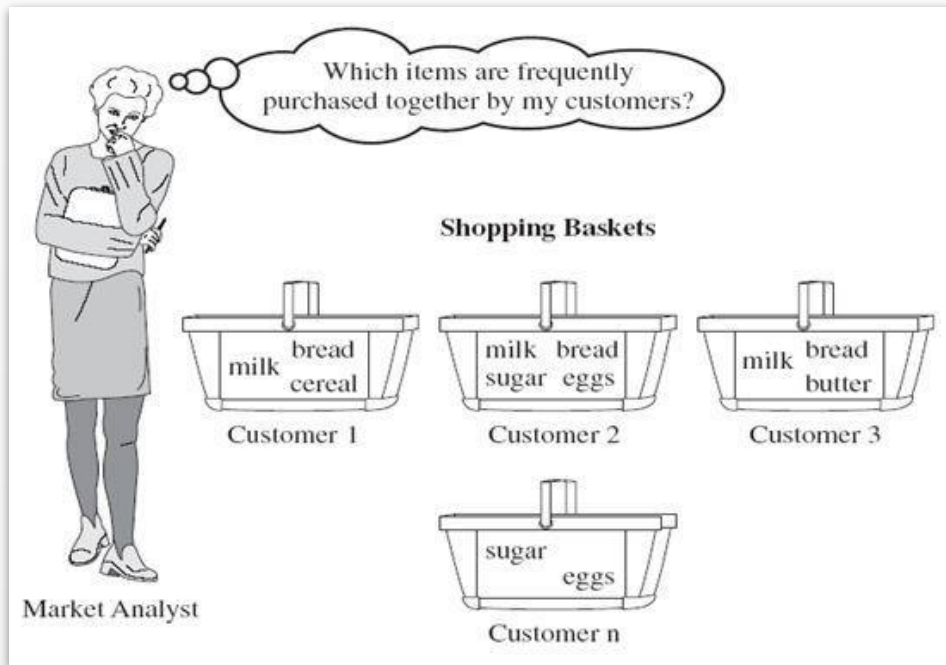


Mining Frequent Itemsets with A-Priori

Market Basket Analysis



Baskets of items

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

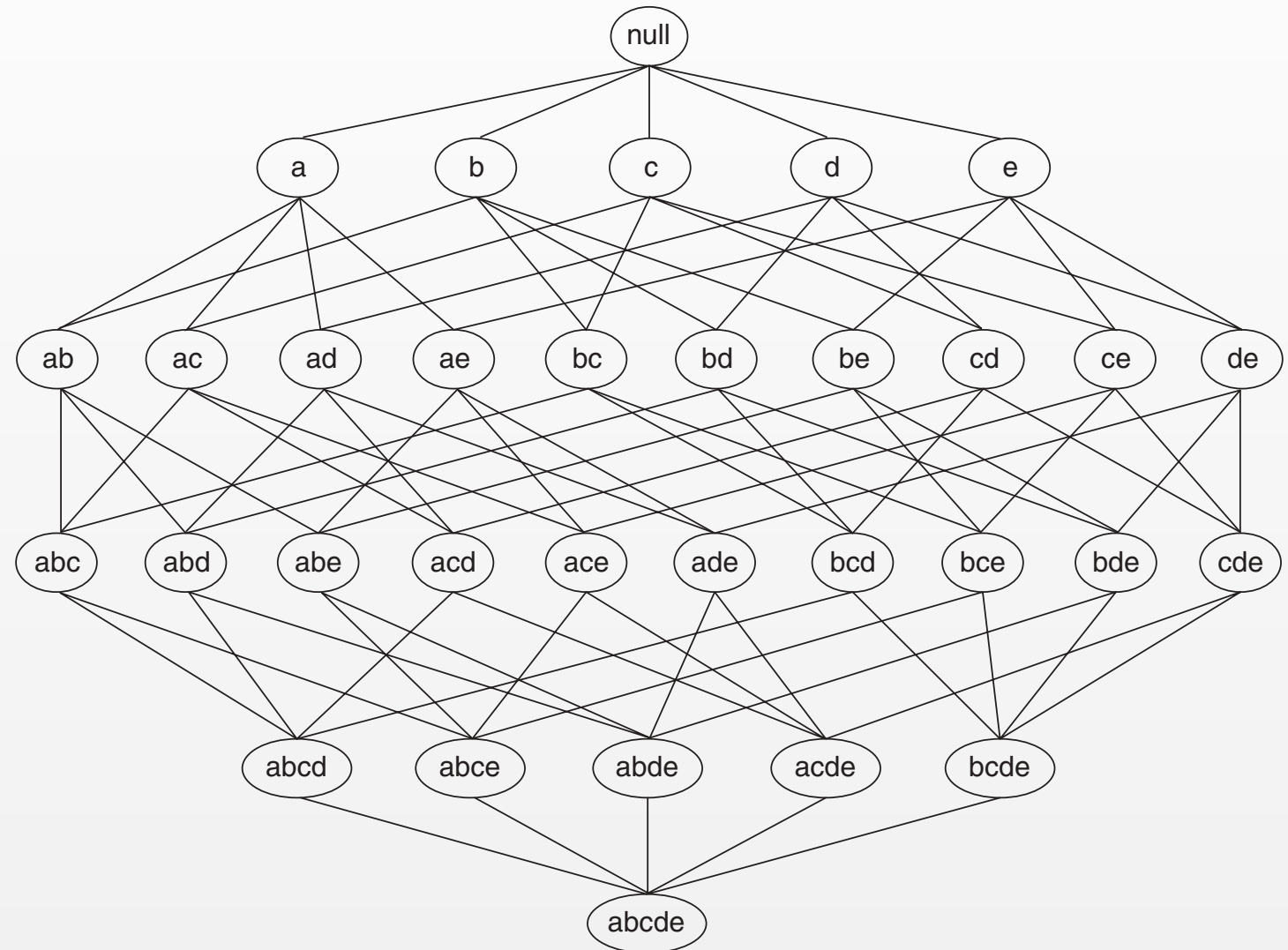
Association Rules

{Milk} --> {Coke}
{Diaper, Milk} --> {Beer}

Finding Frequent Item Sets

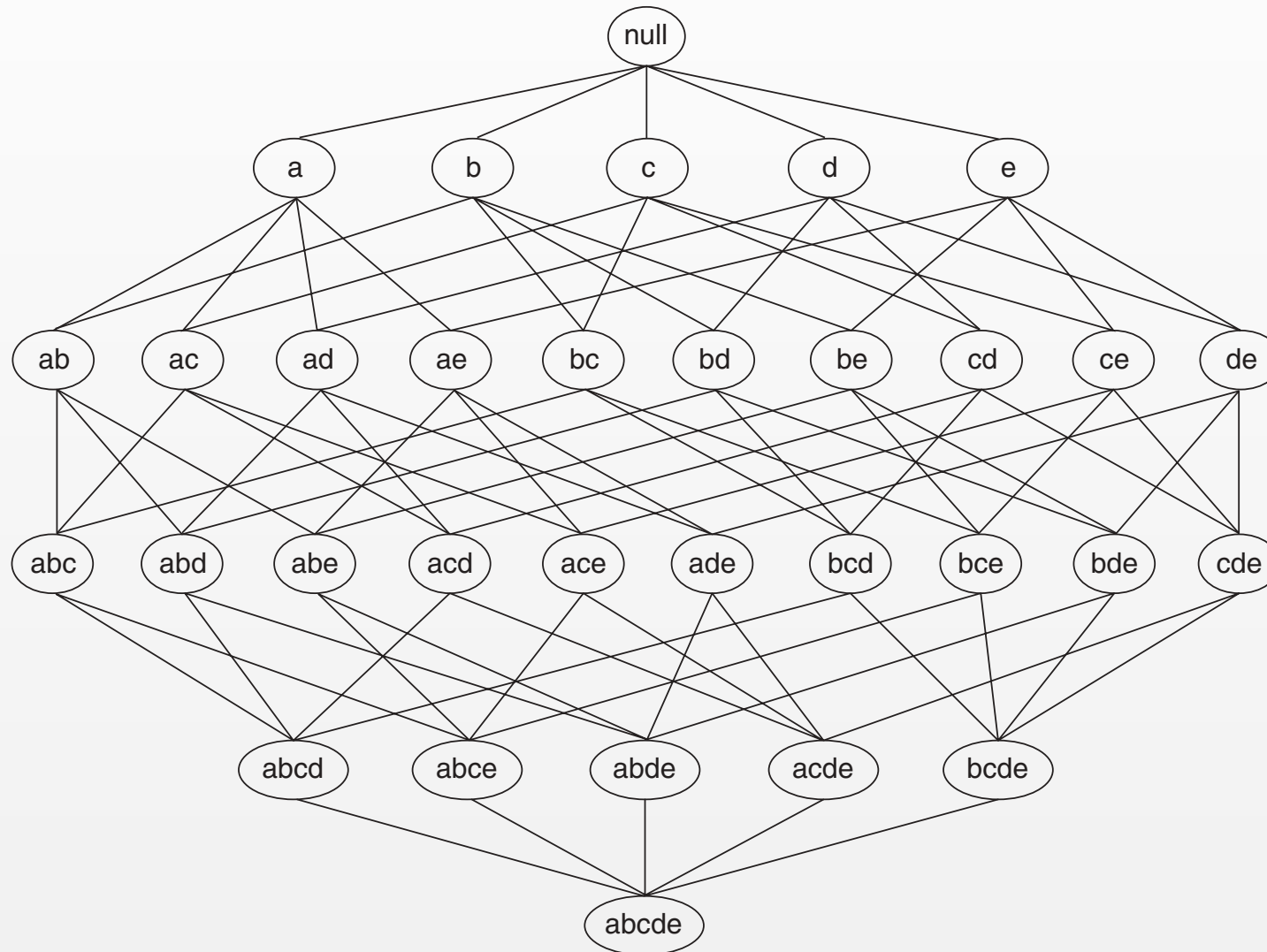
Let I be the set of all items

If $M = |I|$, how many possible itemsets are there?



