# COMP9318: Data Warehousing and Data Mining

## — L6: Association Rule Mining —

- Problem definition and preliminaries
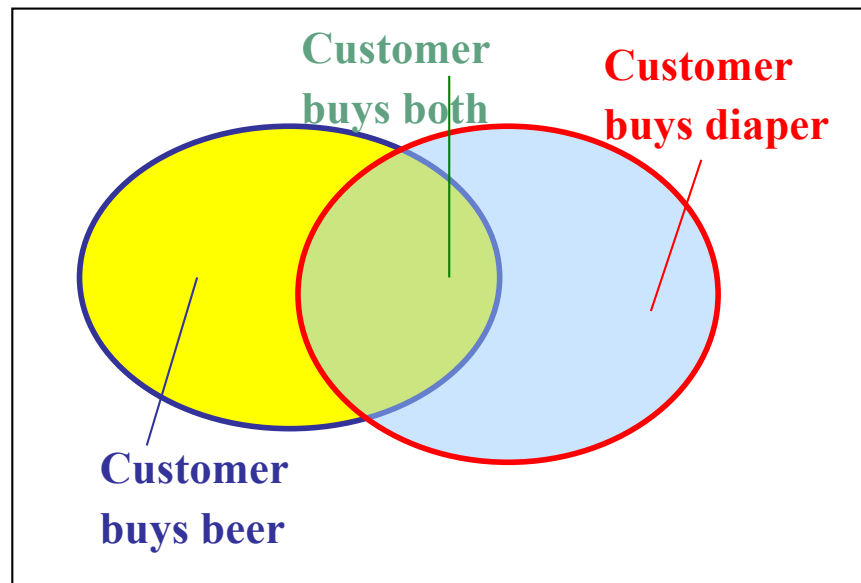
# What Is Association Mining?

- Association rule mining:
  - Finding frequent patterns, associations, correlations, or causal structures among sets of items or objects in transaction databases, relational databases, and other information repositories.
  - Frequent pattern: pattern (set of items, sequence, etc.) that occurs frequently in a database [AIS93]
- Motivation: finding regularities in data
  - What products were often purchased together? — Beer and diapers?!
  - What are the subsequent purchases after buying a PC?
  - What kinds of DNA are sensitive to this new drug?
  - Can we automatically classify web documents?

# Why Is Frequent Pattern or Assoiciation Mining an Essential Task in Data Mining?

- Foundation for many essential data mining tasks
  - Association, correlation, causality
  - Sequential patterns, temporal or cyclic association, partial periodicity, spatial and multimedia association
  - Associative classification, cluster analysis, iceberg cube, fascicles (semantic data compression)
- Broad applications
  - Basket data analysis, cross-marketing, catalog design, sale campaign analysis
  - **Web log** (click stream) **analysis**, DNA sequence analysis, etc. | c.f., google's spelling suggestion |

# Basic Concepts: Frequent Patterns and Association Rules

| Transaction-id | Items bought |
|---|---|
| 10 | { A, B, C } |
| 20 | { A, C } |
| 30 | { A, D } |
| 40 | { B, E, F } |

- Itemset $X=\{x_1, ..., x_k\}$
  - **Shorthand**: $x_1\ x_2\ ...\ x_k$
- Find all the rules $X \rightarrow Y$ with min confidence and support
  - support, $s$, probability that a transaction contains $X \cup Y$
  - confidence, $c$, conditional probability that a transaction having $X$ also contains $Y$.

Customer buys both

Customer buys diaper

Customer buys beer

Let $min\_support = 50\%$,
$min\_conf = 70\%$:
$sup(AC) = 2$

frequent itemset

association rule

$A \rightarrow C$ (50%, 66.7%)
$C \rightarrow A$ (50%, 100%)

# Mining Association Rules—an Example

| Transaction-id | Items bought |
|:---:|:---:|
| 10 | A, B, C |
| 20 | A, C |
| 30 | A, D |
| 40 | B, E, F |

Min. support 50%
Min. confidence 50%

| Frequent pattern | Support |
|:---:|:---:|
| {A} | 75% |
| {B} | 50% |
| {C} | 50% |
| {A, C} | 50% |

For rule $A \rightarrow C$:

