

極大2部クリーク列挙法による遺伝子発現モジュールの抽出

岡田吉史, 藤渕航, ホートン・ポール
独立行政法人 産業技術総合研究所 生命情報科学研究センター

遺伝子発現データから、ある特定の条件下において類似の発現動態を示す“遺伝子発現モジュール”を抽出する方法として、Biclusteringが注目されつつある。本研究では、極大2部クリークの全列挙法に基づくBiclustering法 (BiModule)を開発した。我々は、これを人工データおよび *S. Cerevisiae* の遺伝子発現データに適用し、既存の代表的なBiclustering法によるモジュール抽出結果との比較実験を行った。本報告では、BiModuleが他手法に比較して、人工的に埋め込まれたBiclusterをより高い精度で検出し、さらには、Gene Ontologyにおける機能アノテーションや既知のタンパク質相互作用を良く反映したモジュールを抽出可能であることを示す。

A biclustering method for finding gene expression modules based on a maximal biclique enumeration

Yoshifumi Okada, Wataru Fujibuchi, Paul Horton
Computational Biology Research Center, International Institute of Advanced Industrial Science and Technology (AIST)

In recent years, biclustering methods have been suggested to discover gene expression modules with shared expression behavior under certain experimental conditions. In this report, we propose a new biclustering method, BiModule, based on a maximal biclique enumeration algorithm. Comparative experiments to existing salient biclustering methods are performed to test the validity of biclusters extracted by BiModule using synthetic data and real expression data. We show that BiModule provides high performance compared to the other methods in extracting artificially-embedded modules as well as modules strongly related to GO annotations and protein-protein interactions.

1 はじめに

最近になって、遺伝子発現データにおける局所的なパターン(発現モジュール)を同定するBiclustering法が注目されつつある。従来のクラスタリング法では、全ての実験条件(以下、サンプル)にわたって類似の発現パターンを示す遺伝子群をクラスタとみなすのに対し、Biclustering法では特定のサンプルで類似の発現動態を示す遺伝子サブセットをモジュール(Bicluster)として抽出する。これまで、遺伝子の発現モジュールを発見するためのBiclustering法がいくつか提案されている[1-7]。それらは、特定のサンプルで共発現する遺伝子を組合せ論的に探索するNP完全な問題を扱うため、準最適なBiclusterをグリーディあるいは確率的に探索する近似アルゴリズムに基づく。これに対し、我々は、Biclusteringを2部グラフからの極大完全グラフ(極大2部クリーク)の抽出問題としてとらえ、準最適解としてではなく、全ての可能なBiclusterを列挙することで遺伝子発現モジュールを抽出する方法(BiModule)を

開発した。Biclusterは、遺伝子の頂点集合とサンプルの頂点集合間に枝が張られた2部グラフ上のクリークとみなせる。BiModuleは、極大2部クリークと等価なことで知られる飽和アイテム集合(以下、飽和集合)の高速列挙アルゴリズムが実装されている。これにより、全ての可能なBiclusterを現実的な時間内で列挙できるだけでなく、特定の評価関数やクラスタ妥当性指標を用いて、目的に沿ったBiclusterを選別できるようになる。また、BiModuleでは、一定の発現レベルを持つBiclusterだけでなく、発現変化の傾向をもつBiclusterを生成する工夫もなされている。

本研究の目的は、人工データおよび実際の発現データを用いて、BiModuleと既存手法の性能を統計的に比較評価することである。したがって、抽出されたBiclusterに関する生物学的意味についてはここでは議論しない。まず、人工データを用いた実験では、人為的に埋め込まれた正解モジュールの抽出精度を評価する。また、実際の発現データを用いた実

