

# Sample-Wise Enumeration Methods for Mining Microarray Datasets

Anthony K. H. Tung

Department of Computer Science

National University of Singapore

# A Microarray Dataset

← 1000 - 100,000 columns →

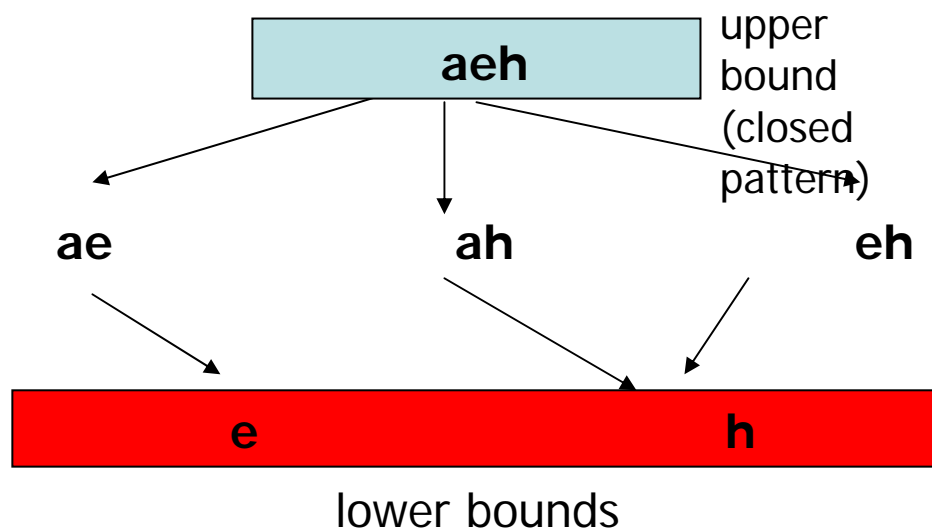
↑ 100-500 rows ↓

	Class	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6	Gene7
Sample1	Cancer							
Sample2	Cancer							
⋮								
SampleN-1	~Cancer							
SampleN	~Cancer							

- Find closed patterns which occur frequently among genes.
- Find rules which associate certain combination of the columns that affect the class of the rows
  - Gene1, Gene10, Gene1001 -> Cancer

# Challenge I

- Large number of patterns/rules
  - number of possible column combinations is extremely high
- Solution: Concept of a **closed pattern**
  - Patterns are found in exactly the same set of rows are grouped together and represented by their upper bound
- Example: the following patterns are found in row 2,3 and 4 and 4



<i>i</i>	<i>ri</i>	<i>Class</i>
1	<b>a,b,c,l,o,s</b>	<b>C</b>
2	<b>a,d,e,h,p,l,r</b>	<b>C</b>
3	<b>a,c,e,h,o,q,t</b>	<b>C</b>
4	<b>a,e,f,h,p,r</b>	<b>~C</b>
5	<b>b,d,f,g,l,q,s,t</b>	<b>~C</b>

"a" however not part of the group

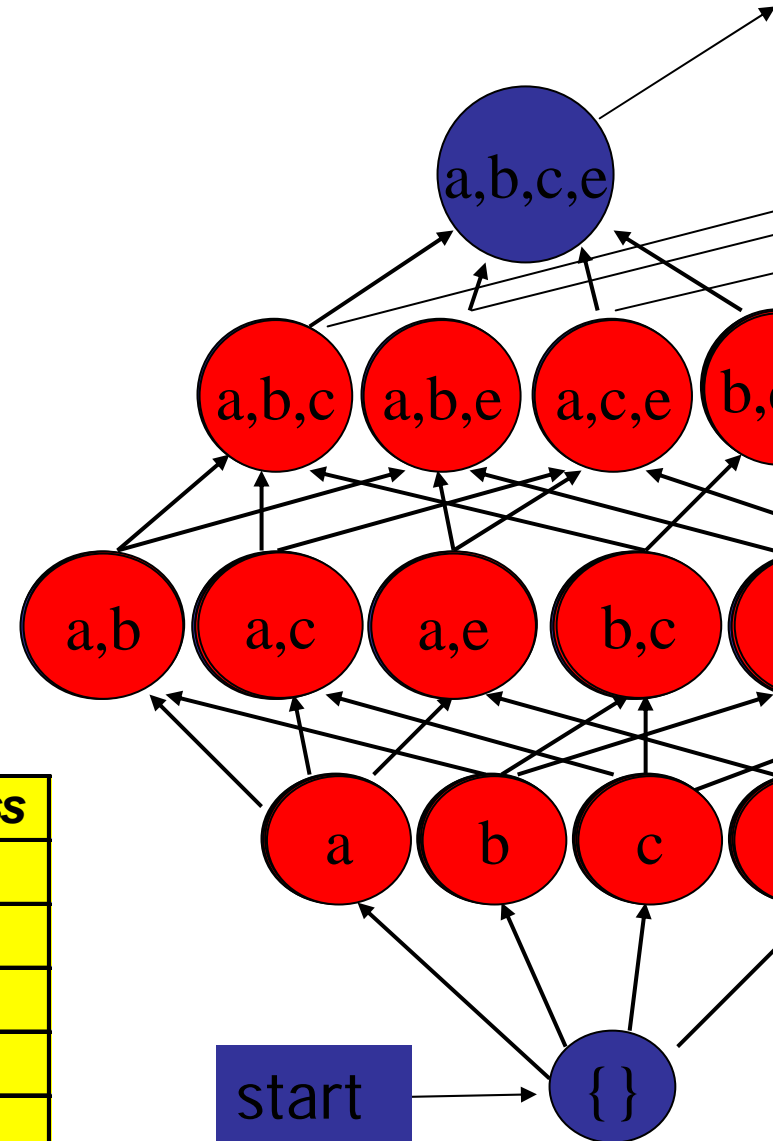
# Challenge II

- Most existing frequent pattern discovery algorithms perform searches in the **column/item enumeration space** i.e. systematically testing various combination of columns/items
- For datasets with 1000-100,000 columns, this search space is enormous
- Instead we adopt a novel row/sample enumeration algorithm for this purpose. **CARPENTER (SIGKDD'03)** is the **FIRST** algorithm which adopt this approach

# Column/Item Enumeration Lattice

- Each nodes in the lattice represent a combination of columns/items
- An edge exists from node A to B if A is subset of B and A differ from B by only 1 column/item
- Search can be done **breadth first**

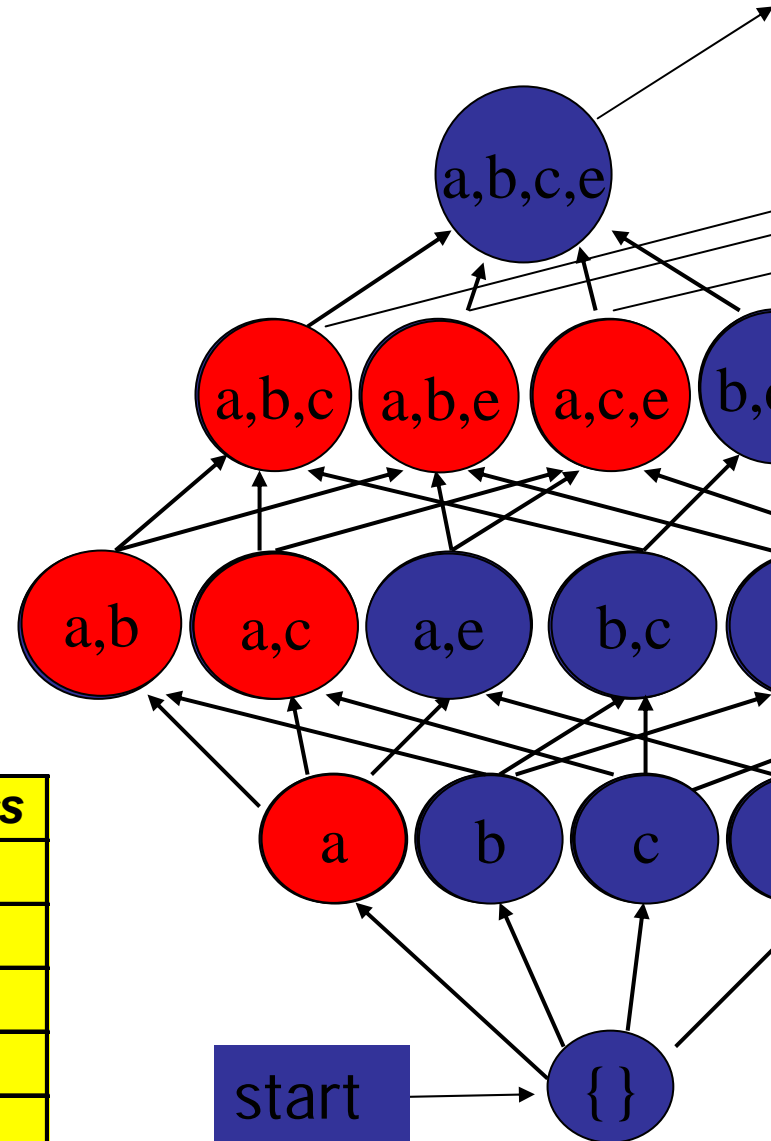
<i>i</i>	<i>ri</i>	<i>Class</i>
1	a,b,c,l,o,s	C
2	a,d,e,h,p,l,r	C
3	a,c,e,h,o,q,t	C
4	a,e,f,h,p,r	~C
5	b,d,f,g,l,q,s,t	~C



# Column/Item Enumeration Lattice

- Each nodes in the lattice represent a combination of columns/items
- An edge exists from node A to B if A is subset of B and A differ from B by only 1 column/item
- Search can be done **depth first**
- Keep edges from parent to child only if child is the prefix of parent