support = support($\{A\} \cup \{C\}$) = 50%

confidence = support($\{A\} \cup \{C\}$)/support($\{A\}$) = 66.6%

major computation challenge: calculate the support of itemsets
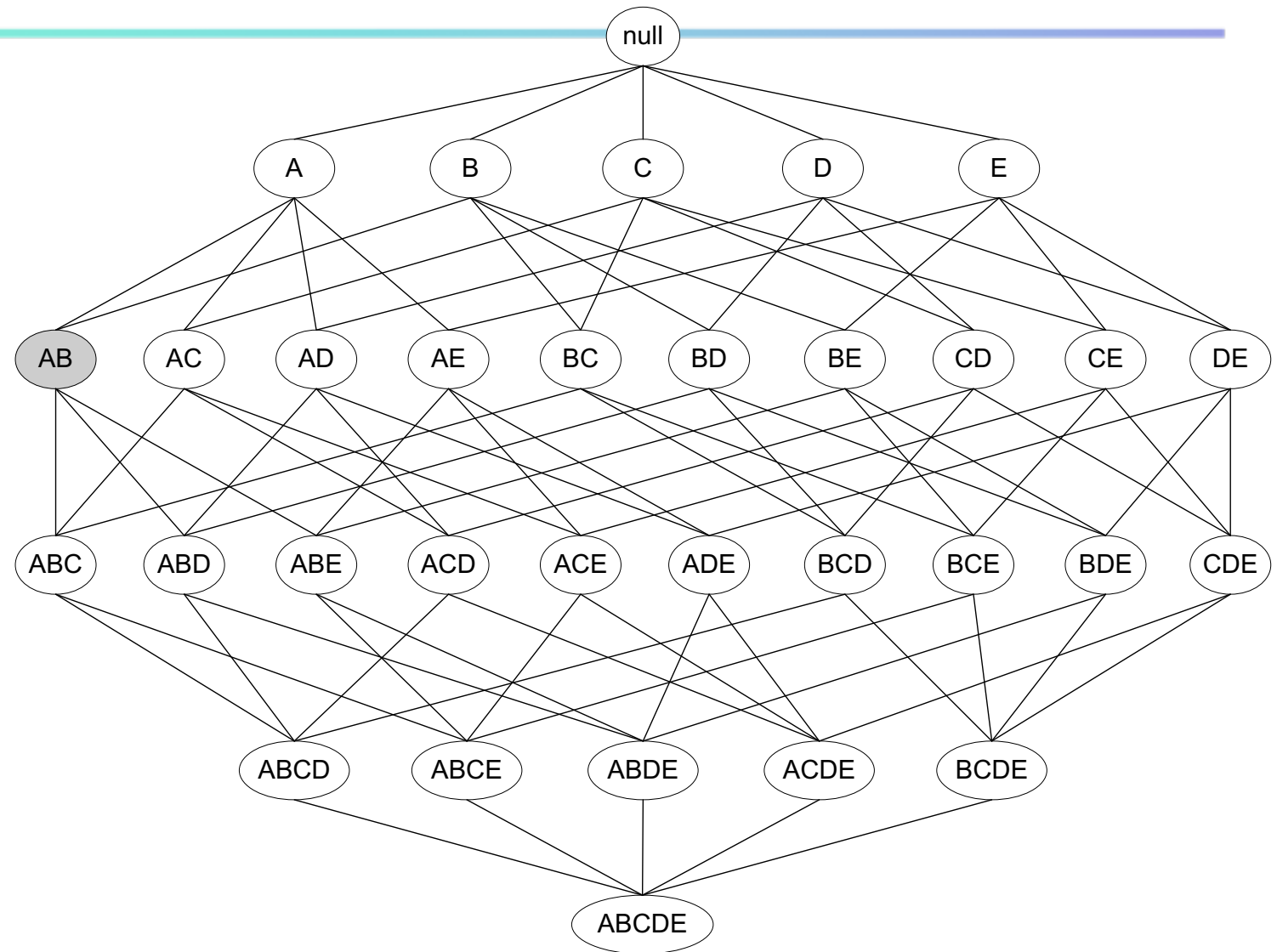← The *frequent itemset mining* problem

- Algorithms for scalable mining of (single-dimensional Boolean) association rules in transactional databases

# Association Rule Mining Algorithms

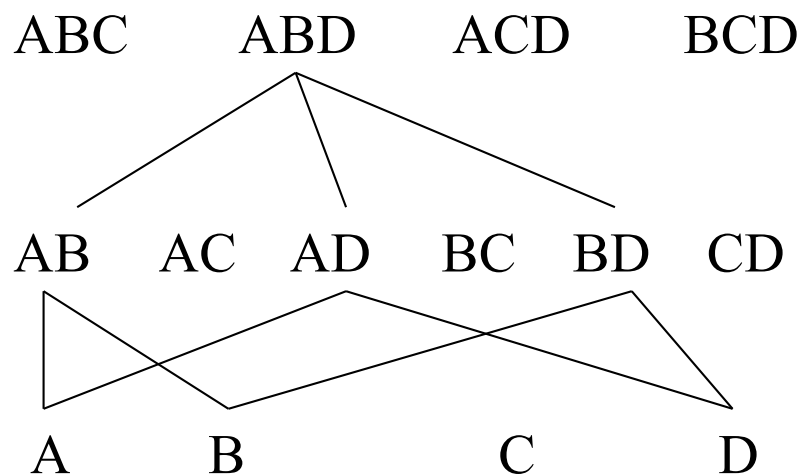Candidate Generation & Verification

- **Naïve algorithm**
    - Enumerate all possible itemsets and check their support against *min_sup*
    - Generate all association rules and check their confidence against *min_conf*
- **The Apriori property**
    - Apriori Algorithm
    - FP-growth Algorithm

