

# BicFinder: A Biclustering Algorithm for Microarray Data Analysis

Wassim Ayadi<sup>1,2</sup>, Mourad Elloumi<sup>1</sup>, and Jin-Kao Hao<sup>2</sup>

<sup>1</sup> UTIC, Higher School of Sciences and Technologies of Tunis, 1008 Tunis, Tunisia

<sup>2</sup> LERIA, University of Angers, 2 Boulevard Lavoisier, 49045 Angers, France

To appear in *Knowledge and Information Systems* DOI: 10.1007/s10115-011-0383-7

**Abstract.** In the context of microarray data analysis, biclustering allows the simultaneous identification of a maximum group of genes that show highly correlated expression patterns through a maximum group of experimental conditions (samples). This paper introduces a heuristic algorithm called BICFINDER<sup>3</sup> for extracting biclusters from microarray data. BICFINDER relies on a new evaluation function called *Average Correspondence Similarity Index* (ACSI) to assess the coherence of a given bicluster and utilizes a directed acyclic graph to construct its biclusters. The performance of BICFINDER is evaluated on synthetic and three DNA microarray datasets. We test the biological significance using a gene annotation web-tool to show that our proposed algorithm is able to produce biologically relevant biclusters. Experimental results show that BICFINDER is able to identify coherent and overlapping biclusters.

**Key words:** biclustering, heuristics, evaluation function, data mining, analysis of DNA microarray data

## 1 Introduction

Microarray data analysis can be carried out according to at least two different and complementary perspectives. On the one hand, researchers on cancer studies are interested in categorical phenotypes like cancer occurrences, specific tumor subtypes or cancer survivals, which naturally leads to supervised classification of data [44]. *Supervised classification*, also called *class prediction* or *class discrimination*, aims to assign samples to pre-defined categories [2]. On the other hand, *data clustering*, i.e., *unsupervised classification*, aims to identify groups of genes, or groups of experimental conditions, that exhibit similar expression patterns. In such a context, it is particularly interesting to be able to identify simultaneously a group of genes that show similar expression trends across specific groups of experimental conditions, also called samples, [19, 30]. This task is typically achieved by a particular type of clustering technique known as *biclustering* and constitutes the focus of this work. A *bicluster* is a subset of genes

---

<sup>3</sup> The BICFINDER software is available at: <http://www.info.univ-angers.fr/pub/hao/BicFinder.html>

associated with a subset of conditions in which these genes are co-expressed. The *biclustering* problem concerns the identification of the *best* biclusters of a given dataset. Given the practical importance of the problem, many algorithms have been devised to extract good and close-to-optimal biclusters [8, 15, 28, 31, 47].

A microarray data is represented by a data matrix where each cell represents the gene expression level of a gene under a particular experimental condition. Formally, a *bicluster* can be defined as follows: Let  $I=\{1, 2, \dots, n\}$  be a set of indices of  $n$  genes,  $J=\{1, 2, \dots, m\}$  be a set of indices of  $m$  conditions and  $M(I, J)$  be a data matrix associated with  $I$  and  $J$ . A *bicluster* associated with the data matrix  $M(I, J)$  is a couple  $(I, J)$  such that  $I \subseteq I$  and  $J \subseteq J$ . Biclustering is known to be NP-hard [15]. Indeed, it is a highly combinatorial problem with a search space of order of  $O(2^{|I|+|J|})$ .

Existing biclustering algorithms can be grouped into two large classes [?]: Those that adopt a *systematic search* approach and those that adopt a *stochastic search* one, also called *metaheuristic* approach. Algorithms that adopt a systematic search approach include greedy algorithms [8, 14, 15, 28, 43], divide-and-conquer algorithms [22, 38] and enumeration algorithms [5, 23, 26, 35, 41]. Those that adopt a metaheuristic approach include neighbourhood-based algorithms [12], GRASP [18] and evolutionary algorithms [11, 20, 21, 33].

In this paper, we present a greedy algorithm, called BICFINDER, for biclustering of DNA microarray data. The main features of this algorithm include the introduction of a new evaluation function called *Average Correspondence Similarity Index* (ACSI) and the utilization of a directed acyclic graph for biclusters extraction. When applied on both synthetic and real datasets, our algorithm shows better performances compared to other biclustering algorithms.

## 2 Our evaluation function: ACSI

To assess a given bicluster, an evaluation function is required. One of the most popular evaluation functions is called *Mean Squared Residue* (MSR) [15], which has been largely used by several biclustering algorithms [4, 11, 14, 18, 33, 45, 46]. However, MSR is known to be deficient to correctly assess the quality of certain types of biclusters [1, 37, 43]. Recently, another evaluation function called *Average Correlation Value* (ACV) has been proposed in [43]. However, the performance of ACV is known to be sensitive to noises [14]. In the following, we briefly present *Average Spearman's rho* (ASR) proposed in [5] and our new evaluation function *Average Correspondence Similarity Index* (ACSI).

### 2.1 Average Spearman's Rho

In [5], we have proposed an evaluation function called *Average Spearman's Rho* (ASR) based on Spearman's rank correlation [24].

Let  $(I, J)$  be a bicluster in data matrix  $M(I, J)$ , the ASR evaluation function is then defined by:

$$ASR(I, J) = 2 \max \left\{ \frac{\sum_{i \in I'} \sum_{j \in I'; j \geq i+1} \rho_{ij}}{|I'|(|I'|-1)}, \frac{\sum_{k \in J'} \sum_{l \in J'; l \geq k+1} \rho_{kl}}{|J'|(|J'|-1)} \right\} \quad (1)$$

where  $\rho_{ij}$  ( $i \neq j$ ) is the Spearman's rank correlation associated with the row indices  $i$  and  $j$  in the bicluster  $(I, J)$  [24],  $\rho_{kl}$  ( $k \neq l$ ) is the Spearman's rank correlation associated with the column indices  $k$  and  $l$  in the bicluster  $(I, J)$  and  $ASR(I, J) \in [-1..1]$ . A high (resp. low) ASR value, *close* to 1 (resp. *close* to -1), indicates that the genes/conditions of the bicluster are strongly (resp. weakly) correlated. The ASR function can be computed in  $O(n^2m)$ .

In [6], it has been shown that Spearman's rank correlation is robust to the presence of noise in the data and does not require any normalization of the expression data matrix. Since the evaluation function ASR is based entirely on this correlation function, ASR is robust to the presence of noise.

## 2.2 Average Correspondence Similarity Index

Recently, in the area of clustering, Son and Baek [40] showed that that correlation coefficients like Spearman's rank correlation [24] or Pearson correlation [34] are not reliable measures when conditions are few. Furthermore, a high correlation coefficient does not necessarily imply a homogeneous cluster, nor does a low correlation coefficient necessarily imply a heterogeneous cluster [36]. To avoid this difficulty, we propose a new evaluation function called *Average Correspondence Similarity Index* (ACSI) for bicluster evaluation. The proposed ACSI is based on *Concordance Index* (CI) [40] which is designed to measure clusters.

In fact, in microarray data analysis, genes are considered to belong to the same cluster if their trajectory patterns of expression levels are similar (see [29, 36, 39]), i.e., if they have the same profile shape (which may be either monotone increasing, or monotone decreasing, or up-down, or down-up, etc.). Our proposed ACSI is designed to keep track of the profile shape over conditions and preserves the similarity information of trajectory patterns of the expression levels.

