

# FREQUENT PATTERN MINING—ESSENTIALS

---

*Hari Sundaram*

[hs1@illinois.edu](mailto:hs1@illinois.edu)

<http://sundaram.cs.illinois.edu>

adapted from slides by Jiawei Han and Kevin Chang

# BASIC CONCEPTS

---

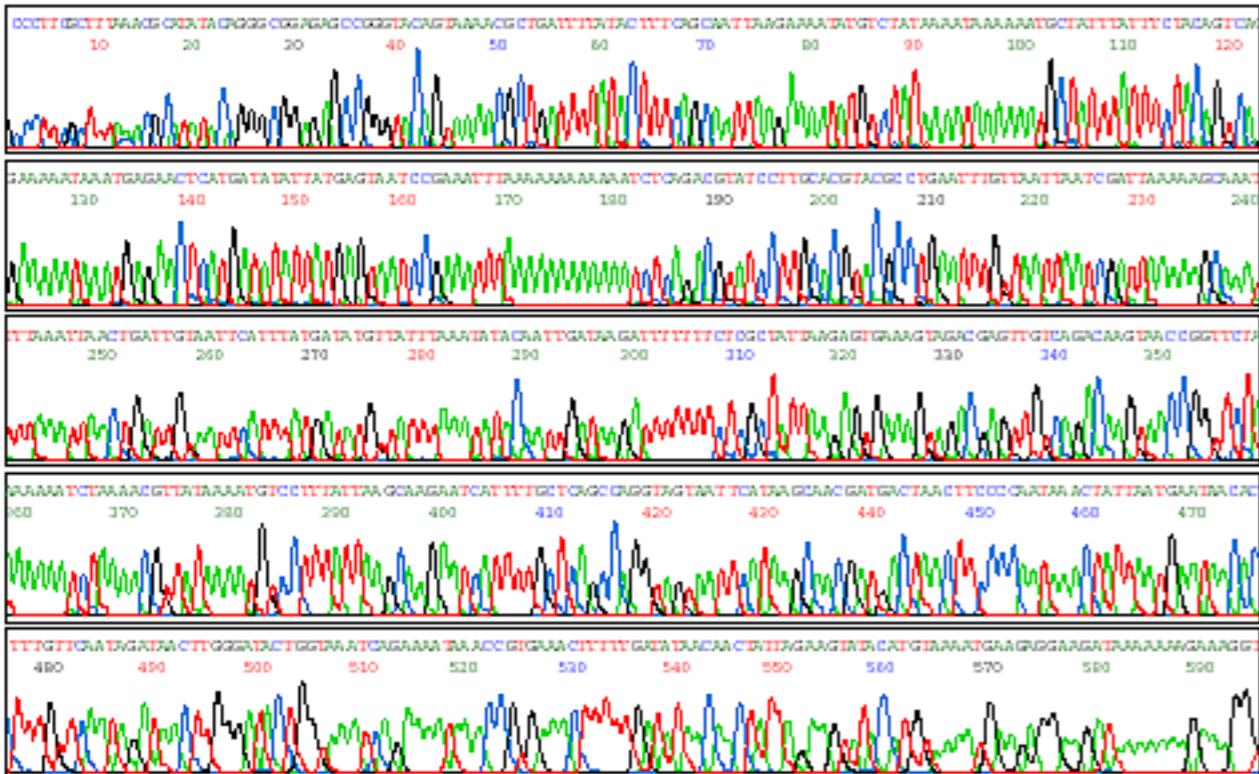
Methods

Pattern Evaluation

Summary



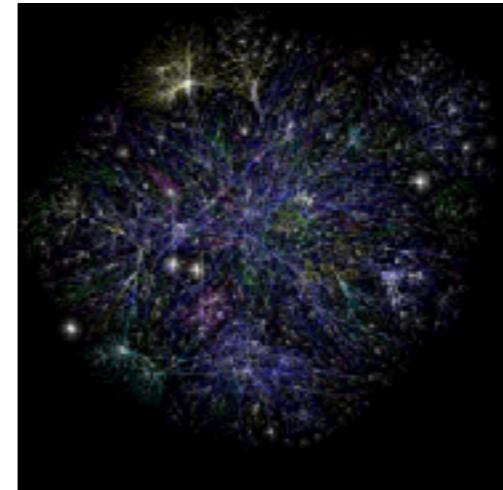
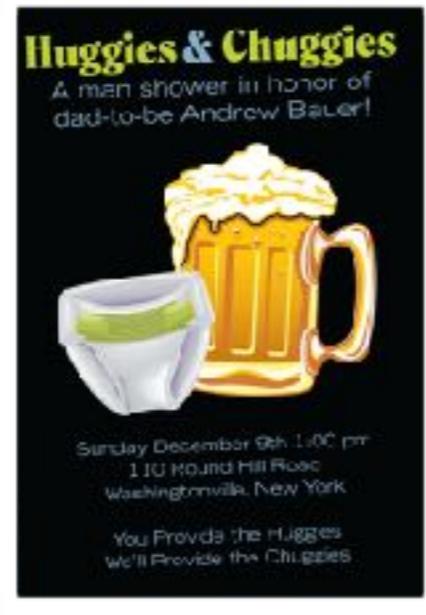
# frequent patterns



a pattern (a set of items, subsequences, substructures, etc.) that occurs frequently in a data set

Agrawal, Rakesh, Tomasz Imieliński, and Arun Swami. "Mining association rules between sets of items in large databases." ACM SIGMOD Record 22.2 (1993): 207-216.

What products are often purchased together?

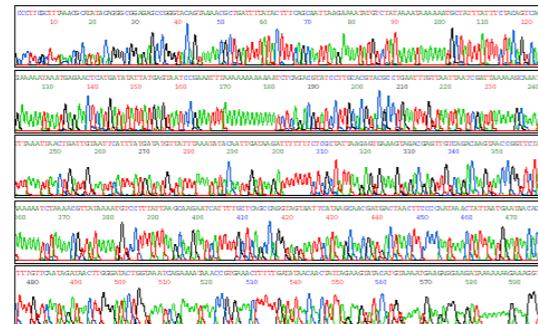


Can we automatically classify web documents?

lots of applications!

# motivation

What kinds of DNA are sensitive to this new drug?



What are the subsequent purchases after buying a PC?



why is it  
important?

# frequent mining importance

Classification

Semantic data compression: fascicles

Jagadish, H. V., Jason Madar, and Raymond T. Ng.  
"Semantic compression and pattern extraction with fascicles." VLDB. Vol. 99. 1999.

Data warehousing:  
iceberg cube and  
cube-gradient

Foundation for  
many essential  
data mining tasks

Cluster analysis

Pattern analysis in  
spatiotemporal,  
multimedia, time-series,  
and stream data

Broad applications!

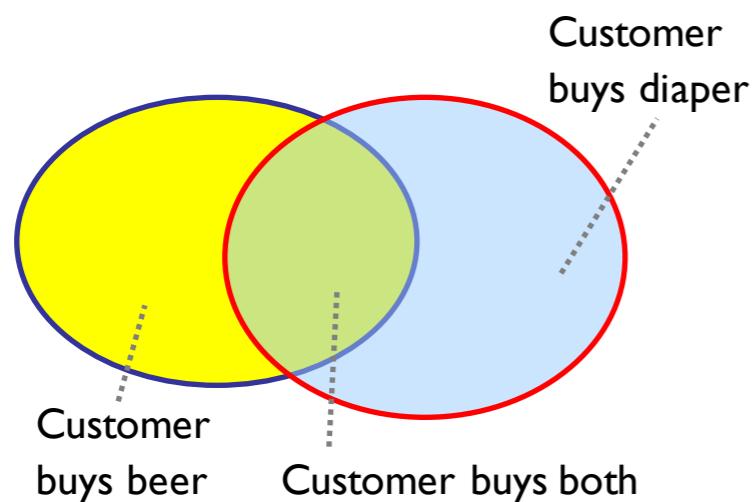
Association,  
correlation, and  
causality analysis

Sequential, structural  
(e.g., sub-graph) patterns

## BASIC IDEAS: FREQUENT PATTERNS

**itemset:** A set of one or more items

Tid	Items bought
10	{Beer, Nuts, Diaper}
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk



**k-itemset  $X = \{x_1, \dots, x_k\}$**

**(absolute)** support, or, support count of  $X$ : Frequency or occurrence of an itemset  $X$

**(relative)** support,  $s$ , is the fraction of transactions that contains  $X$  (i.e., the probability that a transaction contains  $X$ )

An itemset  $X$  is **frequent** if  $X$ 's support is no less than a **minsup threshold**

**itemset**: A set of one or more items

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

## BASIC IDEAS: ASSOCIATION RULES

Find all the rules  $X \rightarrow Y$  with minimum support and confidence

**support**, s, probability that a transaction contains  $X \cup Y$

**confidence**, c, conditional probability that a transaction having X also contains Y

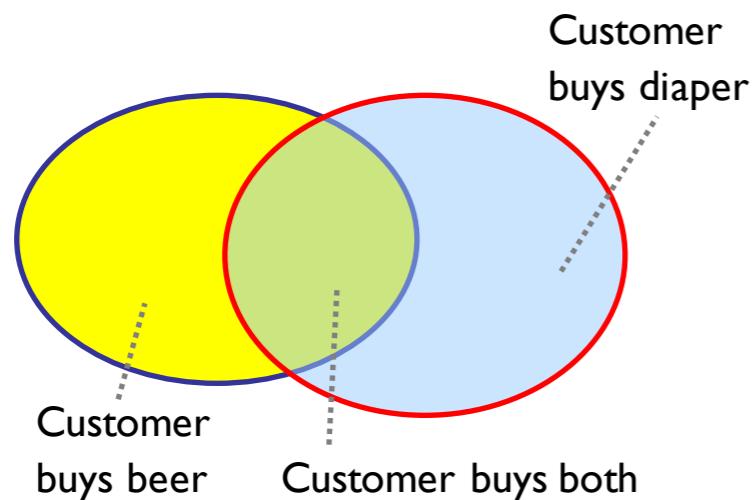
Let  $\text{minsup} = 50\%$ ,  $\text{minconf} = 50\%$

Frequent Patterns: Beer:3, Nuts:3, Diaper:4, Eggs:3, {Beer, Diaper}:3

**Association** rules: (many more!)

$\text{Beer} \rightarrow \text{Diaper}$  (60%, 100%)

$\text{Diaper} \rightarrow \text{Beer}$  (60%, 75%)



joint probability



$$P(X \cup Y) \neq P(X \vee B)$$

X, Y are set elements, not events

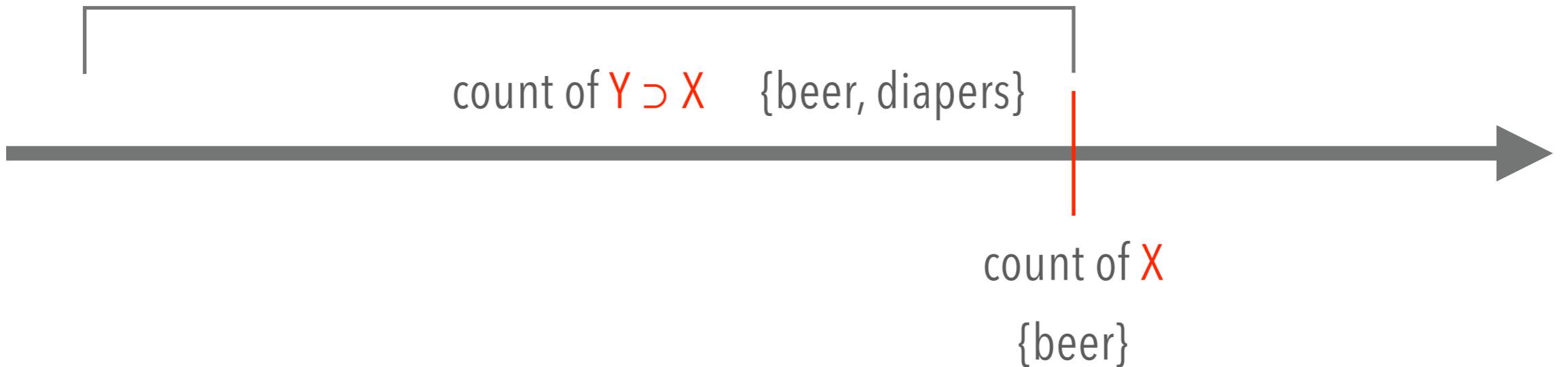
closed

A long pattern  $\{a_1, \dots, a_{100}\}$  contains  
contains a combinatorial number of  
sub-patterns:

$2^{100-1} \approx 10^{30.1}$  sub-patterns!

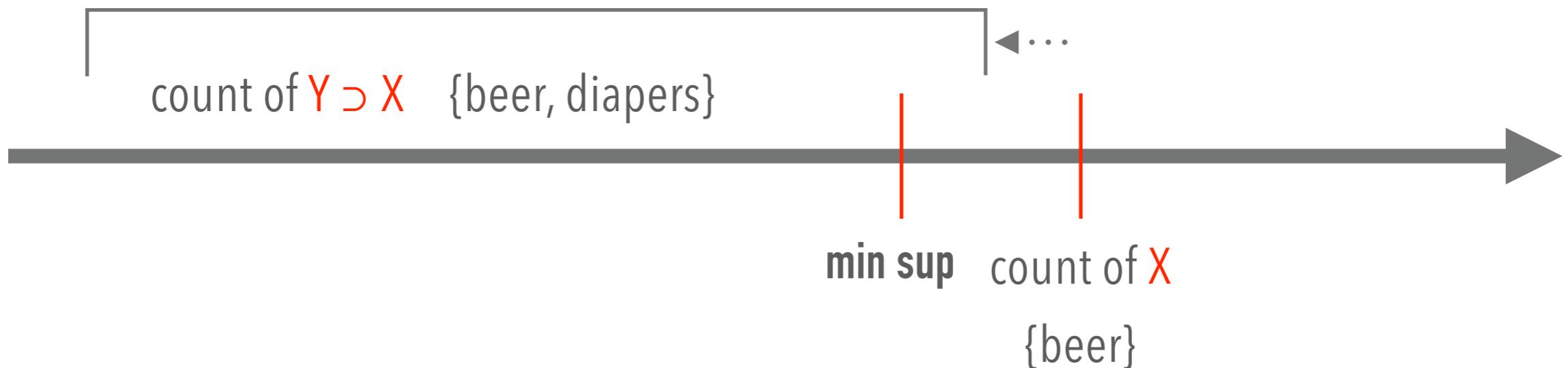
max patterns

the count of the super set  $Y$  is **always** less or equal to the set count  $X$ !



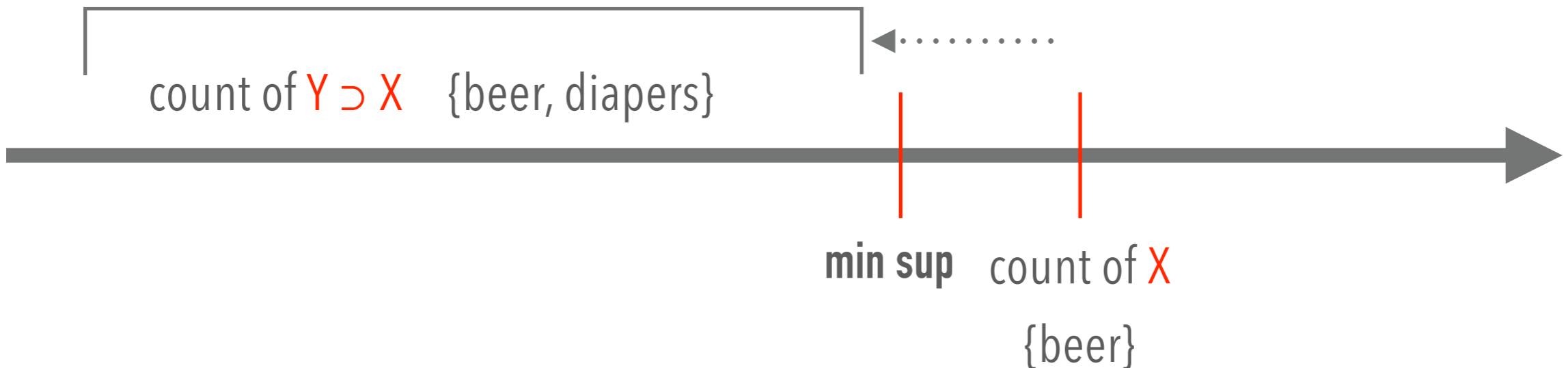


An itemset  $X$  is **closed in dataset D** if there exists **no** proper super-pattern  $Y \supset X$ , with the **same support** as  $X$  in D

