# CS 484
# Data Mining

Association Rule Mining 3

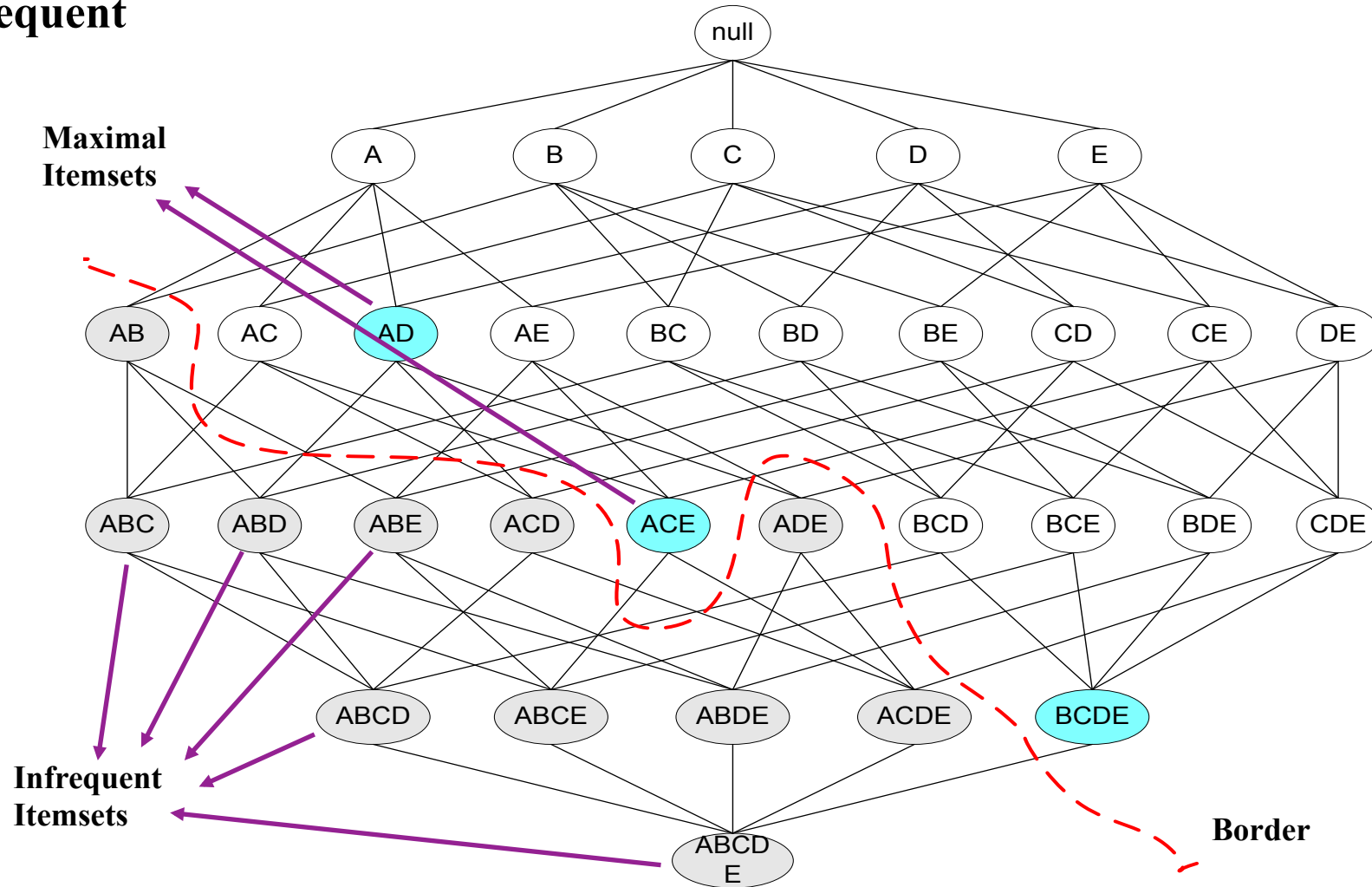# Compact Representation of Frequent Itemsets

- 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 | 0 | 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 | 0 | 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 | 0 | 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 | 0 | 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 | 0 | 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 | 1 | 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 | 1 | 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 | 1 | 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 | 1 | 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 | 1 | 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**