# All Candidate Itemsets for {A, B, C, D, E}

# Apriori Property

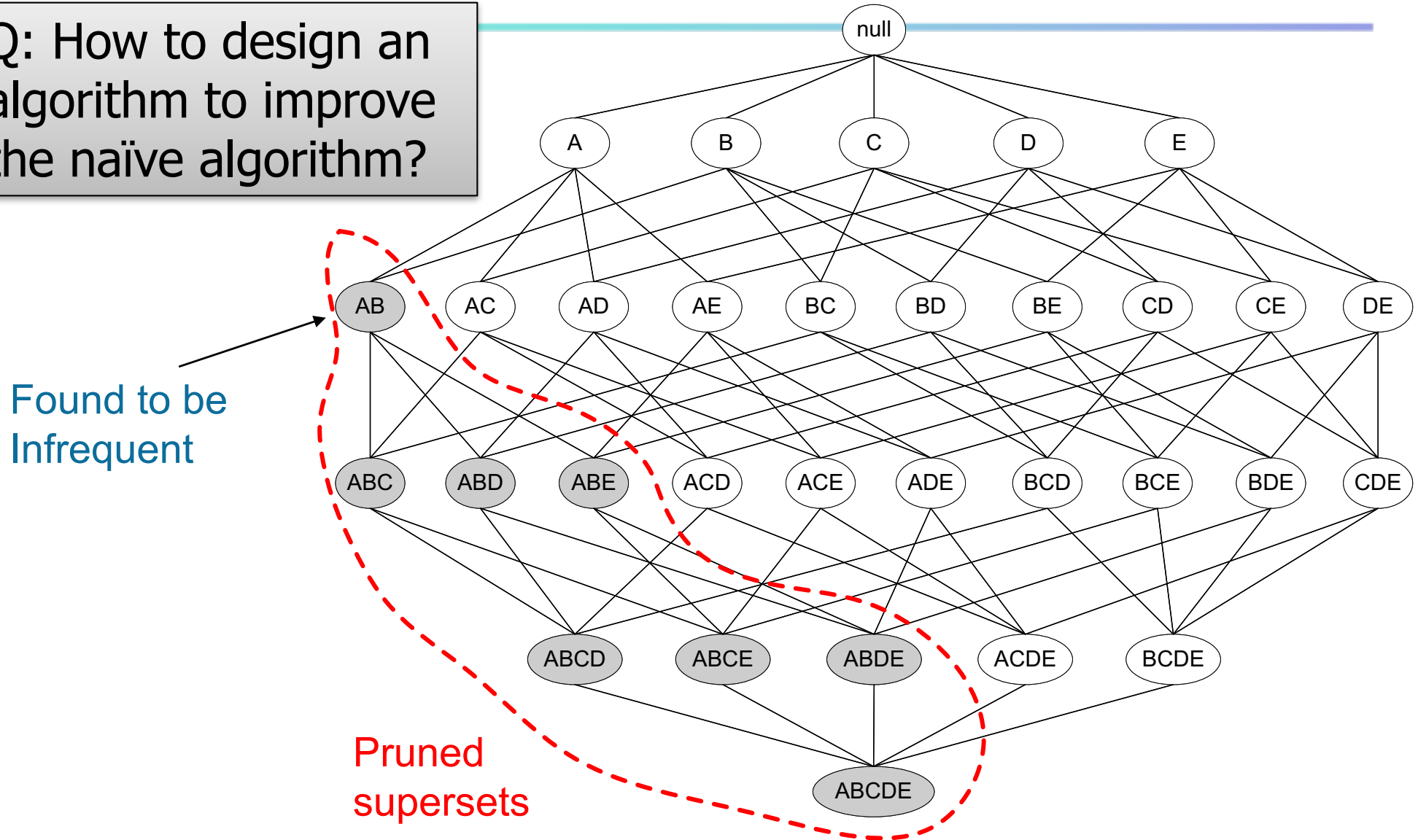- A *frequent* (used to be called *large*) *itemset* is an itemset whose support is ≥ min_sup.

- Apriori property (downward closure): any <span style="color:red">sub</span>sets of a frequent itemset are also frequent itemsets

- Aka the <span style="color:red">anti-monotone</span> property of support

ABC    ABD    ACD    BCD

AB   AC   AD   BC   BD   CD

A        B        C        D

"any <span style="color:red">super</span>sets of an <span style="color:red">in</span>frequent itemset are also <span style="color:red">in</span>frequent itemsets"

# Illustrating Apriori Principle

Q: How to design an algorithm to improve the naïve algorithm?



Found to be Infrequent

Pruned supersets

# Apriori: A Candidate Generation-and-test Approach

- **<u>Apriori pruning principle</u>: If there is any itemset which is infrequent, its superset should not be generated/tested!**

- Algorithm [Agrawal & Srikant 1994]
    1. $C_k \leftarrow$ Perform level-wise candidate generation (from singleton itemsets)
    2. $L_k \leftarrow$ Verify $C_k$ against $L_k$
    3. $C_{k+1} \leftarrow$ generated from $L_k$
    4. Goto 2 if $C_{k+1}$ is not empty

# The Apriori Algorithm

- **Pseudo-code**:

$C_k$: Candidate itemset of size $k$

$L_k$ : frequent itemset of size $k$

```
L₁ = {frequent items};
for (k = 1; Lₖ !=∅; k++) do begin
    C₍ₖ₊₁₎ = candidates generated from Lₖ;
    for each transaction t in database do begin
        increment the count of all candidates in C₍ₖ₊₁₎
        that are contained in t
    end
    L₍ₖ₊₁₎ = candidates in C₍ₖ₊₁₎ with min_support
end
return ∪ₖLₖ;
```

# The Apriori Algorithm—An Example

minsup = 50%

Database TDB

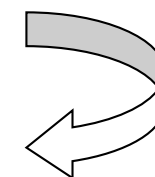| Tid | Items |
|-----|-------|
| 10 | A, C, D |
| 20 | B, C, E |
| 30 | A, B, C, E |
| 40 | B, E |

$1^{st}$ scan

$C_1$

| Itemset | sup |
|---------|-----|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {D} | 1 |
| {E} | 3 |

$L_1$

| Itemset | sup |
|---------|-----|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {E} | 3 |

$C_2$

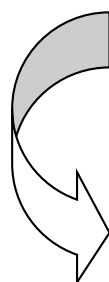| Itemset | sup |
|---------|-----|
| {A, B} | 1 |
| {A, C} | 2 |
| {A, E} | 1 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$2^{nd}$ scan

$C_2$

| Itemset |
|---------|
| {A, B} |
| {A, C} |
| {A, E} |
| {B, C} |
| {B, E} |
| {C, E} |

$L_2$

| Itemset | sup |
|---------|-----|
| {A, C} | 2 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$C_3$

| Itemset |
|---------|
| {B, C, E} |

$3^{rd}$ scan

$L_3$

| Itemset | sup |
|---------|-----|
| {B, C, E} | 2 |

# Important Details of Apriori

1.  How to generate candidates?

    - Step 1: self-joining $L_k$  (what's the join condition? why?)

