# Mining Frequent Itemsets with A-Priori

# Market Basket Analysis

Baskets of items

| TID | Items |
|---|---|
| 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}**



Which items are frequently purchased together by my customers?

Shopping Baskets

Customer 1: milk bread cereal
Customer 2: milk bread sugar eggs
Customer 3: milk bread butter
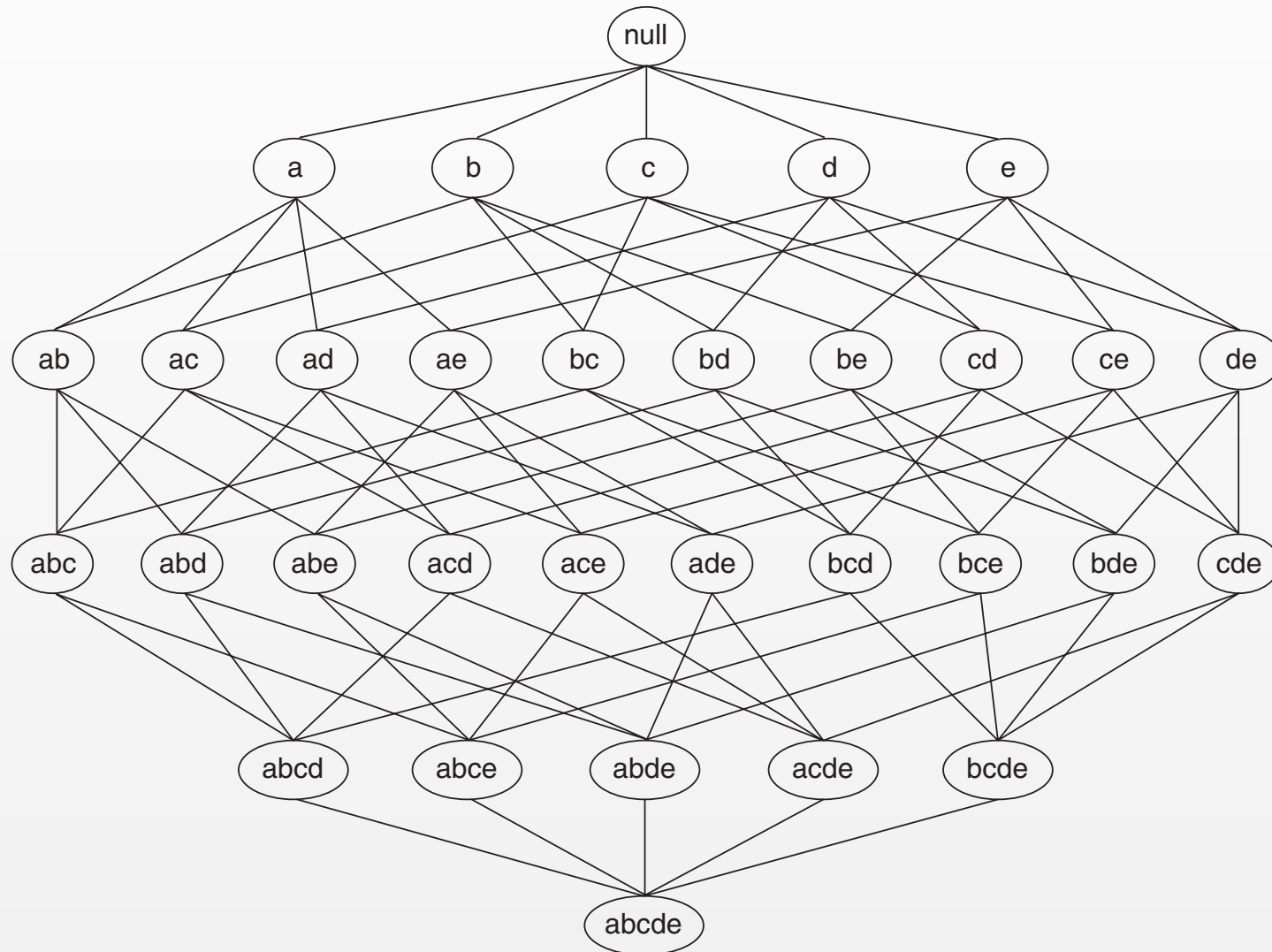Customer n: sugar eggs

Market Analyst

# Finding Frequent Item Sets

e the set of all items

$M = |I|$, how
any possible
:emsets 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