# Closed Itemset

- An itemset is closed if none of its immediate supersets has the same support as the itemset. Using the closed itemset support, we can find the support for the non-closed itemsets.
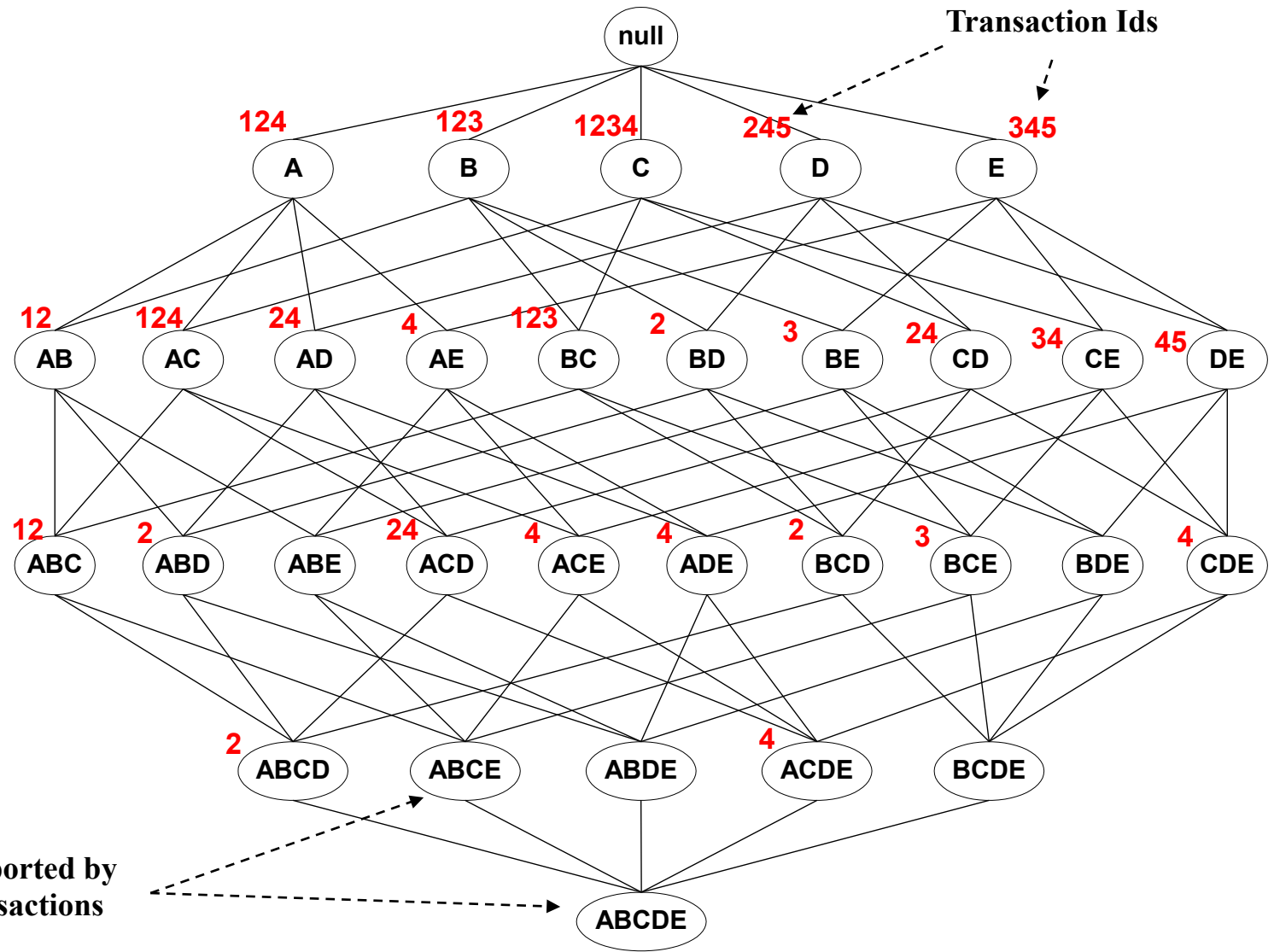
| TID | Items |
|-----|-------|
| 1 | {A,B} |
| 2 | {B,C,D} |
| 3 | {A,B,C,D} |
| 4 | {A,B,D} |
| 5 | {A,B,C,D} |

