

Artificial Intelligence and Machine Learning

Association Analysis

What Is Association?

■ Association rule:

- Finding frequent patterns, associations, correlations, or causal structures among sets of items or objects in transaction databases, relational databases, and other information repositories.
- Frequent patterns and Associations: 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

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. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

■ 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

Definition: Association Rule

■ Association Rule

- An implication expression of the form $X \rightarrow Y$, where X and Y are itemsets
- Meaning: if a basket contains X then it is likely to contain Y .

■ Examples.

- Rule form:
“**Body** \rightarrow **Head** [support, confidence]”.
- $\text{buys}(x, \text{“diapers”})$
 $\rightarrow \text{buys}(x, \text{“beers”}) [0.5\%, 60\%]$
- $\text{major}(x, \text{“CS”}) \wedge \text{takes}(x, \text{“DB”})$
 $\rightarrow \text{grade}(x, \text{“A”}) [1\%, 75\%]$

Rule Evaluation Metrics : Support and Confidence

■ 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

■ Interest (I)

- The interest of an association rule $X \rightarrow Y$ is the absolute value of the amount by which the confidence differs from the probability of Y.

Relationships Among Measures

- Rules with high support and confidence may be useful even if they are not “interesting.”
 - We don’t care if buying bread causes people to buy milk, or whether simply a lot of people buy both bread and milk.
- But high interest suggests a cause that might be worth investigating.

Example

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

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

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

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

$$I = |C - \sigma(\text{Beer})/|T| | \\ = |0.67 - 3/5| \\ = 0.07$$

Association Task

- Given:
 - A large set of transactions
 - Each transaction is a list of items (purchased by a customer in a transaction)
- Find: all rules having
 - support \geq *minsup* threshold
 - confidence \geq *minconf* threshold

The Market-Basket Model

- A large set of items, e.g., things sold in a supermarket.
- A large set of baskets, each of which is a small set of the items, e.g., the things one customer buys in one transaction.

Applications --- (1)

- Real market baskets: chain stores keep terabytes of information about what customers buy together.
 - Tells how typical customers navigate stores, lets them position tempting items.
 - Suggests tie-in “tricks,” e.g., run sale on diapers and raise the price of beer.

Applications --- (2)

- “Baskets” = documents; “items” = words in those documents.
 - Lets us find words that appear together unusually frequently, i.e., linked concepts.
- “Baskets” = sentences, “items” = documents containing those sentences.
 - Items that appear together too often could represent plagiarism.

Applications --- (3)

- “Baskets” = Web pages; “items” = linked pages.
 - Pairs of pages with many common references may be about the same topic.
- “Baskets” = Web pages p ; “items” = pages that link to p .
 - Pages with many of the same links may be mirrors or about the same topic.

Important Point

- “Market Baskets” is an abstraction that models any many-many relationship between two concepts: “items” and “baskets.”
 - Items need not be “contained” in baskets.
- The only difference is that we count co-occurrences of items related to a basket, not vice-versa.

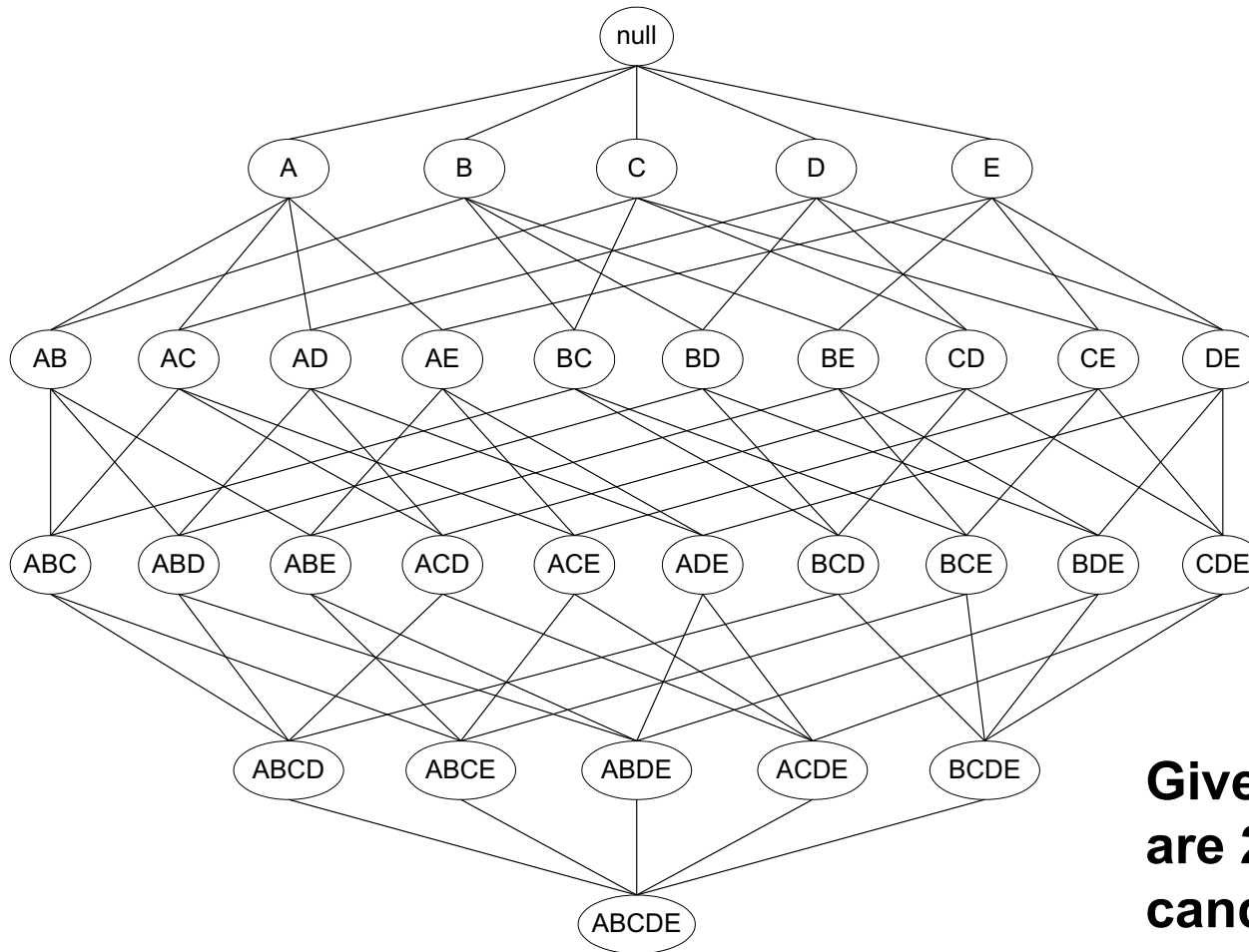
Association Approaches

- 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!

Association Approaches

- Two-step approach:
 - Frequent Itemset Generation
 - Generate all itemsets whose support \geq minsupport
 - 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

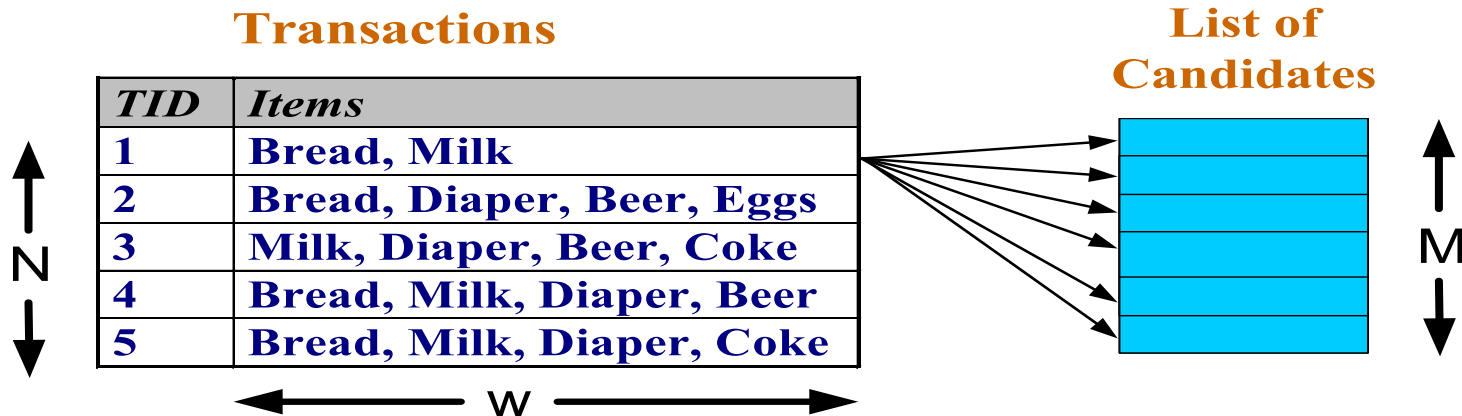
Frequent Itemset Generation



Given d items, there are 2^d possible candidate itemsets

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 or the files



- Match each transaction against every candidate

Computation Model

- Typically, data is kept in a “flat file” rather than a database system.
 - Stored on disk.
 - Stored basket-by-basket.
 - Expand baskets into pairs, triples, etc. as you read baskets.
- The true cost of using disk-resident data is usually the number of disk I/O's.
- In practice, association-rule algorithms read the data in passes --- all baskets read in turn.
- Thus, we measure the cost by the number of passes an algorithm takes.

Main-Memory Bottleneck

- For many frequent-itemset algorithms, main memory is the critical resource.
 - ❑ As we read baskets, we need to count something, e.g., occurrences of pairs.
 - ❑ The number of different things we can count is limited by main memory.
 - ❑ Swapping counts in/out is a disaster.

Naive Algorithm (of Counting Pairs)

- Read file once, counting in main memory the occurrences of each pair.
 - Expand each basket of n items into its $n(n-1)/2$ pairs.
- Fails if $(\text{\#items})^2$ exceeds main memory.
 - Remember: \#items can be 100K (Wal-Mart) or 10B (Web pages).

Details of Main-Memory Counting

- Two approaches:
 1. Count all item pairs, using a triangular matrix.
 2. Keep a table of triples $[i, j, c]$ = the count of the pair of items $\{i, j\}$ is c .
- 1 requires only (say) 4 bytes/pair.
- 2 requires 12 bytes, but only for those pairs with count > 0 .

Frequent Itemset Generation Strategies

- Complexity $\sim O(NMw) \Rightarrow$ Expensive since $M = 2^d$!!!
- 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
 - Data sampling
- 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

■ A-Priori Algorithm

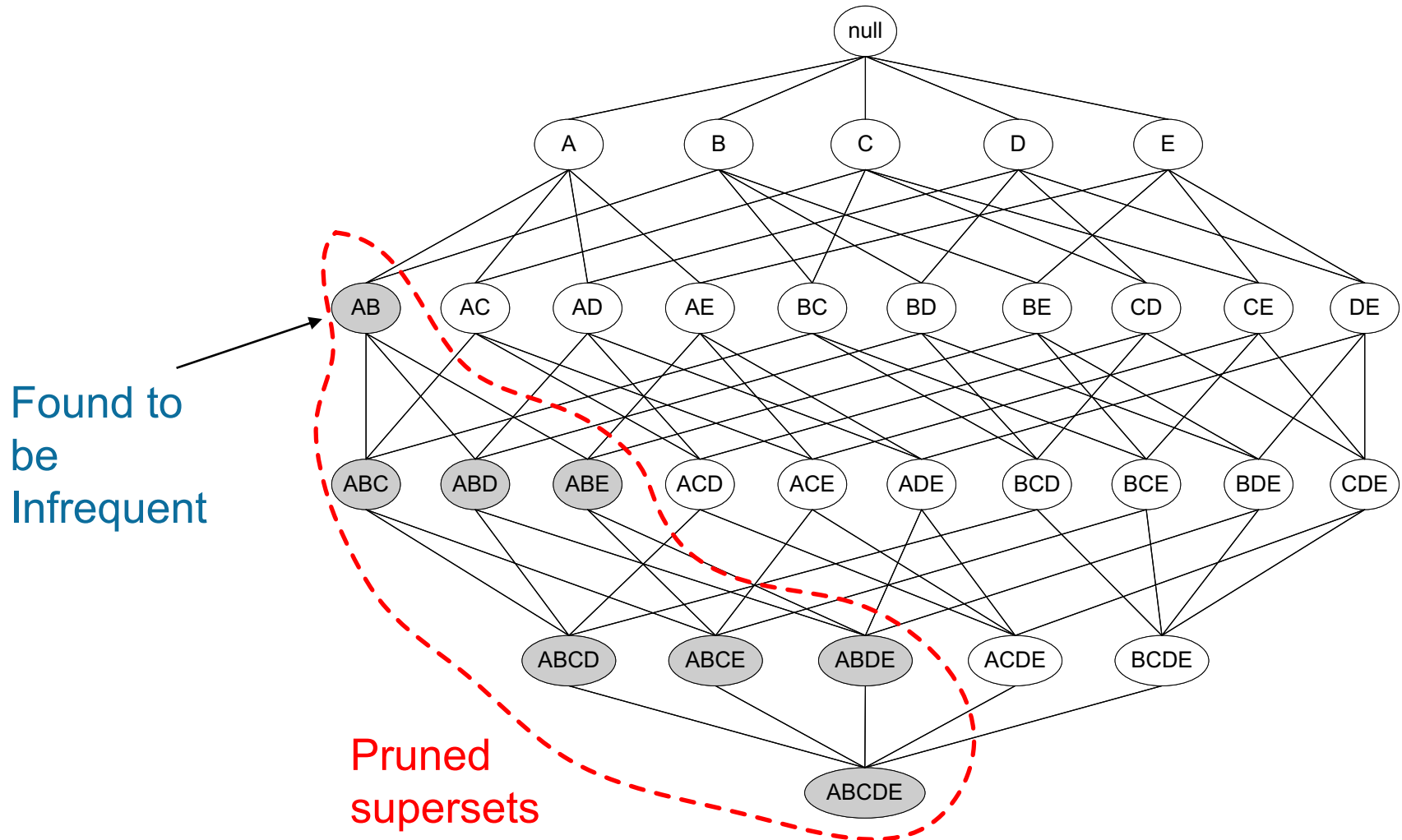
- If an itemset is frequent, then all of its subsets must also be frequent

■ A-Priori 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 A-Priori Principle



A-Priori Algorithm

- A two-pass approach.
- Pass 1: Read baskets and count in main memory the occurrences of each item.
 - Requires only memory proportional to #items.
- Pass 2: Read baskets again and count in main memory only those pairs both of which were found in Pass 1 to be frequent.
 - Requires memory proportional to square of frequent items only.

