

International Journal on Artificial Intelligence Tools
© World Scientific Publishing Company

An Improved Biclustering Method for Analyzing Gene Expression Profiles

Jiong Yang

*EECS Department, Case Western Reserve University
Cleveland, Ohio 44166, USA
jiong@eecs.cwru.edu*

Haixun Wang

*IBM T. J. Watson Research Centers 19 Skyline Drive
Hawthorne New York 10532, USA
haixun@us.ibm.com*

Wei Wang

*CS Department, UNC-Chapel Hill
Chapel Hill, North Carolina 27599, USA
weiwang@cs.unc.edu*

Philip S. Yu

*IBM T. J. Watson Research Centers 19 Skyline Drive
Hawthorne New York 10532, USA
psyu@us.ibm.com*

Microarrays are one of the latest breakthroughs in experimental molecular biology, which provide a powerful tool by which the expression patterns of thousands of genes can be monitored simultaneously and are already producing huge amount of valuable data. The concept of *bicluster* was introduced by Cheng and Church (2000) to capture the coherence of a subset of genes and a subset of conditions. A set of heuristic algorithms were also designed to either find one bicluster or a set of biclusters, which consist of iterations of masking null values and discovered biclusters, coarse and fine node deletion, node addition, and the inclusion of inverted data. These heuristics inevitably suffer from some serious drawback. The masking of null values and discovered biclusters with random numbers may result in the phenomenon of *random interference* which in turn impacts the discovery of high quality biclusters. To address this issue and to further accelerate the biclustering process, we generalize the model of bicluster to incorporate null values and propose a probabilistic algorithm (FLOC) that can discover a set of k possibly overlapping biclusters simultaneously. Furthermore, this algorithm can easily be extended to support additional features that suit different requirements at virtually little cost. Experimental study on the yeast gene expression data shows that the FLOC algorithm can offer substantial improvements over the previously proposed algorithm.

1. Introduction

Microarrays are one of the latest breakthroughs in experimental molecular biology, which provide a powerful tool by which the expression patterns of thousands of genes can be monitored simultaneously and are already producing huge amount of valuable data. Analysis of such data is becoming one of the major bottlenecks in the utilization of the technology. The gene expression data are organized as matrices — tables where rows represent genes, columns represent various samples such as tissues or experimental conditions, and numbers in each cell characterize the expression level of the particular gene in the particular sample. Investigations show that more often than not, several genes contribute to the same pathway, which motivates researchers to identify a subset of genes whose expression levels rise and fall coherently under a subset of conditions, that is, they exhibit fluctuation of a similar shape when conditions change. Discovery of such clusters of genes is essential in revealing the significant connections in gene regulatory networks.

The concept of *bicluster* was introduced by Cheng and Church (2000) to capture the coherence of a subset of genes and a subset of conditions. Unlike previous methods that treat similarity as a function of pairs of genes or pairs of conditions, the bicluster model measures coherence within the subset of genes and conditions. This model may be particularly useful to disclose the involvement of a gene or a condition in multiple pathways, some of which can only be discovered under the dominance of more consistent ones. The coherence score is defined as a symmetric function of genes and conditions involved and thereby the biclustering is a process of simultaneous grouping of genes and conditions. The so called *mean squared residue* was employed and was applied to expression data transformed by a logarithm and augmented by the additive inverse. While the mean squared residue represents the variance of the selected genes and conditions with respect to the coherence, the goal of biclustering is to find biclusters with low mean squared residue. In gene expression data analysis, this goal is often accompanied with an additional requirement of reasonably large row variance. The rationale is that a low mean squared residue only indicates that the gene expression levels fluctuate approximately in unison, which also includes the constant biclusters where there is no or little fluctuation at all. These trivial biclusters may not be as interesting as the biclusters where the set of genes show strikingly similar up-regulation and down-regulation under the set of conditions. It has been proven that the problem of finding biclusters satisfying these criteria is NP-hard in general. Therefore, a set of heuristic algorithms were designed by Cheng and Church (2000) to either find one bicluster or a set of biclusters, which consist of iterations of masking null values and discovered biclusters, coarse and fine node deletion, node addition, and the inclusion of inverted data. The computational complexities are in the order of $O(MN \times (M + N) \times k)$ for discovering k biclusters where M and N are the number of conditions and the number of genes, respectively. The proposed heuristics, which have been demonstrated to be able to produce good quality biclusters, inevitably suffer from some serious drawback. The masking of

null values and discovered biclusters was performed by replacing the relevant cells with random numbers. The rationale of replacing missing values with random numbers was that these random values only have a mathematical chance to form any recognizable pattern and therefore would not result in distorted biclustering. The original intention of masking discovered bicluster was to ensure that each successive run of the (deterministic) algorithm outputs a different bicluster in the case where multiple biclusters are preferred. In both cases, even though the random data is unlikely to form any fictitious pattern, there exists a substantial risk that these random numbers will interfere with the future discovery of biclusters, especially those ones that have overlap with the discovered ones. We call this phenomenon the *random interference*. Our experimental study has confirmed that this random interference will impact the biclustering result.

To address this issue and to further accelerate the biclustering process, we generalize the model of bicluster to incorporate null values and propose a probabilistic algorithm (FLOC^a) that can discover a set of k possibly overlapping biclusters simultaneously. Furthermore, this algorithm can easily be extended to support additional features that suit different application needs at virtually little cost. Typical features include the maximum amount of overlap allowed between biclusters, the maximum/minimum size of each bicluster, the minimum overall coverage of the biclusters, and so on. The general process of FLOC consists of iterations of series of gene and condition moves (i.e., selections or deselections) aiming at achieving the best potential residue reduction. During the course of biclustering, certain move may be “blocked” temporarily if performing such move would lead to an unfavorable situation such as producing a trivial bicluster or violating one or more feature constraints. We implemented FLOC to find 100 biclusters on the same yeast data containing 2884 genes and 17 conditions with the same parameter setting as in Cheng and Church (2000) and found that the biclusters returned by FLOC, on average, have a comparable mean squared residue but a larger size than that reported by Cheng and Church (2000). In addition, FLOC is able to locate these biclusters much faster than the algorithms proposed in Cheng and Church (2000).

