Data Mining Association Analysis: Basic Concepts and Algorithms

Lecture Notes for Chapter 6

Association Analysis

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

| TID | Items |
|-----|---------------------------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

Example of Association Rules

```
{Diaper} \rightarrow {Beer},

{Milk, Bread} \rightarrow {Eggs,Coke},

{Beer, Bread} \rightarrow {Milk},
```

Implication means co-occurrence, not causality!

Association Rule Discovery: Application 1

- Marketing and Sales Promotion:
 - Let the rule discovered be {Bagels, ... } --> {Potato Chips}
 - Potato Chips as consequent => Can be used to determine what should be done to boost its sales.
 - Bagels in the antecedent => Can be used to see which products would be affected if the store discontinues selling bagels.
 - Bagels in antecedent and Potato chips in consequent => Can be used to see what products should be sold with Bagels to promote sale of Potato chips!

Association Rule Discovery: Application 2

- Supermarket shelf management.
 - Goal: To identify items that are bought together by sufficiently many customers.
 - Approach: Process the point-of-sale data collected with barcode scanners to find dependencies among items.
 - A classic rule --
 - If a customer buys diaper and milk, then he is very likely to buy beer.
 - So, don't be surprised if you find six-packs stacked next to diapers!

Association Rule Discovery: Application 3

• Inventory Management:

- Goal: A consumer appliance repair company wants to anticipate the nature of repairs on its consumer products and keep the service vehicles equipped with right parts to reduce on number of visits to consumer households.
- Approach: Process the data on tools and parts required in previous repairs at different consumer locations and discover the co-occurrence patterns.

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 (σ)

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

Support

- Fraction of transactions that contain an itemset
- E.g. $s(\{Milk, Bread, Diaper\}) = 2/5$

Frequent Itemset

 An itemset whose support is greater than or equal to a *minsup* threshold

| TID | Items |
|-----|---------------------------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

Definition: Association Rule

Association Rule

- An implication expression of the form
 X → Y, where X and Y are itemsets
- Example:{Milk, Diaper} → {Beer}

| TID | Items |
|-----|---------------------------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 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:

 $\{Milk, Diaper\} \Rightarrow Beer$

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

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

Association Rule Mining Task

- Given a set of transactions T, the goal of association rule mining is to find all rules having
 - support ≥ minsup threshold
 - confidence ≥ 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, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

Example of Rules:

```
{Milk, Diaper} \rightarrow {Beer} (s=0.4, c=0.67)
{Milk, Beer} \rightarrow {Diaper} (s=0.4, c=1.0)
{Diaper, Beer} \rightarrow {Milk} (s=0.4, c=0.67)
{Beer} \rightarrow {Milk, Diaper} (s=0.4, c=0.67)
{Diaper} \rightarrow {Milk, Beer} (s=0.4, c=0.5)
{Milk} \rightarrow {Diaper, Beer} (s=0.4, c=0.5)
```

Observations:

- All the above rules are binary partitions of the same itemset: {Milk, Diaper, Beer}
- Rules originating from the same itemset have identical support but can have different confidence
- Thus, we may decouple the support and confidence requirements

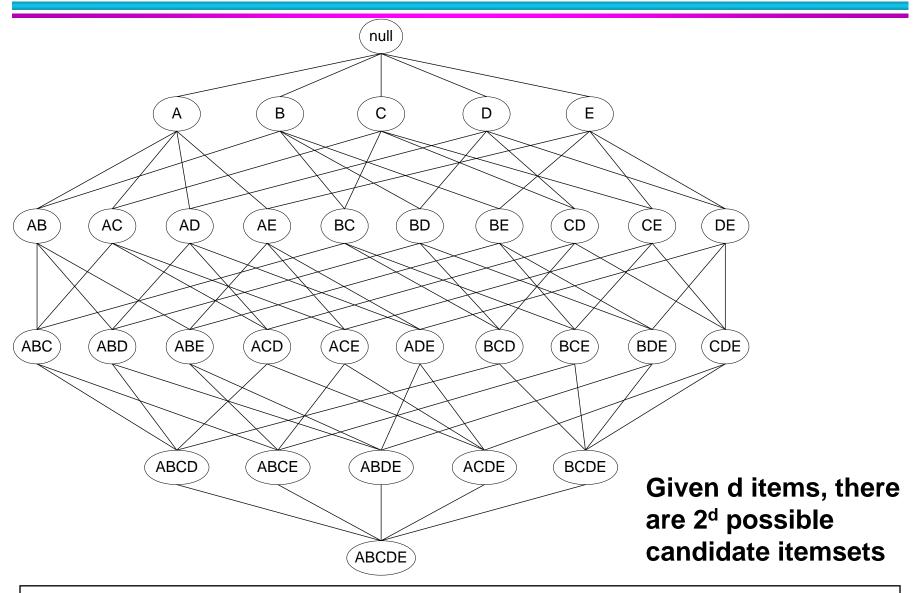
Mining Association Rules

- Two-step approach:
 - 1. Frequent Itemset Generation
 - Generate all itemsets whose support ≥ minsup

2. Rule Generation

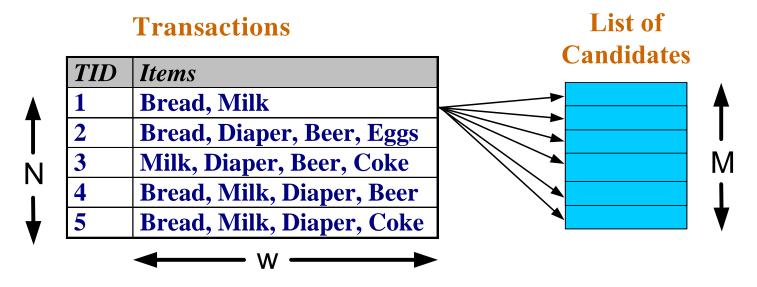
- Generate high confidence rules from each frequent itemset,
 where each rule is a binary partitioning of a frequent itemset
- Frequent itemset generation is still computationally expensive

Frequent Itemset Generation



Frequent Itemset Generation

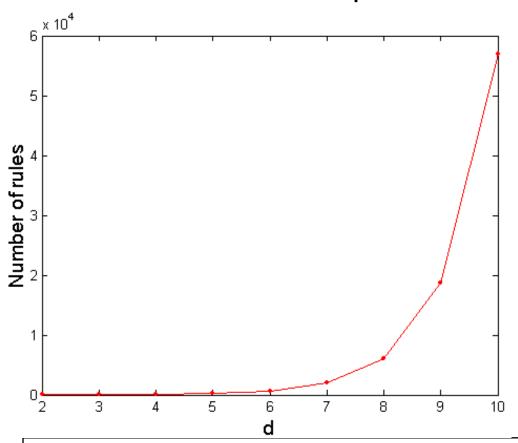
- Brute-force approach:
 - Each itemset in the lattice is a candidate frequent itemset
 - Count the support of each candidate by scanning the database



- Match each transaction against every candidate
- Complexity ~ O(NMw) => Expensive since M = 2^d !!!