    - Step 2: pruning

2.  How to count supports of candidates?

Example of Candidate-generation

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

- 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 $L_3$

- $C_4 = \{abcd\}$

# Generating Candidates in SQL

- Suppose the items in $L_{k-1}$ are listed in an order
- Step 1: self-joining $L_{k-1}$

> insert into $C_k$
>
> select $p.item_1, p.item_2, ..., p.item_{k-1}, q.item_{k-1}$
>
> from $L_{k-1}\ p,\ L_{k-1}\ q$
>
> where $p.item_1=q.item_1, ..., p.item_{k-2}=q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$

- Step 2: pruning

> forall *itemsets c in $C_k$* do
>
>      forall *(k-1)-subsets s of c* do
>
>          **if** *(s is not in $L_{k-1}$)* **then delete** *c* **from** $C_k$

# Derive rules from frequent itemsets

- Frequent itemsets != association rules
- One more step is required to find association rules
- For each frequent itemset $X$,

  For each <span style="color:red">proper nonempty</span> subset $A$ of $X$,

  - Let $B = X - A$

  - $A \rightarrow B$ is an association rule if

    - Confidence $(A \rightarrow B) \geq$ *min_conf*,
      where support $(A \rightarrow B) =$ support $(AB)$, and
      confidence $(A \rightarrow B) =$ support $(AB)$ / support $(A)$

# Example – deriving rules from frequent itemsets

- Suppose 234 is frequent, with supp=50%

    - Proper nonempty subsets: 23, 24, 34, 2, 3, 4, with supp=50%, 50%, 75%, 75%, 75%, 75% respectively

    - These generate these association rules:

        - 23 => 4,      confidence=100%
        - 24 => 3,      confidence=100%
        - 34 => 2,      confidence=67%      $= (N* 50\%)/(N*75\%)$
        - 2 => 34,      confidence=67%
        - 3 => 24,      confidence=67%
        - 4 => 23,      confidence=67%
        - All rules have support = 50%

Q: is there any optimization (e.g., pruning) for this step?

# Deriving rules

- To recap, in order to obtain A → B, we need to have Support(AB) and Support(A)
- This step is not as time-consuming as frequent itemsets generation
  - Why?
- It's also easy to speedup using techniques such as parallel processing.
  - How?
- Do we really need candidate generation for deriving association rules?
  - Frequent-Pattern Growth (FP-Tree)

# Bottleneck of Frequent-pattern Mining

- Multiple database scans are <span style="color:red">costly</span>

- Mining long patterns needs many passes of scanning and generates lots of candidates

  - To find frequent itemset $i_1 i_2 \ldots i_{100}$

    - \# of scans: <span style="color:red">100</span>
    - \# of Candidates: $\binom{100}{1} + \binom{100}{2} + \ldots + \binom{100}{100} = 2^{100} - 1$

- Bottleneck: candidate-generation-and-test

*Can we avoid candidate generation **altogether**?*

- **FP-growth**

# No Pain, No Gain

|         | Java | Lisp | Scheme | Python | Ruby |
|---------|------|------|--------|--------|------|
| Alice   | X    |      |        |        | X    |
| Bob     |      |      |        | X      | X    |
| Charlie | X    |      |        | X      | X    |
| Dora    |      | X    | X      |        |      |

minsup = 1

- Apriori:
  - L1 = {J, L, S, P, R}
  - C2 = all the $\binom{5}{2}$ combinations
    - Most of C2 do not contribute to the result
    - There is no way to tell because

# No Pain, No Gain

| | Java | Lisp | Scheme | Python | Ruby |
|---|---|---|---|---|---|
| Alice | X | | | | X |
| Bob | | | | X | X |
| Charlie | X | | | X | X |
| Dora | | X | X | | |

minsup = 1

**Ideas**:
- Keep the support set for each frequent itemset
- DFS

J ➔ JL?
J ➔ ???
Only need to look at support set for J

{A, C}

J

φ

# No Pain, No Gain
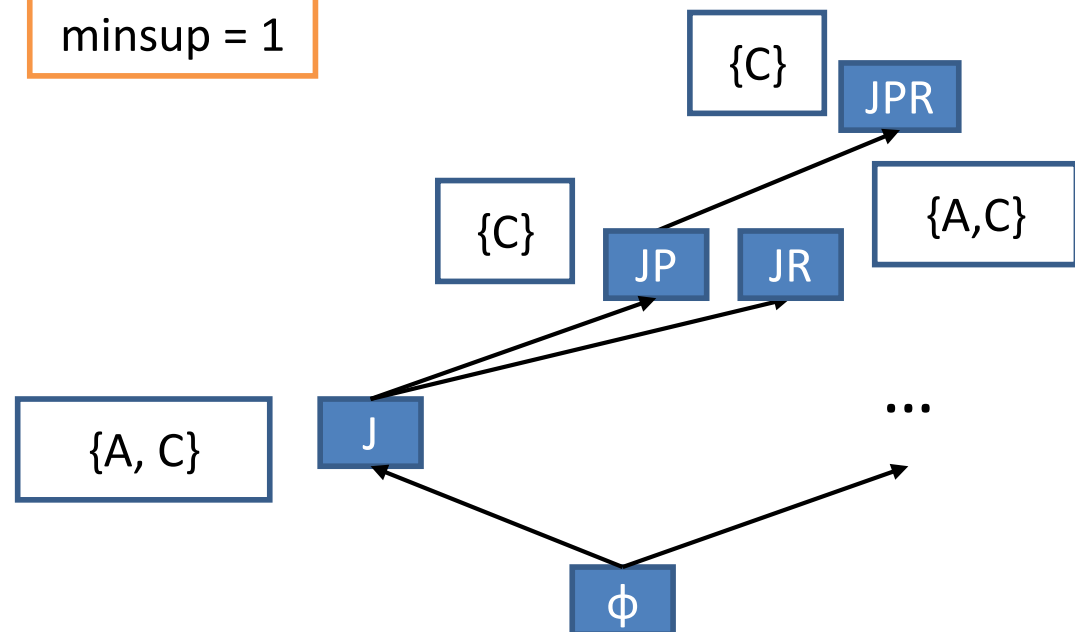
|          | Java | Lisp | Scheme | Python | Ruby |
|----------|------|------|--------|--------|------|
| Alice    | X    |      |        |        | X    |
| Bob      |      |      |        | X      | X    |
| Charlie  | X    |      |        | X      | X    |
| Dora     |      | X    | X      |        |      |