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

$I$ is the set of all items
$2^I$ denotes the power set of $I$

Support follows the anti-monotone property

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

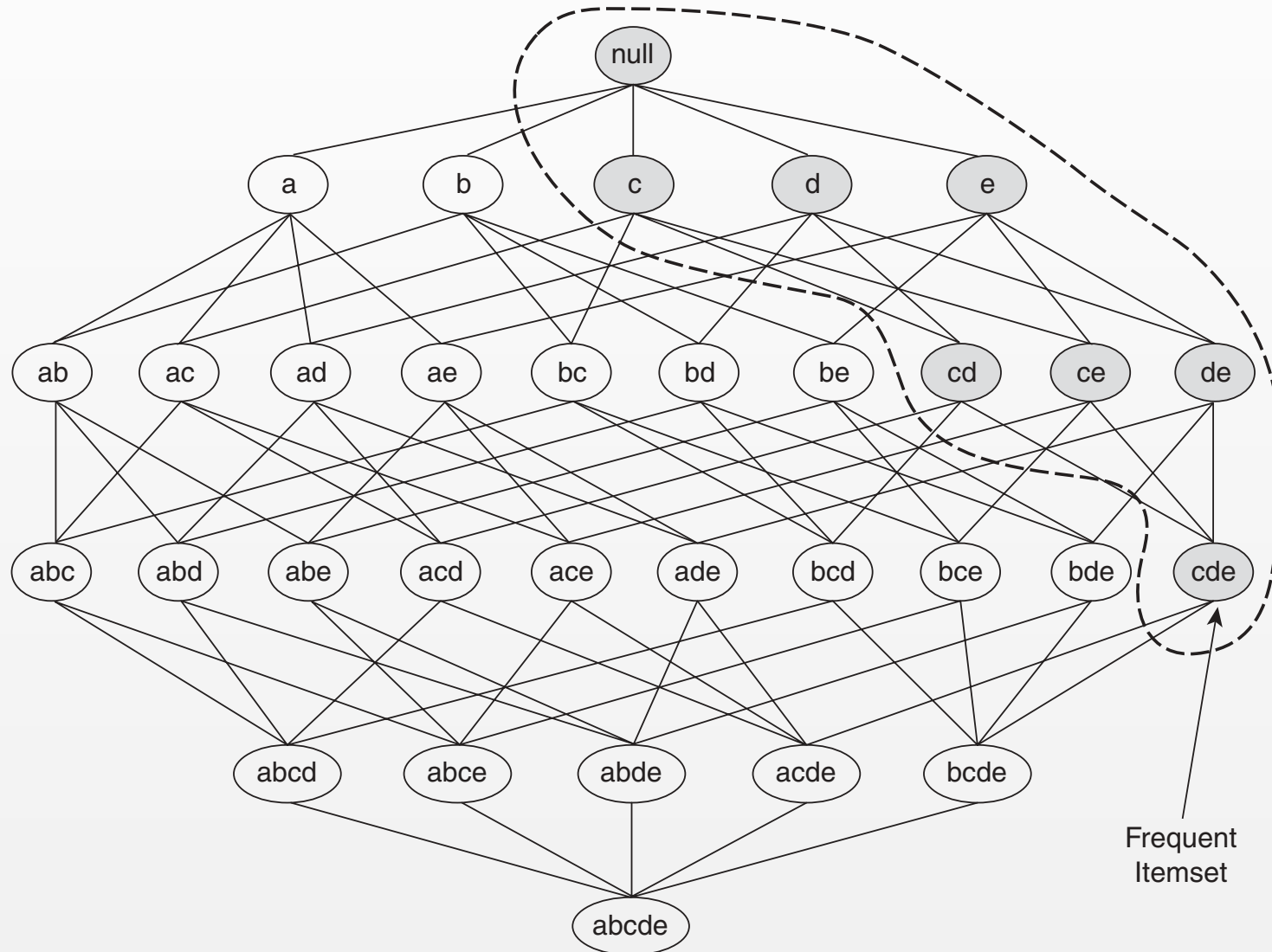$\sigma : 2^I \to \mathbf{N}$
$N = \{0, 1, \ldots, \infty$
set of natural nur
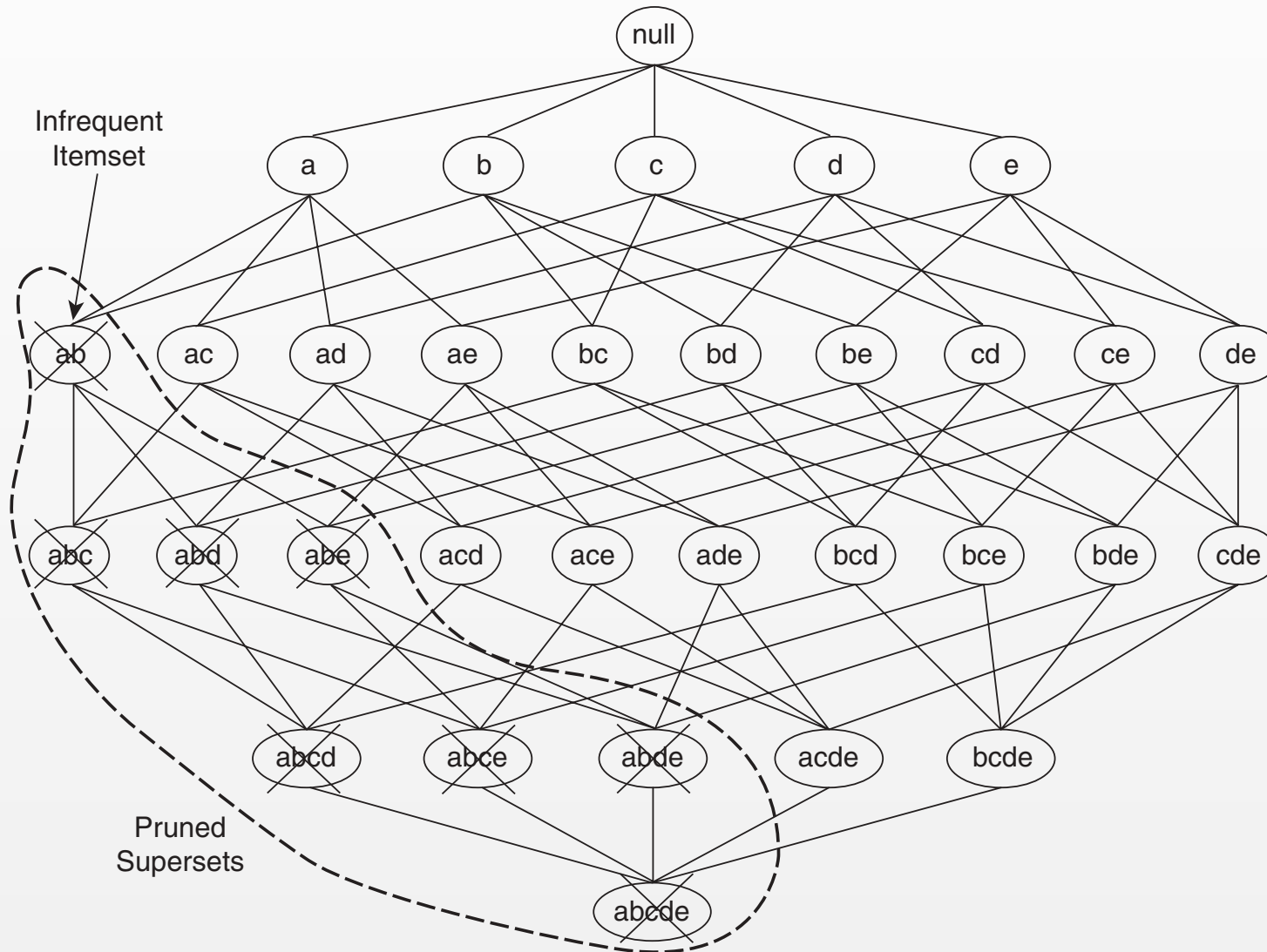$T$: the set of all tra