Computational Complexity

- Given d unique items:
 - Total number of itemsets = 2^d
 - Total number of possible association rules:



$$R = \sum_{k=1}^{d-1} \begin{bmatrix} d \\ k \end{bmatrix} \times \sum_{j=1}^{d-k} \begin{pmatrix} d-k \\ j \end{bmatrix}$$
$$= 3^{d} - 2^{d+1} + 1$$

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
 - Used by DHP and vertical-based 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

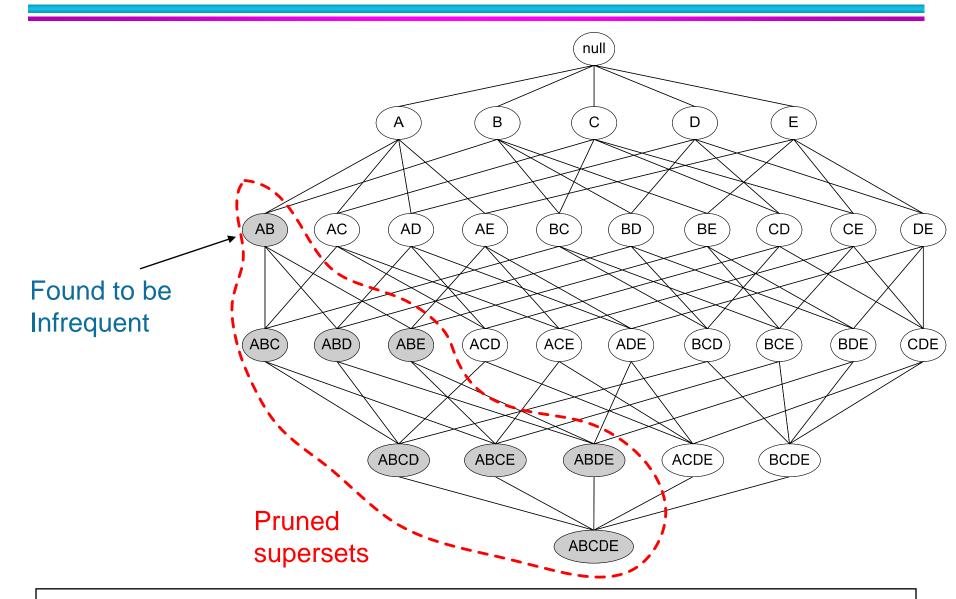
Reducing Number of Candidates

- 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 |
| Beer | 3 |
| Diaper | 4 |
| Eggs | 1 |

Items (1-itemsets)



| Itemset | Count |
|----------------|-------|
| {Bread,Milk} | 3 |
| {Bread,Beer} | 2 |
| {Bread,Diaper} | 3 |
| {Milk,Beer} | 2 |
| (Milk,Diaper) | 3 |
| {Beer,Diaper} | 3 |

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3



Triplets (3-itemsets)

| If every subset is considered, |
|--|
| ${}^{6}C_{1} + {}^{6}C_{2} + {}^{6}C_{3} = 41$ |
| With support-based pruning, |
| 6 + 6 + 1 = 13 |

| Itemset | Count |
|---------------------|-------|
| {Bread,Milk,Diaper} | 3 |



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

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

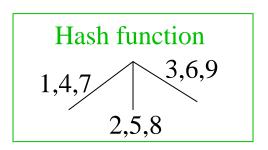
Generate Hash Tree

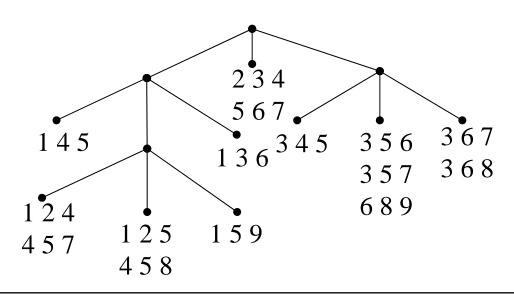
Suppose you have 15 candidate itemsets of length 3:

{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4}, {5 6 7}, {3 4 5}, {3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}

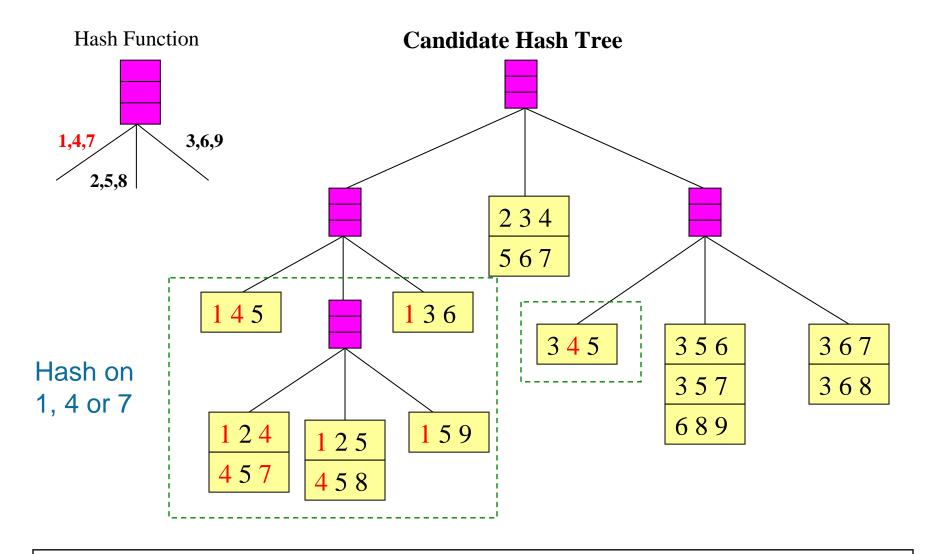
You need:

- Hash function
- Max leaf size: max number of itemsets stored in a leaf node (if number of candidate itemsets exceeds max leaf size, split the node)