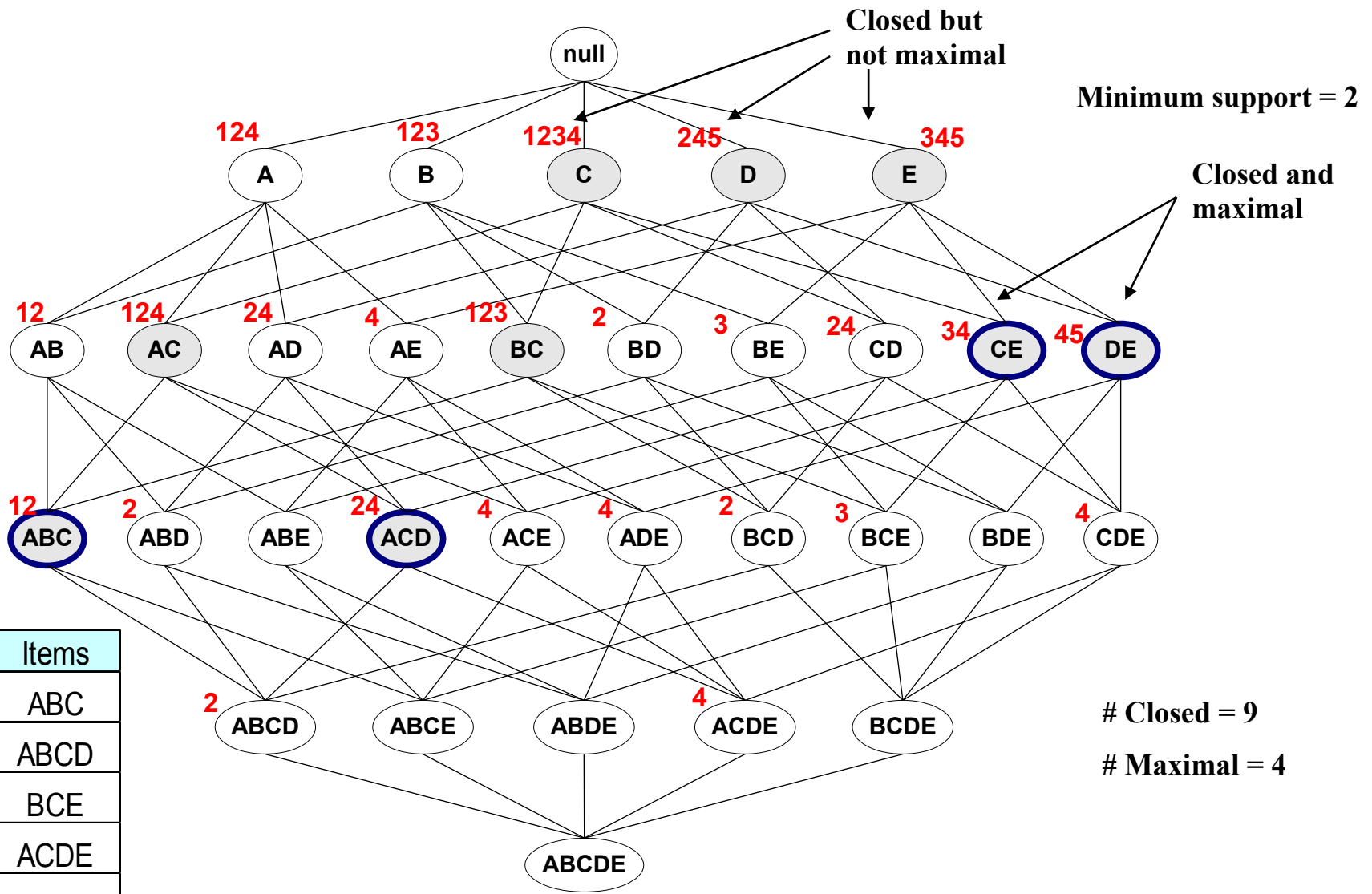| Itemset | Support |
|---------|---------|
| {A} | 4 |
| {B} | 5 |
| {C} | 3 |
| {D} | 4 |
| {A,B} | 4 |
| {A,C} | 2 |
| {A,D} | 3 |
| {B,C} | 3 |
| {B,D} | 4 |
| {C,D} | 3 |

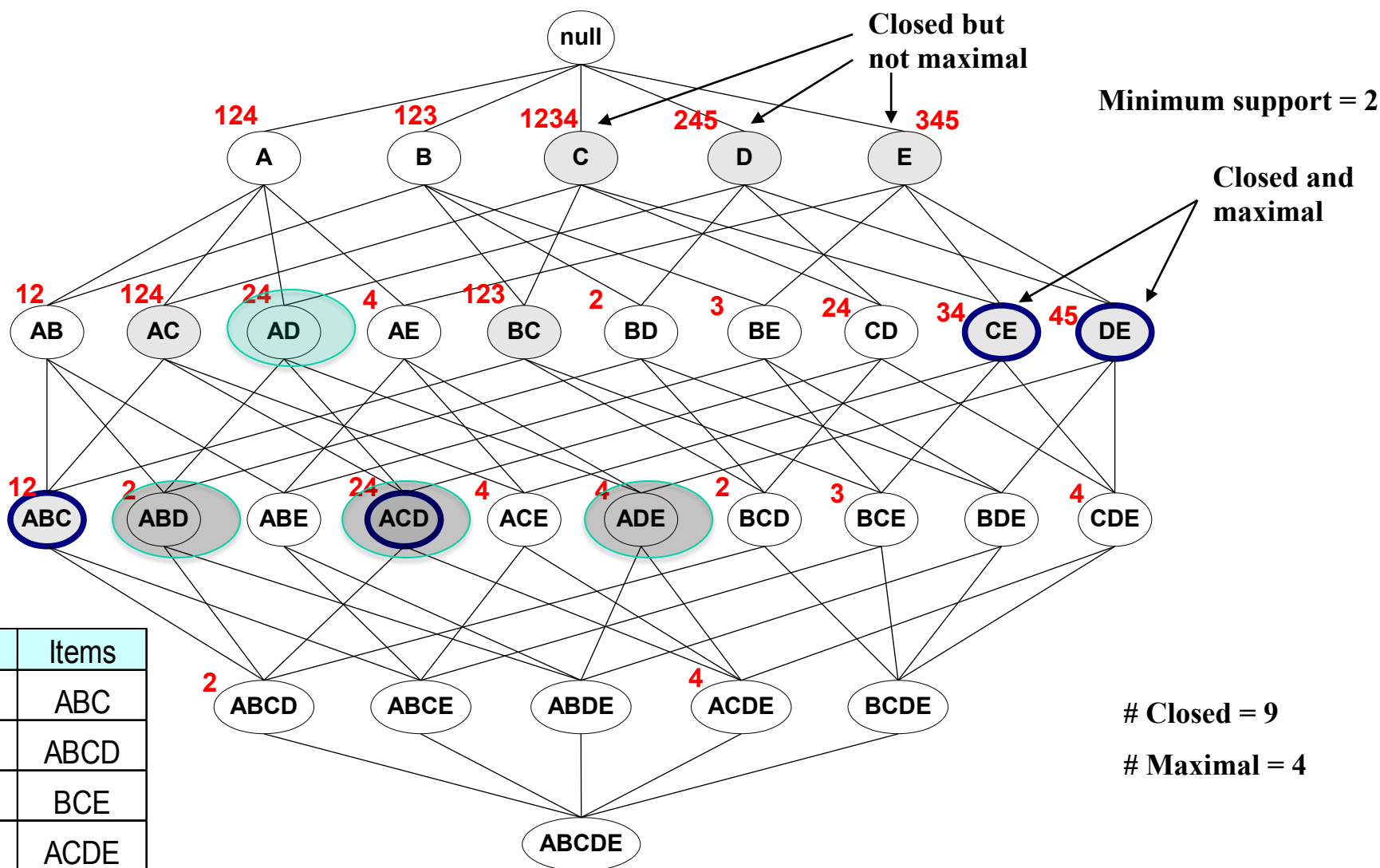| Itemset | Support |
|---------|---------|
| {A,B,C} | 2 |
| {A,B,D} | 3 |
| {A,C,D} | 2 |
| {B,C,D} | 3 |
| {A,B,C,D} | 2 |

