

Association Rule Mining

CS277

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

Slides are taken from the book “Introduction to Data Mining” by Tan, Steinbach, Karpatne, Kumar

THE TASK

- Two ways of defining the task
- General
 - **Input:** A collection of instances
 - **Output:** rules to predict the values of any attribute(s) (not just the class attribute) from values of other attributes
 - E.g. if temperature = cool then humidity =normal
 - If the right hand side of a rule has only the class attribute, then the rule is a **classification rule**
 - Distinction: Classification rules are applied together as sets of rules
- Specific - Market-basket analysis
 - **Input:** a collection of transactions
 - **Output:** rules to predict the occurrence of any item(s) from the occurrence of other items in a transaction
 - E.g. {Milk, Diaper} -> {Beer}
- General rule structure:
 - Antecedents -> Consequents

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 items one customer buys on one day

Market-Baskets – (2)

- Really a general many-many mapping (association) between two kinds of items
 - But we ask about connections among “items,” not “baskets”
- The technology focuses on **common events**, not rare events

EXAMPLE

Market-Basket transactions

<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

Example of Association Rules

$\{\text{Diaper}\} \rightarrow \{\text{Beer}\},$
 $\{\text{Milk, Bread}\} \rightarrow \{\text{Eggs, Coke}\},$
 $\{\text{Beer, Bread}\} \rightarrow \{\text{Milk}\},$

Implication means co-occurrence,
not causality!

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

<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

Definition: Association Rule

- Association Rule
 - An implication expression of the form $X \rightarrow Y$, where X and Y are itemsets
 - Example:
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

<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

- 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:

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

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

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

Association Rule Mining Task

- Given a set of transactions T , the goal of association rule mining is to find all rules having
 - support \geq *minsup* threshold
 - confidence \geq *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

Example of Rules:

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

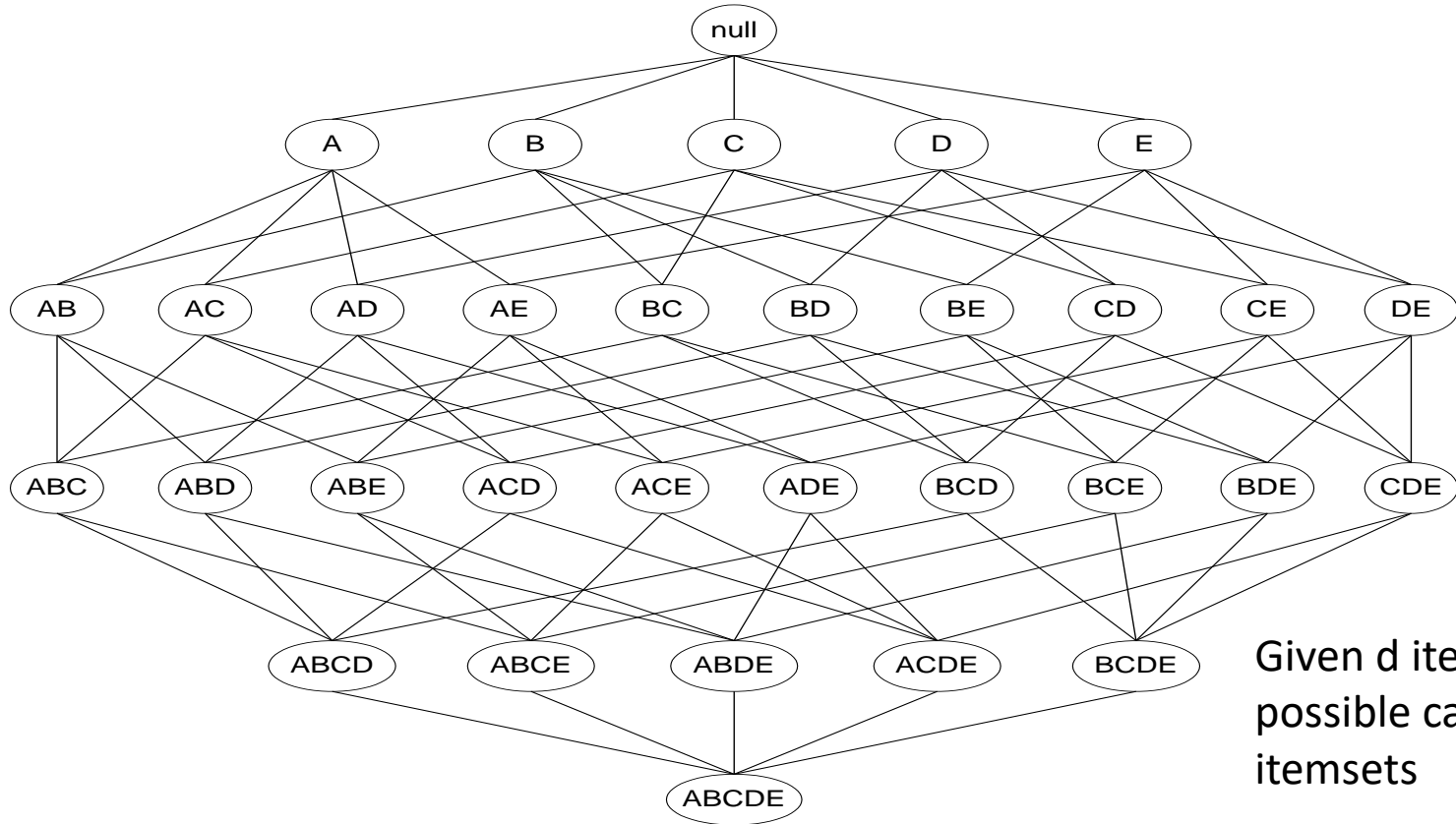
Observations:

- 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** measures
- Thus, we may decouple the support and confidence requirements

Mining Association Rules

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

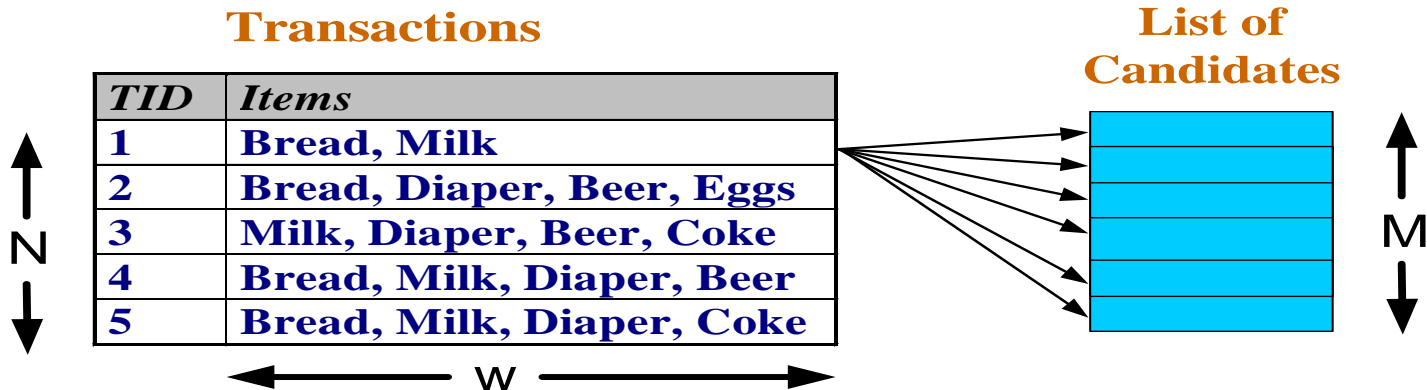
Frequent Itemset Generation



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

Frequent Itemset Generation

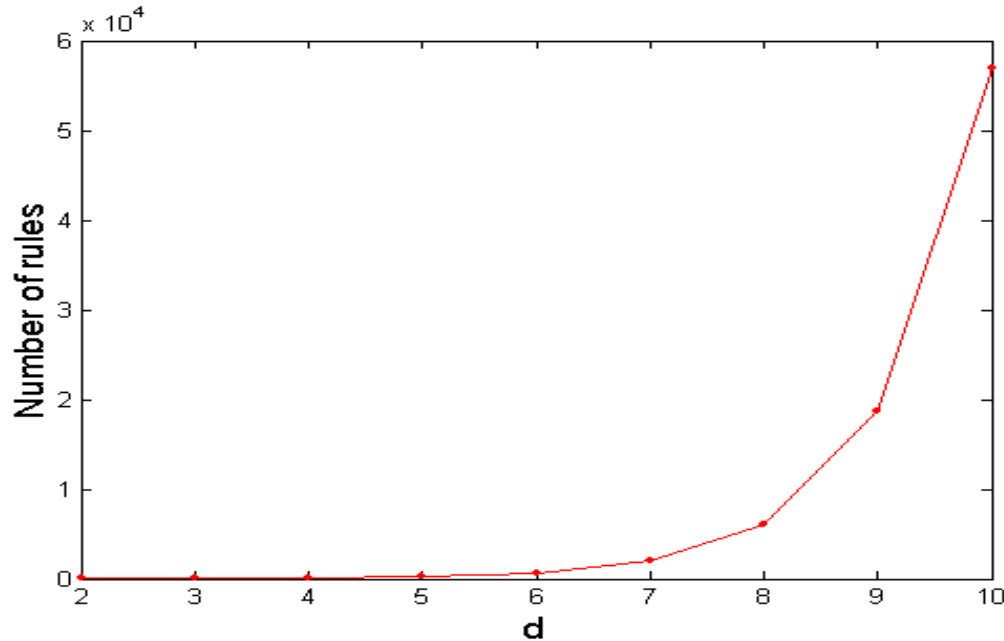
- Brute-force approach:
 - Each itemset in the lattice is a **candidate** frequent itemset
 - Count the support of each candidate by scanning the database