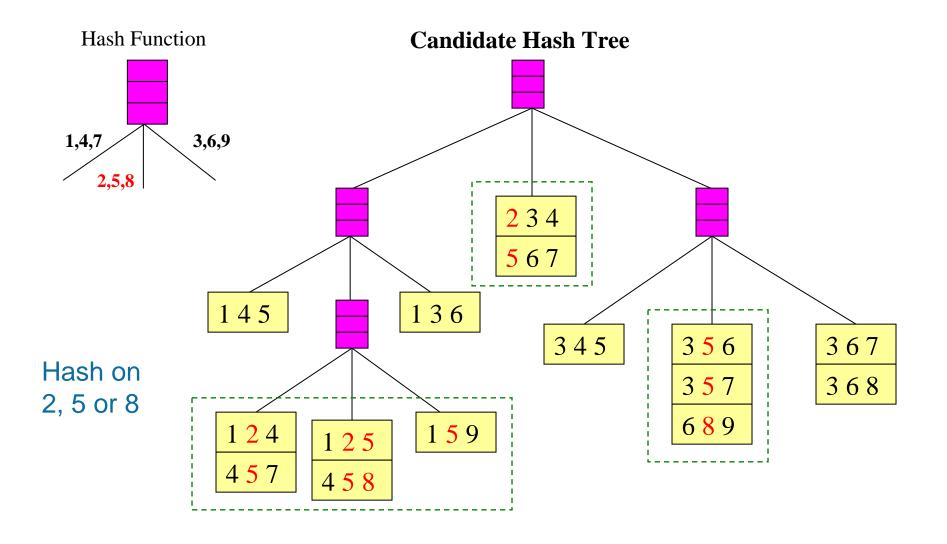




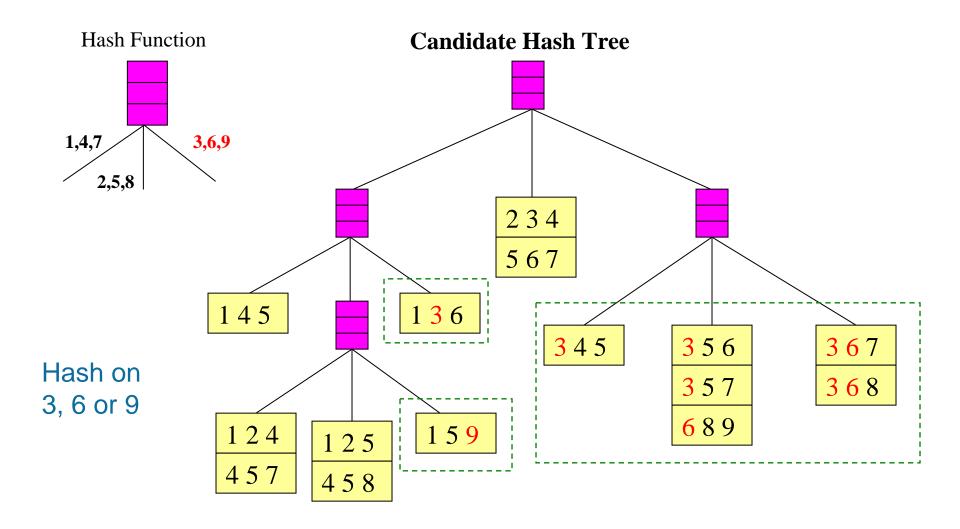
Association Rule Discovery: Hash tree



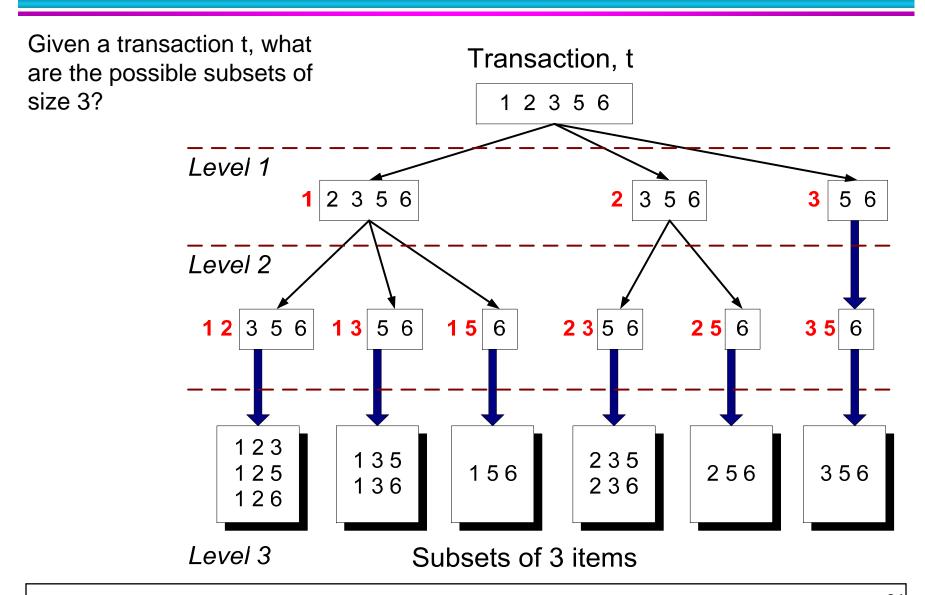
Association Rule Discovery: Hash tree



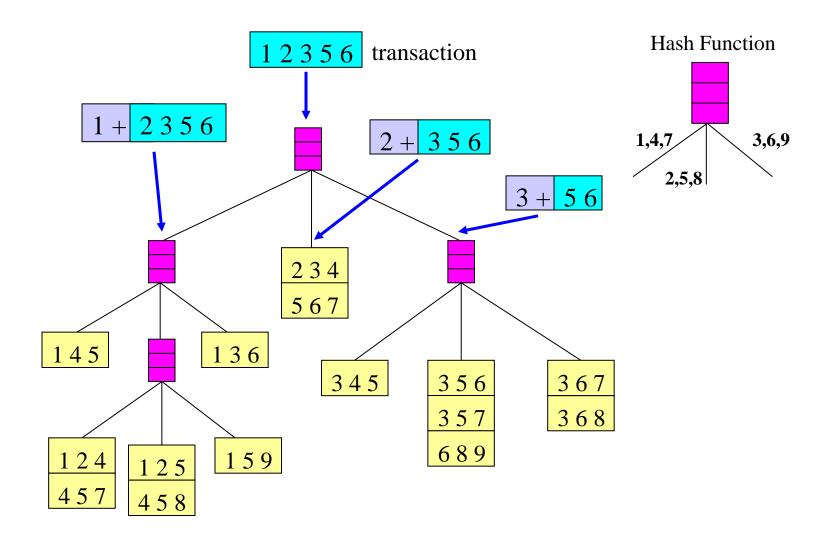
Association Rule Discovery: Hash tree



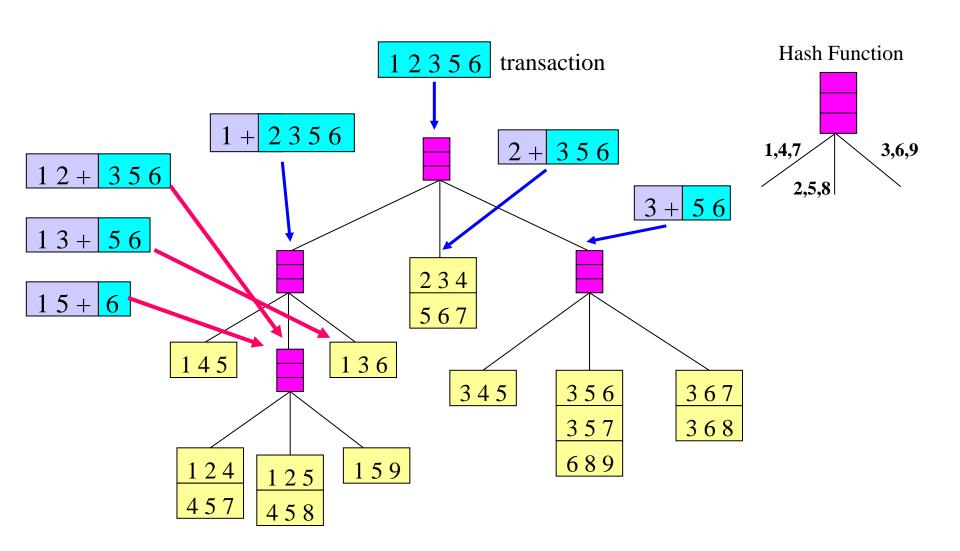
Subset Operation



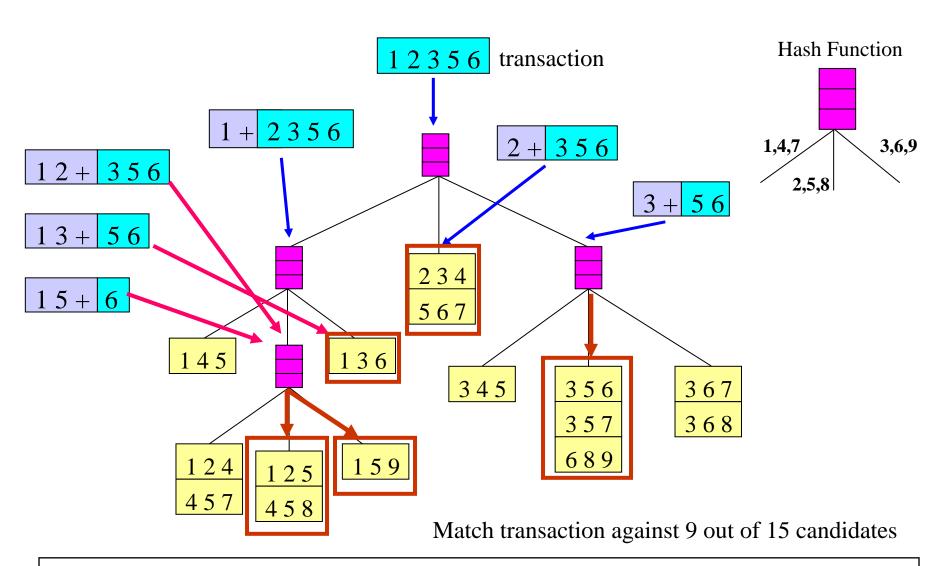
Subset Operation Using Hash Tree



Subset Operation Using Hash Tree



Subset Operation Using Hash Tree



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