Finding Frequent Item Sets

Answer: $2^M - 1$; Cannot enumerate all possible sets



Anti-monotone Property

A function f (defined on sets) is said to follow the **anti-monotone** property if

$$\forall A, B \in 2^I : A \subseteq B \Rightarrow f(A) \geq f(B)$$

I is the set of all items

2^I denotes the power set of I

Support follows the anti-monotone property

$$\sigma(A) = | \{ t \in T : A \subset t \} |$$

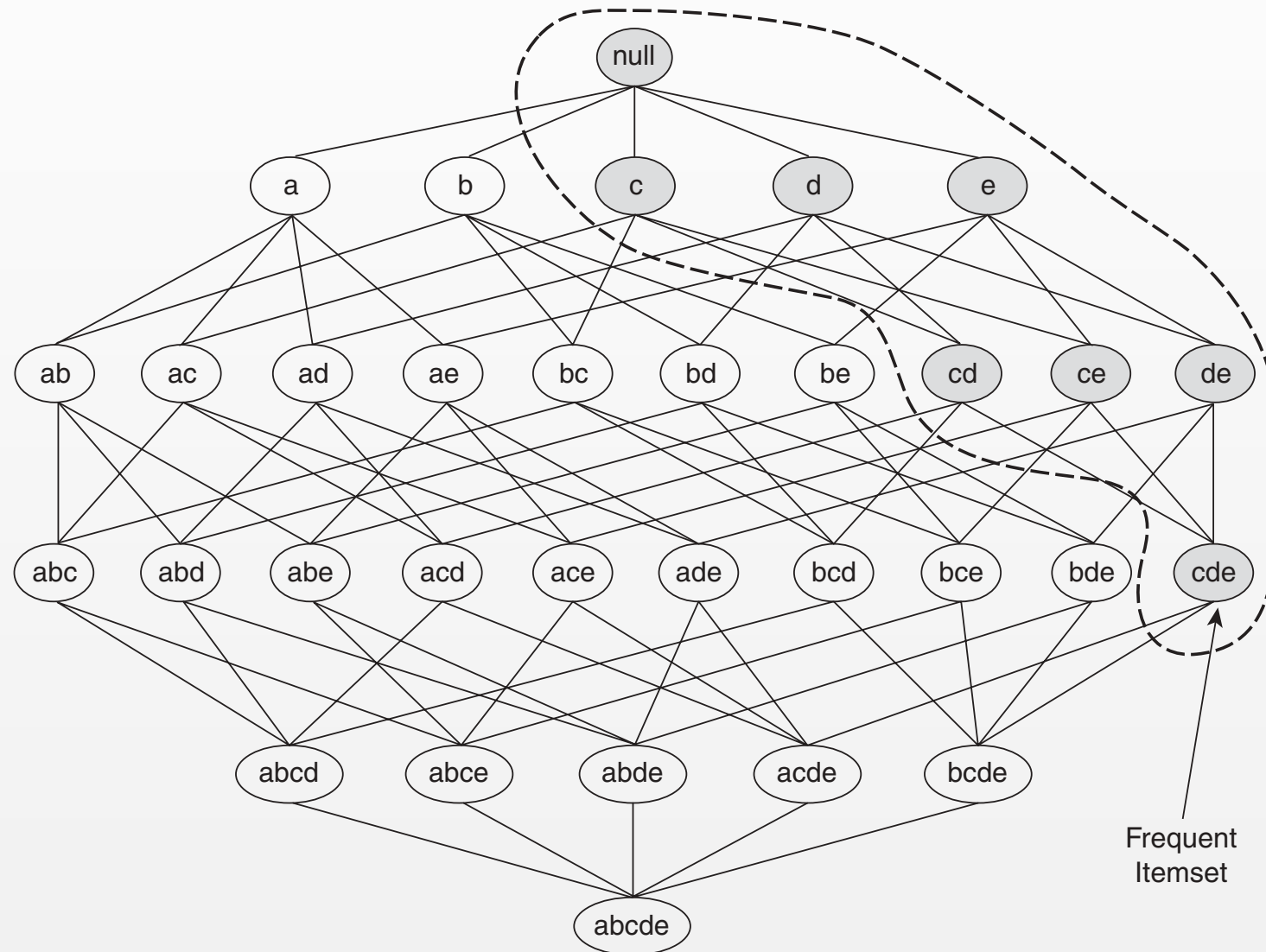
$$\sigma : 2^I \rightarrow \mathbf{N}$$

$\mathbf{N} = \{0, 1, \dots, \infty\}$, the set of natural numbers

T : the set of all transactions

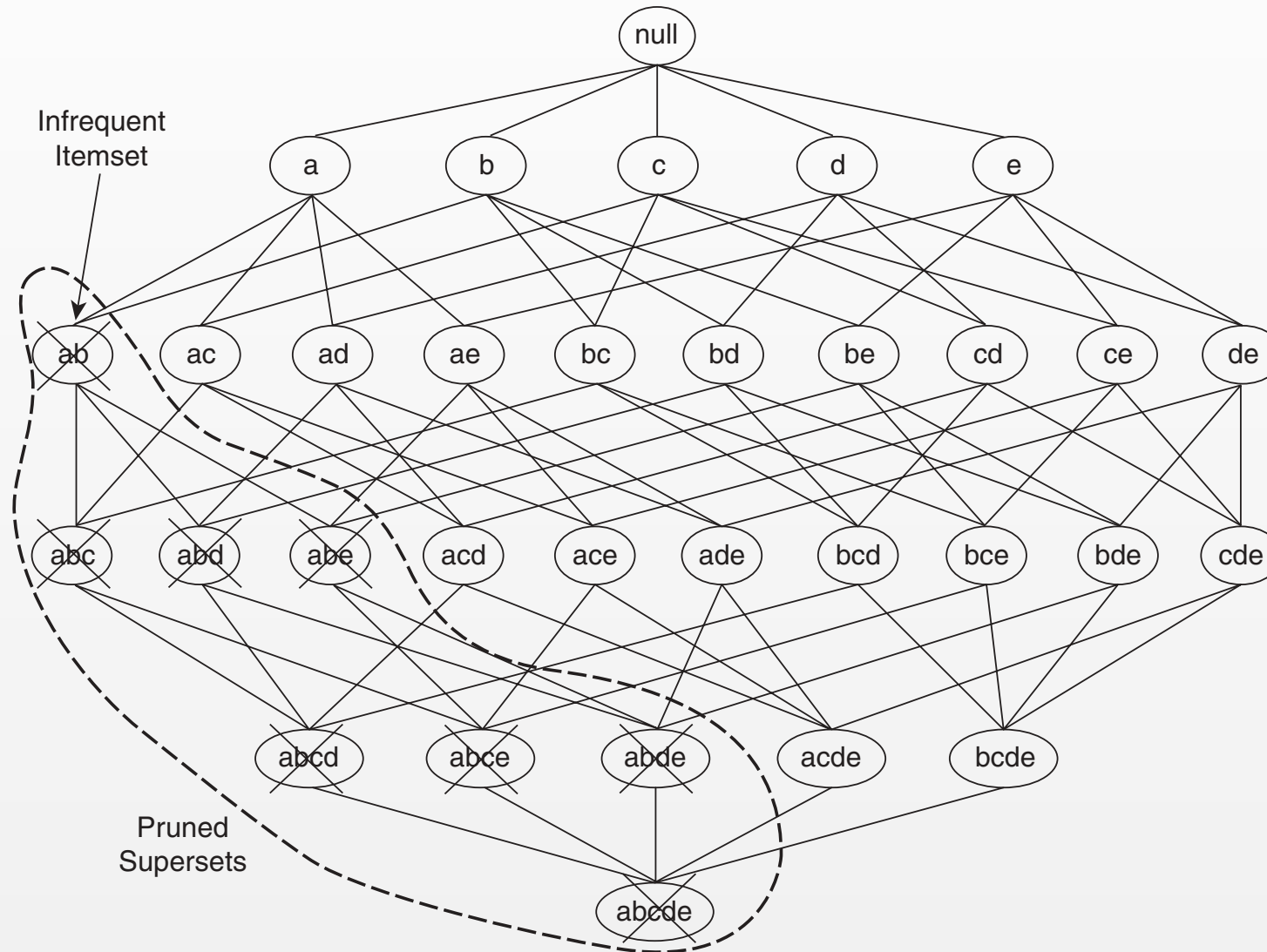
Intuition: A-priori Principle

Observation: Subsets of a frequent item set are **also** frequent



Intuition: A-priori Principle

Corollary: If a set is not frequent, then its supersets are also not frequent



A-priori Algorithm

1. Find all frequent itemsets of size 1
(*only have to check $M = |I|$ possible sets*)
2. For $k = 1, 2, \dots, M$
 - Extend frequent itemsets of size $k - 1$ to create *candidate* itemsets of size k
 - Find candidate sets that are frequent