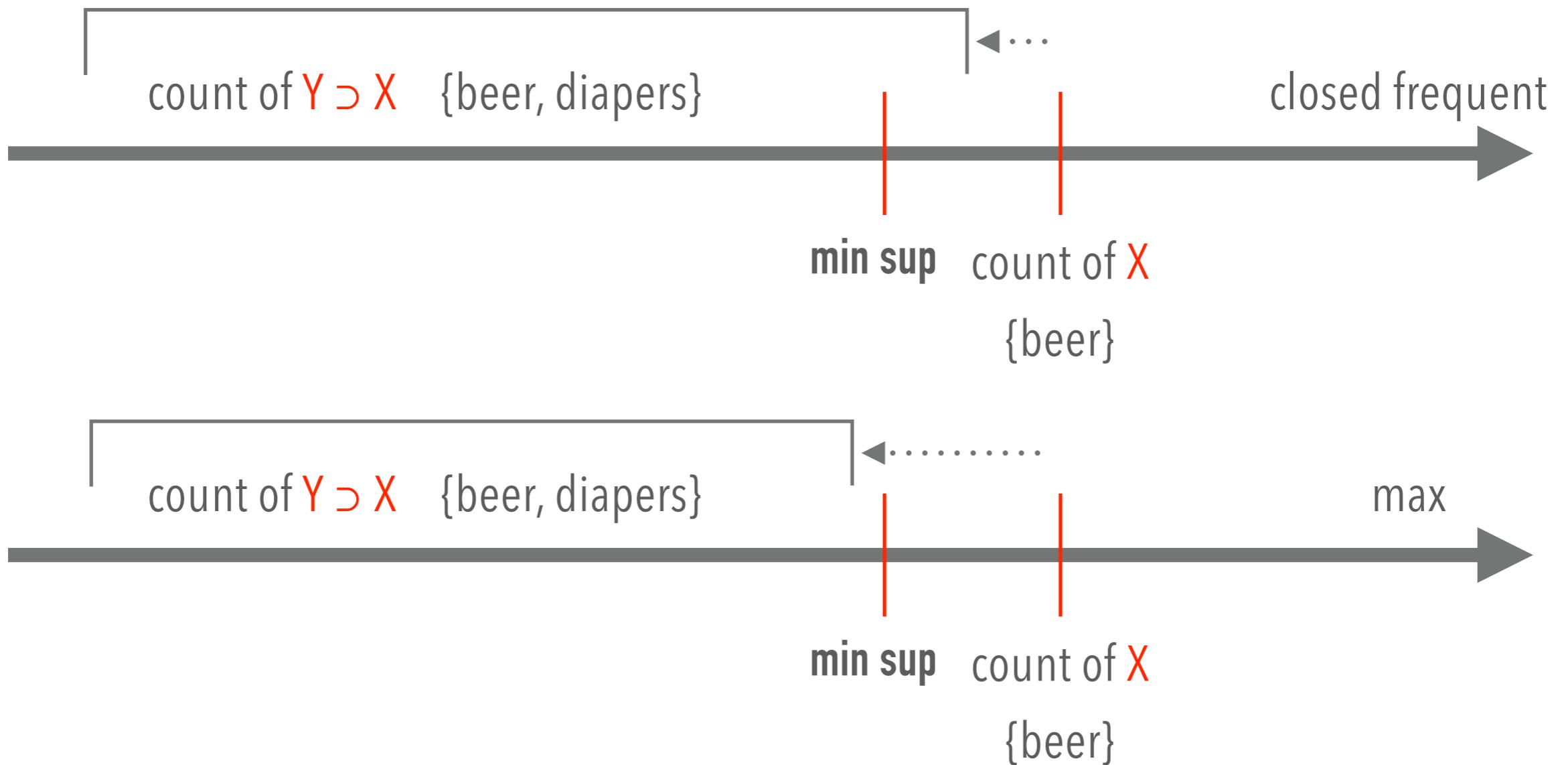


An itemset  $X$  is a closed frequent itemset in set  $D$  if  $X$  is both **closed** and **frequent** in  $D$

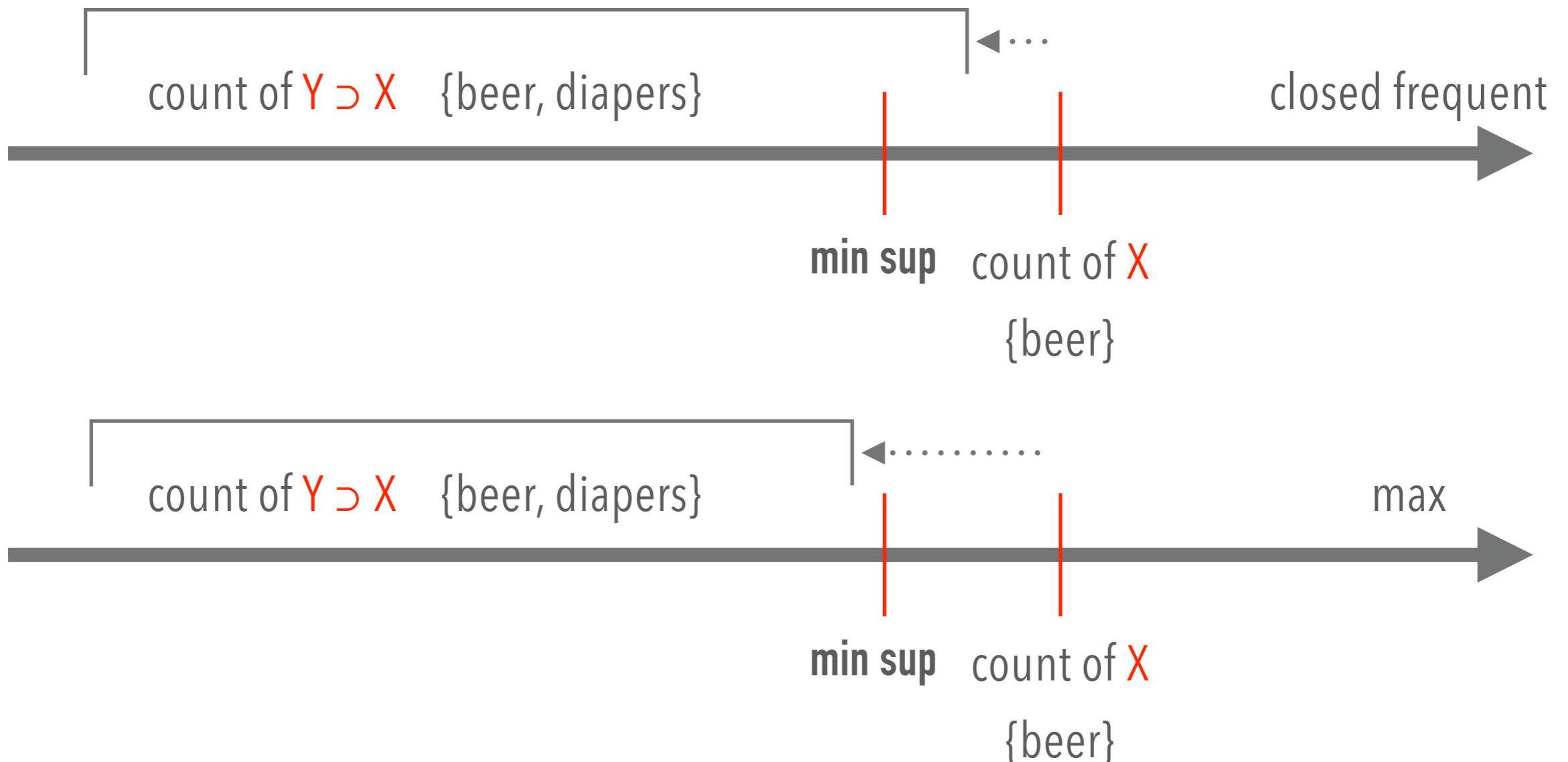
The definition does **not** say anything about  $Y$



An itemset  $X$  is a maximal frequent itemset (or **max-itemset**) in a data set  $D$  if  $X$  is frequent, and there exists no super-itemset  $Y$  such that  $Y \supset X$  and  $Y$  is frequent in  $D$ .



what is a key difference between the set of all  
**max patterns** and the set of all **closed  
frequent** patterns?



You can use the set of closed frequent patterns to enumerate all frequent item sets

assume the following  
four transactions:

{beer, nuts, diapers},  
{beer, nuts, diapers},  
{beer, nuts}, {beer, nuts}

what are closed item sets?

min\_sup=2

$$\{a_1, \dots, a_{100}\}, \{a_1, \dots, a_{50}\}$$

the DB contains only these two transactions

**min sup=1**

$$\{a_1, \dots, a_{100}\}, \{a_1, \dots, a_{50}\}$$

min sup=1

what is the set of closed patterns?

$$\{a_1, \dots, a_{100}\}: 1, \{a_1, \dots, a_{50}\}: 2$$

$$\{a_1, \dots, a_{100}\}, \{a_1, \dots, a_{50}\}$$

$\min \text{ sup} = 1$

what is the set of **max** patterns?

$$\{a_1, \dots, a_{100}\}: I$$

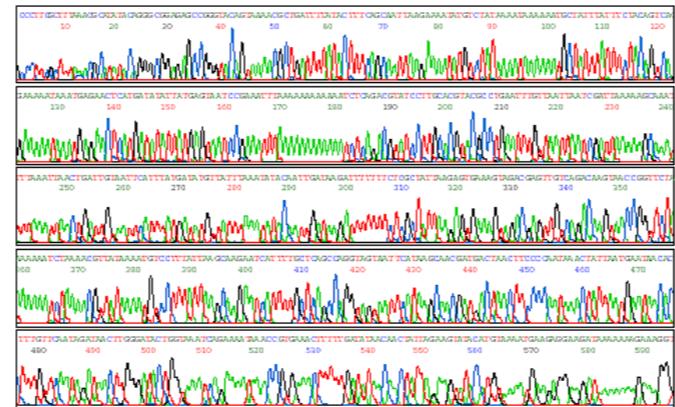
$$\{a_1, \dots, a_{100}\}, \{a_1, \dots, a_{50}\}$$

min sup=1

how many patterns in all?

$$\{a_1\}, \{a_2\}, \{a_1, a_2\}, \dots, \{a_1, \dots, a_{100}\}$$

$2^{100-1}$



a pattern (a set of items, subsequences, substructures, etc.) that occurs frequently in a data set

# BASIC CONCEPTS SUMMARY

Methods

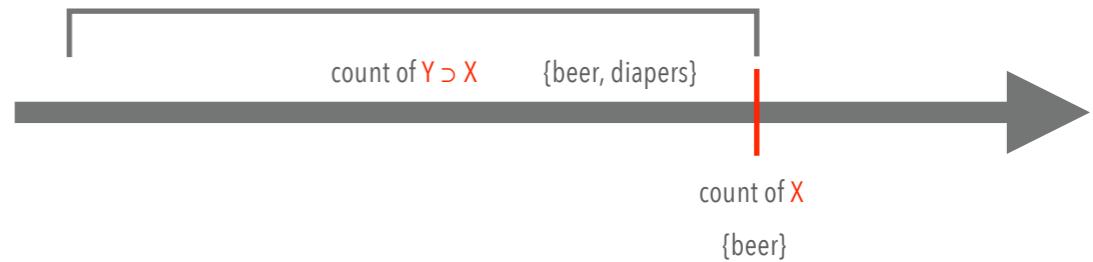
Pattern Evaluation

Summary

itemset, support, frequent

confidence

closed vs. max-pattern



# FREQUENT ITEMSET METHODS

---

Basic Concepts

Pattern Evaluation

Summary

Apriori

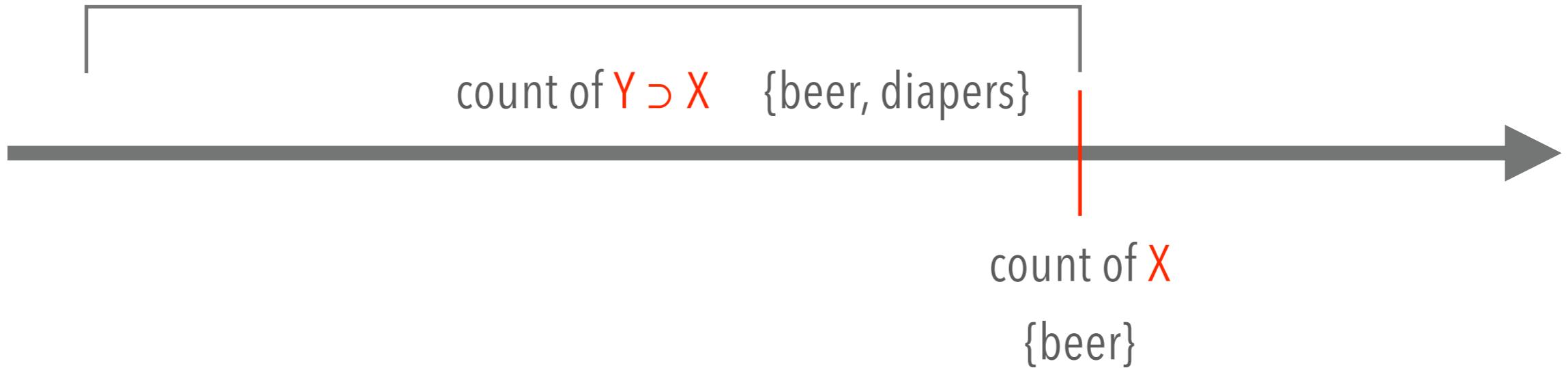
FPGrowth

ECLAT

Max & Closed



the count of a set  $X$  is **always** greater than or equal to the superset  $Y$  count!



If  $\{beer, diaper, nuts\}$  is frequent, so is  $\{beer, diaper\}$

# downward closure property

Any subset of a frequent itemset must be frequent

Rakesh Agrawal and Ramakrishnan Srikant. 1994. [Fast Algorithms for Mining Association Rules in Large Databases](#). In Proceedings of the 20th International Conference on Very Large Data Bases (VLDB '94), Jorge B. Bocca, Matthias Jarke, and Carlo Zaniolo (Eds.). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 487-499.

**Apriori pruning principle:** If there is any itemset which is infrequent, its superset should not be generated or tested!

Rakesh Agrawal and Ramakrishnan Srikant. 1994. [Fast Algorithms for Mining Association Rules in Large Databases](#). In Proceedings of the 20th International Conference on Very Large Data Bases (VLDB '94), Jorge B. Bocca, Matthias Jarke, and Carlo Zaniolo (Eds.). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 487-499.

**Apriori pruning principle:** If there is any itemset which is infrequent, its superset should not be generated or tested!

1. Initially, scan DB once to get frequent 1-itemset
2. Generate length ( $k+1$ ) candidate itemsets from length  $k$  frequent itemsets
3. Test the candidates against DB
4. Terminate when no frequent or candidate set can be generated

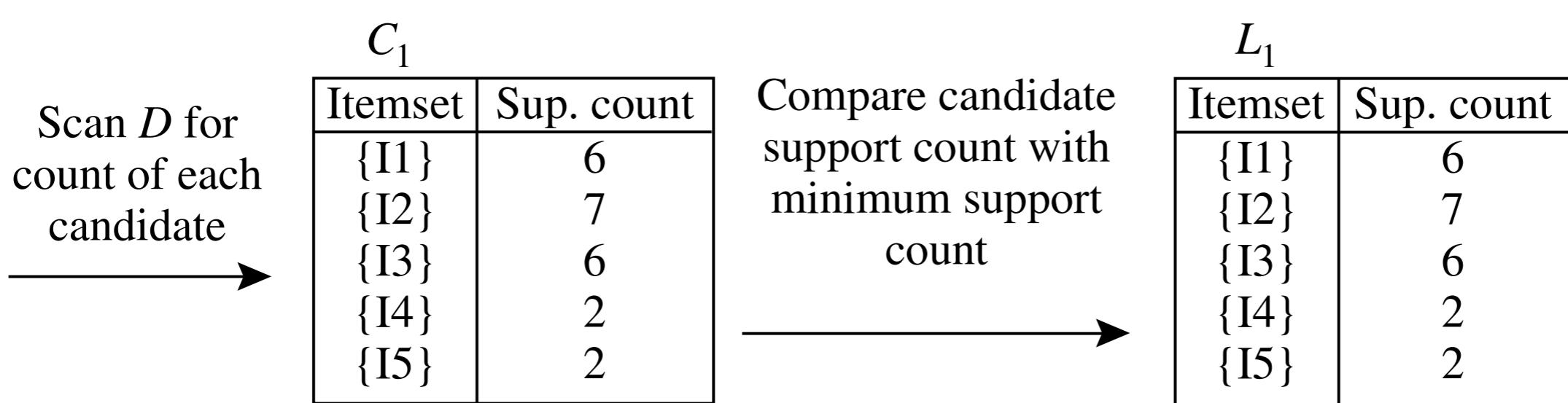
$\text{min sup}=2$

---

<i>TID</i>	<i>List of item_IDs</i>
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

---

<i>TID</i>	<i>List of item_IDs</i>
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3



<i>TID</i>	<i>List of item_IDs</i>
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

Generate  $C_2$   
candidates  
from  $L_1$

Itemset
{I1, I2}
{I1, I3}
{I1, I4}
{I1, I5}
{I2, I3}
{I2, I4}
{I2, I5}
{I3, I4}
{I3, I5}
{I4, I5}

Scan  $D$  for  
count of each  
candidate

$C_2$

Itemset	Sup. count
{I1, I2}	4
{I1, I3}	4
{I1, I4}	1
{I1, I5}	2
{I2, I3}	4
{I2, I4}	2
{I2, I5}	2
{I3, I4}	0
{I3, I5}	1
{I4, I5}	0

Compare candidate  
support count with  
minimum support  
count

Itemset	Sup. count
{I1, I2}	4
{I1, I3}	4
{I1, I5}	2
{I2, I3}	4
{I2, I4}	2
{I2, I5}	2

$L_2$ 

Itemset	Sup. count
{I1, I2}	4
{I1, I3}	4
{I1, I5}	2
{I2, I3}	4
{I2, I4}	2
{I2, I5}	2

TID	List of item_IDs
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

$C_3$   
Generate  $C_3$   
candidates  
from  $L_2$   
 $\longrightarrow$

Scan  $D$  for  
count of each  
candidate  
 $\longrightarrow$

Itemset	Sup. count
{I1, I2, I3}	2
{I1, I2, I5}	2

Compare candidate  
support count  
with minimum  
support count

Itemset	Sup. count
{I1, I2, I3}	2
{I1, I2, I5}	2

# APRIORI IMPLEMENTATION

---

**Algorithm: Apriori.** Find frequent itemsets using an iterative level-wise approach based on candidate generation.

**Input:**

- $D$ , a database of transactions;
- $\text{min\_sup}$ , the minimum support count threshold.