The remainder of this paper is organized as follows. We present the generalized bicluster model in Section 2. Section 3 and 4 present our basic FLOC algorithm and some additional improvements, respectively. Experimental results are shown in Section 5 and we draw the conclusion in Section 6.

2. The General Model of Bicluster

In this section, we formally present the *generalized* bicluster model that can handle null values in a seamless manner. (In the remaining of this paper, we use the term of biclusters to refer to the generalized biclusters.) A bicluster is defined on a gene-expression matrix. Let $\mathfrak{S} = \{A_1, A_2, \dots, A_M\}$ be the set of conditions and $\mathfrak{R} =$

^aFLOC stands for FLexible Overlapped biClustering

$\{O_1, O_2, \dots, O_N\}$ be the set of genes. The data can be viewed as an $M \times N$ matrix D of real numbers. Each entry d_{ij} in this matrix corresponds to the logarithm of the relative abundance of the mRNA of a gene O_i under a specific condition A_j , and may have a null value. Figure 1 illustrates the general format of the data matrix.

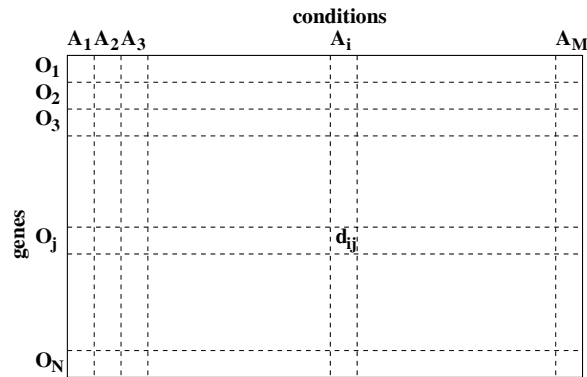


Fig. 1. The Data Matrix

A *bicluster* essentially corresponds to a submatrix that exhibits some coherent tendency. Formally, each bicluster can be uniquely identified by the set of relevant genes and conditions. Even though allowing missing values brings great flexibility to the bicluster model, the amount of missing entries in a bicluster should be limited to some extent to avoid trivial cases. The rule of the thumb is that, despite the missing values, there should still be sufficient evidence to demonstrate the coherency. In Figure 2 (a), many values are missing, which prevents any potential coherence from being observed, even though there is no sign that contradicts the existence of coherence either. To exclude this kind of situation from being considered as a meaningful bicluster, we introduce a parameter α (which is a positive number less than or equal to 1) to limit the amount of missing values for each gene and each condition in a bicluster.

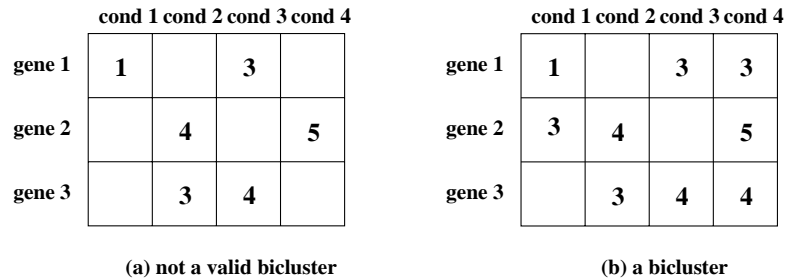


Fig. 2. Missing Values in Biclusters

Definition 2.1. For a given matrix $\mathfrak{S} \times \mathfrak{R}$ and an occupancy threshold α , a **bicluster** (of α occupancy) can be represented by a pair (I, J) where $I \subseteq \{1, \dots, M\}$ is a subset of genes and $J \subseteq \{1, \dots, N\}$ is a subset of conditions. For each gene $i \in I$, $\frac{|J'_i|}{|J|} > \alpha$ where $|J'_i|$ and $|J|$ are the number of specified conditions for gene i in the bicluster and the number of conditions in the bicluster, respectively. Similarly, for each condition $j \in J$, $\frac{|I'_j|}{|I|} > \alpha$ where $|I'_j|$ and $|I|$ are the number of specified genes under condition j in the bicluster and the number of genes in the bicluster, respectively.

Let $\alpha = 0.6$, the submatrix in Figure 2 (a) is not a valid bicluster while the submatrix in Figure 2 (b) is a bicluster. The number of specified (non-missing) entries in the corresponding submatrix is referred to as the *volume* of the bicluster.

Definition 2.2. The **volume** of a bicluster (I, J) (v_{IJ}) is defined as the number of specified entries d_{ij} such that $i \in I$ and $j \in J$.

In the case that all entries are specified, $v_{IJ} = |I| \times |J|$ where $|I|$ and $|J|$ are the number of conditions and the number of genes participating in the bicluster, respectively. Figure 3(a) shows a gene expression matrix with ten genes (one for each rows) under five conditions (one for each column). The bicluster defined by picking $I = \{2, 3, 8\}$ and $J = \{1, 3, 5\}$ is shown in Figure 3(b). The volume of this bicluster is 9.

		conditions				
		1	2	3	4	5
genes	1	4392	284	4108		228
	2	401	281	120	275	298
	3	318	280	37	277	215
	4	401	292	109	580	238
	5	2857	285		271	226
	6	228	290	48	285	224
	7	538	272	266	277	236
	8	322	288	41	278	219
	9	312		40	273	232
	10	329	296	33	274	228

(a) a data matrix

		conditions				
		1	2	3	4	5
genes	1					
	2	401		120		298
	3	318		37		215
	4					
	5					
	6					
	7					
	8	322		41		219
	9					
	10					

(b) a bicluster

Fig. 3. An Example of Bicluster

In order to properly accommodate various expression levels associated with each gene and each condition within a bicluster, we introduce a concept — *base*.