minsup = 1

**Ideas**:
- Keep the support set for each frequent itemset
- DFS

{C}

JPR

{A,C}

{C}

JP   JR

{A, C}

J

...

φ

# Notations and Invariants

- CondiditonalDB:
  - DB|p = {t $\in$ DB | t contains itemset p}
  - DB = DB|$\varnothing$ (i.e., conditioned on nothing)
  - Shorthand: DB|px = DB|(p $\cup$ x)
- SupportSet(p $\cup$ x, DB) = SupportSet(x, DB|p)
  - {x | x mod 6 = 0 $\wedge$ x $\in$ [100] } =

    {x | x mod 3 = 0 $\wedge$ x $\in$ even([100]) }
- A FP-tree is equivalent to a DB|p
  - One can be converted to another
  - Next, we illustrate the alg using conditionalDB

# FP-tree Essential Idea /1

- Recursive algorithm again!

- Freq**Itemsets**(DB|p):

  easy task, as only items (not itemsets) are needed

  all frequent itemsets in DB|p belong to one of the following categories:

  - X = Find**Locally**Frequent**Items**(DB|p)

  output { (x p) | x ∈ X }

  - Foreach $x$ in X

    - DB*|p$x$ = GetConditionalDB$^+$(DB*|p, $x$)

    -

    - Freq**Itemsets**(DB*|p$x$)

  | patterns ~ $x_i p$ |
  | patterns ~ ★$px_1$ |
  | patterns ~ ★$px_2$ |
  | patterns ~ ★$px_i$ |
  | patterns ~ ★$px_n$ |

  *obtained via recursion*

# No Pain, No Gain

DB|J

| | **Java** | **Lisp** | **Scheme** | **Python** | **Ruby** |
|---|---|---|---|---|---|
| Alice | X | | | | X |
| Charlie | X | | | X | X |

minsup = 1

- Freq**Itemsets**(DB|J):
  - {**P**, **R**} ← Find**Locally**Frequent**Items**(DB|J)
  - Output {JP, JR}
  - Get DB*|J**P**; Freq**Itemsets**(DB*|J**P**)
  - Get DB*|J**R**; Freq**Itemsets**(DB*|J**R**)
  - // Guaranteed no other frequent itemset in DB|J

# FP-tree Essential Idea /2

- **Freq*Itemsets*(DB|p):**
  - If boundary condition, then …
  - X = Find**Locally**Frequent**Items**(DB|p)
  - [optional] DB*|p = PruneDB(DB|p, X)

    | output { (x p) | x $\in$ X } |

  - Foreach x in X
    - DB*|px = GetConditionalDB[+](DB*|p, x)
    - [optional] if DB*|px is degenerated, then powerset(DB*|px)
    - Freq**Itemsets**(DB*|px)

Also output each item in X (appended with the conditional pattern)

Remove items not in X; potentially reduce # of transactions ($\varnothing$ or dup). Improves the efficiency.

Also gets rid of items already processed before x → *avoid duplicates*

# Lv 1 Recursion

- minsup = 3

| F C A M P |
|---|
| C B P |
| F C A M P |

DB*|P

DB*|M (sans P)

DB*|B (sans MP)

DB*|A (sans BMP)

DB*|C (sans ABMP)

DB*|F (sans CABMP)

