

EXPLORING FREQUENT ITEMSET GENERATION

- Present by Wei Li

- In lecture talk:
 - Association rule mining and Frequent itemset generation problem
 - ECLAT to generate frequent itemsets in a vertical database
- My project:
 - Implement ECLAT and two algorithms: FP-Growth and Apriori for frequent itemset generation
 - Compare three algorithms in two dataset
- Agenda:
 - Refresh the problem
 - Apriori
 - FP-Grwoth
 - Project Experiments and Results

THE PROBLEM

- Discover association rules
 - if-then relationship based on co-occurrence
 - {Bread, Butter} \rightarrow {Milk}
- Break Association Rules Mining into 2 tasks
 - **Generating frequent itemsets**
 - Generating rules from frequent itemsets
- Three sequential, exhaustive methods
 - **ECLAT**
 - **Apriori**
 - **FP-growth**



| TID | List of item IDs |
|-----|----------------------|
| T10 | Coke, Slice, Kurkure |
| T20 | Coke, Kurkure |
| T30 | Slice, pizza |

- Horizontal database

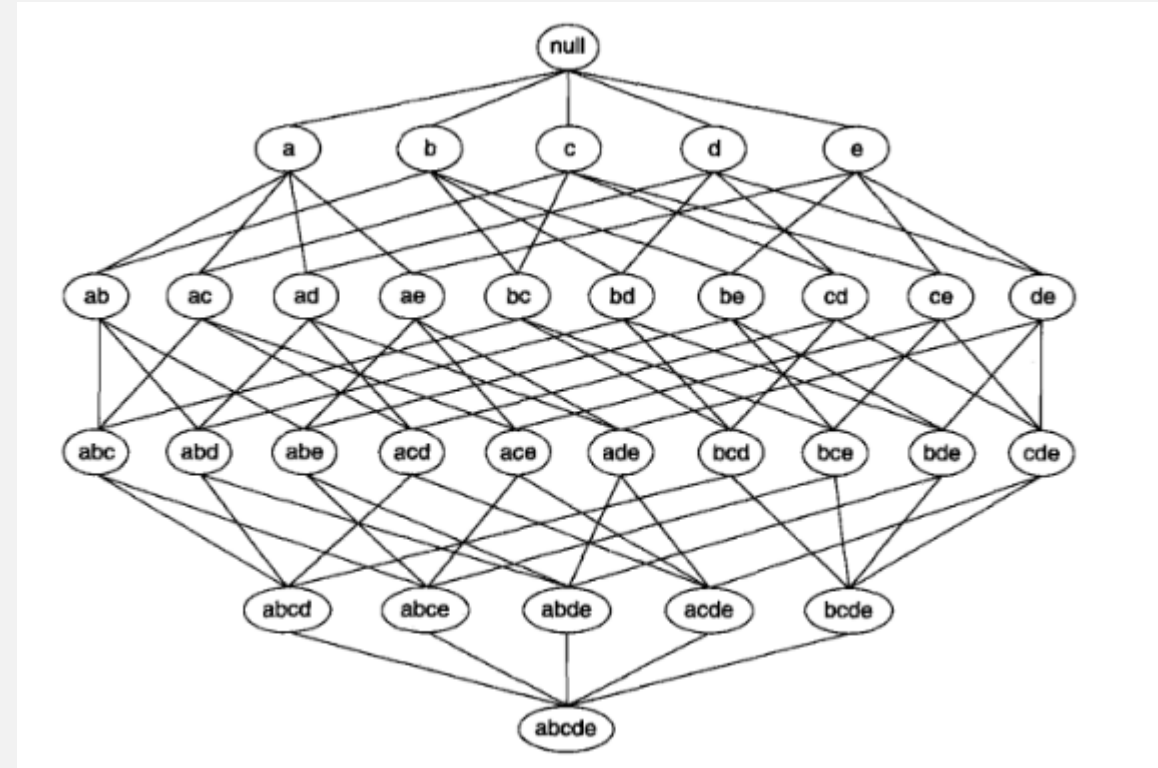
| Itemsets | TID_set |
|----------|----------|
| Coke | T10, T20 |
| Slice | T10, T30 |
| Kurkure | T10, T20 |
| Pizza | T30 |

- Vertical database

TERMS/ GENERATE ALL FREQUENT SET

- **Commodity set** $I = \{i_1, i_2, \dots, i_n\}$
- **transaction set** $\Omega = \{R_1, R_2, \dots, R_m\}$ where $R_i \subseteq I$.
- **Itemset** $X \subseteq I$ **K- Itemset:** $X \subseteq I$ and $|X| = k$
- **Support count :** $\delta(X) = |\{R_i | X \subseteq R_i \& R_i \in \Omega\}|$