**Output:**  $L$ , frequent itemsets in  $D$ .

**Method:**

```
(1)    $L_1 = \text{find\_frequent\_1-itemsets}(D);$ 
(2)   for ( $k = 2; L_{k-1} \neq \emptyset; k++$ ) {
(3)      $C_k = \text{apriori\_gen}(L_{k-1});$ 
(4)     for each transaction  $t \in D$  { // scan  $D$  for counts
(5)        $C_t = \text{subset}(C_k, t);$  // get the subsets of  $t$  that are candidates
(6)       for each candidate  $c \in C_t$ 
(7)          $c.\text{count}++;$ 
(8)     }
(9)      $L_k = \{c \in C_k | c.\text{count} \geq \text{min\_sup}\}$ 
(10)   }
(11)   return  $L = \cup_k L_k;$ 
```

```
procedure apriori_gen( $L_{k-1}$ : frequent  $(k-1)$ -itemsets)
(1)   for each itemset  $l_1 \in L_{k-1}$ 
(2)     for each itemset  $l_2 \in L_{k-1}$ 
(3)       if ( $l_1[1] = l_2[1]$ )  $\wedge (l_1[2] = l_2[2])$ 
(4)          $\dots \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$  then {
(5)            $c = l_1 \bowtie l_2;$  // join step: generate candidates
(6)           if has_infrequent_subset( $c, L_{k-1}$ ) then
(7)             delete  $c;$  // prune step: remove unfruitful candidate
(8)           else add  $c$  to  $C_k;$ 
(9)     }
(9)   return  $C_k;$ 
```

```
procedure has_infrequent_subset( $c$ : candidate  $k$ -itemset;
                                 $L_{k-1}$ : frequent  $(k-1)$ -itemsets); // use prior knowledge
(1)   for each  $(k-1)$ -subset  $s$  of  $c$ 
(2)     if  $s \notin L_{k-1}$  then
(3)       return TRUE;
(4)   return FALSE;
```

How to generate candidates?

Step 1: self-joining  $L_k$

Step 2: pruning

Example of Candidate-generation

$L_3 = \{\text{abc}, \text{abd}, \text{acd}, \text{ace}, \text{bcd}\}$

Self-joining:  $L_3 * L_3$

{abcd} from {abc} and {abd}

{acde} from {acd} and {ace}

Pruning:

{acde} is removed because {ade} is not in  $L_3$

$C_4 = \{\text{abcd}\}$

**Algorithm: Apriori.** Find frequent itemsets using an iterative level-wise approach based on candidate generation.

**Input:**

- $D$ , a database of transactions;
- $min\_sup$ , the minimum support count threshold.

**Output:**  $L$ , frequent itemsets in  $D$ .

**Method:**

```
(1)    $L_1 = \text{find\_frequent\_1-itemsets}(D);$ 
(2)   for ( $k = 2; L_{k-1} \neq \emptyset; k++$ ) {
(3)      $C_k = \text{apriori\_gen}(L_{k-1});$ 
(4)     for each transaction  $t \in D$  { // scan  $D$  for counts
(5)        $C_t = \text{subset}(C_k, t);$  // get the subsets of  $t$  that are candidates
(6)       for each candidate  $c \in C_t$ 
(7)          $c.\text{count}++;$ 
(8)     }
(9)      $L_k = \{c \in C_k | c.\text{count} \geq min\_sup\}$ 
(10)   }
(11)   return  $L = \cup_k L_k;$ 
```

procedure  $\text{apriori\_gen}(L_{k-1}:\text{frequent } (k-1)\text{-itemsets})$

```
(1)   for each itemset  $l_1 \in L_{k-1}$ 
(2)     for each itemset  $l_2 \in L_{k-1}$ 
(3)       if ( $l_1[1] = l_2[1]$ )  $\wedge (l_1[2] = l_2[2])$ 
(4)          $\wedge \dots \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$  then {
(5)          $c = l_1 \bowtie l_2;$  // join step: generate candidates
(6)         if  $\text{has\_infrequent\_subset}(c, L_{k-1})$  then
(7)           delete  $c;$  // prune step: remove unfruitful candidate
(8)           else add  $c$  to  $C_k;$ 
(9)       }
(9)   return  $C_k;$ 
```

procedure  $\text{has\_infrequent\_subset}(c: \text{candidate } k\text{-itemset};$

```
     $L_{k-1}: \text{frequent } (k-1)\text{-itemsets); // use prior knowledge}$ 
(1)   for each  $(k-1)$ -subset  $s$  of  $c$ 
(2)     if  $s \notin L_{k-1}$  then
(3)       return TRUE;
(4)   return FALSE;
```

$$L_3 = \{\text{abc, abd, acd, ace, bcd}\}$$

Self-joining:  $L_3 * L_3$

{abcd} from {abc} and {abd}

{acde} from {acd} and {ace}

{acde} is removed because  
{ade} is **not** in  $L_3$

$$C_4 = \{\text{abcd}\}$$

**apriori is  
inefficient!**



# challenges



Multiple scans of transaction database

reduce passes

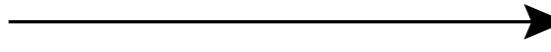
Huge number of candidates

shrink

Tedious workload of  
support counting for  
candidates

make it easy to  
keep count

Create hash table  $H_2$   
using hash function  

$$h(x, y) = ((\text{order of } x) \times 10 + (\text{order of } y)) \bmod 7$$


bucket address	0	1	2	3	4	5	6
bucket count	2	2	4	2	2	4	4
bucket contents	{I1, I4} {I3, I5}	{I1, I5} {I1, I5}	{I2, I3} {I2, I3} {I2, I3}	{I2, I4} {I2, I4}	{I2, I5} {I2, I5}	{I1, I2} {I1, I2}	{I1, I3} {I1, I3} {I1, I3}
							{I1, I2} {I1, I2}

# hashing

A **k**-itemset whose corresponding hashing bucket count is below the threshold cannot be frequent

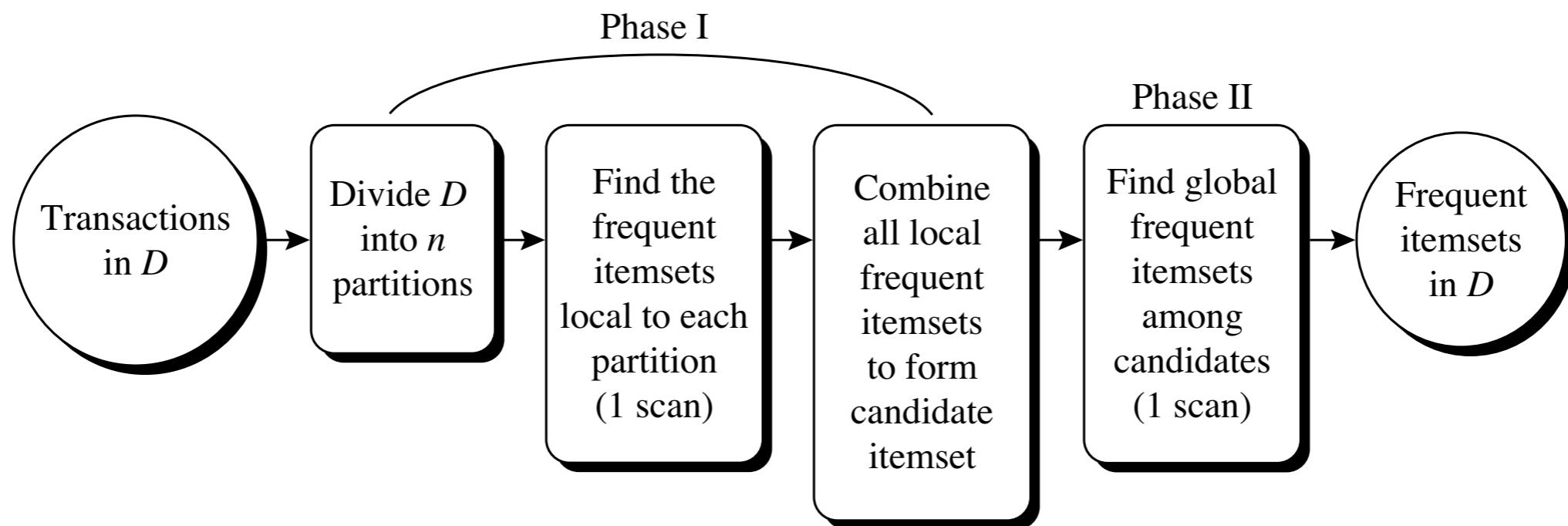
Jong Soo Park, Ming-Syan Chen, and Philip S. Yu. 1995. [An effective hash-based algorithm for mining association rules](#). In Proceedings of the 1995 ACM SIGMOD international conference on Management of data (SIGMOD '95), Michael Carey and Donovan Schneider (Eds.). ACM, New York, NY, USA, 175-186.

Any itemset that is potentially frequent in DB must be frequent in at least one of the partitions of DB

Scan 1: partition database and find local frequent patterns

Scan 2: consolidate global frequent patterns

# partitioning



# sampling

Hannu Toivonen. 1996. **Sampling Large Databases for Association Rules.** In Proceedings of the 22th International Conference on Very Large Data Bases (VLDB '96), T. M. Vijayaraman, Alejandro P. Buchmann, C. Mohan, and Nandlal L. Sarda (Eds.). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 134-145.

# FREQUENT ITEMSET METHODS

---

Basic Concepts

Pattern Evaluation

Summary

Apriori

FPGrowth

ECLAT

Max & Closed



# MINING FREQUENT PATTERNS

Bottlenecks of the Apriori approach

Breadth-first (i.e., level-wise) search

Candidate generation and test (Often generates a huge number of candidates)

The FP-Growth Approach

Depth-first search

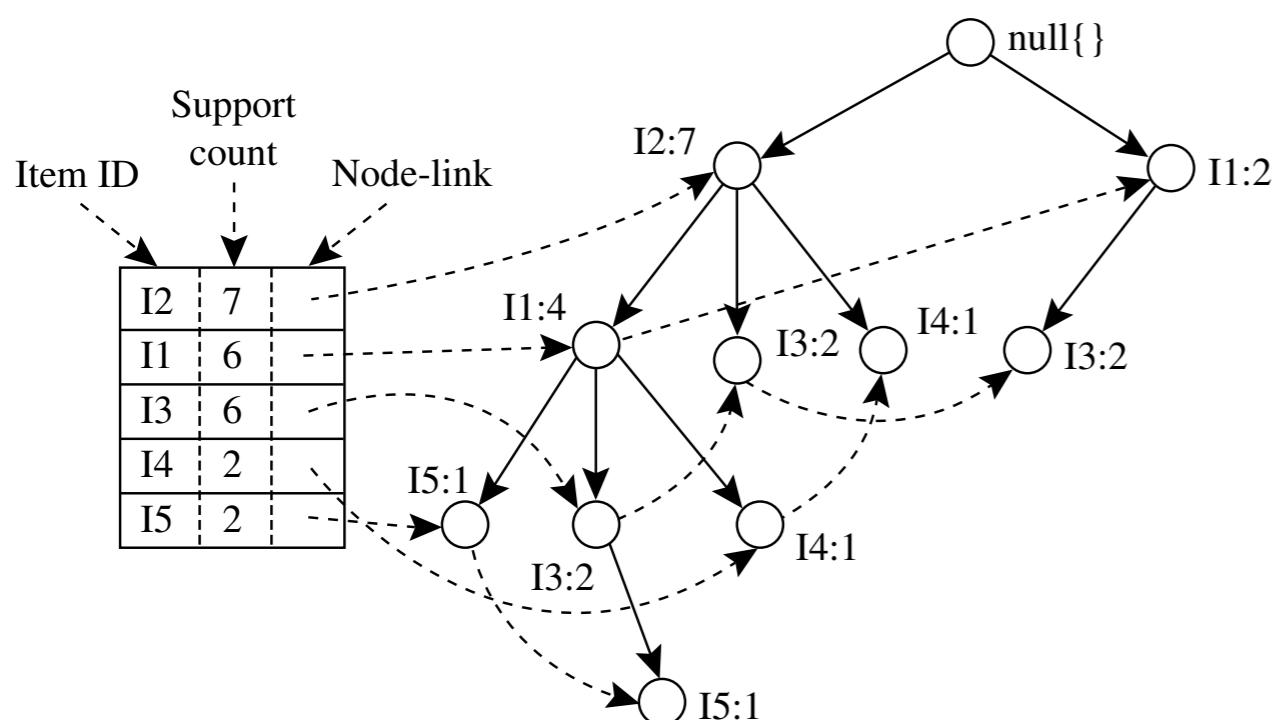
Avoid explicit candidate generation

Major philosophy: Grow long patterns from short ones using local frequent items only

“abc” is a frequent pattern

Get all transactions having “abc”, i.e., project DB on abc: DB|abc

“d” is a local frequent item in DB|abc  
abcd is a frequent pattern



Jiawei Han, Jian Pei, and Yiwen Yin. 2000. Mining frequent patterns without candidate generation. In Proceedings of the 2000 ACM SIGMOD international conference on Management of data (SIGMOD '00). ACM, New York, NY, USA, 1-12.

Prof. Han's most cited paper! ~**7,800** citations

## Branch

---

<i>TID</i>	<i>List of item_IDs</i>
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

---

Branch

<i>TID</i>	<i>List of item IDs</i>
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

most frequent

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

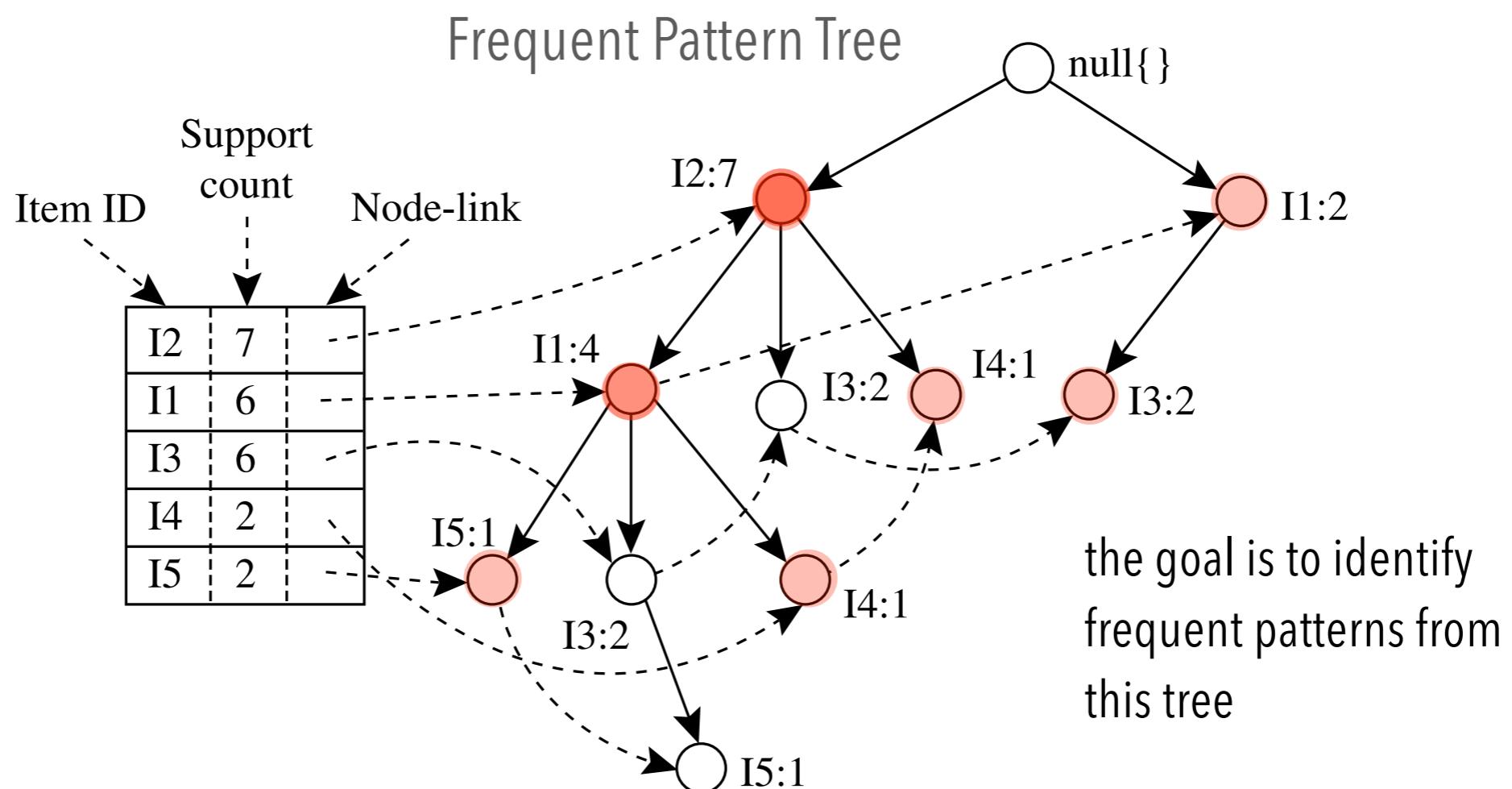
**min sup=2**

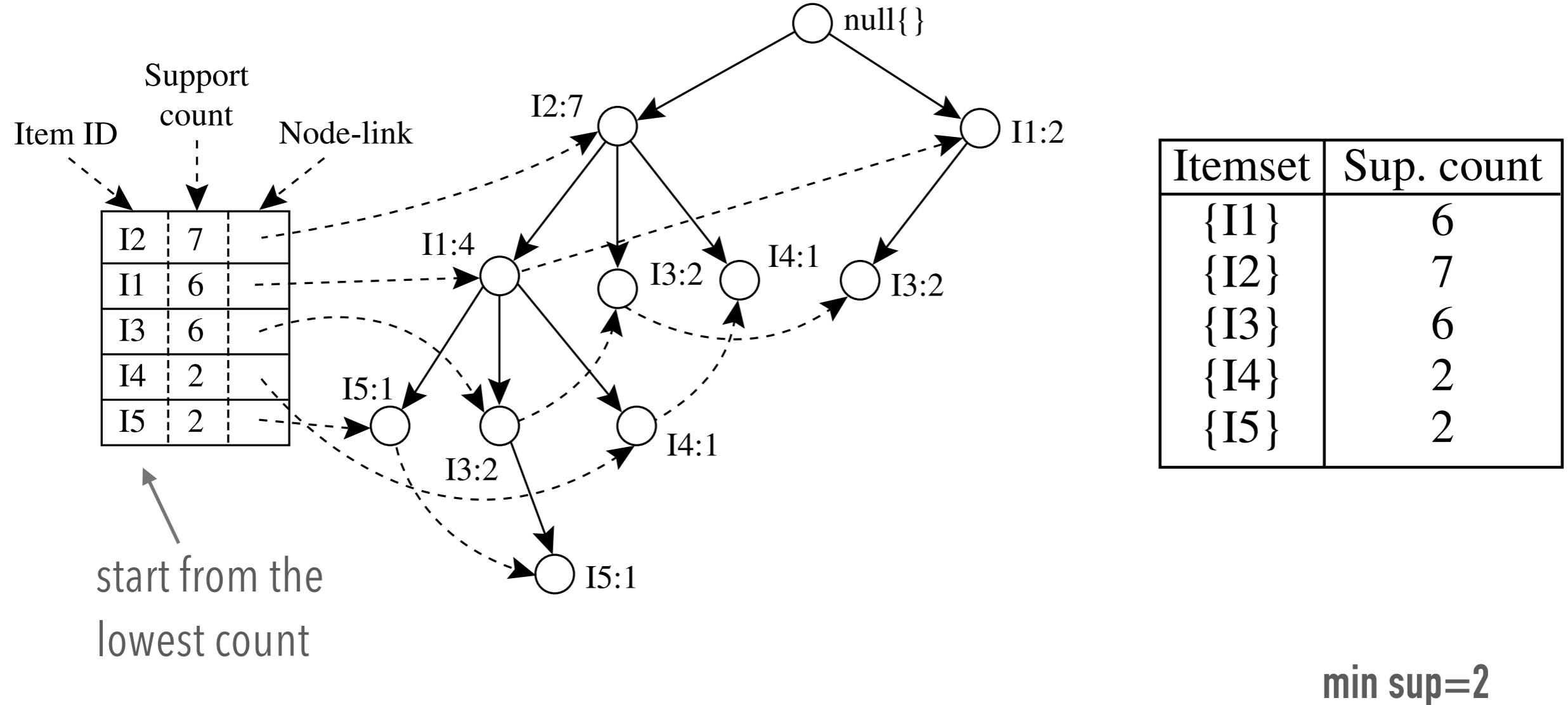
Itemset	Sup. count
{I1}	6
{I2}	7
{I3}	6
{I4}	2
{I5}	2