# Maximal vs Closed Itemsets

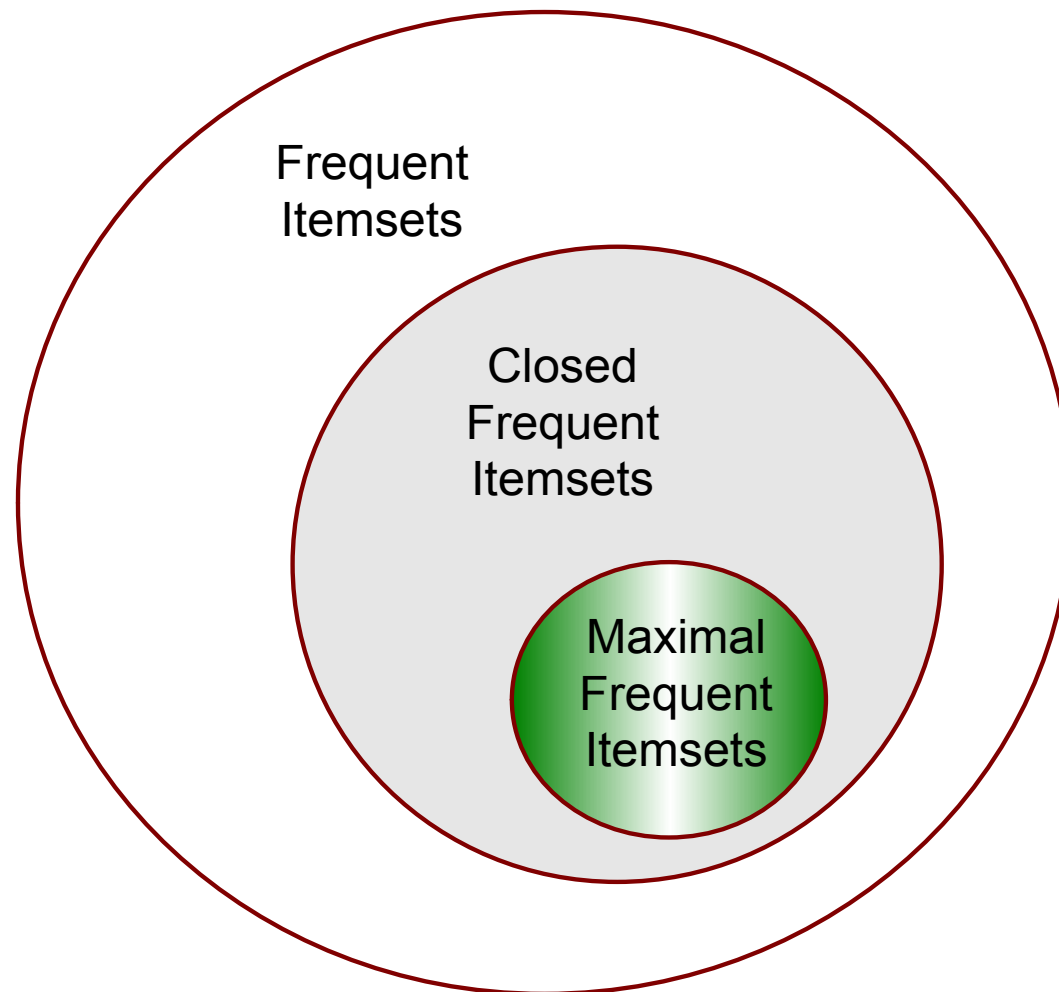| TID | Items |
|-----|-------|
| 1 | ABC |
| 2 | ABCD |
| 3 | BCE |
| 4 | ACDE |
| 5 | DE |

