



University of Alberta

Database Laboratory



# Relevance of Counting in Data Mining Tasks

Osmar R. Zaiane

蔡頤安

<http://www.cs.ualberta.ca/~zaiane/>  
zaiane@cs.ualberta.ca

The First International Conference on  
Advanced Data Mining and Applications  
Wuhan, China, July 22-24, 2005

Calculateur/Ordinateur

计算机 Rechner

Вычислительная машина

آلة الحاسبة

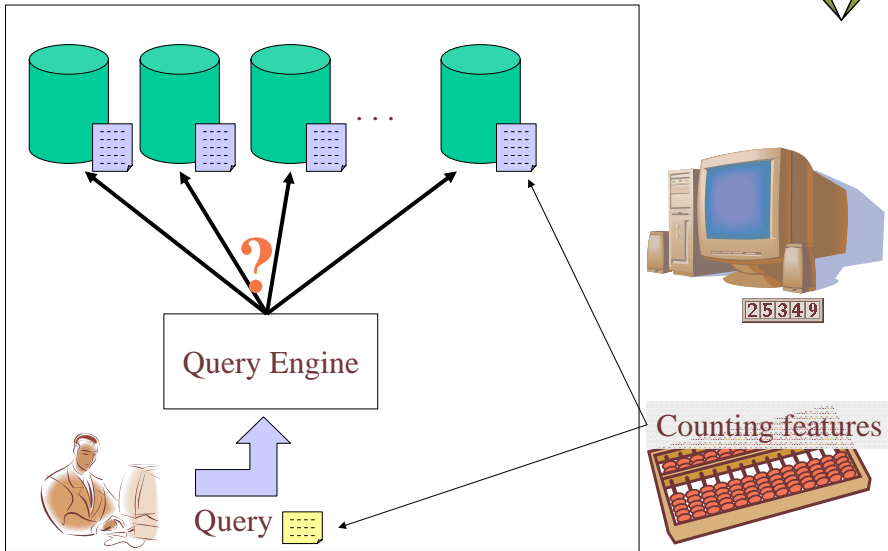
Computer

الحسوب Calcolatore

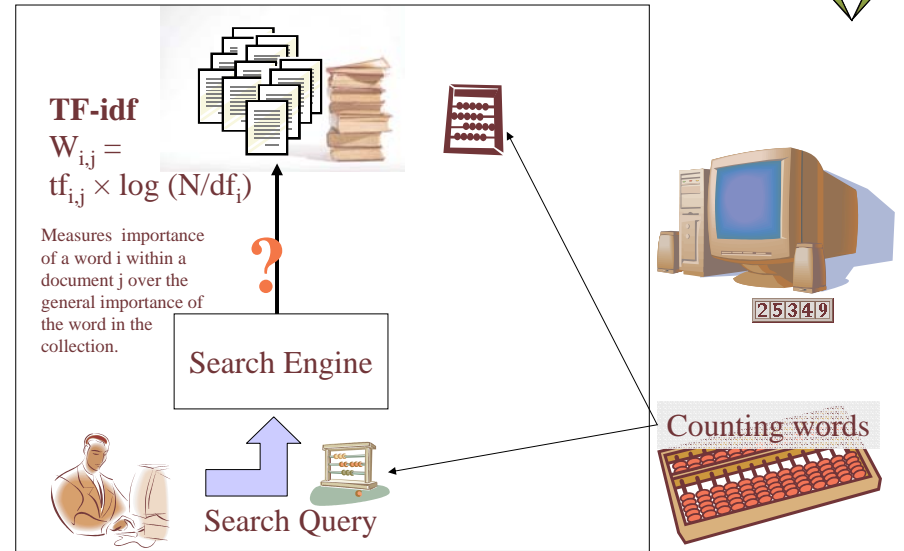
123 45678901234 567890123 4567890 123 45678901 23 4567 890  
1234 56789 012345 67890 1234 5678901 234 567890 12 34567 8901  
2345 678 901 234 5678 90 12 3456 789 0123 456 78901 2345 678 901  
23456 789 012 345 678 90123 4567 89012 34567890 123456789012  
34 567 890 12 345 678 901 23456 7890 1 23 45 67890 123 456 7890  
12345 678 90 12 345 678 901 2345 67890 123 4567 890 1234 56 78901  
23 4567 890 12 3456 78 9012 3456 78 9012 345 678 90 1234 567 8901  
234 567 8901 23456 78 90123 4567 890 1234 5678 90 12 34 5678  
9012 3456 7890 123 4567 890 12 3456 78 901 2345 678 90 1234 56  
7890 123 456 7890 12 34567 890 123 456 7890 1234 5678 9012 34 56  
78 90 12 345 678 901 234 5678 9012 34 56 7890 1234 567 890 1234  
5678 901 234 567 890 1234 5678 9012 34 56 78 90 12 34 567 890 123  
456 78 90 12 34 56 7890 1234 5678 90 12 34 56 789 12 34 567 890 123  
45 67 89 12 3456 7890 12 345 678 9012 34 567 8901 23 456 7890 12  
345 6789 12345 67 890 1234 56789 12345 67890 1234 5678 901 234  
56 78 901 234 567 8901 2345 67890 12345 678 901 234 5678 9012 345  
678 901 23 45 678 9012 34 567 8901 23 45 678 90123 45 678 9012 345  
678 901 2345 67 890 12345 67 890 12 345 678 901 234 56789 1 23 56