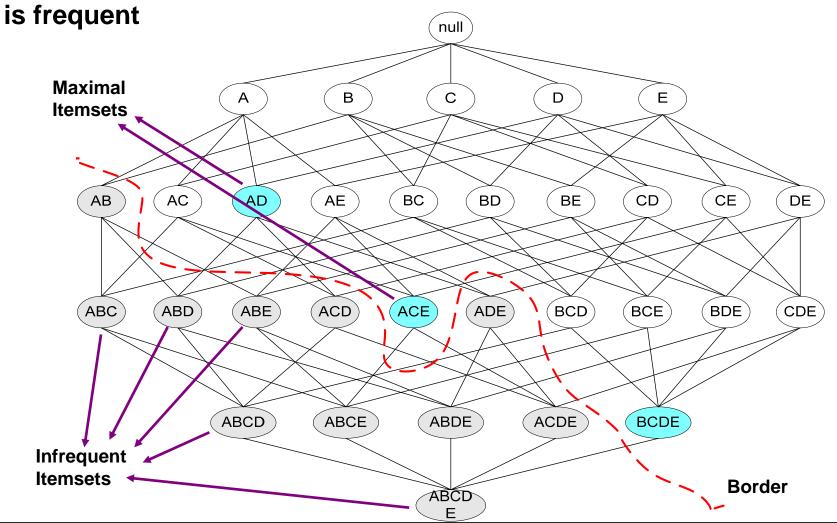
| TID | A 1 | A2 | A3 | A4 | A5 | A6 | A7 | A8 | A9 | A10 | B1 | B2 | В3 | B4 | B5 | B6 | B7 | B8 | B9 | B10 | C1 | C2 | C3 | C4 | C 5 | C6 | C7 | C8 | C9 | C10 |
|-----|------------|----|-----------|-----------|-----------|-----------|----|-----------|-----------|-----|----|----|----|-----------|----|----|----|----|----|-----|-----------|----|----|----|------------|----|-----------|----|----|-----|
| 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 =
$$3 \times \sum_{k=1}^{10} {10 \choose k}$$