Definition 2.3. For a given bicluster (I, J) , the **base** of a gene O_i is defined as the average value of O_i for all specified conditions in J , $d_{iJ} = \frac{\sum_{j \in J'_i} d_{ij}}{|J'_i|}$ where $J'_i \subseteq J$ is the set of specified conditions in J for gene O_i . Similarly, the **base** of a condition A_j is the average specified value of A_j taken by all genes in I , i.e., $d_{Ij} = \frac{\sum_{i \in I'_j} d_{ij}}{|I'_j|}$ where $I'_j \subseteq I$ is the set of genes whose value is specified in condition A_j . The **base** of the bicluster is the average value of all specified entries of the submatrix defined by (I, J) , i.e., $d_{IJ} = \frac{\sum_{i \in I, j \in J} d_{ij}}{v_{IJ}}$ where v_{IJ} is the volume of the bicluster.

For example, we have $d_{2,J} = 273$, $d_{3,J} = 190$, $d_{8,J} = 194$, and $d_{I,1} = 347$, $d_{I,3} = 66$, $d_{I,5} = 244$, and $d_{IJ} = 219$ in Figure 3(b). While d_{iJ} and d_{Ij} take care of the potential tendency that may associate with each individual gene or condition, the value of d_{IJ} set the base point of the entire bicluster. In a perfect bicluster where each gene and condition exhibits an absolutely consistent tendency^b, the value of each entry d_{ij} can be uniquely determined by its gene base d_{iJ} , its condition base d_{Ij} , and the bicluster base d_{IJ} . The difference $d_{iJ} - d_{IJ}$ is essentially the relative tendency held by gene O_i in contrast to other genes in the bicluster. This tendency should hold exactly on the entry d_{ij} as well in a perfect bicluster. That is, $d_{ij} - d_{Ij} = d_{iJ} - d_{IJ}$. Consequently, we have $d_{ij} = d_{iJ} + d_{Ij} - d_{IJ}$. Figure 3(b) is a perfect bicluster even though the values are quite far apart (which may produce poor quality bicluster(s) of the traditional meaning). For example, the entry $d_{2,1} = d_{2,J} - d_{I,1} + d_{IJ} = 273 - 347 + 219 = 401$ and this property holds for every entry in Figure 3(b).

In practice, the bicluster may not always be perfect. The concept of *residue* is thus introduced to quantify the difference between the actual value of an entry and the expected value of an entry predicted from the corresponding gene base, condition base, and the bicluster base.

Definition 2.4. The **residue** of an entry d_{ij} in a bicluster is $r_{ij} = d_{ij} - d_{iJ} - d_{Ij} + d_{IJ}$ if d_{ij} is specified. Otherwise, $r_{ij} = 0$.

It is obvious that every entry in Figure 3(b) has a zero residue. The residue indeed serves as an indicator of the degree of coherence of an entry with the remaining entries in the bicluster given the tendency of the relevant gene and the relevant condition. The lower the residue, the stronger the coherence. To assess the overall quality of a bicluster, the **residue** of the bicluster can be defined as the mean residue of all specified entries. The mean can be in the form of either arithmetic, geometric, or square mean. In this paper, we use the square mean in the assessment of the bicluster residue as in Cheng and Church (2000).

Definition 2.5. The **residue** of a bicluster (I, J) is $r_{IJ} = \frac{\sum_{i \in I, j \in J} r_{ij}^2}{v_{IJ}}$ where r_{ij} is the residue of the entry d_{ij} and v_{IJ} is the volume of the bicluster.

^bThe entries of each gene (or condition) can be exactly generated by shifting the entries of other genes (or conditions) by a common offset.

In the above example, the residue of the bicluster in Figure 3(b) is 0. The lower the residue, the stronger the coherence exhibited by the bicluster, and the better the quality of the bicluster. We also refer to a bicluster (I, J) as a **r -residue bicluster** if its residue $r_{IJ} \leq r$ where r is a constant number. In addition, we may prefer the *row variance* to be relatively large to reject trivial biclusters.

Definition 2.6. The **row variance** of a bicluster (I, J) is defined as $var_{I,J} = \frac{\sum_{i \in I, j \in J} (d_{ij} - d_{iJ})^2}{v_{IJ}}$.

This accompanying score would warrant the bicluster to capture genes exhibiting fluctuating yet coherent trends under some set of conditions. The basic bicluster model can also be easily extended to support some additional features that may be very useful in many applications.

- The amount of overlap allowed between a pair of biclusters $Cons_o$: Some application may require mutually exclusive biclusters while others may prefer some degree of overlap. The user can control the amount of overlap by specifying some threshold.
- The coverage of the biclusters $Cons_c$: the number of genes/conditions should be covered by any of the biclusters. In some case, the user may want every gene to be covered by some bicluster.
- The balance between number of genes and conditions of the bicluster $Cons_b$: the desirable ratio (or its range) between the number of genes and the number of conditions of each bicluster can be also be specified if the user prefers to find “balanced” biclusters.
- The volume of the final biclusters $Cons_v$: The user can also control the volume of the final biclusters. This can be useful in the application where certain statistical significance needs to be warranted.

We shall see later in this paper that our proposed algorithm can be applied with minor modification to suit the above purposes and to produce promising results. A list of notations and conventions that will be used throughout this paper is provided in Table 1.