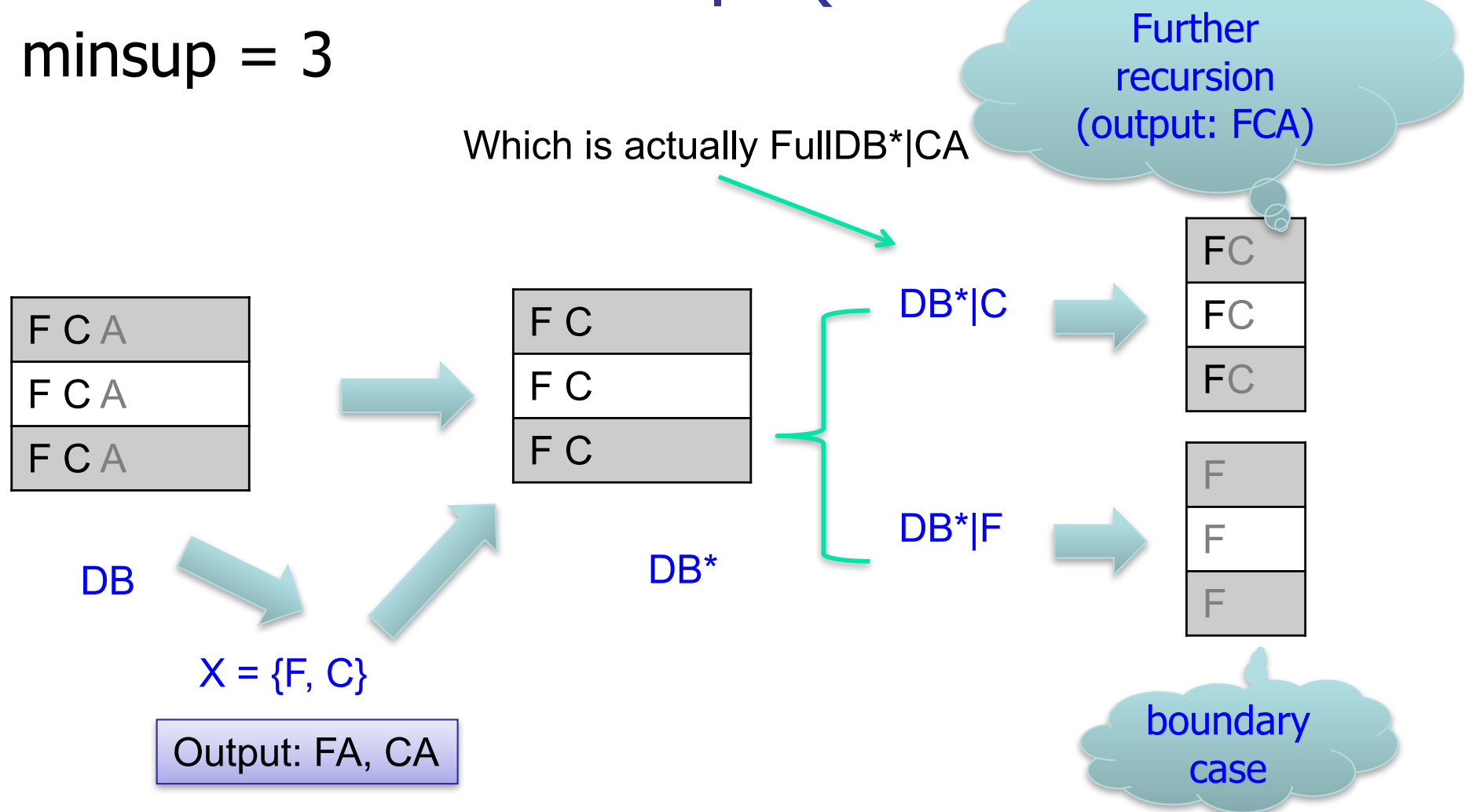| F C A D G I M P |
|---|
| A B C F L M O |
| B F H J O W |
| B C K S P |
| A F C E L P M N |

DB

| F C A M P |
|---|
| F C A B M |
| F B |
| C B P |
| F C A M P |

DB*

X = {F, C, A, B, M, P}

Output: F, C, A, B, M, P

| F C A |
|---|
| F C A |
| F C A |

# Lv 2 Recursion on DB*|P

- minsup = 3

Which is actually FullDB*|CP

| F C A M P |
|-----------|
| C B P     |
| F C A M P |

DB

X = {C}

Output: CP

| C |
|---|
| C |
| C |

DB*

DB*|C

| C |
|---|
| C |
| C |

Context = Lv 3 recursion on DB*|CP: DB has only empty sets or X = {} ➔ immediately returns

# Lv 2 Recursion on DB*|A (sans ...)

- minsup = 3

Further recursion (output: FCA)

Which is actually FullDB*|CA

| F C A |
|-------|
| F C A |
| F C A |

DB

X = {F, C}

Output: FA, CA

| F C |
|-----|
| F C |
| F C |

DB*

DB*|C

| FC |
|----|
| FC |
| FC |

DB*|F

| F |
|---|
| F |
| F |

boundary case

# Different Example:
## Lv 2 Recursion on DB*|P

Output: FAP

X = {F}

| F |
|---|
| F |

- minsup = 2

Which is actually FullDB*|AP

| F C A |
|---|
| F C |
| F A |

DB*|A

| F C |
|---|
| F |

DB*|C

| F |
|---|
| F |

DB*|F

|   |
|---|
|   |
|   |

| F C A M P |
|---|
| F C B P |
| F A P |

DB*

DB

X = {F, C, A}

Output: FP, CP, AP

# I will give you back the FP-tree

- An FP-tree tree of DB consists of:
  - A fixed **order** among items in DB
  - A prefix, **threaded** tree of **sorted** transactions in DB
  - Header table: (item, freq, ptr)
- When used in the algorithm, the input DB is always pruned (c.f., PruneDB())
  - Remove infequent items
  - Remove infrequent items in every transaction

# FP-tree Example

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

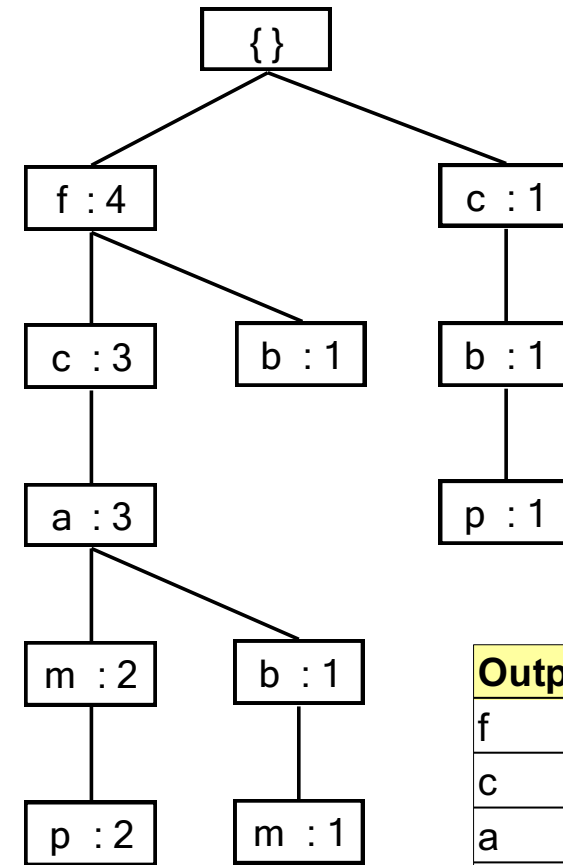| TID | Items bought | (ordered) frequent items |
|-----|--------------|--------------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o, w} | {f, b} |
| 400 | {b, c, k, s, p} | {c, b, p} |
| 500 | {a, f, c, e, l,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