In order to calculate ACSI, we first discretize the initial data matrix  $M(I, J)$ ,  $I = \{1, 2, \dots, n\}$  and  $J = \{1, 2, \dots, m\}$ , into a matrix  $M$  defined as follows:

$$M[i, l] = \begin{cases} 1 & \text{if } M[i, l] < M[i, l+1] \\ -1 & \text{if } M[i, l] > M[i, l+1] \\ 0 & \text{if } M[i, l] = M[i, l+1] \end{cases} \quad (2)$$

with  $i \in [1..n]$  and  $l \in [1..m-1]$ .

Let the *correspondence similarity list* between genes  $g_i$  and  $g_j$  ( $i < j$ ), denoted by  $CSL_{i,j}$ , be the list whose each element is represented by  $T(M[i, l] = M[j, l])$  where  $T(Func)$  is true, if and only if,  $Func$  is true, and  $T(Func)$  is false otherwise. Let  $NumCSL_{i,j}$  be the number of times where we have a true value in  $CSL_{i,j}$  and  $MaxCSL_i = \max \{NumCSL_{i,i+1}, NumCSL_{i,i+2}, \dots,$

$NumCSL_{i,n}$ . We define the *correspondence similarity index* as follows:

$$CSI(i, j, k) = \frac{\sum_{l=1}^{m-1} T(M[i, l] = M[j, l] = M[k, l])}{MaxCSL_i} \quad (3)$$

with  $i \in [1..n - 2]$ ,  $j \in [2..n - 1]$ ,  $k \in [3..n]$ ,  $l \in [1..m - 1]$  and  $i < j < k$ .

**Proposition 1:** Let  $i, j$  and  $k$ ,  $i < j < k$ , be indices of rows in a data matrix  $M(I, J)$ , we have (see Proof 1 in appendix):

$$0 \leq CSI(i, j, k) \leq 1.$$

$CSI(i, j, k)$  indicates the proportion of true's, i.e., the change in the same direction, that exists between rows  $i, j$  and  $k$  in the same set of columns. This enables to see how genes  $g_i, g_j$  and  $g_k$  behave over a subset of conditions.

Finally, for the whole bicluster, we define the *Average Correspondence Similarity Index* (ACSI) for the row  $i$  ( $i \in I$  and  $i < j < k$ ):

$$ACSI_i(I, J) = 2 * \frac{\sum_{j \in I; j \neq i+1} \sum_{k \in I; k \neq j+1} CSI(i, j, k)}{|I| (|I| - 1)} \quad (4)$$

where  $|I|$  is the number of rows (genes) in the bicluster without the row index  $i$ .

**Proposition 2:** Let  $(I, J)$  be a bicluster in a data matrix  $M(I, J)$ . We have (see Proof 2 in appendix):

$$0 \leq ACSI_i(I, J) \leq 1.$$

**Proposition 3:** The ACSI function can be computed in a time  $O(n^2m)$ , where  $n$  is the number of the rows and  $m$  is the number of the columns of the data matrix (see Proof 3 in appendix).

A high (resp. low) ACSI value, *close* to 1 (resp. *close* to 0), indicates that the genes of the bicluster are strongly (resp. weakly) correlated.

To illustrate the computation of ACSI, we use the following example: Let  $M$  be a data matrix (Table 1(Left)). When we discretize  $M$ , using Equation 2, we obtain the data matrix  $M'$  (Table 1(Right)).

Suppose that  $Bic = (\{g_1, g_2, g_4, g_5\}; \{c_1, c_2, c_3, c_4\})$  is a bicluster. The ACSI of  $Bic$  is equal to  $ACSI_1 = \frac{CSI(1,2,4) + CSI(1,2,5) + CSI(1,4,5)}{3(3-1)/2} = \frac{3/4 + 3/4 + 3/4}{3} = 0.75$ .

The difference between CI [40] and ACSI is that CI can only measure the coherence between two genes with all conditions while ACSI measures the coherence of the whole bicluster, i.e., subset of rows under a subset of columns.

With the discretization of the data matrix  $M$ , both ACSI and CI can only count the number of conditions that indicate the same sign of change between

	$c'_1$	$c'_2$	$c'_3$	$c'_4$	$c'_5$	$c'_6$		$c_1$	$c_2$	$c_3$	$c_4$	$c_5$
$g_1$	10	20	5	15	40	18	$g_1$	1	-1	1	1	-1
$g_2$	20	40	10	30	24	20	$g_2$	1	-1	1	-1	-1
$g_3$	23	12	8	15	29	50	$g_3$	-1	-1	1	1	1
$g_4$	4	8	2	6	5	5	$g_4$	1	-1	1	-1	0
$g_5$	12	13	8	15	31	50	$g_5$	1	-1	1	1	1
$g_6$	23	12	8	15	29	50	$g_6$	-1	-1	1	1	1

**Table 1. Left:** Data matrix  $M$ . **Right:** discretized data matrix  $M'$ .

rows. Likewise, both ACSI and CI, may lose the information concerning the change of the bicluster size [40]. That is why, we also compute the ASR evaluation function to take into account the information concerning the change of the bicluster size. Indeed, the existing evaluation functions adopted for biclustering can roughly be classified into two families: *numerical* measures and *qualitative* measures. *Numerical* measures, like Pearson’s correlation or Euclidean distance, are easy to compute but they are quite sensitive toward outliers and noise. *Qualitative* measures, like measures that consider only *ups*, *downs* and *no change* of conditions, are very sensitive to precise the values of changes. So, ASR, which is based on Spearman’s rank correlation, can be considered as a good compromise between numerical measures and qualitative ones.

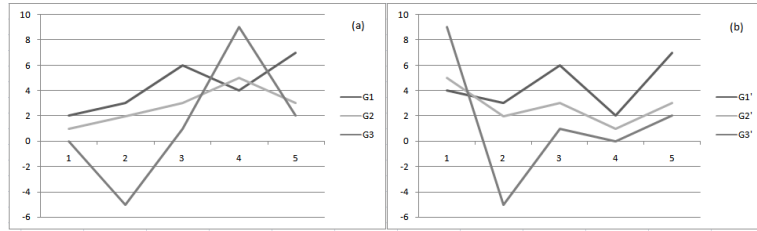
Finally, the example of Figure 1 shows how ACSI can better discriminate biclusters having the same profile shape than other evaluation functions. We consider two biclusters B1 and B2 composed of three genes and five conditions.  $B1 = \{G1(2, 3, 6, 4, 7), G2(1, 2, 3, 5, 3), G3(0, -5, 1, 9, 2)\}$  and  $B2 = \{G1(4, 3, 6, 2, 7), G2(5, 2, 3, 1, 3), G3(9, -5, 1, 0, 2)\}$ . B2 is thus the same as B1, with the first and the fourth conditions of each gene permuted. According to these profile shapes shown in Figure 1, B2 is intrinsically better than B1. Yet the ACV score [43] (based on Pearson correlation coefficient) of B1 and B2 is equally 0.498 just like ASR [5] (based on Spearman correlation coefficient) that leads to the same value of 0.667 for both biclusters. When we apply our ACSI evaluation function<sup>4</sup>, we get an ACSI score of 0.334 (i.e., weakly correlated) for B1 and an ACSI score of 1.0 (i.e., strongly correlated) for B2, showing clearly B2 is better than B1.