0101000100101110101001110100100111010100111001110101011101  
01001110001110011101010011100111001001010110010101000110  
0100111010101001110010100011001010011101010100101010001101  
00100100111001110000100100010001010100011010010101001110  
00111010101110111010101001000111101111110011111100111011  
00111000011001001110011010110001001010101010100011001010100  
1000100100100001001110101001010101011101100010101001011101  
0100010010010101010100100100101010100010001000100101100  
1001001101011011101001010100100011110111111001111110011101  
1001110000110001001110011010110001001010110101010010101001  
0100101010101010100101010011010101001011001110011100111001  
0011100010101101110100101010010001111011111100111111001110  
1100111000011000100111001101011000100111110111001110001011  
10000111010101101001010100100011110111111001111110011101  
10011100001100010011100110101100010011001101010100011000  
11100101101010110111010010101001000111101111100111110011  
10110011100001100100101001100110011001100010011000100101010  
101101011010110111010010101001000111101111100111111001110

## Database Selection Problem



## Information Retrieval



## Discretization & Concept Hierarchies

- Handling continuous data
- Automatic Concept hierarchy building
- Numerosity Reduction
- Smoothing Noise

Discretization is used to reduce the number of values for a given continuous attribute, by dividing the range of the attribute into intervals.

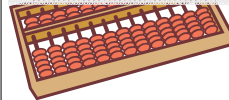
Binning  
Histogram analysis  
Entropy-based  
3-4-5 data segmentation



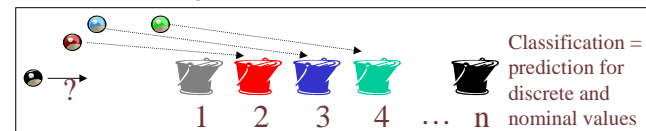
Statistics



Counting values



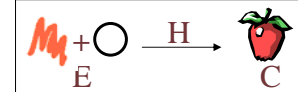
## Naïve/Full Bayesian Classifier



### Bayesian Learning (Bayes Theorem)

Given a hypothesis  $H$  that some data belongs to a class  $C$  and some evidence  $E$  about the data, the posteriori probability of  $H$  given  $E$  is:

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)}$$

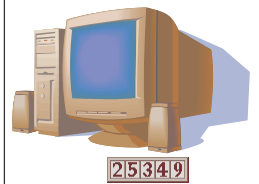


$$P(H) = P(\text{Mushroom})$$

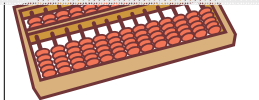
$$P(E) = P(M + \bigcirc)$$

$$P(E|H) = P(M + \bigcirc \text{ if Mushroom})$$

Learning = Calculating prior and posterior probabilities.