# Intuition: A-priori Principle

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



Frequent
Itemset

# Intuition: A-priori Principle

**ollary:** If a set is **not** frequent, then its supersets are **also** not freq



*adapted from*: Tan, Steinbach & Kumar, "Introduction to Data Mining", http://www-users.cs.umn.edu/~kumar/dmbook/ch6.pdf

# A-priori Algorithm

1. Find all frequent itemsets of size $1$
   (*only have to check $M = |I|$ possible sets*)

2. For $k = 1, 2, \ldots . M$

   - 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: $F_k = \{\ i \mid i \in I \wedge \sigma(\{i\}) \geq N \times minsup\}$.     {Find all frequent 1-itemsets}
3: **repeat**       Interpret minsup as a fraction here
4:     $k = k + 1$.
5:     $C_k = $ apriori-gen$(F_{k-1})$.     {Generate candidate itemsets}
6:     **for** each transaction $t \in T$ **do**
7:         $C_t = $ 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:     $F_k = \{\ c \mid c \in C_k \wedge \sigma(c) \geq N \times minsup\}$.     {Extract the frequent $k$-itemsets}
13: **until** $F_k = \emptyset$
14: Result $= \bigcup F_k$.

---

---

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

---

$k = 1$.

$F_k = \{\, i \mid i \in I \wedge \sigma(\{i\}) \geq N \times minsup \,\}$.     {Find all frequent 1-itemsets}

**repeat**

    $k = k + 1$.

    $C_k = \text{apriori-gen}(F_{k-1})$.     {Generate candidate itemsets}

    **for** each transaction $t \in T$ **do**

        $C_t = \text{subset}(C_k, t)$.     {Identify all candidates that belong to $t$}

        **for** each candidate itemset $c \in C_t$ **do**

            $\sigma(c) = \sigma(c) + 1$.     {Increment support count}

        **end for**

    **end for**

    $F_k = \{\, c \mid c \in C_k \wedge \sigma(c) \geq N \times minsup \,\}$.     {Extract the frequent $k$-itemsets}

**until** $F_k = \emptyset$

$\text{Result} = \bigcup F_k$.

---

# 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 $C_k$

*stions*

*ow many k-itemsets are there?*

*ow to reduce number of candidates using the computations eady performed?*

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

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

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

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

# Generating Candidates $C_k$

*uplicates*

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 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 siz could be generate

$$\binom{k}{k-2}$$ times

# Generating Candidates $C_k$

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. $a_1 < a_2 < \ldots < a_k$ in $A = \{a_1, a_2, \ldots, a_k\}$

ID based
have co
advan
compar
is che
compa

# Generating Candidates $C_k$

*ion to duplicates*

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

ombine $\{a,c,e\}$ with $\{f\}$ to give candidate $\{a,c,e,f\}$.

o not combine $\{a,c,e\}$ with $\{d\}$.

o duplicates: Each candidate has only one way of being enerated. $\{a,c,e,f\}$ can only be generated by combining $\{a,c,e$ nd $\{f\}$.

ompleteness: If $\{a,c,e,f\}$ is indeed frequent, $\{a,c,e\}$ and $\{f\}$ hav o be present in $F_3$ and $F_1$. And they would get combined to enerate $\{a,c,e,f\}$.

ubsets of generated candidates might still be infrequent

# Generating Candidates $C_k$

*tion to duplicates*

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

ombine $\{a, c, e\}$ with $\{a, c, f\}$ to give candidate $\{a, c, e, f\}$.
o not combine $\{a, c, e\}$ with $\{a, b, e\}$.

o duplicates: Each candidate has only one way of being
enerated. $\{a, c, e, f\}$ can only be generated by combining $\{a, c, e$
nd $\{a, c, f\}$.