A-Priori 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 (Candidate generation step)
 - Prune candidate itemsets containing subsets of length k that are infrequent (Candidate pruning step)
 - Count the support of each candidate by scanning the DB
 - Eliminate candidates that are infrequent, leaving only those that are frequent

A-Priori 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 more than min_support

end

return $\cup_k L_k$;

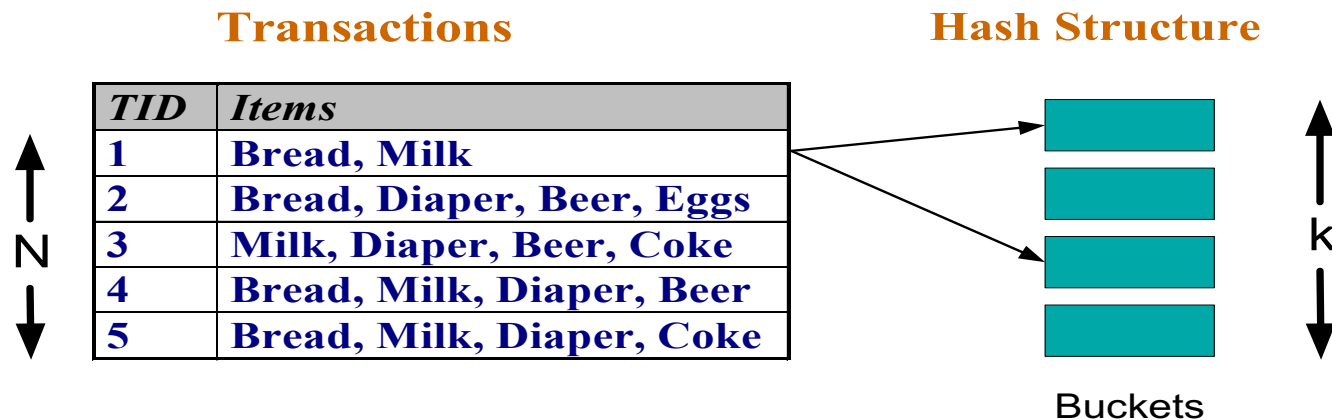
The bottleneck of *A-Priori*

- candidate generation and support counting
 - Use frequent $(k - 1)$ -itemsets to generate candidate frequent k -itemsets
 - Use database scan and pattern matching to collect counts for the candidate itemsets
 - Huge candidate sets:
 - 10^4 frequent 1-itemset will generate 10^7 candidate 2-itemsets
 - To discover a frequent pattern of size 100, e.g., $\{a_1, a_2, \dots, a_{100}\}$, one needs to generate $2^{100} \approx 10^{30}$ candidates.
 - Multiple scans of database:
 - Needs $(n + 1)$ scans, n is the length of the longest pattern

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

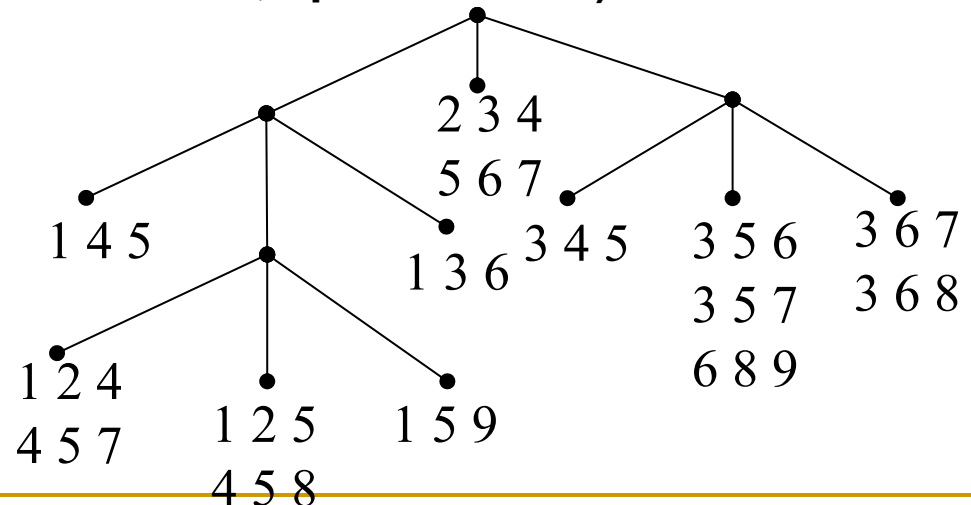
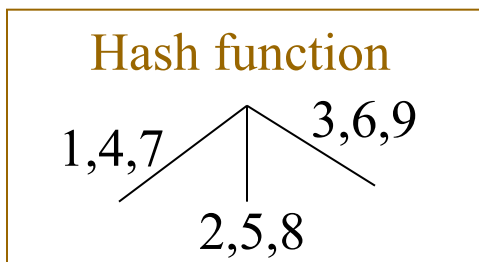
Suppose we have 15 candidate itemsets of length 3:

{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}

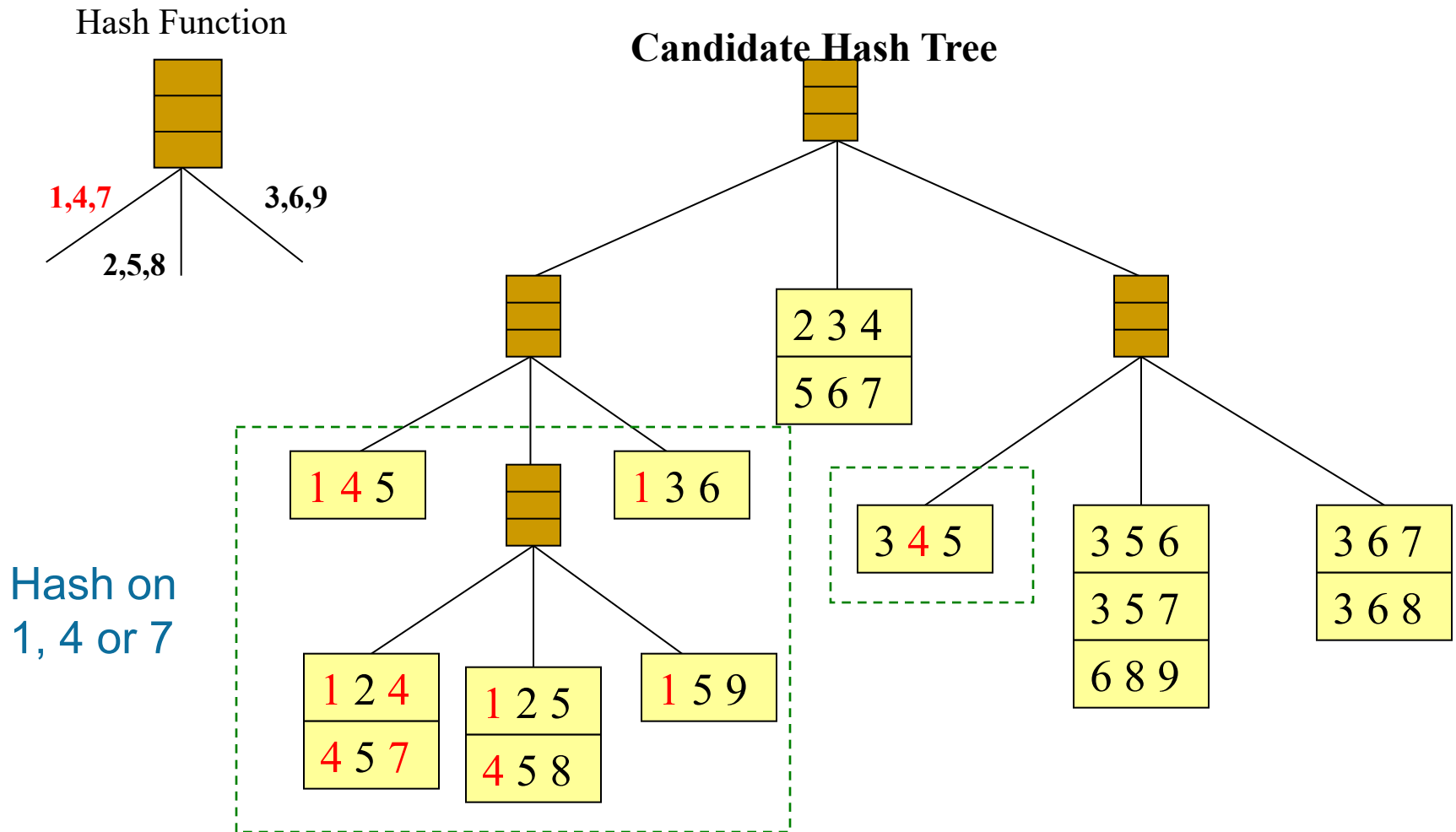
Suppose we have one transaction of length 5 contains: {1, 2, 4, 6, 8}

You need:

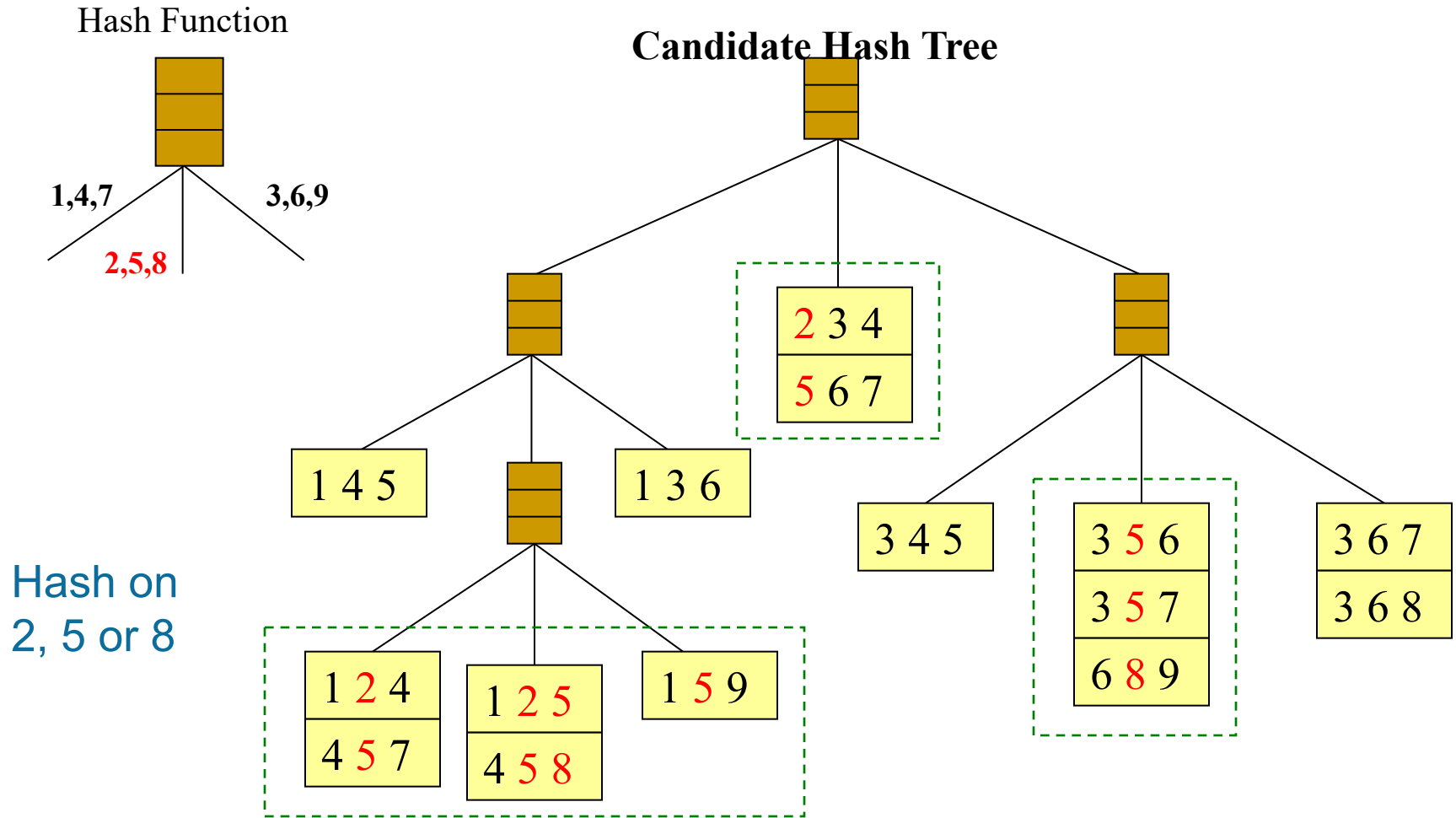
- Hash function
- Max leaf size: max number of itemsets stored in a leaf node (if number of candidate itemsets exceeds max leaf size, split the node)