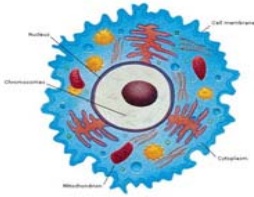

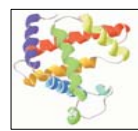


Counting occurrences



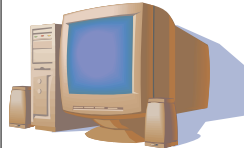
# Genomics and Proteomics



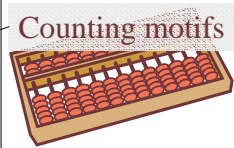




Predicting function or location of a Protein

Investigate frequent sequences and sub-sequences








25349



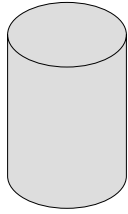
Counting motifs

# Market Basket Analysis

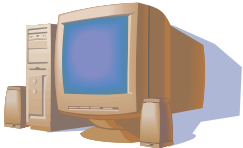


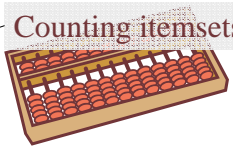
Store



Bread, → Milk  
Coke, Chips → Hot dogs

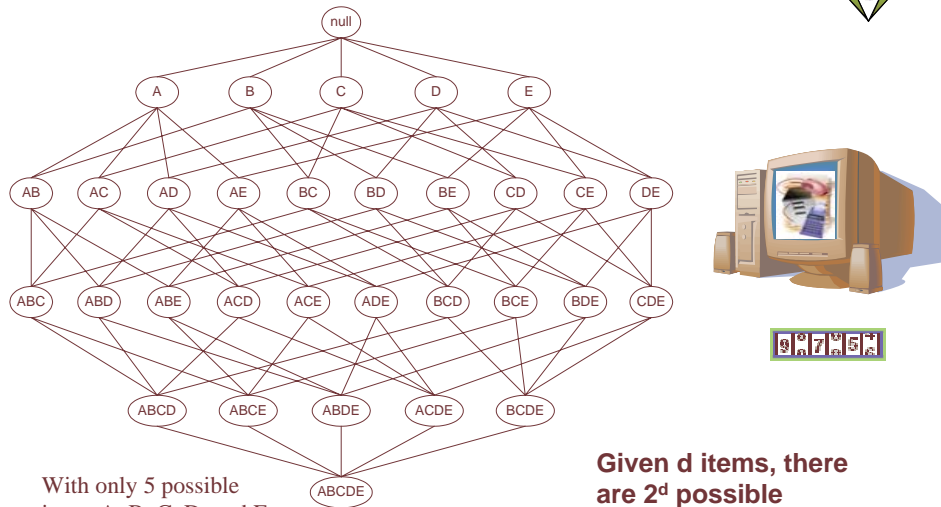


25349



Counting itemsets

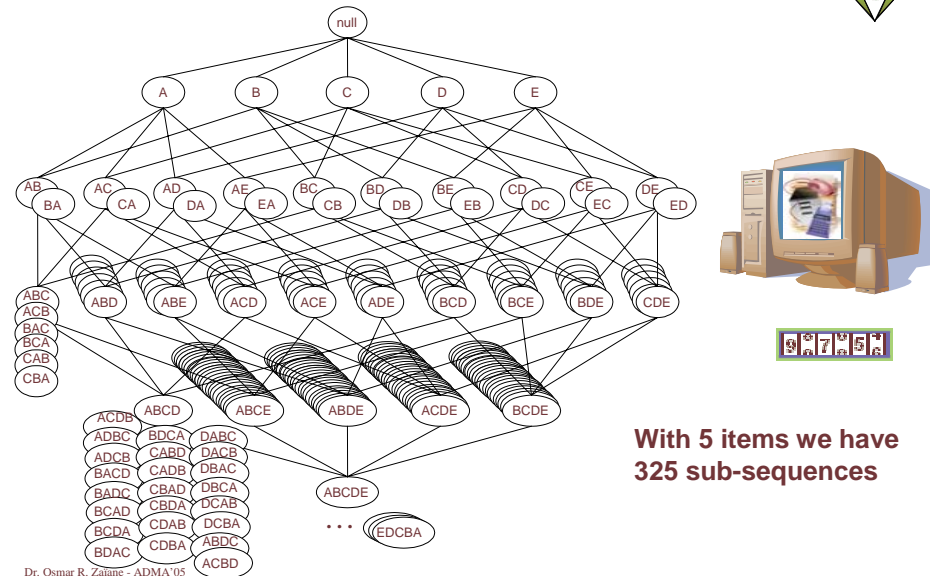
# Frequent Itemset Search Space



With only 5 possible items A, B, C, D, and E, there are  $2^5=32$  cases.