3. The Basic FLOC Algorithm

3.1. Algorithm Description

In general, the bicluster problem is NP-hard as proven by Cheng and Church (2000). Thus, finding an exact solution could be time consuming. In this section, we present a new probabilistic *move-based* algorithm called FLOC, which can efficiently and accurately approximate the k biclusters with low mean squared residues. The data is represented in the form of a matrix as shown in Figure 3(a) where the rows correspond to the genes and the columns correspond to the conditions. The FLOC biclustering algorithm starts from a set of seeds (initial biclusters) and carries out an

Table 1. Notation and Convention

D	the data matrix
\mathfrak{S}	the set of conditions
\mathfrak{R}	the set of genes
M	the number of conditions (i.e., $M = \mathfrak{S} $)
N	the number of genes (i.e., $N = \mathfrak{R} $)
A_j	the j th condition in \mathfrak{S}
O_i	the i th gene in \mathfrak{R}
d_{ij}	the entry in the matrix which corresponds to the value of gene O_i on condition A_j
I (possibly with subscript)	the subset of genes that participate in a bicluster
J (possibly with subscript)	the subset of conditions that participate in a bicluster
(I, J)	the bicluster defined by genes in I and conditions in J
v_{IJ}	the volume of the bicluster (I, J)
d_{iJ}	the average value of entries corresponding to gene O_i of all conditions in J
d_{Ij}	the average value of entries corresponding to condition A_j on all genes in I
d_{IJ}	the average value of entries corresponding to genes in I and conditions in J
r_{ij}	the residue of the entry d_{ij} with respect to some bicluster
r_{IJ}	the residue of the bicluster (I, J)
r	residue threshold
i (possibly with subscript)	index term for gene
j (possibly with subscript)	index term for condition
α	the occupancy threshold
k	number of biclusters
$Cons_o$	the threshold of overlap between a pair of biclusters
$Cons_c$	the percentage threshold of genes (or conditions) covered by at least one bicluster
$Cons_b$	the balance ratio threshold of each bicluster
$Cons_v$	the volume threshold of each bicluster
ρ	the probability that a row or column is assigned to a bicluster at the initial phase
x	a row or a column in the data matrix
c	a bicluster
$Action(x, c)$	the action of changing membership of x with respect to c
a (possibly with subscript)	an action
g (possibly with subscript)	the gain of some action

iterative process to improve the overall quality of the biclustering. At each iteration, each row and column is moved among biclusters to produce a better biclustering in terms of lower mean squared residues. The best biclustering obtained during each iteration will serve as the initial biclustering for the next iteration. The algorithm terminates when the current iteration fails to improve the overall biclustering quality.

The FLOC algorithm has two phases (Figure 4). In the first phase, k initial biclusters are constructed. As we presented in the previous section, a bicluster con-

tains a set of genes (rows) and a set of conditions (columns). A parameter ρ is introduced to control the size of a bicluster. For each initial bicluster, a random switch is employed to determine whether a row or column should be included. Each row and column is included in the bicluster with probability ρ . Consequently, each initial bicluster is expected to contain $M \times \rho$ rows and $N \times \rho$ columns. If the percentage of specified values in an initial cluster falls below the α threshold, then we keep generating new clusters until the percentage of specified values of all columns and rows satisfy the α threshold. In this paper, the α is chosen as follows. Let D be the original matrix, which has N rows and M columns. α is set to $\frac{|D|}{N \times M}$ where $|D|$ is the volume of D . In this case, we can guarantee that any biclusters of D has at most the same percentage of unspecified values as D itself. (We will discuss how to choose ρ in the next section.)

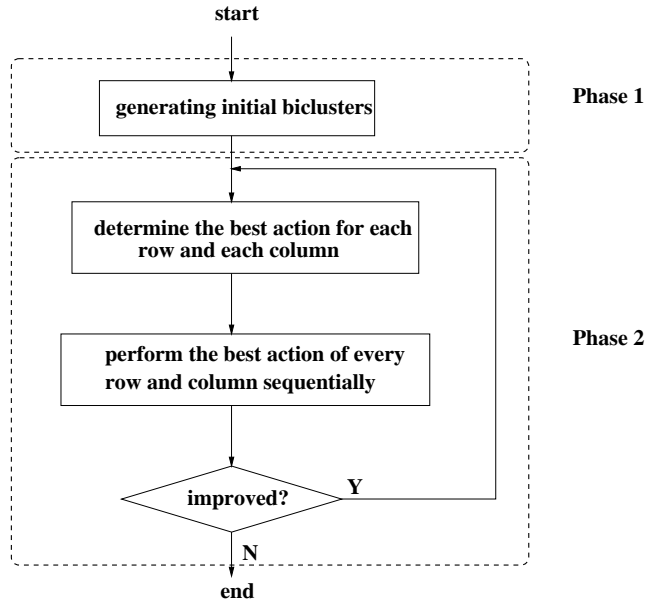


Fig. 4. The Flowchart of the FLOC Algorithm

The second phase is an iterative process to improve the quality of the biclusters continuously. During each iteration in the second phase, each row and each column are examined to determine its best action towards reducing the overall mean squared residue. These actions are then performed successively to improve the biclustering. An **action** is defined with respect to a row (or column) and a bicluster. There are k actions associated with each row (or column), one for each bicluster. For a given row (or column) x and a bicluster c , the **action** $Action(x, c)$ is defined as the change of membership of x with respect to c . Note that this action is uniquely defined at any stage. If x is already included in c , then $Action(x, c)$ represents the removal x from

the bicluster c . Otherwise, $Action(x, c)$ denotes the addition of x to the bicluster c . Figure 5 shows a data matrix with 3 rows and 4 columns. Assume that we want to find two biclusters and their current status is indicated by the dash lines. Bicluster 1 contains row 1, 2 and column 1, 2; whereas bicluster 2 contains row 2, 3 and column 1, 2, 3. Each row (or column) in Figure 5 is then associated with two actions, one for each bicluster. For example, the actions associated with column 3 are (1) inserting into bicluster 1, and (2) deleting from bicluster 2. The better action among these two need to be identified and performed. In general, if the data matrix contains N rows and M columns, then $N + M$ actions will be performed during each iteration, one for each row (or column). We will discuss shortly that sometimes an action may be blocked temporarily during an iteration due to the violation of some constraint (by the action).

	column 1	2	3	4
row 1	3	4	2	2
2	1	3	2	3
3	4		0	4

Fig. 5. An Example of the Actions

Since there are k biclusters, the number of potential actions associated with the row (or column) x is k . Among these k actions, the action that brings most improvement needs to be identified. To assess the amount of improvement that can be brought by an action, we introduce a new concept called **gain**. Since our objective is to find biclusters with low residue ($< r$), the **gain** of an action $Action(x, c)$ is defined as a function of the *relative reduction* of c 's residue and the *relative enlargement* of c 's volume as a consequence of performing $Action(x, c)$.