- Match each transaction against every candidate
- Complexity $\sim O(NMw) \Rightarrow$ **Expensive since $M = 2^d$!!!**

Computational Complexity

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



$$R = 3^d - 2^{d+1} + 1$$

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
- 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

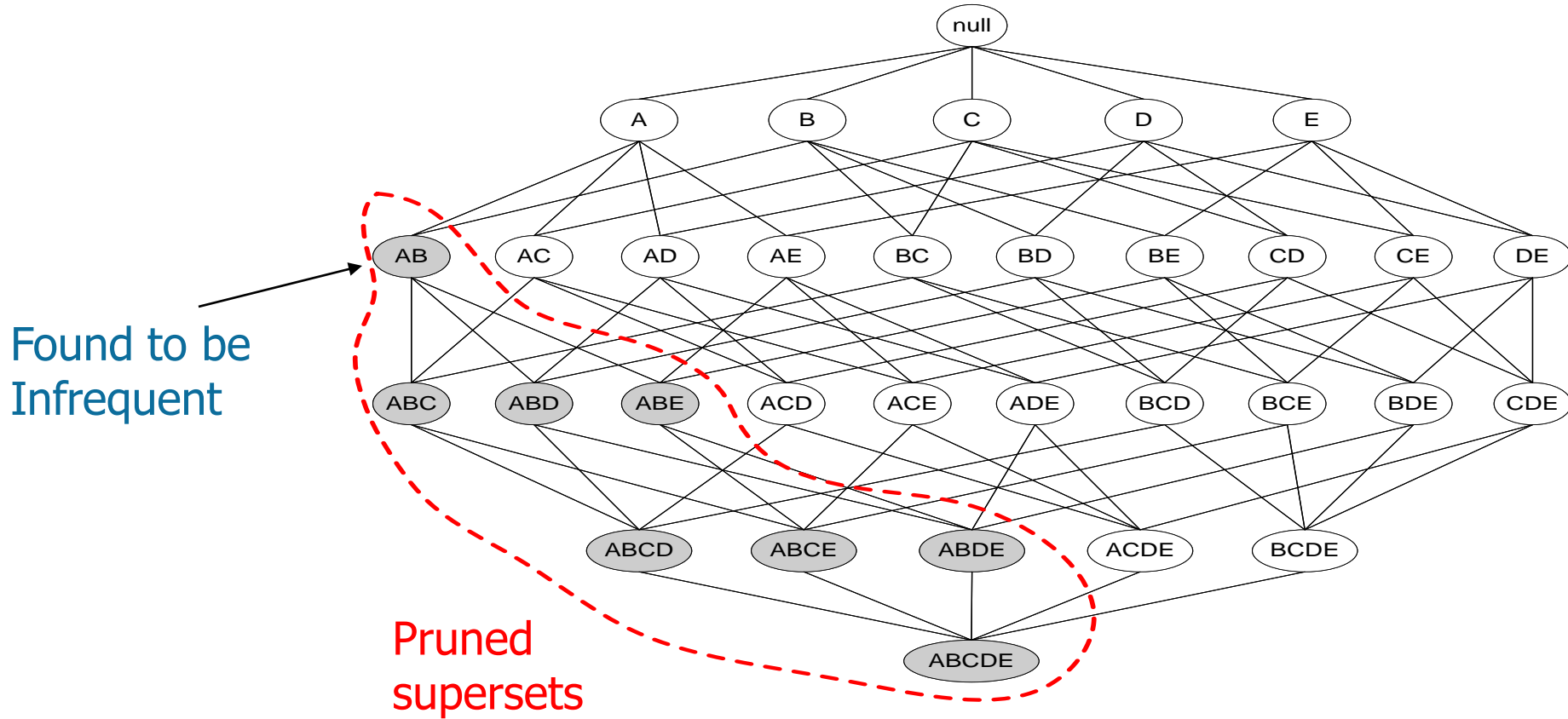
Reducing Number of Candidates

- **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
- 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)



Triplets (3-itemsets)

Itemset	Count
{Bread,Milk,Diaper}	3



Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3 = 41$
With support-based pruning,
 $6 + 6 + 1 = 13$

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

THE APRIORI ALGORITHM: BASIC IDEA

- **Join Step:** C_k is generated by joining L_{k-1} with itself
- **Prune Step:** Any $(k-1)$ -itemset that is not frequent cannot be a subset of a frequent k -itemset

- 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 min_support

end

return $\cup_k L_k$;

THE APRIORI ALGORITHM — EXAMPLE

MinSup=2

Database D

TID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

Scan D

C_1	itemset	sup.
	{1}	2
	{2}	3
	{3}	3
	{4}	1
	{5}	3

L_1

itemset	sup.
{1}	2
{2}	3
{3}	3
{5}	3

L_2

itemset	sup
{1 3}	2
{2 3}	2
{2 5}	3
{3 5}	2

C_2

itemset	sup
{1 2}	1
{1 3}	2
{1 5}	1
{2 3}	2
{2 5}	3
{3 5}	2

Scan D

C_2

itemset
{1 2}
{1 3}
{1 5}
{2 3}
{2 5}
{3 5}

C_3

itemset
{2 3 5}

Scan D

L_3

itemset	sup
{2 3 5}	2

HOW TO GENERATE CANDIDATES?

- Suppose the items in L_{k-1} are listed in an order

- Step 1: self-joining L_{k-1}

insert into C_k

select **$p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$**

from L_{k-1} **p, L_{k-1} q**

where **$p.item_1=q.item_1, \dots, p.item_{k-2}=q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$**