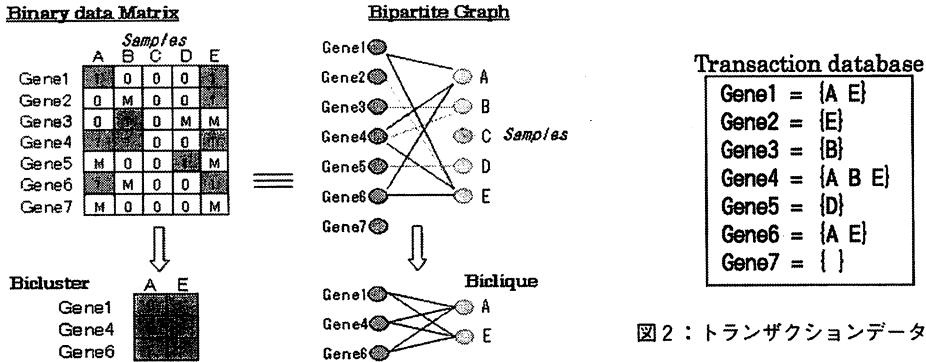


図2：トランザクションデータベース

1: up-regulation, 0: no change, M: missing value

図1：Bicluster と 2部クリークの関係

験では、*Saccharomyces Cerevisiae* (*S. Cerevisiae*) のデータセットを用いて、抽出された Bicluster が GO アノテーションおよびタンパク質間相互作用をどの程度良く反映しているかを評価する。

本稿の構成は以下のとおりである。まず、2章では Bicluster と 2部クリークとの関係について述べ、3章では極大 2部クリークと等価な飽和集合の定義について述べる。4章では BiModule のアルゴリズムについて説明し、5章で他手法との比較評価実験の方法と結果、6章を本研究のまとめとする。

2 Bicluster と 2部クリーク

Biclustering への入力、一般的に用いられる発現データ行列、すなわち行と列がそれぞれ遺伝子とサンプル、各セルが発現値をとる行列である。既存の Biclustering 法のいくつかは、発現値を離散化し、それらを単なる記号データとして扱う[1, 3, 7]。特に、発現値を 2 値化したデータ行列の Biclustering は、2部グラフから 2部クリークを抽出することと同義になる。図1は、2値データ行列における Bicluster と、2部グラフから抽出した 2部クリークの関係を示している。ここでの Bicluster は、特定のサンプル集合で 1 をとる遺伝子の集合で構成される部分行列であり、これは 2部グラフにおける遺伝子の頂点集合とサンプルの頂点集合の間で任意の 2 点間に辺を持つ 2部クリークに対応する。他のクリークに含まれない極大 2部クリークを列挙することで、極大の Bicluster を列挙できる。極大 2部クリークを多項式時間で列挙するアルゴリズムがいくつか提案されている[8,9]。本研究では、実装プログラムが公開されている LCM (Linear time Closed itemset Miner) を用いる[9,10]。これはトランザクションデータベースから飽和集合と呼ばれるパターンを列挙するプログラムとして実装されている。しかしながら、飽和集

合の列挙問題は、極大 2部クリークを見つけ出す問題と等価であるため[11]、飽和集合を列挙すれば、極大の Bicluster を見つけ出すことができる。次章では、飽和集合について説明する。

3 飽和集合と極大 2部クリーク

全アイテムの集合を $I = \{1, \dots, n\}$ 、その部分集合をアイテム集合 (itemset) と呼ぶ。D は全トランザクションの集合であり、トランザクションデータベースと呼ぶ。ここで各トランザクション T は I の部分集合である。あるアイテム集合 P を含む D のトランザクションを P の出現という。P が出現している集合を P の出現集合、その大きさを P の頻度と呼ぶ。ここでアイテム集合 P が、「P と同じ頻度を持ち、かつ、P を含むような集合が存在しない」とき P を飽和集合と呼ぶ。与えられた定数 θ 以上の頻度を持つ飽和集合を頻出飽和集合という。 θ はサポート値とよばれる。例として、図1の 2値データ行列から作成したトランザクションデータ (図2) を用いて飽和集合を説明する。図2の各行が、1つのトランザクション (遺伝子)、トランザクションに含まれる個々のアルファベットがアイテム (サンプル) に対応する。アイテム名は、数字や記号など、アイテムを識別できるものとする。アイテム A は、アイテム集合 $\{A\}$ と同じ頻度かつ $\{A\}$ を含む集合 $\{A, E\}$ が存在するので飽和集合ではない。アイテム集合 $\{A, E\}$ は、同頻度かつ、これを包含する集合がないので飽和集合である。

頂点がトランザクションとアイテムに対応し、トランザクションがアイテムを含むときにそれらに枝を張ると、2部グラフの極大 2部クリークと飽和集合が 1対1に対応する。図2の $\{A, E\}$ のようなサンプルの飽和集合が見つければ、それらに共通して応答する遺伝子 (Gene1, 4, 6) を特定でき、極大 2部クリーク、すなわち Bicluster を列挙できる。