Definition 3.1. Given a residue threshold r , the **gain** of an action $Action(x, c)$ is defined as $Gain(x, c) = \frac{r_c - r_{c'}}{\frac{r^2}{r_c}} + \frac{v_{c'} - v_c}{v_c}$ where $r_c, r_{c'}$ are the residues of bicluster c and the bicluster, c' , obtained by performing $Action(x, c)$ on c , respectively. Similarly, v_c and $v_{c'}$ are the volumes of c and c' , respectively.

When c has a much smaller residue than the threshold r ($r_c \ll r$), the gain measurement would favor those actions that enlarge c , especially if c has a small volume. This encourages the FLOC algorithm to find large biclusters with tolerable residues. On the other hand, when c 's residue is larger than r ($r_c \gg r$), the

measurement of gain inclines to keep the residue of c under control, especially if c already has a sufficiently large volume. Obviously, a positive gain indicates that performing $Action(x, c)$ has a potential to produce a better bicluster while a negative gain suggests that such an action would be likely to degrade the bicluster quality. Therefore, the intermediate goal during the course of biclustering becomes, for each row (or column) x , to find and perform the action with the highest gain. In the above example, the residue of bicluster 1 and bicluster 2 are $\frac{1}{16}$ and $\frac{7}{6}$, respectively. Let $r = 1$ and consider the two actions associated with column 3: (1) inserting into bicluster 1 and (2) deleting from bicluster 2. The resulting bicluster after inserting column 3 into bicluster 1 would contain row 1,2 and column 1,2,3, and its residue is $\frac{1}{6}$. Therefore, the gain of inserting column 3 into bicluster 1 is $\frac{\frac{1}{6} - \frac{1}{16}}{\frac{1}{16}} + \frac{2}{4} = -\frac{5}{768} + \frac{1}{2} = \frac{379}{768}$. Via similar computation, the gain of deleting column 3 from bicluster 2 is $\frac{\frac{7}{6} - 1}{\frac{12}{776}} - \frac{1}{3} = \frac{7}{36} - \frac{1}{3} = -\frac{5}{36}$. Consequently, the first action is chosen as the best action for column 3. Note that the best action for a row or column might be negative. Such negative action(s) will still be performed. The rationale is that the (temporary) degradation of the bicluster quality may lead to an ultimate (bigger) improvement. We will explain shortly that such action will not take into effect if the bicluster quality fails to improve by the end of the iteration. Nevertheless, the highest gain of any action associated with a given row (or column) is positive in many cases and will directly contribute to the improvement of the bicluster quality. For example, the highest gain of any action associated with column 3 is $\frac{379}{768}$.

To compute the gain of a particular action, the residue of the resulting bicluster (if the action was taken) needs to be computed. The straightforward way to compute the residue after each action is to recompute from scratch. This involves the computation of each gene base, each condition base, and bicluster base, and finally the bicluster residue. A more efficient method is to recompute only those gene and condition bases affected by the action. This can be done efficiently (in an incremental manner) if the gene bases and condition bases of the bicluster are maintained along with the bicluster base throughout the course. This technique effectively reduces the time complexity from $O(N \times M)$ to $O(N + M)$ where N and M are the number of rows and the number of columns of the data matrix, respectively.

After the best action is identified for every row (or column), these $N + M$ actions are then performed sequentially. The best biclustering obtained during the last iteration, denoted by *best_biclustering*, is used as the initial biclustering of the current iteration. Let *Biclustering_i* be the set of biclusters after applying the first i actions. After applying all actions, we would obtain $M + N$ sets of biclusterings. Among them, if any biclustering with all r -biclusters^c has a larger aggregated volume than that of *best_biclustering*, then there is an improvement in the current iteration.

^cIt is possible that the biclustering at some stage contains some bicluster with residue larger than r . The biclustering at this stage will not be considered even the aggregated volume is larger than that of *best_biclustering*.

The biclustering with the minimum average residue is stored in *best_biclustering* and the process continues to the next iteration. Otherwise, there is no improvement in the current iteration and the process terminates. The biclustering stored in *best_biclustering* is then returned as the final result. It is obvious that the order to perform these actions plays an important role in this process. The simplest way to decide the order is to assume a fixed order among all actions, e.g., row 1 to row N followed by column 1 to column M . We will explain in the next section that the fixed ordering suffers from some inherent drawback that may be overcome by employing a dynamic ordering.

3.2. Complexity Analysis

In the first phase, a set of k biclusters (seeds) are generated. Thus, the time complexity of the first phase is $O(k \times (N + M))$ where N and M are the number of rows and columns of the matrix D while k is the number of the biclusters. The second phase is a series of iterations. During each iteration, each possible action of each row (or column) needs to be considered. There are k possible actions for a given row (or column). Thus, $(N + M) \times k$ actions have to be considered. In turn, the overall time complexity to evaluate all of these actions is $O((N + M)^2 \times k)$. The time complexity to perform an action is the same as to compute the gain of that action which is $O(N + M)$. There are $(N + M)$ actions to perform. Thus, the overall time complexity for an iteration is $O((N + M)^2 \times k)$, which implies that the complexity of the FLOC algorithm is $O((N + M)^2 \times k \times p)$ where p is the number of iterations till termination. Note that FLOC has less computational complexity than the Cheng-Church algorithm as it is typically the case where $p \ll N + M$.