- Step 2: pruning

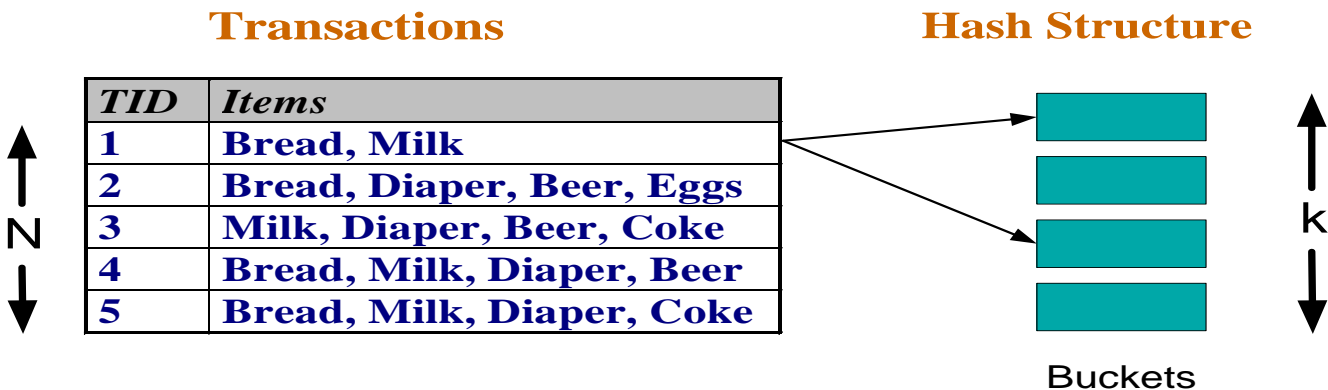
for all ***itemsets* c in C_k** do

for all ***(k-1)-subsets* s of c** do

if (c is not in L_{k-1}) then delete c from C_k

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



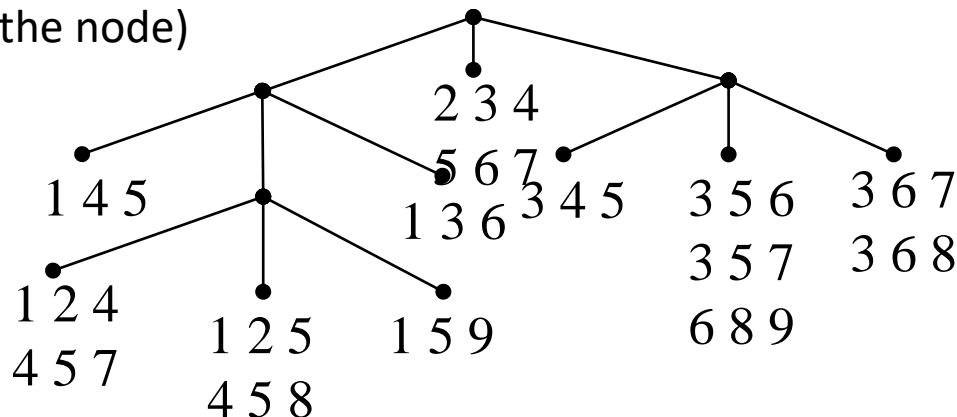
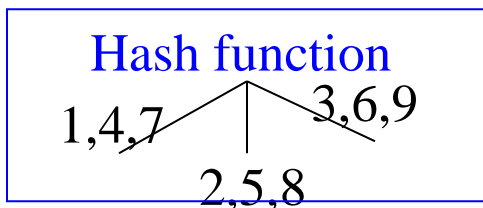
Generate Hash Tree

Suppose you have 15 candidate itemsets of length 3:

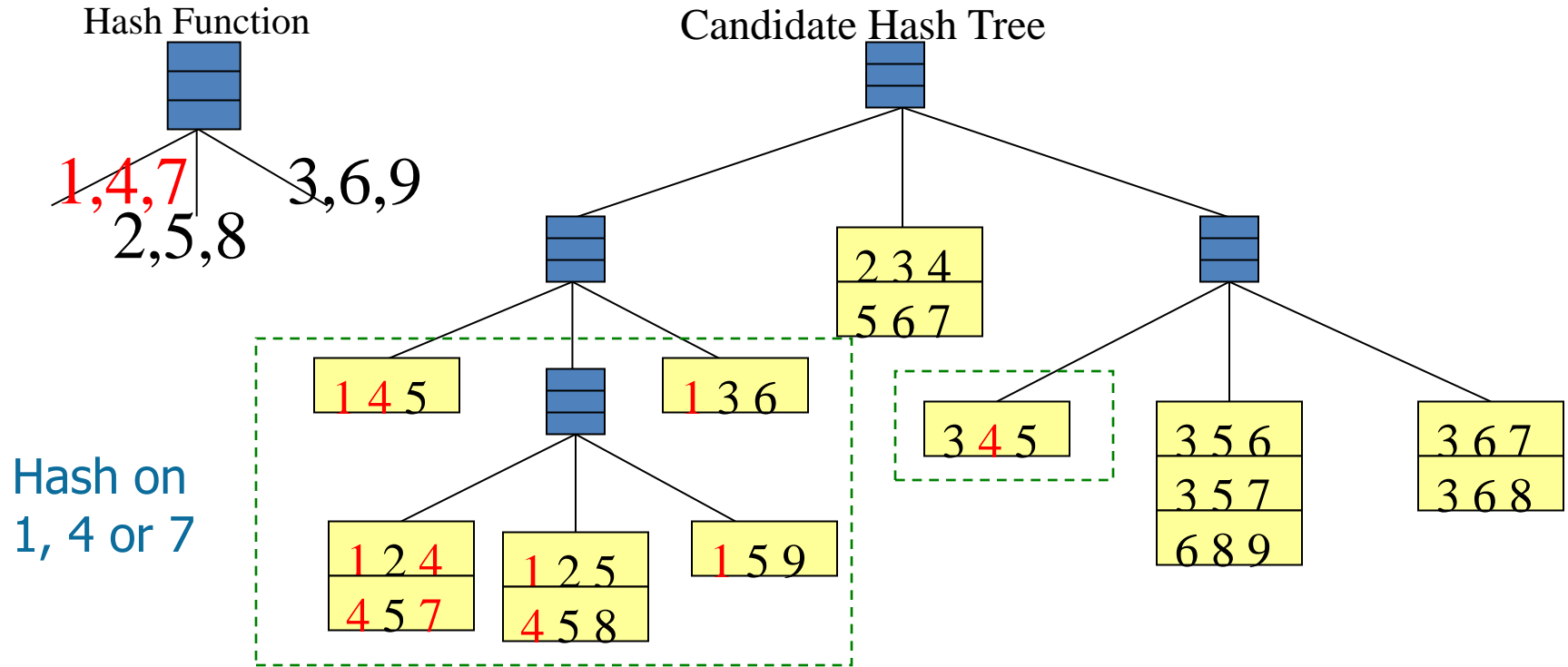
{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4}, {5 6 7}, {3 4 5}, {3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}

You need:

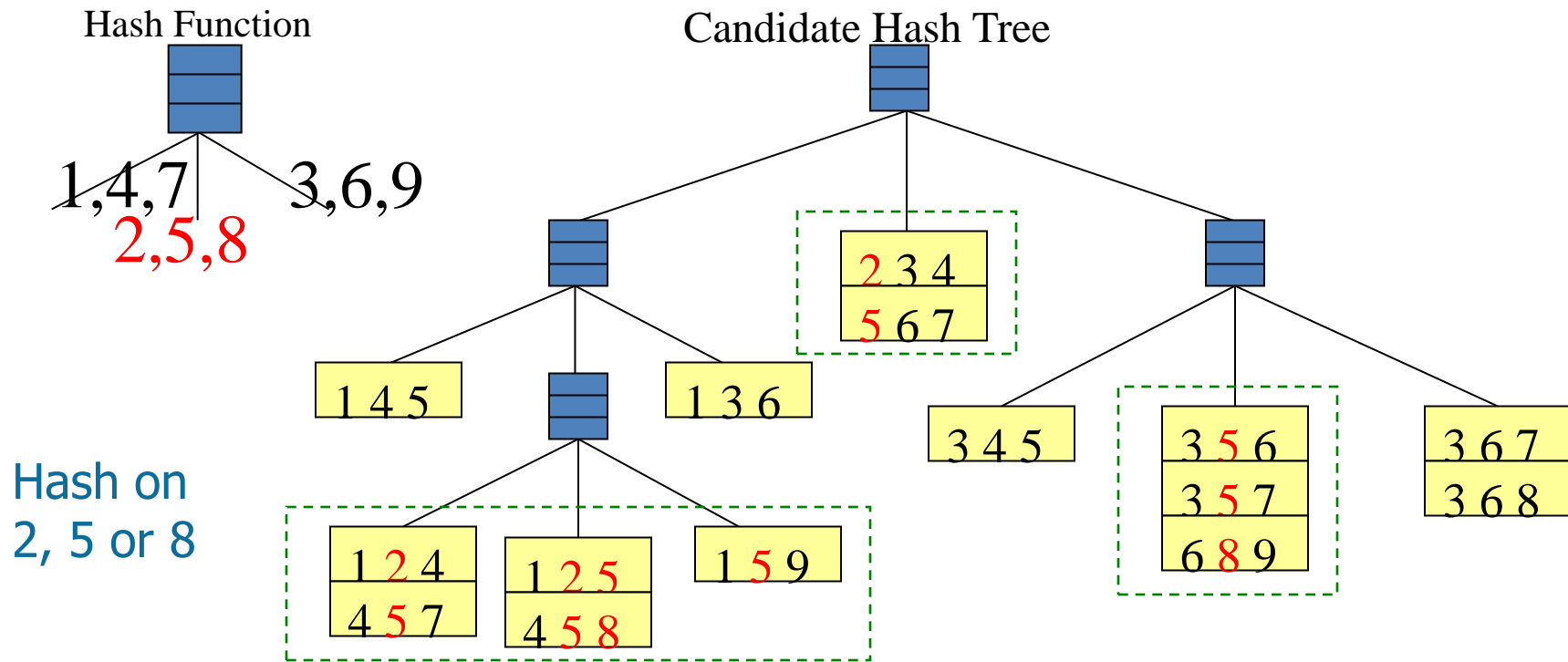
- Hash function: e.g. $h(p) = (p-1) \bmod 3$
- Max leaf size: max number of itemsets stored in a leaf node (if number of candidate itemsets exceeds max leaf size, split the node)