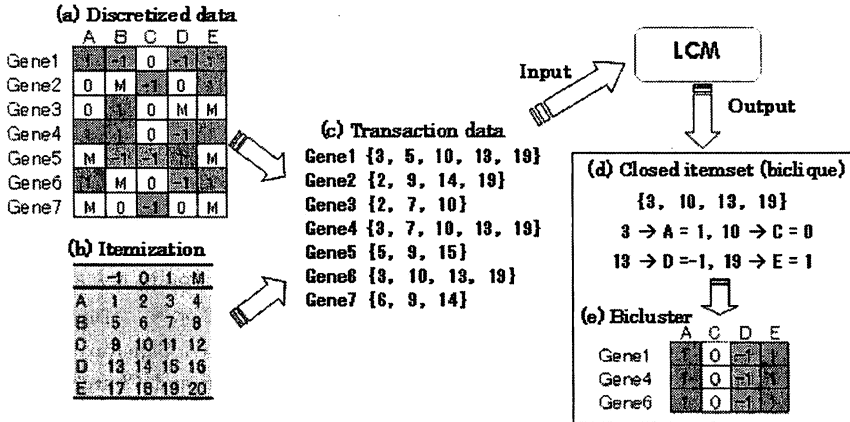


図3: BiModule の処理手順

4 方法

図3は, BiModule で行う処理の概略図である。まず, 与えられた発現データを正規化し, それらを多段階の階級に離散化する(4.1節)。次に離散データをトランザクションデータに変換する。ここでは, 2章で述べたような一定の値をとる Bicluste r だけでなく, サンプルにわたって遺伝子が発現傾向を示す Bicluste r を抽出できるような変換操作を行う(4.2節)。続いて, LCM アルゴリズムを用いて飽和集合をなすサンプルを列挙する(4.3節)。ここで, 他の飽和集合に含まれる飽和集合を除外する。最終的に残った飽和集合に関して応答する遺伝子セットを求めて Bicluste r を列挙する(4.4節)。以下, それぞれの処理の詳細を説明し, 最後に BiModule におけるパラメータ設定や実装について述べる(4.5節)。

4.1 データの正規化と離散化

まず, 各行に対して発現値が平均 0, 分散 1 になるように正規化する。次に, それらの正規化データをいくつかの階級に離散化する。図3aでは簡単のため, 階級数を3 (-1, 0, 1) とした場合を例示している。ここでデータ行列内のMは欠損値を表す。なお, 本研究では7階級 (-3, -2, -1, 0, 1, 2, 3) で離散化を行っている。各階級の値を階級値と呼ぶ。階級幅は, データの最大値と最小値の差を階級数で割った値とする。しかしこのとき, 最大値や最小値が極端な値(外れ値)をとる場合には, データが多く分布する値域において不当に大きな階級幅がとられてしまい, データに対する離散化が適切に行われぬ。そこで, 離散化の前処理として外れ値の置換処理を行った。まず, 反復切断補正法により, 正規化データの $\pm 2SD$ または $\pm 3SD$ 以上のデータを繰り返し除外し, 最終的な最大値・最小値を得る。除外された外れ値については, $+2(3)SD$ 以上, $-2(3)SD$ 以下の値をそれぞれ

最終的な最大値, 最小値に置き換える。離散化は, 発現データに対する正規化, 反復切断補正法および外れ値の最大値・最小値への置換の後に行う。

4.2 トランザクションデータの作成

サンプルにわたり遺伝子が発現変化を示す Bicluste r を生成可能にするため, Itemization という処理を行う(図3b)。これは, 各サンプルの全ての階級値にIDをつけ, どのサンプルでどの階級値を持つかをアイテム化する処理である。ここで, 行はサンプル名, 列は階級値および欠損値である。例えば, アイテム7は, サンプルBで階級値が1であることを表している。続いて, 各遺伝子をトランザクション, 各サンプルの階級値に付されたIDをアイテムとして, データ行列をもとにトランザクションデータベースを作成する(図3c)。このとき, 欠損値(M)を示すアイテムはトランザクションデータベースには含まない。Itemization 処理によって, 多段階の発現レベルをトランザクションデータに反映できるため, 一定の値をとる Bicluste r だけでなく, 発現変化を示す Bicluste r の抽出が可能となる。