A-priori Algorithm

Algorithm 6.1 Frequent itemset generation of the *Apriori* algorithm.

```
1:  $k = 1$ .
2:  $F_k = \{ i \mid i \in I \wedge \sigma(\{i\}) \geq N \times \text{minsup} \}$ .    {Find all frequent 1-itemsets}
3: repeat                                     Interpret minsup as a fraction here
4:    $k = k + 1$ .
5:    $C_k = \text{apriori-gen}(F_{k-1})$ .    {Generate candidate itemsets}
6:   for each transaction  $t \in T$  do
7:      $C_t = \text{subset}(C_k, t)$ .    {Identify all candidates that belong to  $t$ }
8:     for each candidate itemset  $c \in C_t$  do
9:        $\sigma(c) = \sigma(c) + 1$ .    {Increment support count}
10:    end for
11:  end for
12:   $F_k = \{ c \mid c \in C_k \wedge \sigma(c) \geq N \times \text{minsup} \}$ .    {Extract the frequent  $k$ -itemsets}
13: until  $F_k = \emptyset$ 
14: Result =  $\bigcup F_k$ .
```

Matching Transactions to Candidate Sets

Algorithm 6.1 Frequent itemset generation of the *Apriori* algorithm.

```
1:  $k = 1$ .
2:  $F_k = \{ i \mid i \in I \wedge \sigma(\{i\}) \geq N \times \text{minsup} \}$ .    {Find all frequent 1-itemsets}
3: repeat
4:    $k = k + 1$ .
5:    $C_k = \text{apriori-gen}(F_{k-1})$ .    {Generate candidate itemsets}
6:   for each transaction  $t \in T$  do
7:      $C_t = \text{subset}(C_k, t)$ .    {Identify all candidates that belong to  $t$ }
8:     for each candidate itemset  $c \in C_t$  do
9:        $\sigma(c) = \sigma(c) + 1$ .    {Increment support count}
10:    end for
11:  end for
12:   $F_k = \{ c \mid c \in C_k \wedge \sigma(c) \geq N \times \text{minsup} \}$ .    {Extract the frequent  $k$ -itemsets}
13: until  $F_k = \emptyset$ 
14: Result =  $\bigcup F_k$ .
```

Generating Candidates C_k

Objectives

1. ***No Duplicates:*** A candidate itemsets must be unique.
2. ***Completeness:*** At least, all frequent k-itemsets should be included.
3. ***No infrequent subsets:*** A candidate should not have any infrequent subset.

Generating Candidates C_k

Questions

1. *How many k -itemsets are there?*
2. *How to reduce number of candidates using the computations already performed?*

1. $F_{k-1} \times F_1$: *combine frequent $(k-1)$ -itemsets with frequent 1-itemsets to get candidates of size k .*

$$\{a, b, c\} \cup \{d\} = \{a, b, c, d\}$$

2. $F_{k-1} \times F_{k-1}$: *combine frequent $(k-1)$ -itemsets with other frequent $(k-1)$ -itemsets that differs in only 1 item to get candidates of size k .*

$$\{a, b, c\} \cup \{a, b, d\} = \{a, b, c, d\}$$

Generating Candidates C_k

Duplicates

Combining sets arbitrary pairs of sets from F_{k-1} and F_1 will lead to duplicate candidates

$$\{a, b, c\} \cup \{d\} = \{a, b, c, d\}$$

$$\{a, b, d\} \cup \{c\} = \{a, b, c, d\}$$

Each candidate of size k
could be generated k
times

Combining sets arbitrary pairs of sets from F_{k-1} will lead to duplicate candidates

$$\{a, b, c\} \cup \{a, b, d\} = \{a, b, c, d\}$$

$$\{a, b, c\} \cup \{a, c, d\} = \{a, b, c, d\}$$

Each candidate of size k
could be generated
 $\binom{k}{k-2}$ times

Generating Candidates C_k

Solution to duplicates

Sort the items

- **Item ordering:** Define an ordering on all items
 - Either by assigning a unique id to each item. The items are ordered based on their ID.
 - Or by a lexicographic ordering on the item string; e.g.
'coke' < 'cookie'
- Assume that the items in the itemsets in F_{k-1} are sorted. $a_1 < a_2 < \dots < a_k$ in $A = \{a_1, a_2, \dots, a_k\}$

ID based ordering will have computational advantage since comparing numbers is cheaper than comparing strings

Generating Candidates C_k

Solution to duplicates

$F_{k-1} \times F_1$: Combine $A \in F_{k-1}$ and $B = \{b\} \in F_1$ only if $\forall a \in A \quad a < b$.

- Combine $\{a, c, e\}$ with $\{f\}$ to give candidate $\{a, c, e, f\}$.
- Do not combine $\{a, c, e\}$ with $\{d\}$.
- **No duplicates:** Each candidate has only one way of being generated. $\{a, c, e, f\}$ can only be generated by combining $\{a, c, e\}$ and $\{f\}$.
- **Completeness:** If $\{a, c, e, f\}$ is indeed frequent, $\{a, c, e\}$ and $\{f\}$ have to be present in F_3 and F_1 . And they would get combined to generate $\{a, c, e, f\}$.
- Subsets of generated candidates might still be infrequent

Combine if all items in A are less than the only item in B

Generating Candidates C_k

Solution to duplicates

$F_{k-1} \times F_{k-1}$: Combine $A \in F_{k-1}$ and $B \in F_{k-1}$ only if $a_i = b_i$, for $i = 1, 2, \dots, k-2$ and $a_{k-1} < b_{k-1}$.

- Combine $\{a, c, e\}$ with $\{a, c, f\}$ to give candidate $\{a, c, e, f\}$.
- Do not combine $\{a, c, e\}$ with $\{a, b, e\}$.
- **No duplicates:** Each candidate has only one way of being generated. $\{a, c, e, f\}$ can only be generated by combining $\{a, c, e\}$ and $\{a, c, f\}$.
- **Completeness:** If $\{a, c, e, f\}$ is indeed frequent, $\{a, c, e\}$ and $\{a, c, f\}$ have to be present in F_3 . And they would get combined to generate $\{a, c, e, f\}$.
- Subsets of generated candidates might still be infrequent

Combine if the first $k-2$ items in A and B are the same and the last element of A is less than that of B .

Generating Candidates C_k

How to efficiently find itemsets that could be combined?

- **Itemset Ordering:** Use the ordering on items to define an ordering of itemsets
 - Assume that the items in each itemset are presorted.
 $a_1 < a_2 < \dots < a_k$ in $A = \{a_1, a_2, \dots, a_k\}$
 - $A < B$ if $a_i < b_i$, where i is the index of the first item differing in A and B .
 $\{\text{apple, bread, coke, sauce}\} < \{\text{apple, bread, cookie, milk}\}$ or $\{4, 7, 21, 50\} < \{4, 7, 25, 40\}$

Elements of F_{k-1} sorted with the itemset ordering

$\{a, b, c\}, \{a, b, e\}, \{a, b, g\}, \{a, c, d\}, \{a, c, g\} \dots$

$F_{k-1} \times F_{k-1}$

$\{a, b, c\}$ can't be combined with any itemset beyond $\{a, b, g\}$. So no need to compare beyond $\{a, c, d\}$

Without exploiting the itemset ordering
 $|F_{k-1}|(|F_{k-1}| - 1)/2$ comparisons
need to be made

Generating Candidates C_k

Pruning candidates with infrequent subsets

- For a candidate of size k , one only needs to check subsets of size $k - 1$.
- Enumerate subsets of size $k - 1$ by removing one element at a time from the candidate.
- Search for the subsets one after the other in F_{k-1} until a subset is not found or the list of subsets is exhausted
- Binary search could be performed if F_{k-1} is sorted under the itemset ordering for an efficient search.
- Alternatively a hash tree could be build to store the itemsets of F_{k-1} for an efficient search.
- If a subset wasn't found the candidate should be discarded.

If all size $k - 1$ subsets are frequent so are subsets of smaller sizes.

Each candidate of size k will give k subsets of size $k - 1$

For each candidate $O(k |F_{k-1}|)$ comparisons might be needed in the worst case, if searching naively.

This can be improved by binary search $O(k \log(|F_{k-1}|))$ or constructing hash table

Generating Candidates C_k

1. *Self-joining*: Find pairs of sets in F_{k-1} that have first $k - 2$ items in common and differ by **one** element.
2. *Pruning*: Remove all candidates with infrequent subsets

Example: Generating Candidates C_k

- *Frequent itemsets of size 2:*
 $\{b,c\}:5, \{b,m\}:4, \{c,j\}:3 \{c,m\}:3$
- *Self-joining:*
 $\{b,c,m\}, \{c,j,m\}$
- *Pruning:*
 ~~$\{c,j,m\}$~~ since ~~$\{j,m\}$~~ not frequent
- Frequent items of size 3:
 $\{b,c,m\}$

$$B_1 = \{b, c, m\} \quad B_2 = \{j, m, p\}$$

$$B_3 = \{b, m\} \quad B_4 = \{c, j\}$$

$$B_5 = \{b, c, m\} \quad B_6 = \{b, c, j, m\}$$

$$B_7 = \{b, c, j\} \quad B_8 = \{b, c\}$$

Matching Transactions to Candidate Sets

Algorithm 6.1 Frequent itemset generation of the *Apriori* algorithm.

```
1:  $k = 1$ .
2:  $F_k = \{ i \mid i \in I \wedge \sigma(\{i\}) \geq N \times \text{minsup} \}$ .    {Find all frequent 1-itemsets}
3: repeat
4:    $k = k + 1$ .
5:    $C_k = \text{apriori-gen}(F_{k-1})$ .    {Generate candidate itemsets}
6:   for each transaction  $t \in T$  do
7:      $C_t = \text{subset}(C_k, t)$ .    {Identify all candidates that belong to  $t$ }
8:     for each candidate itemset  $c \in C_t$  do
9:        $\sigma(c) = \sigma(c) + 1$ .    {Increment support count}
10:    end for
11:  end for
12:   $F_k = \{ c \mid c \in C_k \wedge \sigma(c) \geq N \times \text{minsup} \}$ .    {Extract the frequent  $k$ -itemsets}
13: until  $F_k = \emptyset$ 
14:  $\text{Result} = \bigcup F_k$ .
```