The aim of the next section is to assess the quality of the proposed ACSI evaluation function in comparison with ASR and two popular functions that are the *Mean Squared Residue* (MSR) and the *Average Correlation Value* (ACV).

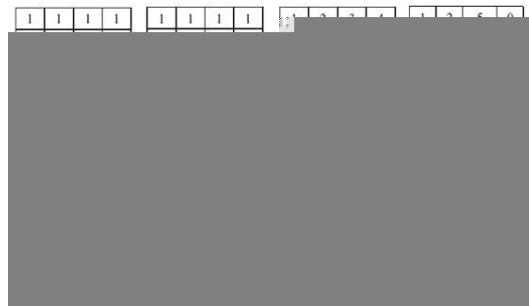
### 2.3 Study of the ACSI evaluation function

As in [43], we compare the ACSI evaluation function with the most used evaluation ones: MSR [15] and ACV [43]. We also compare ACSI with ASR [5].

<sup>4</sup> The values of ACSI and ACV (between 0 and 1) are normalized to drop in [-1 .. 1] to be compared with the scores of ASR which are between -1 and 1.



**Fig. 1.** The expression profiles of two biclusters B1 (a) and B2 (b) with 3 genes and 5 conditions. The genes of B1 show different profiles while those of B2 show similar patterns. B2 is a better bicluster than B1.



**Fig. 2.** Different typical Biclusters.

To compare evaluation functions mentioned above, we use seven data matrices, named  $M_1, M_2, \dots, M_7$ , representing all typical biclusters [30, 43]: A constant bicluster ( $M_1$ ), a constant rows bicluster ( $M_2$ ), a constant column bicluster ( $M_3$ ), a bicluster with coherent values using an additive model ( $M_4$ ), a bicluster with coherent values using a multiplicative model ( $M_5$ ), a bicluster with coherent values using a multiplicative model where the first row of  $M_5$  is multiplied by 10 ( $M_6$ ) and a coherent evolution bicluster ( $M_7$ ).

For each bicluster  $M_1, M_2, \dots, M_7$ , we calculate the coherence score with our evaluation functions ACSI and ASR. The values of MSR and ACV were taken from Table 2 in [43]. The results are summarized in Table 2. Concerning MSR, a low (resp. high) value, *close* to 0 (resp. higher than a fixed threshold), indicates that the genes/conditions of the bicluster are strongly (resp. weakly) correlated. Concerning ACV, a high (resp. low) value, *close* to 1 (resp. *close* to 0), indicates that the genes/conditions of the bicluster are strongly (resp. weakly) correlated. According to Table 2, the ASR, ACSI, ACV and MSR functions are perfect to assess the quality of biclusters  $M_1, M_2, M_3$  and  $M_4$ . However, MSR is deficient on  $M_6$  and  $M_7$ , confirming the claim that MSR is not appropriate on certain types of biclusters [1, 37, 43]. On the other hand, ASR and ACV are perfect to assess the quality of biclusters  $M_5$  and  $M_6$ . ASR is slightly better than ACV

when applied on  $M_7$ . Let us notice that ACSI is perfect when applied on in all biclusters  $M_1, M_2, \dots, M_7$ .

	$M_1$	$M_2$	$M_3$	$M_4$	$M_5$	$M_6$	$M_7$
<i>MSR</i>	0.00	0.00	0.00	0.00	0.62	2.425	131.87
<i>ACV</i>	1.00	1.00	1.00	1.00	1.00	1.00	0.84
<i>ASR</i>	1.00	1.00	1.00	1.00	1.00	1.00	0.99
<i>ACSI</i>	1.00	1.00	1.00	1.00	1.00	1.00	1.00

**Table 2.** ACSI versus ASR, MSR and ACV.

### 3 BicFinder: our proposed biclustering algorithm

The BICFINDER biclustering algorithm is based on the construction of a *Directed Acyclic Graph* (DAG) to represent the different correspondence similarity lists  $CSL_{i,j}$ ,  $1 \leq i < j \leq n$ , between rows. ASR and ACSI are used as evaluation functions. BICFINDER operates in four main steps: discretization of the data matrix, construction of the associated DAG and, extraction and selection of biclusters. The discretization of the data matrix  $M$  is obtained thanks to Equation 2.

#### 3.1 Construction of DAG and, extraction and selection of biclusters

A DAG associated with a data matrix  $M$  is represented as follows: a node  $n_i$  represents a gene  $g_i$  and an arc connecting a node  $n_i$  to a node  $n_j$  if and only if  $i < j$ . To each arc  $(n_i, n_j)$  we assign  $CSL_{i,j}$ .

The next step extracts coherent biclusters. For each node  $n_i$ , first we initialize the associated bicluster  $B_i=(I_i, J_i)$  to  $(\emptyset, \emptyset)$ . Then, we sort the arcs leaving  $n_i$  in a decreasing way according to the number of true's associated with each arc. We consider then the sorted arcs successively. Let  $(n_i, n_k)$  the current arc, if the evaluation function ACSI associated with the bicluster  $(I_i \cup \{g_i, g_k\}, J_i \cup \{c_l, c_{l+1}$  such that  $T(M [i, l] = M [k, l]) = true$  } is greater than or equal to a fixed threshold  $\alpha$  then we set  $B_i=(I_i \cup \{g_i, g_k\}, J_i \cup \{c_l, c_{l+1}$  such that  $T(M [i, l] = M [k, l]) = true$  }. This process is repeated until all the arcs leaving  $n_i$  are processed.

Finally, we consider only the obtained biclusters for which the ASR evaluation function is greater than or equal to another fixed threshold  $\beta$ . The set of such biclusters represents a solution to our problem.

To describe formally the BICFINDER algorithm, let us define some variables:

$M$  (resp.  $M$ ): data matrix (resp. discretized data matrix),

$\mathcal{B}$ : set of biclusters,

$n_i$  (resp.  $n_k$ ): node that represents a gene  $g_i$  (resp.  $g_k$ ),

$B_i = (I_i, J_i)$ : current bicluster,

$I_c, J_c$ : current subset of genes and current subset of conditions,  
 $\alpha, \beta$ : quality thresholds, respectively, according to ACSI and ASR.

---

**Algorithm 1** BICFINDER
 

---

```

1: Input:  $M, \alpha, \beta$  ; Output:  $\mathcal{B}$ 
2: Discretize  $M$  using Equation 2 to obtain  $M'$  // Discretization step
3: Construct the DAG associated with  $M'$  // Construction step
4:  $\mathcal{B} = \emptyset$  // Extraction step
5: for each  $n_i$  in the DAG do
6:    $I'_i = \emptyset; J'_i = \emptyset;$  //  $B_i = (I'_i, J'_i)$ 
7:   Sort the arcs leaving  $n_i$  in a decreasing way according to the number of true's
8:   for each arc  $(n_i, n_k)$  do
9:      $I_c = I'_i \cup \{g_i, g_k\}; J_c = J'_i \cup \{c_l, c_{l+1}$  such that  $T(M'[i, l] = M'[k, l]) = true\};$ 
10:    if  $ACSI_i(I_c, J_c) \geq \alpha$  then  $B_i = (I_c, J_c)$ 
11:  endfor
12:   $\mathcal{B} = \mathcal{B} \cup B_i$ 
13: endfor
14: for each bicluster  $B_i = (I'_i, J'_i)$  in  $\mathcal{B}$  do // Selection step
15:   if  $ASR(I'_i, J'_i) < \beta$  then  $\mathcal{B} = \mathcal{B} \setminus B_i$ 
16: endfor
17: Return  $\mathcal{B}$ 

```

---

**Proposition 4:** Time complexity of BICFINDER is  $O(n^5m)$ , where  $n$  is the number of the rows and  $m$  is the number of the columns of the data matrix (see Proof 4 in appendix).