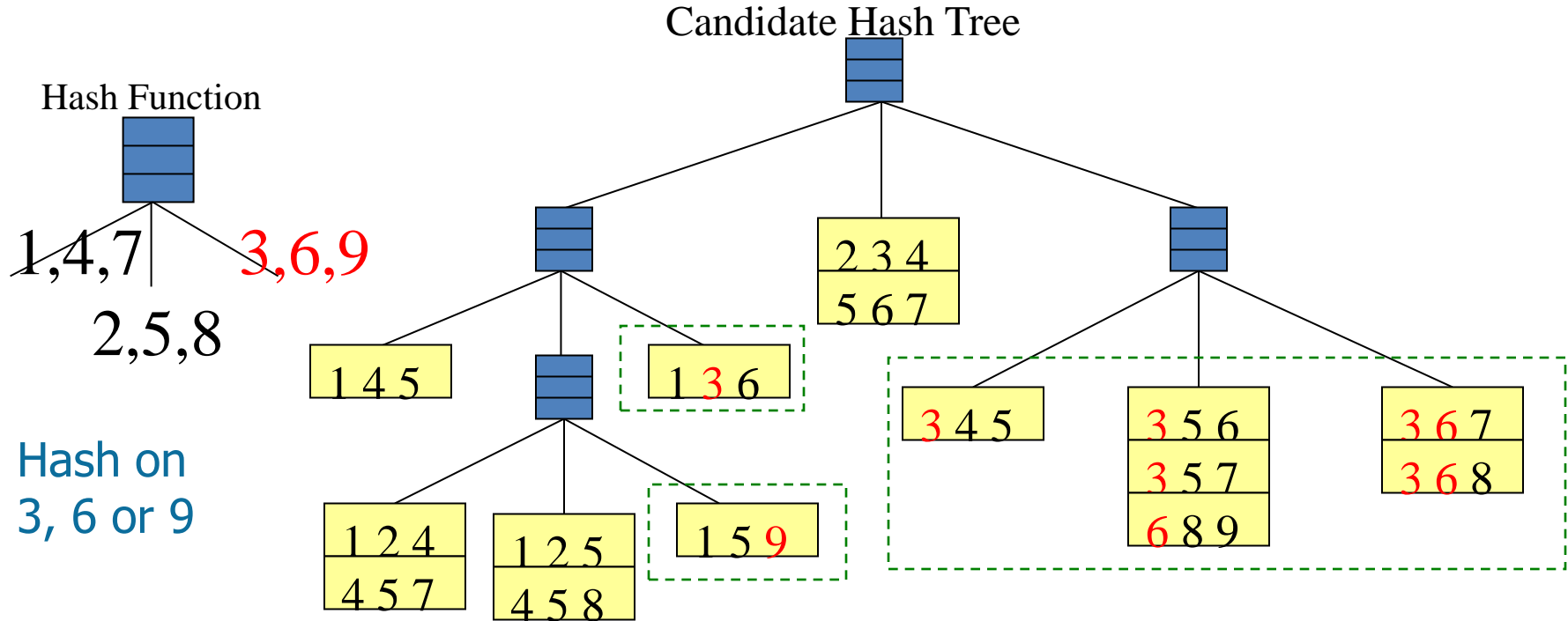
Association Rule Discovery: Hash tree



Association Rule Discovery: Hash tree

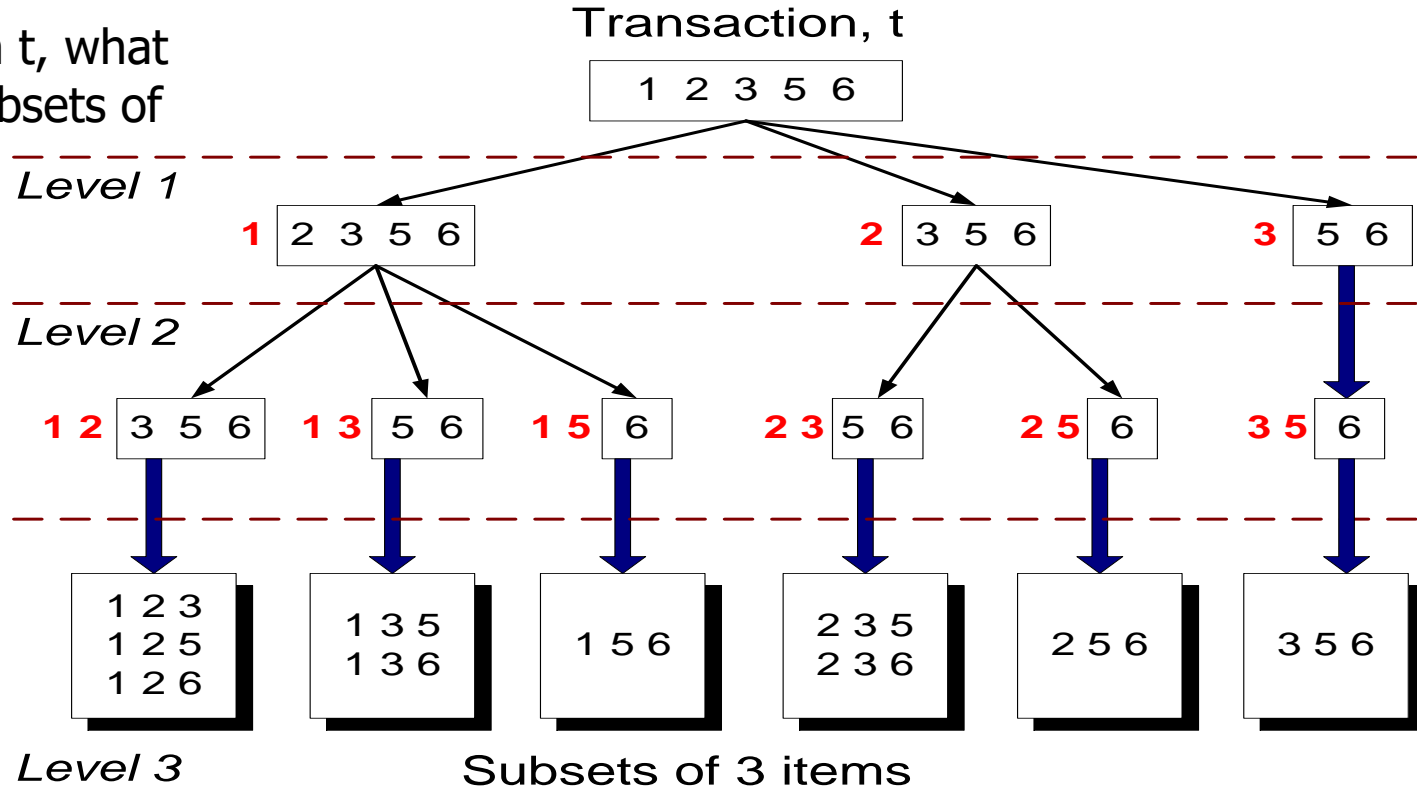


Association Rule Discovery: Hash tree



Subset Operation

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



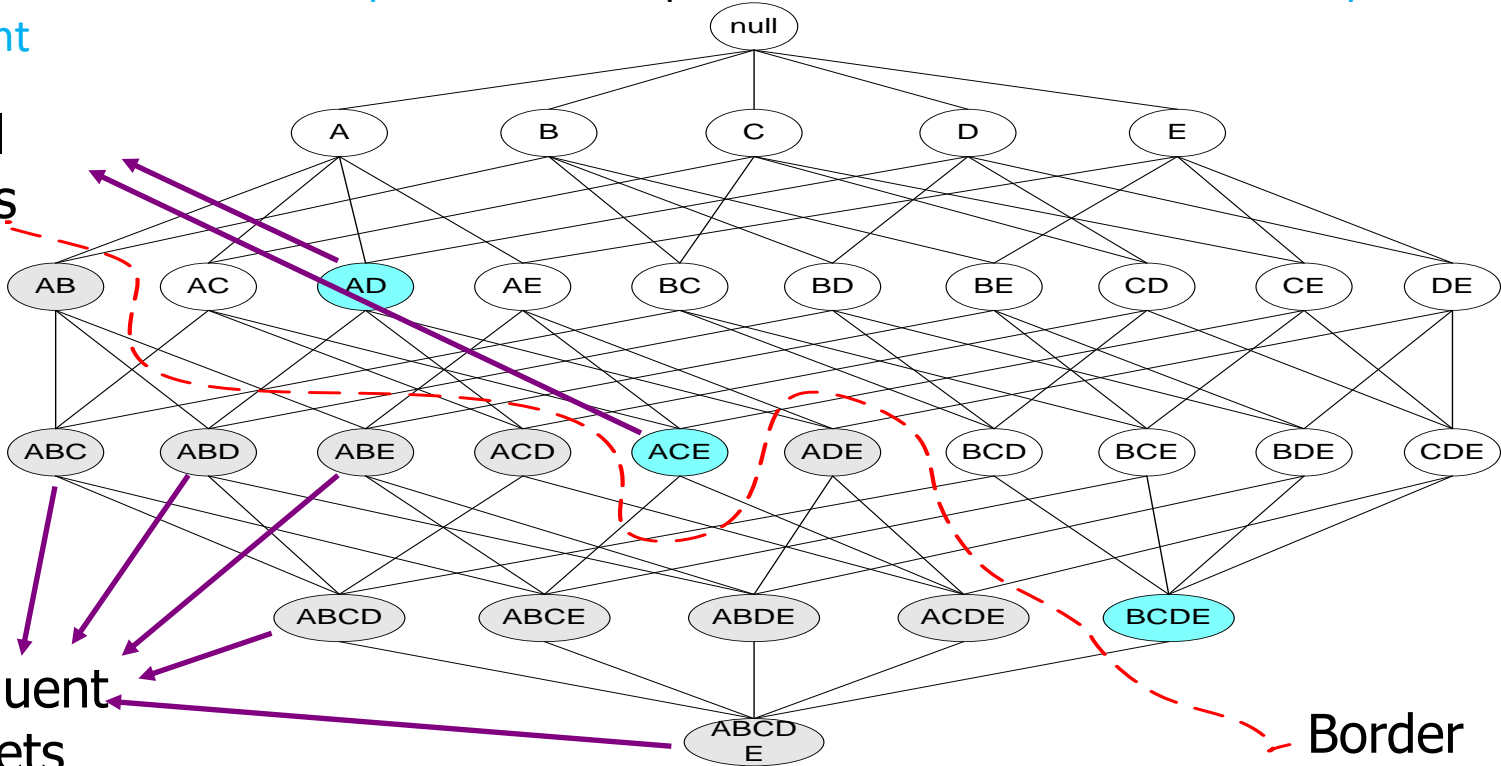
Factors Affecting Complexity