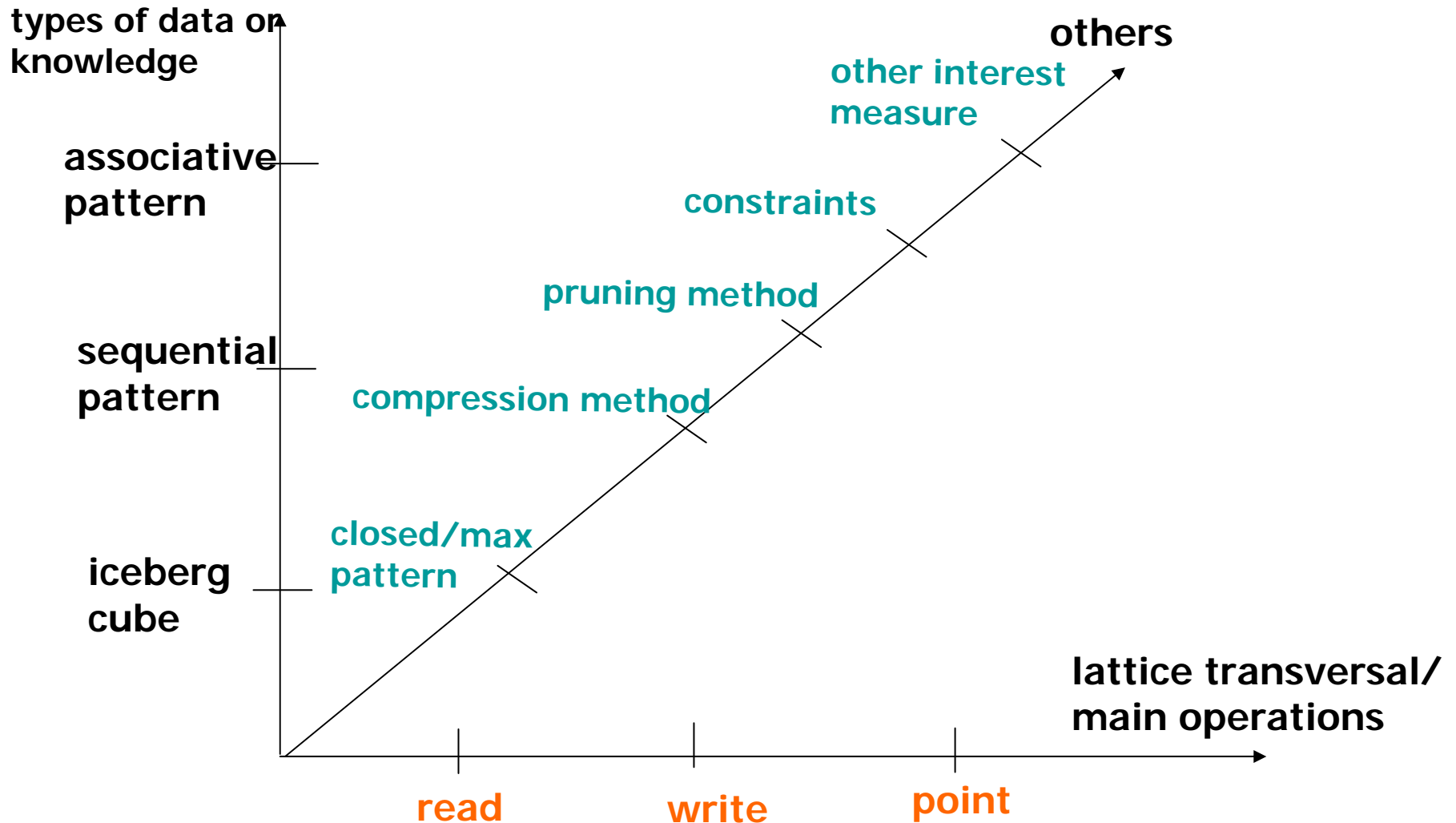
<i>i</i>	<i>ri</i>	<i>Class</i>
1	<i>a,b,c,l,o,s</i>	<i>C</i>
2	<i>a,d,e,h,p,l,r</i>	<i>C</i>
3	<i>a,c,e,h,o,q,t</i>	<i>C</i>
4	<i>a,e,f,h,p,r</i>	<i>~C</i>
5	<i>b,d,f,g,l,q,s,t</i>	<i>~C</i>



# General Framework for Column/Item Enumeration

	Read-based	Write-based	Point-based
Association Mining	Apriori[AgSr94], DIC	Eclat, MaxClique[Zaki01], FPGrowth [HaPe00]	Hmine
Sequential Pattern Discovery	GSP[AgSr96]	SPADE [Zaki98,Zaki01], PrefixSpan [PHPC01]	
Iceberg Cube	Apriori[AgSr94]		BUC[BeRa99], H- Cubing [HPDW01]

# A Multidimensional View





# Sample/Row Enumeration Algorithms

- To avoid searching the large column/item enumeration space, our mining algorithm search for patterns/rules in the **sample/row enumeration space**
- Our algorithms does not fitted into the column/item enumeration algorithms
- They are not YAARMA (Yet Another Association Rules Mining Algorithm)
- Column/item enumeration algorithms simply does not scale for microarray datasets

# Existing Row/Sample Enumeration Algorithms

- CARPENTER(SIGKDD'03)
  - Find closed patterns using row enumeration
- FARMER(SIGMOD'04)
  - Find interesting rule groups and building classifiers based on them
- COBBLER(SSDBM'04)
  - Combined row and column enumeration for tables with large number of rows and columns
- FARMER's demo (VLDB'04)
- Balance the scale: 3 row enumeration algorithms vs >50 column enumeration algorithms

# Concepts of CARPENTER

<i>i</i>	<i>ri</i>	Class
1	a,b,c,l,o,s	C
2	a,d,e,h,p,l,r	C
3	a,c,e,h,o,q,t	C
4	a,e,f,h,p,r	~C
5	b,d,f,g,l,q,s,t	~C

Example Table

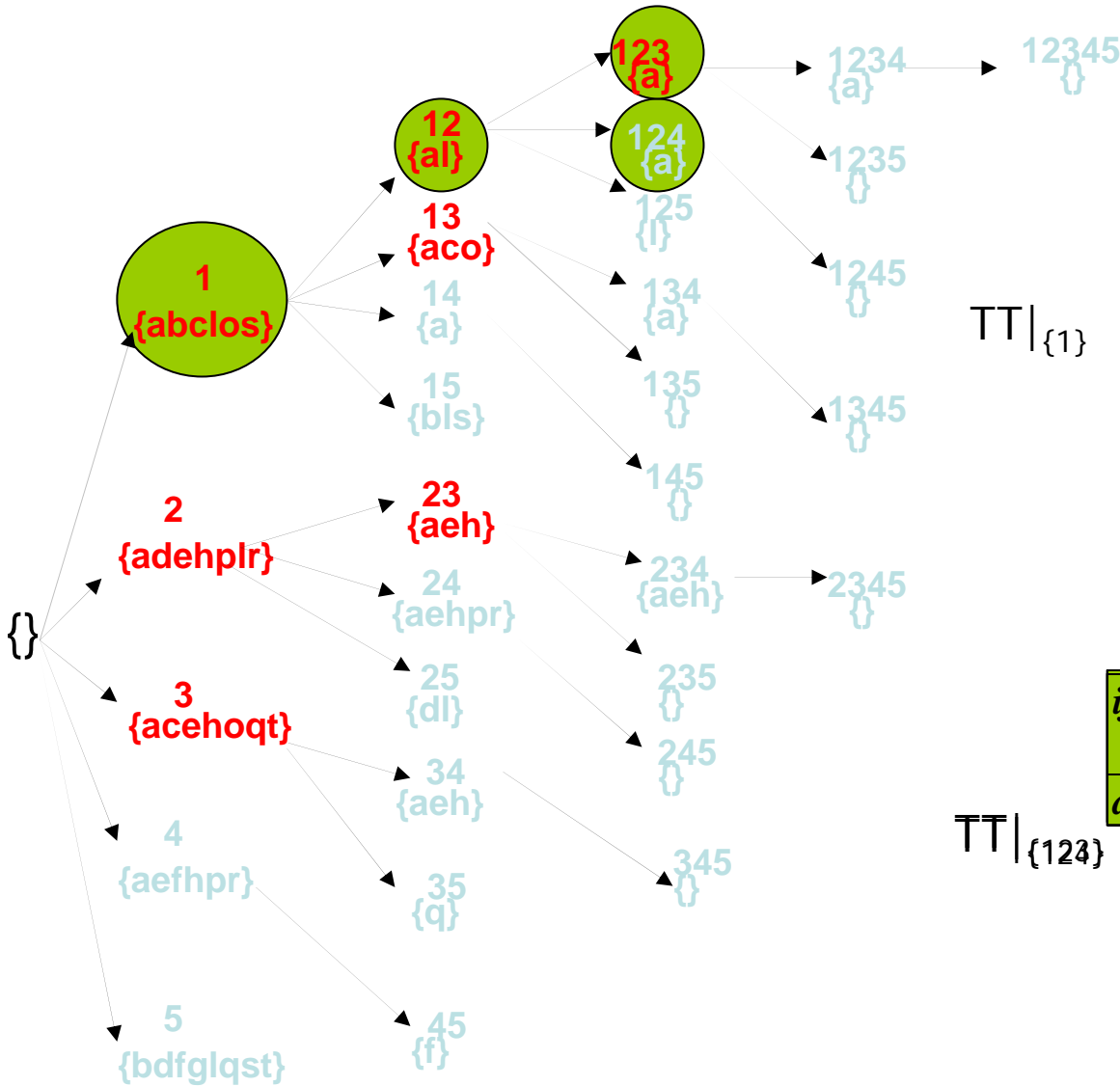
<i>ij</i>	<i>R (ij )</i>	
	<i>C</i>	<i>~C</i>
<i>a</i>	1,2,3	4
<i>b</i>	1	5
<i>c</i>	1,3	
<i>d</i>	2	5
<i>e</i>	2,3	4
<i>f</i>		4,5
<i>g</i>		5
<i>h</i>	2,3	4
<i>l</i>	1,2	5
<i>o</i>	1,3	
<i>p</i>	2	4
<i>q</i>	3	5
<i>r</i>	2	4
<i>s</i>	1	5
<i>t</i>	3	5

Transposed Table, TT

	<i>C</i>	<i>~C</i>
<i>a</i>	1,2,3	4
<i>e</i>	2,3	4
<i>h</i>	2,3	4

TT|<sub>{2,3}</sub>

# Row Enumeration



$TT|_{\{1\}}$

$ij$	$R(ij)$	
	$C$	$\sim C$
$a$	1,2,3	4
$b$	1	5
$c$	1,3	
$l$	1,2	5
$o$	1,3	
$s$	1	5

$TT|_{\{123\}}$

$ij$	$R(ij)$	
	$C$	$\sim C$
$a$	1,2,3	4

$ij$	$R(ij)$	
	$C$	$\sim C$
$a$	1,2,3	4
$b$	1	5
$c$	1,3	
$d$	2	5
$e$	2,3	4
$f$		4,5
$g$		5
$h$	2,3	4
$l$	1,2	5
$o$	1,3	
$p$	2	4
$q$	3	5
$r$	2	4
$s$	1	5
$t$	3	5

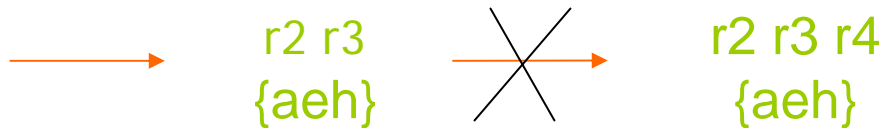
$ij$	$R(ij)$	
	$C$	$\sim C$
$a$	1,2,3	4
$l$	1,2	5

$TT|_{\{12\}}$

# Pruning Method 1

- Removing rows that appear in all tuples of transposed table will not affect results