4.3 サンプルに関する飽和集合列挙

LCM への入力は, トランザクションデータファイルとサポート値である。サポート値は Bicluste r を構成する遺伝子の最小数である。出力は, 図3dに示されるようなアイテムの飽和集合である。ここでは, {3, 10, 13, 19}が飽和集合として出力されている。これをサンプル名に直すには, Itemization 表を参照する。すると, このサンプル集合は{A, C, D, E}であり, それらの階級値は{1, 0, -1, 1}であることがわかる。

以下, LCM アルゴリズムの概要を説明する。アルゴリズムに関する詳しい説明や性能については, 文献[9]を参照されたい。

LCM では、prefix 保存飽和拡張という手法を用いて、飽和集合の探索を高速化する。飽和集合 P に対し、そのコアアイテム $\text{core}(P_i)$ を $\text{clo}(P_{\leq i})=P$ が成り立つような最小添え字を持つアイテムとする。ここで、 $\text{clo}(P_{\leq i})$ は、 i 以下の添え字を持つ P の閉包（出現集合の共通部分）である。このとき、 P_i が P の prefix 保存飽和拡張であるとは、あるアイテム $e \in \text{core}(P_i)$ に対して $P_i = \text{core}(P \cup \{e\})$ 、 $P_{\leq e-1} = P_{\leq e-1}$ が成り立つことである。 P_0 を、全てのトランザクションに含まれる集合の中で極大なものとする。ただし、 P_0 は空集合にもなりえる。 P_0 でない任意の飽和集合は、自身より頻出度大きい飽和集合の prefix 保存飽和拡張であり、さらに、そのような飽和集合は唯一に決まる。したがって、 P_0 から出発して再帰的に prefix 保存飽和拡張を生成することで、全ての頻出飽和集合を深さ優先で発見できる。

LCM を実行することにより、発見されたサンプルの飽和集合が逐次的に列挙される。

4.4 フィルタリング処理と Biclusters 出力

LCM に入力するデータサイズが大きい場合には、非常に膨大な数の飽和集合が生成される。しかしながら、そのほとんどは他の飽和集合に包含されているか、あるいは大部分が重複したものである。そこで、出力された飽和集合を以下のスコア S に基づいてソートした後で、上位の飽和集合に 80% 以上包含される飽和集合を除外する。

$$S = A \times \log_2(g) \times \log_2(s)$$

ここで、 A は飽和集合のサンプルが持つ階級値の絶対値の平均値、 g は飽和集合のサンプルに回答する遺伝子数、 s は飽和集合に含まれるサンプル数である。これにより、サイズがある程度大きく、他に大部分を包含されない Biclusters を取り出せる。また、階級値の絶対値の平均を考慮することで、平均的に発現レベルの低い Biclusters を排除できる。最終的に残った飽和集合に関して回答する遺伝子セットを求めて Biclusters を列挙する（図 3 e）。

4.5 実装とパラメータ設定

BiModule は、LCM による飽和集合列挙の部分を除き、全て Java 言語で開発されている。なお、LCM は C 言語で実装されている。BiModule の入力は、正規化前の発現データ行列および、抽出する Biclusters サイズに関するパラメータ M_g と M_s である。ここで M_g と M_s はそれぞれ、出力される Biclusters の遺伝子数およびサンプル数に関する最小値である。

5 評価方法および結果

5.1 比較に用いた Biclustering 手法

Prelic らは、6 つの代表的な Biclustering 法を、1) 人為的に正解モジュールを埋め込んだデータセット、および 2) 実際の発現データセットに適用し、生成される Biclusters の妥当性評価を行った[1]。ここで用いられたデータセットおよび Biclustering 結果は[12]で公開されている。Prelic らが比較に利用した 6 つの手法を以下に示す：Divide-and-conquer Algorithm[1], Iterative Signature Algorithm[2], SAMBA[3,4], Cheng and Church's Algorithm[5] Order Preserving Submatrix Algorithm[6], xMotif[7]。以下本稿では、それぞれ BiMax, ISA, SAMBA, CC, OPSM, xMotif と呼ぶ。本研究では、Prelic らによって提供されている人工データと実際の発現データに BiModule と上記の 6 手法を適用した結果を比較する[12]。なお、BiModule に入力したパラメータは $M_g=7$, $M_s=4$ である。