Need a compact representation

Maximal Frequent Itemset

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



Closed Itemset

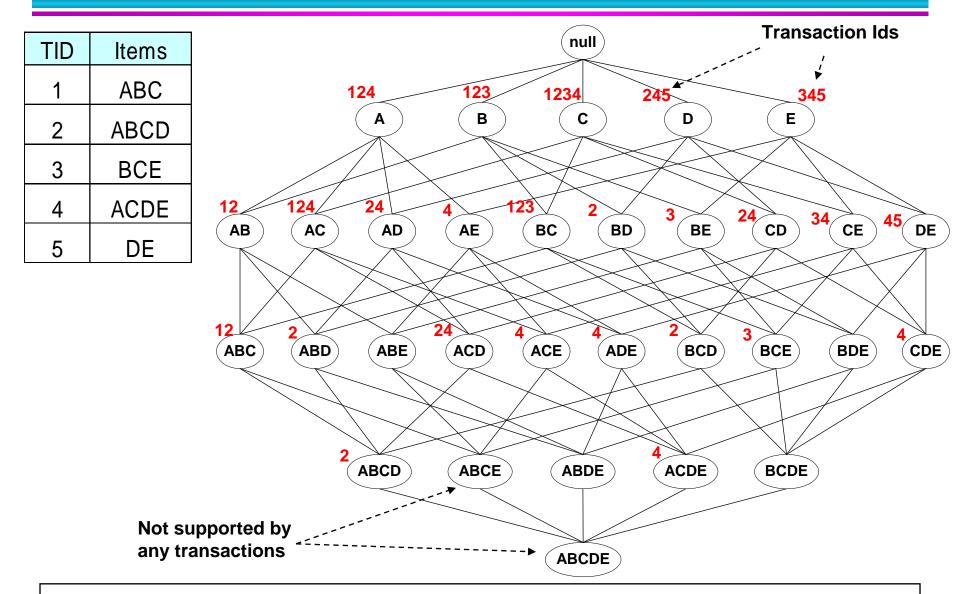
 An itemset is closed if none of its immediate supersets has the same support as the itemset

| TID | Items |
|-----|---------------|
| 1 | {A,B} |
| 2 | $\{B,C,D\}$ |
| 3 | $\{A,B,C,D\}$ |
| 4 | $\{A,B,D\}$ |
| 5 | $\{A,B,C,D\}$ |

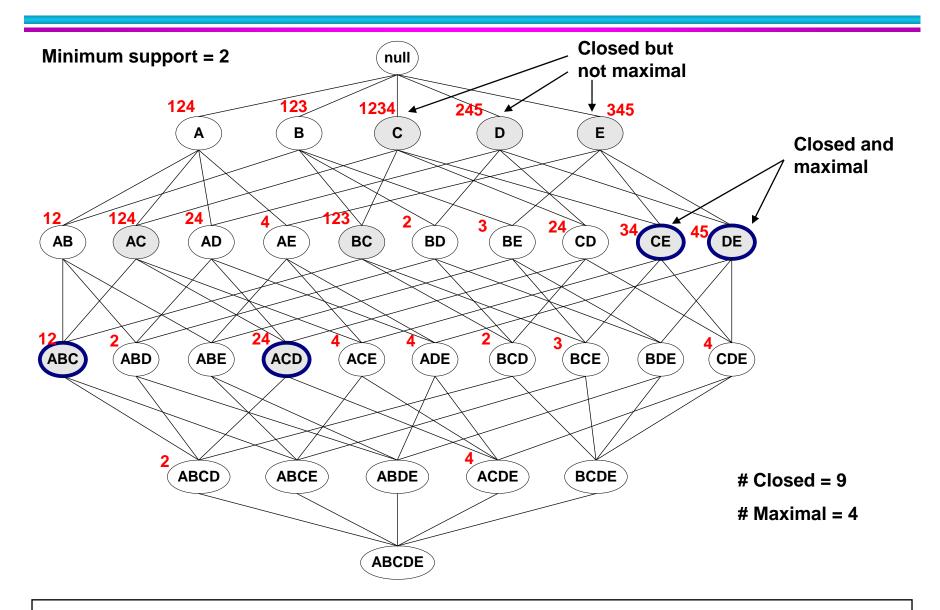
| Itemset | Support |
|---------|---------|
| {A} | 4 |
| {B} | 5 |
| {C} | 3 |
| {D} | 4 |
| {A,B} | 4 |
| {A,C} | 2 |
| {A,D} | 3 |
| {B,C} | 3 |
| {B,D} | 4 |
| {C,D} | 3 |