ompleteness: If$\{a, c, e, f\}$ is indeed frequent, $\{a, c, e\}$ and $\{a, c, f\}$
ave to be present in $F_3$. And they would get combined to
enerate $\{a, c, e, f\}$.
ubsets of generated candidates might still be infrequent

# Generating Candidates $C_k$

*to efficiently find itemsets that could be combined*

**set Ordering:** Use the ordering on items to define an ordering sets

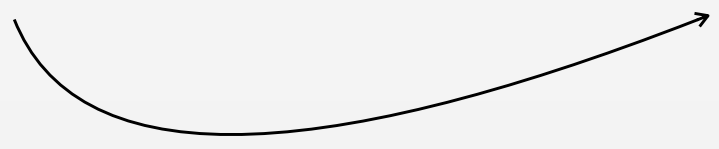sume that the items in each itemset are presorted.

$< a_2 < \ldots < a_k$ in $A = \{a_1, a_2, \ldots, a_k\}$

$B$ if $a_i < b_i$, where $i$ is the index of the first item differing in $A$ an

ple, bread, coke, sauce$\} < \{$apple, bread, cookie, milk$\}$ or $\{4,7,21,50\} < \{4,7$

ements of $F_{k-1}$ sorted with the itemset ordering
$$\{a,b,c\}, \{a,b,e\}, \{a,b,g\}, \{a,c,d\}, \{a,c,g\} \ldots$$

$k-1$

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

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

# Generating Candidates $C_k$

*...ng candidates* with infrequent subsets

...or a candidate of size $k$, one only needs to check ...ubsets of size $k - 1$.
...numerate subsets of size $k - 1$ by removing one ...ement at a time from the candidate.
...earch for the subsets one after the other in $F_{k-1}$ until ...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.
...a subset wasn't found the candidate should be ...iscarded.

If all size $k - 1$
are frequent s...
subsets of sm...

Each candida...
will give $k$ sub...
size $k - 1$

For each ...
$O(k|F_{k-1}|)$
might be ne...
worst case,
naiv...
This can be
binary
$O(k \log(|...$
constructing...

# Generating Candidates $C_k$

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

*Pruning:* Remove all candidates
with infrequent subsets

# Example: Generating Candidates $C_k$

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

*quent itemsets of size 2:*
*,c}:5, {b,m}:4,{c,j}:3 {c,m}:3*

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

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

*elf-joining:*
*,c,m}, {c,j,m}*

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

*uning:*

~~{j,m}~~ since ~~{j,m}~~ not frequent

*quent items of size 3:*

*b,c,m}*

---

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

---

$k = 1$.
$F_k = \{\ i \mid i \in I \wedge \sigma(\{i\}) \geq N \times minsup\}$.     {Find all frequent 1-itemsets}
**repeat**
   $k = k + 1$.
   $C_k = \text{apriori-gen}(F_{k-1})$.     {Generate candidate itemsets}
   **for** each transaction $t \in T$ **do**
      $C_t = \text{subset}(C_k, t)$.     {Identify all candidates that belong to $t$}
      **for** each candidate itemset $c \in C_t$ **do**
         $\sigma(c) = \sigma(c) + 1$.     {Increment support count}
      **end for**
   **end for**
   $F_k = \{\ c \mid c \in C_k \wedge \sigma(c) \geq N \times minsup\}$.     {Extract the frequent $k$-itemsets}
**until** $F_k = \emptyset$
$\text{Result} = \bigcup F_k$.

---

# A-priori Algorithm

---

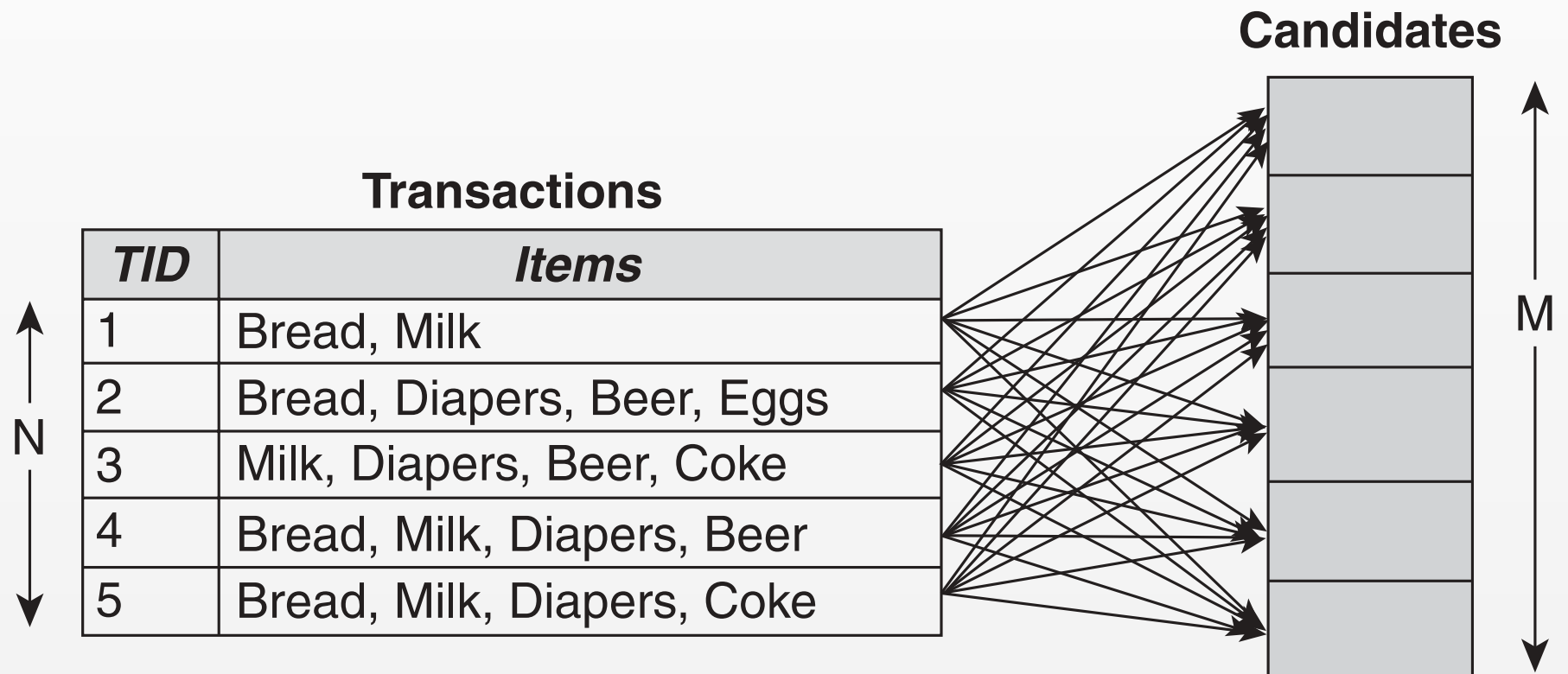**Algorithm 6.1** Frequent itemset generation of the *Apriori* algorithm.