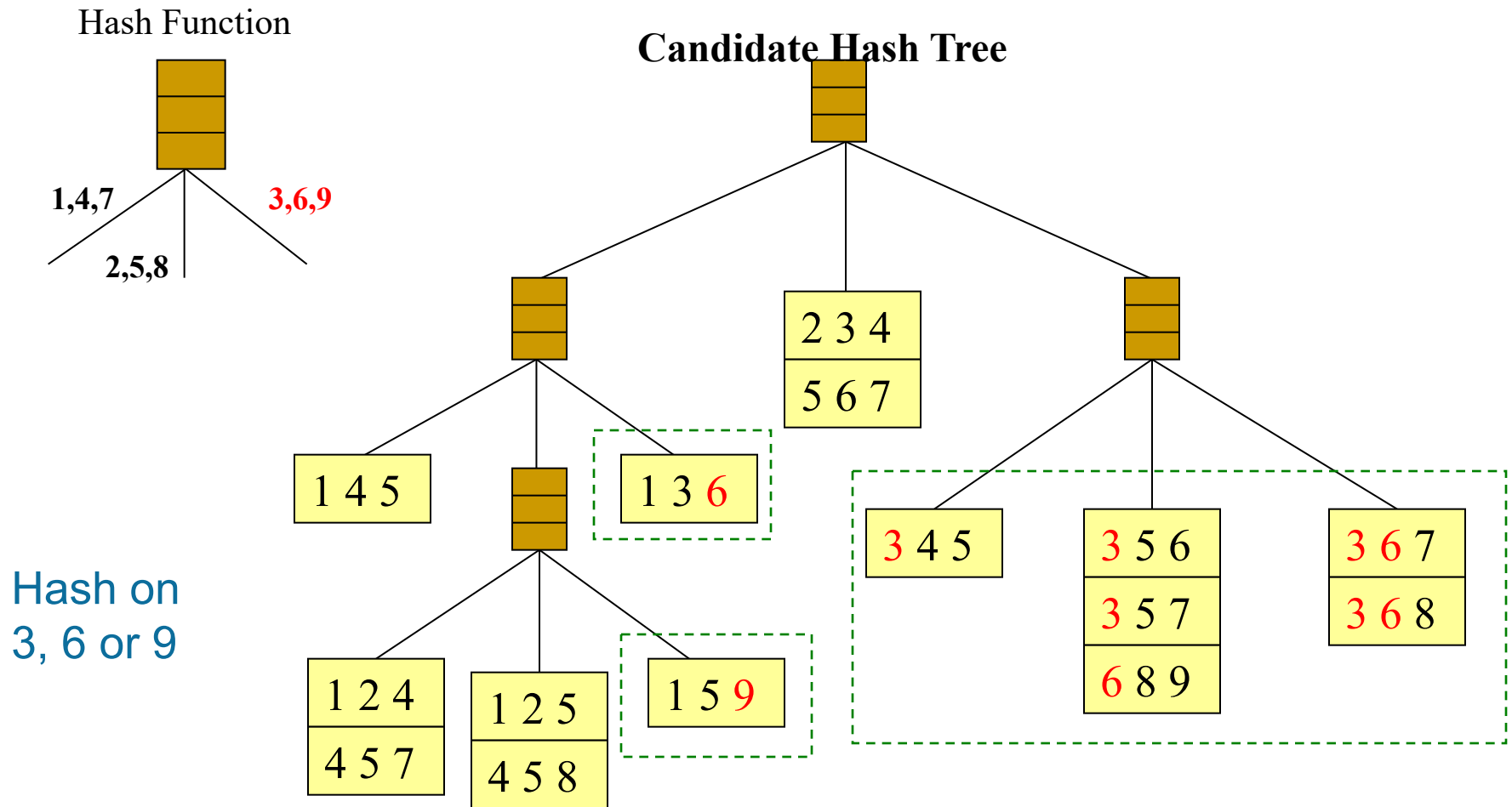
Association Rule Discovery: Hash tree



Association Rule Discovery: Hash tree

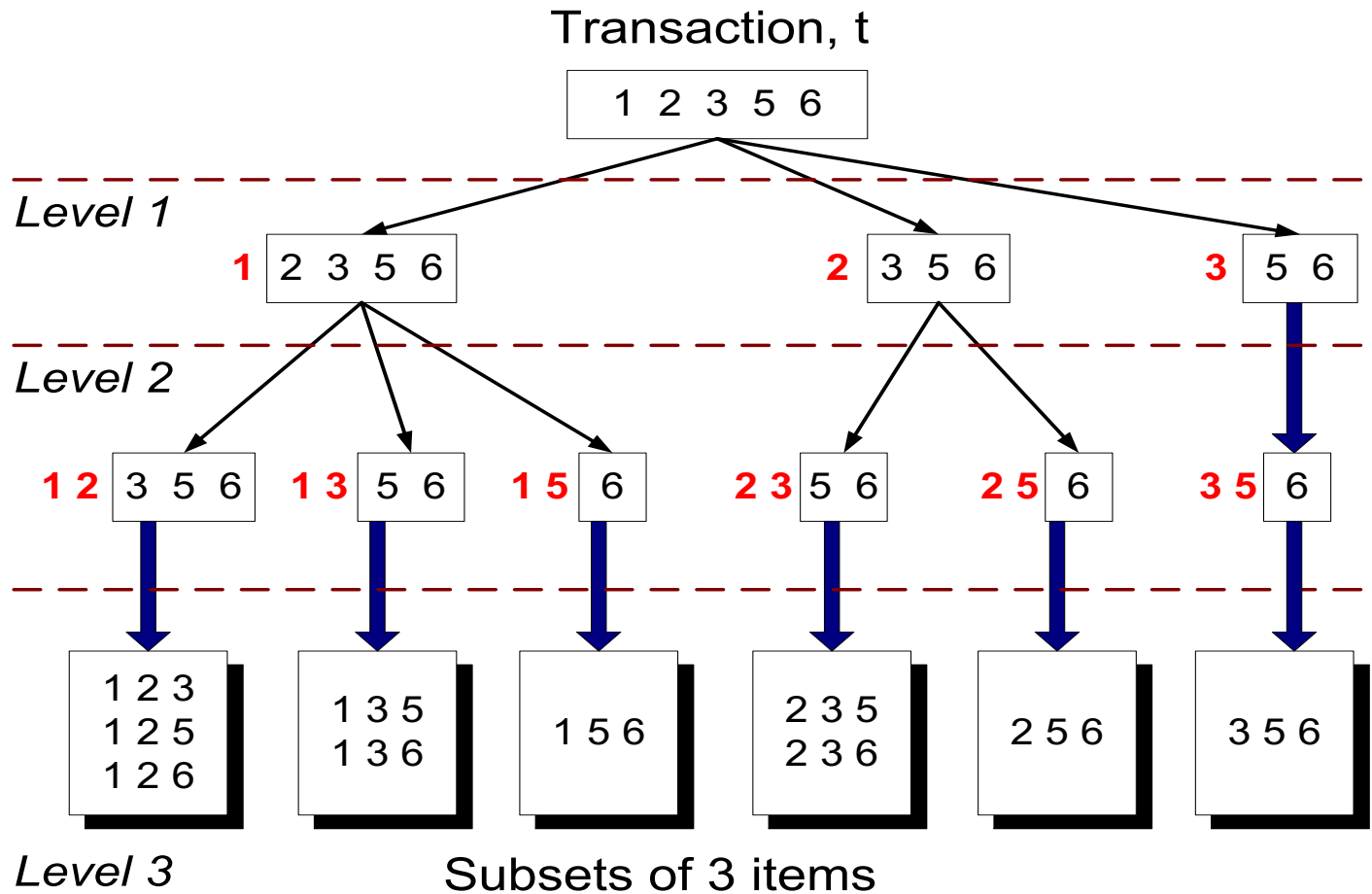


Association Rule Discovery: Hash tree

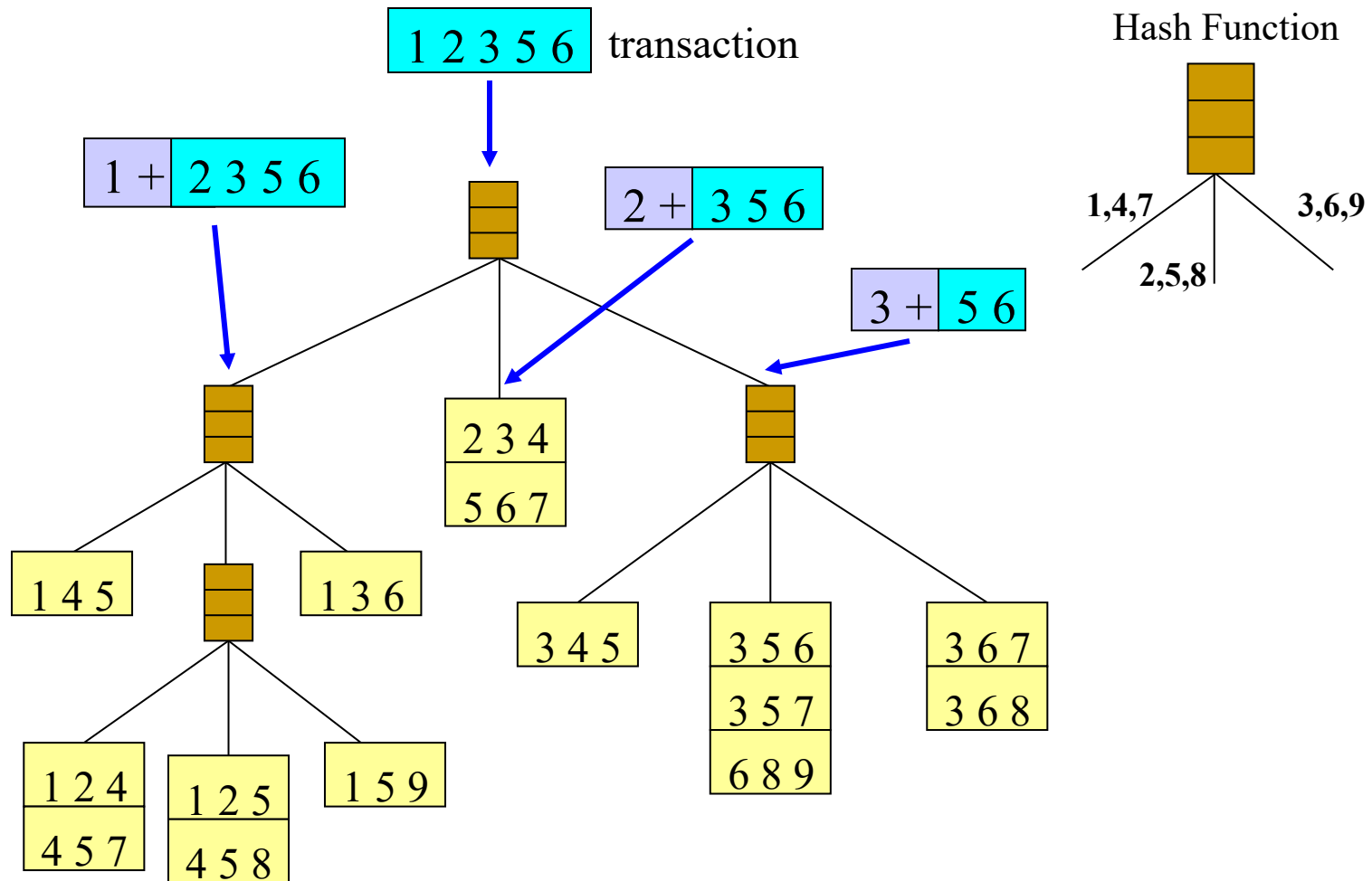


Subset Operation

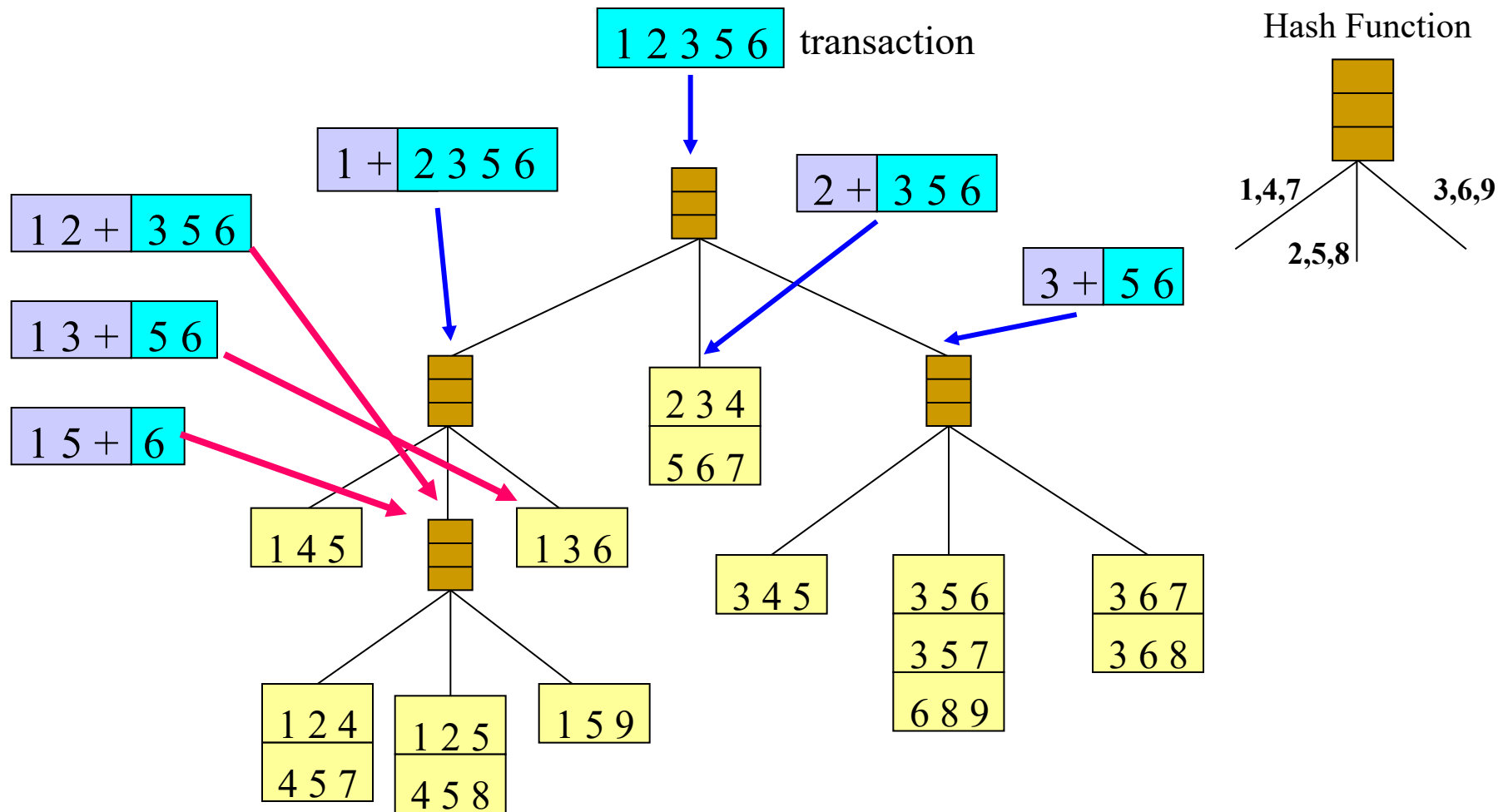
Given a transaction t , what are the possible subsets of size 3?



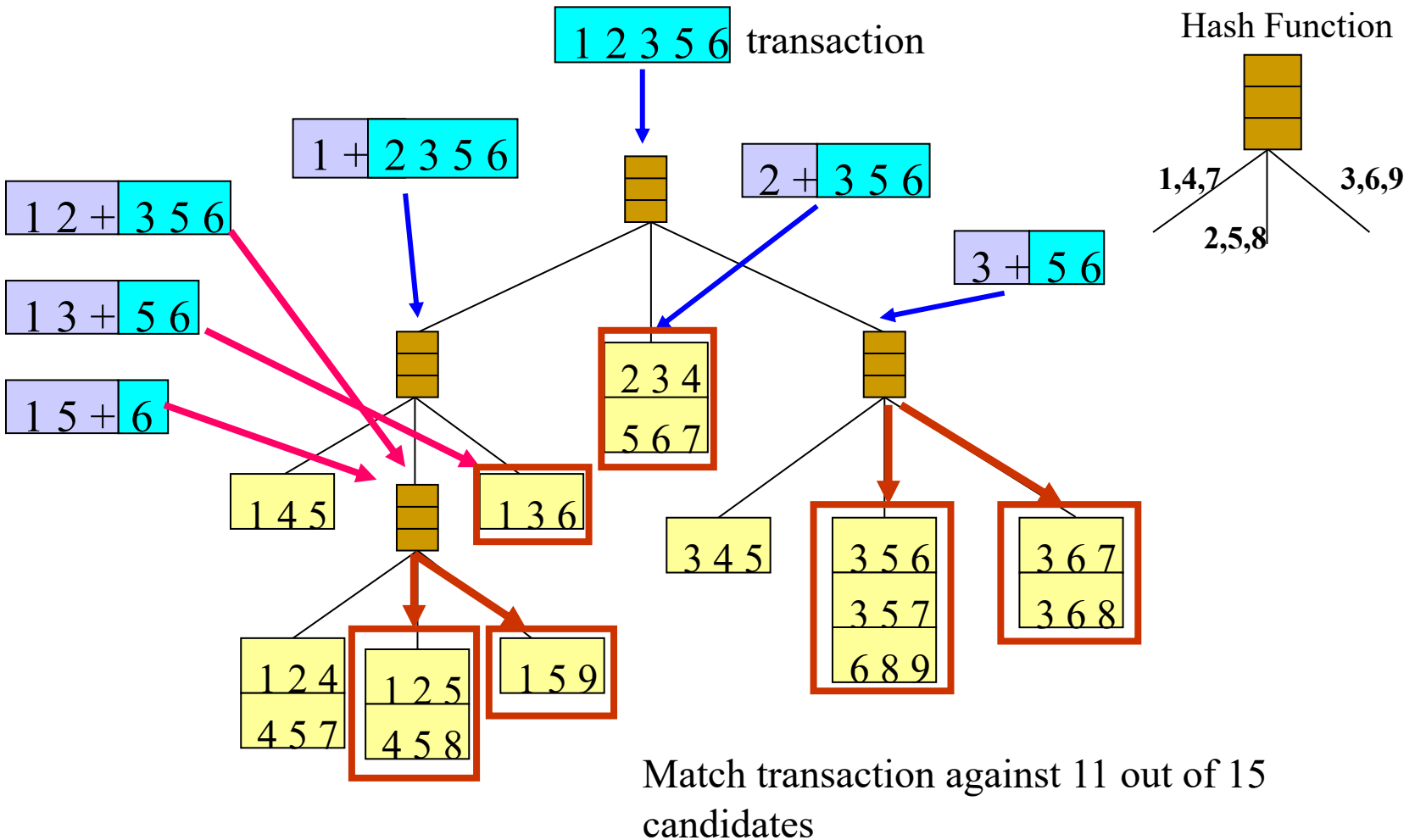
Subset Operation Using Hash Tree



Subset Operation Using Hash Tree



Subset Operation Using Hash Tree



Factors Affecting Complexity (I)

- 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

Factors Affecting Complexity (II)

- Size of database