5.2 人工データを用いた評価

人為的にモジュールを埋め込んだデータを用いて、BiModule および既存手法がどの程度正確にそれらのモジュールを抽出できるかを、「遺伝子適合スコア」と呼ばれる指標に基づいて評価する。ここでは、以下の 3 つのタイプのモジュールが埋め込まれたデータ行列を扱う。

- ・ Constant モジュールを含むデータ行列
- ・ Coherent モジュールを含むデータ行列
- ・ Overlapping モジュールを含むデータ行列

Constant モジュールとは全てのセルが同一の値を持つモジュールであり、Coherent モジュールとはサンプルにわたって値が変化し、全ての遺伝子が同じ発現傾向を示すモジュールである。Constant および Coherent タイプではともに、100 遺伝子×50 サンプルのデータ行列内に、10 遺伝子×5 サンプルのモジュールが 10 個含まれている。なお、10 個のモジュールは互いに重複しないものとする。

Overlapping モジュールは、データ行列内の対角線上に並ぶ 10 個の Coherent モジュールが重複度 d で重なり合っているモジュールである。 $d=0$ はモジュール間で全く重なりが無い場合、 $d=10$ は隣接する 2 つのモジュールが遺伝子側に 10、サンプル側に 10 の重複（すなわち $10 \times 10 = 100$ セル分）を持つことを意味する。

抽出された Biclusters の妥当性を評価するため、以下の遺伝子適合スコアを用いる。

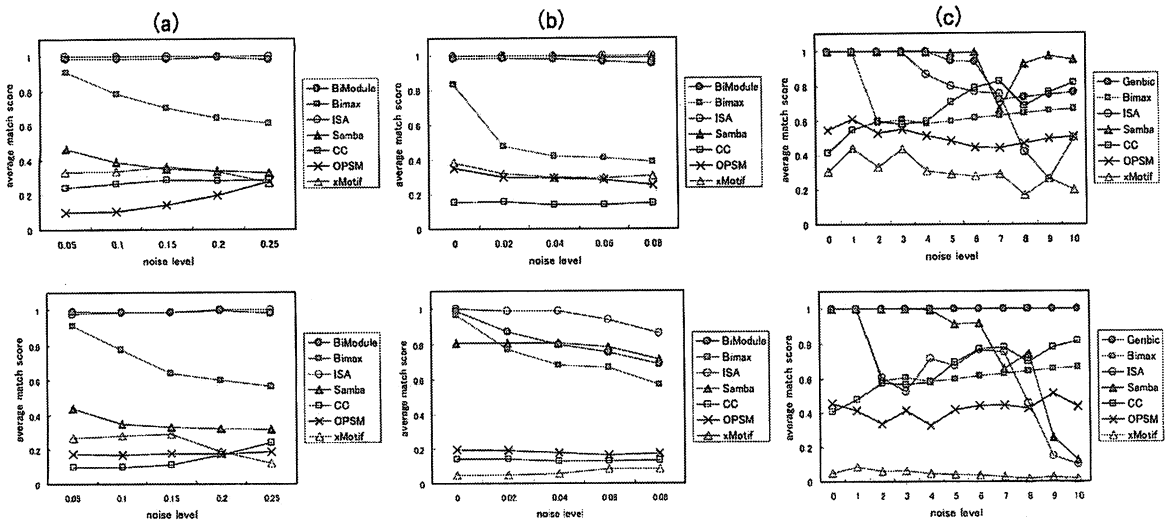


図4：遺伝子適合スコア

$$GMS_G(M_1, M_2) = \frac{1}{|M_1|} \sum_{(G_1, S_1) \in M_1} \max_{(G_2, S_2) \in M_2} \frac{|G_1 \cap G_2|}{|G_1 \cup G_2|}$$

ここで、 M_1 , M_2 はモジュール (または Bicluster) のセット、 G , S はそれぞれ1つのモジュールに含まれる遺伝子セットおよびサンプルセットである。このスコアは、 M_1 の各モジュールについて、 M_2 に含まれるモジュールとの最大適合スコアの平均値を表わしている。今、 M_{opt} を人為的に埋め込まれたモジュールセット (正解セット)、 B をある Biclustering 法から出力された Bicluster セットとすると、 $GMS_G(B, M_{opt})$ は出力された Bicluster と正解モジュールセットとの一致度を表し、 $GMS_G(M_{opt}, B)$ は正解セットをもれなく抽出する再現度を意味する。どちらのスコアも $M_{opt} = B$ のときに最大値1をとる。