- 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)

Maximal Frequent Itemset

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

Maximal
Itemsets



Infrequent
Itemsets

Border

Closed Itemset

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

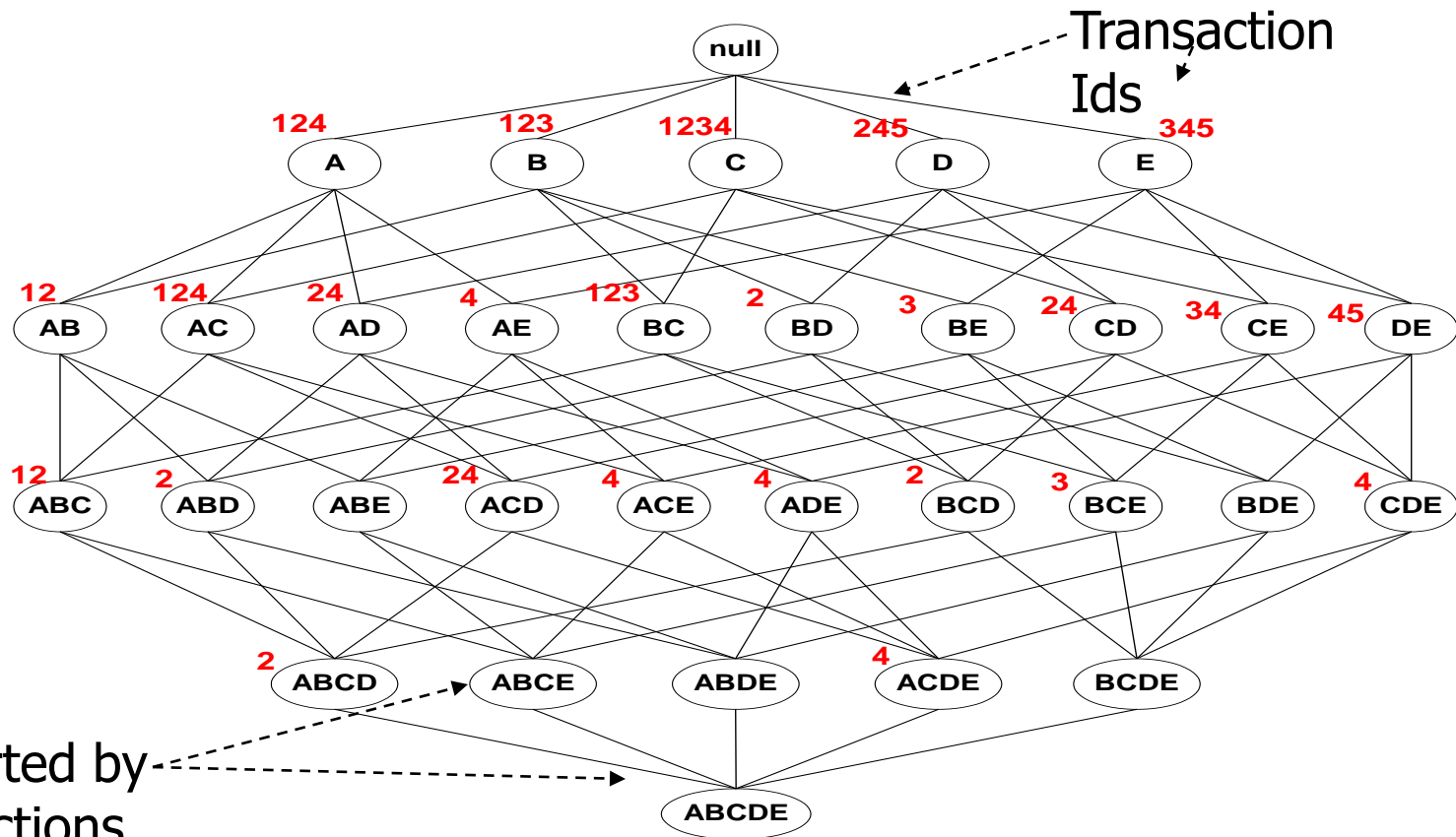
TID	Items
1	{A,B}
2	{B,C,D}
3	{A,B,C,D}
4	{A,B,D}
5	{A,B,C,D}