Given  $d$  items, there are  $2^d$  possible candidate itemsets

# Frequent subsequence Search Space



With 5 items we have 325 sub-sequences

# Association Rule Mining



## Frequent Itemset Mining

FIM

1

abc

Bound by a **support** threshold

## Association Rules Generation

2

$ab \rightarrow c$   
 $b \rightarrow ac$

Bound by a **confidence** threshold

- Frequent itemset generation is still computationally expensive

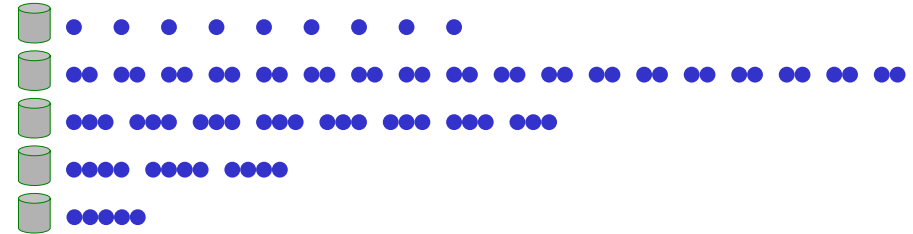
# Apriori Algorithm



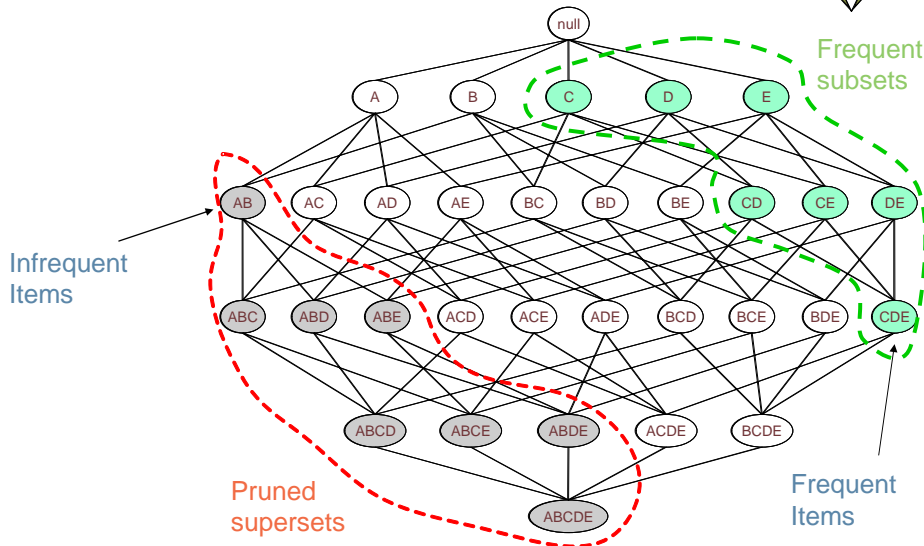
## Apriori (Agrawal et al. 1994)