- Since A-priori makes multiple passes, run time of algorithm may increase with number of transactions

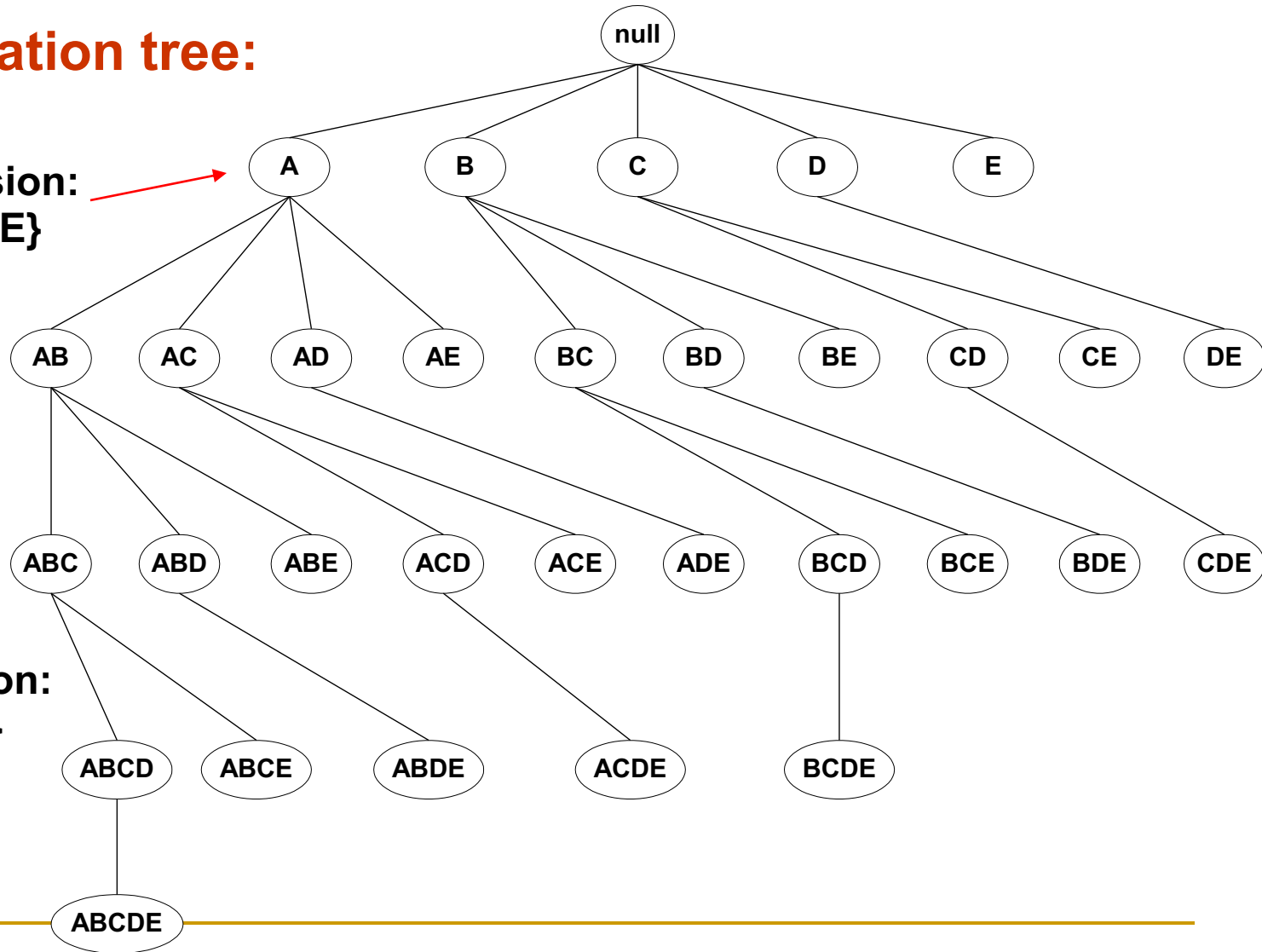
- Average transaction width

- 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)

Alternative: Tree Projection

Set enumeration tree:

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



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 $T \cap E(A)$

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

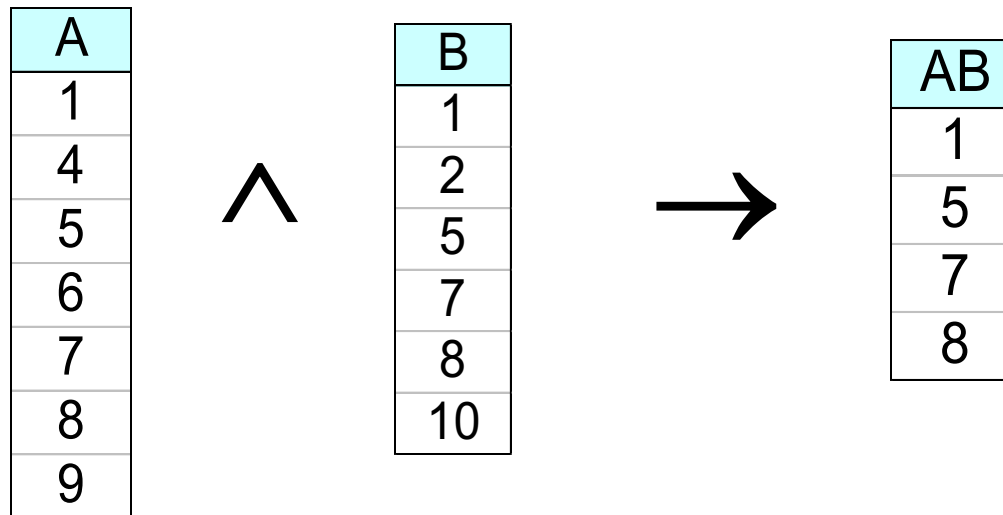
A	B	C	D	E
1	1	2	2	1
4	2	3	4	3
5	5	4	5	6
6	7	8	9	
7	8	9		
8	10			
9				



TID-list

ECLAT

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



- Advantage: very fast support counting
- Disadvantage: intermediate tid-lists may become too large for memory