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

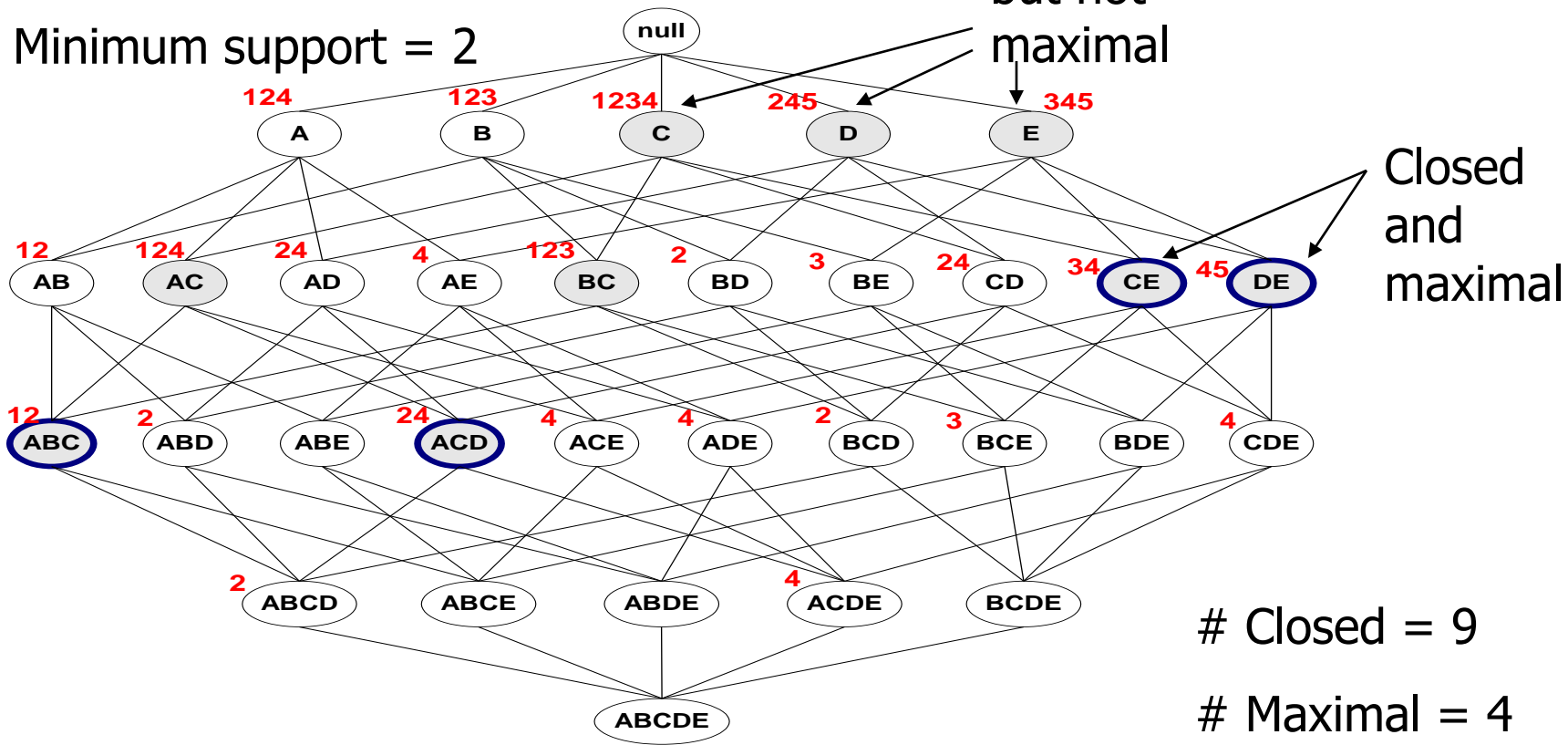
Maximal vs Closed Itemsets

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



Maximal vs Closed Frequent Itemsets

Minimum support = 2



Closed = 9
Maximal = 4

Compact Representation of Frequent Itemsets

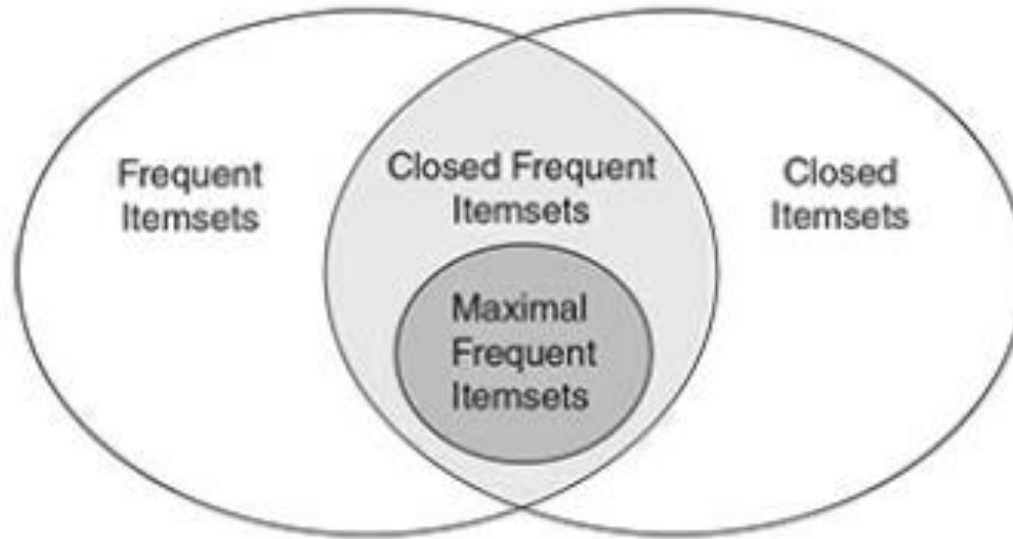
- Some itemsets are redundant because they have identical support as their supersets

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	1	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	1	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	1	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	1	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	1	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	0	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	0	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	0	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	0	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	0	1	1	1	1	1	1	1	1	1

- Number of frequent itemsets
- Need a compact representation

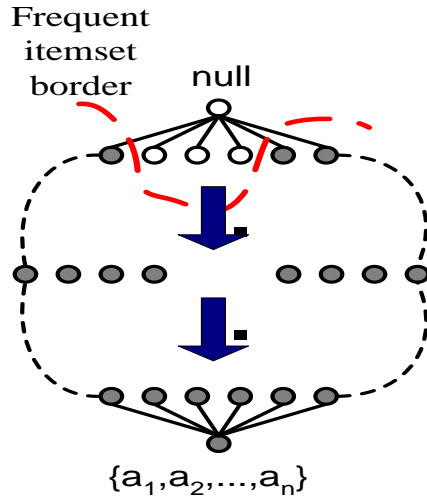
$$= 3 \times \sum_{k=1}^{10} \binom{10}{k}$$

Maximal vs Closed Itemsets

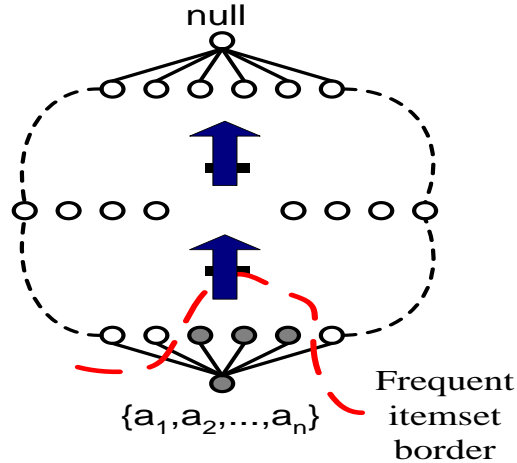


Alternative Methods for Frequent Itemset Generation

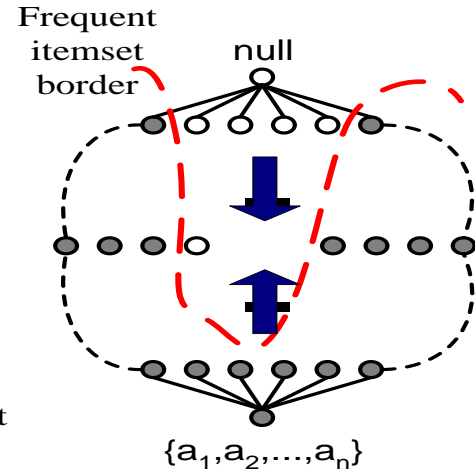
- Traversal of Itemset Lattice
 - General-to-specific vs Specific-to-general



(a) General-to-specific



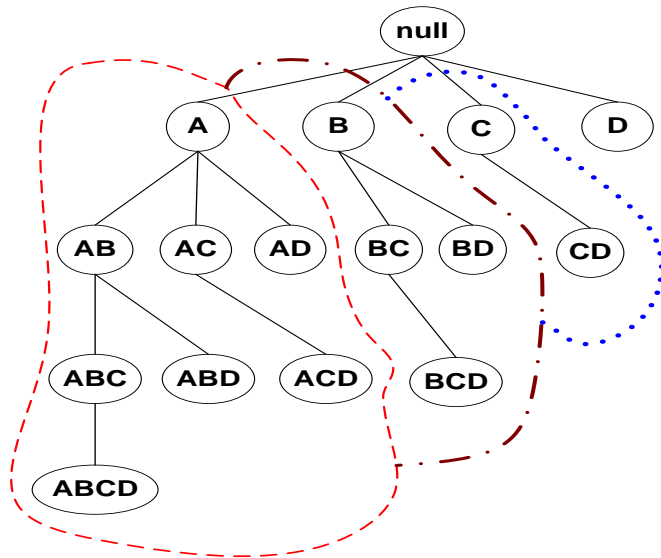
(b) Specific-to-general



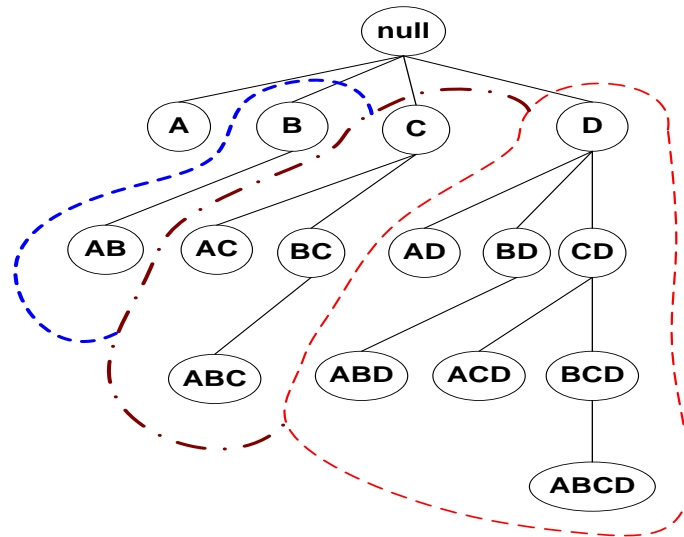
(c) Bidirectional

Alternative Methods for Frequent Itemset Generation

- Traversal of Itemset Lattice
 - Equivalent Classes (two itemsets belong to the same class if they share same common prefix or suffix)