**Transaction Ids**

null

124 **A**   123 **B**   1234 **C**   245 **D**   345 **E**

12 **AB**   124 **AC**   24 **AD**   4 **AE**   123 **BC**   2 **BD**   3 **BE**   24 **CD**   34 **CE**   45 **DE**

12 **ABC**   2 **ABD**   **ABE**   24 **ACD**   4 **ACE**   4 **ADE**   2 **BCD**   3 **BCE**   **BDE**   4 **CDE**

2 **ABCD**   **ABCE**   **ABDE**   4 **ACDE**   **BCDE**

**ABCDE**

**Not supported by
any transactions**

# Maximal vs Closed Frequent Itemsets



Closed but not maximal

Closed and maximal

Minimum support = 2

null

124  A
123  B
1234 C
245  D
345  E

12  AB
124 AC
24  AD
4   AE
123 BC
2   BD
3   BE
24  CD
34  CE
45  DE

12 ABC
2  ABD
24 ACD
4  ACE
4  ADE
2  BCD
3  BCE
BDE
4  CDE

2 ABCD
ABCE
ABDE
4 ACDE
BCDE

ABCDE

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

# Closed = 9

# Maximal = 4

# Determining support for non-closed itemsets

# Closed Frequent Itemset

- An itemset is closed frequent itemset if it is closed and it support is greater than or equal to "minsup".


- Useful for removing redundant rules
  - A rules X -> Y is redundant if there exists another rule X' -> Y' where X is a subset of X' and Y is a subset of Y', such that the support/confidence for both rules are identical
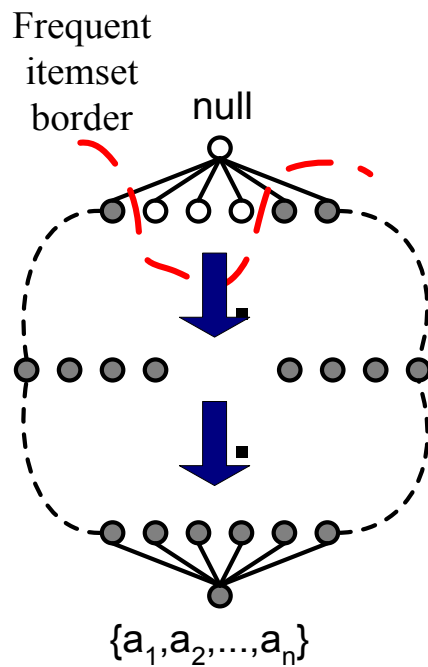
# Maximal vs Closed Itemsets

# Apriori Problems

- High I/O

- Poor performance for dense datasets because of increasing width of dimensions.
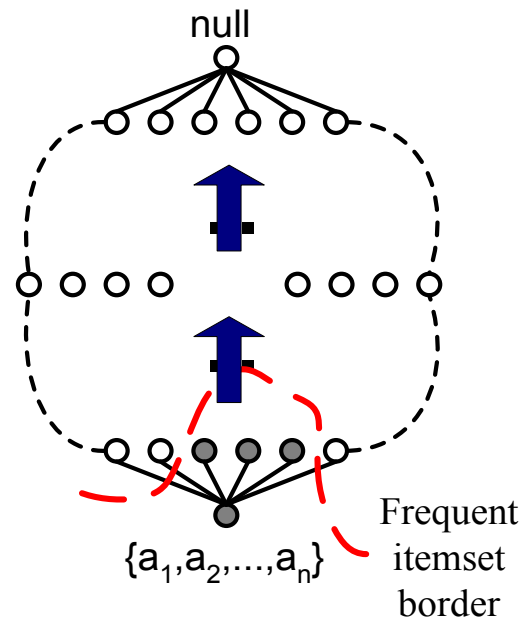
# Alternative Methods for Frequent Itemset Generation
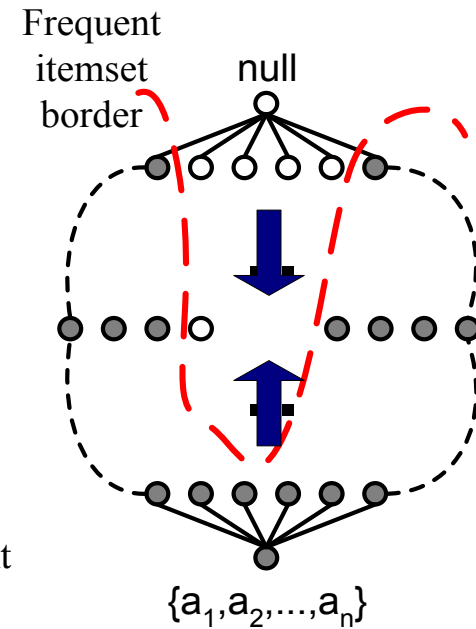
- Traversal of Itemset Lattice
  - General-to-specific vs Specific-to-general
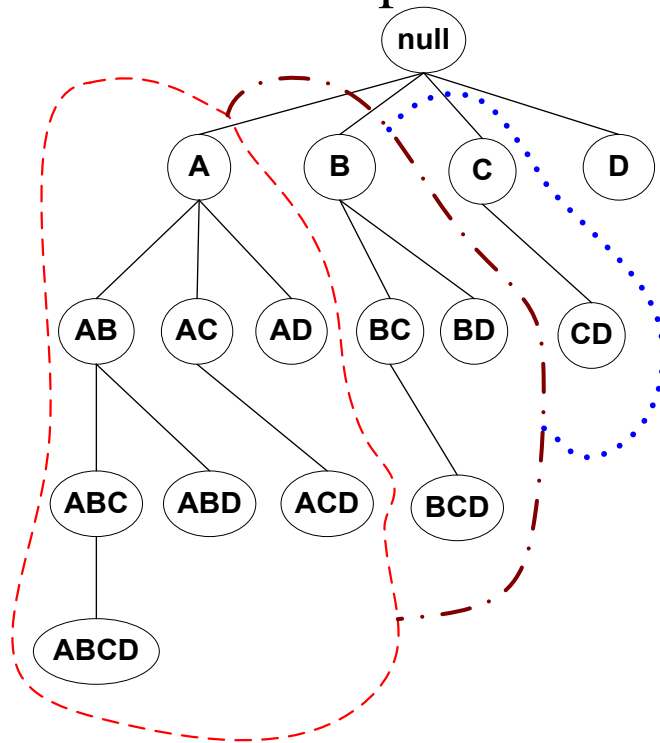


