



# Data Mining

## -- Association Rules

Instructor: Jen-Wei Huang

Office: 92528 in the EE building  
jwhuang@mail.ncku

## FP-Growth [1]

- ▶ Mining frequent patterns without candidate generation
  - Depth-first search approach
- ▶ 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

# Construct FP-tree

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

$min\_support = 3$

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

## Header Table

<i>Item</i>	<i>frequency</i>	<i>head</i>
f	4	
c	4	
a	3	
b	3	
m	3	
p	3	

**F-list** = f-c-a-b-m-p

# Construct FP-tree

<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>
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}

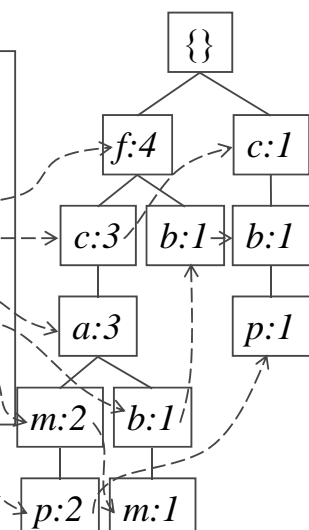
$min\_support = 3$

3. Scan DB again, sort items in the transaction by frequency and construct FP-tree

## Header Table

<i>Item</i>	<i>frequency</i>	<i>head</i>
f	4	
c	4	
a	3	
b	3	
m	3	
p	3	

**F-list** = f-c-a-b-m-p

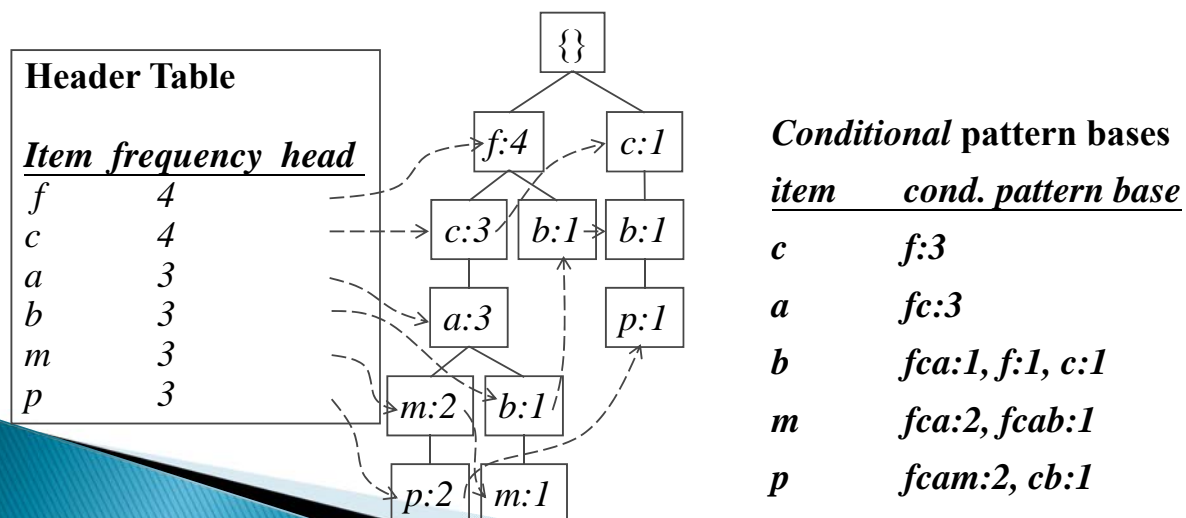


# Partition Database

- ▶ Frequent patterns can be partitioned into subsets according to f-list
  - F-list = f-c-a-b-m-p
  - Patterns containing p
  - Patterns having m but no p
  - ...
  - Patterns having c but no a nor b, m, p
  - Pattern f
- ▶ Completeness and non-redundancy