- **Frequent Itemset:** itemset whose $\delta(X)$ is greater than a threshold (eg, 3)
- **Maximal frequent Itemset:** An frequent itemset is maximal frequent Itemset if none of its supersets are frequent.
- **Problem: frequent itemset generation**
 - **Input:** commodity set I , transaction set Ω , frequent threshold $minsup$
 - **Output:** all frequent itemset $\delta(X) > minsup$



APRIORI

- Apriori Property:
 - If an itemset is frequent, all of its subset are frequent
 - If an itemset is not frequent, all of its superset must not be frequent

Algorithm 1: Apriori(Ω , $minsup$)

Result: F

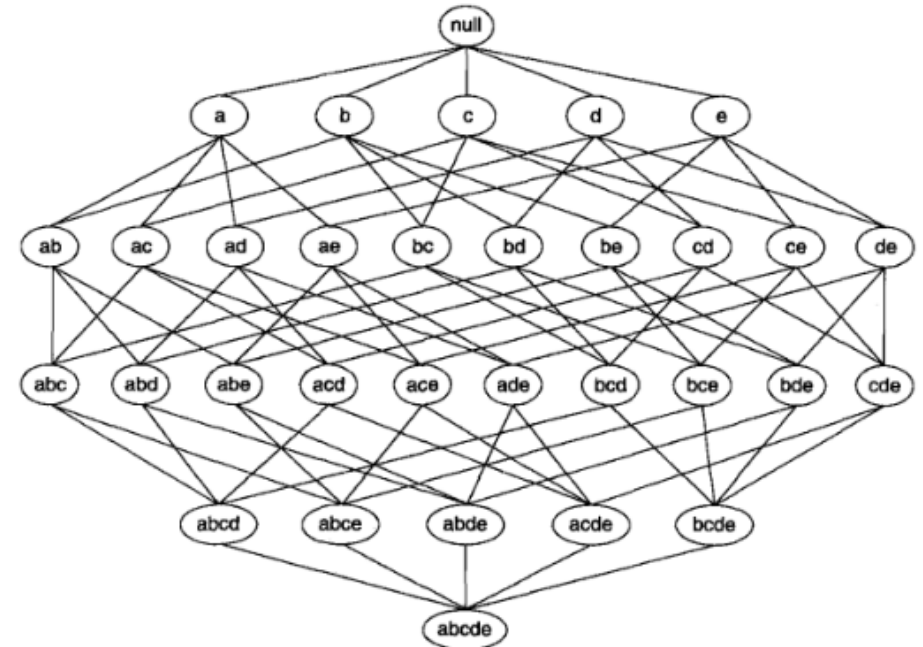
```
1 GLOBAL VAR, estimate;  
2  $k \leftarrow 1$ ;  
3  $F_1 \leftarrow$  all frequent 1-itemsets;  
4  $F \leftarrow \{(f, \delta(f)) | f \in F_1\}$ ;  
5 while  $F_k \neq \phi$  do  
6    $k \leftarrow k + 1$ ;  
7    $C_k \leftarrow candidate(F_{k-1}, \dots)$ ;  
8   for  $R_i \in \Omega$  do  
9      $D_i \leftarrow \{c | c \in C_k \text{ and } c \subseteq R_i\}$ ;  
10    for  $d \in D_i$  do  
11       $\delta(d) \leftarrow \delta(d) + 1$ ;  
12    end  
13  end  
14   $F_k = \{c | c \in C_k \text{ and } \delta(c) \geq minsup\}$ ;  
15   $F \leftarrow F \cup \{(f, \delta(f)) | f \in F_k\}$ ;  
16 end  
17 return  $F$ ;
```

K=1

K=2

K=3

K=4



FP-GROWTH

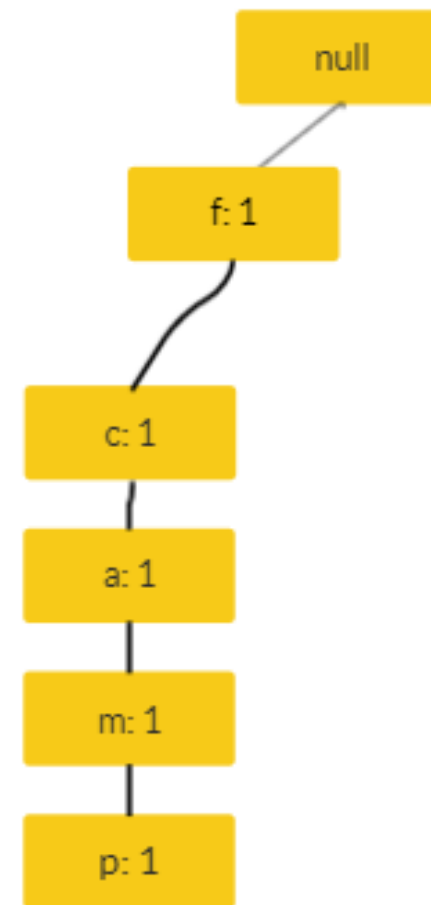
Minsup = 3

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

| item | freq | rank id |
|------|------|---------|
| f | 4 | 1 |
| c | 4 | 2 |
| a | 3 | 3 |
| b | 3 | 4 |
| m | 3 | 5 |
| p | 3 | 6 |