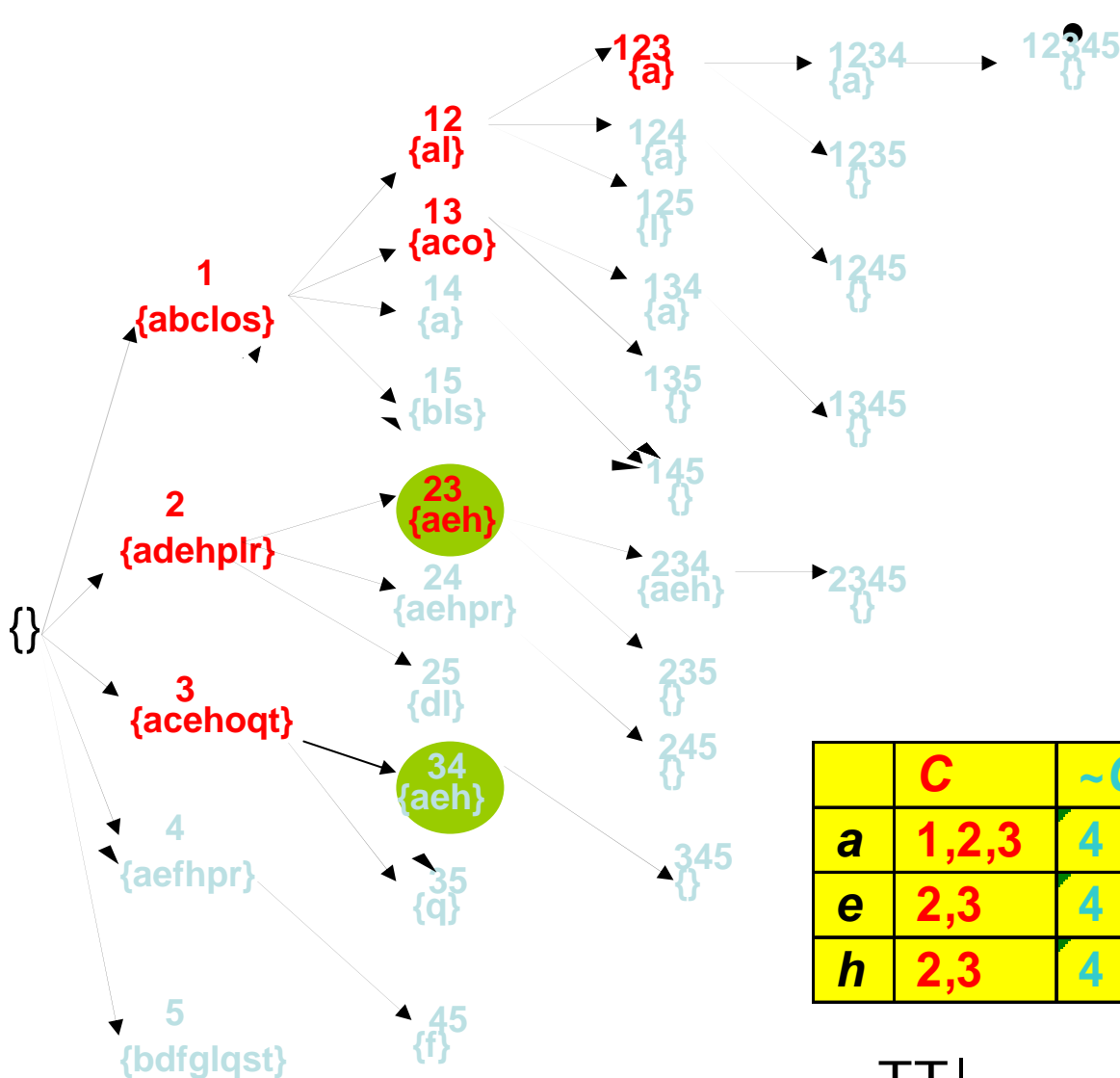
	C	~C
a	1,2,3	4
e	2,3	4
h	2,3	4



TT|<sub>{2,3}</sub>

r4 has 100% support in the conditional table of “r2r3”, therefore branch “r2 r3r4” will be pruned.

# Pruning method 2



if a rule is discovered before, we can prune enumeration below this node

- Because all rules below this node has been discovered before
- For example, at node 34, if we found that {aeh} has been found, we can prune off all branches below it

	C	~C
a	1,2,3	4
e	2,3	4
h	2,3	4

TT|<sub>{3,4}</sub>

# Pruning Method 3: Minimum Support

- Example: From  $TT|_{\{1\}}$ , we can see that the support of all possible pattern below node  $\{1\}$  will be at most 5 rows.

$TT|_{\{1\}}$

$ij$	$R(ij)$	
	$C$	$\sim C$
$a$	1,2,3	4
$b$	1	5
$c$	1,3	
$l$	1,2	5
$o$	1,3	
$s$	1	5

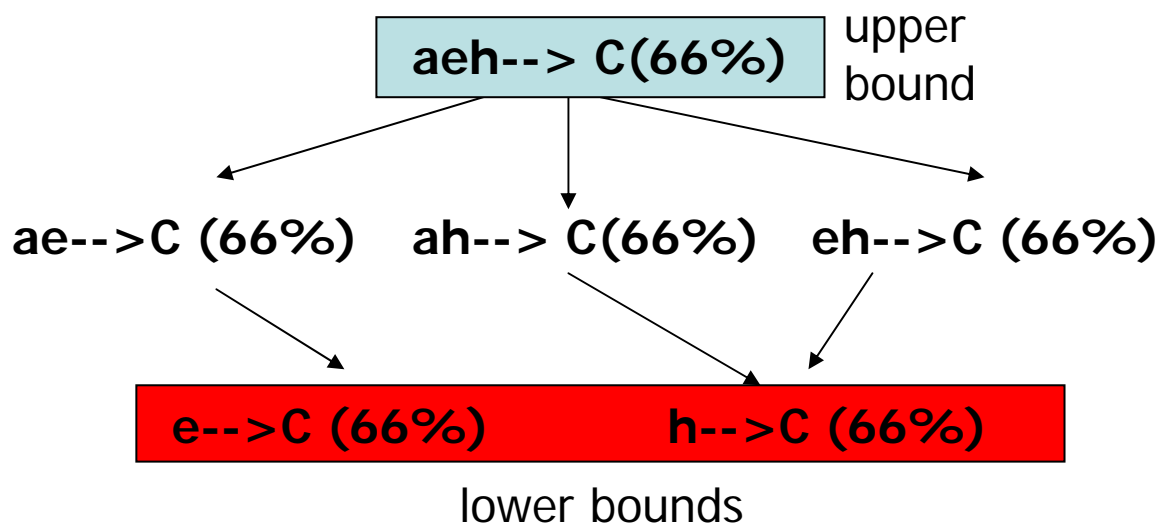
# From CARPENTER to FARMER

- What if classes exists ? What more can we do ?
- Pruning with Interestingness Measure
  - Minimum confidence
  - Minimum chi-square
- Generate lower bounds for classification/prediction



# Interesting Rule Groups

- Concept of a **rule group/equivalent class**
  - rules supported by exactly the same set of rows are grouped together
- Example: the following rules are derived from row 2,3 and 4 with 66% confidence



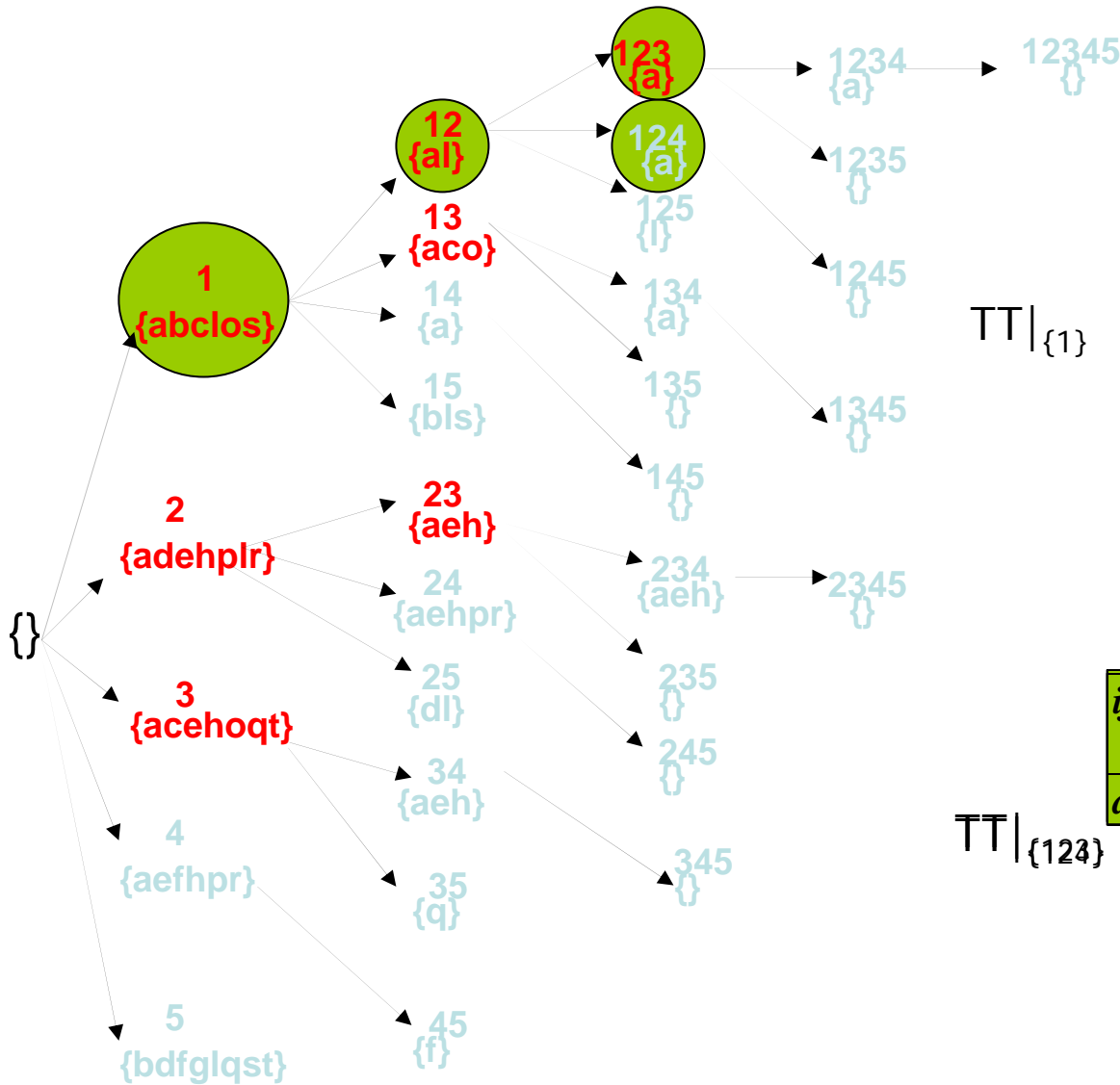
<i>i</i>	<i>ri</i>	<i>Class</i>
1	<b>a,b,c,l,o,s</b>	<b>C</b>
2	<b>a,d,e,h,p,l,r</b>	<b>C</b>
3	<b>a,c,e,h,o,q,t</b>	<b>C</b>
4	<b>a,e,f,h,p,r</b>	<b>~C</b>
5	<b>b,d,f,g,l,q,s,t</b>	<b>~C</b>

a-->C however is not in the group

# Pruning by Interestingness Measure

- In addition, find only interesting rule groups (IRGs) based on some measures:
  - **minconf**: the rules in the rule group can predict the class on the RHS with high confidence
  - **minchi**: there is high correlation between LHS and RHS of the rules based on chi-square test
- Other measures like lift, entropy gain, conviction etc. can be handle similarly

