# High Confidence Rule Mining for Microarray Analysis

Kang Deng

University of Alberta

UNIVERSITY OF ALBERTA

---

## Outline

- Introduction
- Row Enumeration
- Confidence-based Prune Strategy
- MAXCONF Algorithm
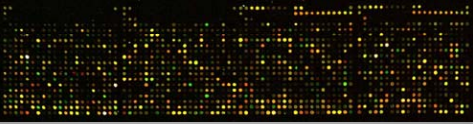- Evaluation
- References

---

"High Confidence Rule Mining for Microarray Analysis", by Tara McIntosh, Sanjay Chawla, 2006

- 27 pages
- 14 definitions
- 2 lemmas
- 4 tables
- Figures, formulas, etc

Confidence-based Strategy



---

## What is Microarrays?

- A DNA microarray is a collection of microscopic DNA spots, commonly representing single genes, arrayed on a solid surface by covalent attachment to a chemical matrix.

~~genes~~   items

~~samples~~
transactions

## Our Task

- **One main objective of molecular biology is to develop a deeper understanding of how genes are functionally related.**

~~minimum support~~          minimum confidence

**Minimum Support = 30%, Minimum Confidence = 80%**

$\text{GENE1} \Rightarrow \text{GENE2 (support 10\%, confidence 90\%)}$

**We do not mine association rules, but Confidence Rules.**

---

## Row Enumeration   Explosive increase of candidates

- Traditional Dataset
- Microarray Dataset



items

1-length
2-length
3-length

transactions

The length of the patterns is much less than the average number of items in One transaction

transactions

items

Width: 12    Length: 10000

Width: <500    Length: >>6000

How can we make the right rectangle like the left one?

EXPLOSION!!!

---

## Outline

- Introduction
- **Row Enumeration**
- Confidence-based Prune Strategy
- MAXCONF Algorithm
- Evaluation
- References

---

## Row Enumeration   Transposed table & Tree

| Transaction | Items |
|---|---|
| 1 | A B C D E G |
| 2 | A C D E G |
| 3 | C D E F G H I |
| 4 | B C D E G |
| 5 | A C E G I |
| 6 | A D I |
| 7 | D I J |
| 8 | A B C D G |

| Items | Transactions |
|---|---|
| A | 1,2,5,6 |
| B | 1,4,8 |
| C | 1,2,3,4,5,8 |
| D | 1,2,3,4,6,7,8 |
| E | 1,2,3,4,5 |
| F | 3 |
| G | 1,2,3,4,8 |
| H | 3 |
| I | 3,5,6,7 |
| J | 7 |

**Row Enumeration**

## Row Enumeration **Tree**



**If the current parent node n, is completely contained within a sibling node, a child node is not constructed.
For Example, node 2.**

## Confidence-based Strategy



**RER II, "Mining frequent closed patterns in microarray data." by G. Cong, K.-L. Tan, A. Tung, and F. Pan, 2004**   **Support-based pruning strategy**

**Minimum Support = 30%, Minimum Confidence = 80%**

$$GENE1 \Rightarrow GENE2 \ (support\ 10\%,\ confidence\ 90\%)$$

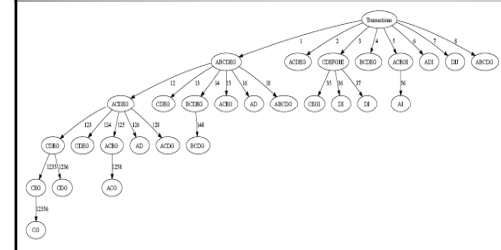**In Biology, we care confidence rules, but not support**

## Outline

- Introduction
- Row Enumeration
- **Confidence-based Prune Strategy**
- MAXCONF Algorithm
- Evaluation
- References