TID = 100

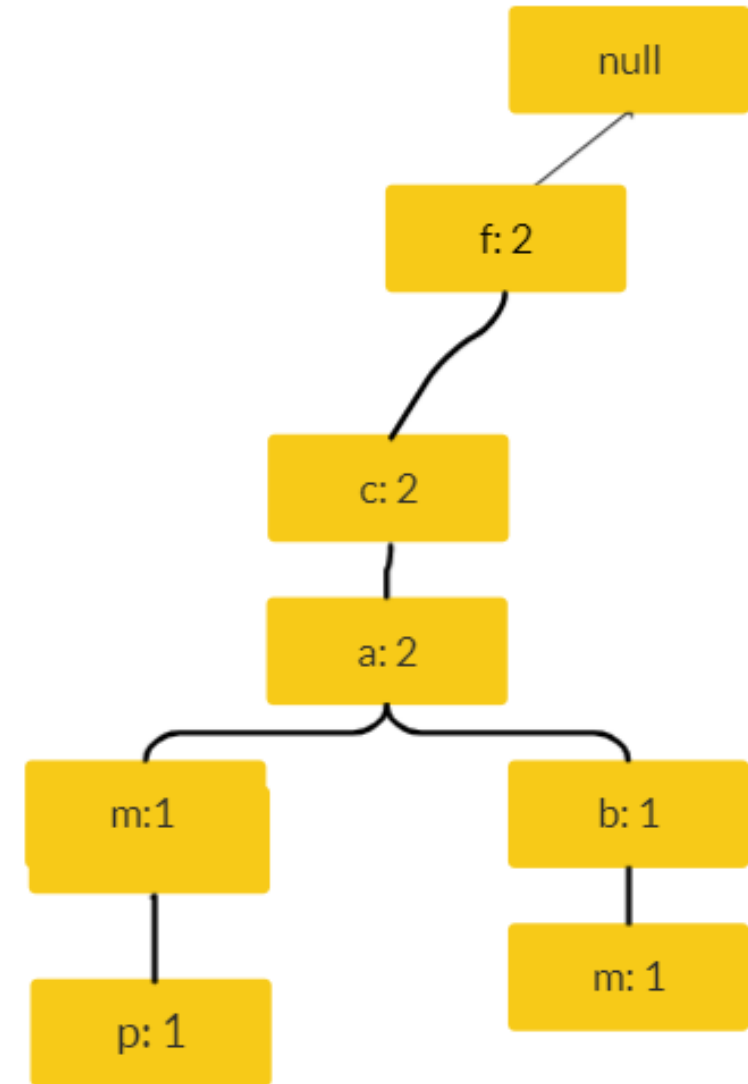
| TID | Item bought |
|-----|-----------------|
| 100 | {f, c, a, m, p} |
| 200 | {f, c, a, b, m} |
| 300 | {f, b} |
| 400 | {c, b, p} |
| 500 | {f, c, a, m, p} |



FP-GROWTH

TID = 200

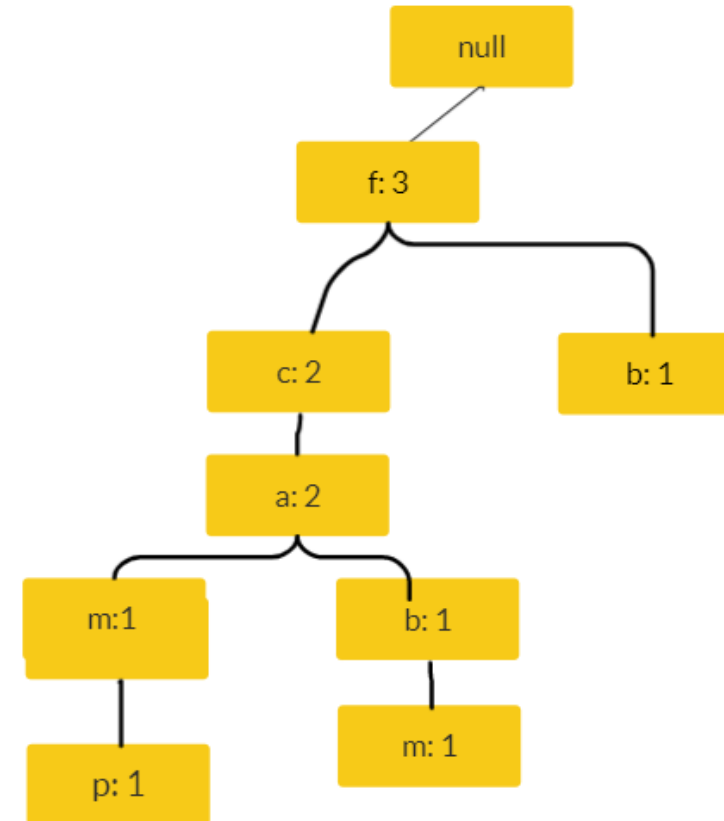
| TID | Item bought |
|-----|-----------------|
| 100 | {f, c, a, m, p} |
| 200 | {f, c, a, b, m} |
| 300 | {f, b} |
| 400 | {c, b, p} |
| 500 | {f, c, a, m, p} |