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

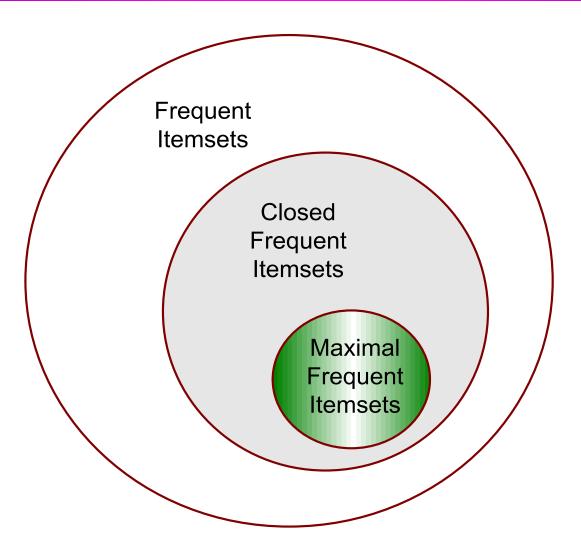
Maximal vs Closed Itemsets



Maximal vs Closed Frequent Itemsets



Maximal vs Closed Itemsets



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 \rightarrowD, ABD \rightarrowC, ACD \rightarrowB, BCD \rightarrowA, A \rightarrowBCD, B \rightarrowACD, C \rightarrowABD, D \rightarrowABC AB \rightarrowCD, AC \rightarrow BD, AD \rightarrow BC, BC \rightarrowAD, BD \rightarrowAC, CD \rightarrowAB,
```

 If |L| = k, then there are 2^k − 2 candidate association rules (ignoring L → Ø and Ø → L)

Rule Generation

- How to efficiently generate rules from frequent itemsets?
 - In general, confidence does not have an antimonotone property

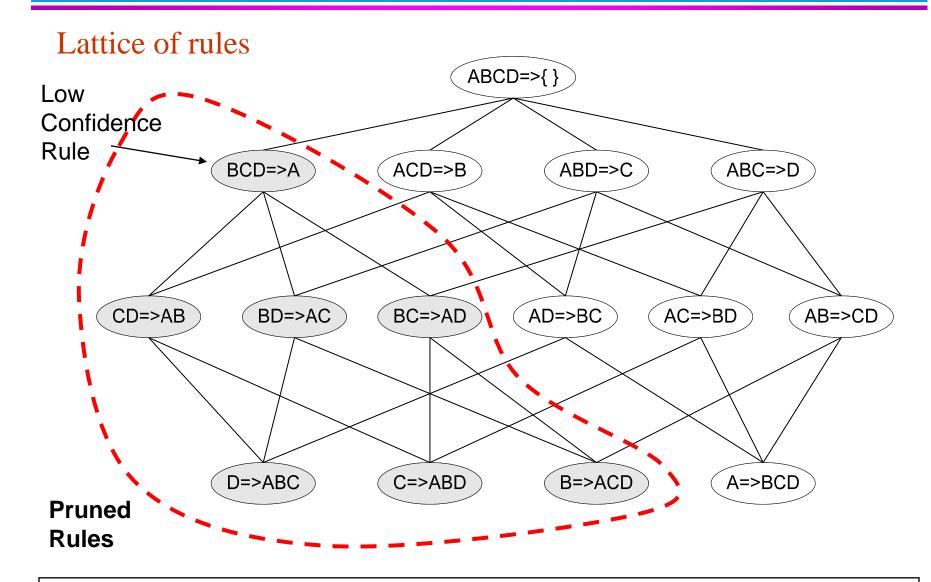
 $c(ABC \rightarrow D)$ can be larger or smaller than $c(AB \rightarrow D)$

- But confidence of rules generated from the same itemset has an anti-monotone property
- e.g., L = {A,B,C,D}:

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

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

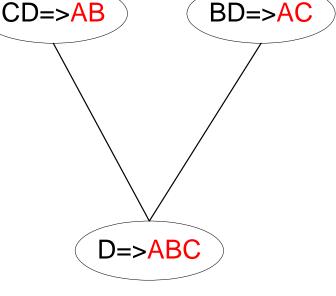
Rule Generation for Apriori Algorithm



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

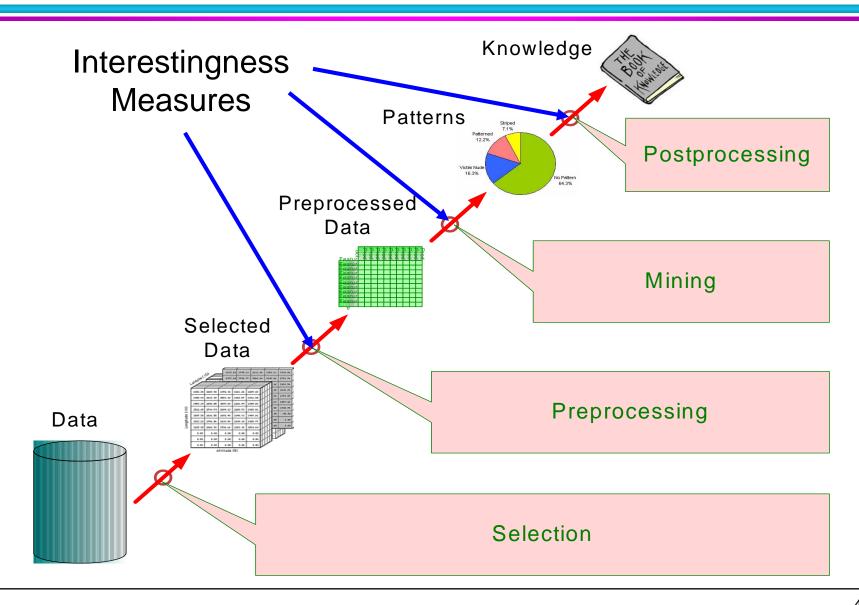


 Prune rule D=>ABC if its subset AD=>BC does not have high confidence

Pattern Evaluation

- Association rule algorithms tend to produce too many rules as the size and dimensionality of real commercial databases can be very large
- Easily end up with thousands or even millions of patterns
 - 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

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 $X \to Y$

| | Υ | 7 | |
|---|-----------------|-----------------|-----------------|
| X | f ₁₁ | f ₁₀ | f ₁₊ |
| X | f ₀₁ | f ₀₀ | f _{o+} |
| | f ₊₁ | f ₊₀ | Τ |

f₁₁: support of X and Y

 f_{10} : support of X and \overline{Y}

f₀₁: support of 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

Suppose we are interested in analyzing the relationship between people who drink tea and coffee

| | Coffee | Coffee | |
|-----|--------|--------|-----|
| Tea | 15 | 5 | 20 |
| 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

$$\Rightarrow$$
 P(Coffee|Tea) = 0.9375

A person is a tea drinker actually decreases her probability of being a coffee drinker from 0.9 to 0.75

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)
 - $P(S \land B) = 420/1000 = 0.42$
 - $P(S) \times P(B) = 0.6 \times 0.7 = 0.42$
 - $P(S \land B) = P(S) \times P(B) => Statistical independence$
 - $P(S \land B) > P(S) \times P(B) => Positively correlated$
 - P(S∧B) < P(S) × P(B) => 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 | Coffee | |
|-----|--------|--------|-----|
| Tea | 15 | 5 | 20 |
| Tea | 75 | 5 | 80 |
| | 90 | 10 | 100 |

Association Rule: Tea → Coffee

Confidence= P(Coffee|Tea) = 0.75but P(Coffee) = 0.9

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

Drawback of Lift & Interest

- We illustrate the limitation of interest factor with an example from the text mining
- Reasonable to assume that the association between a pair of words depends on the number of documents that contain both words
- Expect the words "data" and 'mining" to appear together more frequently than the words "complier" and "mining" in a collection of computer science articles

X=complier and Y=mining

| | Υ | Y | |
|---|----|----|-----|
| X | 10 | 0 | 10 |
| X | 0 | 90 | 90 |
| | 10 | 90 | 100 |

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

X=data and Y=mining

| | Υ | Y | |
|---|----|----|-----|
| X | 90 | 0 | 90 |
| X | 0 | 10 | 10 |
| | 90 | 10 | 100 |

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

Statistical independence:

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

| | # | Measure | Formula |
|------------------------------------|----|-------------------------------|---|
| There are lots of | 1 | ϕ -coefficient | $\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1 - P(A))(1 - P(B))}}$ $\sum_{j} \max_{k} P(A_{j}, B_{k}) + \sum_{k} \max_{j} P(A_{j}, B_{k}) - \max_{j} P(A_{j}) - \max_{k} P(B_{k})$ |