### 3.2 An illustrative example

By using the previous example in Section 2.2, we represent the data matrix  $M$  (Table 1(Right)) by the DAG that represents all the genes with their  $CSL$  (see Figure 3). Let us set  $\alpha = 0.75$  and  $\beta = 0.9$ .

Figure 3 shows the first node  $g_1$  with its  $CSL$ , i.e., arcs  $\{(a),(b),(c),(d),(e)\}$ . We have two  $MaxCSL_1^5$ : arc (a) with  $NumCSL_{1,2}=4$  and arc (d) with  $NumCSL_{1,5}=4$ . When we compute  $ACSI_1$  of (a) and (d), we find that  $ACSI_1$  is equal to 0.75, i.e.,  $ACSI_1 = \frac{CSL(1,2,5)}{2(2-1)/2} = \frac{3/4}{1} = 0.75$ , then we can merge them in  $B_1$ , giving  $B_1 = (\{g_1, g_2, g_5\}; \{c_1, c_2, c_3, c_4\})$ .

After that, we test ACSI of (a) and (d) with (b) and get  $ACSI_1 = \frac{CSL(1,2,3)+CSL(1,2,5)+CSL(1,3,5)}{3(3-1)/2} = 0.58$ . Hence, this case is discarded because ACSI is lower than 0.75. Then, we test ACSI of (a) and (d) with (c) and we obtain  $ACSI_1 = \frac{CSL(1,2,4)+CSL(1,2,5)+CSL(1,4,5)}{3(3-1)/2} = 0.75$ , then we can merge them. So,  $B_1$  becomes  $B_1 = (\{g_1, g_2, g_4, g_5\}; \{c_1, c_2, c_3, c_4\})$ . Finally, with the node  $g_1$ , we test ACSI of (a), (d) and (c) with (e), then we obtain  $ACSI_1 = \frac{CSL(1,2,4)+CSL(1,2,5)+CSL(1,2,6)+CSL(1,4,5)+CSL(1,4,6)+CSL(1,5,6)}{4(4-1)/2} =$

<sup>5</sup> arcs with stars



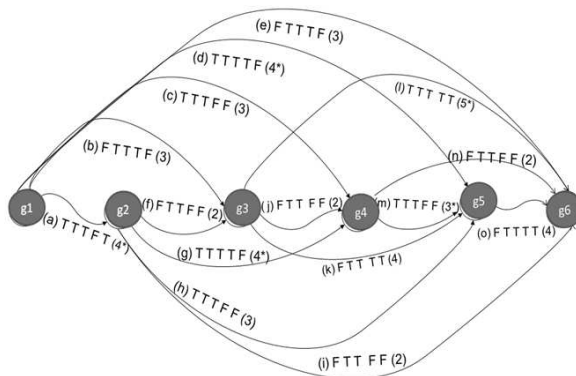


Fig. 3. DAG associated with  $M'$ .

0.625. This case is discarded because ACSI is lower than 0.75. This completes the first node  $g_1$  of DAG, and we repeat the same process for all the other nodes.

After that, we select only the biclusters which are higher than a  $\beta$  threshold, i.e.,  $\beta = 0.9$ . At the end, we obtain two biclusters  $B_1 = (\{g_1, g_2, g_4, g_5\}; \{c_1, c_2, c_3, c_4\})$  with  $ASR=0.9$  and  $B_3 = (\{g_3, g_5, g_6\}; \{c_2, c_3, c_4, c_5, c_6\})$  with  $ASR=1$ . So,  $\mathcal{B} = \{(\{g_1, g_2, g_4, g_5\}; \{c_1, c_2, c_3, c_4\}); (\{g_3, g_5, g_6\}; \{c_2, c_3, c_4, c_5, c_6\})\}$  is the output of our algorithm.

## 4 Results

We assess the BICFINDER algorithm on both synthetic and real DNA microarray datasets. For synthetic data, we compare our results with the results of some prominent biclustering algorithms used by the community, namely, CC [15], OPSM [8], ISA [9] and Bimax [38]. For these references, we have used *Biclustering Analysis Toolbox* (BicAT) which is a recent software platform for clustering-based data analysis that integrates all these biclustering algorithms [7]. For real datasets, in addition to the algorithms mentioned before, we compare our algorithm with the results of Samba [41], RMSBE [28], MOEA [33], MOPSOB [27] and CMOPSOB [25]. For both synthetic and real datasets, the parameters of BICFINDER are fixed after a number of simulations. The parameter settings used for CC, OPSM, ISA and Bimax for the synthetic datasets are the default values as used in [28]. For all the other experiments, we report the results of the compared algorithms from their original papers. The BICFINDER algorithm was implemented in Java and was run on a PC Intel Core 2 Duo T6400 with 2.0GHz CPU and 3.5Gb RAM.

#### 4.1 Computational results of BicFinder on synthetic data

According to [12, 13, 43], we have randomly generated two types of synthetic datasets<sup>6</sup> of size  $(|I|, |J|) = (200, 20)$ . Different types of biclusters are embedded like constant, additive, multiplicative and coherent evolution biclusters. The first (resp. second) dataset contains biclusters without (resp. with) overlapping. To obtain statistically stable results, for each type of datasets, we have generated 10 problem instances by randomly inserting the biclusters at different places in the data matrix.

**Comparison criteria** Following [13], we use the following two ratios to evaluate our biclustering algorithm:

$$\theta_{Shared} = \frac{S_{cb}}{Tot_{size}} * 100 \quad (5)$$

where  $S_{cb}$  is the portion size of biclusters correctly extracted and  $Tot_{size}$  is the total size of correct biclusters.

$$\theta_{NotShared} = \frac{S_{ncb}}{Tot_{size}} * 100 \quad (6)$$

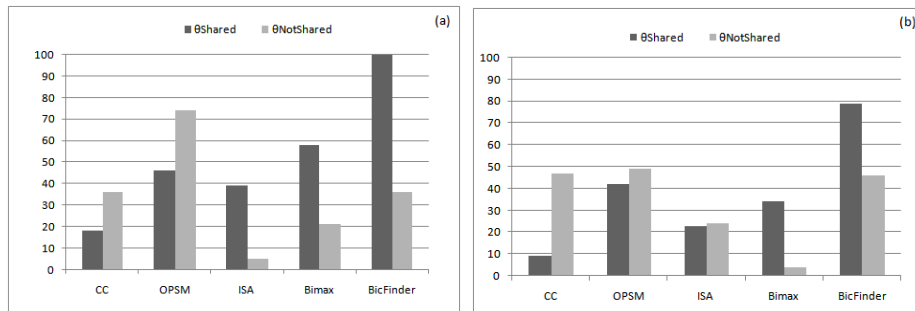
where  $S_{ncb}$  is the portion size of biclusters not correctly extracted and  $Tot_{size}$  is the total size of corrected biclusters.