FP-GROWTH

TID = 300

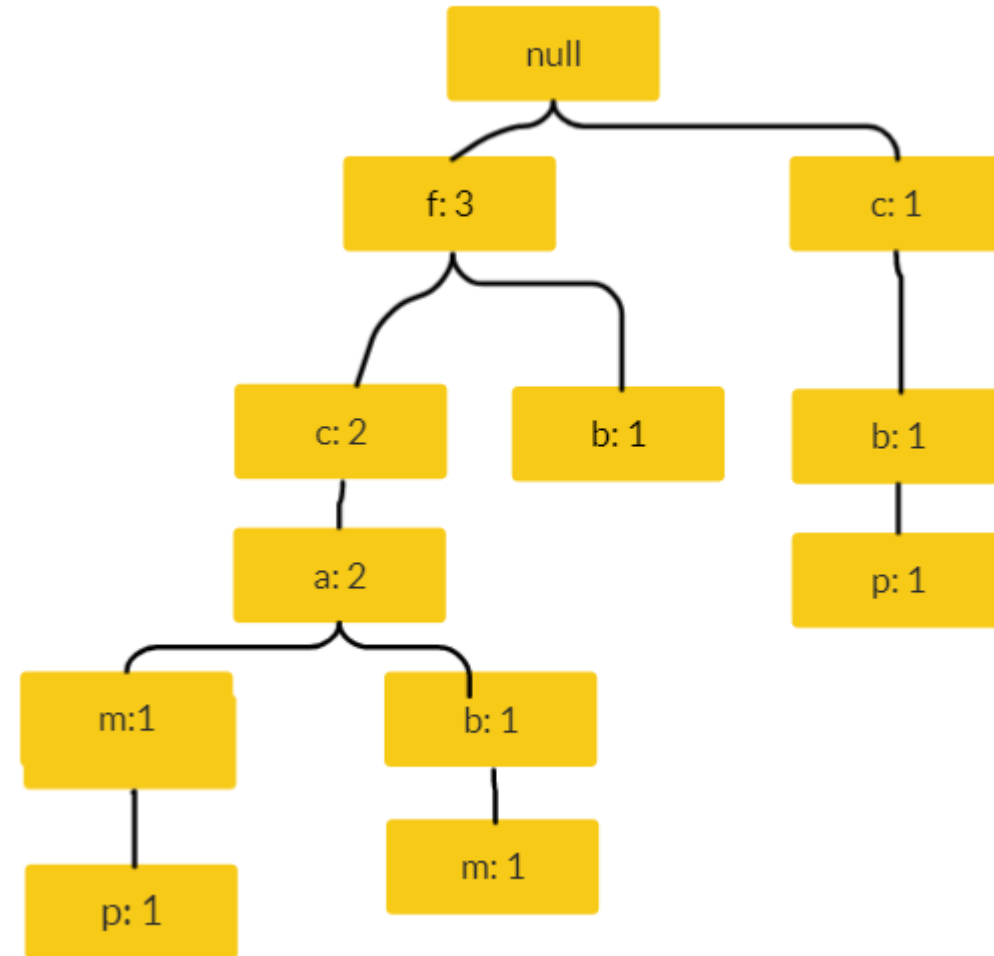
| TID | Item bought |
|-----|-----------------|
| 100 | {f, c, a, m, p} |
| 200 | {f, c, a, b, m} |
| 300 | {f, b} |
| 400 | {c, b, p} |
| 500 | {f, c, a, m, p} |