# Ordering of Rows: All Class C before ~C



<i>ij</i>	<i>R</i> ( <i>ij</i> )	
	<i>C</i>	<i>~C</i>
<i>a</i>	1,2,3	4
<i>b</i>	1	5
<i>c</i>	1,3	
<i>d</i>	2	5
<i>e</i>	2,3	4
<i>f</i>		4,5
<i>g</i>		5
<i>h</i>	2,3	4
<i>l</i>	1,2	5
<i>o</i>	1,3	
<i>p</i>	2	4
<i>q</i>	3	5
<i>r</i>	2	4
<i>s</i>	1	5
<i>t</i>	3	5

TT|<sub>{1}</sub>

<i>ij</i>	<i>R</i> ( <i>ij</i> )	
	<i>C</i>	<i>~C</i>
<i>a</i>	1,2,3	4
<i>b</i>	1	5
<i>c</i>	1,3	
<i>l</i>	1,2	5
<i>o</i>	1,3	
<i>s</i>	1	5

<i>ij</i>	<i>R</i> ( <i>ij</i> )	
	<i>C</i>	<i>~C</i>
<i>a</i>	1,2,3	4
<i>l</i>	1,2	5

TT|<sub>{12}</sub>

<i>ij</i>	<i>R</i> ( <i>ij</i> )	
	<i>C</i>	<i>~C</i>
<i>a</i>	1,2,3	4

TT|<sub>{123}</sub>

<i>ij</i>	<i>R</i> ( <i>ij</i> )	
	<i>C</i>	<i>~C</i>
<i>a</i>	1,2,3	4

# Pruning Method: Minimum Confidence

- Example: In  $TT|_{\{2,3\}}$  on the right, the maximum confidence of all rules below node  $\{2,3\}$  is at most  $4/5$

	<b>C</b>	<b>~C</b>
<b>a</b>	<b>1,2,3,6</b>	<b>4,5</b>
<b>e</b>	<b>2,3,7</b>	<b>4,9</b>
<b>h</b>	<b>2,3</b>	<b>4</b>

$TT|_{\{2,3\}}$

# Pruning method: Minimum chi-square

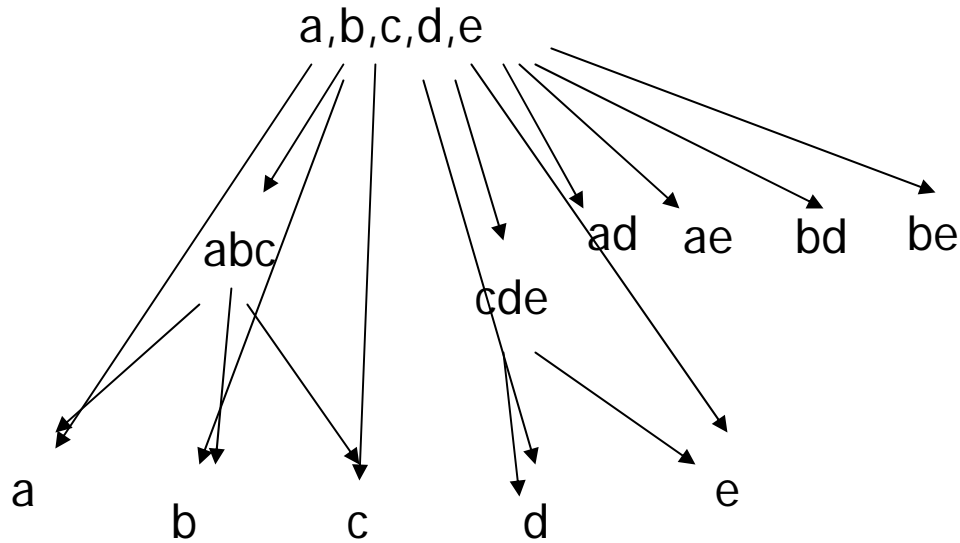
- Same as in computing maximum confidence

	<b>C</b>	<b>~C</b>
<b>a</b>	<b>1,2,3,6</b>	<b>4,5</b>
<b>e</b>	<b>2,3,7</b>	<b>4,9</b>
<b>h</b>	<b>2,3</b>	<b>4</b>

	C	~C	Total
A	<b>max=5</b>	<b>min=1</b>	Computed
~A	Computed	Computed	Computed
	<b>Constant</b>	<b>Constant</b>	<b>Constant</b>

$TT|_{\{2,3\}}$

# Finding Lower Bound, MineLB



- Example: An upper bound rule with antecedent  $A=abcde$  and two rows ( $r1 : abc$ ) and ( $r2 : cde$ )
- Initialize lower bounds  $\{a, b, c, d, e\}$
- add “**abc**f”--- new lower  $\{d, e\}$
- Add “**cde**g”--- new lower bound  $\{ad, bd, ae, be\}$

Candidate lower bound: a, b, c, d, e, cd, ce

Removed since,  $a, b, c, d, e$  is better

# Implementation

- In general, CARPENTER FARMER can be implemented in many ways:
  - FP-tree
  - Vertical format
- For our case, we assume the dataset can be fitted into the main memory and used pointer-based algorithm similar to BUC

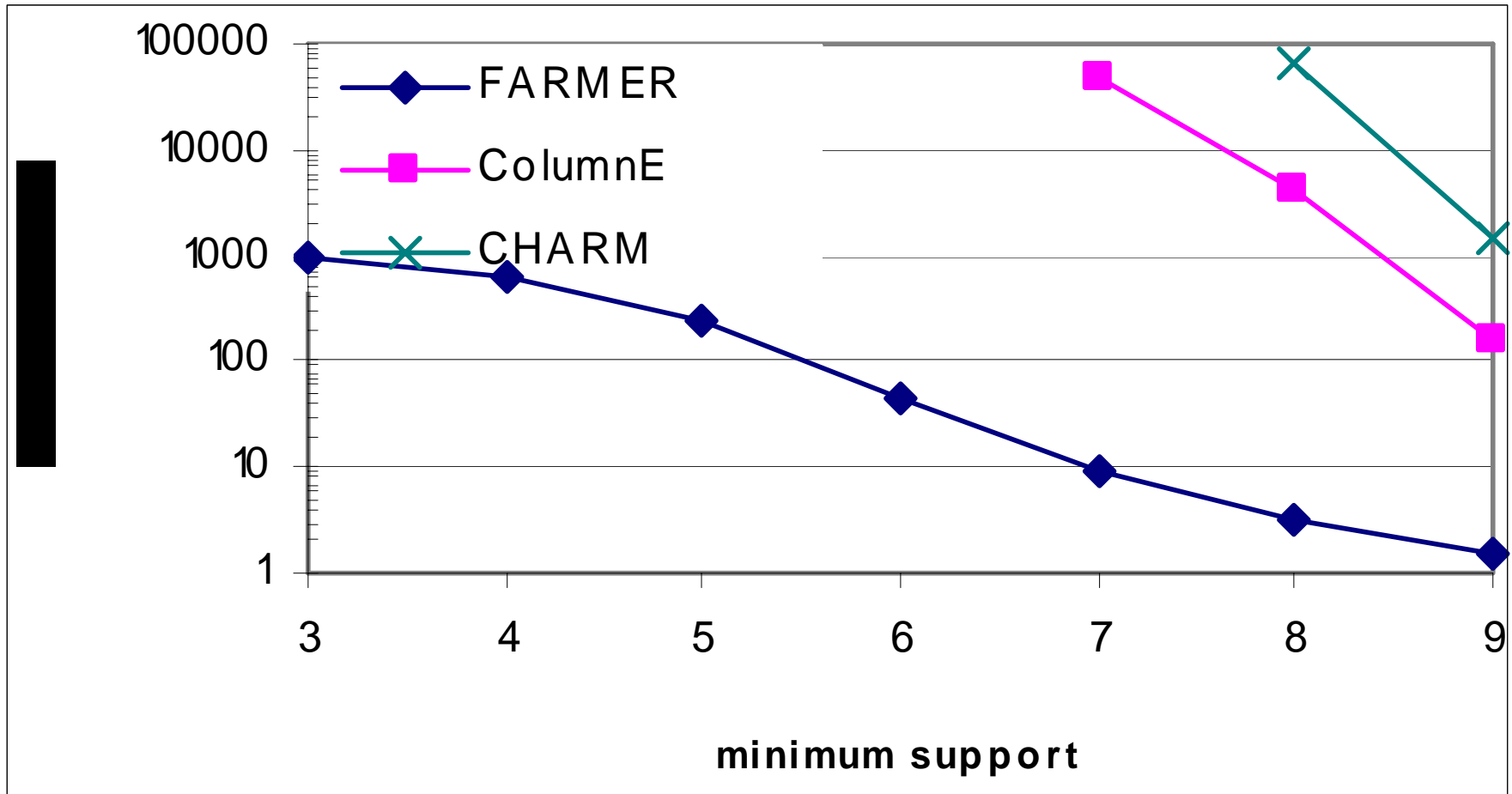
<i>ij</i>	<i>R (ij )</i>	
	<i>C</i>	<i>~C</i>
<i>a</i>	1,2,3	4
<i>b</i>	1	5
<i>c</i>	1,3	
<i>d</i>	2	5
<i>e</i>	2,3	4
<i>f</i>		4,5
<i>g</i>		5
<i>h</i>	2,3	4
<i>l</i>	1,2	5
<i>o</i>	1,3	
<i>p</i>	2	4
<i>q</i>	3	5
<i>r</i>	2	4
<i>s</i>	1	5
<i>t</i>	3	5