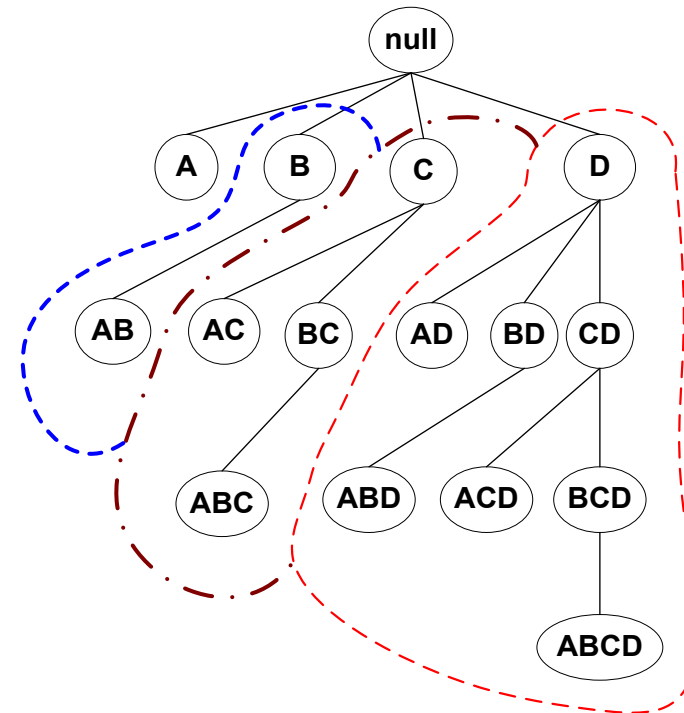(a) General-to-specific     (b) Specific-to-general     (c) Bidirectional

# Alternative Methods for Frequent Itemset Generation

- Traversal of Itemset Lattice
  - Equivalent Classes based on prefix or suffix
  - Consider frequent itemsets from these classes.
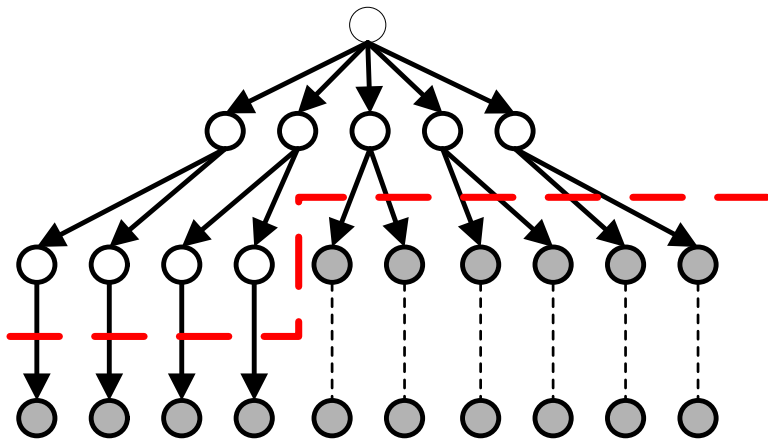