ここでは、各 Biclustering 法の「データに含まれるノイズに対する感度」と「ノイズを含まないが重複しているモジュールの抽出精度」の2つの観点から評価した結果を示す。図4は、ノイズを含む (a) Constant および (b) Coherent モジュール、(c) ノイズを含まない Overlapping モジュールに関する遺伝子適合スコアの一致度 (上段) と再現度 (下段) を示している。(a) と (b) については、正規分布から得られるランダムな値をノイズとして、元のデータ行列に付加した。図(a), (b) および (c) は異なるノイズレベル (標準偏差) および重複度 ($d=0, 1, \dots, 10$) において、それぞれ10個の異なるデータ行列を生成し、それらに対するスコアの平均値をプロットしている。

Constant モジュール (図4 a) では、BiModule と ISA はノイズの増加の影響をほとんど受けず、一致度と再現度ともに98%以上を維持している。続く BiMax はノイズの増加とともに正解モジュールの

抽出精度が、一致度、再現度ともに91% (noise level=0.05) から60%程度 (noise level=0.25) に減少している。SAMBA, CC, OPSM, xMotif に関しては全てのノイズレベルにわたり、一致度、再現度ともに50%以下であった。Coherent モジュール (図4 b) では、一致度については、BiModule, ISA, そして SAMBA が全てのノイズレベルにわたって95%以上の精度を示した。BiMax は、Constant モジュールの場合よりもノイズの影響を強く受けている。再現度については、BiModule, BiMax, ISA, SAMBA が上位グループを形成しているが、中でも ISA が最も高く安定した精度を示した。Overlapping モジュール (図4 c) では、一致度、再現度ともに BiModule が良好な精度を示している。特に、再現度に関して言えば、重複度の増加に関わらず全てのモジュールをもれなく抽出できている。SAMBA は一致度に関しては良好な成績だが、再現性では重複度の増加とともに急激な抽出漏れが見受けられる。ノイズに最も頑健であった ISA は、重なりが少ない ($d=1$ or 3) 場合はどちらのスコアも100%を示すが、重複度が増すにつれ抽出精度は減少し、 $d=9$ では一致度、再現度ともに20%代まで落ち込んでしまう。

全体をとおして、BiModule が最も安定かつ高い精度でモジュールの抽出を行っている。これは、BiModule における以下の操作が有効に働いていると思われる。BiModule は、4章で述べた Itemization によって全てのセル値が同一な Bicluster (Constant モジュール) だけでなく、サンプルにわたる発現変化の傾向を持つ Bicluster (Coherent モジュール) を抽出できる。また、離散化により、微弱なノイズを含む値や外れ値の影響を取り除くため、ノイズに

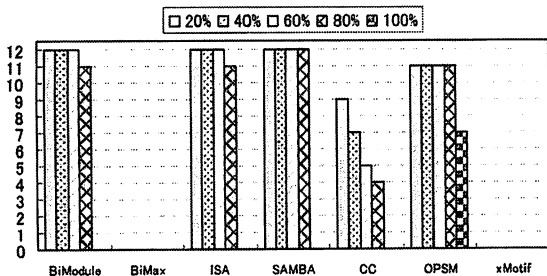


図5：有意なGOカテゴリ ($p < 0.001$) を含む Bicluster 数

対する感度を抑えることができる。さらに、BiModule では、Bicluster に対してサンプル数と遺伝子数のバランスを考慮したフィルタリングを行うため、CC や xMotif のようにどちらか一方の次元のみに拡張した Bicluster は抽出されない。

5.3 遺伝子発現データを用いた評価

ここでは Gasch ら[13]によって提供されている *S. Cerevisiae* の遺伝子発現データセットを用いる。これは、さまざまなストレス条件下で測定された 173 サンプルに関する 2993 遺伝子の発現データを含む。我々は、このデータセットに個々の手法を適用し、得られた Bicluster が、Gene Ontology (GO) アノテーションと既知のタンパク質間相互作用をどの程度反映しているかについての評価を行う。なお、BiModule に入力したパラメータは、 $M_g=80$, $M_s=5$ である。通常、Biclustering のアルゴリズムによって、生成される Bicluster 数は異なる。そこで、出力 Bicluster 数が 100 を超える場合には、それらのサイズが大きい順に 100 個を最適な Bicluster とみなして実験を行った。結果として、各手法が生成した Bicluster 数は、BiModule (20), BiMax (100), ISA (67), SAMBA (100), CC (100), OPSM (12), xMotif (100) であった。