# Experimental studies

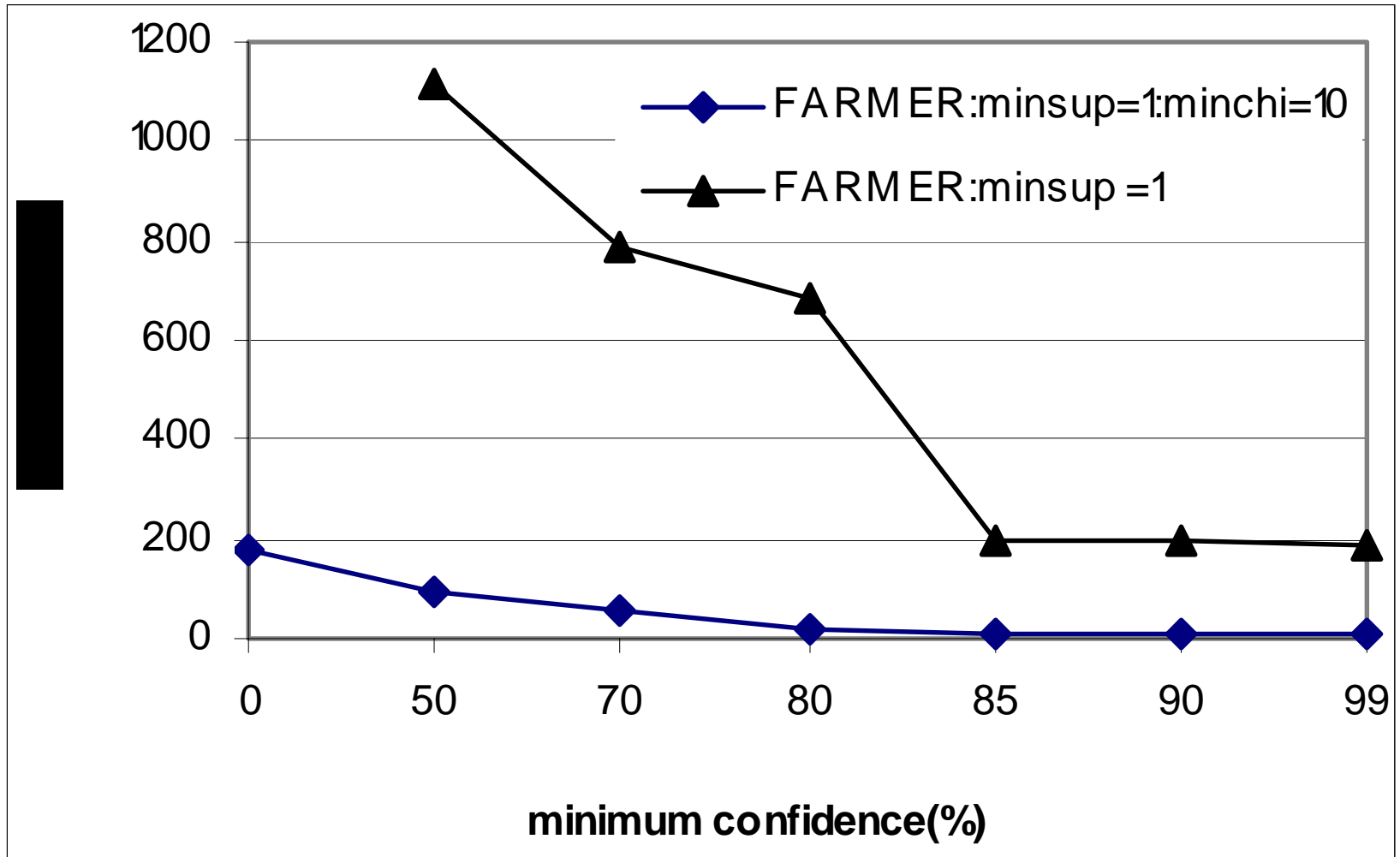
- Efficiency of FARMER
  - On five real-life dataset
    - lung cancer (LC), breast cancer (BC) , prostate cancer (PC), ALL-AML leukemia (ALL), Colon Tumor(CT)
  - Varying minsup, minconf, minchi
  - Benchmark against
    - CHARM [ZaHs02] ICDM'02
    - Bayardo's algorithm (ColumE) [BaAg99] SIGKDD'99
- Usefulness of IRGs
  - Classification



# Example results--Prostate



# Example results--Prostate



# Naive Classification Approach

- Generate the upper bounds of IRGs
- Rank the upper bounds, thus ranking the IRGs;
- Apply coverage pruning on the IRGs;
- Predict the test data based on the IRGs that it covers.

# Classification results

Dataset	# features	# training	# test	IRG Classifier	CBA	SVM
BC	619	78	19	<b>78.95%</b>	57.89%	36.84%
LC	2173	32	149	89.93%	81.88%	<b>96.64%</b>
CT	135	47	15	<b>93.33%</b>	73.33%	73.33%
PC	1554	102	34	<b>88.24%</b>	82.35%	79.41%
ALL	866	38	34	64.71%	91.18%	<b>97.06%</b>
Average Accuracy				<b>83.03%</b>	77.33%	76.66%

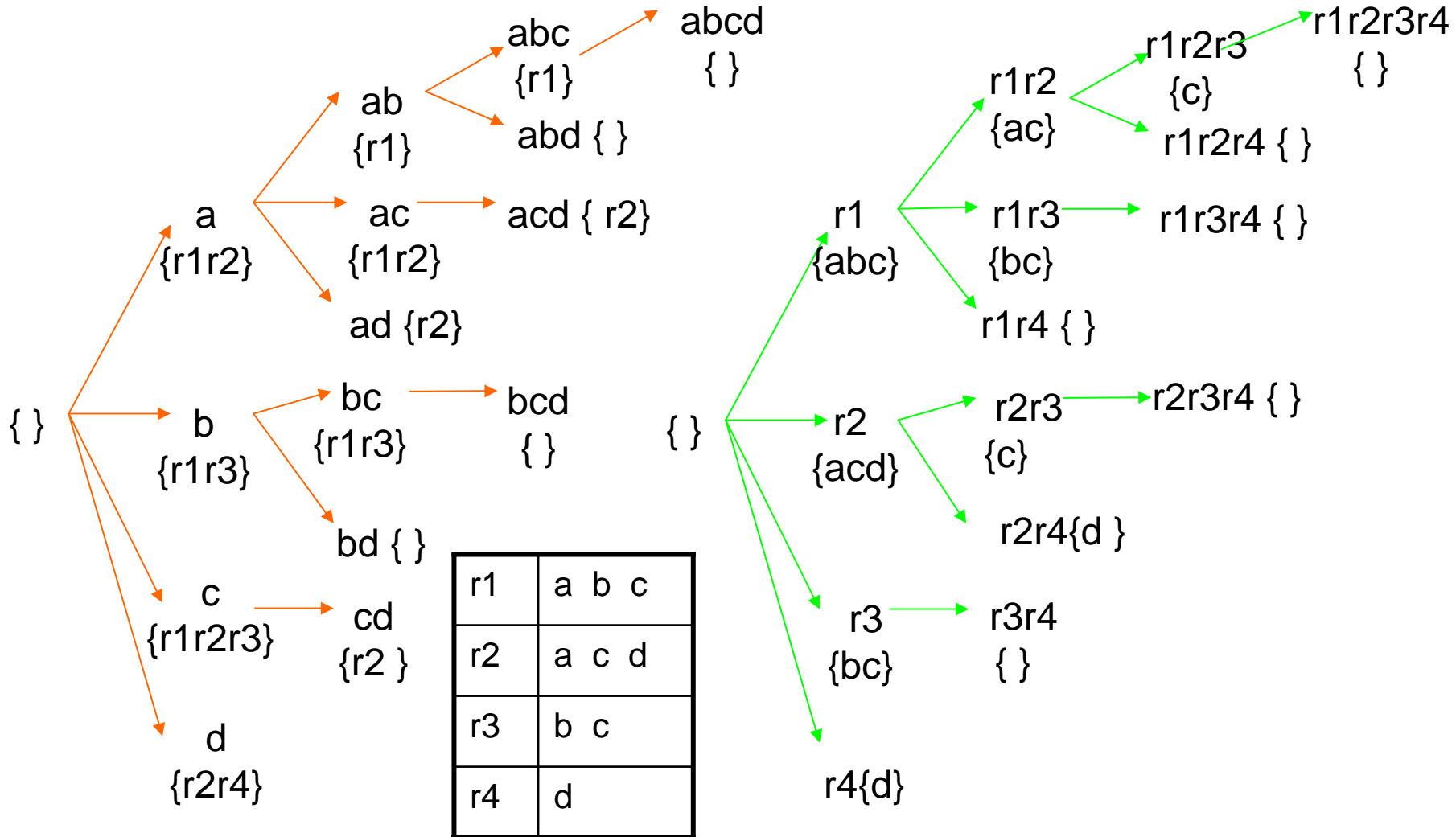
# Summary of Experiments

- FARMER is much more efficient than existing algorithms
- There are evidences to show that IRGs is useful for classification of microarray datasets

# COBBLER: Combining Column and Row Enumeration

- Extend CARPENTER to handle datasets with both large number of columns and rows
- Switch dynamically between column and row enumeration based on estimated cost of processing

# Single Enumeration Tree



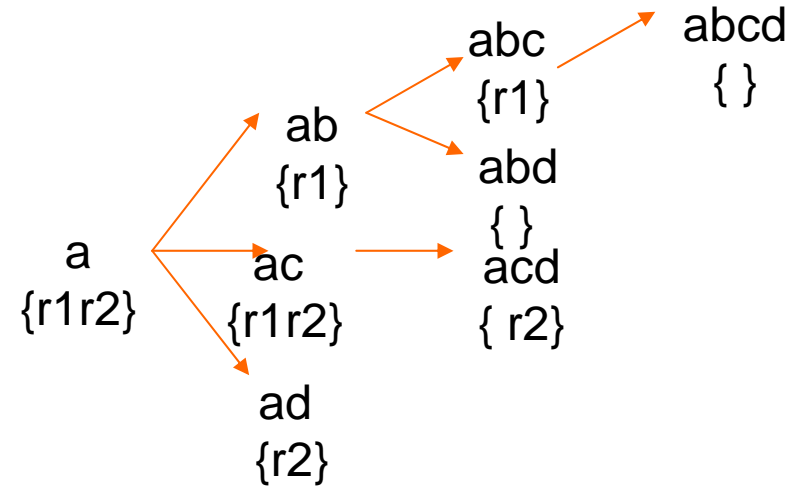
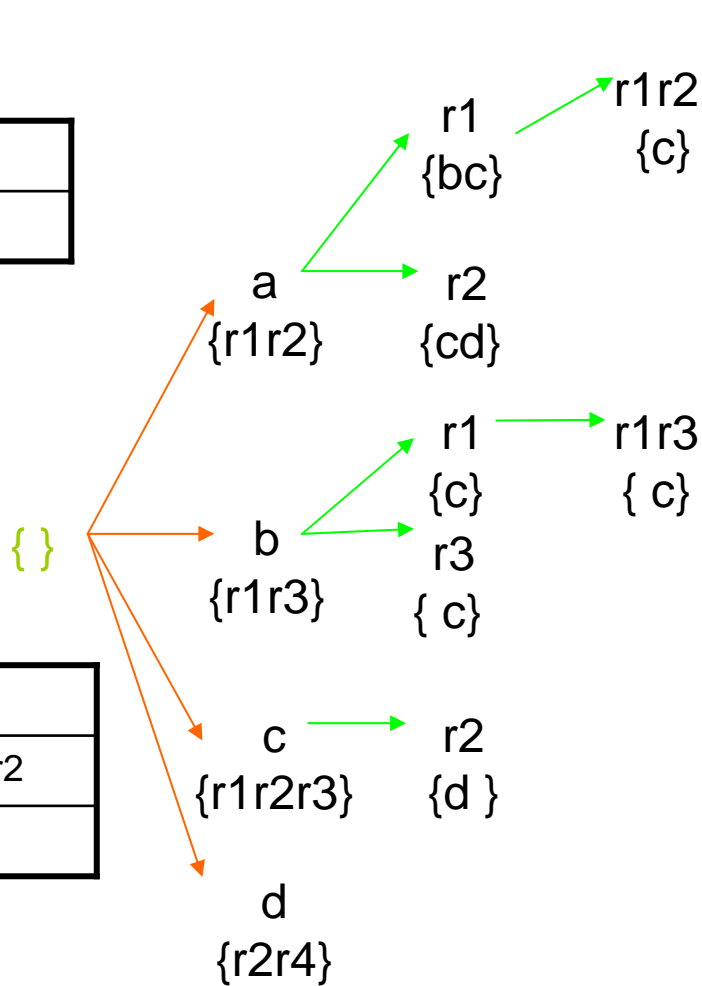
Feature enumeration