The ratio  $\theta_{Shared}$  (resp.  $\theta_{NotShared}$ ) expresses the percent of shared (resp. not shared) biclusters volume which corresponds (resp. does not corresponds) with the real biclusters. In fact, when  $\theta_{Shared}$  (resp.  $\theta_{NotShared}$ ) is equal to 100% the algorithm extracts the correct (resp. not correct) biclusters. A perfect solution has  $\theta_{Shared} = 100\%$  and  $\theta_{NotShared} = 0\%$  respecting, thus, the exact number of genes and conditions of implanted biclusters.

**Results and comparisons** For our algorithm, we fix the threshold of ACSI  $\alpha = 0.85$  and threshold of ASR  $\beta = 0.3$ . We run all the algorithms and we select the 4 biclusters obtained by each algorithm which best fit the 4 real biclusters. We compute the  $\theta_{Shared}$  and the  $\theta_{NotShared}$  for each algorithm to show the averaged percentage of volume of the resulting biclusters which is shared and not shared with the real biclusters. The objective of this experiment is to determine which algorithm is able to extract all types of biclusters. Figure 4 shows the statistics of the best biclusters provided by each algorithm for the two datasets.

As we can see in Figure 4(a), BICFINDER obtains 100% of the volume of the correctly extracted biclusters (last row, second column), with an extra volume that represents 36.18 % (last row, third column). In fact, to enlarge a bicluster, BICFINDER tries to add genes that have the maximum similarity with each other over the same subset of conditions. Hence, it can provide an extra volume only

<sup>6</sup> Datasets available at <http://www.info.univ-angers.fr/pub/hao/BicFinder.html>.



**Fig. 4.** (a) Results of BICFINDER and comparison with other algorithms on synthetic data without overlapped biclusters. (b) Results of BICFINDER and comparison with other algorithms on synthetic data with overlapped biclusters.

on conditions but gives exactly the correct number of genes. This is due to the higher number of extracted conditions. In fact, BICFINDER tries to keep the maximum number of conditions which have the same gene expression profile shape.

On the other hand, the best of the studied algorithms, i.e., Bimax, can extract only 58.18% of implanted biclusters with 21.39% of extra volume. The discretization preprocessing used by Bimax may render it impossible to identify the elements in the coherent biclusters and miss some implanted biclusters. ISA has, despite the low value of  $\theta_{NotShared}$  (5.31%), a poor performance since it extracts only 39.38 % of the implanted biclusters ( $\theta_{Shared} = 39.38\%$ ). In fact, ISA uses only up-regulated and down-regulated constant expression values. When coherent biclusters exist, ISA may miss some rows and columns of the implanted biclusters. As to CC, when the signal of the implanted biclusters is weak, the greedy nature of CC may delete some rows and columns of the implanted biclusters at the beginning of the algorithm and miss definitively the deleted rows and columns in the output biclusters. OPSM seeks only up and down regulation expression values with coherent evolution. Its performance decreases when there exist constant biclusters.

Figure 4(b) illustrates the statistics of the best biclusters provided by each algorithm for the second dataset. As we can see in Figure 4(b), the results with BICFINDER present the highest coverage of the correctly extracted biclusters (79.94%) with an extra volume that represent 46.11% of implanted biclusters. On the contrary, the best of the studied algorithms, i.e., OPSM, can extract only 42.87% of implanted biclusters with 49.31% of extra volume. To find overlapped biclusters in a given matrix, some algorithms, e.g. CC, need to mask the discovered biclusters with random values which is not necessary for BICFINDER. ISA and OPSM are sensitive to overlapping biclusters. They use a normalization step in the first preprocessing step. With overlapping biclusters, the expression value range after normalization becomes narrower. Figure 4(b) shows

that BICFINDER is marginally affected by the implanted overlap biclusters compared to other algorithms. This study tends to show that, contrary to existing algorithms, BICFINDER is able to extract all types of biclusters.

#### 4.2 Computational results of BicFinder on real data

In this section, we show computational results of BICFINDER on three well-known real datasets. For the experiments, the two thresholds of BicFinder, i.e.,  $\alpha$  (ASCI threshold) and  $\beta$  (ASR threshold), are fixed after a number of simulations. To proceed, we fix one threshold and tune the other, and vice-versa. For each experiment, five values are tested between 0.5 and 1 with a stepwise of 0.1 for  $\alpha$  and ten values are tested between 0.1 and 1 with a stepwise of 0.1 for  $\beta$ . For each combination, we compute the  $p$ -values of the obtained biclusters. Note that a smaller  $p$ -value, *close* to 0, is indicative of a better match [42]. When a good combination of  $\alpha$  and  $\beta$  is identified, we test ten more values above and below the fixed thresholds with a stepwise of 0.01 (for both  $\alpha$  and  $\beta$ ) to obtain more precise values which lead to biclusters with the lowest  $p$ -values.

**Human B-cell Lymphoma dataset** The Human B-cell Lymphoma dataset [3] contains 4026 genes and 96 conditions <sup>7</sup>. As in [11, 15, 25, 27, 33], we use the criterion of the coverage which is defined as the total number of cells in microarray data matrix covered by the obtained biclusters. For this experiment, the two parameters of BICFINDER  $\alpha$  and  $\beta$  are experimentally set to 0.85 and 0.27. The running time of BICFINDER on this test was 78 minutes.

We select one hundred biggest biclusters out of 727 with high ASR like [11, 15, 25, 27]. These biclusters cover 55.89% of the genes, 100% of the conditions and in total 44.24% cells of the expression data matrix.

On the contrary, in [25] the authors report an average coverage of 38.3% cells of the dataset matrix, while a coverage of 20.96%, 36.9% and 36.81% cells is reported in [33], [27] and [15], respectively.

Table 3 shows the information of five biclusters out of the one hundred biclusters found by BICFINDER on the Human lymphoma dataset with their ASR. The largest bicluster, in this table, is of size 55814.

Genes	Conditions	size	ASR
649	86	55814	0.3064
472	91	42952	0.3246
588	86	50568	0.3107
596	83	49468	0.3009
1018	41	41738	0.2746

**Table 3.** Five biclusters found on Human B-cell Lymphoma dataset.

<sup>7</sup> Available at <http://arep.med.harvard.edu/biclustering/>

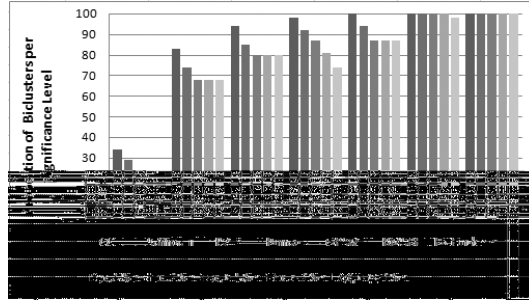


Fig. 5. Proportions of Biclusters significantly enriched by GO annotations.

This implies that BICFINDER can generate maximal biclusters with high coverage of a data matrix.

**Saccharomyces Cerevisiae dataset** The *Saccharomyces Cerevisiae* dataset<sup>8</sup> contains the expression levels of 2993 genes under 173 experimental conditions. In order to evaluate the biological relevance of our proposed biclustering algorithm, we compute the  $p$ -values to indicate the quality of the extracted biclusters. For this experiment, the two parameters of BICFINDER  $\alpha$  (ASCI threshold) and  $\beta$  (ASR threshold) are experimentally set as follows:  $\alpha=0.85$  and  $\beta=0.29$ . The running time of BICFINDER on this test was 484 minutes.