## Branch

<i>TID</i>	<i>List of item_IDs</i>
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

Itemset	Sup. count
{I1}	6
{I2}	7
{I3}	6
{I4}	2
{I5}	2

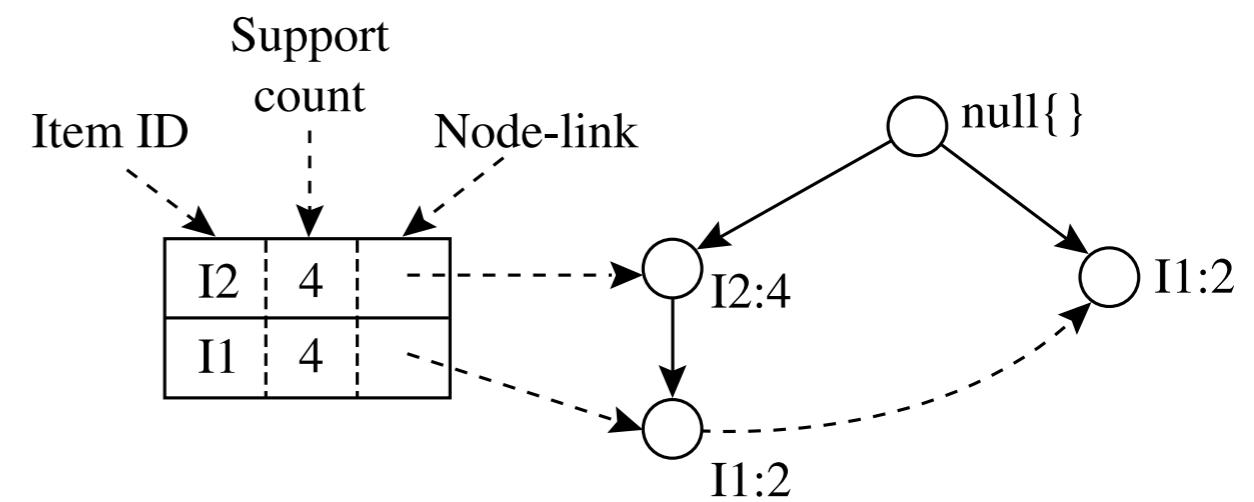
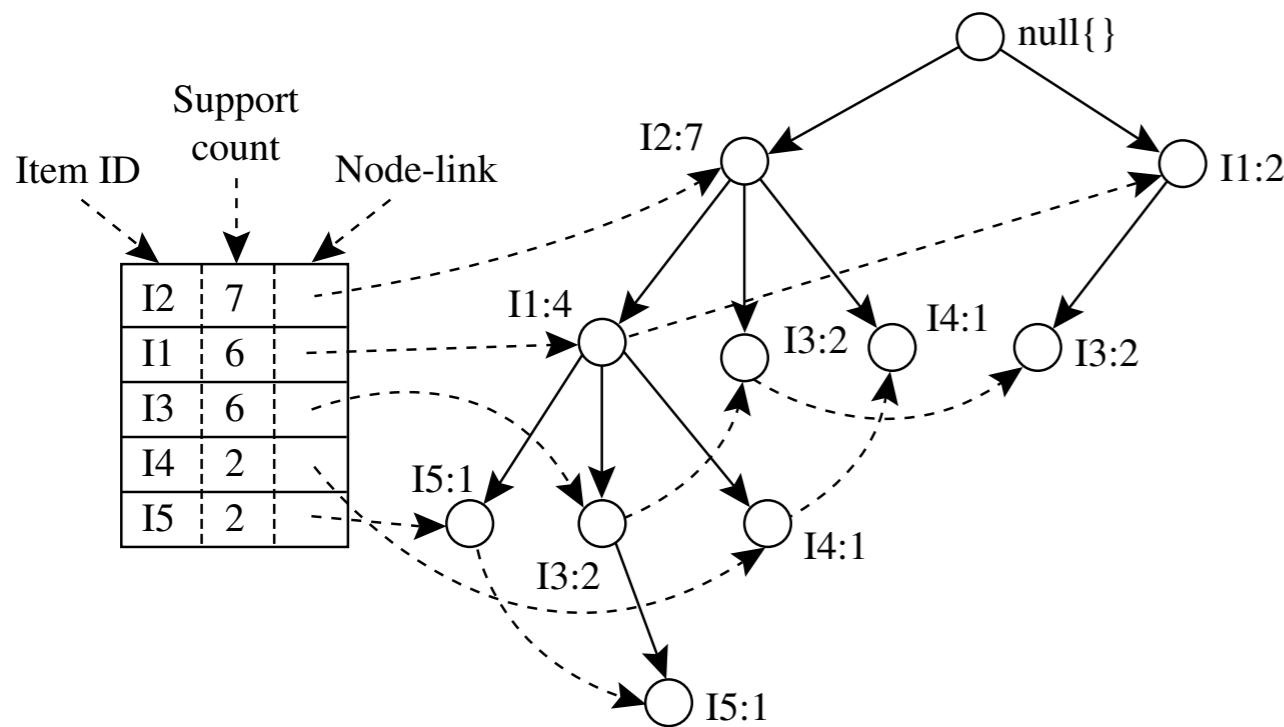




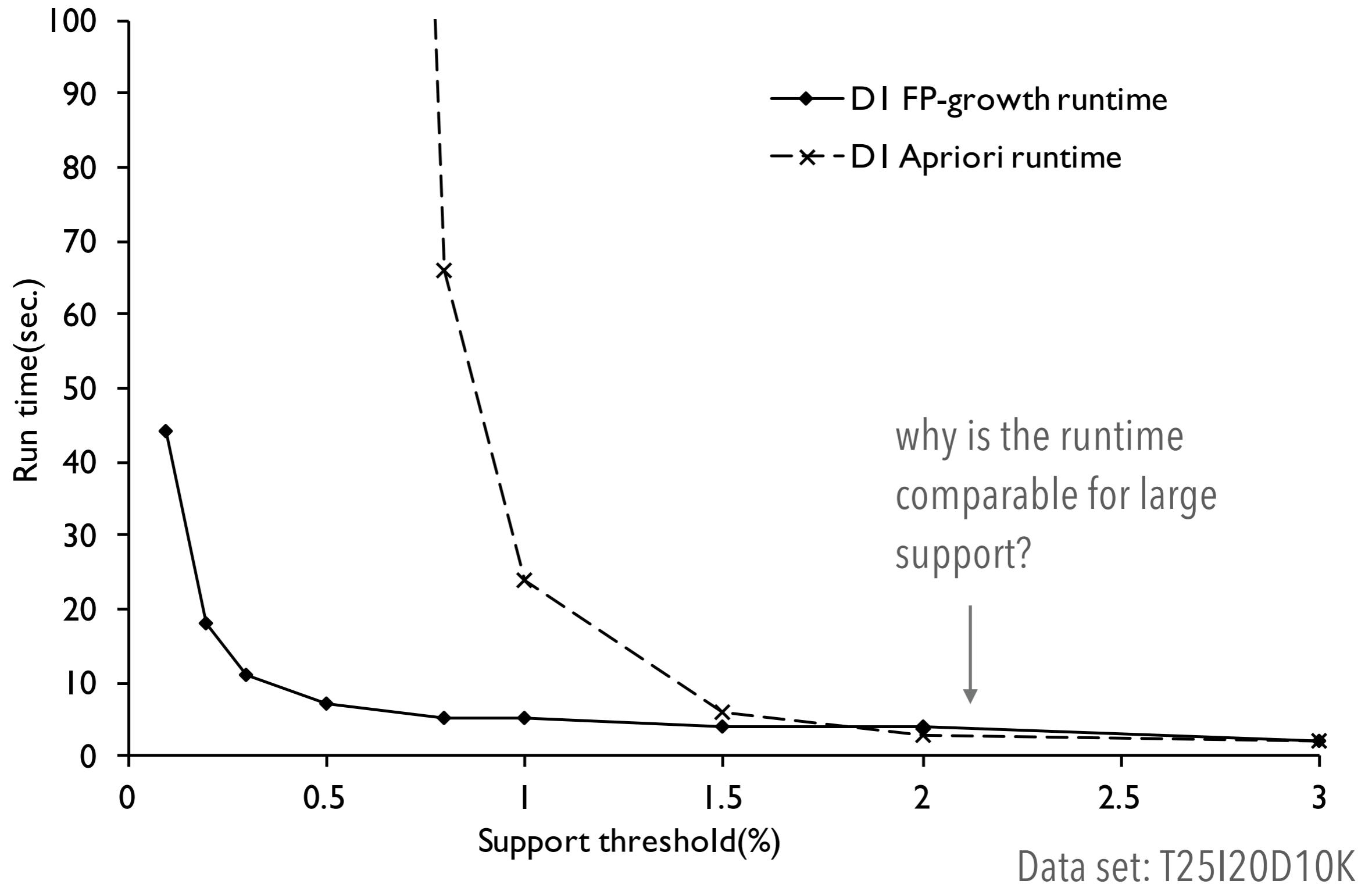
Item	Conditional Pattern Base	Conditional FP-tree	Frequent Patterns Generated
I5	$\{\{I2, I1: 1\}, \{I2, I1, I3: 1\}\}$	$\langle I2: 2, I1: 2 \rangle$	$\{I2, I5: 2\}, \{I1, I5: 2\}, \{I2, I1, I5: 2\}$
I4	$\{\{I2, I1: 1\}, \{I2: 1\}\}$	$\langle I2: 2 \rangle$	$\{I2, I4: 2\}$
I3	$\{\{I2, I1: 2\}, \{I2: 2\}, \{I1: 2\}\}$	$\langle I2: 4, I1: 2 \rangle, \langle I1: 2 \rangle$	$\{I2, I3: 4\}, \{I1, I3: 4\}, \{I2, I1, I3: 2\}$
I1	$\{\{I2: 4\}\}$	$\langle I2: 4 \rangle$	$\{I2, I1: 4\}$

$\text{min sup}=2$

Item	Conditional Pattern Base	Conditional FP-tree	Frequent Patterns Generated
I5	$\{\{I2, I1: 1\}, \{I2, I1, I3: 1\}\}$	$\langle I2: 2, I1: 2 \rangle$	$\{I2, I5: 2\}, \{I1, I5: 2\}, \{I2, I1, I5: 2\}$
I4	$\{\{I2, I1: 1\}, \{I2: 1\}\}$	$\langle I2: 2 \rangle$	$\{I2, I4: 2\}$
I3	$\{\{I2, I1: 2\}, \{I2: 2\}, \{I1: 2\}\}$	$\langle I2: 4, I1: 2 \rangle, \langle I1: 2 \rangle$	$\{I2, I3: 4\}, \{I1, I3: 4\}, \{I2, I1, I3: 2\}$
I1	$\{\{I2: 4\}\}$	$\langle I2: 4 \rangle$	$\{I2, I1: 4\}$



# FP-Growth vs.Apriori: Scalability With the Support Threshold



# FP-GROWTH SUMMARY

**Idea:** Frequent pattern growth

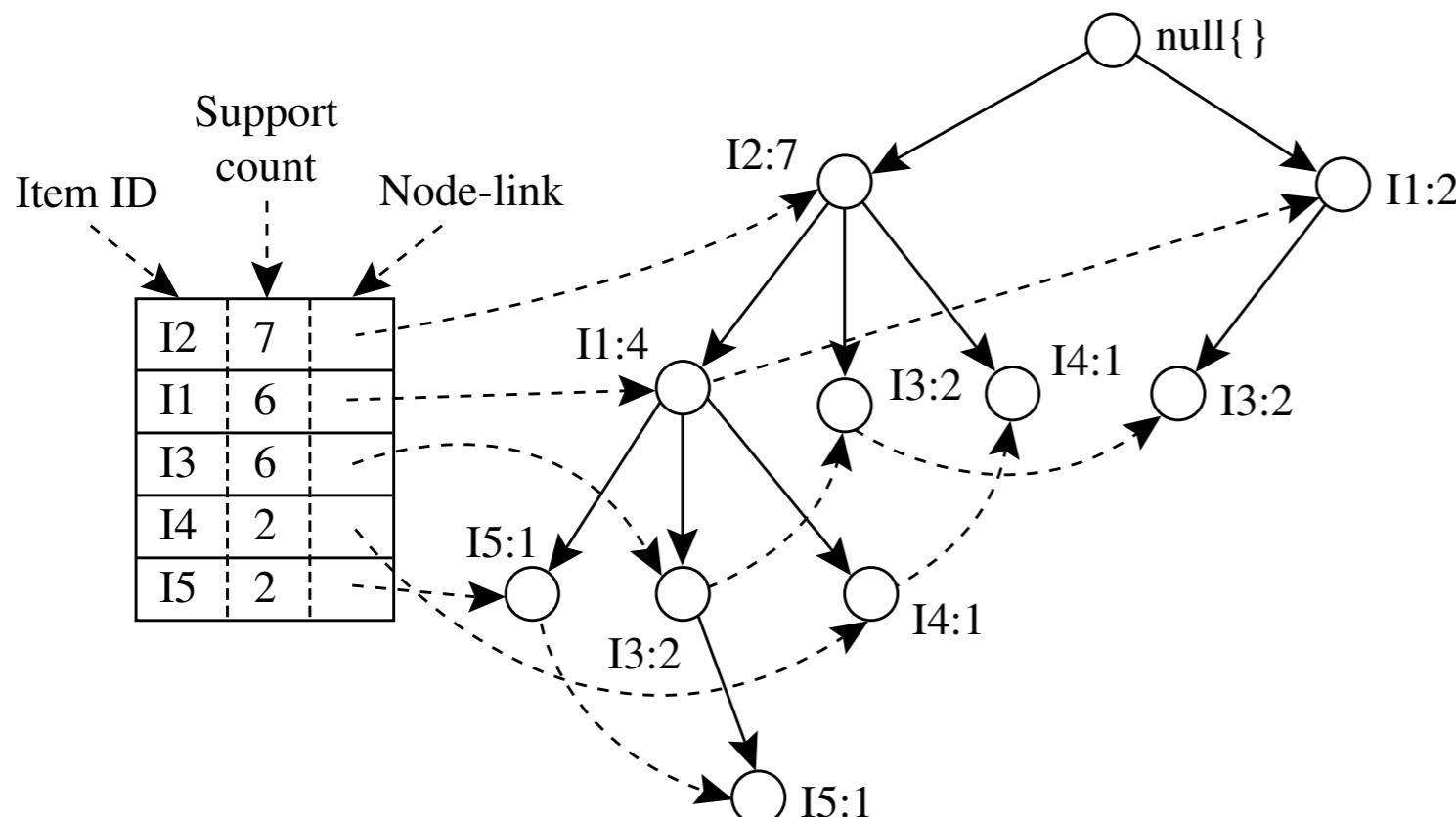
Recursively grow frequent patterns by pattern and database partition

## Method

For each frequent 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



# FP-TREE BENEFITS

---

## Divide-and-conquer:

Decompose both the mining task and DB according to the frequent patterns obtained so far