FP-GROWTH

TID = 400

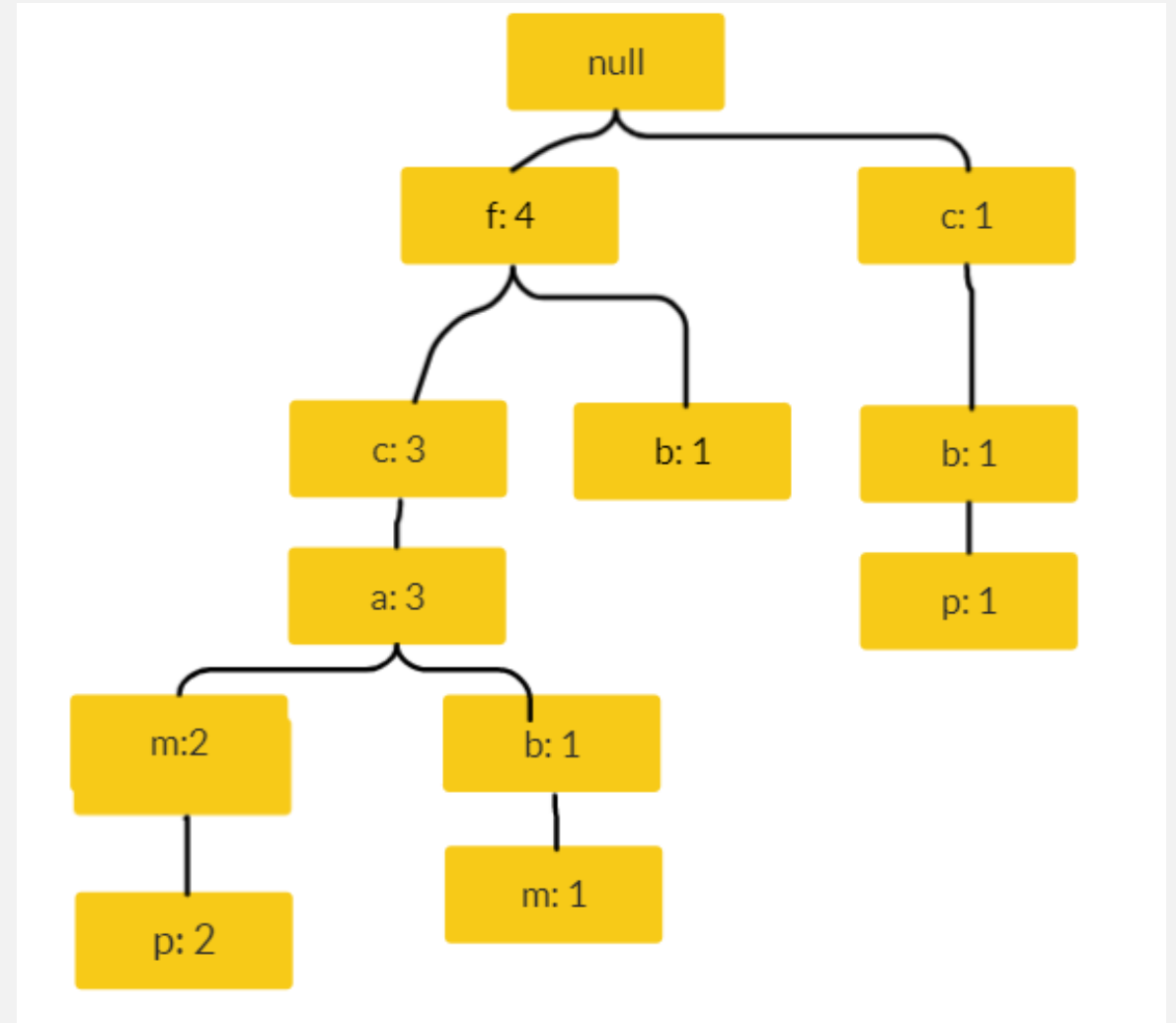
| TID | Item bought |
|-----|-----------------|
| 100 | {f, c, a, m, p} |
| 200 | {f, c, a, b, m} |
| 300 | {f, b} |
| 400 | {c, b, p} |
| 500 | {f, c, a, m, p} |