1: $k = 1$.
2: $F_k = \{ i \mid i \in I \wedge \sigma(\{i\}) \geq N \times minsup\}$.     {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:     $F_k = \{ c \mid c \in C_k \wedge \sigma(c) \geq N \times minsup\}$.     {Extract the frequent $k$-itemsets}
13: **until** $F_k = \emptyset$
14: Result $= \bigcup F_k$.

---

# roblem: Naive Matching is Expensiv

O(N M), where N is number of baskets and M is number of cand

**Candidates**

**Transactions**

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

N

M

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

Transaction, t

| 1 2 3 5 6 |

(*items are sorted*)

*Level 1*

**1** | 2 3 5 6     **2** | 3 5 6     **3** | 5 6

*Level 2*

**1 2** | 3 5 6    **1 3** | 5 6    **1 5** | 6    **2 3** | 5 6    **2 5** | 6    **3 5** | 6

1 2 3
1 2 5
1 2 6

1 3 5
1 3 6

1 5 6

2 3 5
2 3 6

2 5 6

3 5 6

*Level 3*      Subsets of 3 items

# Hash Tree for Itemsets

**Hash Function**

**Candidate Hash Tree**

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

**Hash on 1, 4 or 7**

| 2 3 4 |
|-------|
| 5 6 7 |

| 1 4 5 |

| 1 3 6 |

| 3 4 5 |

| 3 5 6 |
|-------|
| 3 5 7 |
| 6 8 9 |

| 3 6 7 |
|-------|
| 3 6 8 |

| 1 2 4 |
|-------|
| 4 5 7 |

| 1 2 5 |
|-------|
| 4 5 8 |

| 1 5 9 |

15 candidate 3-itemsets, distributed across 9 leaf nodes

*adapted from*: Tan, Steinbach & Kumar, "Introduction to Data Mining", http://www-users.cs.umn.edu/~kumar/dmbook/ch6.pdf

# Hash Tree for Itemsets

**Hash Function**

**Candidate Hash Tree**

1,4,7          3,6,9

2,5,8

**Hash on 2, 5 or 8**

| 2 3 4 |
| 5 6 7 |

| 1 4 5 |

| 1 3 6 |

| 3 4 5 |

| 3 5 6 |
| 3 5 7 |
| 6 8 9 |

| 3 6 7 |
| 3 6 8 |

| 1 2 4 |
| 4 5 7 |

| 1 2 5 |
| 4 5 8 |

| 1 5 9 |

15 candidate 3-itemsets, distributed across 9 leaf nodes

*adapted from*: Tan, Steinbach & Kumar, "Introduction to Data Mining", http://www-users.cs.umn.edu/~kumar/dmbook/ch6.pdf