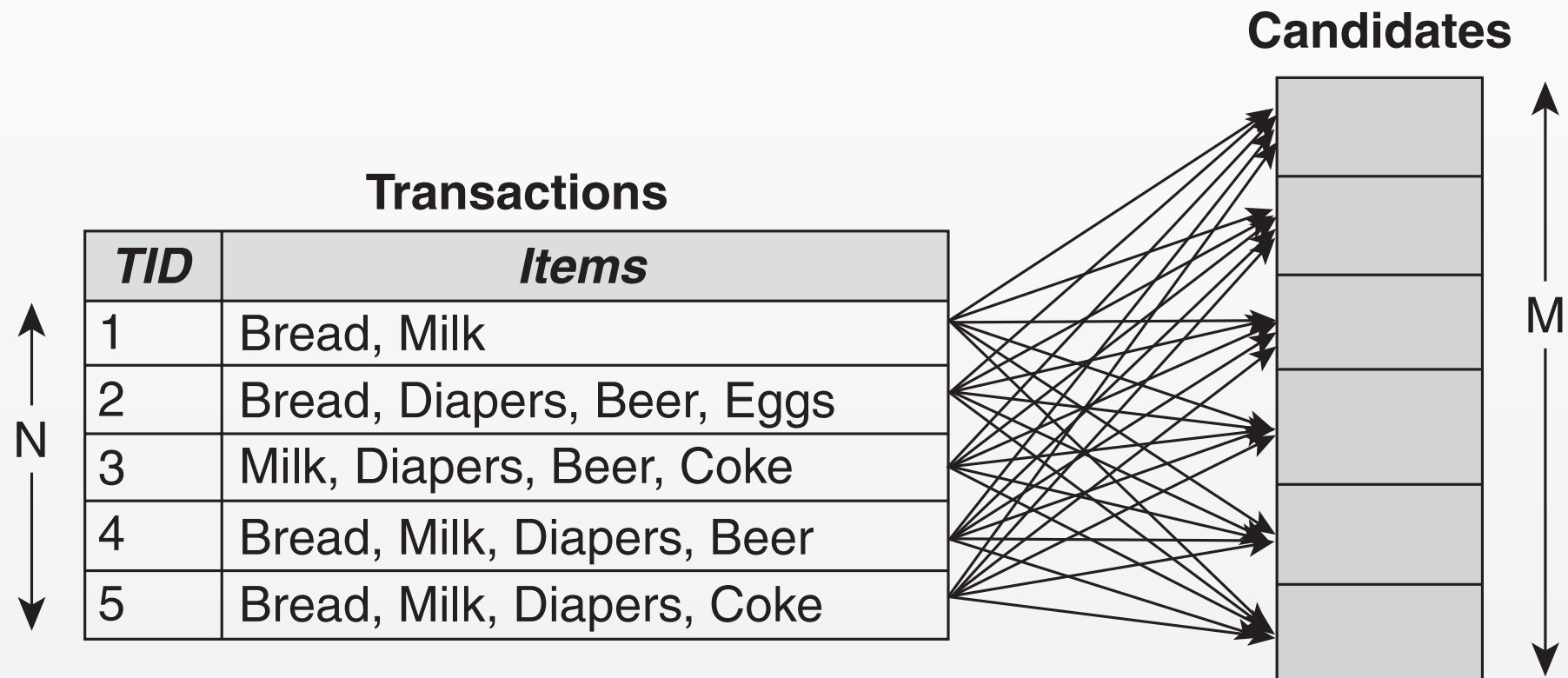
A-priori Algorithm

Algorithm 6.1 Frequent itemset generation of the *Apriori* algorithm.

```
1:  $k = 1$ .
2:  $F_k = \{ i \mid i \in I \wedge \sigma(\{i\}) \geq N \times \text{minsup} \}$ .    {Find all frequent 1-itemsets}
3: repeat                                     Interpret minsup as a fraction here
4:    $k = k + 1$ .
5:    $C_k = \text{apriori-gen}(F_{k-1})$ .    {Generate candidate itemsets}
6:   for each transaction  $t \in T$  do
7:      $C_t = \text{subset}(C_k, t)$ .    {Identify all candidates that belong to  $t$ }
8:     for each candidate itemset  $c \in C_t$  do
9:        $\sigma(c) = \sigma(c) + 1$ .    {Increment support count}
10:    end for
11:  end for
12:   $F_k = \{ c \mid c \in C_k \wedge \sigma(c) \geq N \times \text{minsup} \}$ .    {Extract the frequent  $k$ -itemsets}
13: until  $F_k = \emptyset$ 
14: Result =  $\bigcup F_k$ .
```

