains any infrequent subsets of size k
- Pruning step has to be modified:
  - Prune only if subset contains the first item

  e.g.:  Candidate={Broccoli, Coke, Milk}   (ordered according to minimum support)

  {Broccoli, Coke} and {Broccoli, Milk} are frequent but {Coke, Milk} is infrequent
  - Candidate is not pruned because {Coke,Milk} does not contain the first item, i.e., Broccoli.

# 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

# 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 $\underline{X}$ and $\overline{Y}$
$f_{01}$: support of $\overline{X}$ and $\underline{Y}$
$f_{00}$: 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)
  - $P(S \wedge B) = 420/1000 = 0.42$
  - $P(S) \times P(B) = 0.6 \times 0.7 = 0.42$

  - $P(S \wedge B) = P(S) \times P(B) \Rightarrow$ 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

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

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

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

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

# Example: Lift/Interest

|  | Coffee | $\overline{\text{Coffee}}$ |  |
|---|---|---|---|
| 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

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

# Drawback of Lift/Interest

|   | Y | $\overline{Y}$ |   |
|---|---|---|---|
| X | 10 | 0 | 10 |
| $\overline{X}$ | 0 | 90 | 90 |
|   | 10 | 90 | 100 |

|   | Y | $\overline{Y}$ |   |
|---|---|---|---|
| X | 90 | 0 | 90 |
| $\overline{X}$ | 0 | 10 | 10 |
|   | 90 | 10 | 100 |

$$Lift = \frac{0.1}{(0.1)(0.1)} = 10$$

$$Lift = \frac{0.9}{(0.9)(0.9)} = 1.11$$

**Statistical independence:**

**If P(X,Y)=P(X)P(Y)  => Lift = 1**