# Strategy 2: Hashing Itemsets



15 candidate 3-itemsets, distributed across 9 leaf nodes

*adapted from*: Tan, Steinbach & Kumar, "Introduction to Data Mining", http://www-users.cs.umn.edu/~kumar/dmbook/ch6.pdf
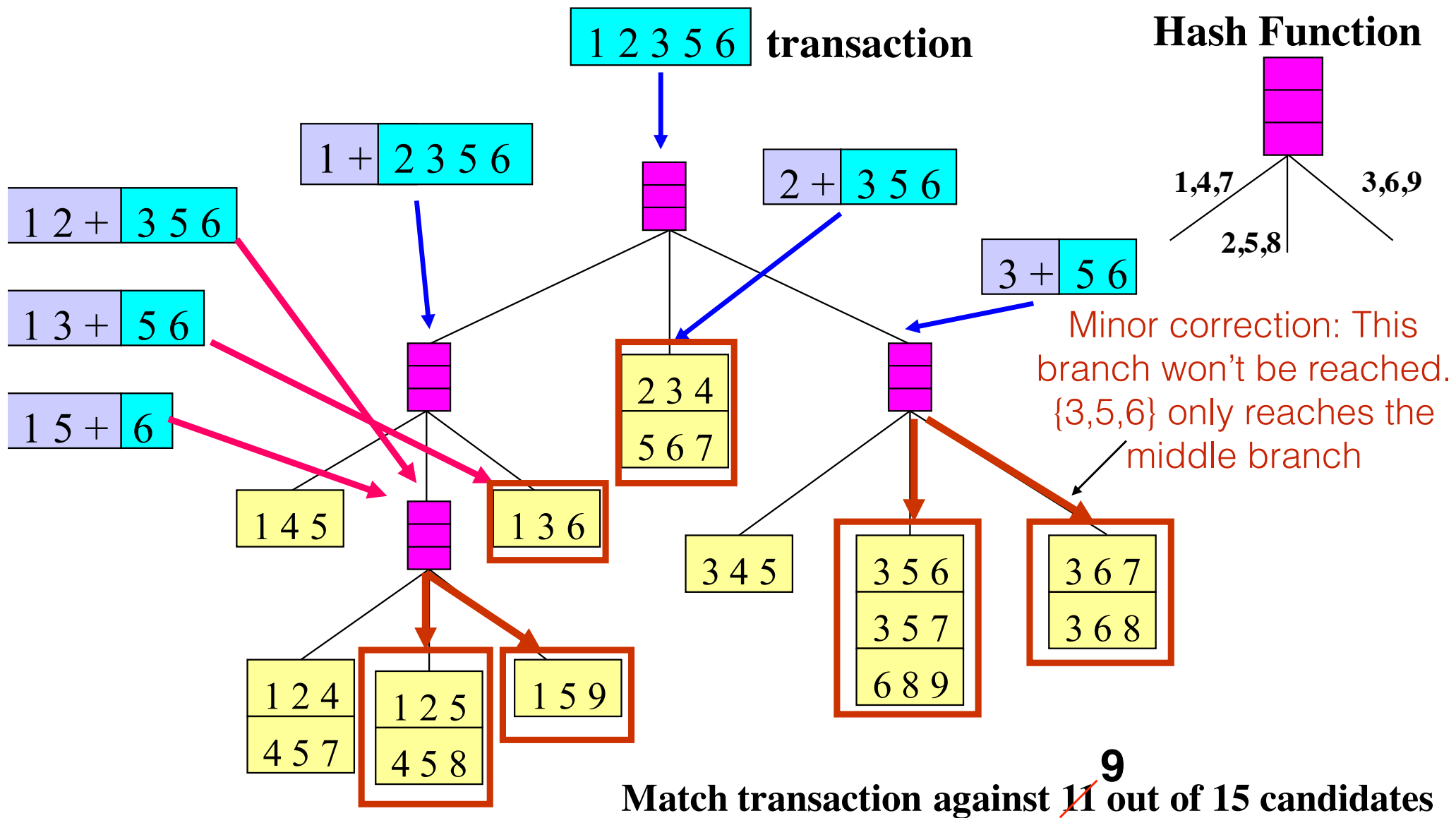
# Strategy 2: Hash Tree for Candidates

# Strategy 2: Hash Tree for Candidates

# Strategy 2: Hash Tree for Candidates



**1 2 3 5 6** transaction

**Hash Function**

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

1 + 2 3 5 6

2 + 3 5 6

3 + 5 6

1 2 + 3 5 6

1 3 + 5 6

1 5 + 6

2 3 4
5 6 7

1 4 5

1 3 6

3 4 5

3 5 6
3 5 7
6 8 9

3 6 7
3 6 8

1 2 4
4 5 7

1 2 5
4 5 8

1 5 9

Minor correction: This branch won't be reached. {3,5,6} only reaches the middle branch

**9**

**Match transaction against ~~11~~ out of 15 candidates**

*adapted from*: Tan, Steinbach & Kumar, "Introduction to Data Mining", http://www-users.cs.umn.edu/~kumar/dmbook/ch6.pdf

# A-priori Algorithm

---

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

1: $k = 1$.
2: $F_k = \{ i \mid i \in I \wedge \sigma(\{i\}) \geq N \times minsup \}$.    {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:    $F_k = \{ c \mid c \in C_k \wedge \sigma(c) \geq N \times minsup \}$.    {Extract the frequent $k$-itemsets}
13: **until** $F_k = \emptyset$
14: Result $= \bigcup F_k$.

---

# Rule Generation

…ems of each frequent itemset $Y$ can be partitioned into the …onsequent and the the antecedent to give a rule. For an $X \subset Y$

$$X \to Y - X$$

…$= \{a, b, c\}$ could give the six rule $\{a, b\} \to \{c\}, \{a, c\} \to \{b\},$
…$, c\} \to \{a\}, \{a\} \to \{b, c\}, \{b\} \to \{a, c\}, \{c\} \to \{a, b\}$.
…frequent k-itmeset can potentially give to $2^k - 2$ rules.
…ot all rules are confident

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

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

…ow to find confident association rule without enumerating them all?