Lead to focused search of smaller databases

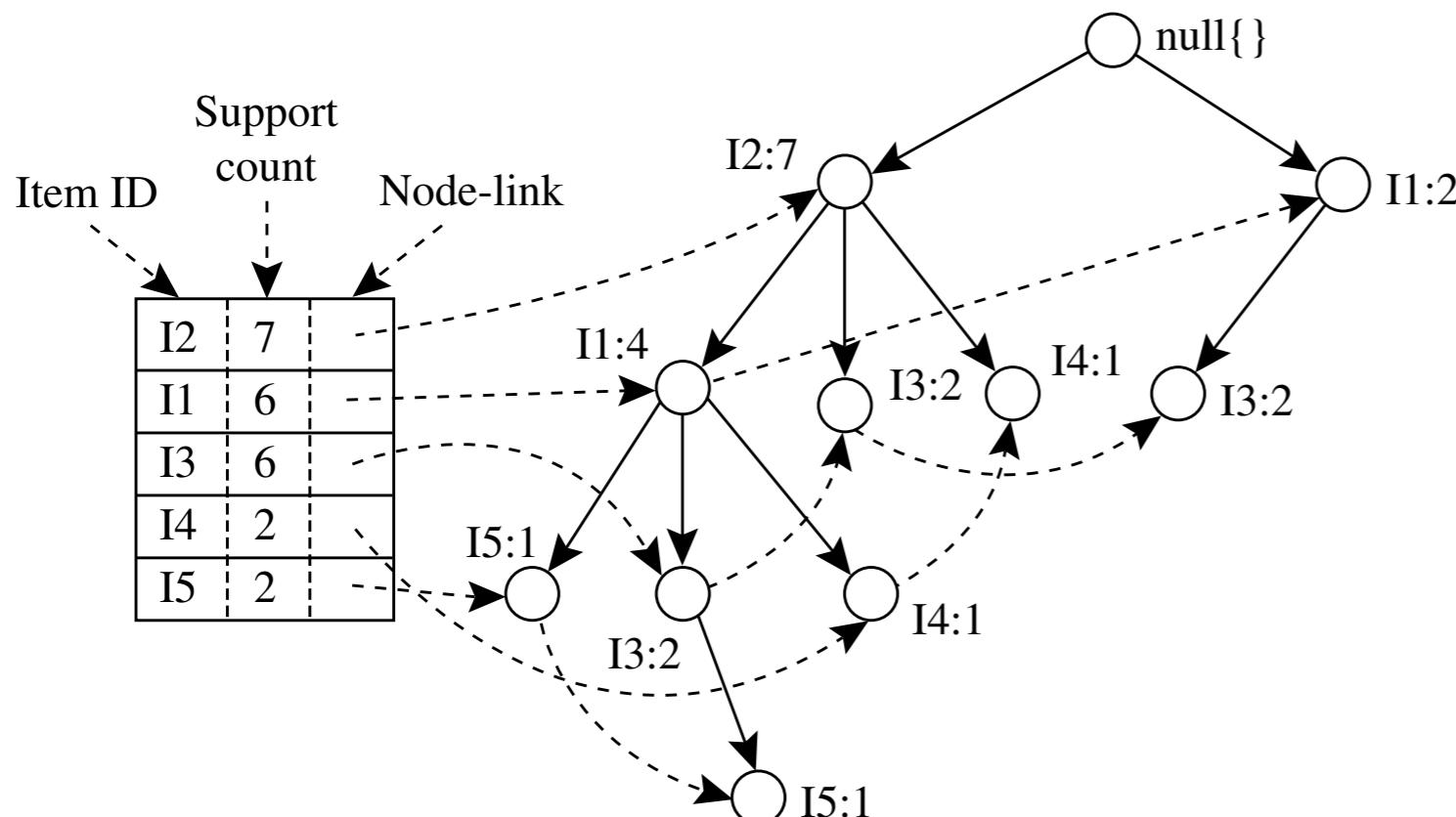
## Other factors

No candidate generation, no candidate test

Compressed database: FP-tree structure

No repeated scan of entire database

Basic operations: counting local frequent items and building sub FP-tree, no pattern search and matching



# FP-TREE BENEFITS

.....

## Completeness

Preserves complete information  
for frequent pattern mining

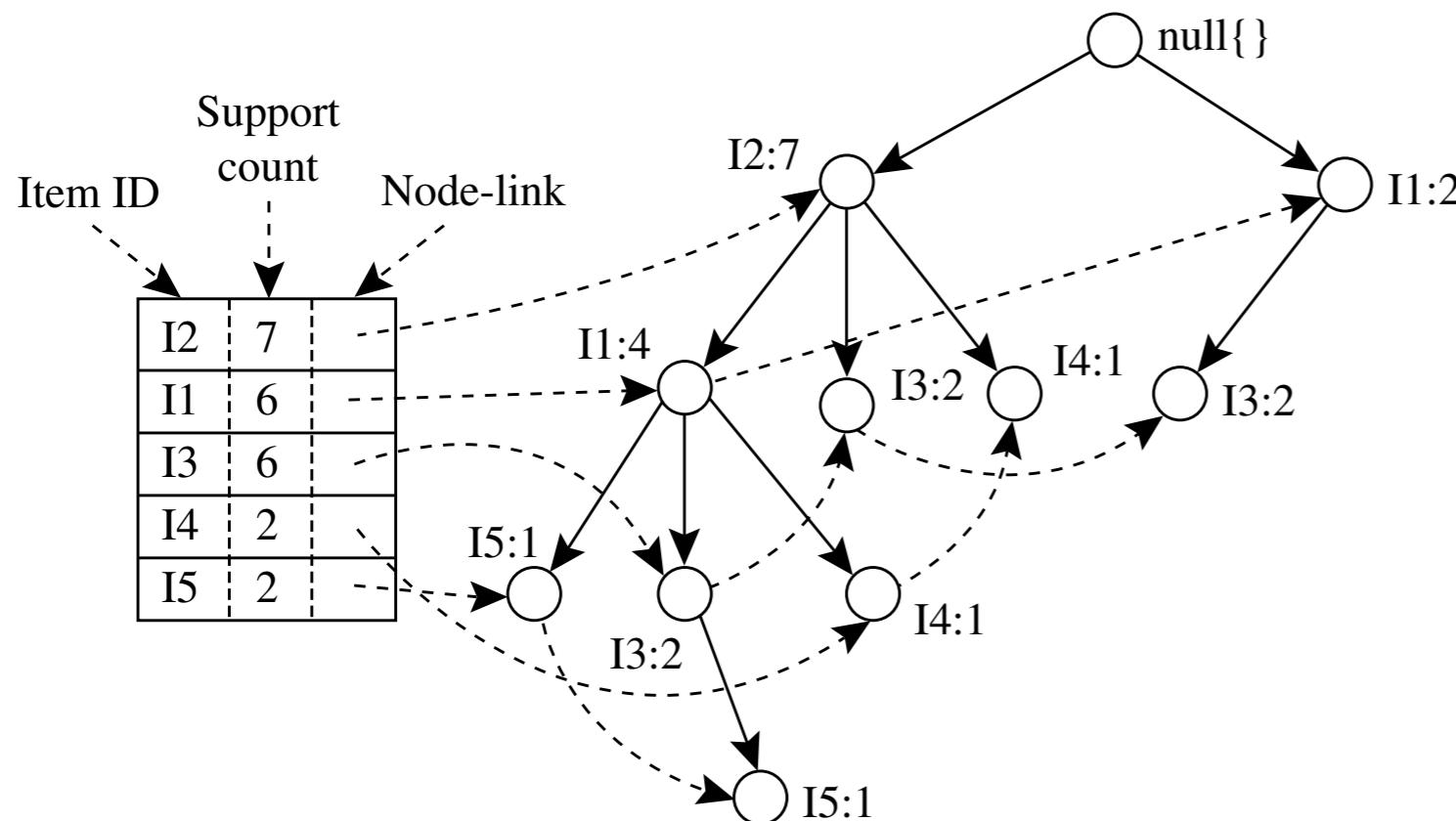
Never breaks a long pattern of  
any transaction

## Compactness

Reduce irrelevant information—  
infrequent items are gone

Items in frequency descending  
order: the more frequently  
occurring pattern, the more likely  
it will be shared

Is never larger than the original  
database (not counting node-  
links and the count field)



what if the tree  
cannot fit in main  
memory?

partition!

# FREQUENT ITEMSET METHODS

---

Basic Concepts

Pattern Evaluation

Summary

Apriori

FPGrowth

ECLAT

Max & Closed



<i>itemset</i>	<i>TID_set</i>
I1	{T100, T400, T500, T700, T800, T900}
I2	{T100, T200, T300, T400, T600, T800, T900}
I3	{T300, T500, T600, T700, T800, T900}
I4	{T200, T400}
I5	{T100, T800}

<i>itemset</i>	<i>TID_set</i>
{I1, I2}	{T100, T400, T800, T900}
{I1, I3}	{T500, T700, T800, T900}
{I1, I4}	{T400}
{I1, I5}	{T100, T800}
{I2, I3}	{T300, T600, T800, T900}
{I2, I4}	{T200, T400}
{I2, I5}	{T100, T800}
{I3, I5}	{T800}

## ECLAT: VERTICAL DATA MINING

.....

Vertical format:  $t(AB) = \{T_{11}, T_{25}, \dots\}$

tid-list: list of transaction ids containing an itemset

Deriving frequent patterns based on vertical intersections

$t(X) = t(Y)$ : X and Y always happen together

$t(X) \subset t(Y)$ : transaction having X always has Y

Using diffset to accelerate mining

Only keep track of differences of tids

$t(X) = \{T_1, T_2, T_3\}$ ,  $t(XY) = \{T_1, T_3\}$

Diffset (XY, X) = {T<sub>2</sub>}

Zaki, M.J., Parthasarathy, S., Ogihsara, M. and Li, W., 1997, August. New Algorithms for Fast Discovery of Association Rules. In KDD (Vol. 97, pp. 283-286).

# FREQUENT ITEMSET METHODS

---

Basic Concepts

Pattern Evaluation

Summary

Apriori

FPGrowth

ECLAT

Max & Closed

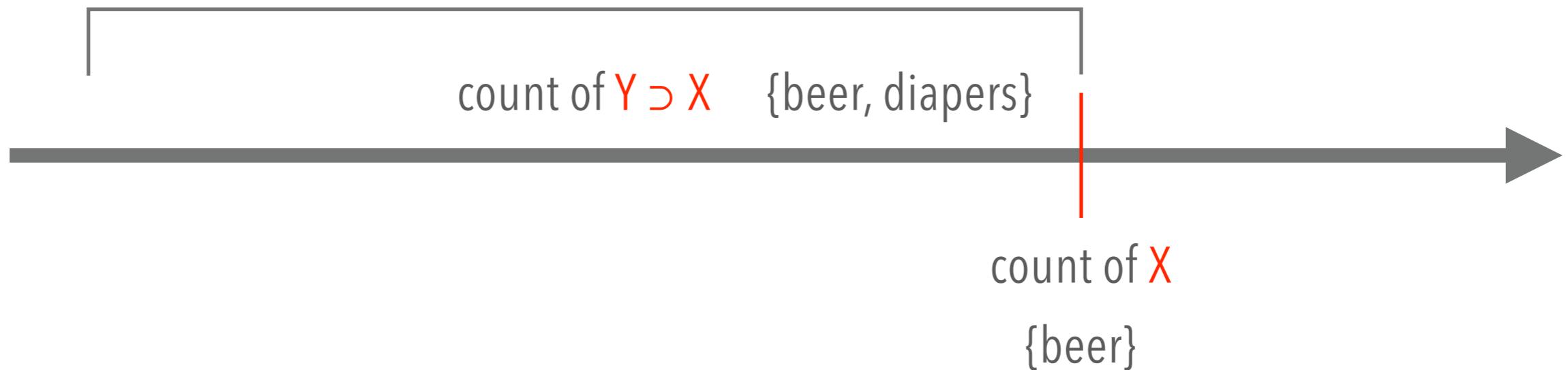


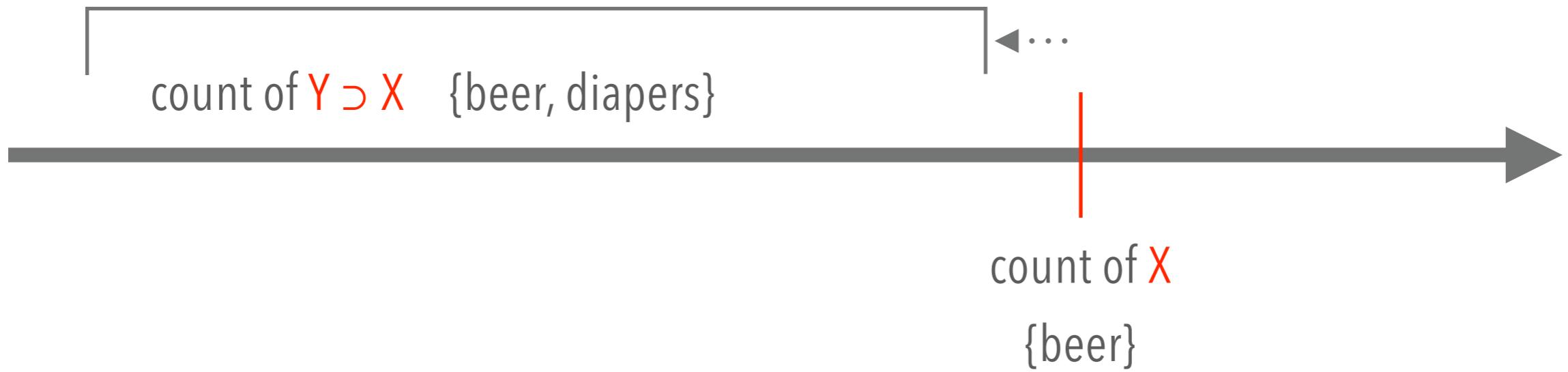
why be interested  
in **closed** and **max-**  
patterns?

A long pattern  $\{a_1, \dots, a_{100}\}$  contains  
contains a **combinatorial** number of  
sub-patterns:

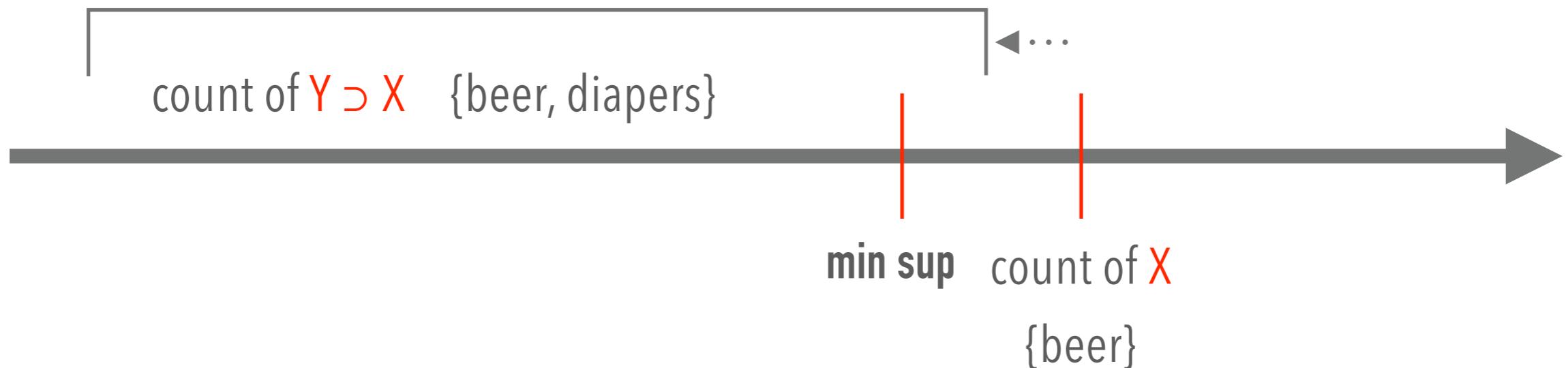
$2^{100-1} \approx 10^{30.1}$  sub-patterns!

the count of the super set  $Y$  is **always** less or equal to the count of set  $X$ !