Following the same process as in [17, 28, 38], we extract the largest 100 biclusters out of 1163. The results of BICFINDER are compared against the reported scores of RMSBE, Bimax, OPSM, ISA, Samba and CC from [28, 38]. The idea is to determine whether the set of genes discovered by biclustering algorithms shows significant enrichment with respect to a specific *Gene Ontology* (GO) annotation. We use the web-tool *FuncAssociate* [10] for this purpose. *FuncAssociate* computes the adjusted significance scores for each bicluster, i.e., adjusted  $p$ -values ( $p$ ). Indeed, the adjusted significance scores assess genes in each bicluster by computing , which indicates how well they match with the different GO categories.

Figure 5 presents different significant scores  $p$  for each algorithm over the percentage of total extracted biclusters. On the one hand, BICFINDER and RMSBE seem to outperform other algorithms. The BICFINDER (resp. RMSBE) result shows that 100% (resp. 98%) of discovered biclusters are statistically significant with  $p < 0.001\%$ . On the other hand, apart from CC, other algorithms have reasonably good performance, i.e., the best of the other compared algorithms, OPSM, has 87% of biclusters with  $p < 0.001\%$ . CC under-performs because it is unable to find coherent biclusters and its lack of robustness against noise.

<sup>8</sup> Available at <http://www.tik.ethz.ch/sop/bimax/>

**Yeast Cell-Cycle dataset** The *Yeast Cell-Cycle* dataset is described in [42]. This dataset is processed in [15] and publicly available from [16]. It contains the expression profiles of more than 6000 yeast genes measured at 17 conditions over two complete cell cycles. In our experiments we use 2884 genes selected by [15]. To assess the quality of the extracted biclusters, we use a well-known web-tool on yeast cell-cycle to search for the significant shared Gene Ontology terms of the groups of selected genes. For this experiment, the two parameters of BICFINDER  $\alpha$  (ASCI threshold) and  $\beta$  (ASR threshold) are experimentally set as follows:  $\alpha=0.80$  and  $\beta=0.45$ . The running time of BICFINDER on this test was 5 minutes.

In order to identify the biological annotations for the biclusters, we use *GOTermFinder*<sup>9</sup> which is a tool available in the *Saccharomyces Genome Database* (SGD). *GOTermFinder* is designed to search for the significant shared GO terms of the groups of genes and provides users with the means to identify the characteristics that the genes may have in common. We present the significant shared GO terms (or parent of GO terms) used to describe the two selected set of genes with 92 genes  $\times$  16 conditions and 50 genes  $\times$  17 conditions in each bicluster with ASR equal to 0.4543 and 0.7844 respectively, for biological process, molecular function and cellular component. As [32], we report the most significant GO terms shared by these biclusters. For example, with the first bicluster (Table 4), the genes (YEL034W, YER074W, YER117W, YER131W, YGR214W, YHL001W, YIL069C, YJL111W, YJL136C, YJL177W, YJL189W, YJL190C, YJR123W, YKL056C, YKL156W, YKL180W, YKR057W, YKR094C, YLR029C, YLR048W, YLR075W, YLR150W, YLR167W, YLR185W, YLR248W, YLR249W, YLR325C, YLR344W, YLR367W, YLR380W, YLR406C, YLR441C, YLR448W, YML026C, YML063W, YML073C, YMR143W, YMR225C, YNL030W, YNL067W, YNL096C, YNL162W, YNL301C, YNL302C, YOL039W, YOL040C, YOL127W, YOL139C, YOR167C, YOR234C, YOR293W, YOR312C, YOR369C, YPL037C, YPL081W, YPL090C, YPL143W, YPR043W, YPR102C, YPR163C) are particularly involved in the cellular protein metabolic process and protein metabolic process. The values within parentheses after each GO term in Table 4, such as (59.7%, 4.33e-30) for *translation* in the first bicluster, indicate the cluster frequency and the statistical significance. The cluster frequency (59.7%) shows that out of 92 genes in the first bicluster 55 belong to this process, and the statistical significance is provided by a *p*-value of 4.33e-30 (highly significant).

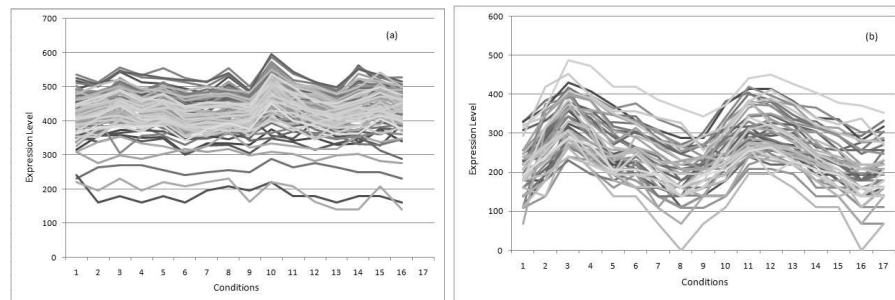
Figure 6 shows the two biclusters of Table 4 found by BICFINDER algorithm on the yeast dataset. From a visual inspection of the biclusters presented, we can notice that the genes do present a similar behavior under the selected conditions. This fact confirms the claim that in microarray data analysis, genes are considered to be in the same cluster if their trajectory patterns of expression levels are similar across a set of conditions [29, 36, 39].

The experiments on these real datasets show that our proposed algorithms can identify biclusters with a high biological relevance.

<sup>9</sup> Available at <http://db.yeastgenome.org/cgi-bin/GO/goTermFinder>

Biclusters	Biological Process	Molecular function	Cellular component
92 genes × 16 conditions	translation (59.7%, 4.33e-30) cellular protein metabolic process (65.2%, 1.48e-18) protein metabolic process (65.2%, 1.43e-17)	structural constituent of ribosome (52.1%, 3.67e-48) structural molecule activity (53.2%, 3.76e-40)	cytosolic ribosome (53.2%, 6.22e-56) cytosolic part (54.3%, 8.15e-52)
50 genes × 17 conditions	cellular response to DNA damage stimulus (37.3% ,3.24e-14) response to DNA damage stimulus (37.3% , 5.27e-13)	double-stranded DNA binding (9.8%, 0.00026) cyclin-dependent protein kinase regulator activity (7.8% , 0.00124)	replication fork (21.6% , 3.36e-12) chromosome (35.3% , 1.94e-10) chromosomal part (31.4%, 5.28e-09)

**Table 4.** Most significant shared GO terms (process, function, component) for two biclusters on yeast data extracted by BICFINDER



**Fig. 6.** Two Biclusters found by BICFINDER on yeast data: (a) 92 genes × 16 conditions with ASR = 0.4543 (b) 50 genes × 17 conditions with ASR = 0.7844.

## 5 Conclusion

We have proposed a novel greedy biclustering algorithm called BICFINDER. BICFINDER, do not require fixing a minimum or a maximum number of genes or conditions, enabling a generation of diversified biclusters.

BICFINDER is based on several new features including a more informative evaluation function, *Average Correspondence Similarity Index* (ACSI), and an effective bicluster extraction mechanism using a *Directed Acyclic Graph* (DAG).