Alternative Method for Generating Frequent Itemset --- FP-Growth

- Compress a large database into a compact, Frequent-Pattern tree (FP-tree) structure
 - highly condensed, but complete for frequent pattern finding
 - avoid costly database scans
- Develop an efficient, FP-tree-based frequent pattern finding method
 - A divide-and-conquer methodology
 - Avoid candidate generation

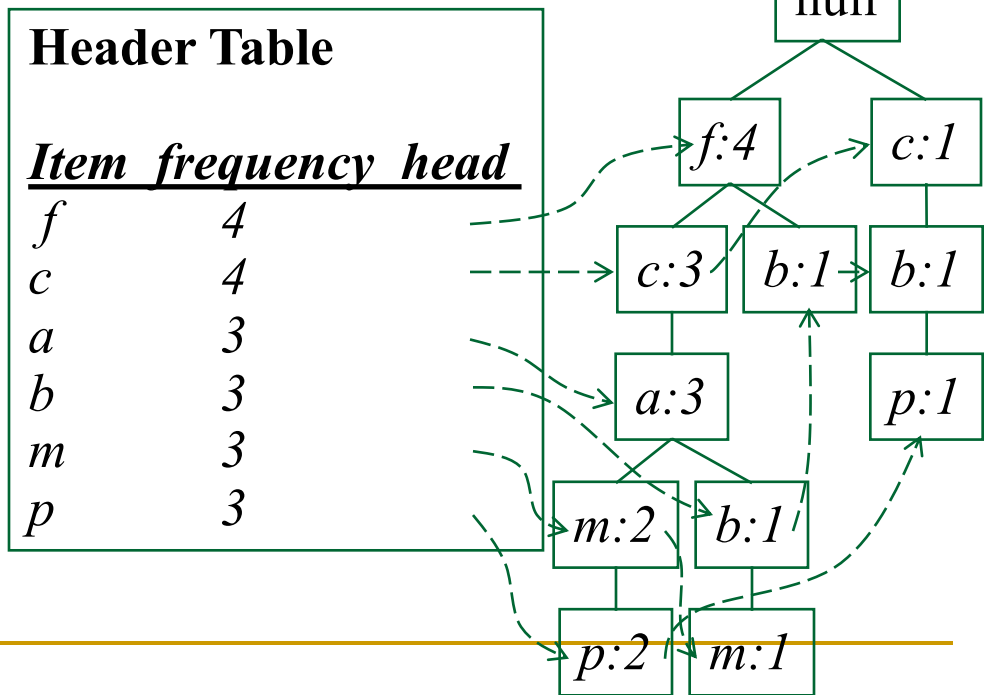
FP-tree construction

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

min_support = 0.5

Steps:

1. Scan DB once, find frequent 1-itemset (single item pattern)
2. Order frequent items in frequency descending order
3. Scan DB again, construct FP-tree



Benefits of the FP-tree Structure

■ Completeness:

- ❑ never breaks a long pattern of any transaction
- ❑ preserves complete information for frequent pattern finding

■ Compactness

- ❑ reduce irrelevant information—infrequent items are gone
- ❑ frequency descending ordering: more frequent items are more likely to be shared
- ❑ never be larger than the original database (if not count node-links and counts)

Finding Frequent Patterns Using FP-tree

- General idea (divide-and-conquer)
 - Recursively grow frequent pattern path using the FP-tree
- Method
 - For each item, construct its **conditional pattern-base**, and then its **conditional FP-tree**
 - Repeat the process on each newly created conditional FP-tree
 - Until the resulting FP-tree is **empty**, or it contains **only one path** (single path will generate all the combinations of its sub-paths, each of which is a frequent pattern)

Major Steps to Construct FP-tree

- 1) Construct conditional pattern base for each node in the FP-tree
- 2) Construct conditional FP-tree from each conditional pattern-base
- 3) Recursively construct conditional FP-trees and grow frequent patterns obtained so far
 - If the conditional FP-tree contains a single path, simply enumerate all the patterns

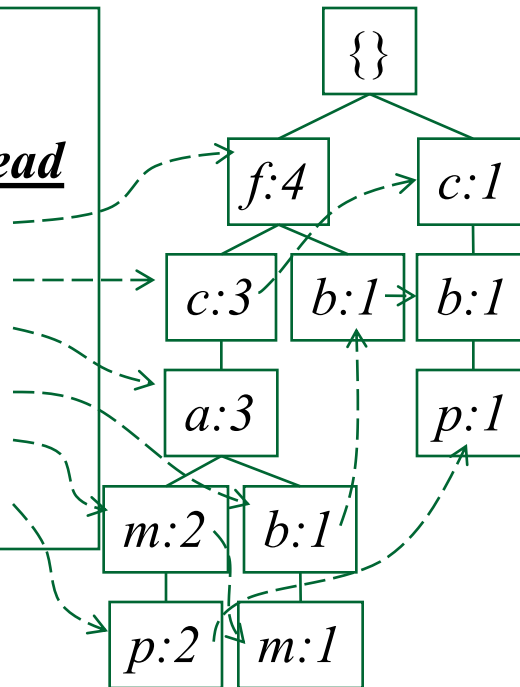
Step 1: From FP-tree to Conditional Pattern Base

- Starting at the frequent header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item
- Accumulate all of transformed prefix paths of that item to form a conditional pattern base

Header Table

Item frequency head

<i>f</i>	4
<i>c</i>	4
<i>a</i>	3
<i>b</i>	3
<i>m</i>	3
<i>p</i>	3



Conditional pattern bases

item cond. pattern base

<i>c</i>	<i>f:3</i>
<i>a</i>	<i>fc:3</i>
<i>b</i>	<i>fca:1, f:1, c:1</i>
<i>m</i>	<i>fca:2, fcab:1</i>
<i>p</i>	<i>fcam:2, cb:1</i>

Properties of FP-tree for Conditional Pattern Base Construction

- Node-link property

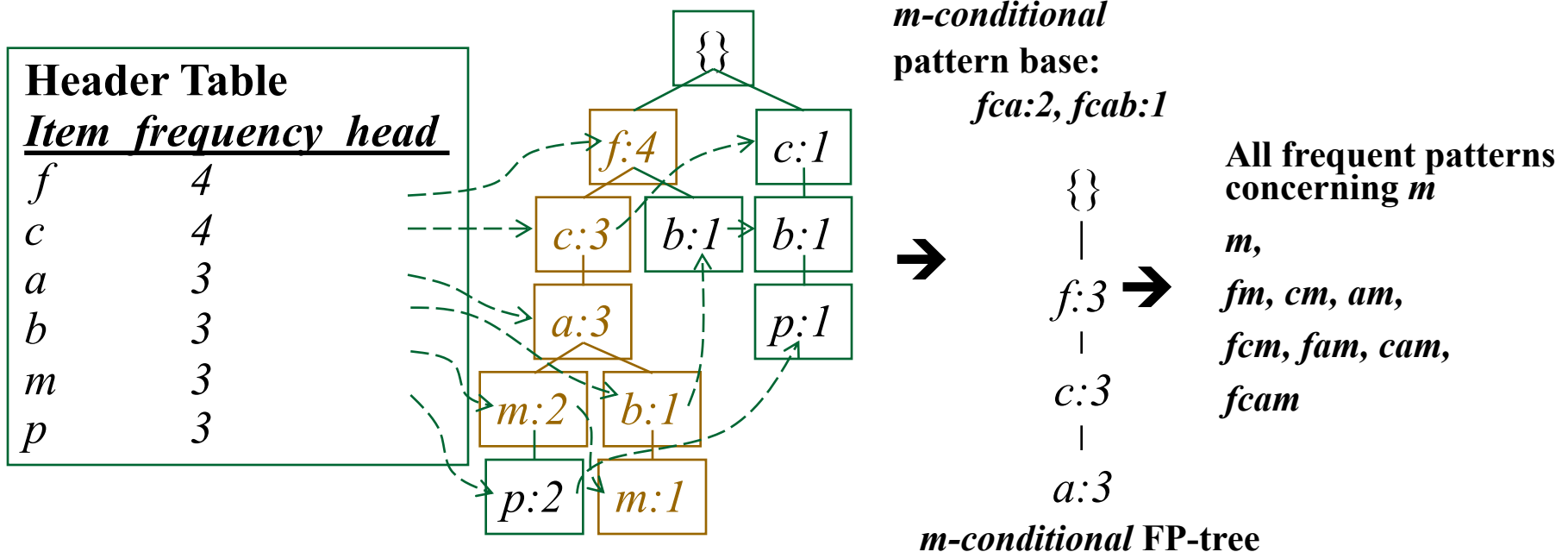
- For any frequent item a_i , all the possible frequent patterns that contain a_i can be obtained by following a_i 's node-links, starting from a_i 's head in the FP-tree header

- Prefix path property

- To calculate the frequent patterns for a node a_i in a path P , only the prefix sub-path of a_i in P need to be accumulated, and its frequency count should carry the same count as node a_i .

Step 2: Construct Conditional FP-tree

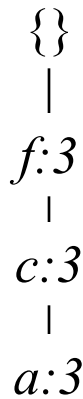
- For each pattern-base
 - Accumulate the count for each item in the base
 - Construct the FP-tree for the frequent items of the pattern base



Finding Frequent Patterns by Creating Conditional Pattern-Bases

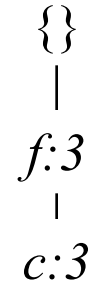
Item	Conditional pattern-base	Conditional FP-tree
p	{(fcam:2), (cb:1)}	{(c:3)} p
m	{(fca:2), (fcab:1)}	{(f:3, c:3, a:3)} m
b	{(fca:1), (f:1), (c:1)}	Empty
a	{(fc:3)}	{(f:3, c:3)} a
c	{(f:3)}	{(f:3)} c
f	Empty	Empty

Step 3: Recursively Construct the conditional FP-tree



m-conditional FP-tree

Cond. pattern base of "am": (fc:3)



am-conditional FP-tree

Cond. pattern base of "cm": (f:3)



cm-conditional FP-tree

Cond. pattern base of "cam": (f:3)



cam-conditional FP-tree

Single FP-tree Path Generation

- Suppose an FP-tree T has a single path P
- The complete set of frequent pattern of T can be generated by enumeration of all the combinations of the sub-paths of P

$\{\}$
|
 $f:3$
|
 $c:3$
|
 $a:3$



**All frequent patterns
concerning m**

$m,$
 $fm, cm, am,$
 $fcm, fam, cam,$
 $fcam$

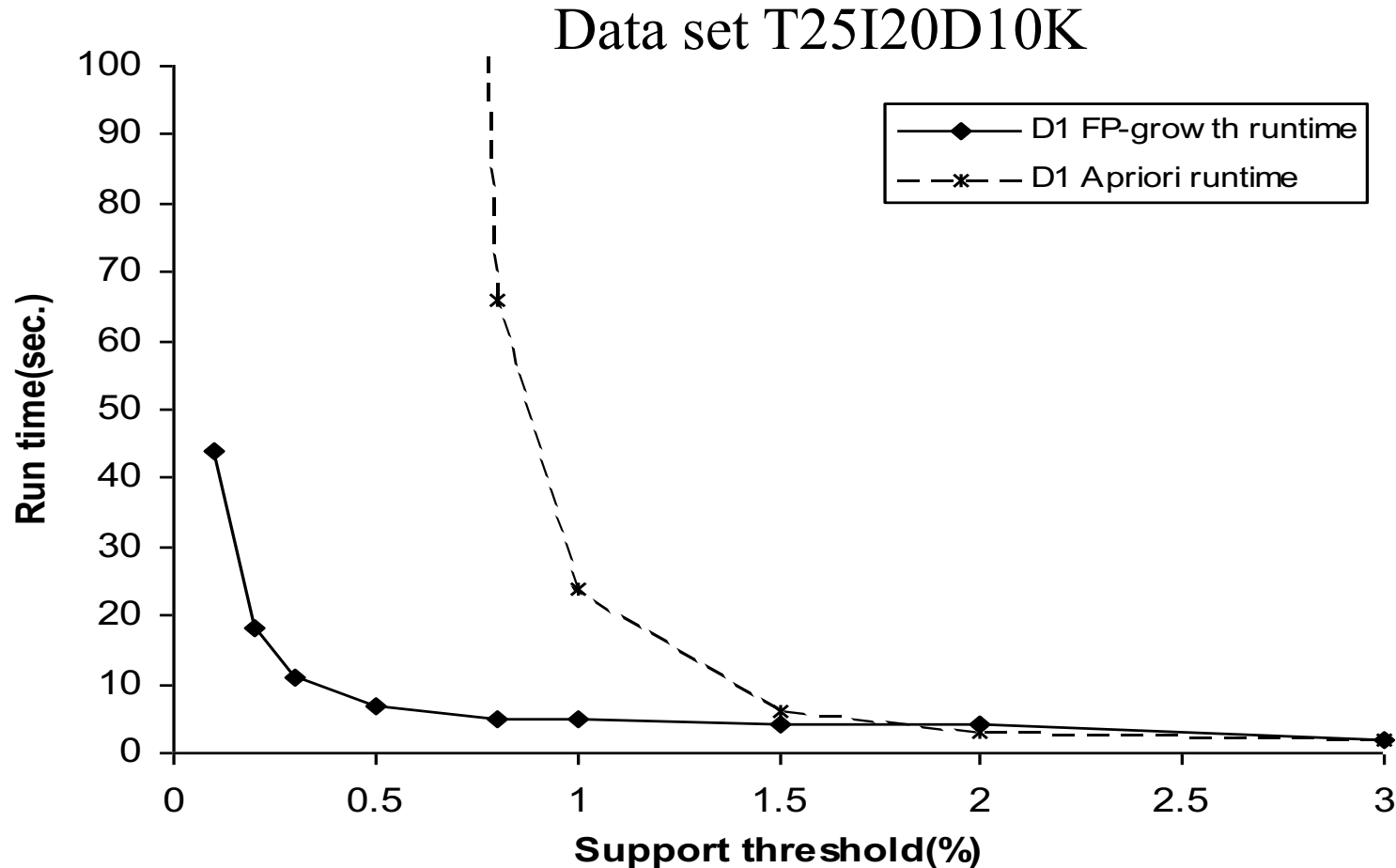
Principles of Frequent Pattern Growth

- Pattern growth property
 - Let α be a frequent itemset in DB, B be α 's conditional pattern base, and β be an itemset in B . Then $\alpha \cup \beta$ is a frequent itemset in DB iff β is frequent in B .
- “*abcdef*” is a frequent pattern, if and only if
 - “*abcde*” is a frequent pattern, and
 - “*f*” is frequent in the set of transactions containing “*abcde*”

Why Is FP Growth Fast?