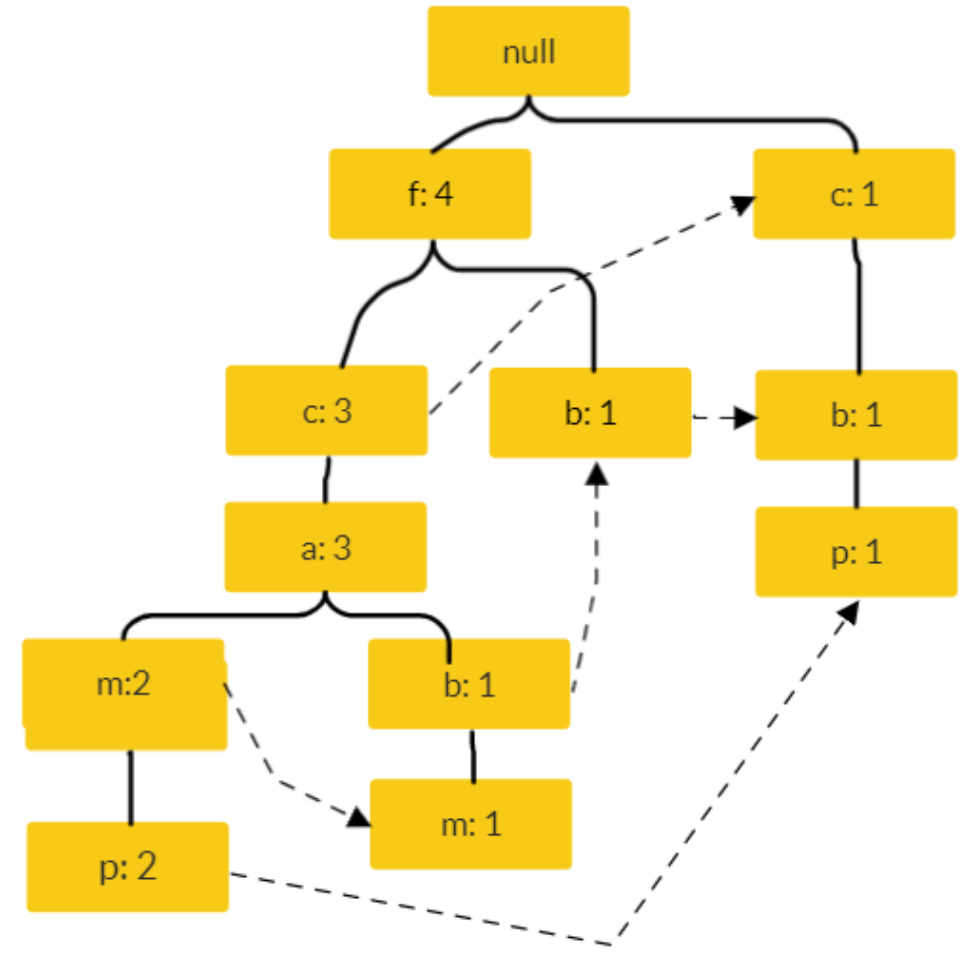
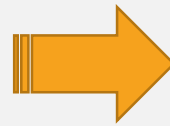
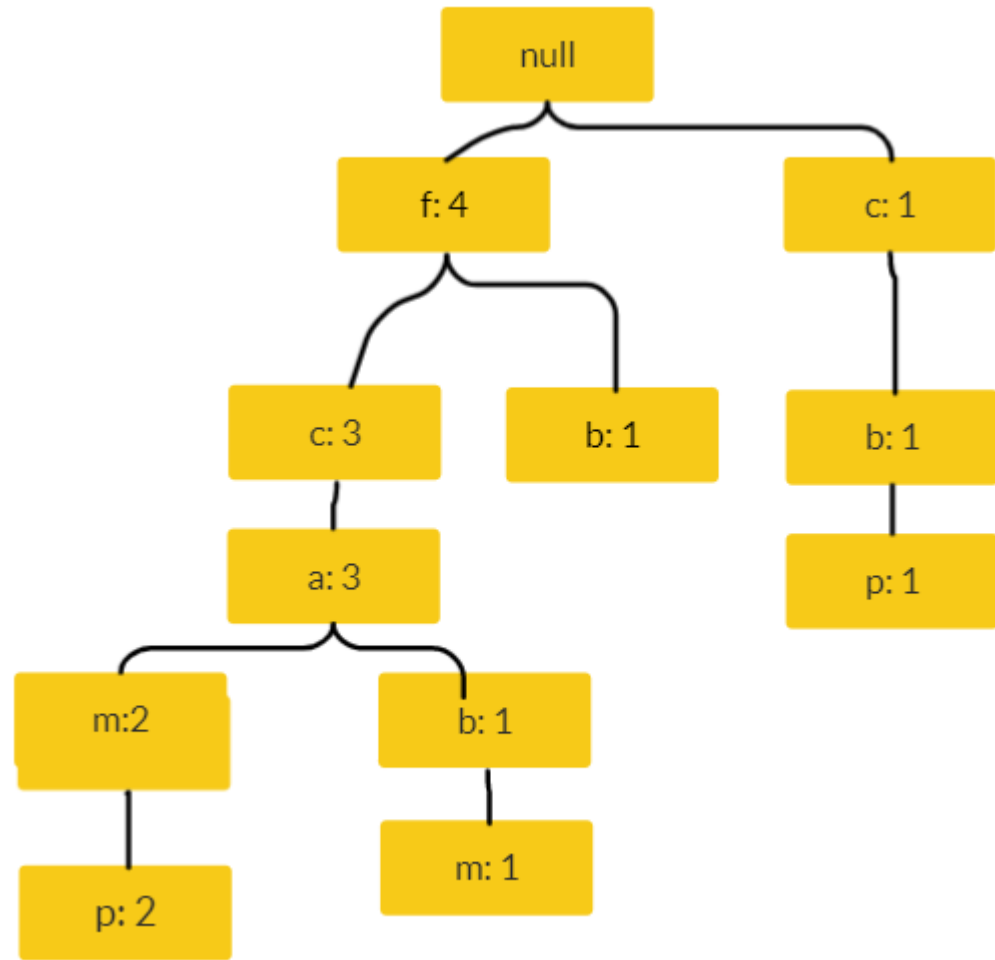
FP-GROWTH