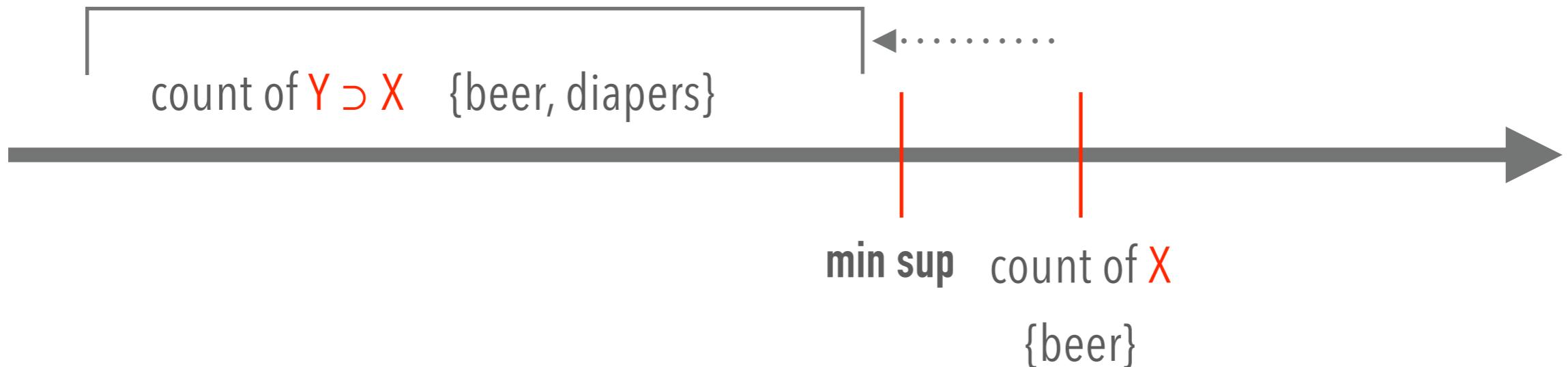


An itemset  $X$  is **closed in dataset D** if there exists **no** proper super-pattern  $Y \supset X$ , with the **same support** as  $X$  in D

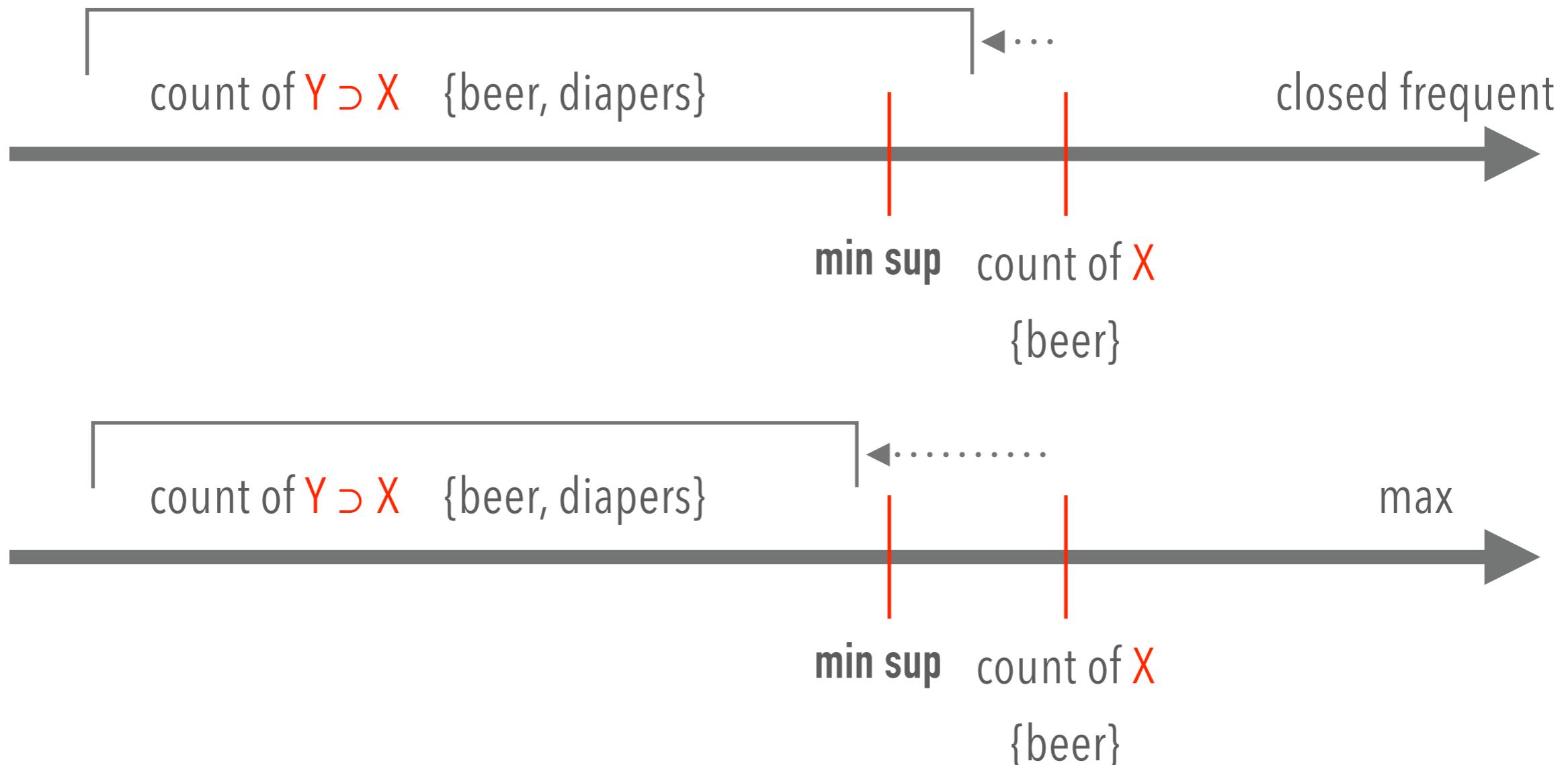


An itemset  $X$  is a closed frequent itemset in set  $D$  if  $X$  is both **closed** and **frequent** in  $D$

The definition does **not** say anything about  $Y$



An itemset  $X$  is a maximal frequent itemset (or **max-itemset**) in a data set  $D$  if  $X$  is frequent, and there exists no super-itemset  $Y$  such that  $Y \supset X$  and  $Y$  is frequent in  $D$ .



what is a key difference between the set of all **max patterns** and the set of all **closed frequent** patterns?

TID	Items
10	a, c, d, e, f
20	a, b, e
30	c, e, f
40	a, c, d, f
50	c, e, f

Pei, Jian, Jiawei Han, and Runying Mao. "CLOSET: An Efficient Algorithm for Mining Frequent Closed Itemsets." In ACM SIGMOD workshop on research issues in data mining and knowledge discovery, vol. 4, no. 2, pp. 21-30. 2000.

## CLOSET: CLOSED PATTERN MINING

**Flist:** list of all frequent items in support descending order

→  $f\_list: \langle c:4, e:4, f:4, a:3, d:2 \rangle$

**Divide** search space

Patterns having **d**

Patterns having **a** but no **d**, etc.

Find frequent closed pattern recursively

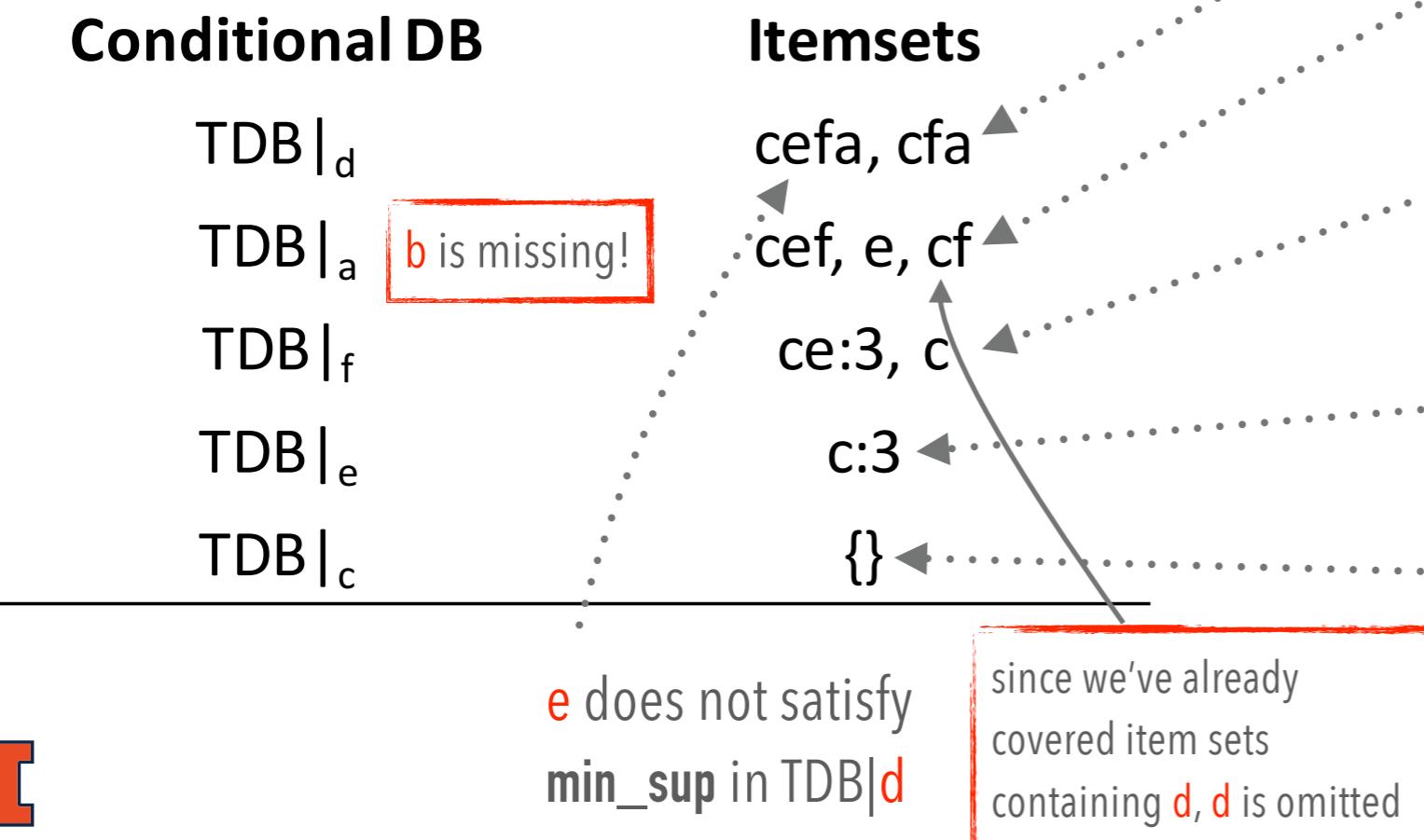
item set merging

Every transaction having **d** also has  $cfa \rightarrow cfad$  is a frequent closed pattern

## CLOSET: DIVIDE SEARCH SPACE

TID	Items	
10	a, c, d, e, f	b does not satisfy min_sup
20	a, <u>b</u> , e	
30	c, e, f	d has the lowest count satisfying min_sup
40	a, c, d, f	
50	c, e, f	

f\_list:<c:4, e:4, f:4, a:3, d:2>



All possible frequent closed patterns can be divided into 5 non-overlapping search spaces based on f\_list, i.e. **x**-conditional database ( $TDB|x$ )

$TDB|_d$ : The ones containing **d**

$TDB|_a$ : The ones containing **a** but no **d**

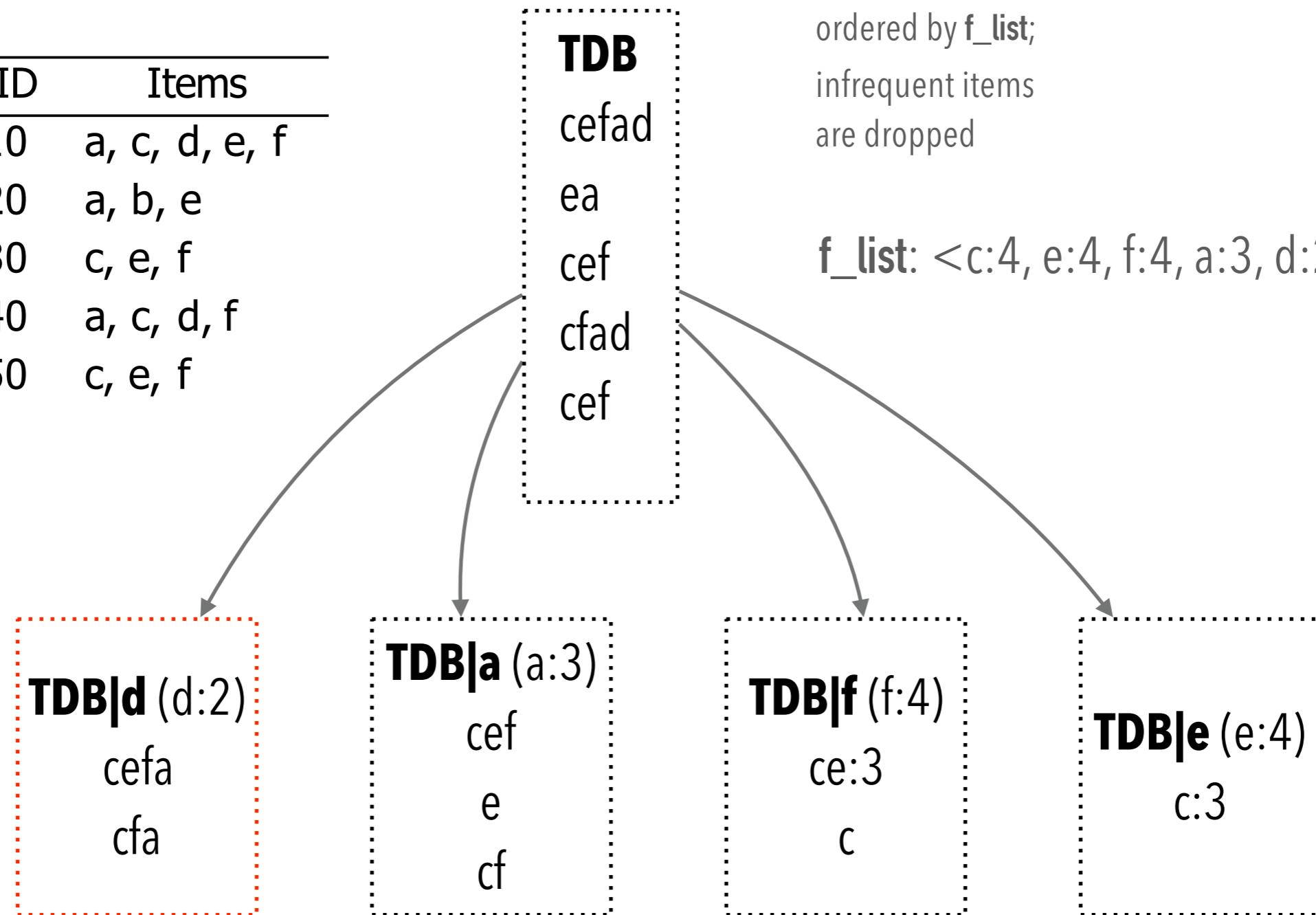
$TDB|_f$ : The ones containing **f** but no **a** nor **d**

$TDB|_e$ : The ones containing **e** but no **f**, **a** nor **d**

$TDB|_c$ : The ones containing only **c**

# find closed patterns containing d

TID	Items
10	a, c, d, e, f
20	a, b, e
30	c, e, f
40	a, c, d, f
50	c, e, f



$\langle c:2, f:2, a:2 \rangle$

F.C.P: cfad

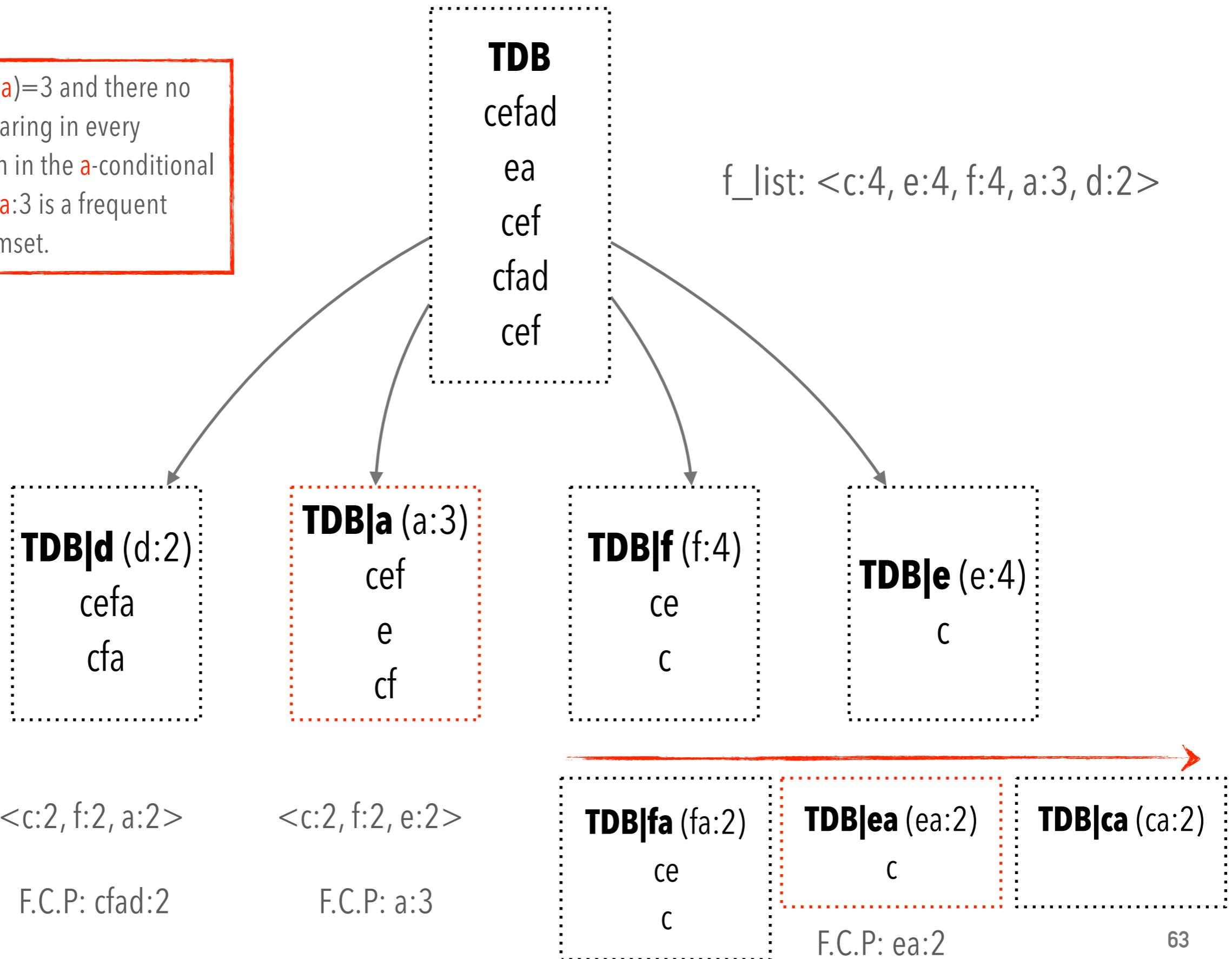
ordered by `f_list`;  
infrequent items  
are dropped