(a) Prefix tree                    (b) Suffix tree
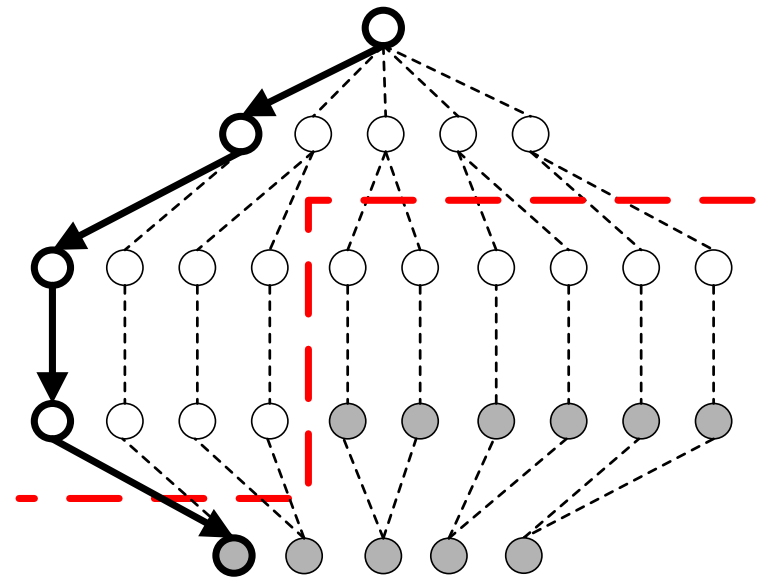
# Alternative Methods for Frequent Itemset Generation

- Traversal of Itemset Lattice
  - Breadth-first vs Depth-first
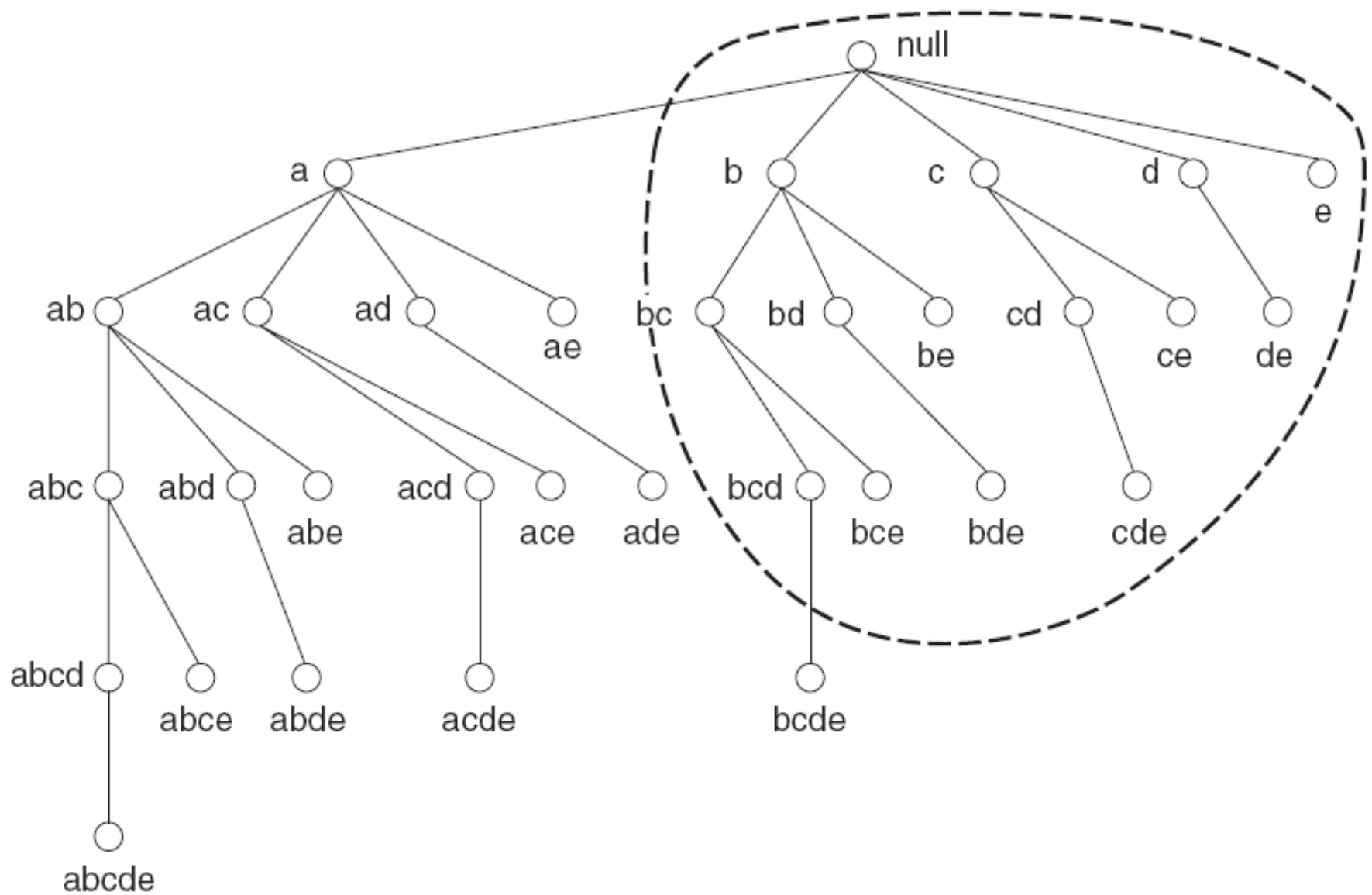


(a) Breadth first                    (b) Depth first

**Figure 6.22.** Generating candidate itemsets using the depth-first approach.

# 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 | | | | |

# 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

# Subjective Interestingness Measure

- ## Objective measure:
  - Rank patterns based on statistics computed from data
  - e.g., 21 measures of association (support, confidence, Laplace, Gini, mutual information, Jaccard, etc).