Repetitive I/O scans

Huge Computation to generate candidate items



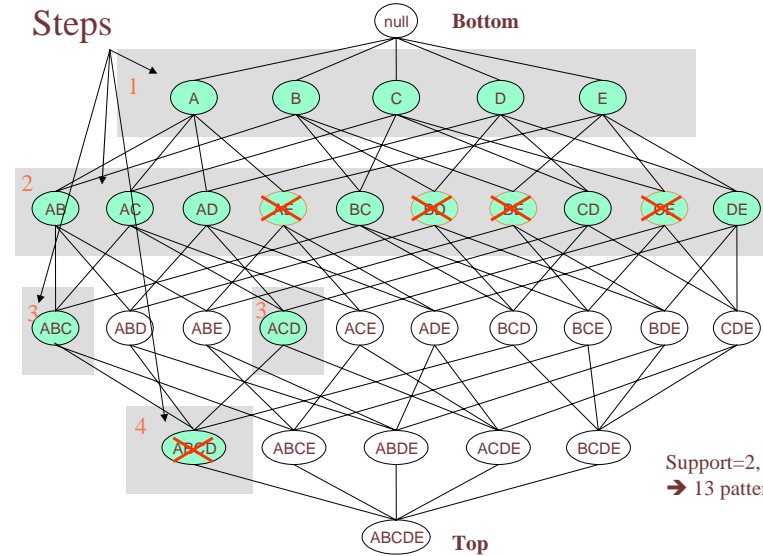
# Illustrating Apriori Principle



# Bottom - Up Example



## Steps

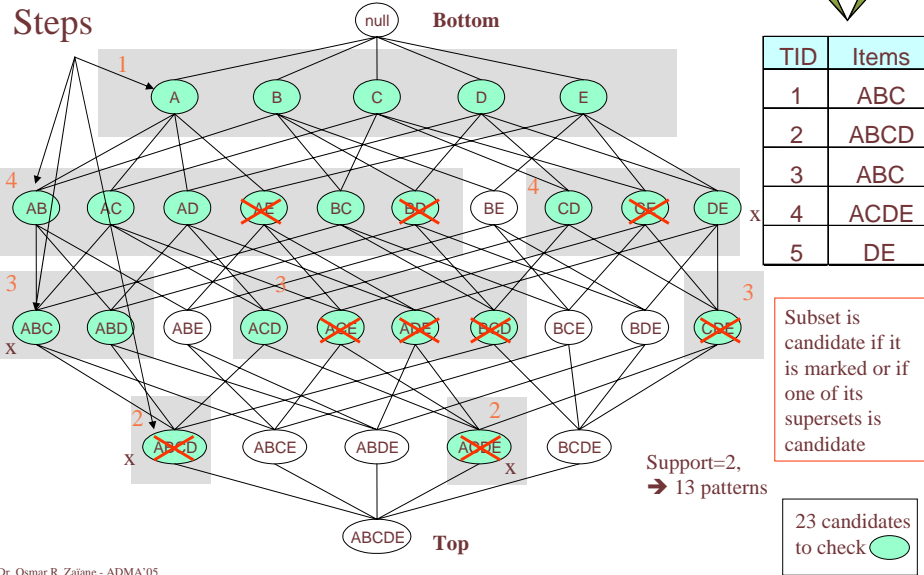


TID	Items
1	ABC
2	ABCD
3	ABC
4	ACDE
5	DE

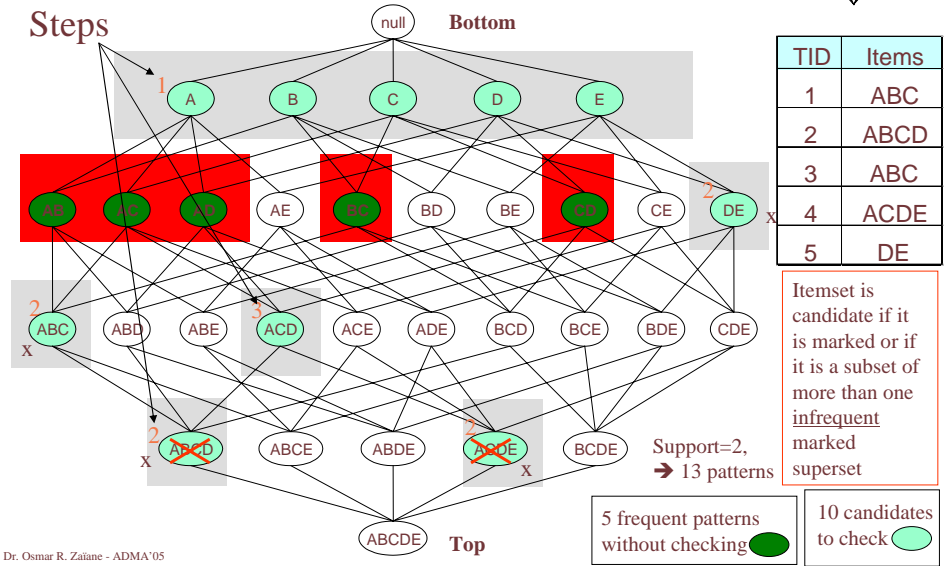
Superset is candidate if ALL its subsets are frequent

18 candidates to check

## Top - Down Example



## Leap Traversal Example



## Many Candidates – Many Patterns

Not only there are too many candidate itemsets but there are also too many frequent ones.

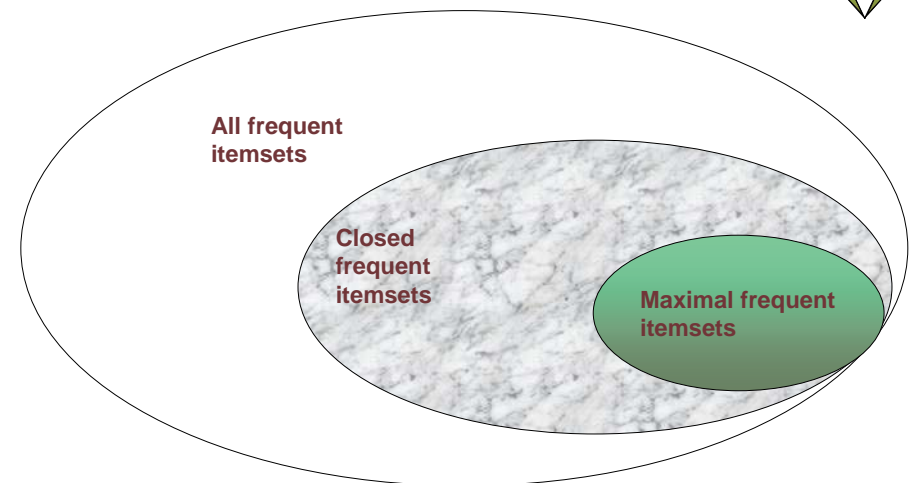
Frequent pattern  $\{a_1, \dots, a_{100}\}$

$$\rightarrow (100^1) + (100^2) + \dots + (100^{100})$$

$$= 2^{100} - 1$$

$$= 1.27 * 10^{30} \text{ frequent sub-patterns!}$$

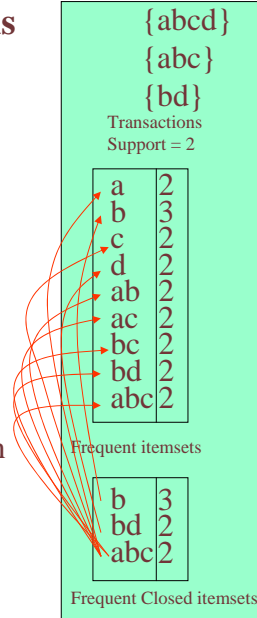
## Compressed Representation



Maximal frequent itemsets  $\subseteq$  Closed frequent itemsets  $\subseteq$  All frequent itemsets

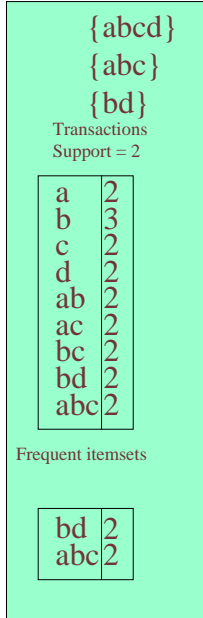
## Frequent Closed Patterns

- N. Pasquier et al. In ICDT'99
- For frequent itemset X, if there exists no item y such that every transaction containing X also contains y, then X is a frequent closed pattern
- In other words, frequent itemset X is closed if there is no item y, not already in X, that always accompanies X in all transactions where X occurs.
- Concise representation of frequent patterns. Can generate all frequent patterns with their support from frequent closed ones.
- Reduce number of patterns and rules