`f_list`:  $\langle c:4, e:4, f:4, a:3, d:2 \rangle$

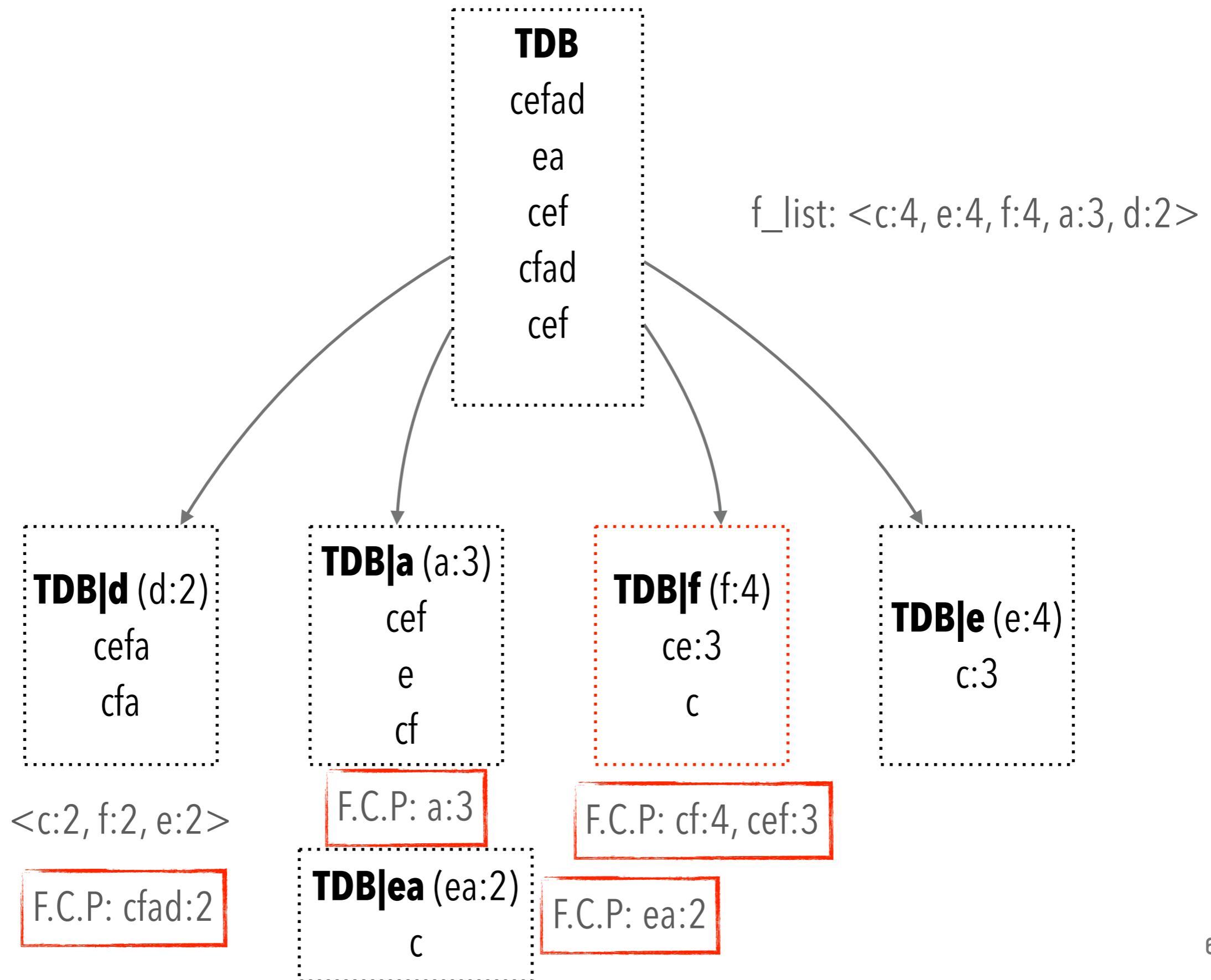
The support of  $d=2$ ; Items c, f and a appear twice respectively in  $TDB|d$ . That is every transaction containing d also contains c, f and a. Moreover e is infrequent since it appears only once in  $TDB|d$ . Therefore, cfad:2 is a frequent closed itemset. Since cfad covers all frequent itemsets containing d, the mining of  $TDB|d$  is complete.

# find closed patterns containing a but not d

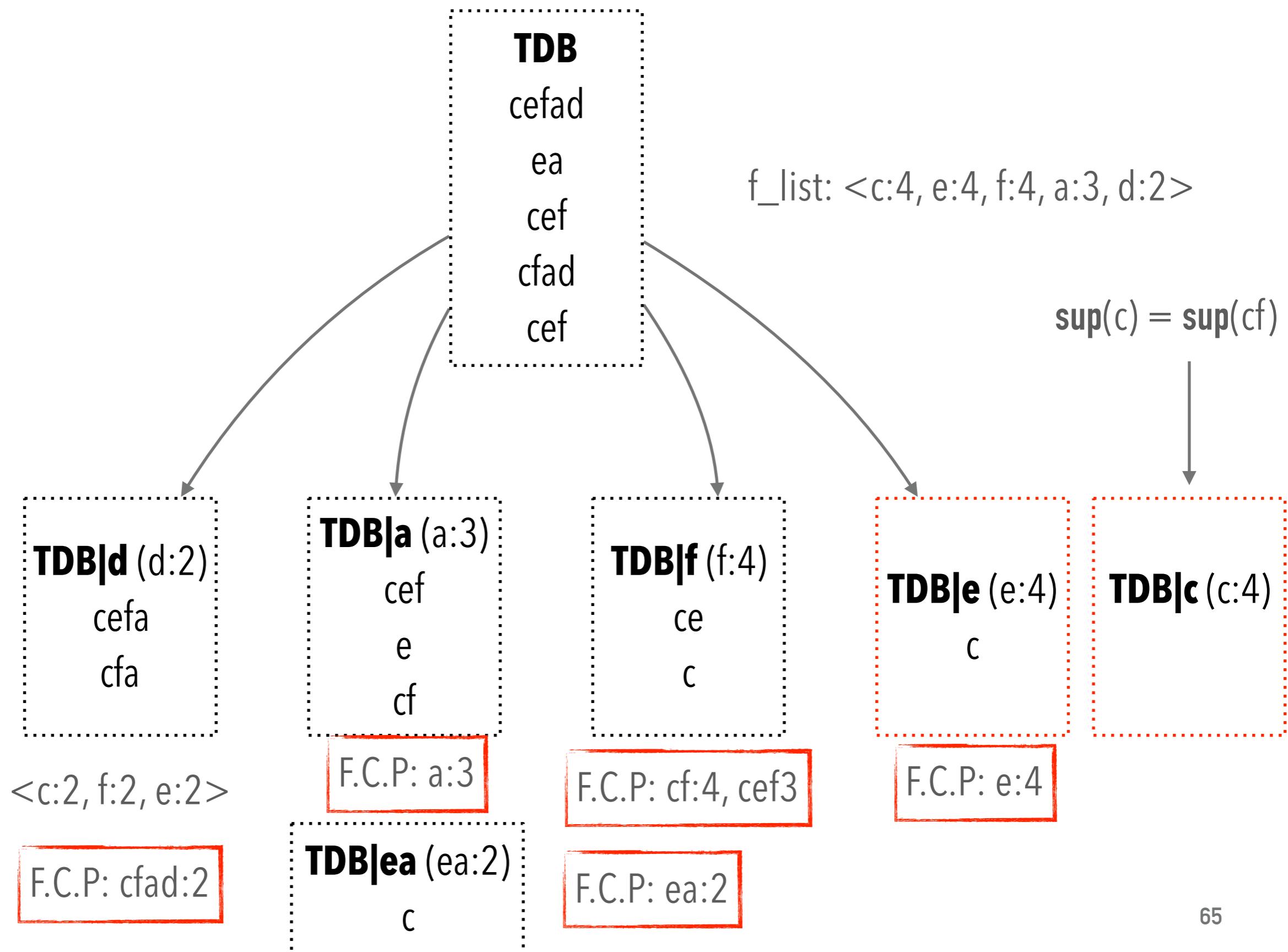
Since  $\text{sup}(a)=3$  and there no item appearing in every transaction in the  $a$ -conditional database,  $a:3$  is a frequent closed itemset.



# find closed patterns containing f but not a,d



# find closed patterns containing e but **not** a,f,d



# MAX MINER

.....  
1st scan: find frequent items

potential  
max  
patterns

A, B, C, D, E

2nd scan: find support for

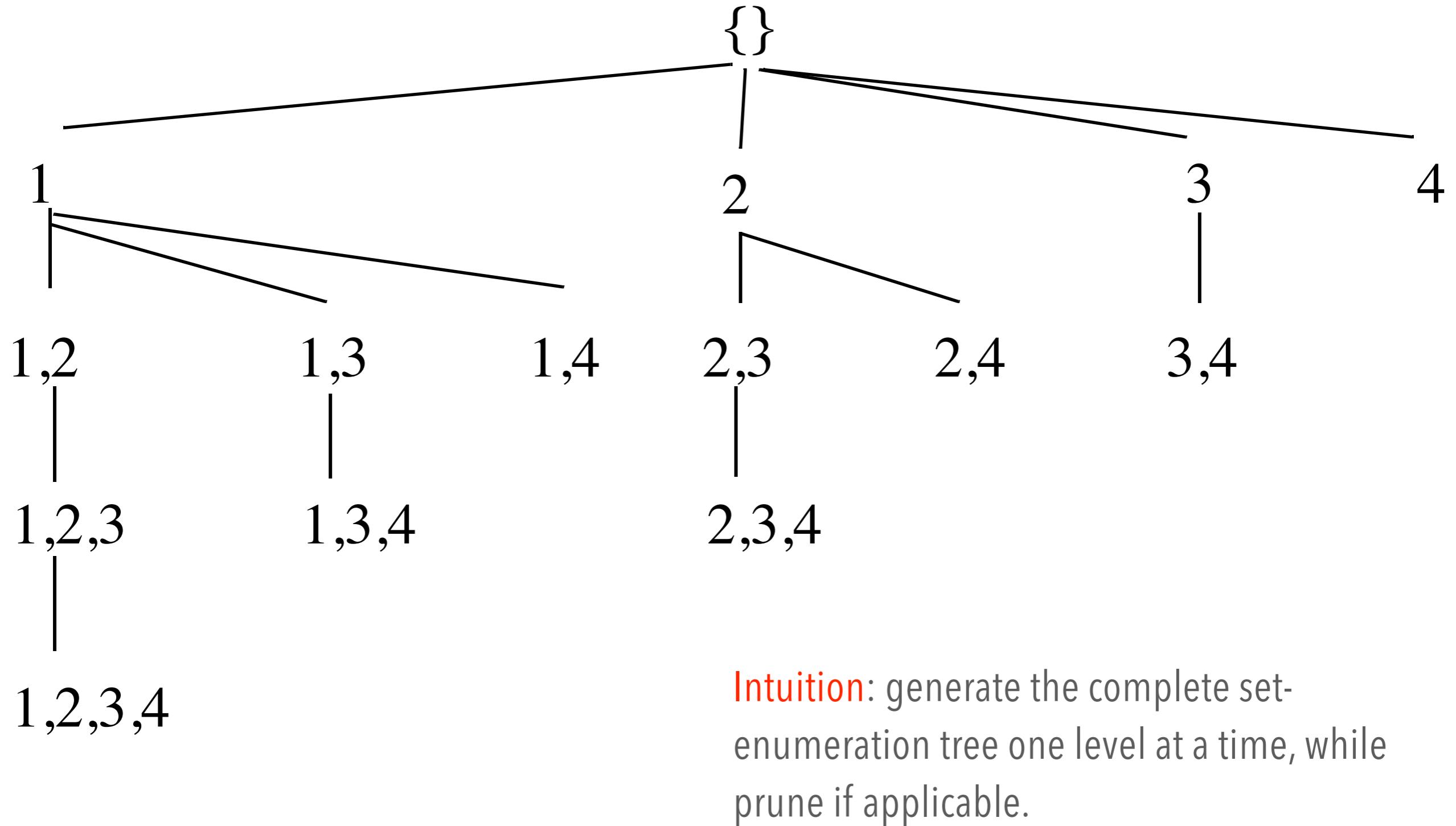
AB, AC, AD, AE, **ABCDE**

BC, BD, BE, **BCDE**

CD, CE, **CDE**, DE

Since BCDE is a max-pattern,  
no need to check BCD, BDE,  
CDE in later scan

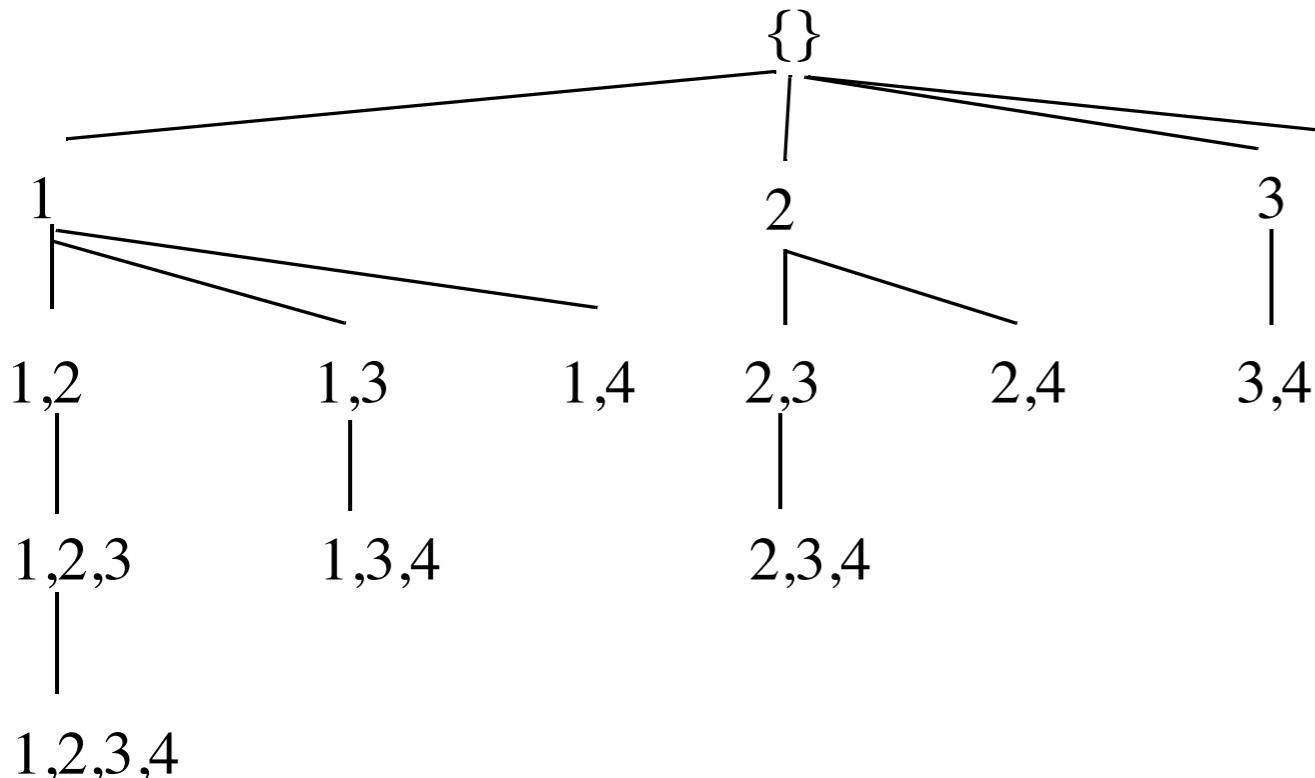
Bayardo Jr, Roberto J. "Efficiently mining long patterns from databases." In ACM Sigmod Record, vol. 27, no. 2, pp. 85-93. ACM, 1998.



**Intuition:** generate the complete set-enumeration tree one level at a time, while prune if applicable.

# MAXMINER

.....



Initially, generate one node  $N=(1234)$ , where head of  $N$  is  $h(N)=\Phi$  and tail of  $N$  is  $t(N)=\{1,2,3,4\}$ .

Consider expanding  $N$ ,

4 If  $h(N) \cup t(N)$  is frequent, do not expand  $N$  and prune the whole sub-tree.

Else: do local pruning. If for some  $i \in t(N)$ ,  $h(N) \cup \{i\}$  is not frequent, remove  $i$  from  $t(N)$  before expanding  $N$ :

Apply global pruning: When a max pattern is identified, prune all nodes across sub-tree where  $h(N) \cup t(N)$  is a sub-set of the pattern

# Apriori

$C_2$	
Itemset	
{I1, I2}	
{I1, I3}	
{I1, I4}	
{I1, I5}	
{I2, I3}	
{I2, I4}	
{I2, I5}	
{I3, I4}	
{I3, I5}	
{I4, I5}	

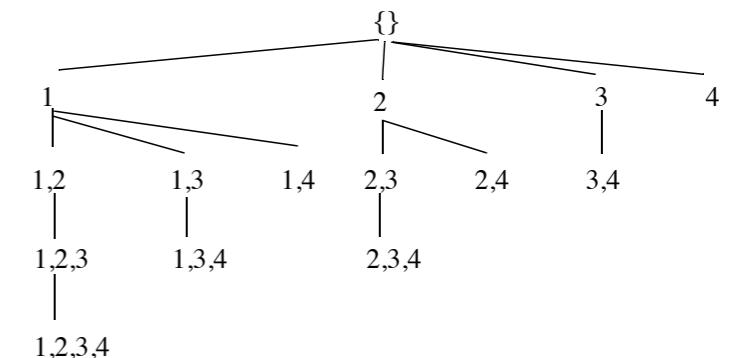
Generate  $C_2$  candidates from  $L_1$  →

$C_2$	
Itemset	Sup. count
{I1, I2}	4
{I1, I3}	4
{I1, I4}	1
{I1, I5}	2
{I2, I3}	4
{I2, I4}	2
{I2, I5}	2
{I3, I4}	0
{I3, I5}	1
{I4, I5}	0