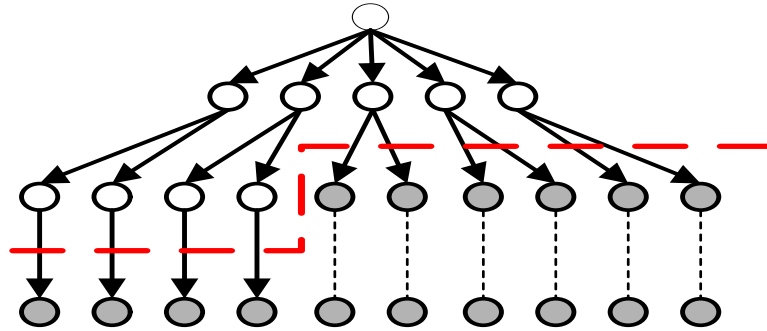
(a) Prefix tree



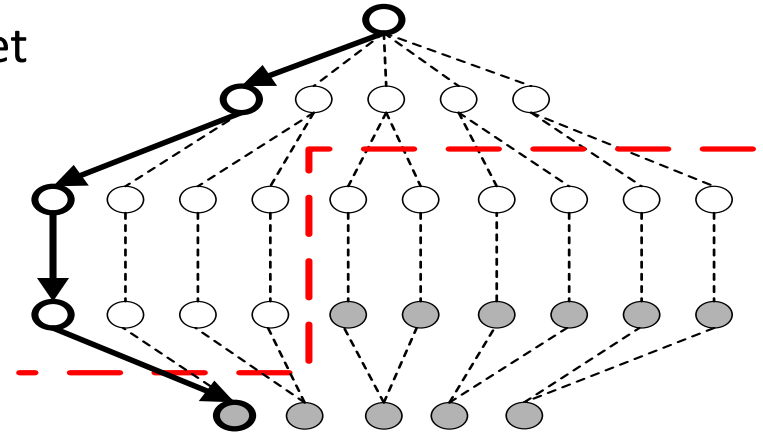
(b) Suffix tree

Alternative Methods for Frequent Itemset Generation

- Traversal of Itemset Lattice
 - Breadth-first vs Depth-first
 - Apriori traverses in BFS manner
 - DFS quickly finds maximal frequent set



(a) Breadth first



(b) Depth first

ECLAT: ANOTHER METHOD FOR FREQUENT ITEMSET GENERATION

- ECLAT: for each item, store a list of transaction ids (tids); 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

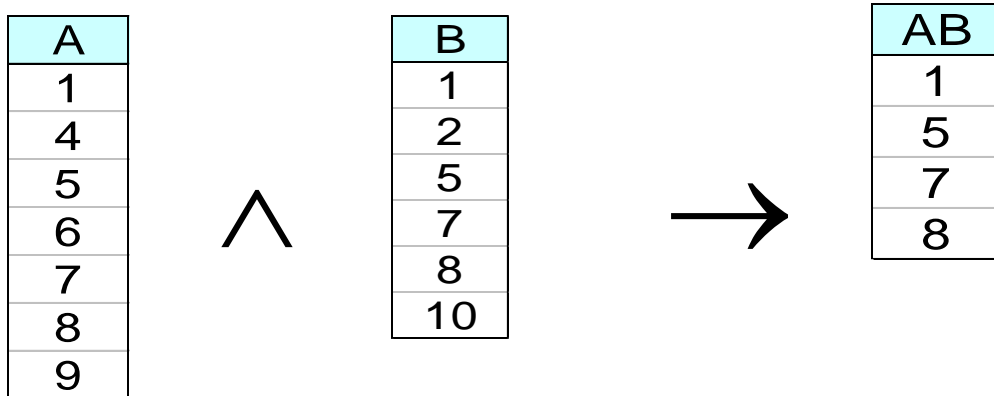
A	B	C	D	E
1	1	2	2	1
4	2	3	4	3
5	5	4	5	6
6	7	8	9	
7	8	9		
8	10			
9				



TID-list

ECLAT: ANOTHER METHOD FOR FREQUENT ITEMSET GENERATION

- 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

FP GROWTH ALGORITHM

- **Apriori**: uses a generate-and-test approach – generates candidate itemsets and tests if they are frequent
 - Generation of candidate itemsets is expensive (in both space and time)
 - Support counting is expensive
 - Subset checking (computationally expensive)
 - Multiple Database scans (I/O)
- **FP-Growth**: allows frequent itemset discovery without candidate itemset generation. Two step approach:
 - **Step 1**: Build a compact data structure called the FP-tree
 - Built using 2 passes over the data-set.
 - **Step 2**: Extracts frequent itemsets directly from the FP-tree

PHASE 1: FP-TREE CONSTRUCTION

- FP-Tree is constructed using 2 passes over the data-set:

Pass 1:

- **Scan** data and find support for each item
- **Discard** infrequent items
- **Sort** frequent items in decreasing order based on their support

Use this order when building the FP-Tree, so common prefixes can be shared.

STEP 1: FP-TREE CONSTRUCTION

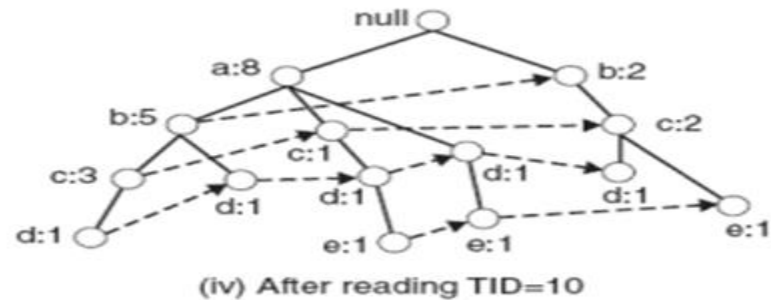
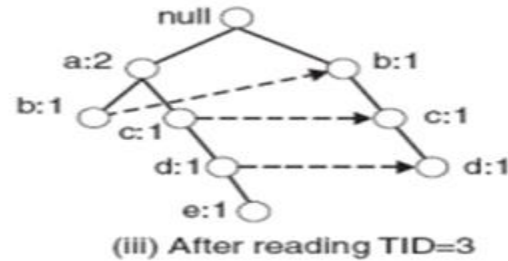
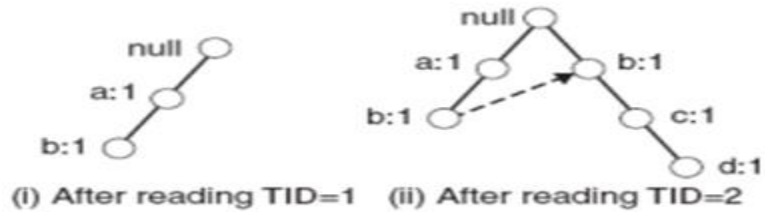
Pass 2:

Nodes correspond to items and have a counter

- ❑ FP-Growth reads 1 transaction at a time and maps it to a path
- ❑ Fixed order is used, so paths can overlap when transactions share items (when they have the same prefix).
- ❑ In this case, counters are incremented
- ❑ Pointers are maintained between nodes containing the same item, creating singly linked lists (dotted lines)
- ❑ The more paths that overlap, the higher the compression. FP-tree may fit in memory.

STEP 1: FP-TREE CONSTRUCTION (EXAMPLE)

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	{a}
8	{a,b,c}
9	{a,b,d}
10	{b,c,e}



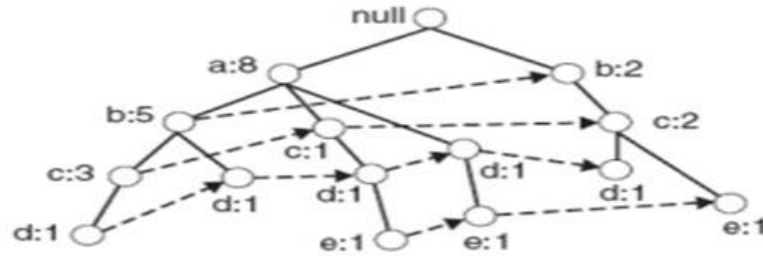
FP-TREE SIZE

- The FP-Tree usually has a **smaller size** than the uncompressed data - typically many transactions share items (and hence prefixes).
 - **Best case scenario:** all transactions contain the same set of items.
 - 1 path in the FP-tree
 - **Worst case scenario:** every transaction has a unique set of items (no items in common)
 - Size of the FP-tree is at least as large as the original data.
 - Storage requirements for the FP-tree are higher - need to store the pointers between the nodes and the counters.
- The size of the FP-tree depends on how the items are ordered
- **Ordering by decreasing support** is typically used but it does not always lead to the smallest tree (it's a heuristic).