TID = 500

| TID | Item bought |
|-----|-----------------|
| 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} |



FP-GROWTH



FP-GROWTH

| item | freq | rank id |
|------|------|---------|
| f | 4 | 1 |
| c | 4 | 2 |
| a | 3 | 3 |
| b | 3 | 4 |
| m | 3 | 5 |
| p | 3 | 6 |

Find frequent itemset {f}

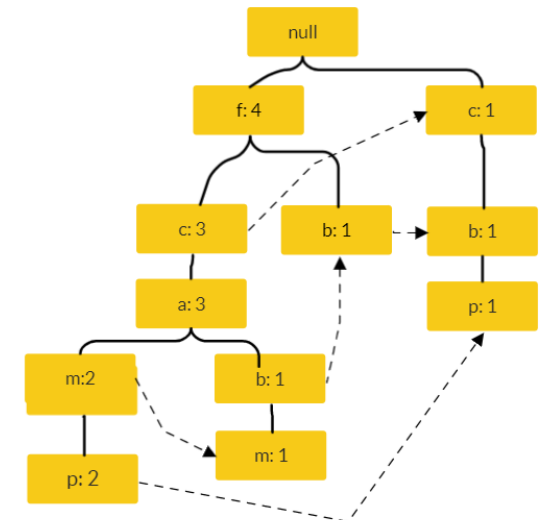
Find frequent itemset {f, c}, {c}

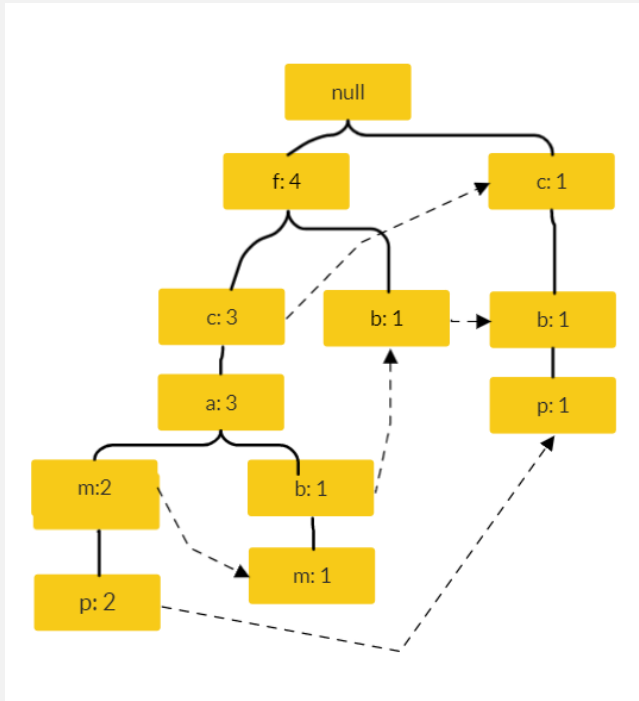
...

Find frequent itemset end with 'm'

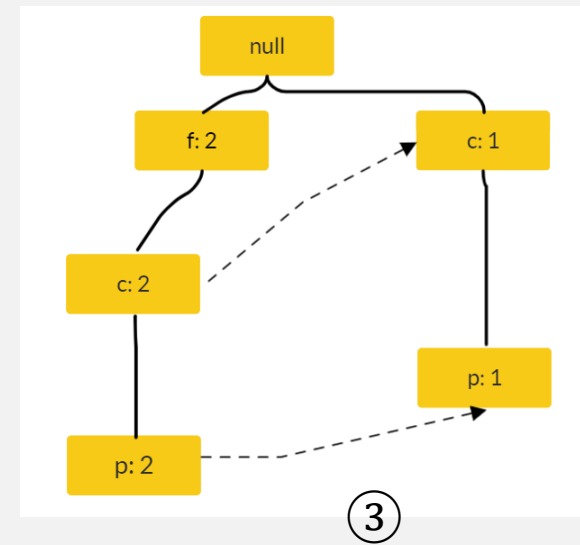
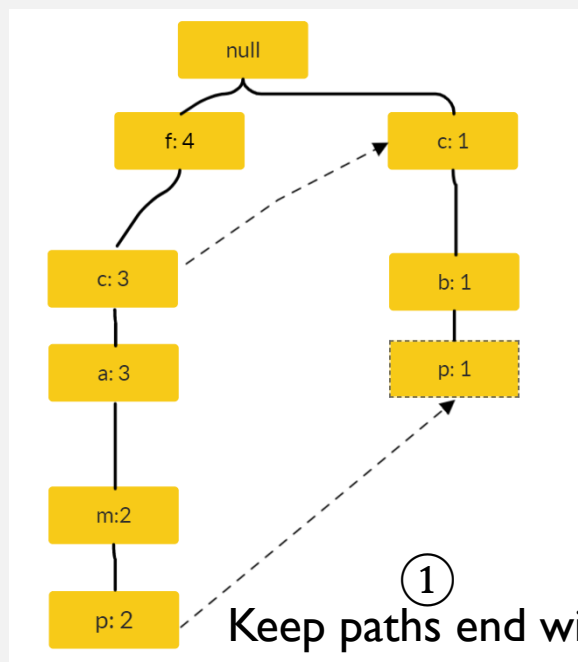
Find frequent itemset end with 'p'