Scan  $D$  for count of each candidate →

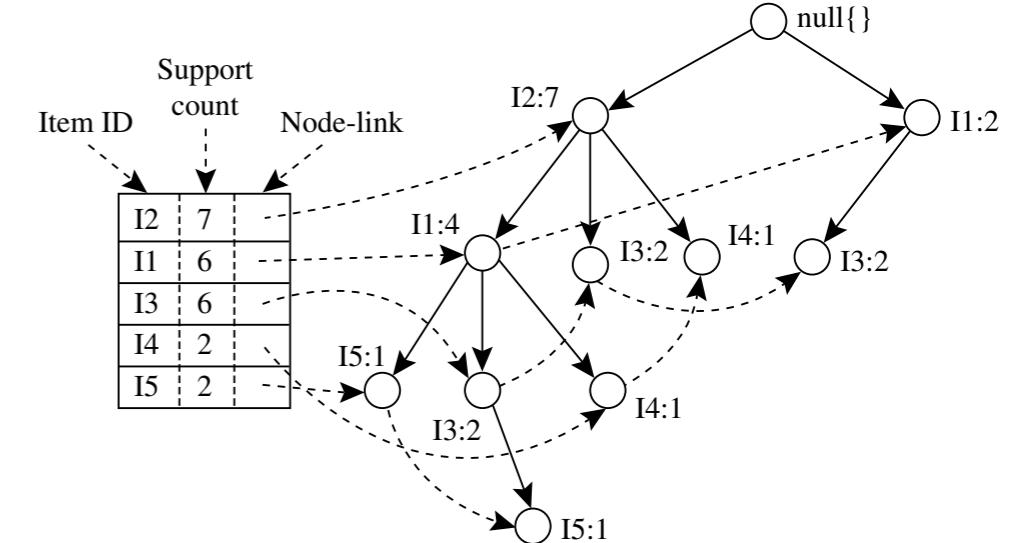
Conditional DB	Itemsets
$TDB _d$	cefa, cfa
$TDB _a$	cef, e, cf
$TDB _f$	ce:3, c
$TDB _e$	c:3
$TDB _c$	{}

## Max & Closed



# SUMMARY

Frequent Pattern Mining Methods



ECLAT

FPGrowth

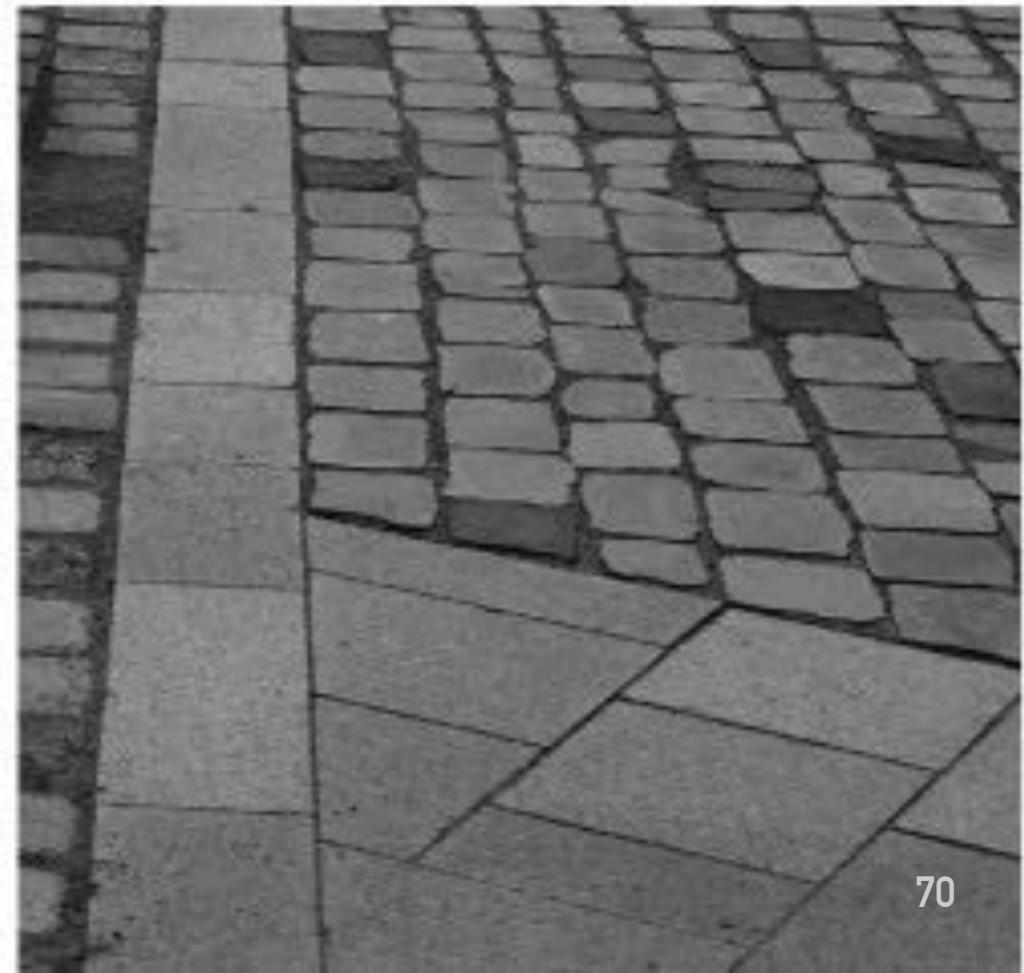
# PATTERN EVALUATION

---

Basic Concepts

Methods

Summary



# CORRELATION

	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

play basketball → eat cereal [40%, 66.7%] is misleading

This is because the overall percentage of students eating cereal is  $75\% > 66.7\%$ .

Why is this happening?  $P(A|B)$

$$lift(A, B) = \frac{P(A|B)}{P(A)P(B)}$$

play basketball → not eat cereal [20%, 33.3%] is more accurate, although with lower support and confidence

**lift:** Measure of dependent or correlated events:

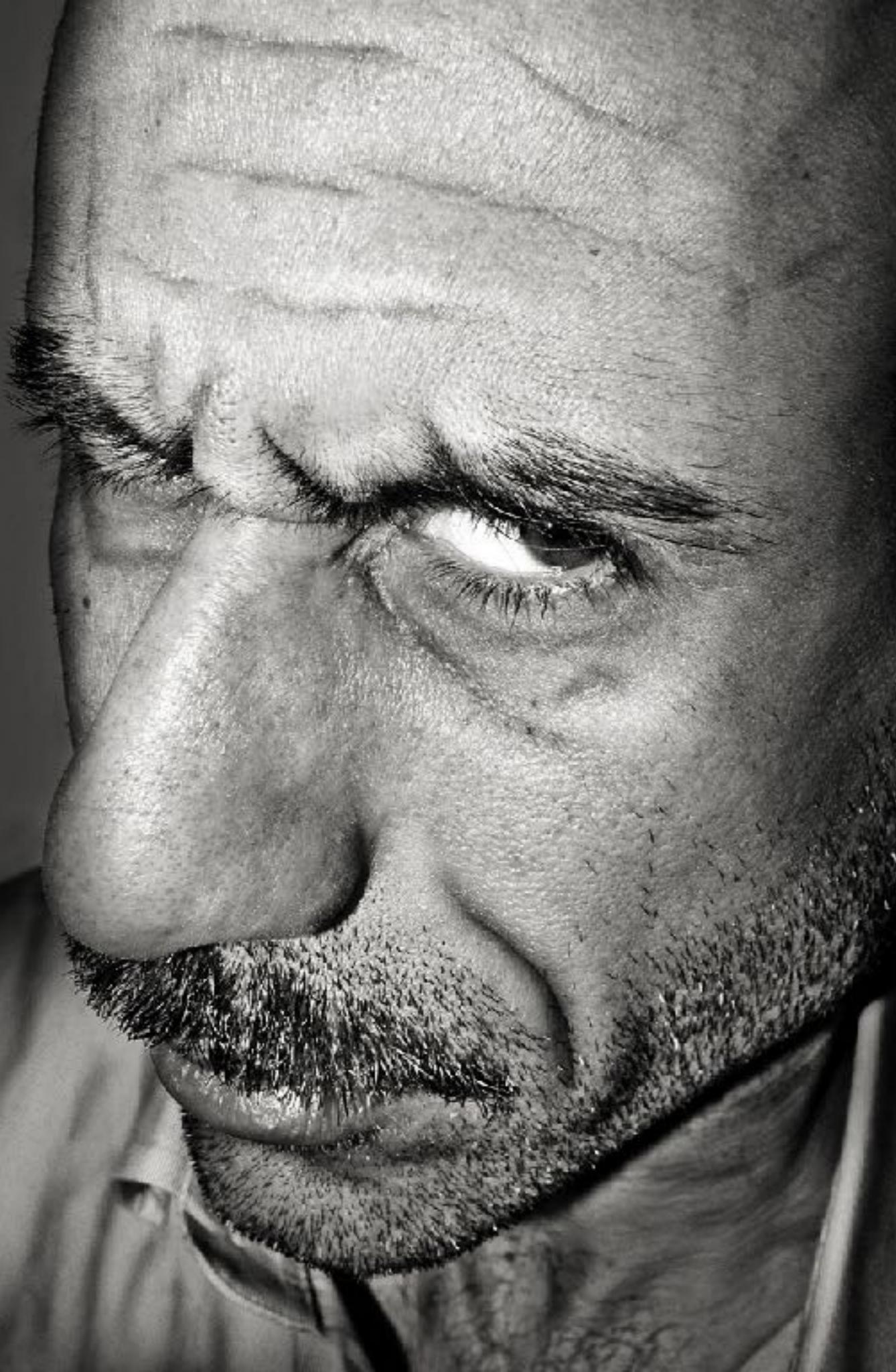
$$lift(B, C) = \frac{2000/5000}{3000/5000 \times 3750/5000}$$

$$lift(B, \bar{C}) = \frac{1000/5000}{3000/5000 \times 1250/5000}$$

## Exercise!

	<i>game</i>	$\overline{\text{game}}$	$\Sigma_{\text{row}}$
<i>video</i>	4000	3500	7500
$\overline{\text{video}}$	2000	500	2500
$\Sigma_{\text{col}}$	6000	4000	10,000

games → videos?



## WHAT COULD BE A PROBLEM WITH LIFT?

.....

$$lift(A, B) = \frac{P(A \cup B)}{P(A)P(B)}$$

video and games are **negatively** correlated!

	game	$\bar{\text{game}}$	$\Sigma_{\text{row}}$
video	4000 (4500)	3500 (3000)	7500
$\bar{\text{video}}$	2000 (1500)	500 (1000)	2500
$\Sigma_{\text{col}}$	6000	4000	10,000

expected values

$$\begin{aligned}\chi^2 &= \sum \frac{(observed - expected)^2}{expected} = \frac{(4000 - 4500)^2}{4500} + \frac{(3500 - 3000)^2}{3000} \\ &\quad + \frac{(2000 - 1500)^2}{1500} + \frac{(500 - 1000)^2}{1000} = 555.6. \quad p < 0.001\end{aligned}$$

$$all\_conf(A, B) = \frac{sup(A \cup B)}{\max\{sup(A), sup(B)\}} = \min\{P(A|B), P(B|A)\}$$

$$max\_conf(A, B) = \max\{P(A|B), P(B|A)\}.$$

# Null invariant measures

$$Kulc(A, B) = \frac{1}{2}(P(A|B) + P(B|A)).$$

$$\begin{aligned} cosine(A, B) &= \frac{P(A \cup B)}{\sqrt{P(A) \times P(B)}} = \frac{sup(A \cup B)}{\sqrt{sup(A) \times sup(B)}} \\ &= \sqrt{P(A|B) \times P(B|A)}. \end{aligned}$$

$$all\_conf(A, B) = \frac{sup(A \cup B)}{\max\{sup(A), sup(B)\}} = \min\{P(A|B), P(B|A)\}$$

$$lift(A, B) = \frac{P(A \cup B)}{P(A)P(B)}$$

$$max\_conf(A, B) = \max\{P(A|B), P(B|A)\}.$$

$$Kulc(A, B) = \frac{1}{2}(P(A|B) + P(B|A)).$$

$$\begin{aligned} cosine(A, B) &= \frac{P(A \cup B)}{\sqrt{P(A) \times P(B)}} = \frac{sup(A \cup B)}{\sqrt{sup(A) \times sup(B)}} \\ &= \sqrt{P(A|B) \times P(B|A)}. \end{aligned}$$

	<i>milk</i>	$\overline{milk}$	$\Sigma_{row}$
<i>coffee</i>	<i>mc</i>	$\overline{mc}$	<i>c</i>
$\overline{coffee}$	$m\bar{c}$	$\overline{m\bar{c}}$	$\bar{c}$
$\Sigma_{col}$	<i>m</i>	$\overline{m}$	$\Sigma$

# canonical contingency table

<i>Data</i>										
<i>Set</i>	<i>mc</i>	$\overline{mc}$	$m\bar{c}$	$\overline{m\bar{c}}$	$\chi^2$	<i>lift</i>	<i>all_conf.</i>	<i>max_conf.</i>	<i>Kulc.</i>	<i>cosine</i>
$D_1$	10,000	1000	1000	100,000	90557	9.26	0.91	0.91	0.91	0.91
$D_2$	10,000	1000	1000	100	0	1	0.91	0.91	0.91	0.91
$D_3$	100	1000	1000	100,000	670	8.44	0.09	0.09	0.09	0.09
$D_4$	1000	1000	1000	100,000	24740	25.75	0.5	0.5	0.5	0.5
$D_5$	1000	100	10,000	100,000	8173	9.18	0.09	0.91	0.5	0.29
$D_6$	1000	10	100,000	100,000	965	1.97	0.01	0.99	0.5	0.10

$$all\_conf(A, B) = \frac{sup(A \cup B)}{\max\{sup(A), sup(B)\}} = \min\{P(A|B), P(B|A)\}$$

$$lift(A, B) = \frac{P(A \cup B)}{P(A)P(B)}$$

$$max\_conf(A, B) = \max\{P(A|B), P(B|A)\}.$$

$$Kulc(A, B) = \frac{1}{2}(P(A|B) + P(B|A)).$$

$$\begin{aligned} cosine(A, B) &= \frac{P(A \cup B)}{\sqrt{P(A) \times P(B)}} = \frac{sup(A \cup B)}{\sqrt{sup(A) \times sup(B)}} \\ &= \sqrt{P(A|B) \times P(B|A)}. \end{aligned}$$

**imbalance ratio**

$$IR(A, B) = \frac{|sup(A) - sup(B)|}{sup(A) + sup(B) - sup(A \cup B)}$$

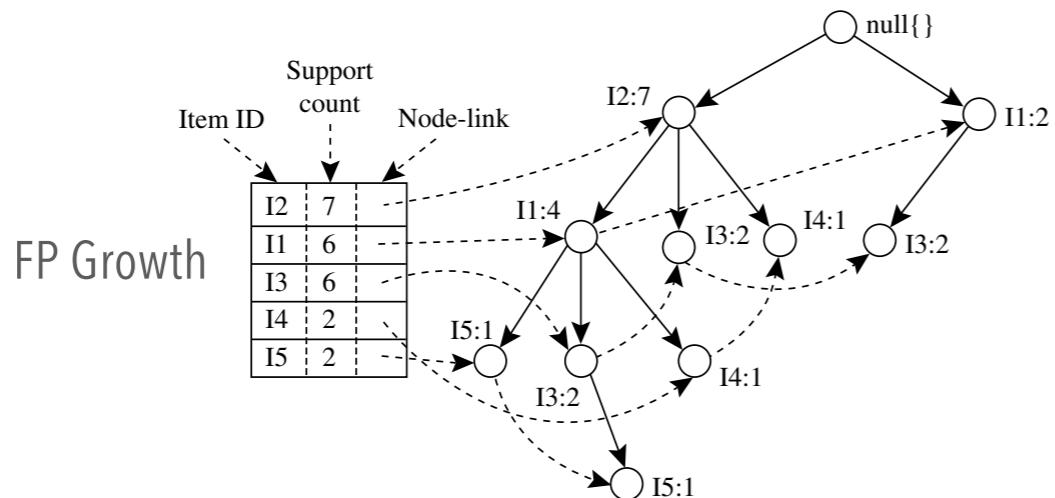
# canonical contingency table

*Data*

Set	$mc$	$\bar{mc}$	$m\bar{c}$	$\bar{m}\bar{c}$	$\chi^2$	$lift$	$all\_conf.$	$max\_conf.$	$Kulc.$	$cosine$
$D_1$	10,000	1000	1000	100,000	90557	9.26	0.91	0.91	0.91	0.91
$D_2$	10,000	1000	1000	100	0	1	0.91	0.91	0.91	0.91
$D_3$	100	1000	1000	100,000	670	8.44	0.09	0.09	0.09	0.09
$D_4$	1000	1000	1000	100,000	24740	25.75	0.5	0.5	0.5	0.5
$D_5$	1000	100	10,000	100,000	8173	9.18	0.09	0.91	0.5	0.29
$D_6$	1000	10	100,000	100,000	965	1.97	0.01	0.99	0.5	0.10

$|R=0.89$

$|R=0.99$



# FREQUENT

## ITEMSET

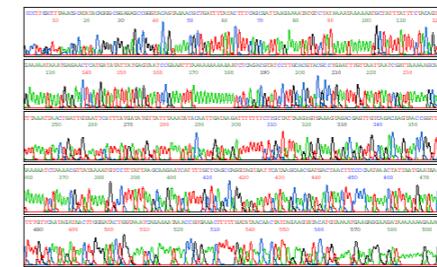
## METHODS

Evaluation

lift,  $\chi^2$ , min, max,  
Kulc, cosine

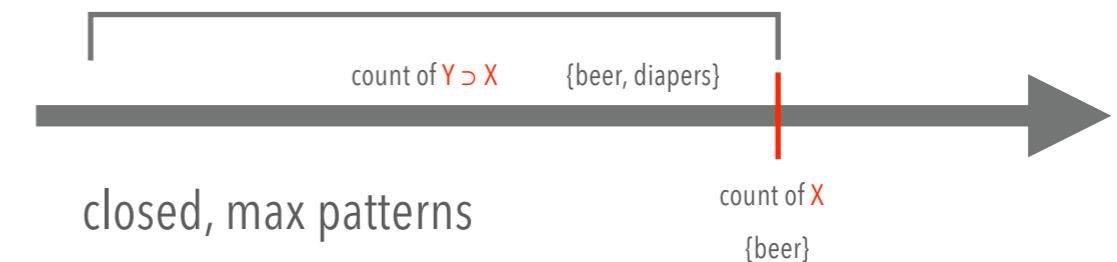
null invariant

Summary



a pattern (a set of items, subsequences, substructures, etc.) that occurs frequently in a data set

itemset, support, frequent, confidence



**Apriori pruning principle:** If there is any itemset which is infrequent, its superset should not be generated or tested!

challenges

hashing,  
partitioning,  
sampling