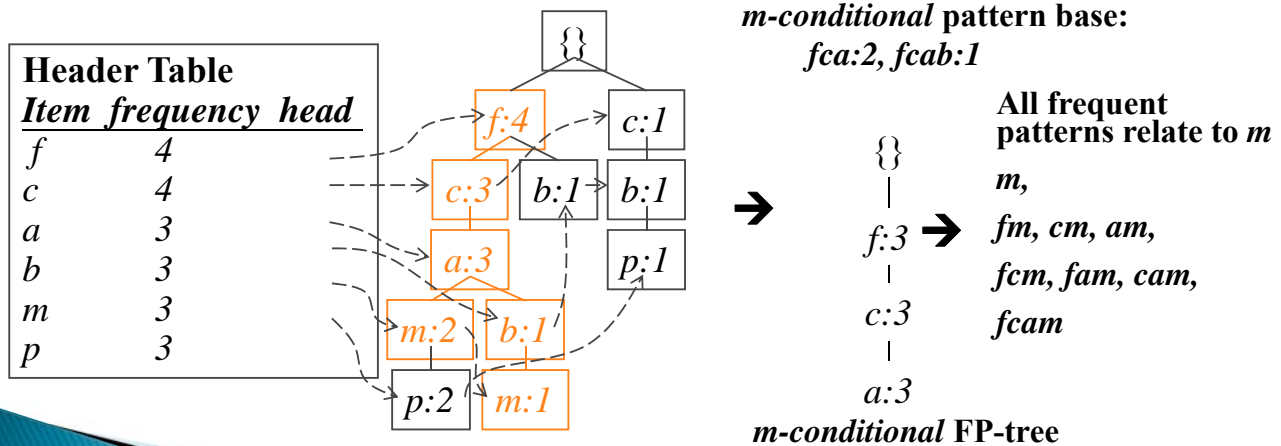
## Conditional Pattern Bases

- ▶ Starting at the least frequent item in the header table
- ▶ Traverse the FP-tree by following the link of each frequent item
- ▶ Accumulate all of *transformed prefix paths* of the item to form its conditional pattern base