## Confidence-based Strategy **Prune #1**



*Definition 5 (Maximum Support [4]):* Given a node $n$ with $k$ sibling nodes, the *maximum support* of the itemset at $n$, represented as $\sigma_{max}(n)$, or any of $n$'s potential child nodes is:

$$\sigma_{max}(n) = n.\text{initial\_support} + k \qquad (5)$$

$$\sigma_{max}(5) = 1 + 2 = 3$$

## Confidence-based Strategy · Prune #1

*Definition 7 (Minimum Feature):* The item $i_1$ in the itemset $I$ is the *minimum feature* if:

$$\sigma(i_1) \leq \sigma(i_2) \mid \forall i_2 \in I$$

**In the itemset {A,B,C}, Support(A)<=Support(B), Support(A)<=Support(C)
So, A is the minimum feature in {A,B,C}**

*Definition 8 (I-Spanning Rule):* Given an itemset $I$, a rule $r$ is an *I-spanning rule* if:

$$antecedent(r) \cup consequent(r) = I \text{ and}$$
$$|antecedent(r)| = 1$$

**(A)->(B, C) is an I-spanning rule; (B, C)->(A) is not**

*Definition 9 (Maximum Confidence):* Given a node $n$ with minimum feature $i$, the *maximum confidence* of any spanning rule of the itemset at $n$ is:

$$conf_{max}(n) = \frac{\sigma_{max}(n)}{\sigma(i)} \tag{11}$$

## Confidence-based Strategy · Prune #1

$$(I) \rightarrow (ACEG) \quad \textbf{This rule has the highest confidence}$$

$$(AI) \rightarrow (CEG) \quad \textbf{What about this one?} \quad \sqrt{}$$

**Itemset becomes larger, the support of it will not change or even become smaller**

$$\sigma(AI) \leq \sigma(I)$$

$$\uparrow confidence = \frac{\sigma(itemset)}{\sigma(antecedent)} \downarrow$$

## Confidence-based Strategy · Prune #1



**Maximum Support of 5** $\quad \sigma_{max}(5) = 1 + 2 = 3$

**Minimum Feature in this itemset is I** $\quad \sigma(I) = 4$

**Maximum Confidence of 5:** $\quad conf_{max}(5) = \frac{\sigma_{max}(5)}{\sigma(I)} = 3/4$

**If minimum confidence is 4/5, the child of node #5 will be pruned**

$$\sigma(antecedent) \downarrow \qquad \sigma(I) \uparrow$$

## Confidence-based Strategy · Prune #2

*Definition 10 (Maximum features):* Given an itemset $I$, let $R_I$ be the set of all confident $I$-spanning rules. The set of *maximum features*, $M_I$, is the set of all antecedents of the spanning rules.

**Itemset: {CDEG}** $\quad C \rightarrow DEG, E \rightarrow CDG, G \rightarrow CDE$

**The maximum feature of CDEG is CEG**

**Prune Strategy #2:
If maximum feature set M of an itemset at node n is not empty, we can prune all child nodes of n whose itemsets are subsets of M.**

## Slide: Confidence-based Strategy — Prune #2

**Confidence-based Strategy**  **Prune #2**



**Itemset: (1234){CDEG}**  $C \rightarrow DEG, E \rightarrow CDG, G \rightarrow CDE$

**The maximum feature of CDEG is CEG**

**Node (12345)generates:**  $C \rightarrow EG, E \rightarrow CG, G \rightarrow CE$  **Sub-rules**

## Slide: MAXCONF Algorithm

**MAXCONF Algorithm**

Algorithm 1: MAXCONF - High Confidence Rule Mining
Input: Transaction database $D$, minimum confidence $minconf$
Output: High confidence spanning rules satisfying $minconf$
Initialisation:
Let $N$ = set of parent nodes corresponding to each transaction in $D$. Let $n.items$ = itemset represented by node $n$ with support $\sigma(n)$. For each transaction node, $\sigma(n) = 1$ initially. Let $R := \emptyset$ be the set of maximal confident rules.
Procedure: MAXCONF_depthfirst($N$)
  foreach node $n_i \in N$ do
1  if $n_i$ has been discovered then delete $n_i$ and return;
2  Level 1 Confidence Pruning;
    if $n_i$ cannot form a confident spanning rule then delete $n_i$ and continue;
3  Expand subtree;
    Calculate $\sigma(n_i)$ and form children of $n_i$;
4  Maximal Rule Generation;
    $M$:= getMaxFeatures($n_i$);
    foreach $m \in M$ do
      if $m \notin n_i.parentMaxFeatures$ then add rule $m \Rightarrow \{n_i.items - m\}$ into $R$
5  Level 2 Confidence Pruning;
    foreach child $c \in n_i.children$ do
      if $c.items \subset M$ then delete $c$
6  if $n_i.children \neq \emptyset$ then MAXCONF_depthfirst($n_i.children$)
Procedure: getMaxFeatures($n$)
7  maxFeatures := $\emptyset$;
  foreach item $i \in n.items$ do
    if $\sigma(n)/\sigma(i) \geq minconf$ then maxFeatures.insert($i$)
  return maxFeatures