Row enumeration

# Dynamic Enumeration Tree

r1	bc
r2	cd

b	r1
c	r1 r2
d	r2



abc: {r1}

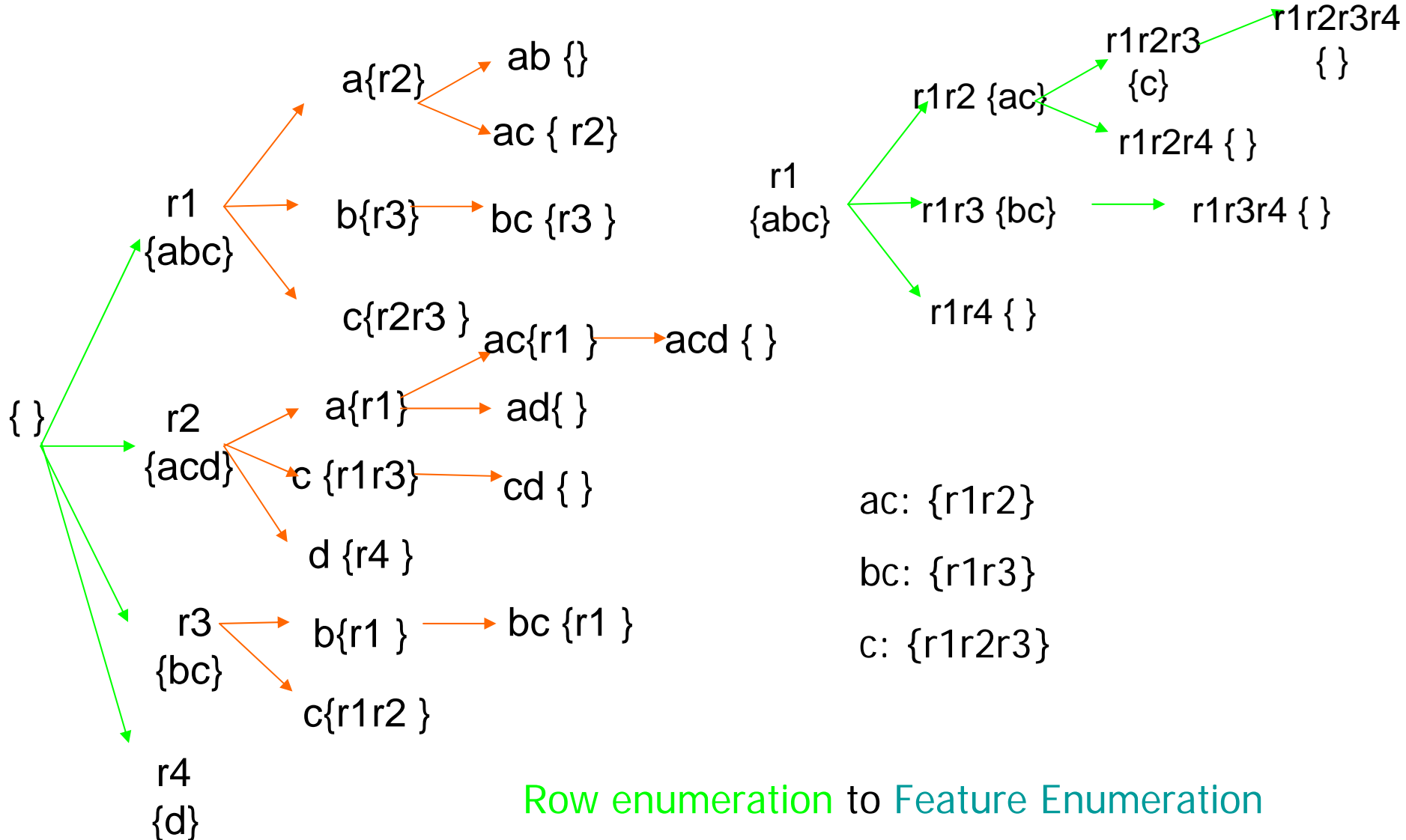
ac: {r1r2}

acd: {r2}

Feature enumeration to Row enumeration



# Dynamic Enumeration Tree

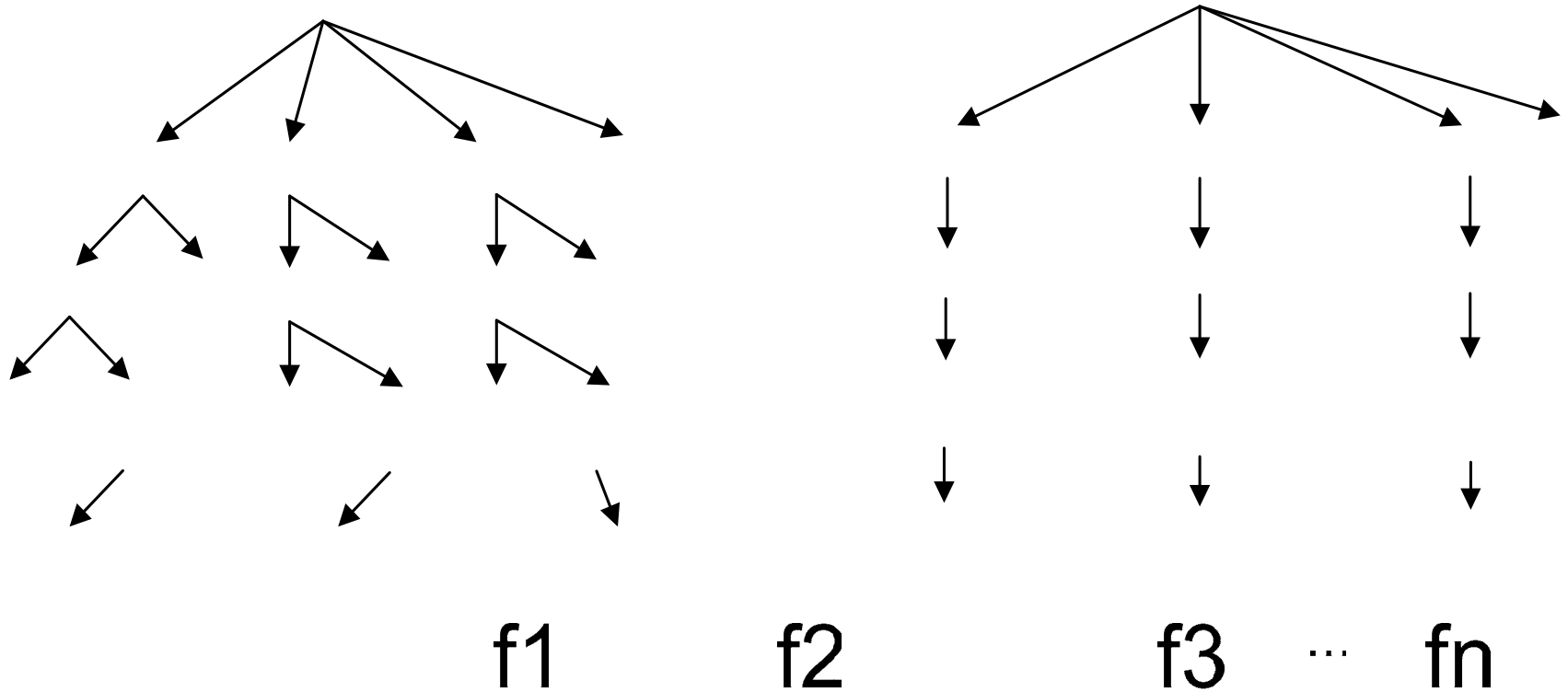


Row enumeration to Feature Enumeration

# Switching Condition

- Naïve idea of switching based on row number and feature number does not work well
- to estimate the required computation of an enumeration sub-tree, i.e., row enumeration sub-tree or feature enumeration sub-tree.
  - Estimate the maximal level of enumeration for each children sub-tree
- Example of estimating the maximal level of enumeration:
  - Suppose  $r=10$ ,  $S(f_1)=0.8$ ,  $S(f_2)=0.5$ ,  $S(f_3)=0.5$ ,  $S(f_4)=0.3$  and  $minsup=2$
  - $S(f_1)*S(f_2)*S(f_3)*r = 2 \geq minsup$
  - $S(f_1)*S(f_2)*S(f_3)*S(f_4)*r = 0.6 < minsup$
  - Then the estimated deepest node under  $f_1$  is  $f_1f_2f_3$

# Switching Condition



$\{f_1, f_2\} \dots \{f_1, f_n\} \{f_2, f_3\} \dots \{f_2, f_n\} \{f_3, f_4\} \dots \{f_3, f_n\}$

$\{f_1, f_2, f_3\} \dots \{f_2, f_3, f_4\} \dots \{f_3, f_4, f_5\} \dots$

# Switching Condition

- $F = \{f_1, f_2, \dots, f_n\}$ .
- $\mathcal{S}(f_j, T|_F)$ , the frequency of feature  $f_j$  in  $T|_F$ .
- $r = |\mathcal{R}(F)|$ , the number of rows conditional table  $T|_F$  contains.

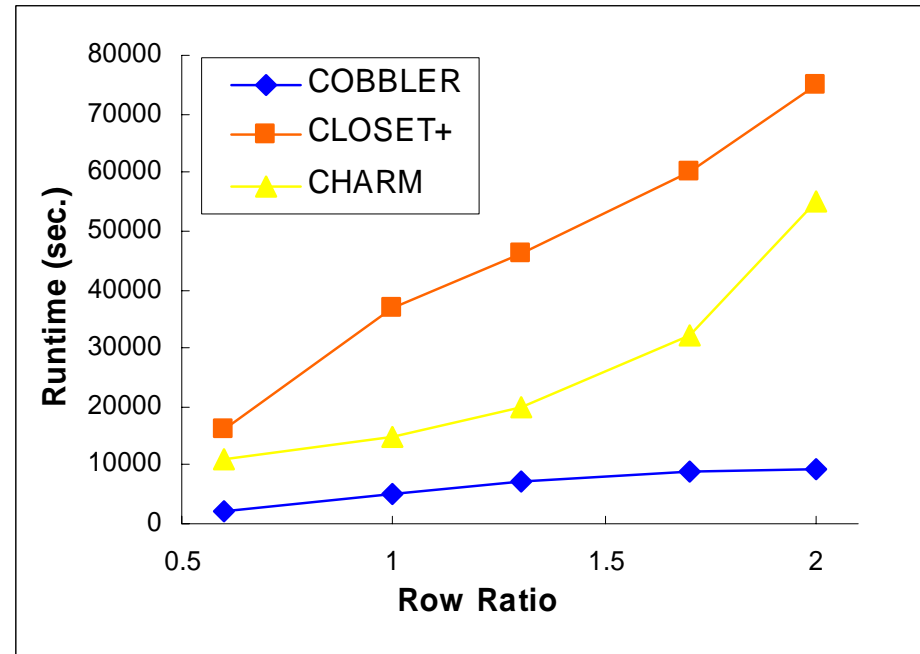
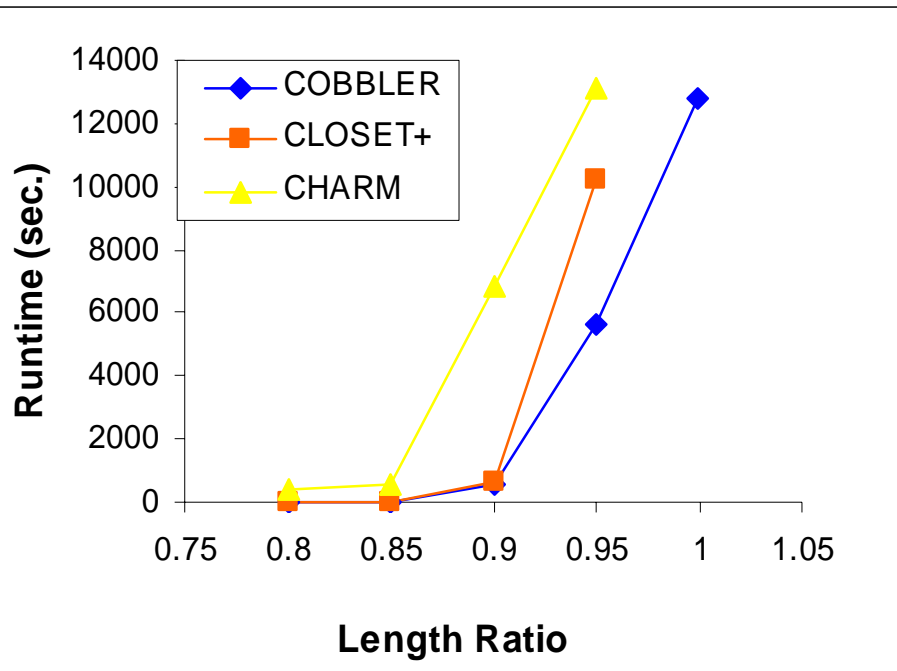
To estimate for a node:

To estimate for a path:

$$\mathcal{L}(N_i) = \sum_{k=1}^h (r \cdot RowProcessTime \cdot \prod_{j=1}^k \mathcal{S}(f_j, T|_{F'_i})).$$

To sum up estimation of all paths as the final estimation

# Length and Row ratio



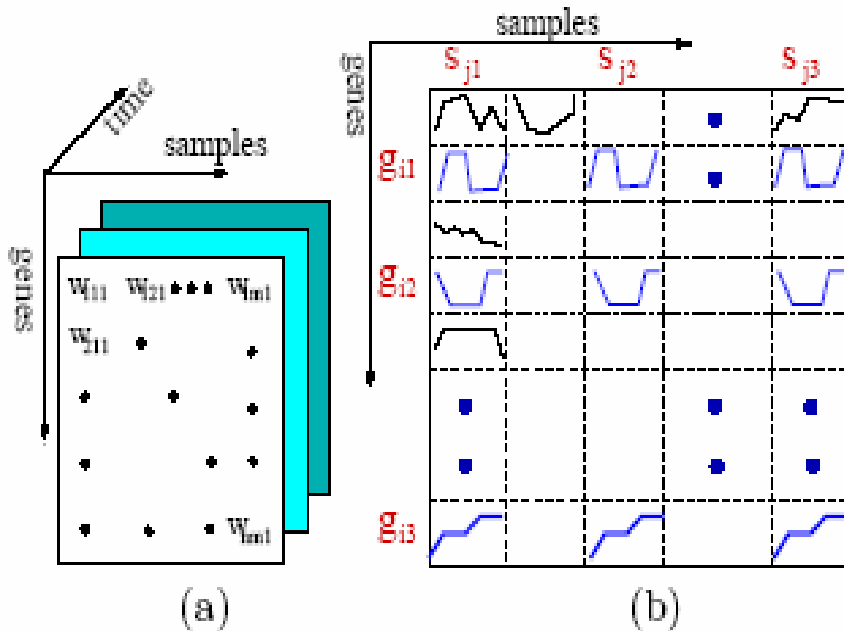
Synthetic data

# Extension of our work by other groups (with or without citation)

- **[1]** [\*Using transposition for pattern discovery from microarray data\*](#), Francois Rioult (GREYC CNRS), Jean-Francois Boulicaut (INSA Lyon), Bruno Cremileux (GREYC CNRS), Jeremy Besson (INSA Lyon)
- See the presence and absence of genes in the sample as a binary matrix. Perform a transposition of the matrix which is essentially our transposed table. Enumeration methods are the same otherwise.

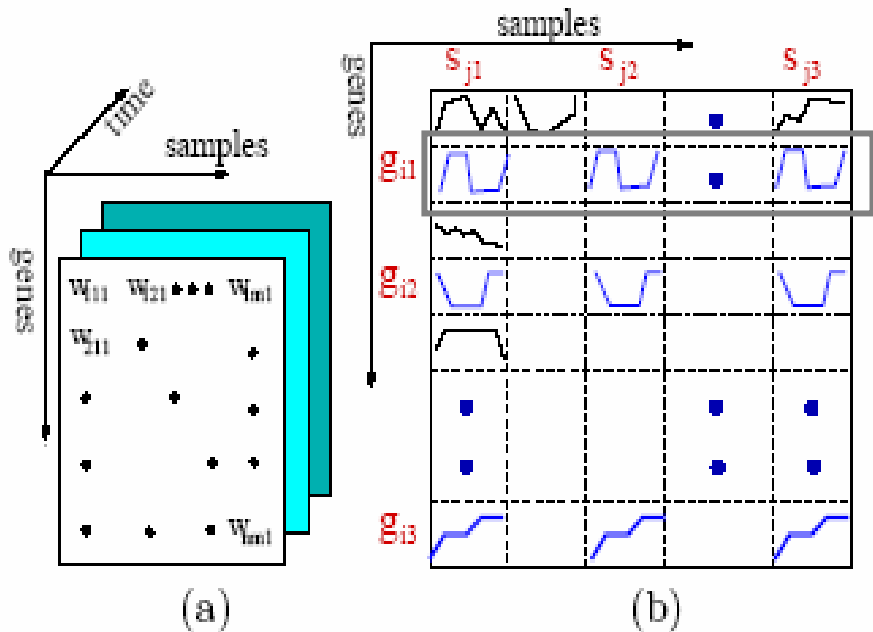
# Extension of our work by other groups (with or without citation) II

- [2] [Mining Coherent Gene Clusters from Gene-Sample-Time Microarray Data](#). D. Jiang, **Jian Pei**, M. Ramanathan, C. Tang and A. Zhang. (**Industrial full paper**, Runner-up for the best application paper award). *SIGKDD'2004*



	Gene1	Gene 2	Gene3	Gene 4
Sample1				
Sample2				
.				
.				
.				
SampleN-1				
SampleN				

# Extension of our work by other groups (with or without citation) III



A gene in two samples are say to be coherent if their time series satisfied a certain matching condition

	Gene 1	Gene 2	Gene 3	Gene 4
S1		1.23		
S2		1.34		
.				
.				
.				
SN-1		1.52		
SN				

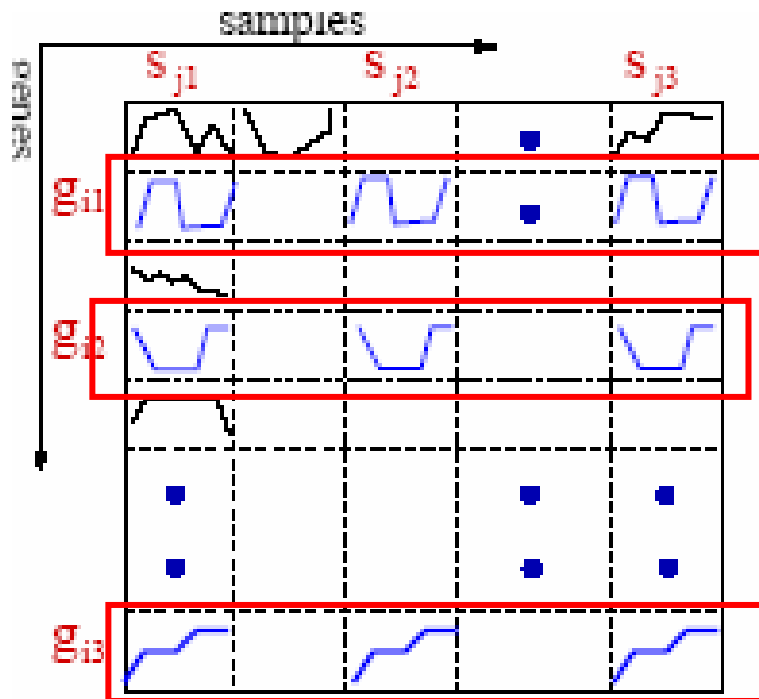
In CARPENTER, a gene in two samples are say to be matching if their expression in the two samples are almost the same



# Extension of our work by other groups (with or without citation) IV

[2] Try to find a subset of samples  $S$  such that a subset of genes  $G$  is coherent for each pair of samples in  $S$ .

$$|S| > \min_s, |G| > \min_g$$

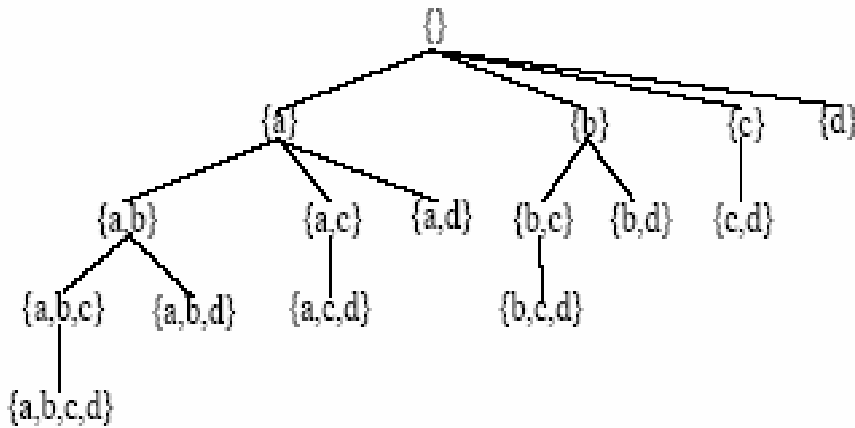


In CARPENTER, we try to find a subset of samples  $S$  in which a subset of genes  $G$  is similar in expression level for each pair of samples in  $S$ .