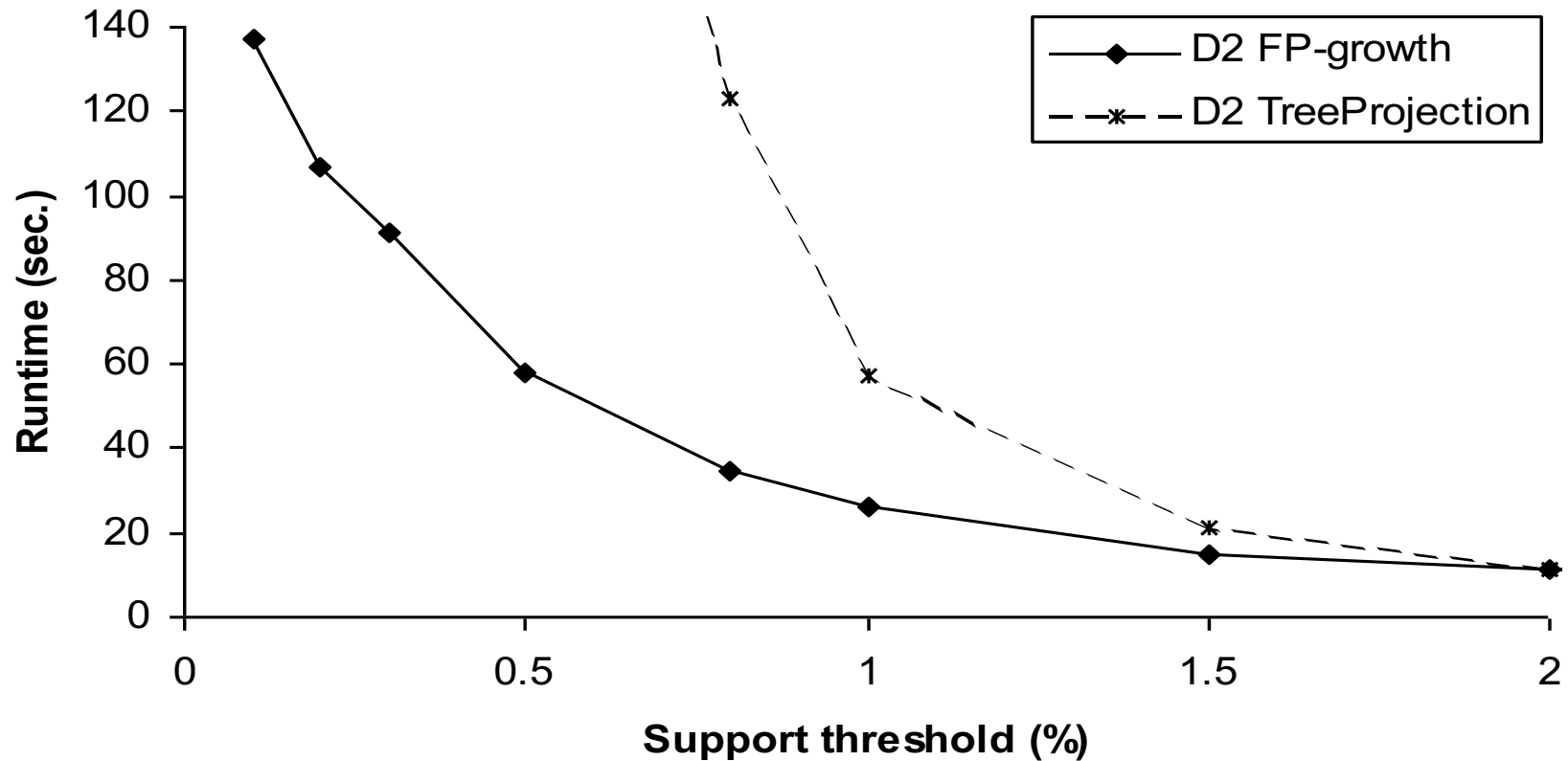
- Our performance study shows
 - FP-growth is an order of magnitude faster than Apriori, and is also faster than tree-projection
- Reasoning
 - No candidate generation, no candidate test
 - Use compact data structure
 - Eliminate repeated database scan
 - Basic operation is counting and FP-tree building

FP-growth vs. Apriori: Scalability With the Support Threshold



FP-growth vs. Tree-Projection: Scalability with Support Threshold

Data set T25I20D100K



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:

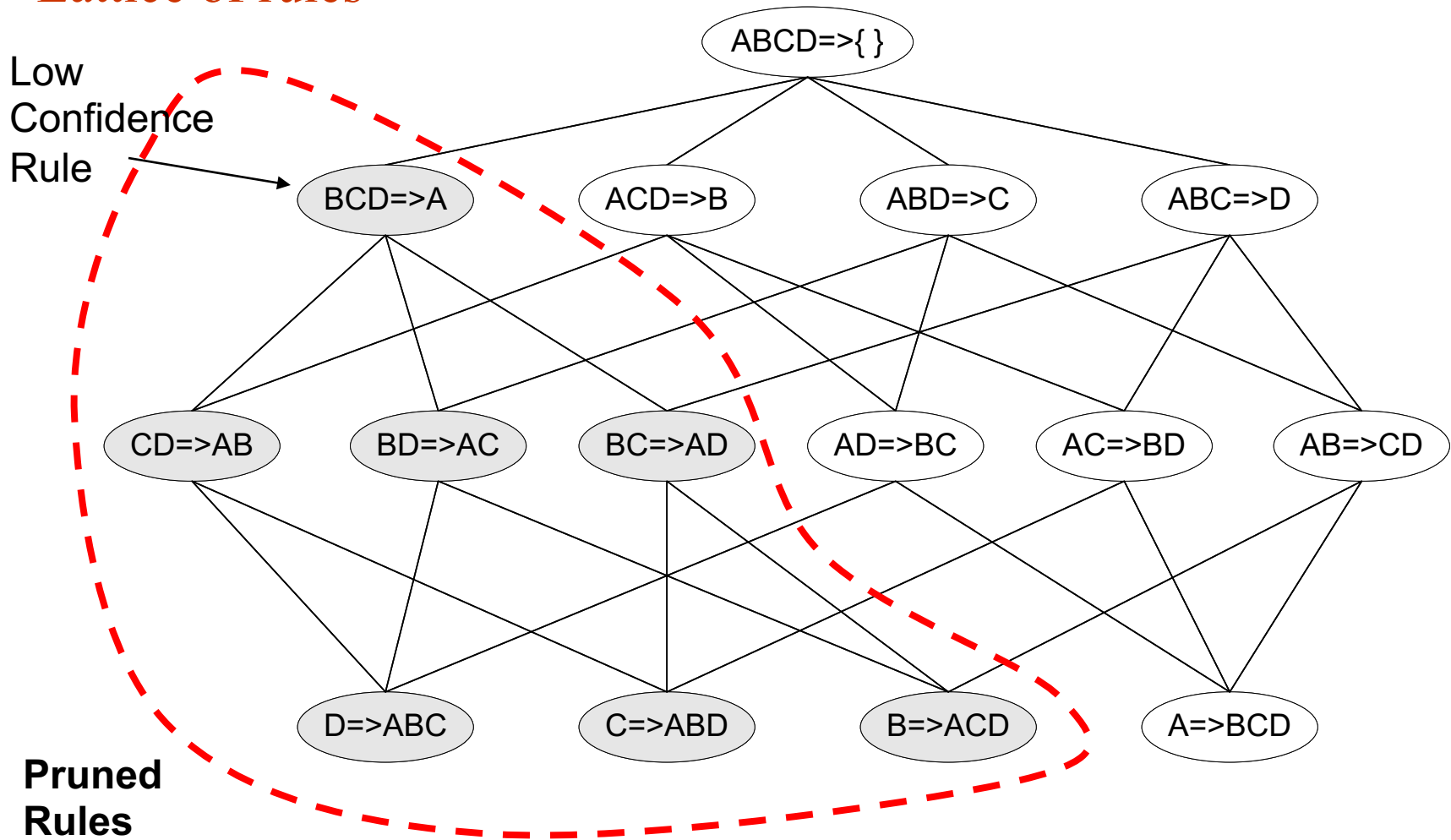
$ABC \rightarrow D,$	$ABD \rightarrow C,$	$ACD \rightarrow B,$	$BCD \rightarrow A,$
$A \rightarrow BCD,$	$B \rightarrow ACD,$	$C \rightarrow ABD,$	$D \rightarrow ABC$
$AB \rightarrow CD,$	$AC \rightarrow BD,$	$AD \rightarrow BC,$	$BC \rightarrow AD,$
$BD \rightarrow AC,$	$CD \rightarrow AB,$		
- If $|L| = k$, then there are $2^k - 2$ candidate association rules (ignoring $L \rightarrow \emptyset$ and $\emptyset \rightarrow L$)

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

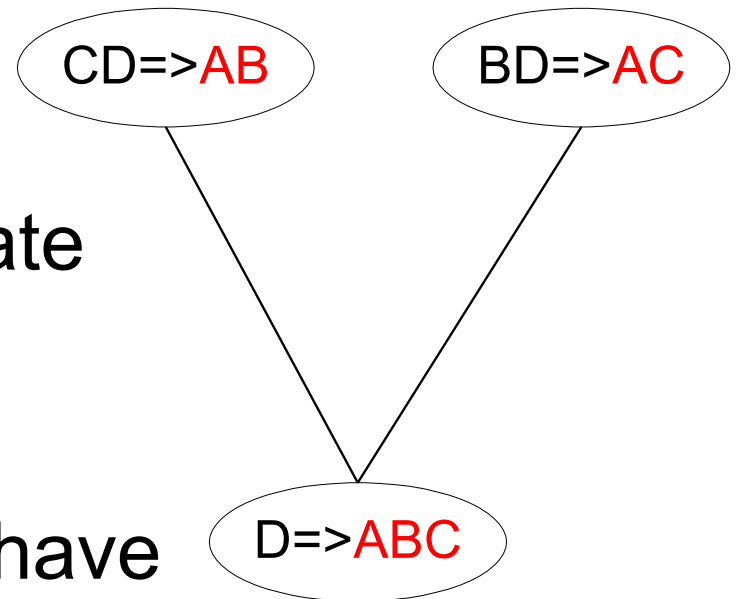
Rule Generation for A-Priori Algorithm

Lattice of rules



Rule Generation for A-Priori 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