| measures proposed | 2 | Goodman-Kruskal's (λ) | $2-\max_{j}P(A_{j})-\max_{k}P(B_{k})$ |
| in the literature | 3 | Odds ratio (α) | $\frac{P(A,B)P(\overline{A},\overline{B})}{P(A,\overline{B})P(\overline{A},B)}$ |
| | 4 | Yule's Q | $\frac{P(A,B)P(\overline{AB}) - P(A,\overline{B})P(\overline{A},B)}{P(A,B)P(\overline{AB}) + P(A,\overline{B})P(\overline{A},B)} = \frac{\alpha - 1}{\alpha + 1}$ |
| | 5 | Yule's Y | $\frac{\sqrt{P(A,B)P(\overline{AB})} - \sqrt{P(A,\overline{B})P(\overline{A},B)}}{\sqrt{P(A,B)P(\overline{AB})} + \sqrt{P(A,\overline{B})P(\overline{A},B)}} = \frac{\sqrt{\alpha} - 1}{\sqrt{\alpha} + 1}$ |
| Some measures are good for certain | 6 | Kappa (κ) | $\frac{P(A,B)+P(\overline{A},\overline{B})-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A)P(B)-P(\overline{A})P(\overline{B})}$ |
| applications, but not | 7 | Mutual Information (M) | $\frac{\sum_{i} \sum_{j} P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\min(-\sum_{i} P(A_i) \log P(A_i), -\sum_{j} P(B_j) \log P(B_j))}$ |
| for others | 8 | J-Measure (J) | $\max \left(P(A,B) \log(\frac{P(B A)}{P(B)}) + P(A\overline{B}) \log(\frac{P(\overline{B} A)}{P(\overline{B})}), \right.$ |
| | | | $P(A,B)\log(rac{P(A B)}{P(A)}) + P(\overline{A}B)\log(rac{P(\overline{A} B)}{P(\overline{A})})\Big)$ |
| | 9 | Gini index (G) | $= \max \left(P(A)[P(B A)^2 + P(\overline{B} A)^2] + P(\overline{A})[P(B \overline{A})^2 + P(\overline{B} \overline{A})^2] \right)$ |
| What criteria should | | | $-P(B)^2-P(\overline{B})^2,$ |
| we use to determine | | | $P(B)[P(A B)^{2} + P(\overline{A} B)^{2}] + P(\overline{B})[P(A \overline{B})^{2} + P(\overline{A} \overline{B})^{2}]$ |
| whether a measure | | | $-P(A)^2 - P(\overline{A})^2$ |
| is good or bad? | 10 | Support (s) | P(A,B) |
| | 11 | Confidence (c) | $\max(P(B A), P(A B))$ |
| | 12 | Laplace (L) | $\max\left(rac{NP(A,B)+1}{NP(A)+2},rac{NP(A,B)+1}{NP(B)+2} ight)$ |
| What about Apriori- | 13 | Conviction (V) | $\max\left(rac{P(A)P(\overline{B})}{P(A\overline{B})},rac{P(B)P(\overline{A})}{P(B\overline{A})} ight)$ |
| style support based | 14 | Interest (I) | $\frac{P(A,B)}{P(A)P(B)}$ |
| pruning? How does | 15 | cosine (IS) | $\frac{P(A,B)}{\sqrt{P(A)P(B)}}$ |
| it affect these | 16 | Piatetsky-Shapiro's (PS) | P(A,B) - P(A)P(B) |
| measures? | 17 | Certainty factor (F) | $\max\left(rac{P(B A)-P(B)}{1-P(B)},rac{P(A B)-P(A)}{1-P(A)} ight)$ |
| | 18 | Added Value (AV) | $\max(P(B A) - P(B), P(A B) - P(A))$ |
| | 19 | Collective strength (S) | $\frac{P(A,B)+P(\overline{AB})}{P(A)P(B)+P(\overline{A})P(\overline{B})} \times \frac{1-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A,B)-P(\overline{AB})}$ |
| | 20 | Jaccard (ζ) | $\frac{P(A,B)}{P(A)+P(B)-P(A,B)}$ |
| | 21 | Klosgen (K) | $\sqrt{P(A,B)}\max(P(B A)-P(B),P(A B)-P(A))$ |

Properties of A Good Measure

- Piatetsky-Shapiro:
 - 3 properties a good measure M must satisfy:
 - -M(A,B) = 0 if A and B are statistically independent
 - M(A,B) increase monotonically with P(A,B) when P(A) and P(B) remain unchanged
 - M(A,B) decreases monotonically with P(A) [or P(B)]
 when P(A,B) and P(B) [or P(A)] remain unchanged

Consistency among objective measures

- Given the wide variety of measures available, it is reasonable to question whether the measures can produce similar ordering results when applied to a set of association patterns
- If the measures are consistent, then we can choose any one of them as our evaluation metric
- Otherwise, it is important to understand what their differences are in order to determine which measure is more suitable

Comparing Different Measures

10 examples of contingency tables:

| Example | f ₁₁ | f ₁₀ | f ₀₁ | f ₀₀ |
|---------|-----------------|-----------------|-----------------|-----------------|
| E1 | 8123 | 83 | 424 | 1370 |
| E2 | 8330 | 2 | 622 | 1046 |
| E3 | 9481 | 94 | 127 | 298 |
| E4 | 3954 | 3080 | 5 | 2961 |
| E5 | 2886 | 1363 | 1320 | 4431 |
| E6 | 1500 | 2000 | 500 | 6000 |
| E7 | 4000 | 2000 | 1000 | 3000 |
| E8 | 4000 | 2000 | 2000 | 2000 |
| E9 | 1720 | 7121 | 5 | 1154 |
| E10 | 61 | 2483 | 4 | 7452 |