5.3.1 Gene Ontology (GO) アノテーション

Bicluster 内の遺伝子について統計的に有意に多く含まれる GO カテゴリを調べるため、フリーで利用可能な GO 解析ツール BiNGO を用いた[14]。図5は、有意と判定された GO カテゴリ ($p < 0.001$) を含む Bicluster の個数を、それらのカテゴリを含む割合別 (20, 40, 60, 80, 100%) に示したものである。ただし、ここでは Bicluster 数が最も少ない OPSM に合わせ、各手法において Bicluster のサイズが大きい順に 12 個とりだした場合を示してある。図から、BiModule, ISA, SAMBA, OPSM の4手法が、GO アノテーションを良く反映した Bicluster を構成していることがわかる。BiModule, ISA については、

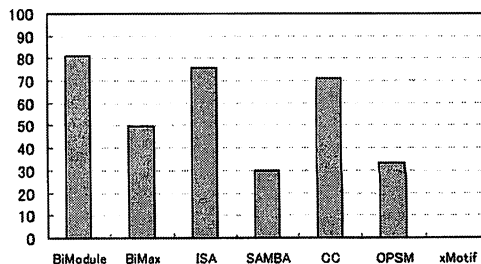


図6：タンパク質間相互作用に関わる遺伝子ペアが有意に含まれる Bicluster の割合

11 個の Bicluster において 80% 以上、SAMBA については、全ての Bicluster で 80% 以上の遺伝子が有意な GO カテゴリを有している。これらの3手法は、人工データからのモジュール抽出においても好成績を示していた。一方、OPSM は、人工データに対してはノイズの影響を受けやすいだけでなく、重複したモジュールの抽出精度も比較的良かった。しかし現実のデータにおいては、11 個の Bicluster において 80% 以上、7 個の Bicluster において全ての遺伝子が有意な GO カテゴリを持っていた。CC に関しては、上記の4手法に比較して有意な GO カテゴリを含む Bicluster が少なく、それらのカテゴリを持つ遺伝子の割合も低かった。これは、CC で得られるサイズの大きい Bicluster は、発現値が 0 のセルを多く含むため、生物学的に意味のある遺伝子が含まれにくいことが考えられる。BiMax と xMotif では、有意な GO カテゴリを持つ Bicluster は存在しなかった。BiMax では発現データを 2 値データに離散化してしまうため、データの持つ情報が失われ、生物学的に関連性の無い Bicluster が多数生成されている可能性がある。また、xMotif では、得られる Bicluster のほとんどが 1 つの遺伝子しか含まなかったため、有意な GO カテゴリは見つからなかった。

5.3.2 タンパク質間相互作用

Bicluster とタンパク質間相互作用との関連性を調べるため、DIP データベース[15]で提供される相互作用データを用いた。このデータには、14,742 組のタンパク質間相互作用が、遺伝子のペアで記述されている。与えられた Bicluster がどの程度タンパク質間の相互作用を反映しているかを見積もるスコアとして、「Bicluster 内において相互作用していない遺伝子ペアの割合」を用いる。このスコアは、ランダムに選ばれた遺伝子グループよりも、Bicluster 内における遺伝子群に関して有意に小さい値をとることが期待される。そこで、各手法で生成された個々の Bicluster と同じサイズのランダムな遺伝子グループを 1000 セットずつ作成し、それぞれの Bicluster に

関するスコアがランダムな遺伝子グループよりも有意に小さいかどうかを Z-検定により判定した。図 6 は、ランダムな遺伝子グループに対し、スコアが有意に小さい BiCluster($p < 0.001$) の割合である。BiModule は 81% の BiCluster においてランダムな遺伝子グループに対して有意にスコアが小さく、7 手法の中で最も良くタンパク質間相互作用を反映していた。これに続き、ISA が 75.8%、CC が 71% と好成績を示している。CC については、サイズの小さい BiCluster で有意差が見られる傾向にあった。