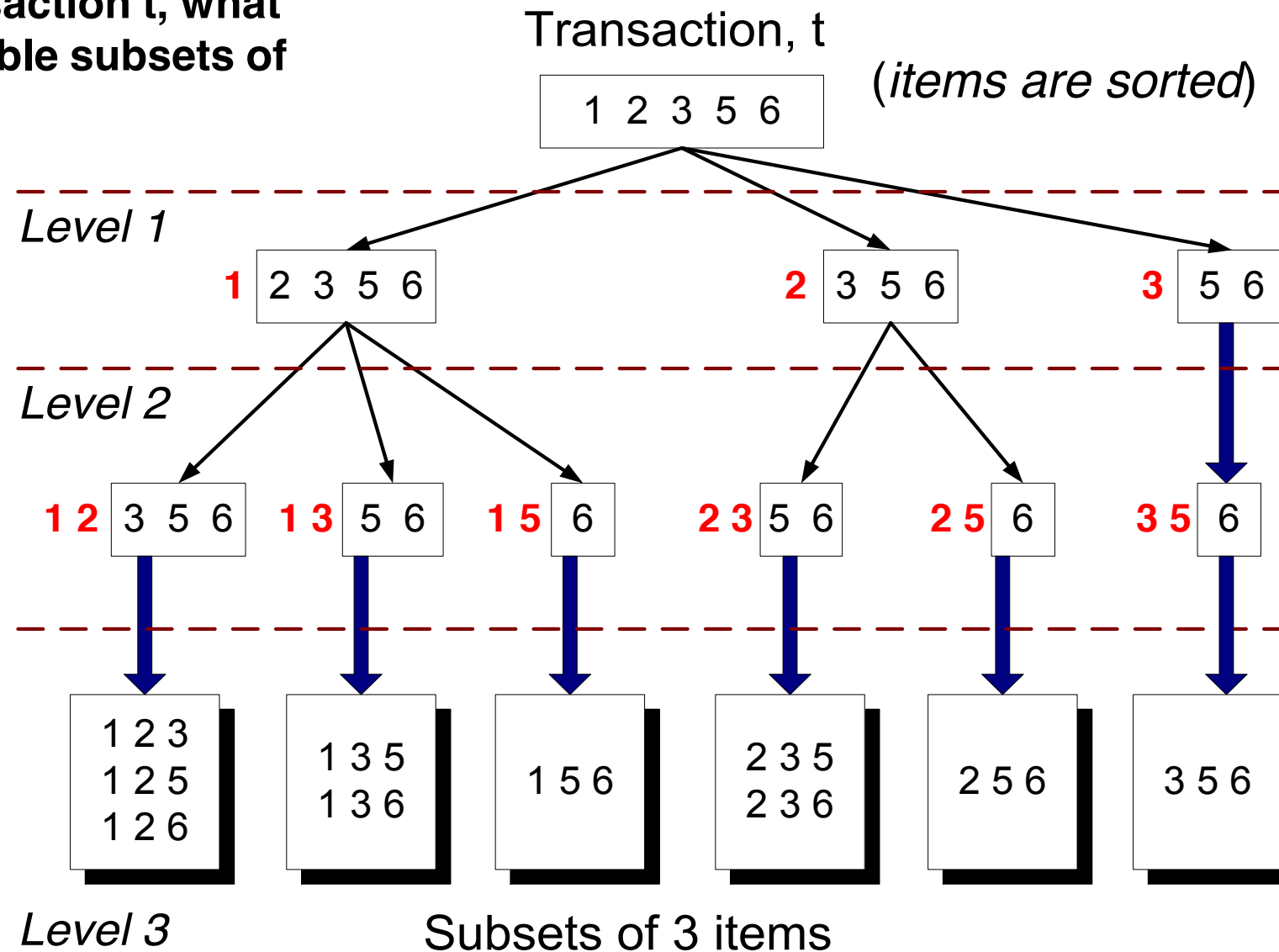
Problem: Naive Matching is Expensive

Cost: $O(N M)$, where N is number of baskets and M is number of candidates

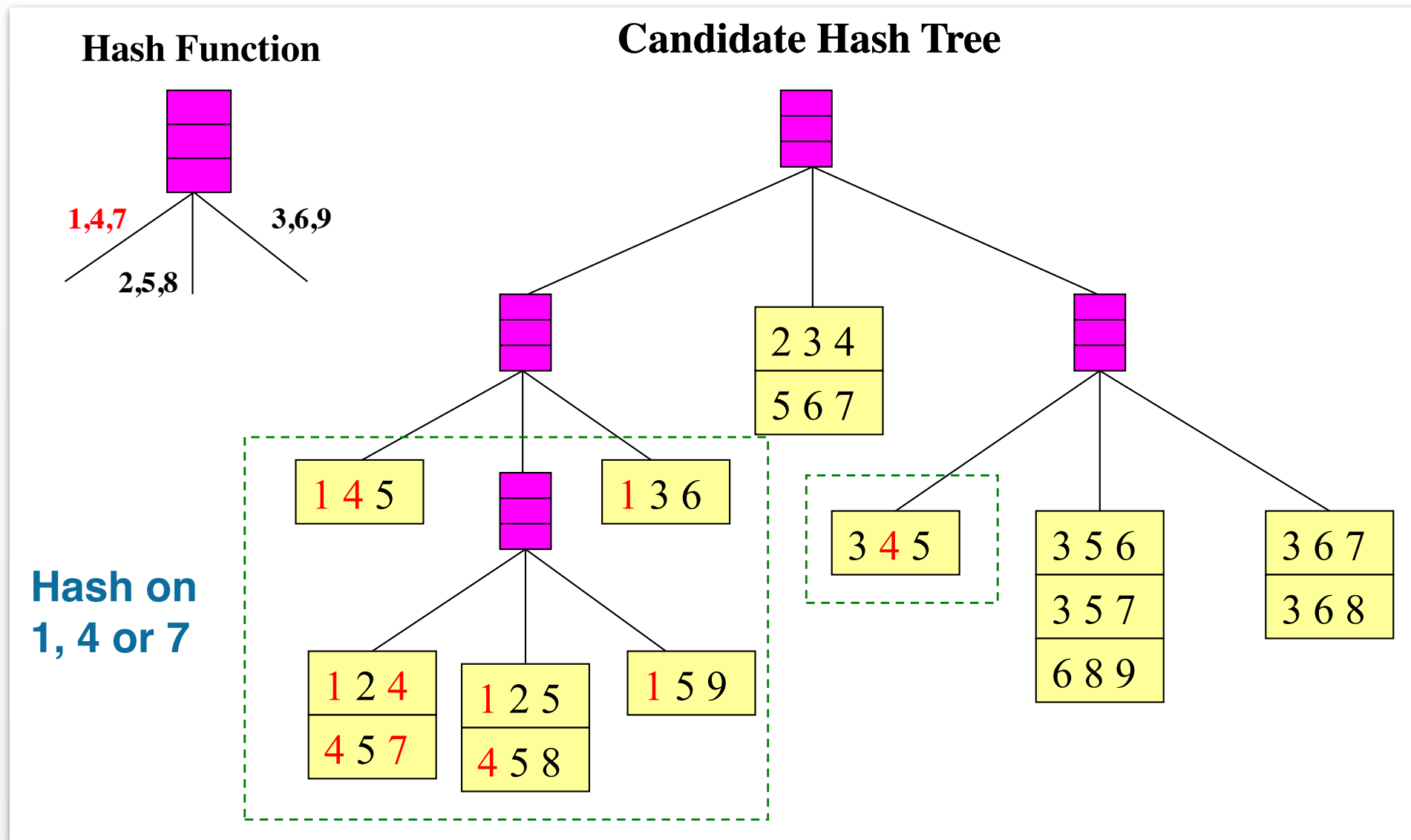


Strategy 1: Enumerating Transaction Subsets

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

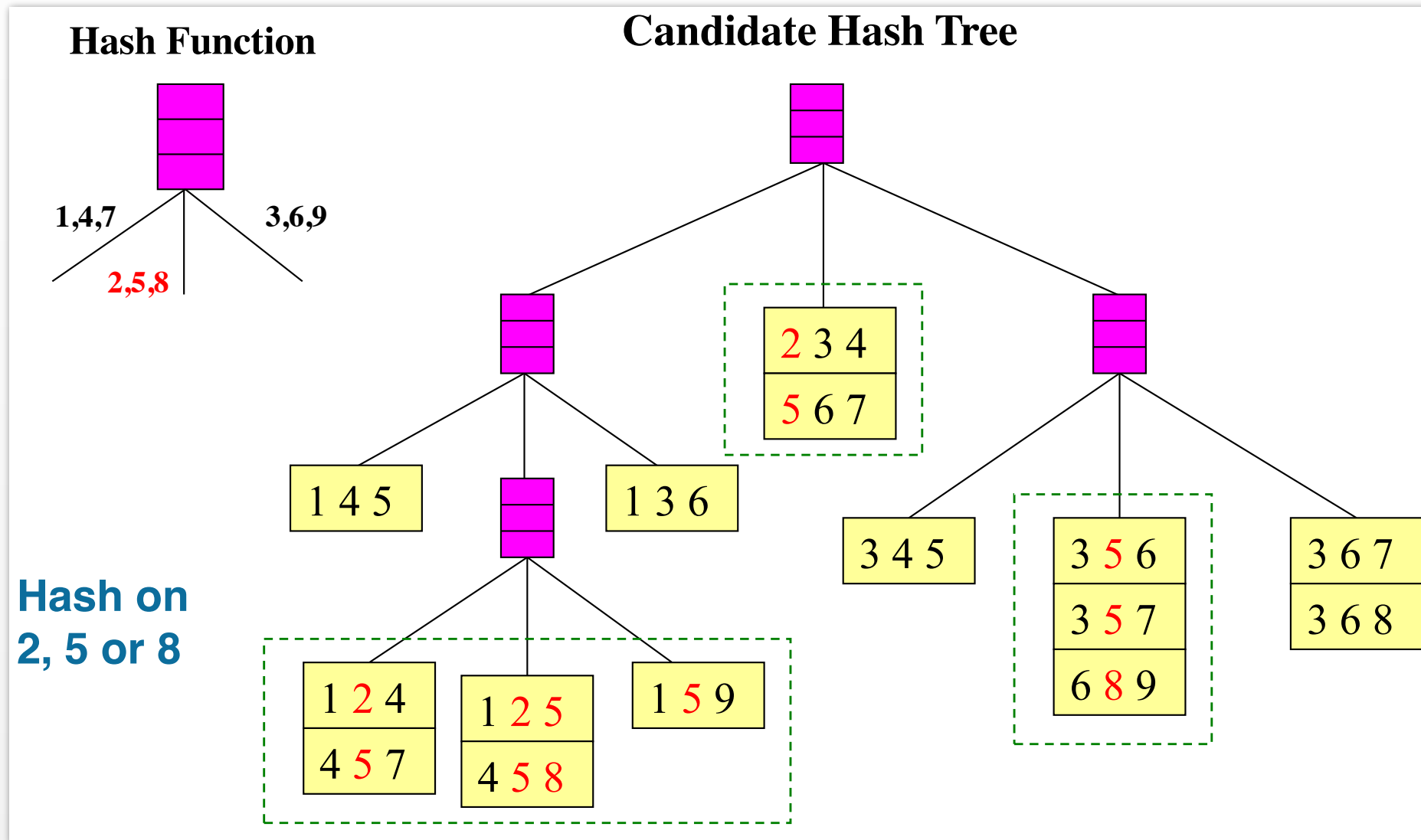


Hash Tree for Itemsets



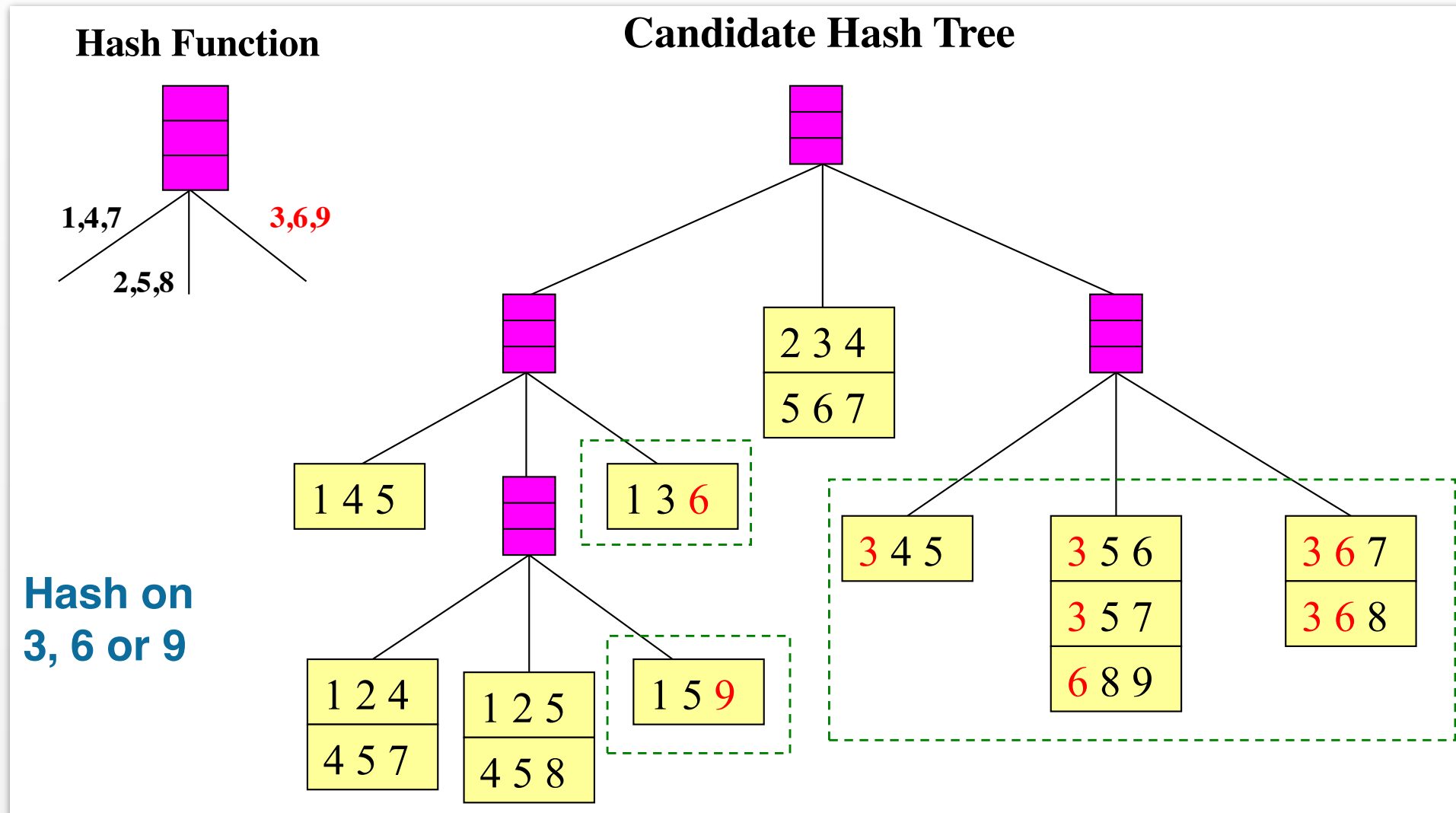
15 candidate 3-itemsets, distributed across 9 leaf nodes