Rankings of contingency tables using various measures:

| # | φ | λ | α | Q | Y | κ | M | J | G | 8 | c | L | V | I | IS | PS | F | AV | S | ζ | K |
|-----|----|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| E1 | 1 | 1 | 3 | 3 | 3 | 1 | 2 | 2 | 1 | 3 | 5 | 5 | 4 | 6 | 2 | 2 | 4 | 6 | 1 | 2 | 5 |
| E2 | 2 | 2 | 1 | 1 | 1 | 2 | 1 | 3 | 2 | 2 | 1 | 1 | 1 | 8 | 3 | 5 | 1 | 8 | 2 | 3 | 6 |
| E3 | 3 | 3 | 4 | 4 | 4 | 3 | 3 | 8 | 7 | 1 | 4 | 4 | 6 | 10 | 1 | 8 | 6 | 10 | 3 | 1 | 10 |
| E4 | 4 | 7 | 2 | 2 | 2 | 5 | 4 | 1 | 3 | 6 | 2 | 2 | 2 | 4 | 4 | 1 | 2 | 3 | 4 | 5 | 1 |
| E5 | 5 | 4 | 8 | 8 | 8 | 4 | 7 | 5 | 4 | 7 | 9 | 9 | 9 | 3 | 6 | 3 | 9 | 4 | 5 | 6 | 3 |
| E6 | 6 | 6 | 7 | 7 | 7 | 7 | 6 | 4 | 6 | 9 | 8 | 8 | 7 | 2 | 8 | 6 | 7 | 2 | 7 | 8 | 2 |
| E7 | 7 | 5 | 9 | 9 | 9 | 6 | 8 | 6 | 5 | 4 | 7 | 7 | 8 | 5 | 5 | 4 | 8 | 5 | 6 | 4 | 4 |
| E8 | 8 | 9 | 10 | 10 | 10 | 8 | 10 | 10 | 8 | 4 | 10 | 10 | 10 | 9 | 7 | 7 | 10 | 9 | 8 | 7 | 9 |
| E9 | 9 | 9 | 5 | 5 | 5 | 9 | 9 | 7 | 9 | 8 | 3 | 3 | 3 | 7 | 9 | 9 | 3 | 7 | 9 | 9 | 8 |
| E10 | 10 | 8 | 6 | 6 | 6 | 10 | 5 | 9 | 10 | 10 | 6 | 6 | 5 | 1 | 10 | 10 | 5 | 1 | 10 | 10 | 7 |

Comparing Different Measures

- The results shown in previous table suggest that a significant number of the measures provide conflicting information about the quality of a pattern
- To understand their differences, we need to examine the properties of these measures

Property under Variable Permutation

| | В | $\overline{\mathbf{B}}$ | | A | $\overline{\mathbf{A}}$ |
|-------------------------|---|-------------------------|-------------------------|---|-------------------------|
| A | p | q | В | р | r |
| $\overline{\mathbf{A}}$ | r | S | $\overline{\mathbf{B}}$ | 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 |
| | 1 | <u> </u> | |

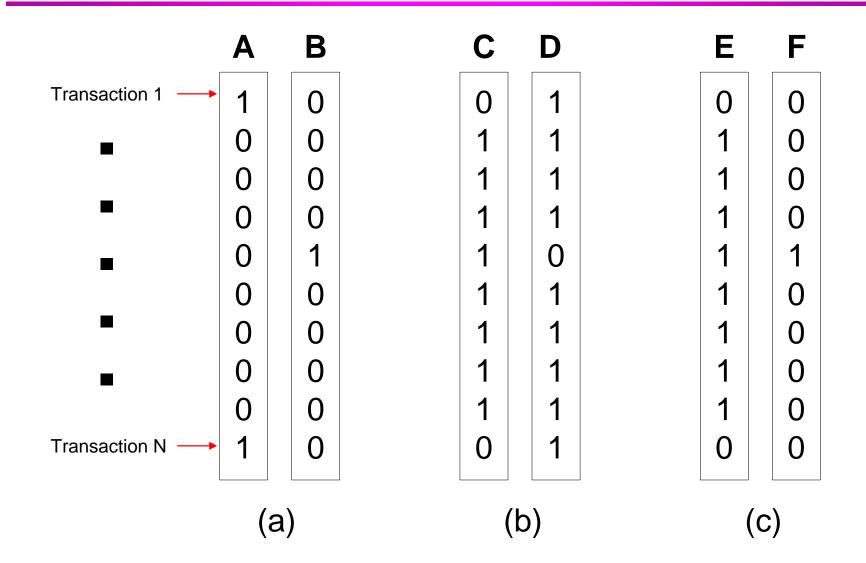
10x

2x

Mosteller:

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

Property under Inversion Operation



Example: ϕ -Coefficient

 φ-coefficient is analogous to correlation coefficient for continuous variables

| | Υ | Y | |
|---|----|----|-----|
| X | 60 | 10 | 70 |
| X | 10 | 20 | 30 |
| | 70 | 30 | 100 |

| | Υ | Y | |
|---|----|----|-----|
| X | 20 | 10 | 30 |
| X | 10 | 60 | 70 |
| | 30 | 70 | 100 |

$$\phi = \frac{0.6 - 0.7 \times 0.7}{\sqrt{0.7 \times 0.3 \times 0.7 \times 0.3}} \qquad \phi = \frac{0.2 - 0.3 \times 0.3}{\sqrt{0.7 \times 0.3 \times 0.7 \times 0.3}}$$
$$= 0.5238 \qquad = 0.5238$$

Coefficient is the same for both tables

Property under Null Addition

- Suppose we are interested in analyzing the relationship between a pair of words, such as "data" and "mining", in a set of documents.
- If a collection of articles about ice fishing is added to the data set
- Should the association between "data" and "mining" be affected?

| | В | $\overline{\mathbf{B}}$ | | | В | $\overline{\mathbf{B}}$ |
|-------------------------|---|-------------------------|---|-------------------------|---|-------------------------|
| A | p | q | | A | р | q |
| $\overline{\mathbf{A}}$ | r | S | / | $\overline{\mathbf{A}}$ | r | s + k |

Invariant measures:

support, cosine, Jaccard, etc

Non-invariant measures:

correlation, Gini, mutual information, odds ratio, etc

Subjective Interestingness Measure

Objective measure:

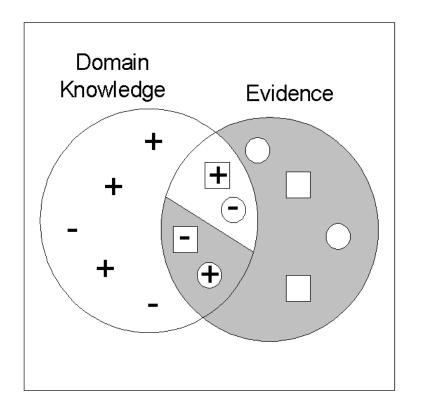
- Rank patterns based on statistics computed from data
- e.g., 21 measures of association (support, confidence, Laplace, Gini, mutual information, Jaccard, etc).

Subjective measure:

- Rank patterns according to user's interpretation
 - A pattern is subjectively interesting if it contradicts the expectation of a user (Silberschatz & Tuzhilin)
 - A pattern is subjectively interesting if it is actionable (Silberschatz & Tuzhilin)

Interestingness via Unexpectedness

Need to model expectation of users (domain knowledge)



- + Pattern expected to be frequent
- Pattern expected to be infrequent
- Pattern found to be frequent
- Pattern found to be infrequent
- **Expected Patterns**
- Unexpected Patterns

 Need to combine expectation of users with evidence from data (i.e., extracted patterns)