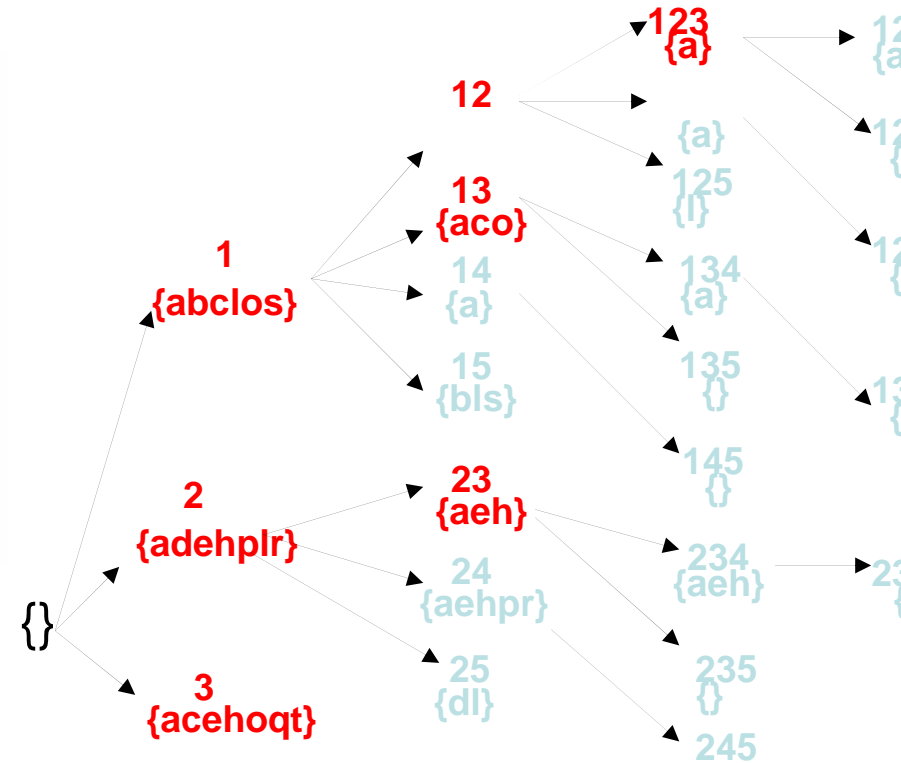
$$|S| > \min_s, |G| > 0$$

	Gene1	Gene2	Gene 3	Gene4
S1		1.23		
S2		1.34		
.				
.				
.				
SN-1		1.52		
SN				

# Extension of our work by other groups (with or without citation) V



[2] Perform sample-wise enumeration and remove genes that are not pairwise coherent across the samples enumerated



CARPENTER: Perform sample-wise enumeration and remove genes that does not have the same expression level across the samples enumerated

# Extension of our work by other groups (with or without citation) VI

From [2]: Pruning Rule 3.1  
(Pruning small sample sets). *At a node  $v = fs1 ; : ; : ; s i k g$ , the subtree of  $v$  can be pruned if  $(k + jTailj) < min_s$*

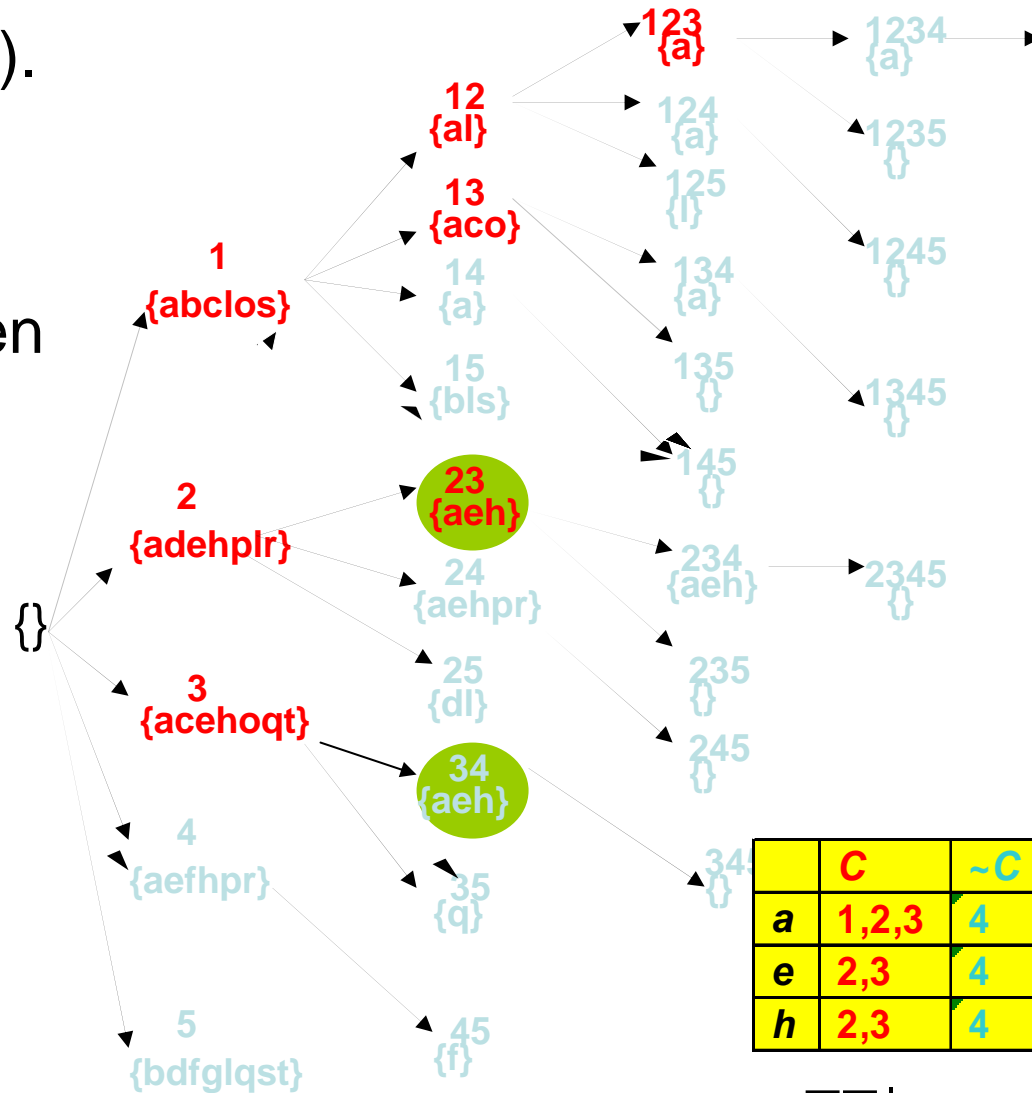
$TT|_{\{1\}}$

$ij$	$R(ij)$	
	$C$	$\sim C$
$a$	1,2,3	4
$b$	1	5
$c$	1,3	
$l$	1,2	5
$o$	1,3	
$s$	1	5

- Pruning Method 3 in CARPENTER:  
From  $TT|_{\{1\}}$ , we can see that the support of all possible pattern below node  $\{1\}$  will be at most 5 rows.

# Extension of our work by other groups (with or without citation) VII

- [2] Pruning Rule 3.2  
(Pruning subsumed sets).  
*At a node  $v = \{s_i \dots s_{ik}\}$  if  $\{s_{i1}, \dots, s_{ik}\} \cup Tail$  is a subset of some maximal coherent sample set, then the subtree of the node can be pruned.*

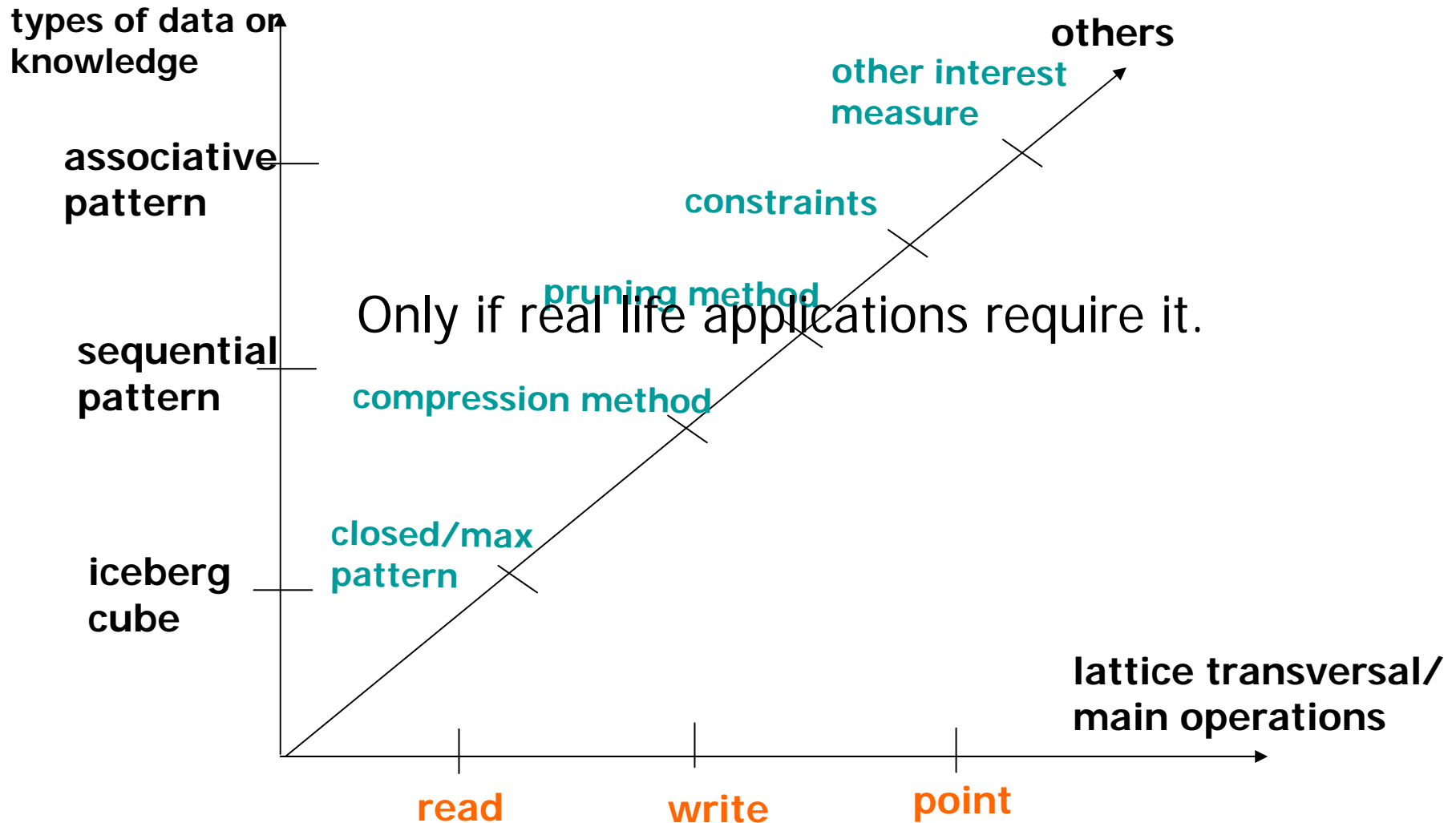


- CARPENTER Pruning Method 2: if a rule is discovered before, we can prune enumeration below this node

# Extension of our work (Conclusion)

- The sample/enumeration framework had been successfully adopted by other groups in mining microarray datasets
- We are proud of our contribution as the group that produced the first row/sample enumeration algorithm CARPENTER and is happy that other groups also find the method useful
- However, citations from these groups would have been nice. After all academic integrity is the most important thing for a researcher.

# Future Work: Generalize Framework for Row Enumeration Algorithms?



# Conclusions

- Many datasets in bioinformatics have very different characteristics compared to those that has been previously studied
- These characteristics can either work against you or for you
- In the case of microarray datasets with large number columns but small number of rows/samples, we turn what is against us to our advantage
  - Row/Sample enumeration
  - Pruning strategy
- We show how our methods have been modified by other groups to produce useful algorithm for mining microarray datasets

Thank you!!!

[atung@comp.nus.edu.sg](mailto:atung@comp.nus.edu.sg)

[www.comp.nus.edu.sg/~atung/sfu\\_talk.pdf](http://www.comp.nus.edu.sg/~atung/sfu_talk.pdf)