Hash Tree for Itemsets



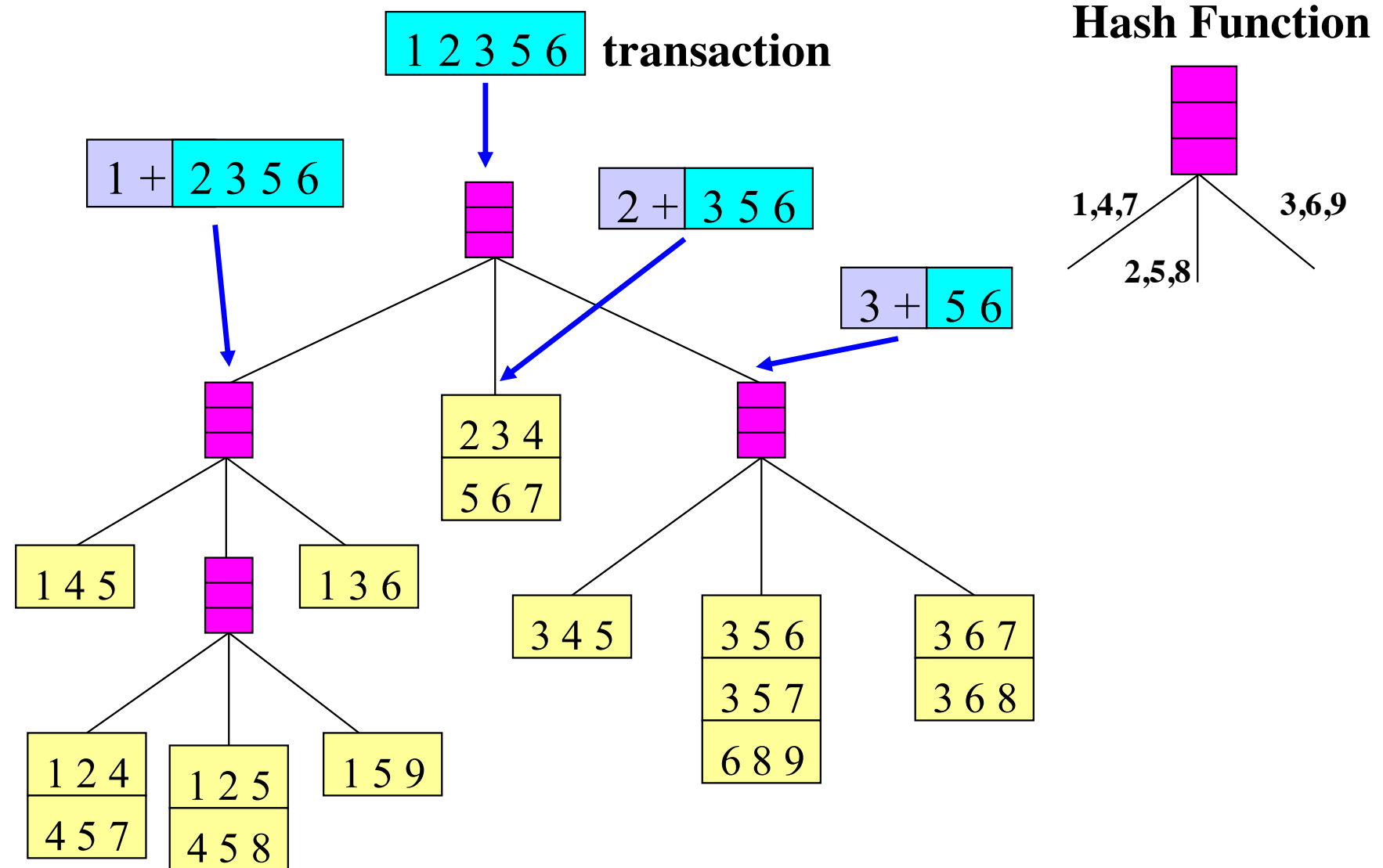
15 candidate 3-itemsets, distributed across 9 leaf nodes