3.3. Additional Features

As mentioned previously, the bicluster model can support many optional constraint specified by the user. In order to enforce the constraint, the basic FLOC algorithm needs to be modified. In the first phase where the set of initial biclusters are generated, the produced biclusters have to comply with the specified constraint. During the second phase where the iterative improvement is carried out, some action may be “blocked” temporarily (e.g., the gain is assigned to $-\infty$) during an iteration if it will result in the violation of some constraint. Only those actions that fully comply with the constraint will be performed. For example, if an action would cause the percentage of specified values for a gene or a condition in a bi-cluster falls below α , then the gain of this action will be assigned to $-\infty$. Furthermore, if all actions for a row or column are associated with $-\infty$ gain, then no action of this row or column will be performed at this iteration. It is obvious that the result produced by this modified version of the FLOC algorithm is guaranteed to satisfy the specified constraint.

3.4. Inclusion of Inverted Rows

Sometimes, it is also interesting to discover genes that are co-regulated but receiving opposite regulation. Their patterns essentially present mirror images. The FLOC algorithm can also find this type of patterns by allowing bicluster to contain inverted row(s). The rationale is that if a row resembles the mirror image of the rest set of rows in a bicluster, the inverse of this row (obtained by augmenting the additive inverse $-$ to each entry of the row) should obey the same regulation as the rest set of rows in the bicluster. The FLOC algorithm can be easily extended to discover this type of patterns by introducing three new actions for each **row**^d x with respect to each bicluster c .

- (1) *inserting the inversion of x in c* if x does not participate in c yet,
- (2) *deleting the inversion of x from c* if the inversion of x is currently in c ,
- (3) *inverting x in c* if either x or the inversion of x is currently in c .

Accordingly, the residue of an entry d_{ij} becomes

$$r_{ij} = \begin{cases} 0 & \text{if } i \notin I'_j \\ d_{ij} - d_{iJ} - d_{Ij} + d_{IJ} & \text{if } i \in I'_j \\ -d_{ij} - d_{iJ} - d_{Ij} + d_{IJ} & \text{if } -i \in I'_j \end{cases}$$

where

$$d_{Ij} = \frac{\sum_{i \in I'_j} d_{ij} - \sum_{-i \in I'_j} d_{ij}}{|I'_j|},$$

$$d_{iJ} = \begin{cases} \frac{\sum_{j \in J'_i} d_{ij}}{|J'_i|} & \text{if } i \in I \\ -\frac{\sum_{j \in J'_i} d_{ij}}{|J'_i|} & \text{if } -i \in I \end{cases},$$

$$d_{IJ} = \frac{\sum_{i \in I, j \in J} d_{ij} - \sum_{-i \in I, j \in J} d_{ij}}{v_{IJ}}.$$

At any time during the biclustering, a row may have two possible actions with respect to a cluster and a total of $2k$ actions are evaluated for each row, from which the best one is chosen and performed during the biclustering process. Beside the enriched set of actions, the rest part of the FLOC algorithm remains intact and so does the computational complexity.

4. Improvements of the Basic FLOC Algorithm

In the previous section, we presented the basic algorithm that still has room for improvement. In this section, we present the further optimization on two aspects of the FLOC algorithm.

^dThe actions for each column remain the same.

4.1. *Initial Size of a Bicluster*

As mentioned before, the parameter ρ is used to control the size of a bicluster in the initial assignment phase. The average number of rows and columns in a bicluster is $\rho \times N$ and $\rho \times M$, respectively. We experiment with different value of ρ . If the optimal bicluster size is very different from the initial (seed) bicluster size, then it would take more iterations to reach the optimal bicluster and the response time could be prolonged. In order to expedite the response time, it is beneficial to make the initial seed as close to the optimal bicluster size as possible. The optimal bicluster size is usually unknown in advance and may be significant disparate for different biclusters. Therefore, we generate initial biclusters of different sizes, i.e., different ρ values for different biclusters. In our experiments, we explore the effects of mixed ρ values and find that our biclustering algorithm can well adjust the bicluster size and discover both large and small biclusters.

4.2. *Order of Actions*

In the previous section, we assume the same order to perform actions at all iterations. This implies that for the set of rows and columns that are at the beginning of the order, their membership always has a higher priority to be changed than the rows and columns at the end of the order. This approach has one drawback. If a large number of actions with negative gain proceed a small number of actions with positive gain at the end of the performance list, then the set of actions with positive gain may never be given a full play. To solve the problem, we need to deploy a dynamic ordering among all actions. In this section, we discuss two possible ways to determine the action order.

4.2.1. *Random Order of Actions*

In this approach, the order of actions are randomly determined at the beginning of each iteration. This means that the membership of every column and row has the same priority to be changed. There are many algorithms to produce a random ordered sequence. Here we present one algorithm that can produce such a sequence. Assume that these actions are stored in an array. A series of action swapping is then performed, where each time two actions are randomly chosen and their positions are swapped. This random swapping procedure repeats g times. We experiment with various value of g and found that the randomness of the list is satisfactory where $g \geq 2 \times (M + N)$. Thus, we chose $g = 2 \times (M + N)$ to generate a random sequence in this paper. An action on a specific row or column may be performed first at one iteration, and performed very late at another iteration. With this technique, the problem mentioned above can be avoided since the positive actions at the end of one iteration may be at a very early position of the next iteration. Table ?? in the experimental results section shows the benefit of this improvement.