6. まとめ

本研究では、飽和集合列挙アルゴリズムに基づいて、遺伝子とサンプル間で形成される極大 2 部クリークを全て列挙し、他に包含されない BiCluster を抽出する BiModule を提案した。これを人工データおよび *S. Cerevisiae* の発現データに適用し、代表的な 6 つの手法との比較を統計的観点から行った。以下は、本研究で得られた結果のまとめである。

1. 人為的に埋め込まれたモジュールの抽出実験では、全体をとおして BiModule、ISA、SAMBA が好成績であった。特に、BiModule はノイズの影響を受けにくいだけでなく、重複度合いによらず、重なり合ったモジュールを適切に抽出した。概して、BiModule が最も安定して高い精度を示した。
2. *S. Cerevisiae* の発現データを用いた実験では、総合的に見て、BiModule と ISA が GO アノテーションおよびタンパク質間相互作用との関連度が高い BiCluster を生成した。

以上の結果から、BiModule が遺伝子発現モジュールの発見に有効であることが示された。

BiModule は、大きく分けて 1) データの正規化と離散化、2) トランザクションデータベースの作成、3) 飽和集合列挙、4) フィルタリングの 4 つの部分に分けられる。しかしながら、本稿では、これらのどの部分がモジュール発見の性能を左右しているかを明らかにしていない。例えば、離散化における階級数、フィルタリングでのスコア関数、あるいは入力パラメータの最小遺伝子数や最小サンプル数によって、モジュール発見精度にどのような違いが見られるかを明確にする必要がある。また、本研究では統計的指標に基づく評価に留まったが、今後は、抽出されたモジュールに含まれる遺伝子の機能や転写制御、あるいは分子間ネットワークといった異なるレベルにおける生物学的知見との照らし合わせから評価を行っていく。

参考文献

1. Prelic, A. et al. (2006) A Systematic comparison and evaluation of biclustering methods for gene expression data. *Bioinformatics*, 22, 1122-1129.
2. Ihmels, J. et al. (2004) Defining transcription modules using large-scale gene expression data. *Bioinformatics*, 20, 1993-2003.
3. Tanay, A. et al. (2002) Discovering statistically significant biclusters in gene expression data. *Bioinformatics*, 18 (Suppl. 1), S136-S144.
4. Tanay, A. et al. (2004) Revealing modularity and organization in the yeast molecular network by integrated analysis of highly heterogeneous genomewide data. *Proc. Natl Acad. Sci. USA*, 101, 2981-2986.
5. Cheng, Y. and Church, G. (2000) Biclustering of expression data. *Proc. Int. Conf. Intell. Syst. Mol. Biol.*, 93-103.
6. Ben-Dor, A., Chor, B., Karp, R. and Yakhini, Z. (2002) Discovering local structure in gene expression data: the order-preserving sub-matrix problem. *Proc. of the 6th Annual Int. Conf. on Computational Biology*, ACM Press, New York, NY, USA, 49-57.
7. Murali, T.M. and Kasif, S. (2003) Extracting conserved gene expression motifs from gene expression data. *Pac. Symp. Biocomput.*, 8, 77-88.
8. Tukiya, S. et al. (1977) A new algorithm for generating all the maximum independent sets, *SIAM Journal on Computing*, 6, 505-517.
9. Uno, T. et al. (2004a) An efficient algorithm for enumerating closed patterns in transaction databases, *Lecture Notes in Artificial Intelligence*, 3245, 16-31
10. <http://research.nii.ac.jp/~uno/codes-j.html>
11. Boros, E. et al. (2002) On the complexity of generating maximal frequent and minimal infrequent sets. *Proc. of STACS 2002*, Springer, Berlin, 133-141.
12. <http://www.tik.ee.ethz.ch/sop/bimax/SupplementMaterials/Biclustering.html>
13. Gasch, A.P. et al. (2000) Genomic expression programs in the response of yeast cells to environmental changes. *Mol. Biol. Cell*, 11, 4241-4257.
14. Maere, S. et al. (2005) BiNGO: a Cytoscape plugin to assess overrepresentation of Gene Ontology categories in biological networks. *Bioinformatics* 21, 3448-3449.
15. <http://dip.doe-mbi.ucla.edu/>