**Insert t₁ tree:**

{} → f : 1 → c : 1 → a : 1 → m : 1 → p : 1

**Insert t₂ tree:**

{} → f : 2 → c : 2 → a : 2 → (m : 1 → p : 1) and (b : 1 → m : 1)

| Item | freq | head |
|------|------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |
| p | 3 | |

**Insert all tᵢ tree:**
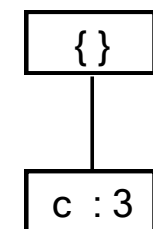
{} branches to f : 4 and c : 1

f : 4 → c : 3 and b : 1

c : 3 → a : 3

a : 3 → (m : 2 → p : 2) and (b : 1 → m : 1)

c : 1 → b : 1 → p : 1

| Output |
|--------|
| f |
| c |
| a |
| b |
| m |
| p |

Insert t₁

Insert t₂

Insert all tᵢ
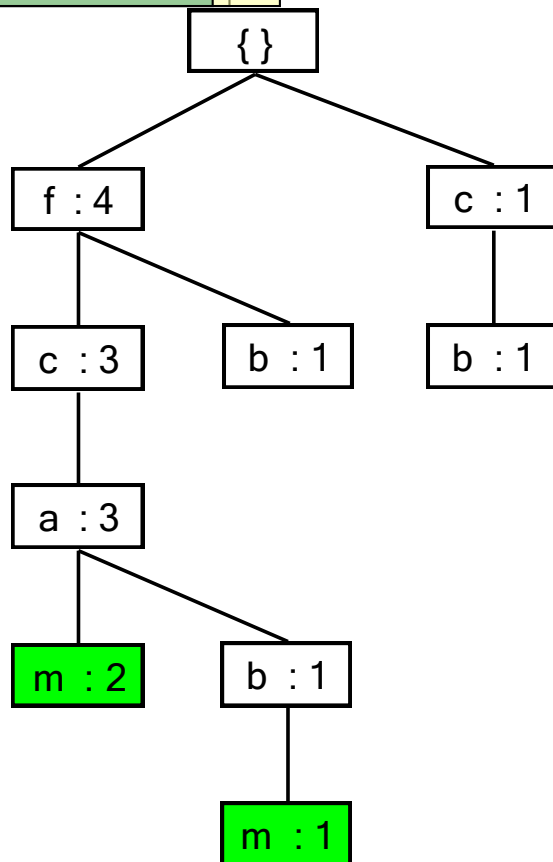
| TID | frequent items |
|-----|----------------|
| 100 | {f, c, a, m, p} |
| 200 | {f, c, a, b, m} |
| 300 | {f, b} |
| 400 | {c, b, p} |
| 500 | {f, c, a, m, p} |

**p's conditional pattern base**

| | |
|---|---|
| f  c  a    m : 2 | |
| c    b    : 1 | |

2 **3** 2 1 2

**Output**

pc

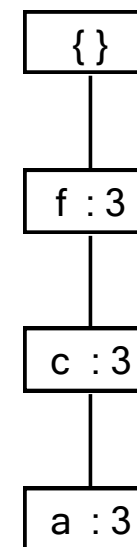| Item | freq | head |
|------|------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |
| p | 3 | |

{}

f : 4          c : 1

c : 3     b : 1     b : 1

a : 3          p : 1

m : 2     b : 1

p : 2     m : 1

**Cleaned p's conditional pattern base**

| C | :2 |
|---|----|
| C | :1 |

STOP      Header Table

{}

c : 3

| TID | frequent items |
|-----|----------------|
| 100 | {f, c, a, m, p} |
| 200 | {f, c, a, b, m} |
| 300 | {f, b} |
| 400 | {c, b, p} |
| 500 | {f, c, a, m, p} |

**m's conditional pattern base**

```
f c a   : 2
f c a b : 1
```

3 3 3 1

**Output**

| |
|---|
| mf |
| mc |
| ma |

| Item | freq | head |
|------|------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |

Tree:

{}
- f : 4
  - c : 3
    - a : 3
      - m : 2
      - b : 1
        - m : 1
  - b : 1
- c : 1
  - b : 1

gen_powerset

Header Table

{}
- f : 3
  - c : 3
    - a : 3

**Output**

| |
|---|
| mac |
| maf |
| mcf |
| macf |

**b's conditional pattern base**

| f c a : 1 |
| f     : 1 |
|   c   : 1 |

2 2 1

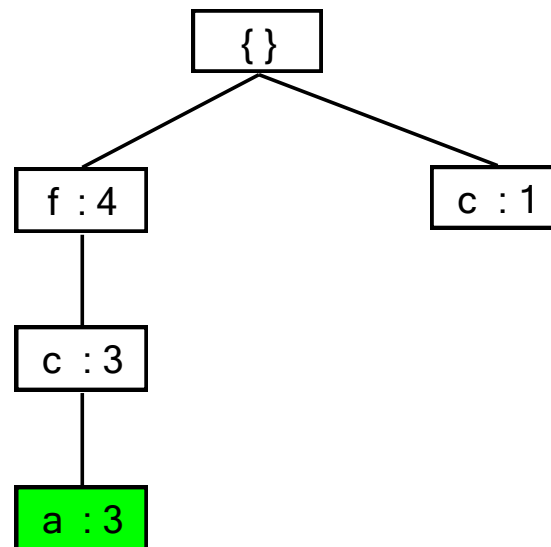| Item | freq | head |
|------|------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |

{}

f : 4

c : 1

c : 3

b : 1

b : 1

a : 3

b : 1

STOP

**a's conditional pattern base**

f c : 3

3  3

**Output**

af

ac

| Item | freq | head |
|------|------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |

{}

f : 4          c : 1

c : 3

a : 3

gen_powerset

Header Table

{}

f : 3

c : 3

**Output**

acf

**c's conditional pattern base**

`f : 3`

3

**Output**

cf

| Item | freq | head |
|------|------|------|
| f | 4 | |
| c | 4 | |

{ }
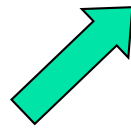
f : 4          c : 1

c : 3

{ }

f : 3

Header
Table

STOP

STOP

{ }

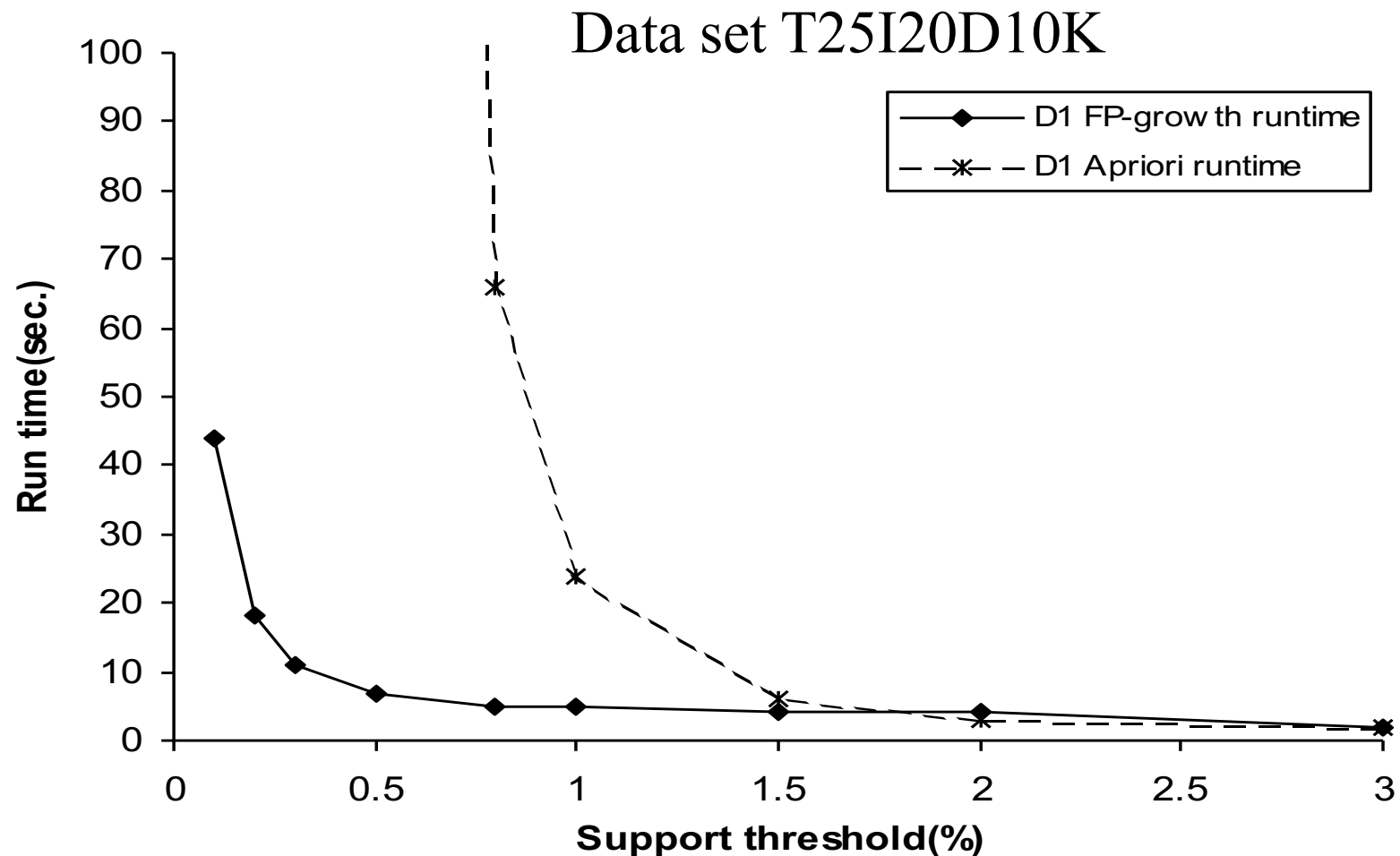| Item | freq | head |
|------|------|------|
| f | 4 | |

f : 4

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



Data set T25I20D10K

# Why Is FP-Growth the Winner?

- Divide-and-conquer:
  - decompose both the mining task and DB according to the frequent patterns obtained so far
  - leads 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 ops—counting local freq items and building sub FP-tree, no pattern search and matching