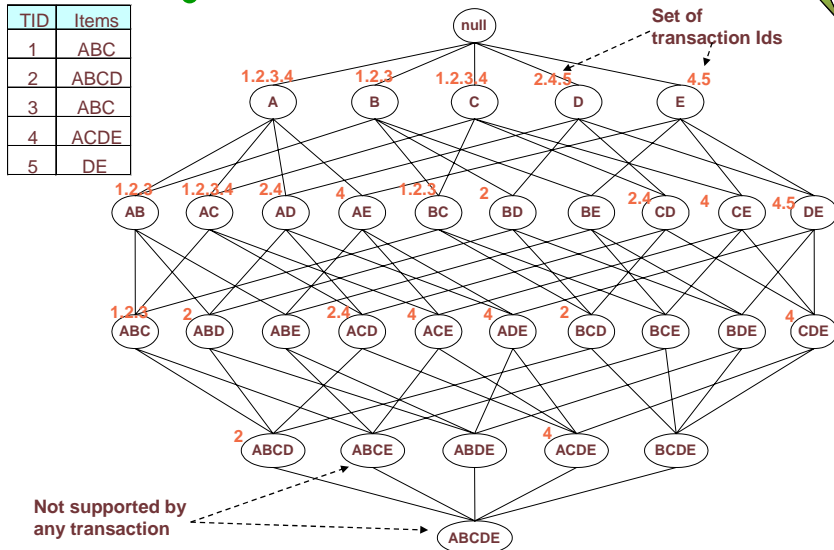


## Frequent Maximal Patterns

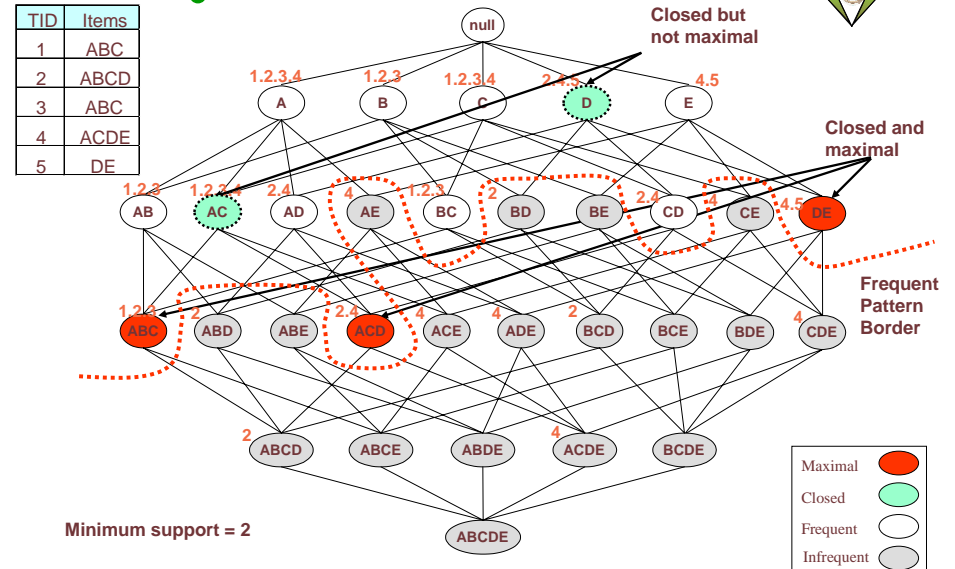
- R. Bayardo. In SIGMOD'98
- Frequent itemset X is maximal if there is no other frequent itemset Y that is superset of X.
- In other words, there is no other frequent pattern that would include a maximal pattern.
- More concise representation of frequent patterns but the information about supports is lost.
- Can generate all frequent patterns from frequent maximal ones but without their respective support.