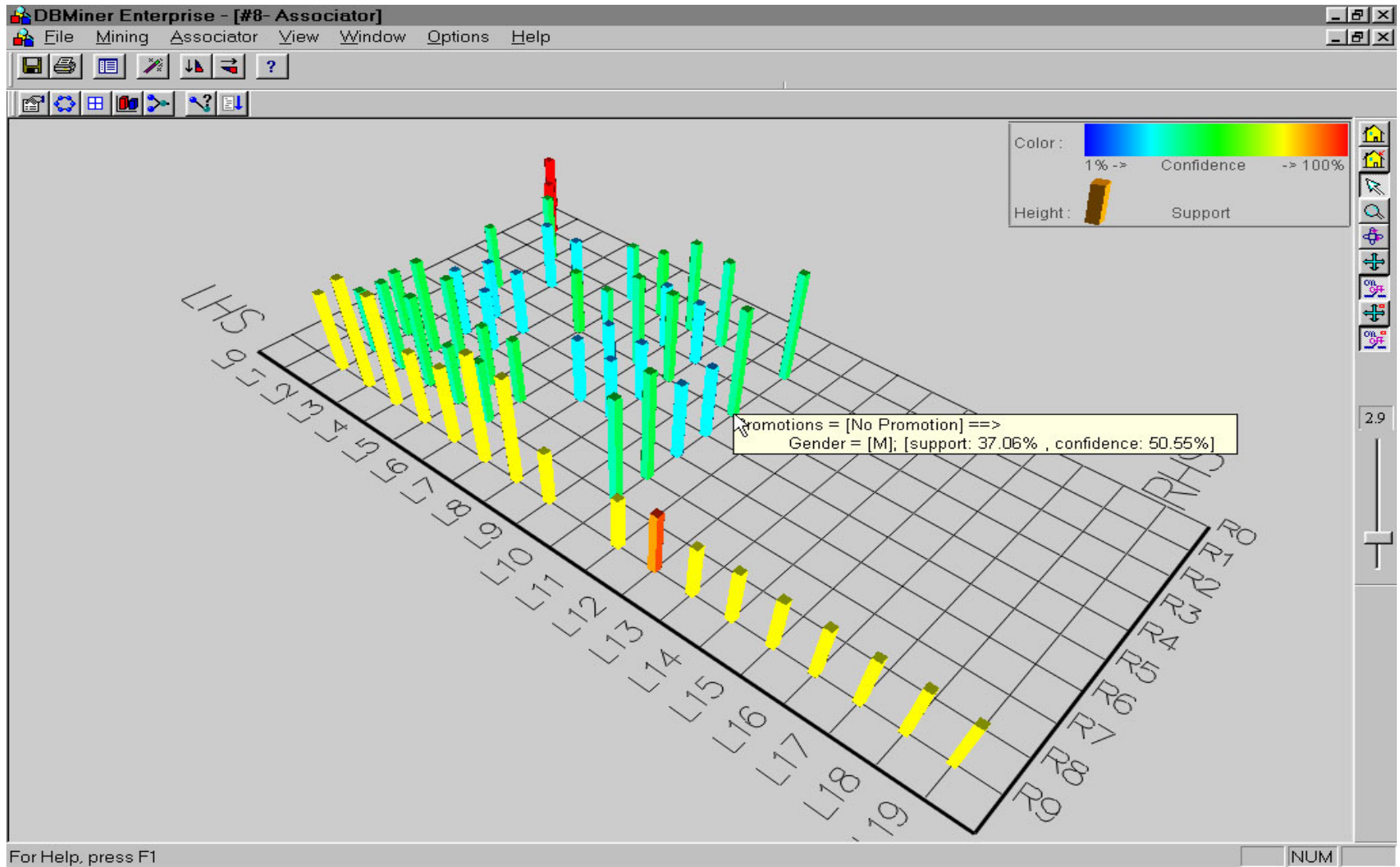


Presentation of Association Rules (Table Form)

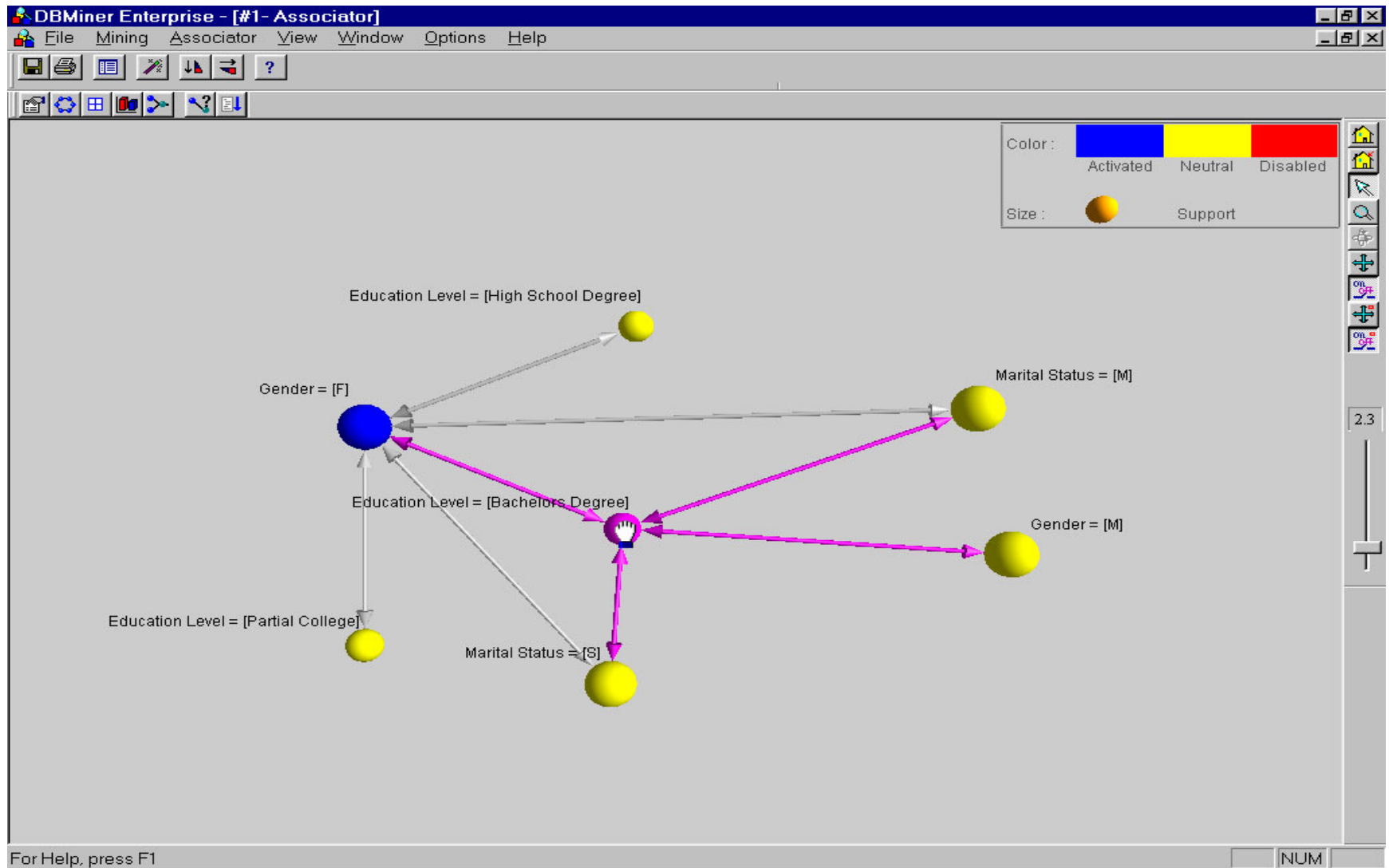
	Body	Implies	Head	Supp (%)	Conf (%)	F	G	H	I
1	cost(x) = '0.00~1000.00'	==>	revenue(x) = '0.00~500.00'	28.45	40.4				
2	cost(x) = '0.00~1000.00'	==>	revenue(x) = '500.00~1000.00'	20.46	29.05				
3	cost(x) = '0.00~1000.00'	==>	order_qty(x) = '0.00~100.00'	59.17	84.04				
4	cost(x) = '0.00~1000.00'	==>	revenue(x) = '1000.00~1500.00'	10.45	14.84				
5	cost(x) = '0.00~1000.00'	==>	region(x) = 'United States'	22.56	32.04				
6	cost(x) = '1000.00~2000.00'	==>	order_qty(x) = '0.00~100.00'	12.91	69.34				
7	order_qty(x) = '0.00~100.00'	==>	revenue(x) = '0.00~500.00'	28.45	34.54				
8	order_qty(x) = '0.00~100.00'	==>	cost(x) = '1000.00~2000.00'	12.91	15.67				
9	order_qty(x) = '0.00~100.00'	==>	region(x) = 'United States'	25.9	31.45				
10	order_qty(x) = '0.00~100.00'	==>	cost(x) = '0.00~1000.00'	59.17	71.86				
11	order_qty(x) = '0.00~100.00'	==>	product_line(x) = 'Tents'	13.52	16.42				
12	order_qty(x) = '0.00~100.00'	==>	revenue(x) = '500.00~1000.00'	19.67	23.88				
13	product_line(x) = 'Tents'	==>	order_qty(x) = '0.00~100.00'	13.52	98.72				
14	region(x) = 'United States'	==>	order_qty(x) = '0.00~100.00'	25.9	81.94				
15	region(x) = 'United States'	==>	cost(x) = '0.00~1000.00'	22.56	71.39				
16	revenue(x) = '0.00~500.00'	==>	cost(x) = '0.00~1000.00'	28.45	100				
17	revenue(x) = '0.00~500.00'	==>	order_qty(x) = '0.00~100.00'	28.45	100				
18	revenue(x) = '1000.00~1500.00'	==>	cost(x) = '0.00~1000.00'	10.45	96.75				
19	revenue(x) = '500.00~1000.00'	==>	cost(x) = '0.00~1000.00'	20.46	100				
20	revenue(x) = '500.00~1000.00'	==>	order_qty(x) = '0.00~100.00'	19.67	96.14				
21									
22									
23	cost(x) = '0.00~1000.00'	==>	revenue(x) = '0.00~500.00' AND order_qty(x) = '0.00~100.00'	28.45	40.4				
24	cost(x) = '0.00~1000.00'	==>	revenue(x) = '0.00~500.00' AND order_qty(x) = '0.00~100.00'	28.45	40.4				
25	cost(x) = '0.00~1000.00'	==>	revenue(x) = '500.00~1000.00' AND order_qty(x) = '0.00~100.00'	19.67	27.93				
26	cost(x) = '0.00~1000.00'	==>	revenue(x) = '500.00~1000.00' AND order_qty(x) = '0.00~100.00'	19.67	27.93				
27	cost(x) = '0.00~1000.00' AND order_qty(x) = '0.00~100.00'	==>	revenue(x) = '500.00~1000.00'	19.67	33.23				

Sheet1

Visualization of Association Rule Using Plane Graph



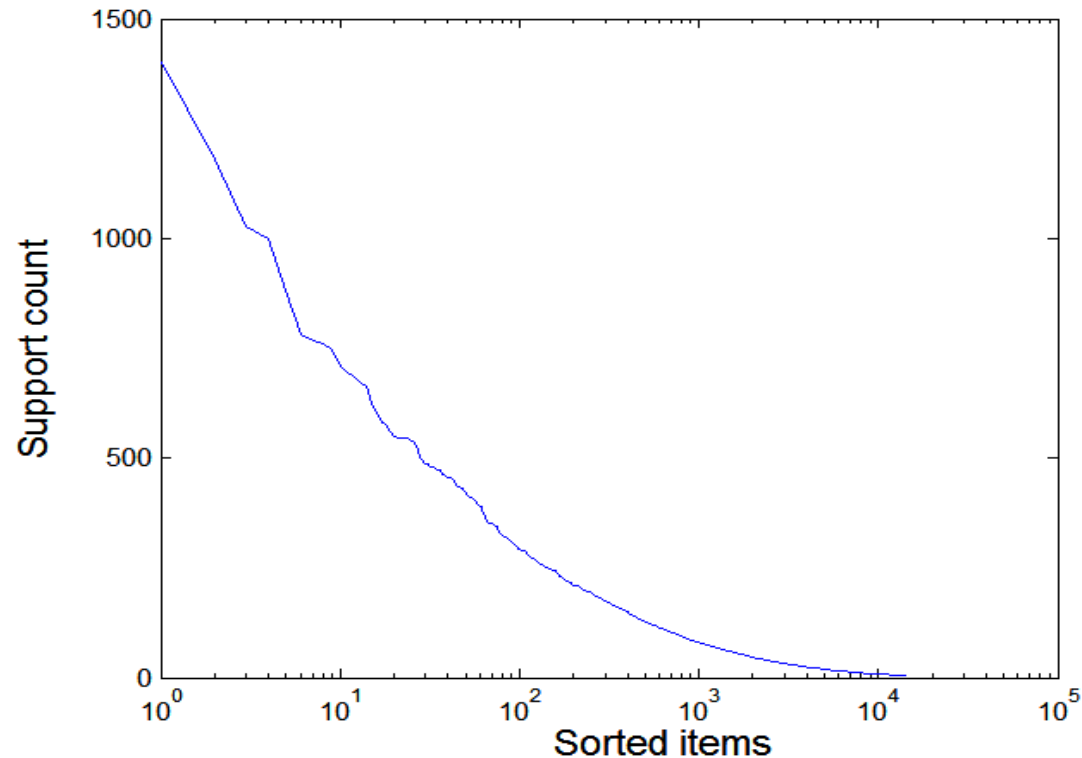
Visualization of Association Rule Using Rule Graph