To assess the performance of BICFINDER algorithm, experiments were done on both synthetic and real DNA microarray datasets. This experimental study that we have achieved showed highly competitive results of BICFINDER in comparison with other popular biclustering algorithms. Biological significance of biclustering results has been verified using *Gene Ontology* (GO) annotations. This study also shows that the BICFINDER algorithm competes very favorably with several existing algorithms in terms of the percentage and the number of functionally enriched biclusters for  $p$ -values.

## Acknowledgement

We are grateful to the reviewers for their careful reviews and highly helpful comments. The work is partially supported by the French *Biogenouest* Network and the Region "Pays de La Loire" *via* the "Bioinformatique Ligrienn" Project.

## References

1. J. S. Aguilar-Ruiz. Shifting and scaling patterns from gene expression data. *Bioinformatics*, 21:3840–3845, 2005.
2. A. Akadi, A. Amine, A. El Ouardighi, and D. Aboutajdine. A two-stage gene selection scheme utilizing MRMR filter and GA wrapper. *Knowledge and Information Systems, Published online: 10 March 2010*.
3. A.A. Alizadeh, M.B. Eisen, R.E. Davis, and et al. Distinct types of diffuse large b-cell lymphoma identified by gene expression profiling. *Nature*, 403:503–511, 2000.
4. F. Angiulli, E. Cesario, and C. Pizzuti. Random walk biclustering for microarray data. *Journal of Information Sciences*, pages 1479–1497, 2008.
5. W. Ayadi, M. Elloumi, and J. K. Hao. A biclustering algorithm based on a bicluster enumeration tree: application to dna microarray data. *BioData Mining*, 2(1):9, 2009.
6. R. Balasubramaniyan, H. Ilermeier, E. Weskamp, and J. Kamper. Clustering of gene expression data using a local shape-based similarity measure. *Bioinformatics*, 21:1069–1077, 2005.
7. S. Barkow, S. Bleuler, A. Prelic, P. Zimmermann, and E. Zitzler. Biccat: a biclustering analysis toolbox. *Bioinformatics*, 22(10):1282–1283, 2006.
8. A. Ben-Dor, B. Chor, R. Karp, and Z. Yakhini. Discovering local structure in gene expression data: the order-preserving submatrix problem. In *RECOMB '02: Proceedings of the sixth annual international conference on Computational biology*, pages 49–57, New York, NY, USA, 2002. ACM.



9. S. Bergmann, J. Ihmels, and N. Barkai. Defining transcription modules using large-scale gene expression data. *Bioinformatics*, 20(13):1993–2003, 2004.
10. G.F Berriz, O.D. King, B. Bryant, C. Sander, and F.P. Roth. Characterizing gene sets with funcassociate. *Bioinformatics*, 19(18):2502–2504, 2003.
11. S. Bleuler, A. Prelic, and E. Zitzler. An ea framework for biclustering of gene expression data. In *Proceedings of Congress on Evolutionary Computation*, pages 166–173, 2004.
12. K. Bryan, P. Cunningham, and N. Bolshakova. Application of simulated annealing to the biclustering of gene expression data. In *IEEE Transactions on Information Technology on Biomedicine*, 10(3), pages 519–525, 2006.
13. C. Cano, L. Adarve, J. Lpez, and A. Blanco. Possibilistic approach for biclustering microarray data. In *Computers in Biology and Medicine*, 37, pages 1426–1436, 2007.
14. K.O. Cheng, N.F. Law, W.C. Siu, and A.W. Liew. Identification of coherent patterns in gene expression data using an efficient biclustering algorithm and parallel coordinate visualization. *BMC Bioinformatics*, 9(210):1282–1283, 2008.
15. Y. Cheng and G. M. Church. Biclustering of expression data. In *Proceedings of the Eighth International Conference on Intelligent Systems for Molecular Biology*, pages 93–103. AAAI Press, 2000.
16. Y. Cheng and G.M. Church. Biclustering of expression data. Technical report, (supplementary information), 2006.
17. Y. Christinat, B. Wachmann, and L. Zhang. Gene expression data analysis using a novel approach to biclustering combining discrete and continuous data. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 5(4):583–593, 2008.
18. A. Dharan and A.S. Nair. Biclustering of gene expression data using reactive greedy randomized adaptive search procedure. *BMC Bioinformatics*, 10(Suppl 1):S27, 2009.
19. P. Dimaggio, S. Mcallister, and C. Floudas. Biclustering via optimal re-ordering of data matrices in systems biology: rigorous methods and comparative studies. *BMC Bioinformatics*, 9(1), 2008.
20. F. Divina and J.S. Aguilar-Ruiz. A multi-objective approach to discover biclusters in microarray data. In *GECCO '07: Proceedings of the 9th annual conference on Genetic and evolutionary computation*, pages 385–392, New York, NY, USA, 2007. ACM.
21. C.A. Gallo, J.A. Carballido, and I. Ponzoni. Microarray biclustering: A novel memetic approach based on the pisa platform. In *EvoBIO '09: Proceedings of the 7th European Conference on Evolutionary Computation, Machine Learning and Data Mining in Bioinformatics*, pages 44–55, Berlin, Heidelberg, 2009. Springer-Verlag.
22. J. A. Hartigan. Direct clustering of a data matrix. *Journal of the American Statistical Association*, 67(337):123–129, 1972.
23. D. Jiang, J. Pei, M. Ramanathan, C. Lin, C. Tang, and A. Zhang. Mining gene-sample-time microarray data: a coherent gene cluster discovery approach. *Knowledge and Information Systems*, 13(3):305–335, 2007.
24. E.L. Lehmann and H.J.M. D’Abrera. *Nonparametrics: Statistical Methods Based on Ranks*, pages 292–323. rev. ed. Englewood Cliffs, NJ: Prentice-Hall, 1998.
25. J. Liu, Z. Li, X. Hu, and Y. Chen. Biclustering of microarray data with MOSPO based on crowding distance. *BMC Bioinformatics*, 10(S-4), 2009.
26. J. Liu and W. Wang. Op-cluster: Clustering by tendency in high dimensional space. *IEEE International Conference on Data Mining*, pages 187–194, 2003.