Pick a
$k$ it
conse
rema
be
antece
$Y \to \varnothing$

# Rule Generation

**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.*
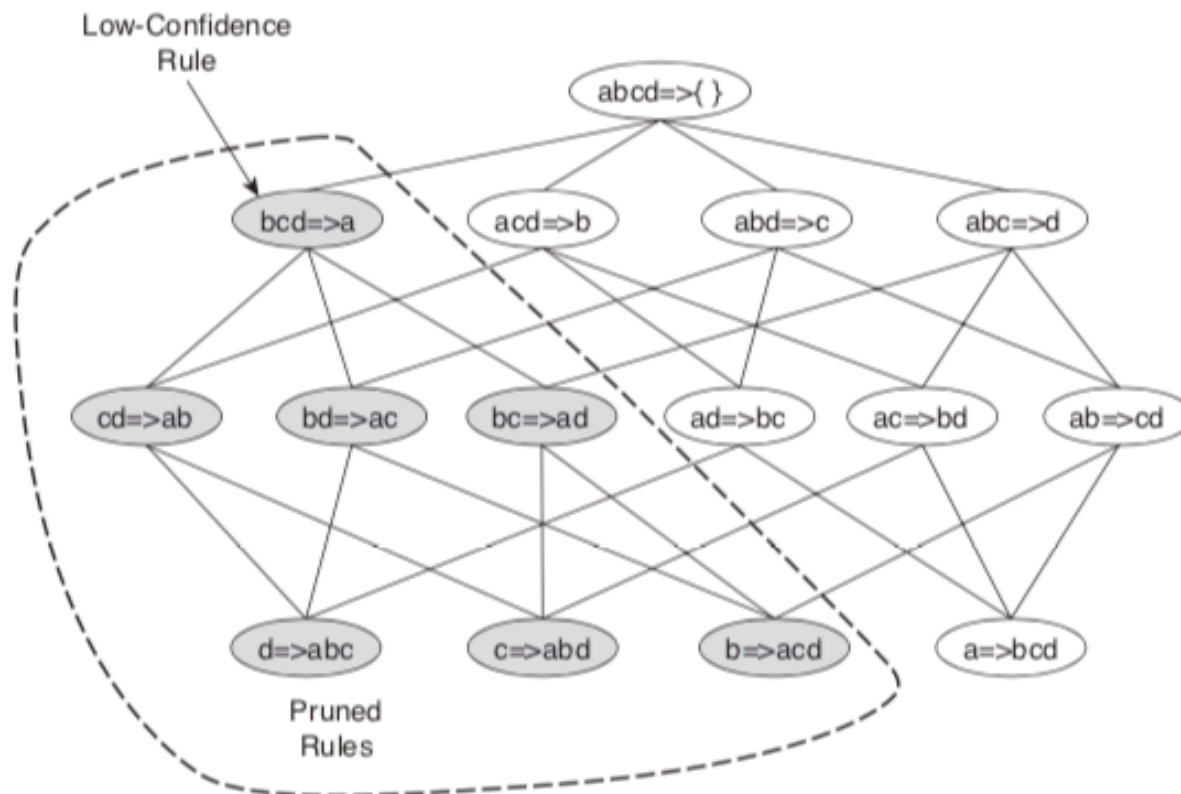


**Figure 6.15.** Pruning of association rules using the confidence measure.

# 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:         $conf = \sigma(f_k)/\sigma(f_k - h_{m+1}).$
7:         **if** $conf \geq minconf$ **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

umber of frequent item
can be exponential in
number of items.

ht be useful to work with
pact representations

imal frequent itemsets:
mmediate superset is
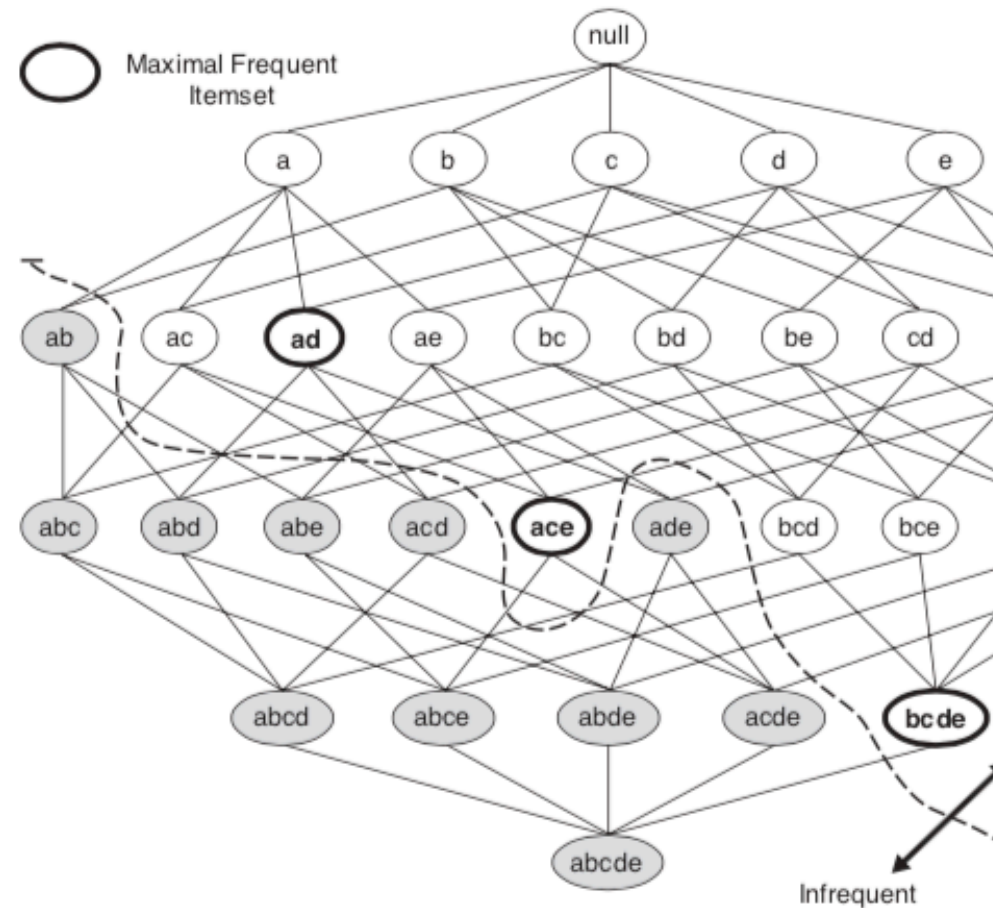uent