Effect of Support Distribution

- Many real data sets have skewed support distribution

**Support
distribution of
a retail data set**



Effect of Support Distribution

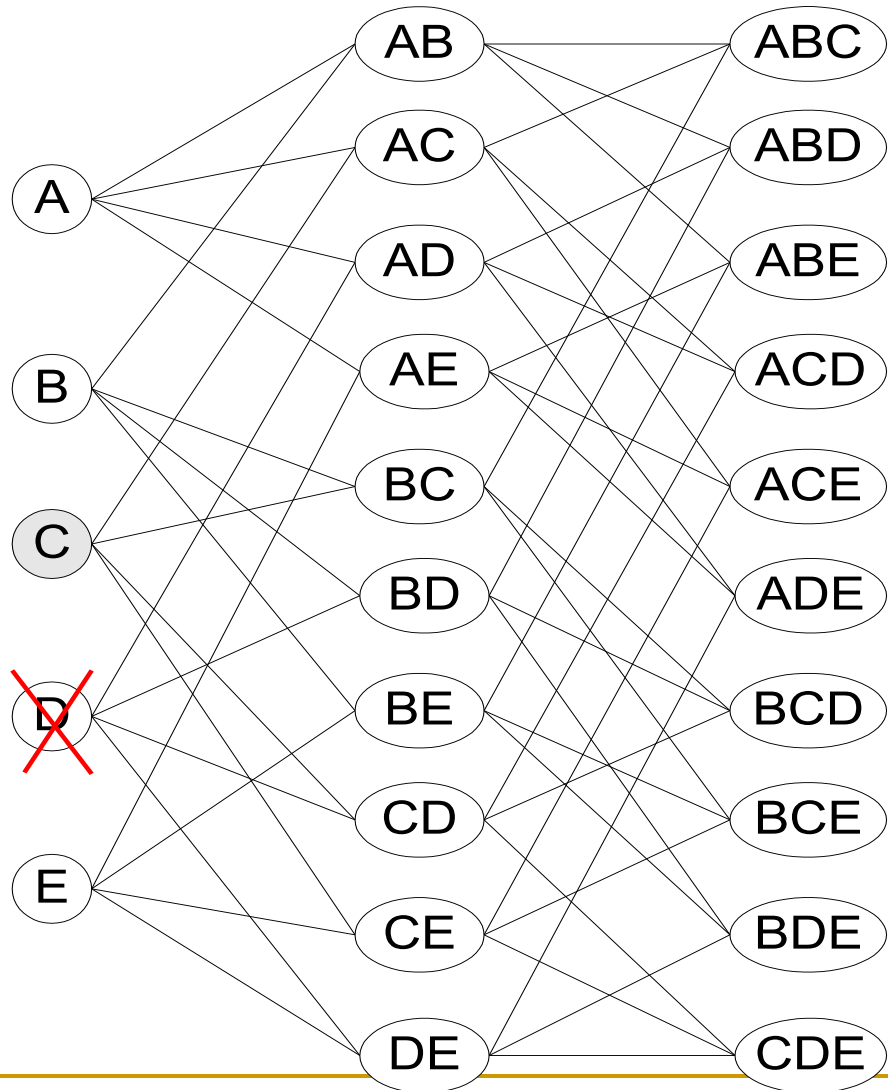
- How to set the appropriate *minsup* threshold?
 - If *minsup* is set too high, we could miss itemsets involving interesting rare items (e.g., expensive products)
 - If *minsup* is set too low, it is computationally expensive and the number of itemsets is very large
- Using a single minimum support threshold may not be effective

Multiple Minimum Support

- How to apply multiple minimum supports?
 - $MS(i)$: minimum support for item i
 - e.g.: $MS(\text{Milk})=5\%$, $MS(\text{Coke}) = 3\%$,
 $MS(\text{Broccoli})=0.1\%$, $MS(\text{Salmon})=0.5\%$
 - $MS(\{\text{Milk}, \text{Broccoli}\}) = \min (MS(\text{Milk}), MS(\text{Broccoli}))$
 $= 0.1\%$
 - 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

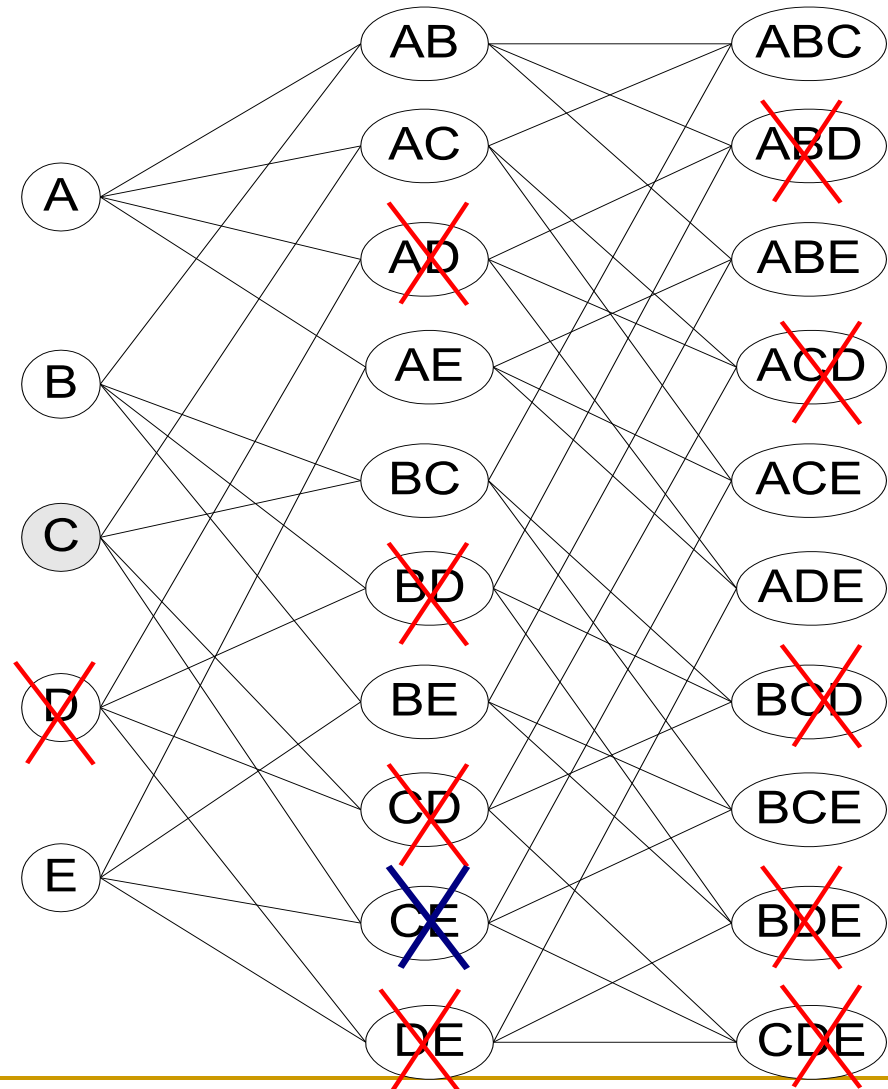
Multiple Minimum Support

Item	MS(I)	Sup(I)
A	0.10%	0.25%
B	0.20%	0.26%
C	0.30%	0.29%
D	0.50%	0.05%
E	3%	4.20%



Multiple Minimum Support

Item	MS(I)	Sup(I)
A	0.10%	0.25%
B	0.20%	0.26%
C	0.30%	0.29%
D	0.50%	0.05%
E	3%	4.20%



Multiple Minimum Support

- 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 A-Priori 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

■ Modifications to A-Priori:

□ In traditional A-Priori,

- 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

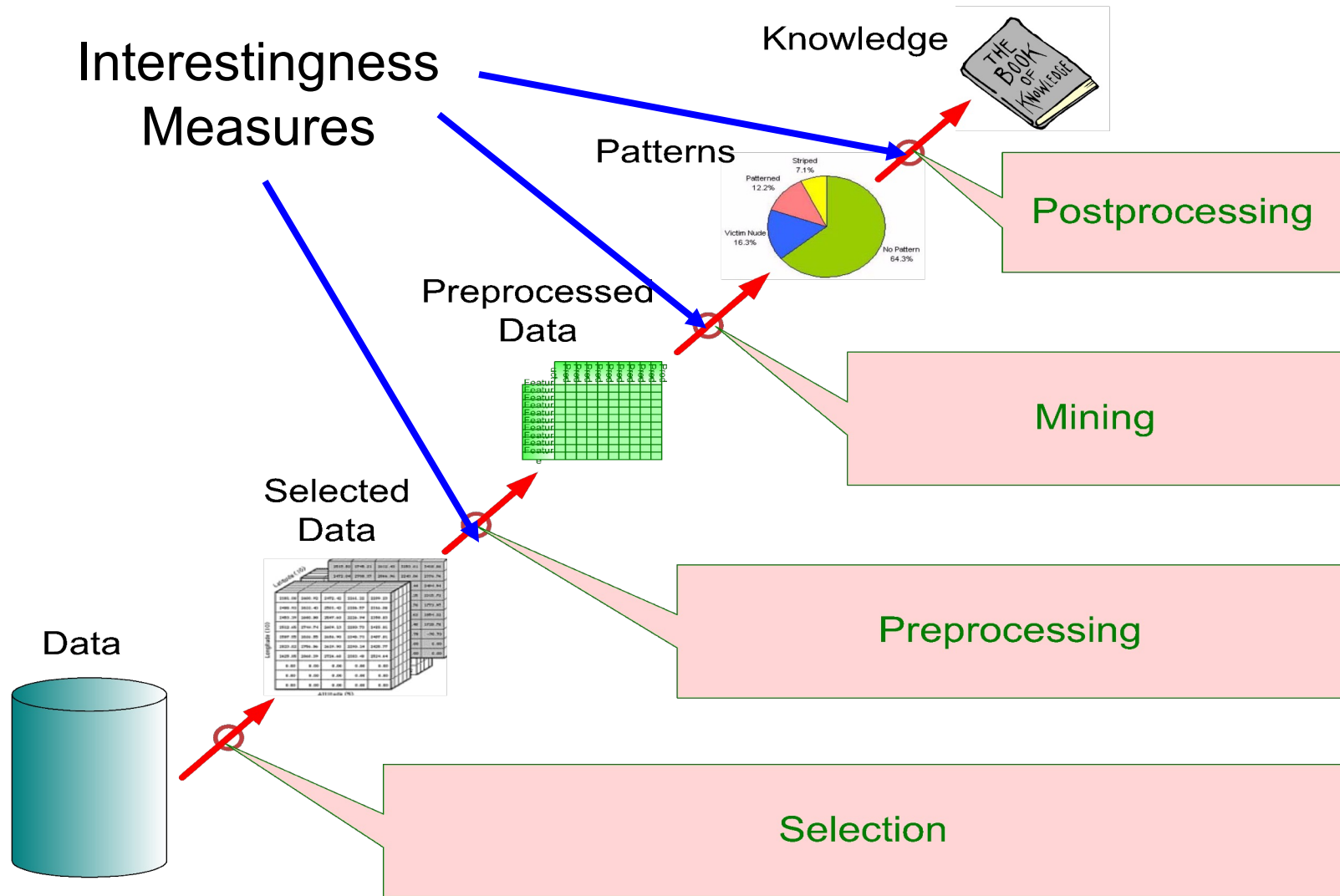
□ Pruning step has to be modified:

- 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

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

f_{11} : support of X and Y

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

f_{01} : support of \overline{X} and Y

f_{00} : support of X and \overline{Y}

Used to define various measures

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

Drawback of Confidence

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

Association Rule: Tea \rightarrow Coffee

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

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

\Rightarrow Although confidence is high, rule is misleading

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

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)

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

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

Statistical-based Measures

- Measures that take into account statistical dependence

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

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

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

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

Example: Lift/Interest

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

Association Rule: Tea \rightarrow Coffee

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

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

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

Drawback of Lift & Interest

	Y	\bar{Y}	
X	10	0	10
\bar{X}	0	90	90
	10	90	100

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

	Y	\bar{Y}	
X	90	0	90
\bar{X}	0	10	10
	90	10	100

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

Statistical independence:

If $P(X,Y)=P(X)P(Y) \Rightarrow Lift = 1$

Properties of A Good Measure

- There are lots of measures proposed in the literature
- **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

Subjective Interestingness Measure

- Objective measure:

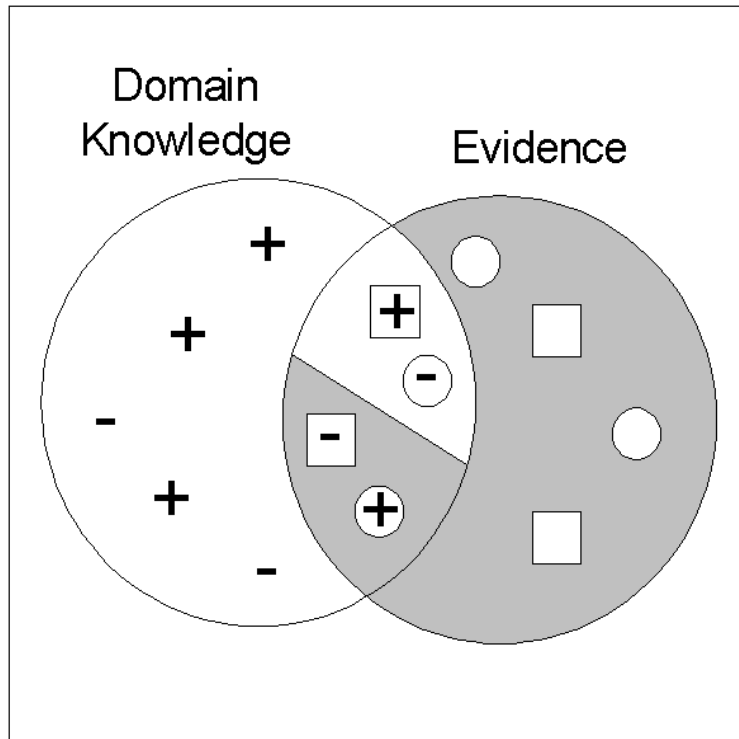
- Rank patterns based on statistics computed from data
- e.g., 21 measures of association (support, confidence, Laplace, Gini, mutual information, Jaccard, etc).

- 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

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)