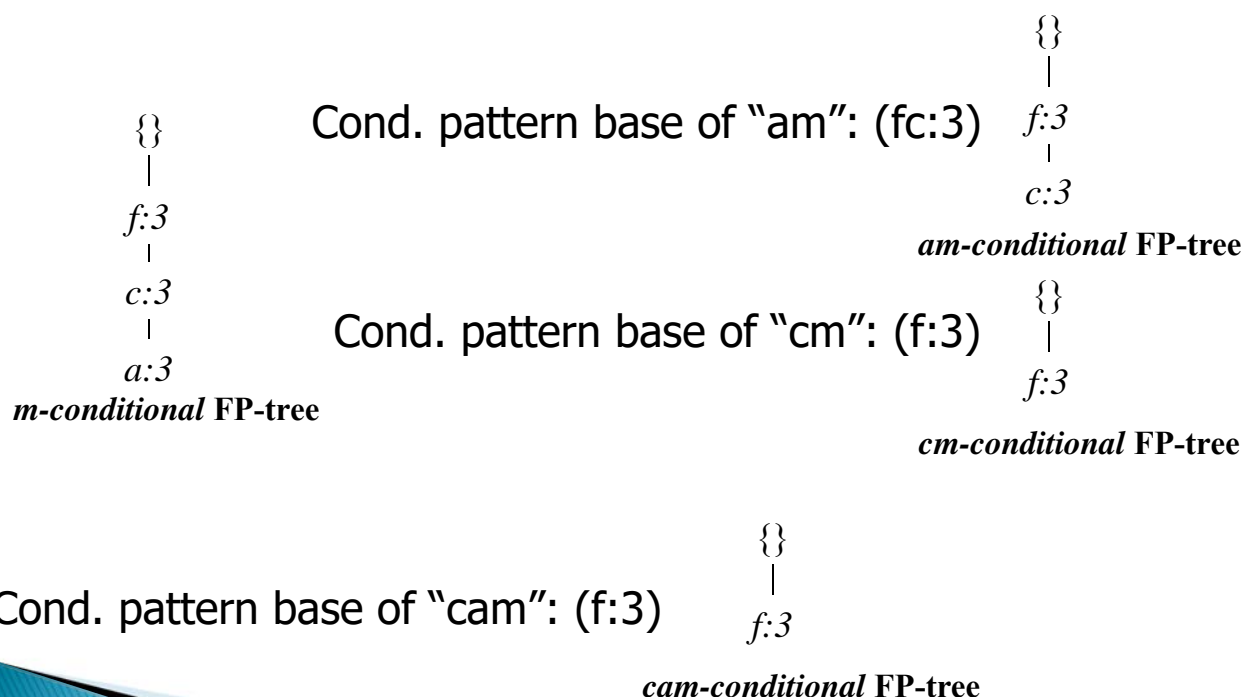


# Conditional FP-trees

- ▶ For each pattern-base
  - Accumulate the count for each item in the base
  - Construct the FP-tree for the frequent items of the pattern base

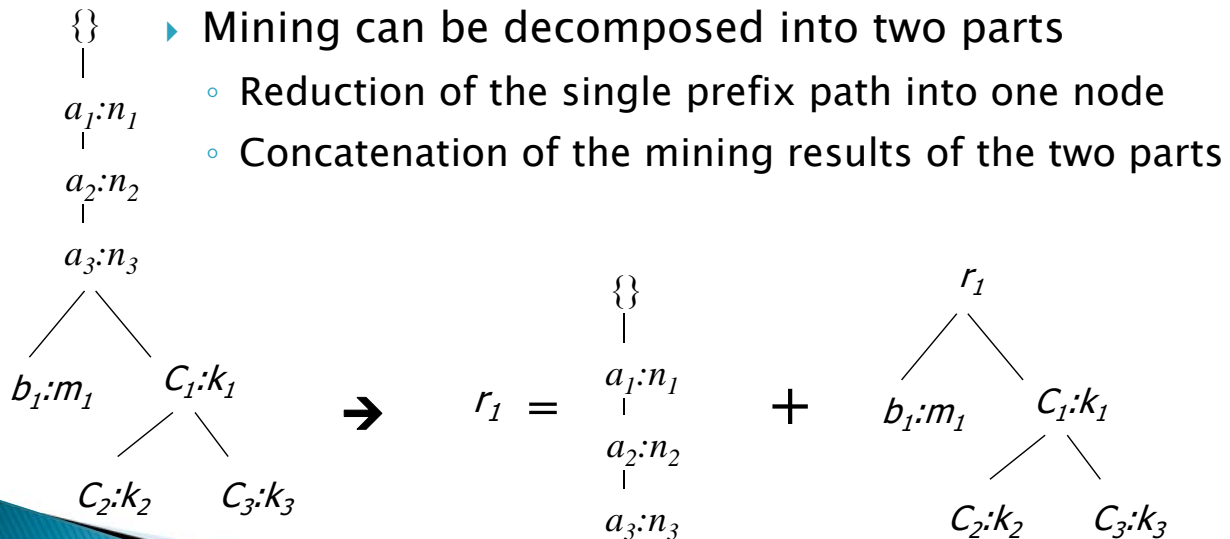


# Mining Conditional FP-tree



# Single Prefix Path

- Suppose a (conditional) FP-tree  $T$  has a shared single prefix-path  $P$



## Benefits of FP-tree

- Completeness
  - Preserve complete information for frequent pattern mining
  - Never break a long pattern of any transaction
- Compactness
  - Reduce irrelevant info—infrequent items are gone
  - Items in frequency descending order: the more frequently occurring, the more likely to be shared
  - Never be larger than the original database (not count node-links and the *count* field)

# FP-Growth Algorithm

- ▶ 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

## Problems of FP-Growth

- ▶ What about if FP-tree cannot fit in memory?
  - DB projection
- ▶ First partition a database into a set of projected DBs
- ▶ Then construct and mine FP-tree for each projected DB