es more pruning



Maximal Frequent
Itemset

Infrequent

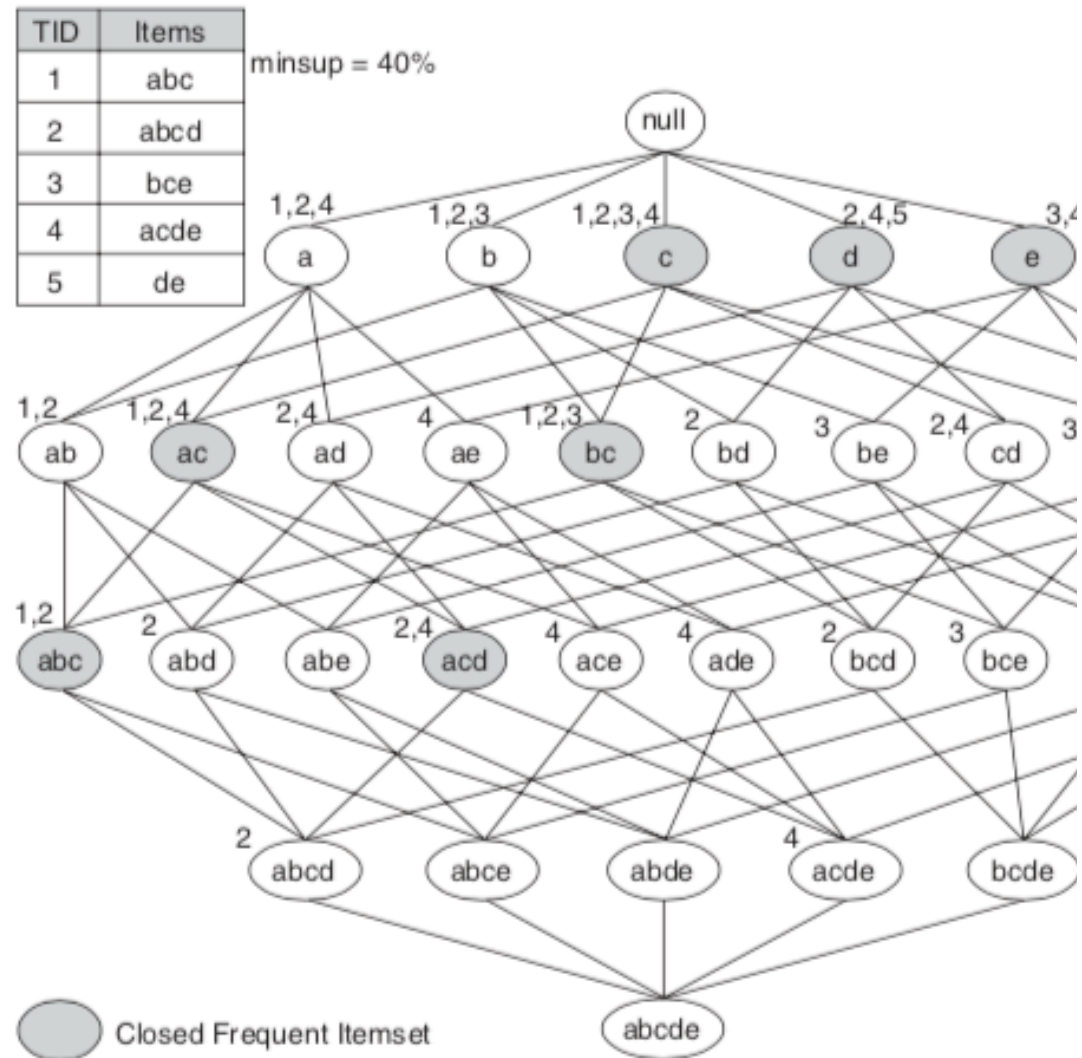**Figure 6.16.** Maximal frequent itemset.

# Compacting the Output

nediate superset has same

not only frequent
ation, but exact counts

ounts of non-closed frequent
can be obtained as the
um of its closed frequent
set

dant association rules are
nerated if using closed
nt itemsets.

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

| TID | Items |
|-----|-------|
| 1 | abc |
| 2 | abcd |
| 3 | bce |
| 4 | acde |
| 5 | de |

minsup = 40%

Closed Frequent Itemset

adapted from: J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

# Example: Maximal vs Closed

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

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

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

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

*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