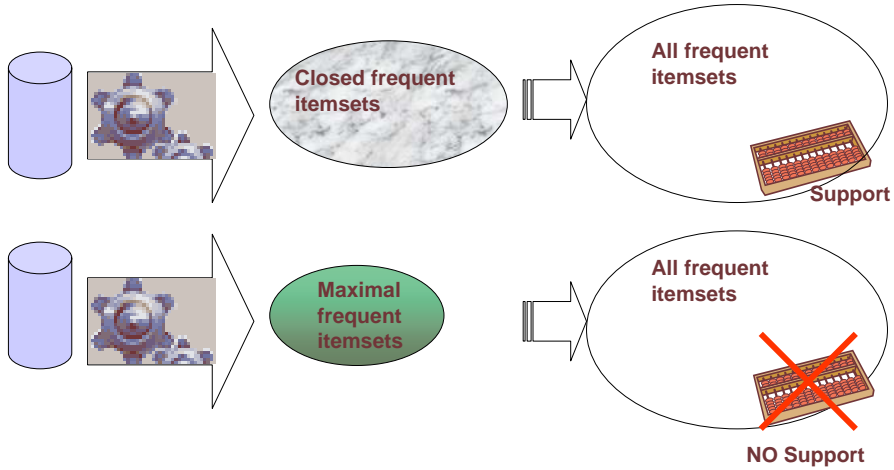
## Maximal Versus Closed Patterns



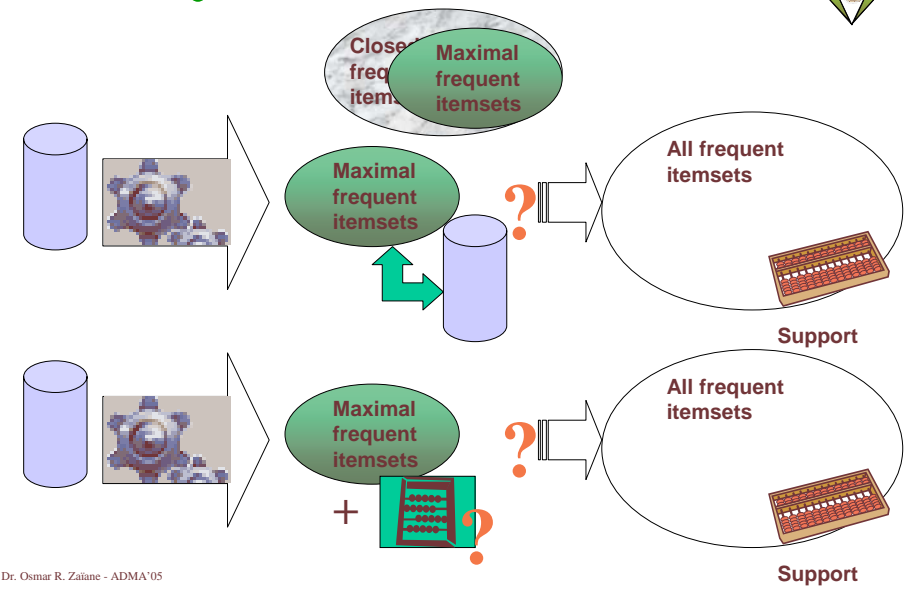
## Maximal Versus Closed Patterns



# Do We Need to Count All?



# What about Maximals?

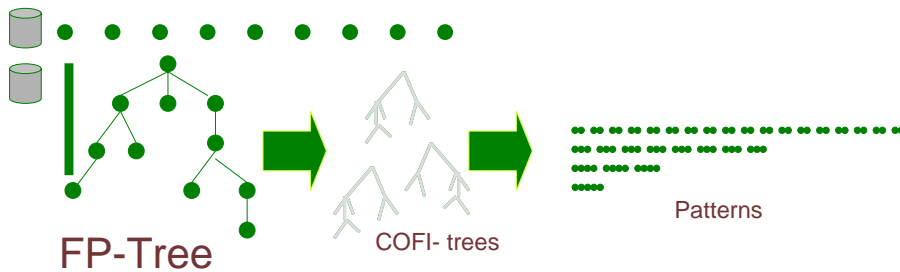


# The COFI Approach

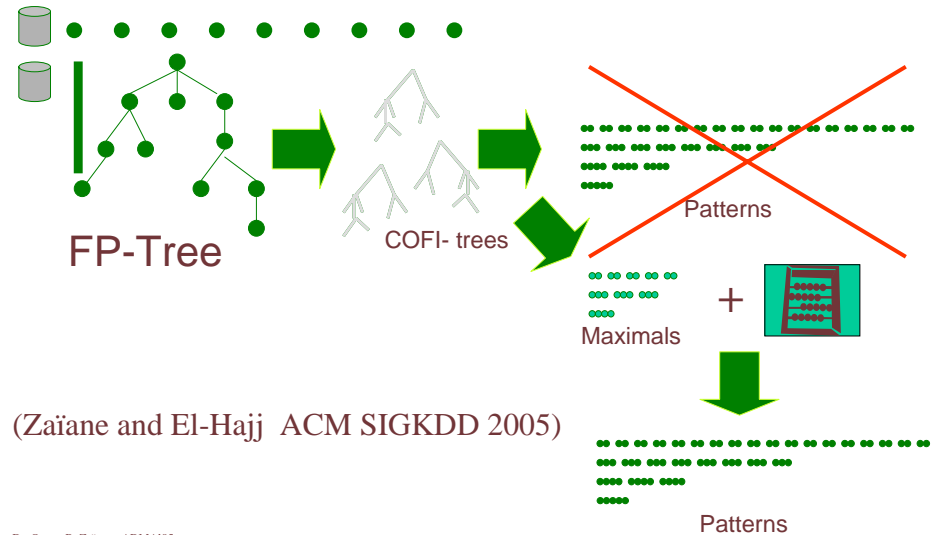


COFI  
(El-Hajj and Zaiane, 2003)