27. J.W. Liu, Z.J. Li, F.F. Liu, and Y.M. Chen. Multi-objective particle swarm optimization biclustering of microarray data. In *IEEE International Conference on Bioinformatics and Biomedicine(BIBM 2008)*, pages 363–366, Washington, DC, USA, 2008. IEEE Computer Society.
28. X. Liu and L. Wang. Computing the maximum similarity bi-clusters of gene expression data. *Bioinformatics*, 23(1):50–56, 2007.
29. Y. Luan and H. Li. Clustering of time-course gene expression data using a mixed-effects model with b-splines. *Bioinformatics*, 19:474–482, 2003.
30. Sara C. Madeira and Arlindo L. Oliveira. Biclustering algorithms for biological data analysis: A survey. *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, 1(1):24–45, 2004.
31. S.C. Madeira and A.L. Oliveira. A polynomial time biclustering algorithm for finding approximate expression patterns in gene expression time series. *Algorithms for Molecular Biology*, 4:8, 2009.
32. U. Maulik, A. Mukhopadhyay, and S. Bandyopadhyay. Combining pareto-optimal clusters using supervised learning for identifying co-expressed genes. *BMC Bioinformatics*, 10:27, 2009.
33. S. Mitra and H. Banka. Multi-objective evolutionary biclustering of gene expression data. *Pattern Recogn.*, 39(12):2464–2477, 2006.
34. J.L. Myers and D.W. Arnold. Research design and statistical analysis. 2003.
35. Y. Okada, K. Okubo, P. Horton, and W. Fujibuchi. Exhaustive search method of gene expression modules and its application to human tissue data. In *IAENG International Journal of Computer Science*, 34, pages 1–16, 2007.
36. S.D. Peddada, E.K. Lobenhofer, L. Li, C.A. Afshari, C.R. Weinberg, and D.M. Umbach. Gene selection and clustering for time-course and dose-response microarray experiments using order-restricted inference. *Bioinformatics*, 19:834–841, 2003.
37. B. Pontes, F. Divina, R. Giráldez, and J.S. Aguilar-Ruiz. Virtual error: A new measure for evolutionary biclustering. In *Evolutionary Computation, Machine Learning and Data Mining in Bioinformatics*, pages 217–226, 2007.
38. A. Prelic, S. Bleuler, P. Zimmermann, P. Buhlmann, W. Gruissem, L. Hennig, L. Thiele, and E. Zitzler. A systematic comparison and evaluation of biclustering methods for gene expression data. *Bioinformatics*, 22(9):1122–1129, 2006.
39. A. Schliep, A. Schonhuth, and C. Steinhoff. Using hidden markov models to analyze gene expression time course data. *Bioinformatics*, 19:i255–i263, 2003.
40. Y.S. Son and J. Baek. A modified correlation coefficient based similarity measure for clustering time-course gene expression data. *Pattern Recognition Letters*, 29(3):232–242, 2008.
41. A. Tanay, R. Sharan, and R. Shamir. Discovering statistically significant biclusters in gene expression data. *Bioinformatics*, 18:S136–S144, 2002.
42. S. Tavazoie, J. D. Hughes, M. J. Campbell, R. J. Cho, and G. M. Church. Systematic determination of genetic network architecture. *Nature Genetics*, 22:281–285, 1999.
43. Li Teng and Laiwan Chan. Discovering biclusters by iteratively sorting with weighted correlation coefficient in gene expression data. *J. Signal Process. Syst.*, 50(3):267–280, 2008.
44. J.M. Wei, S.Q. Wang, and X.J. Yuan. Ensemble rough hypercuboid approach for classifying cancers. *IEEE Transactions on Knowledge and Data Engineering*, 22(3):381–391, 2010.
45. J. Yang, H. Wang, W. Wang, and P. Yu. Enhanced biclustering on expression data. In *BIBE '03: Proceedings of the 3rd IEEE Symposium on Bioinformatics and BioEngineering*, page 321, Washington, DC, USA, 2003. IEEE Computer Society.

46. Z. Zhang, A. Teo, B.C. Ooi, and K.L. Tan. Mining deterministic biclusters in gene expression data. *Bioinformatic and Bioengineering, IEEE International Symposium on*, pages 283–290, 2004.
47. H. Zhao, A. Liew, X. Xie, and H. Yan. A new geometric biclustering algorithm based on the hough transform for analysis of large scale microarray data. *Journal of Theoretical Biology*, 251:264–274, 2008.

## Appendix

**Proof 1:** In the same set of columns,  $\sum_{l=1}^{m-1} T(M'[i, l] = M'[j, l] = M'[k, l])$  represents the number of true's between indices  $i, j$  and  $k$ , and  $MaxCSL_i$  represents the maximum number of true's between  $\{(i, i+1), (i, i+2), \dots, (i, n)\}$ , i.e.,  $MaxCSL_i = \max\{NumCSL_{i,i+1}, NumCSL_{i,i+2}, \dots, NumCSL_{i,n}\}$ . So, the number of true's in  $\sum_{l=1}^{m-1} T(M'[i, l] = M'[j, l] = M'[k, l])$  is usually smaller than or equal to  $MaxCSL_i$ . Since we have:

$$\sum_{l=1}^{m-1} T(M'[i, l] = M'[j, l] = M'[k, l]) \geq 0, \text{ and } MaxCSL_i \geq 0, \text{ then we have:}$$

$$0 \leq \frac{\sum_{l=1}^{m-1} T(M'[i, l] = M'[j, l] = M'[k, l])}{MaxCSL_i} \leq 1, \text{ i.e., } 0 \leq CSI(i, j, k) \leq 1. \quad \bullet$$

**Proof 2:** Indeed, we have  $\frac{|I''|(|I''|-1)}{2}$  values of  $CSI$  to calculate. For  $CSI(i, j, k)$ ,  $i < j < k$ , we have then:

$$0 \leq CSI(i, j, k) \leq 1, \text{ then, } 0 \leq \sum_{j \in I'; j \geq i+1} \sum_{k \in I'; k \geq j+1} CSI(i, j, k) \leq \frac{|I''|(|I''|-1)}{2},$$

$$\sum_{j \in I'; j \geq i+1} \sum_{k \in I'; k \geq j+1} CSI(i, j, k)$$

hence,  $0 \leq \frac{\sum_{j \in I'; j \geq i+1} \sum_{k \in I'; k \geq j+1} CSI(i, j, k)}{\frac{|I''|(|I''|-1)}{2}} \leq 1, \text{ i.e., } 0 \leq ACSI_i(I', J') \leq 1. \quad \bullet$

**Proof 3:** The computing of  $T(M'[i, l] = M'[j, l] = M'[k, l])$  can be achieved in a time  $O(1)$ . We repeat this function  $m-1$  times. We have  $\frac{|I''|(|I''|-1)}{2}$  values of *Correspondence Similarity Index* to calculate, i.e.,  $O(n^2)$ . Hence, the ACSI function can be computed in a time  $O(n^2m)$ .  $\bullet$

**Proof 4:** Time complexity of the discretization step is  $O(nm)$ . Indeed, this step is achieved *via* a scanning of the whole data matrix  $M$  that is of size  $nm$ .

Time complexity of the construction step is  $O(n^2m)$ . In fact, the construction of an arc and the counting of the number of the associated true's in the DAG is made in  $O(m)$ . Since the total number of arcs is  $\frac{n(n-1)}{2}$  then time complexity of the construction step is  $O(n^2m)$ .

Time complexity of the extraction step is  $O(n^5m)$ . Indeed, in the worst case, a node  $n_i$  has  $(n-1)$  outgoing arcs. So, the determination of the outgoing arcs of a node  $n_i$  that have the maximum number of true's is achieved in  $O(n)$ . Each one of these arcs will be compared to a bicluster by computing the corresponding ACSI evaluation function. According to Proposition 3, this is done in a time  $O(n^2m)$ . In the worst case, the comparison between a bicluster and all the outgoing arcs, that have the maximum number of true's, is achieved in a time  $O(n^3m)$ . In the worst case, we have  $O(n)$  biclusters, so the comparison between all the existing biclusters and all the outgoing

arcs, that have the maximum number of true's, is achieved in  $O(n^4m)$ . We have  $n$  nodes, so time complexity of the extraction step is  $O(n^5m)$ .

Time complexity of the selection step is  $O(n^4m)$ . In fact, for a bicluster, the ASR evaluation function is computed in a time  $O(n^2m)$ . In the worst case, we have  $\frac{n(n-1)}{2}$  extracted biclusters. Thus, this step is achieved in a time  $O(n^4m)$ .

Hence, time complexity of BICFINDER is  $O(n^5m)$ . •