- ## Subjective measure:
  - Rank patterns according to user's interpretation
    - A pattern is subjectively interesting if it contradicts the expectation of a user
    - A pattern is subjectively interesting if it is actionable

# 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)
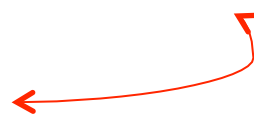  - 420 students know how to swim and bike (S,B)

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

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

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

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

Lift is equivalent to Interest Factor for binary variables.

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

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

Correlation for binary variables

# Interestingness Measure: Lift

- *play basketball* $\Rightarrow$ *eat cereal* [40%, 66.7%] is misleading

  - The overall % of students eating cereal is 75% > 66.7%.

- *play basketball* $\Rightarrow$ *not eat cereal* [20%, 33.3%] is more accurate, although with lower support and confidence

- Measure of dependent/correlated events: lift (= Interest Factor)

$$lift = \frac{P(A \cup B)}{P(A)P(B)}$$

|  | Basketball | Not basketball | Sum (row) |
|---|---|---|---|
| Cereal | 2000 | 1750 | 3750 |
| Not cereal | 1000 | 250 | 1250 |
| Sum(col.) | 3000 | 2000 | 5000 |

$$lift(B,C) = \frac{2000/5000}{3000/5000 * 3750/5000} = 0.89 \qquad lift(B, \neg C) = \frac{1000/5000}{3000/5000 * 1250/5000} = 1.33$$

# Example: Lift/Interest Factor

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

Association Rule: Tea $\rightarrow$ 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 Factor

|  | 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

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

# Properties of Objective Measures

- Symmetric/Asymmetric
- Scaling Property
- Inversion property
- Null Addition Property

# Property under Variable Permutation

|   | **B** | **$\overline{B}$** |
|---|---|---|
| **A** | p | q |
| **$\overline{A}$** | r | s |

$\Longrightarrow$

|   | **A** | **$\overline{A}$** |
|---|---|---|
| **B** | p | r |
| **$\overline{B}$** | q | s |

Does M(A,B) = M(B,A)?

Symmetric measures:

- support, lift, collective strength, cosine, Jaccard, etc

Asymmetric measures:

- confidence, conviction, Laplace, J-measure, etc

# Property under Row/Column Scaling

Grade-Gender Example (Mosteller, 1968):

|      | Male | Female |    |
|------|------|--------|----|
| High | 2    | 3      | 5  |
| Low  | 1    | 4      | 5  |
|      | 3    | 7      | 10 |

|      | Male | Female |    |
|------|------|--------|----|
| High | 4    | 30     | 34 |
| Low  | 2    | 40     | 42 |
|      | 6    | 70     | 76 |

2x      10x

Mosteller:

Underlying association should be independent of the relative number of male and female students in the samples

# Property under Inversion Operation

|  | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| Transaction 1 → | 1 | 0 | 0 | 1 | 0 | 0 |
| | 0 | 0 | 1 | 1 | 1 | 0 |
| | 0 | 0 | 1 | 1 | 1 | 0 |
| | 0 | 0 | 1 | 1 | 1 | 0 |
| | 0 | 1 | 1 | 0 | 1 | 1 |
| | 0 | 0 | 1 | 1 | 1 | 0 |
| | 0 | 0 | 1 | 1 | 1 | 0 |
| | 0 | 0 | 1 | 1 | 1 | 0 |
| | 0 | 0 | 1 | 1 | 1 | 0 |
| Transaction N → | 1 | 0 | 0 | 1 | 0 | 0 |
| | (a) | | (b) | | (c) | |

# Example: φ-Coefficient

- φ-coefficient is analogous to correlation coefficient for continuous variables

|   | Y | $\overline{Y}$ |   |
|---|---|---|---|
| X | 60 | 10 | 70 |
| $\overline{X}$ | 10 | 20 | 30 |
|   | 70 | 30 | 100 |

|   | Y | $\overline{Y}$ |   |
|---|---|---|---|
| X | 20 | 10 | 30 |
| $\overline{X}$ | 10 | 60 | 70 |
|   | 30 | 70 | 100 |

$$\phi = \frac{0.6 - 0.7 \times 0.7}{\sqrt{0.7 \times 0.3 \times 0.7 \times 0.3}}$$
$$= 0.5238$$

$$\phi = \frac{0.2 - 0.3 \times 0.3}{\sqrt{0.7 \times 0.3 \times 0.7 \times 0.3}}$$
$$= 0.5238$$

φ Coefficient is the same for both tables

# Property under Null Addition

|   | **B** | **$\overline{B}$** |
|---|---|---|
| **A** | p | q |
| **$\overline{A}$** | r | s |

$\Longrightarrow$

|   | **B** | **$\overline{B}$** |
|---|---|---|
| **A** | p | q |
| **$\overline{A}$** | r | s + k |

Invariant measures:

- ◆ support, cosine, Jaccard, etc

Non-invariant measures:

- ◆ correlation, Gini, mutual information, odds ratio, etc

# Resources

- Good summary of interestingness measures:

http://michael.hahsler.net/research/
association_rules/measures.html