Strategy 2: Hashing Itemsets

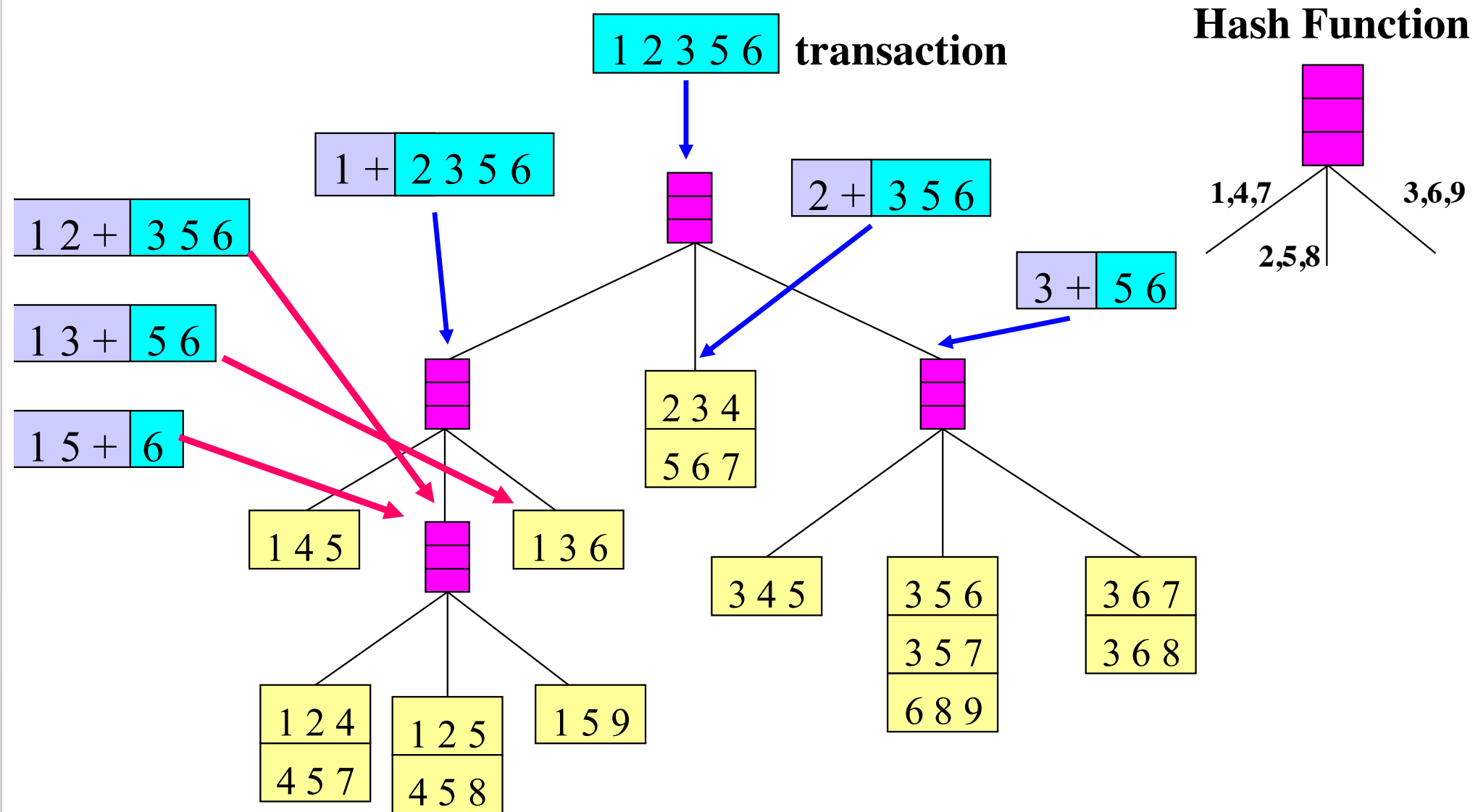


15 candidate 3-itemsets, distributed across 9 leaf nodes

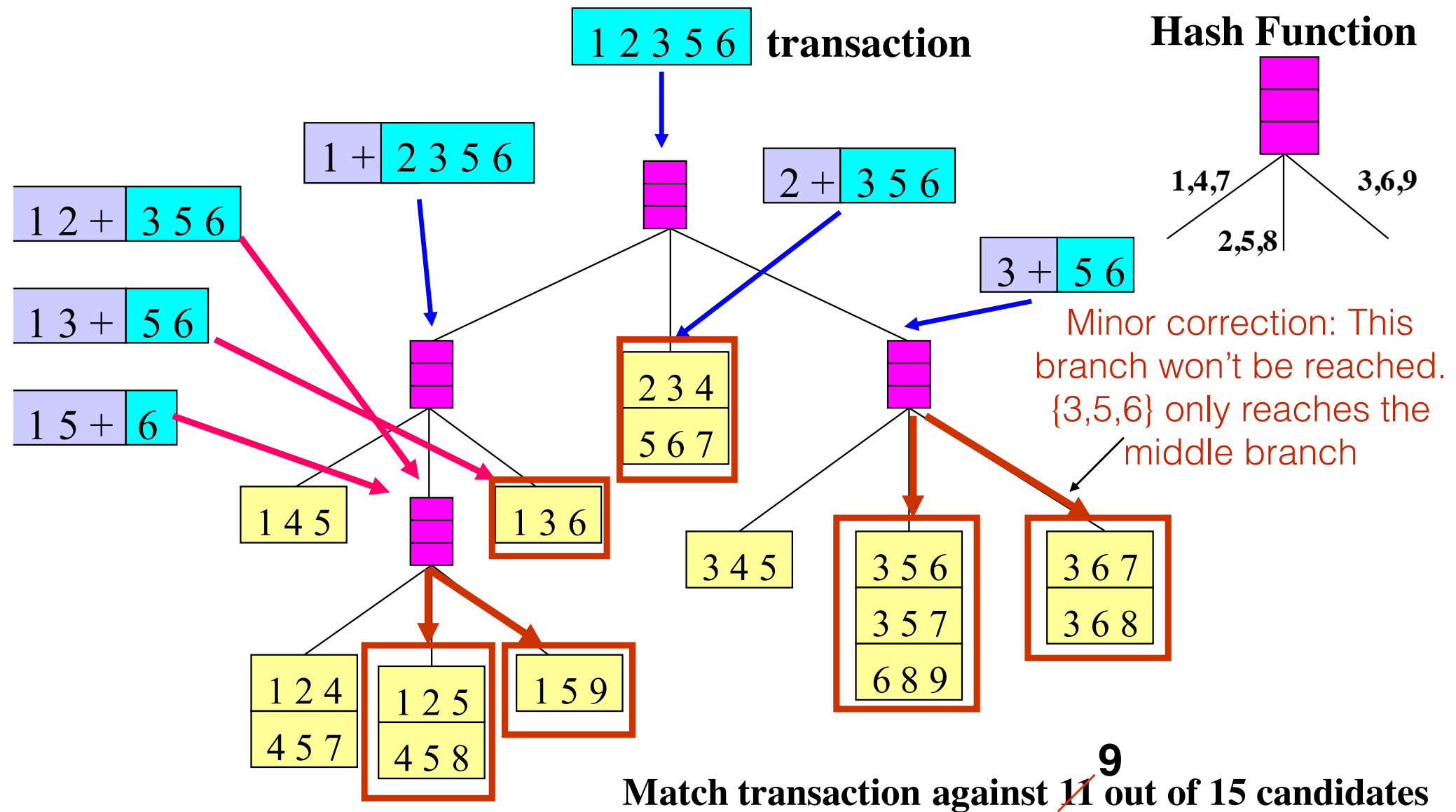
Strategy 2: Hash Tree for Candidates



Strategy 2: Hash Tree for Candidates



Strategy 2: Hash Tree for Candidates



A-priori Algorithm

Algorithm 6.1 Frequent itemset generation of the *Apriori* algorithm.

```
1:  $k = 1$ .
2:  $F_k = \{ i \mid i \in I \wedge \sigma(\{i\}) \geq N \times \text{minsup} \}$ .    {Find all frequent 1-itemsets}
3: repeat                                     Interpret minsup as a fraction here
4:    $k = k + 1$ .
5:    $C_k = \text{apriori-gen}(F_{k-1})$ .    {Generate candidate itemsets}
6:   for each transaction  $t \in T$  do
7:      $C_t = \text{subset}(C_k, t)$ .    {Identify all candidates that belong to  $t$ }
8:     for each candidate itemset  $c \in C_t$  do
9:        $\sigma(c) = \sigma(c) + 1$ .    {Increment support count}
10:    end for
11:  end for
12:   $F_k = \{ c \mid c \in C_k \wedge \sigma(c) \geq N \times \text{minsup} \}$ .    {Extract the frequent  $k$ -itemsets}
13: until  $F_k = \emptyset$ 
14: Result =  $\bigcup F_k$ .
```

Rule Generation

- Items of each frequent itemset Y can be partitioned into the consequent and the antecedent to give a rule. For an $X \subset Y$
$$X \rightarrow Y - X$$
- $Y = \{a, b, c\}$ could give the six rules $\{a, b\} \rightarrow \{c\}$, $\{a, c\} \rightarrow \{b\}$, $\{b, c\} \rightarrow \{a\}$, $\{a\} \rightarrow \{b, c\}$, $\{b\} \rightarrow \{a, c\}$, $\{c\} \rightarrow \{a, b\}$.
- A frequent k -itemset can potentially give to $2^k - 2$ rules.
- Not all rules are confident

$$C(X \rightarrow Y - X) = \sigma(Y) / \sigma(X) < \text{minconf}$$

X is also frequent by anti-monotonicity. However, the rule might not meet the minimum confidence threshold.

- How to find confident association rule without enumerating them all?

Pick a subset of the k items as a consequent. The remaining items become the antecedent. Remove $Y \rightarrow \emptyset$ and $\emptyset \rightarrow Y$.

Rule Generation

Rule Pruning

- **Theorem 6.2.** *If a rule $X \longrightarrow Y - X$ does not satisfy the confidence threshold, then any rule $X' \longrightarrow Y - X'$, where X' is a subset of X , must not satisfy the confidence threshold as well.*

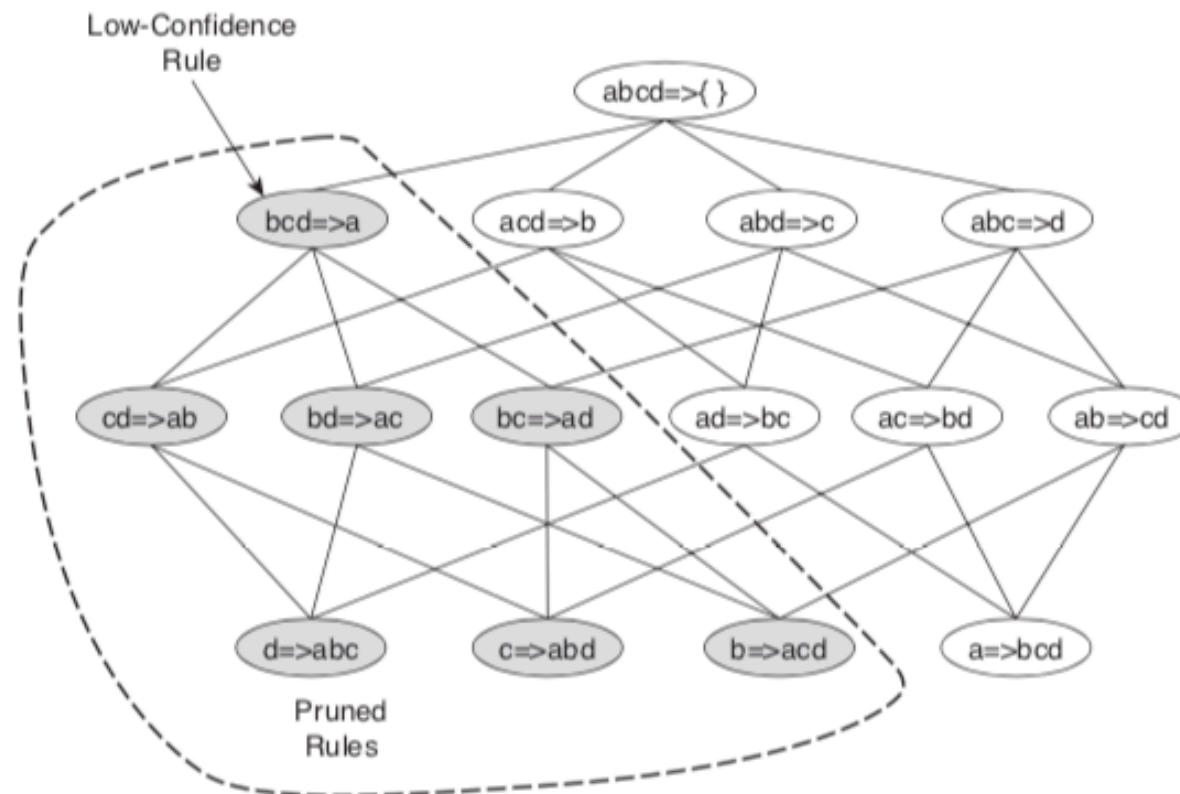


Figure 6.15. Pruning of association rules using the confidence measure.

Rule Generation

Algorithm 6.2 Rule generation of the *Apriori* algorithm.