- Suffix based Divide-and-Conquer
- Find frequent itemset end with 'p'
 - frequent itemset end with 'mp' /(bp, ap, cp, fp)
 - frequent itemset end with 'bmp'

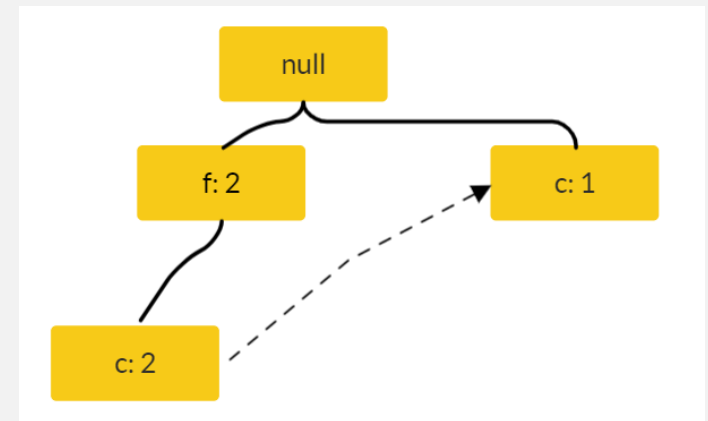
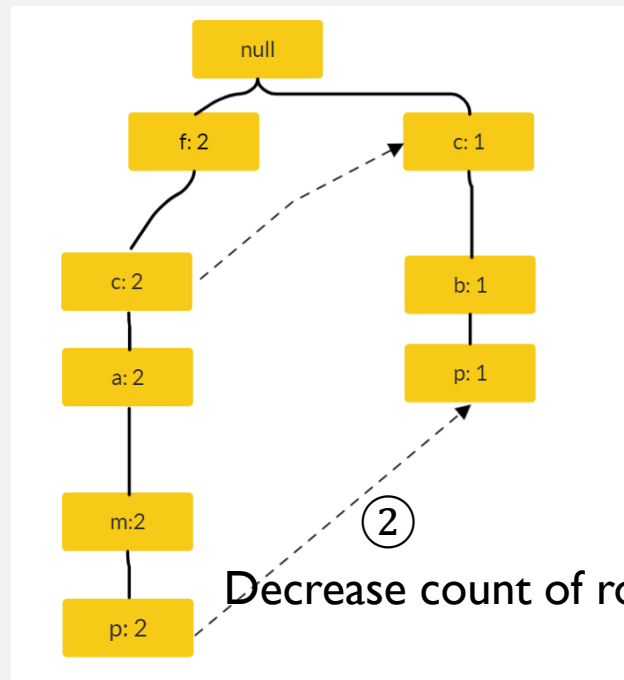




generate frequent itemset {P}



Remove non-frequent items



Recursively solve the conditional FP-tree

EXPERIMENT

- Data set
 - datasetA: realistic groceries data set, 14963 transactions, 167 unique commodities, average transaction length is 2.54
 - datasetB: synthetic data set with 540455 transactions and 2603 unique items. Average transaction length is 4.37
- Minsup:
 - 5%, 2%, 1.5%, 1%, 0.75%, 0.5%

$$support = \frac{\delta(X)}{|\Omega|}$$

RESULT



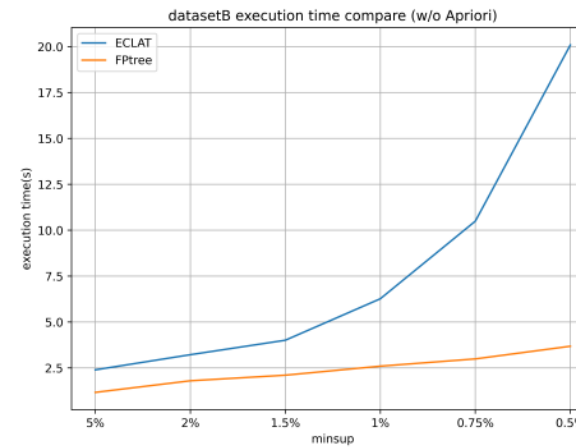
(a) Comparing three algorithms on DatasetA



(b) Comparing three algorithms on DatasetB



(a) Comparing ECLAT and FP-tree on DatasetA



(b) Comparing ECLAT and FP-tree on DatasetB

- Thanks

[\[Scalable Algorithms for Association Mining\]](#)