2 I/O scans  
reduced candidacy generation  
Small memory footprint

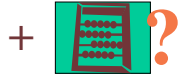


# COFI-MAX and the extra info



(Zaiane and El-Hajj ACM SIGKDD 2005)

## Ordered Partitioning Bases



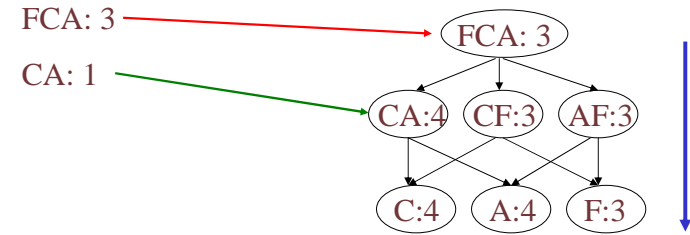
What is this extra information?

- A data structure containing *frequent pattern bases* and their *branch support*;
- The data structure is a “free” bonus since it is used to mine for maximals;
- *Frequent pattern bases* are those marked sub-transactions in the leap-approach and their descendants if not frequent.
- The *branch support* is the number of times the *frequent pattern base* occurs alone (not subsumed by another pattern)

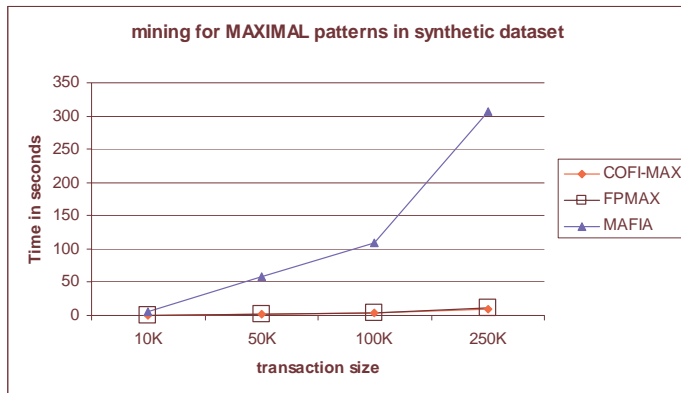
## Counting the Support



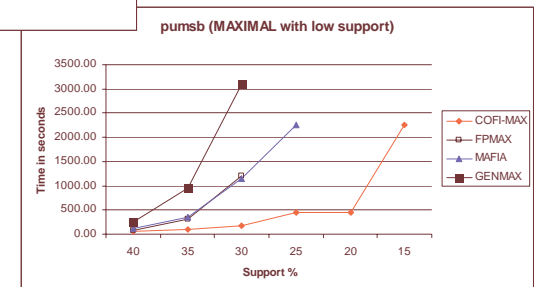
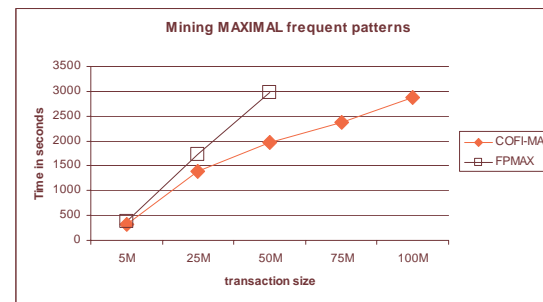
Support of any pattern is the summation of the supports of its supersets of frequent-path-bases



## Some Results (synthetic data)



## Some Other Selected Results



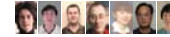


## In Conclusion ...



- Computers are machines that count and compute
- Many data mining tasks consist in counting
- The task of enumerating and counting is essential but not necessarily easy.
- We do not need to count all possibilities or even all patterns of direct interest
- The challenge is to reduce the enumeration without losing effectiveness (loss-less compression)
- There is no winner / no best way to count

## My Students



- Maria-Luiza Antonie
- Jiyang Chen
- Mohammad El-Hajj
- Andrew Foss
- Yan Jin
- Yi Li
- Yaling Pei



- Stanley Oliveira
- Yang Wang
- Lisheng Sun
- Jia Li
- Alex Strilets
- William Cheung
- Yue Zhang
- Chi-Hoon Lee
- Weinan Wang
- Ayman Ammoura
- Hang Cui
- Jun Luo
- Yuan Ji

Without my students this research work wouldn't have been possible.

# 谢谢



蔡頤安



**Osmar R. Zaiane, Ph.D.**  
Associate Professor  
Department of Computing Science

221 Athabasca Hall  
Edmonton, Alberta  
Canada T6G 2E8

Telephone: Office +1 (780) 492 2860  
Fax +1 (780) 492 1071  
E-mail: [zaiane@cs.ualberta.ca](mailto:zaiane@cs.ualberta.ca)  
<http://www.cs.ualberta.ca/~zaiane/>

<http://www.cs.ualberta.ca/~zaiane/>