4.2.2. *Weighted Random Order of Actions*

In the above random order scheme, each action has the same probability to be assigned as the first action, the second action, and so on, regardless its gain. It means that a positive action has the same chance to be scheduled as the first action as a negative action. Intuitively, it is more desirable to perform actions with greater positive gain early so that its effect can be brought into play early. However, if we sort the actions in descending order of their gains and perform them accordingly, then we may only find the local optimal biclustering, but not the global optimal biclustering. As a result, we propose a weighted random order scheme to provide better ordering. Informally, this algorithm is very similar to the random order algorithm. The only difference is whether a swap would occur for two randomly picked actions. The rule of thumb is that if the action in the front has a larger gain than the one in the back, then the swap is less likely to occur. Let's consider two actions a_i and a_j and a_i is in front of a_j . The gain of the two actions are g_i and g_j , respectively. The probability $p(i, j)$ of swapping of a_i and a_j is a function of to $g_j - g_i$. Let R be the difference between the maximum gain and minimum gain of all actions. Then $p(i, j) = 0.5 + \frac{g_j - g_i}{2 \times R}$. When a_j has the maximum gain and a_i has the minimum gain, then the probability of swapping a_i and a_j is 1. On the other hand, if g_j is the minimum gain and g_i is the maximum gain, then $p(i, j) = 0$. In the case that $g_i = g_j$, $p(i, j) = 0.5$. Table ?? shows the improvement of the weighted random order algorithm. The weighted approach produces about 5% improvement in the quality of the final biclustering over the random ordering due to the fact that more positive actions are performed. In other words, the priority of membership migration favors the row or column whose gain is high. As a result, the weighted ordering can provide the best final biclustering.

5. Experimental Results

The FLOC algorithm is implemented with C programming language and is executed on an IBM AIX machine. We compare our FLOC algorithm with the Cheng-Church algorithm (CC algorithm) proposed by Cheng and Church (2000) on the yeast micro array containing 2884 genes under 17 conditions. Table 2 shows the performance of the two algorithms. Both algorithms are used to find 100 largest biclusters whose residue is less than 300.

Table 2. Performance comparison of FLOC and CC algorithms

	CC algorithm	FLOC algorithm
avg. residue	204.293	187.543
avg. volume	1576.98	1825.78
avg. gene num.	167	195
avg. cond. num.	12	12.8
time	12 min	6.7 min

At a glance, the FLOC algorithm is able to locate larger biclusters with smaller residue in substantially less amount of computation time. This is due to the fact that the FLOC algorithm is able to find highly coherent genes and condition which will lead to less residue. On the other hand, since random data is used to replace the discovered bicluster, the Cheng-Church algorithm tends to find smaller cluster with less coherence, i.e., larger residue. For almost all biclusters discovered by the Cheng-Church algorithm, they are subclusters of some bicluster discovered by the FLOC algorithm. Due to space limitations, we show two examples in this paper. Figure 6 shows two biclusters output by the FLOC algorithm, which entirely encompass some bicluster reported by the Cheng-Church algorithm. In Figure 6(a), the fine curves represent the expression levels of genes in bicluster 66 in Cheng and Church (2000) while the bold curve is the additional gene (YOR074C) discovered by FLOC, which also exhibits a similar behavior as the rest. It is interesting to know that the residue *decreases* with the inclusion of this gene because the addition of highly coherent matches will reduce the residue. Figure 6(b) presents another bicluster that contains two more conditions (condition 0 and condition 9) and six more genes (YAR007C, YAR008W, YBR089W, YDR097C, YJL187C, and YKL113C) than bicluster 95 reported by the Cheng-Church algorithm. The residue of this bicluster is 279.85 comparing to residue 311.7 of the bicluster reported by Cheng-Church algorithm.

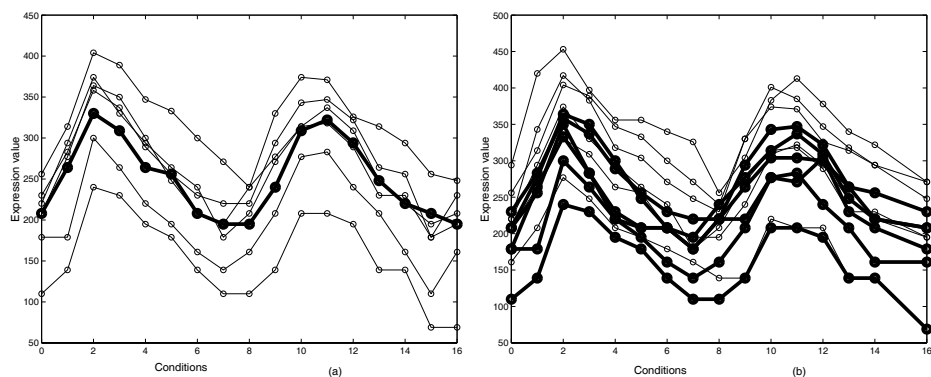


Fig. 6. Discovered Biclusters

After a thorough study, we found that the deficiency of the Cheng-Church algorithm results from a combined effect of the random interference and the vulnerability of the Cheng-Church algorithm to local optimum. As the Cheng-Church algorithm proceeds, more biclusters are located and are masked with random data, and the phenomenon of random interference becomes more severe, which impacts the discovery of large biclusters. For instance, the missed genes by Cheng-Church algorithm in Figure 6 is due to the fact that all (or majority) of the conditions on the missing genes have been included in some previous discovered biclusters. As a result, the exact value is replaced by some random numbers, and in turn, they are not included

in the new bicluster. On the other hand, the order of bicluster discovery has little importance in our FLOC algorithm. Therefore, FLOC is able to include these genes in the discovered bicluster. This can be further confirmed by Figure 7 which show the trends (after smoothing) of residues and volumes presented by the 100 biclusters from the Cheng-Church algorithm. Note that this results in smaller average volume over the discovered biclusters comparing to FLOC. Nevertheless the average residue under FLOC is still smaller as the additional genes are high quality matches. We believe that the random interference is the prime reason to cause this phenomenon.