Pruning #2

**Itemset: (1234){CDEG}**  $conf_{max}(5) = \frac{1}{1} = 3/4$

$C \rightarrow DEG, E \rightarrow CDG, G \rightarrow CDE$

Maximum Feature {CEG}
Itemset of child node {CEG}



## Slide: Outline (MAXCONF Algorithm highlighted)

# Outline

## Slide: Outline (Evaluation highlighted)

# Outline

## Evaluation
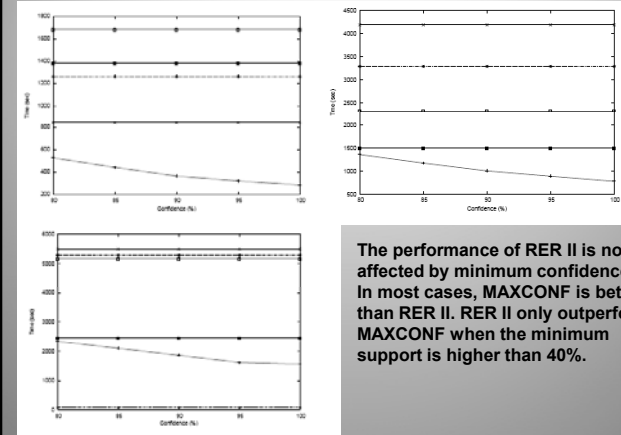
MICROARRAY DATASETS USED IN EXPERIMENTS

| Dataset | # Genes | # Items | #Trans. | Mean trans. size | Min. trans. size | Max. trans. size |
|---|---|---|---|---|---|---|
| Hughes *et al.* (2000) [19] | 6316 | 10044 | 300 | 198 | 2 | 2339 |
| Mnaimneh *et al.* (2004) [20] | 6316 | 8330 | 215 | 228 | 7 | 1111 |
| Spellman *et al.* (1999) [9] | 6178 | 6179 | 82 | 1397 | 205 | 2613 |

**MAXCONF vs RER II**
**Two Aspect:**
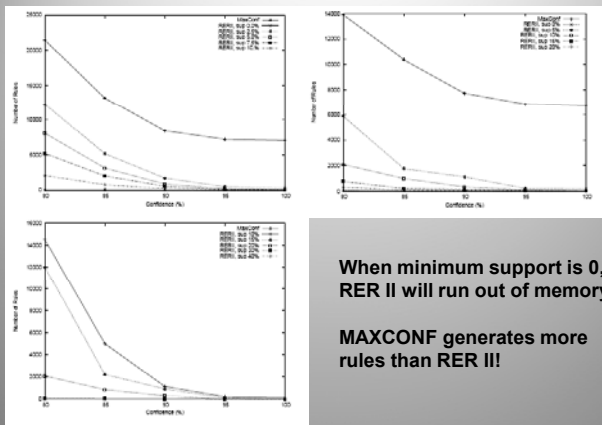1. **Rule Generation**
2. **Scalability**

## Evaluation   Scalability



The performance of RER II is not affected by minimum confidence In most cases, MAXCONF is better than RER II. RER II only outperforms MAXCONF when the minimum support is higher than 40%.

## Evaluation   Rule Generation



When minimum support is 0, RER II will run out of memory

MAXCONF generates more rules than RER II!

## References

- **MAXCONF, "High Confidence Rule Mining for Microarray Analysis",** by Tara McIntosh, Sanjay Chawla, 2006
- **RER II, "Mining frequent closed patterns in microarray data." by G. Cong, K.-L. Tan, A. Tung, and F. Pan, 2004**

Any Question?

Thanks for your Attentions