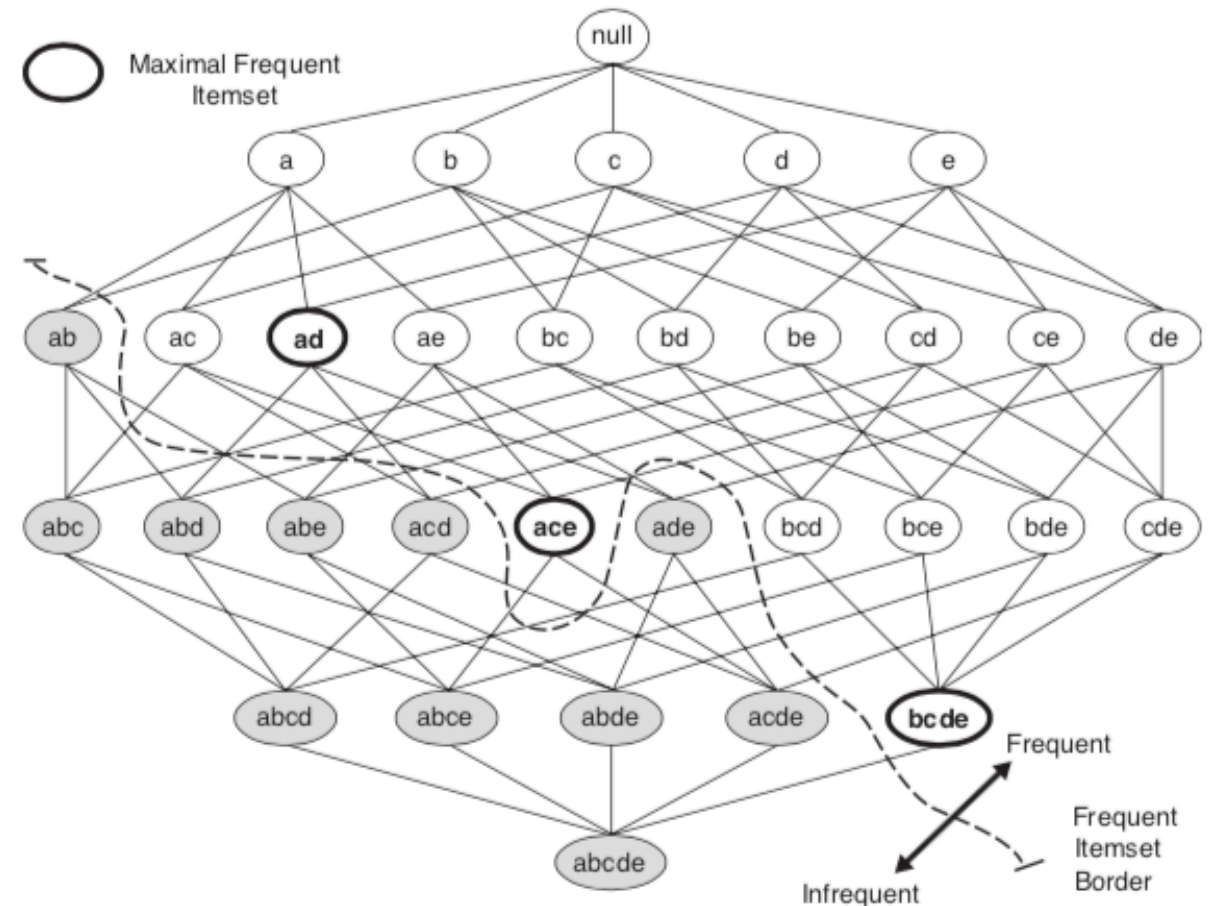
```
1: for each frequent  $k$ -itemset  $f_k, k \geq 2$  do
2:    $H_1 = \{i \mid i \in f_k\}$       {1-item consequents of the rule.}
3:   call ap-genrules( $f_k, H_1.$ )
4: end for
```

Algorithm 6.3 Procedure ap-genrules(f_k, H_m).

```
1:  $k = |f_k|$     {size of frequent itemset.}
2:  $m = |H_m|$     {size of rule consequent.}
3: if  $k > m + 1$  then
4:    $H_{m+1} = \text{apriori-gen}(H_m).$ 
5:   for each  $h_{m+1} \in H_{m+1}$  do
6:      $\text{conf} = \sigma(f_k) / \sigma(f_k - h_{m+1}).$ 
7:     if  $\text{conf} \geq \text{minconf}$  then
8:       output the rule  $(f_k - h_{m+1}) \longrightarrow h_{m+1}.$ 
9:     else
10:      delete  $h_{m+1}$  from  $H_{m+1}.$ 
11:    end if
12:  end for
13:  call ap-genrules( $f_k, H_{m+1}.$ )
14: end if
```

Compacting the Output

- To number of frequent item sets can be exponential in the number of items.
- Might be useful to work with compact representations
- **Maximal frequent itemsets:**
No immediate superset is frequent
 - Gives more pruning

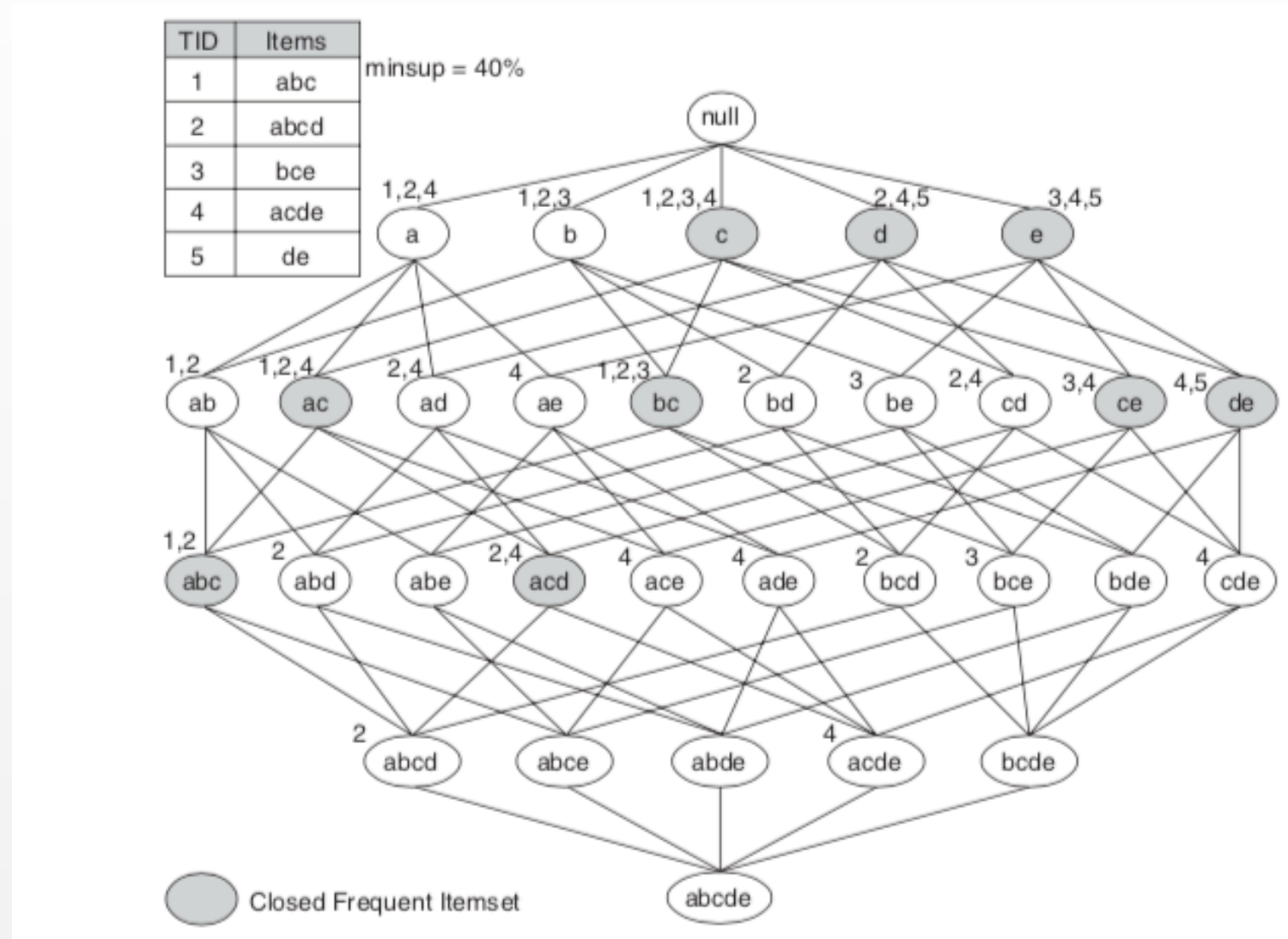


Compacting the Output

Closed frequent itemsets:

- No immediate superset has same count
- Stores not only frequent information, but exact counts
- The counts of non-closed frequent items can be obtained as the maximum of its closed frequent superset
- Redundant association rules are not generated if using closed frequent itemsets.

$\{b\} \rightarrow \{a\}$ and $\{b,c\} \rightarrow \{a\}$ will have the same support and confidence because $\{b\}$ is not closed, but $\{b,c\}$ is



Example: Maximal vs Closed

$$B_1 = \{m, c, b\} \quad B_2 = \{m, p, j\}$$

$$B_3 = \{m, b\} \quad B_4 = \{c, j\}$$

$$B_5 = \{m, c, b\} \quad B_6 = \{m, c, b, j\}$$

$$B_7 = \{c, b, j\} \quad B_8 = \{b, c\}$$

Frequent itemsets:

{m}:5, {c}:6, {b}:6, {j}:4,	Closed
{m,c}:3, {m,b}:4, {c,b}:5, {c,j}:3,	Maximal
{m,c,b}:3	

Example: Maximal vs Closed