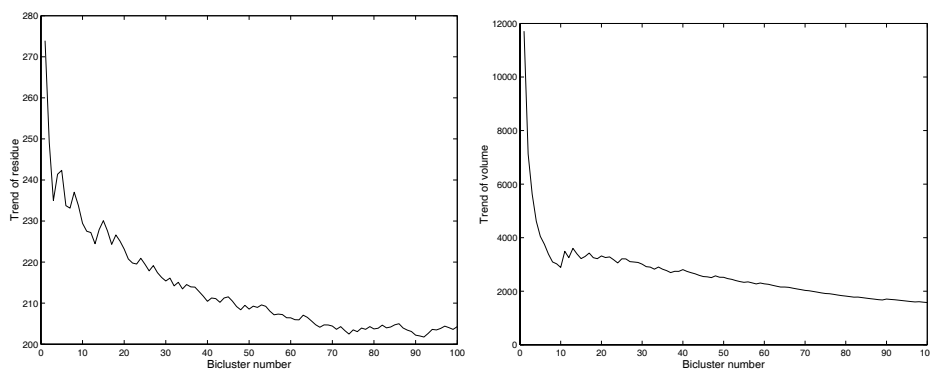


Fig. 7. Trend of Cheng-Church algorithm

5.1. Order of Actions

One of the improvements that we made to the FLOC algorithm is the action reordering at the beginning of each iteration. Table 2 shows the effects of action reordering. There are three ordering methods: fixed, random, and weighted random. Since the FLOC algorithms using random order and weighted random order are not deterministic algorithms (i.e., each run may produce different results), we also include the range in addition to the mean value for each measure we took. It is clear that the weighted random ordering is the winner. We believe that this is due to the fact that the weighted random order favors actions with large gains while still allow enough room for the algorithm to surpass local optimum.

5.2. Initial Bicluster Volumes

We also studied the effect of initial bicluster volume to the performance of FLOC. In all the test, the volume of initial bicluster follows the Erlang distribution with various mean and variance. We found that the performance of FLOC is robust and shows little dependency on the varying volume of initial biclusters. Even if the initial

Table 3. Performance comparison of different action orders in FLOC

	fixed order	random order	weighted random order
avg. residue	199.393	192.67 \pm 10	187.543 \pm 10
avg. volume	1696.11	1753.26 \pm	1825.78 \pm 100
avg. gene num.	177	188 \pm 10	195 \pm 10
avg. cond. num.	12	13 \pm 1	12.8 \pm 1
time	8.3 min	8.5 \pm 0.4 min	6.7 \pm 0.3 min

biclusters are far from the optimal ones, our FLOC algorithm can still “reshape” the biclusters into good quality.

6. Conclusions

In this paper, we proposed a generalized model of bicluster to address potential problems when dealing with gene expression data with missing values. Our observation of the random interference phenomenon, which is inherit to the “masking” operation in the Cheng-Church algorithm, motivates us to devise and deploy a novel algorithm, FLOC, to efficiently discover biclusters on gene expression data. The key feature of the FLOC algorithm is to trigger temporary “blocking” of certain action if such action has a potential to violate any criterion. This mechanism not only enables the FLOC algorithm to deliver better result faster than the Cheng-Church algorithm, but also enriches the bicluster model by supporting additional feature constraint at nearly no extra cost.

References

- Aach, J., Rindone, W., and Church, G. (2000) Systematic management and analysis of yeast gene expression data. *Genome Research*, 10, 431-445.
- Alizadeh, A. et al. (2000) Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling. *Nature*, 403, 503-510.
- Ben-Dor, A. and Yakhini, Z. (1999) Clustering gene expression patterns. *Proc. ACM RECOMB*, 33-43.
- Ben-Dor, A., Friedman, N., and Yakhini, Z. (2001) Class discovery in gene expression data. *Proc. ACM RECOMB*, 31-38.
- Bussemaker, H., Li, H., and Siggia, E. (2001) Regulatory element detection using correlation with expression. *Nature Genetics*, 27, 167-171.
- Califano, A., Stolovitzky, G., and Tu, Y. (2000) Analysis of gene expression microarrays for phenotype classification. *Proc. ISMB*.
- Cheng, Y. and Church, G. (2000) Biclustering of expression data, *Proc. ISMB*.
- Cho, R., Campbell, M., et al. (1998) A genome-wide transcriptional analysis of the mitotic cell cycle. *Mol. Cell*, 2, 65-73.

Eisen M. and Brown P. (1999) DNA arrays for analysis of gene expression. *Methods in Enzymology*, 303, 179-205.

Eisen, M., Spellman, P., et al. (1998) Clustering analysis and display of genome-wide expression patterns. *Proc. Natl. Acad. Sci. USA*, 96, 14863-14868.

Hartuv, E., Schmitt, A., Lange, J., Meier-Ewert, S., Lehrach, H., and Shamir, R. (1999) An algorithm for clustering cDNAs for gene expression analysis using short oligonucleotide fingerprints. *Proc. ACM RECOMB*, 188-197.

Kaufmann, L. and Rousseeuw, P. (1990) Finding groups in data – an introduction to cluster analysis, *Wiley series in Probability and Mathematical Statistics*.

Mirkin, B. (1996) *Mathematical Classification and Clustering*, Kluwer.

Segal, E., Taskar, B., Gasch, A., Friedman, N., and Koller, D. (2001) Rich probabilistic models for gene expression. *Bioinformatics*, 17, 243-252.

Sharan, R. and Shamir, R. (2000) CLICK: a clustering algorithm with applications to gene expression analysis. *Proc. ISMB*.

Tavazoie, S., Hughes, J., Campbell, M., Cho, R., and Church, G. (1999) Systematic determination of genetic network architecture. *Nature Genetics*, 22, 281-285.

Xing, E. and Karp, R. (2001) CLIFF: clustering of high-dimensional microarray data via iterative feature filtering using normalized cuts. *Bioinformatics*, 17, 306-315.

Yeast Micro Data Set, available at http://arep.med.harvard.edu/network_discovery.

Zien, A., Aigner, T., Zimmer, R., and Lengauer, T. (2001) Centralization: a new method for the normalization of gene expression data. *Bioinformatics*, 17, 323-331.