# ECLAT [3]

- ▶ Mining by exploring vertical data format
- ▶ Vertical format:  $t(AB) = \{T_{11}, T_{25}, \dots\}$ 
  - tid-list: list of trans.-ids containing an itemset
- ▶ Deriving frequent patterns based on vertical intersections
- ▶ 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  $(XY, X) = \{T_2\}$

## Basic Extensions

- ▶ **Max-pattern** [5]
  - R. J. Bayardo. Efficiently mining long patterns from databases. SIGMOD'98.
- ▶ **Closed-pattern** [6]
  - N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Discovering frequent closed itemsets for association rules. ICDT'99.
- ▶ **Sequential pattern** [7]
  - R. Agrawal and R. Srikant. Mining sequential patterns. ICDE'95



# Closed Patterns and Max-Patterns

- ▶ A long pattern contains a combinatorial number of sub-patterns, e.g.,  $\{a_1, \dots, a_{100}\}$  contains  $\binom{100}{1} + \binom{100}{2} + \dots + \binom{100}{100} = 2^{100} - 1 = 1.27 \cdot 10^{30}$  sub-patterns!
- ▶ Solution: Mine *closed patterns* and *max-patterns* instead
- ▶ An itemset  $X$  is **closed** if  $X$  is *frequent* and there exists *no super-pattern*  $Y \supset X$ , *with the same support* as  $X$ 
  - Closed pattern is a lossless compression of freq. patterns
  - Reducing the # of patterns and rules
- ▶ An itemset  $X$  is a **max-pattern** if  $X$  is frequent and there exists no frequent super-pattern  $Y \supset X$

## Examples

- ▶ Exercise.  $DB = \{\langle a_1, \dots, a_{100} \rangle, \langle a_1, \dots, a_{50} \rangle\}$ 
  - $\text{Min\_sup} = 1$ .
- ▶ What is the set of **closed itemset**?
  - $\langle a_1, \dots, a_{100} \rangle: 1$
  - $\langle a_1, \dots, a_{50} \rangle: 2$
- ▶ What is the set of **max-pattern**?
  - $\langle a_1, \dots, a_{100} \rangle: 1$
- ▶ What is the set of **all patterns**?
  - !!



# Computational Complexity

- ▶ How many itemsets are potentially to be generated in the worst case?
  - The number of frequent itemsets to be generated is sensitive to the min\_sup threshold
  - When min\_sup is low, there exist potentially an exponential number of frequent itemsets
  - The worst case:  $M^N$  where M: # distinct items, and N: max length of transactions
- ▶ The worst case complexity vs. the expected probability
  - Ex. Suppose Walmart has  $10^4$  kinds of products
    - The chance to pick up one product  $10^{-4}$
    - The chance to pick up a particular set of 10 products:  $\sim 10^{-40}$
    - What is the chance this particular set of 10 products to be frequent  $10^3$  times in  $10^9$  transactions?

## References

- ▶ [1] J. Han, J. Pei, and Y. Yin. Mining frequent patterns without candidate generation. SIGMOD'00
- ▶ [2] G. Grahne and J. Zhu, Efficiently Using Prefix-Trees in Mining Frequent Itemsets, Proc. ICDM'03 Int. Workshop on Frequent Itemset Mining Implementations (FIMI'03), Melbourne, FL, Nov. 2003
- ▶ [3] M. J. Zaki, Scalable Algorithms for Association Mining. IEEE Transactions on Knowledge and Data Engineering, 12(3):372–390. May/June 2000
- ▶ [4] M. J. Zaki and Karam Gouda, Fast Vertical Mining Using Diffsets. In 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. August 2003.
- ▶ [5] R. J. Bayardo. Efficiently mining long patterns from databases. SIGMOD'98.
- ▶ [6] N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Discovering frequent closed itemsets for association rules. ICDT'99.
- ▶ [7] R. Agrawal and R. Srikant. Mining sequential patterns. ICDE'95

# References

- ▶ Slides from Prof. J.-W. Han, UIUC
- ▶ Slides from Prof. M.-S. Chen, NTU
- ▶ Slides from Prof. W.-Z. Peng, NCTU