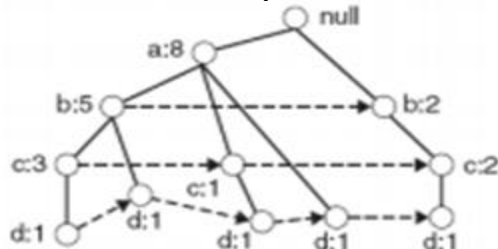
STEP 2: FREQUENT ITEMSET GENERATION

- FP-Growth **extracts frequent itemsets from the FP-tree**
- **Bottom-up algorithm:** from the leaves towards the root
- **Divide and conquer:** first look for frequent itemsets ending in e, then de, etc. . . then d, then cd, etc. . .
- First, **extract prefix path sub-trees ending in an item(set)**. (hint: use the linked lists)

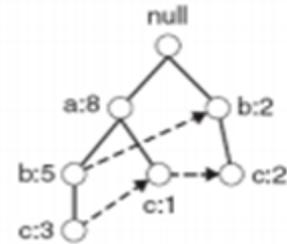
PREFIX PATH SUB-TREES (EXAMPLE)



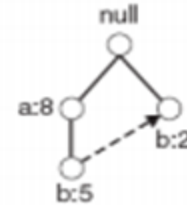
Complete FP-Tree



(b) Paths containing node d



(c) Paths containing node c



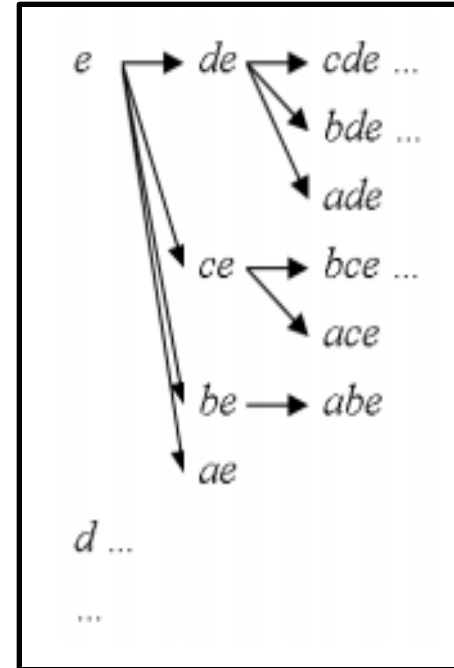
(d) Paths containing node b



(e) Paths containing node a

STEP 2: FREQUENT ITEMSET GENERATION

- Each prefix path sub-tree is now processed **recursively to extract the frequent itemsets**
- Solutions are then merged.
 - E.g. the prefix path sub-tree for *e* will be used to extract frequent itemsets ending in 'e', then in "de", "ce", "be" and "ae", then in "cde", "bde", "cde", etc.
 - Divide and conquer approach

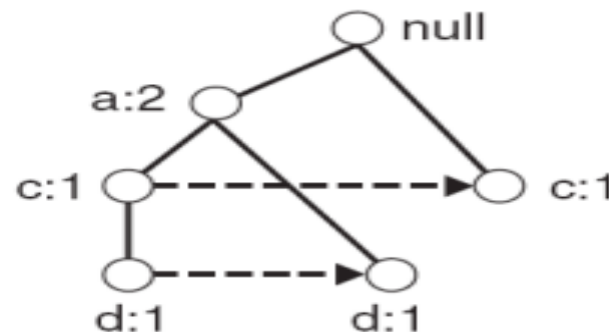


CONDITIONAL FP-TREE

- The FP-Tree that would be built if we only consider transactions containing a particular itemset (and then removing that itemset from all transactions).
- Example: FP-Tree conditional on *e*.

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	{a}
8	{a,b,c}
9	{a,b,d}
10	{b,c,e}

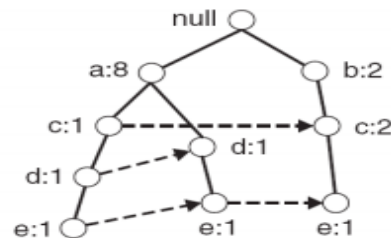
TID	Items
1	{a,b}
2	{b,c,d}
3	{a,c,d, e }
4	{a,d, e }
5	{a,b,c}
6	{a,b,c,d}
7	{a}
8	{a,b,c}
9	{a,b,d}
10	{b,c, e }



EXAMPLE

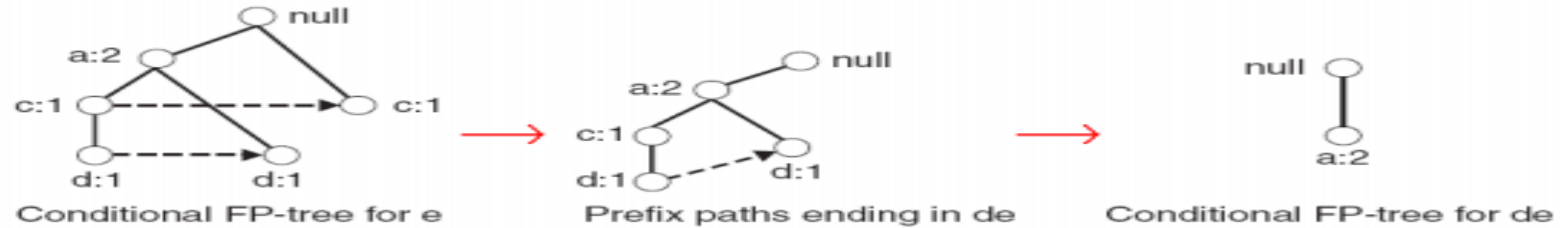
Let $\text{minSup} = 2$ and extract all frequent itemsets containing e

- Obtain the prefix path sub-tree for e :
- Check if e is a frequent item by adding the counts along the linked list (dotted line). If so, extract it.
 - Yes, count = 3 so $\{e\}$ is extracted as a frequent itemset.
- As e is frequent, find frequent itemsets ending in e . i.e. de , ce , be and ae .
- Use the the conditional FP-tree for e to find frequent itemsets ending in de , ce and ae
 - Note that be is not considered as b is not in the conditional FP-tree for e .
- For each of them (e.g. de), find the prefix paths from the conditional tree for e , extract frequent itemsets, generate conditional FP-tree, etc... (recursive)

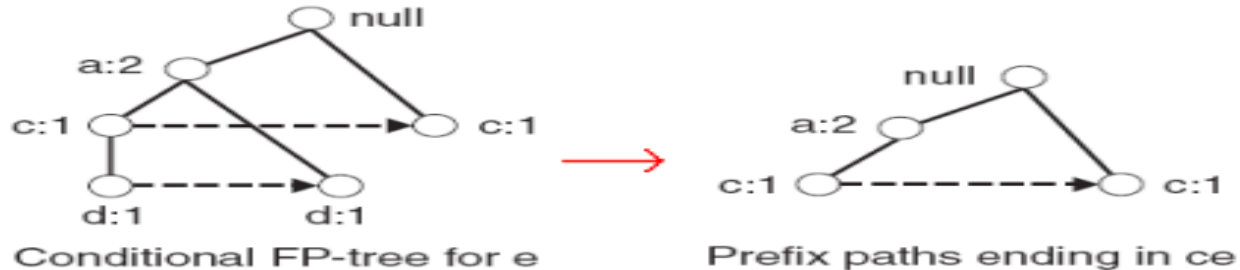


EXAMPLE

- Example: $e \rightarrow de \rightarrow ade$ ($\{d,e\}$, $\{a,d,e\}$ are found to be frequent)



- Example: $e \rightarrow ce$ ($\{c,e\}$ is found to be frequent)



RESULT

Frequent itemsets found (ordered by suffix and order in which they are found):

Transaction
Data Set

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	{a}
8	{a,b,c}
9	{a,b,d}
10	{b,c,e}

Suffix	Frequent Itemsets
e	{e}, {d,e}, {a,d,e}, {c,e}, {a,e}
d	{d}, {c,d}, {b,c,d}, {a,c,d}, {b,d}, {a,b,d}, {a,d}
c	{c}, {b,c}, {a,b,c}, {a,c}
b	{b}, {a,b}
a	{a}

ADVANTAGES AND DISADVANTAGES

- **Advantages** of FP-Growth
 - only 2 passes over data-set
 - “compresses” data-set
 - no candidate generation
 - much faster than Apriori
- **Disadvantages** of FP-Growth
 - FP-Tree may not fit in memory!!
 - FP-Tree is expensive to build

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:
 $ABC \rightarrow D, \quad ABD \rightarrow C, \quad ACD \rightarrow B, \quad BCD \rightarrow A,$
 $A \rightarrow BCD, \quad B \rightarrow ACD, \quad C \rightarrow ABD, \quad D \rightarrow ABC$
 $AB \rightarrow CD, \quad AC \rightarrow BD, \quad AD \rightarrow BC, \quad BC \rightarrow AD, \quad BD \rightarrow AC, \quad CD \rightarrow AB,$
- 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 APRIORI 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 (i.e. *confidence below threshold*)

