

最尤推定法を利用した異数倍数体の遺伝子発現量定量化手法の開発

理学専攻 情報科学コース 山田 恵

1 はじめに

網羅的に遺伝子発現量を得る方法として、RNA-seqがある(図1)。RNA-seqは遺伝子の転写物であるmRNAを大量に読み、既知の遺伝子配列やゲノム配列と対応(アラインメント)させる。対応する配列の本数を数える事で遺伝子の発現量を計測する。近年のシーケンサでは、連続して読める配列(リード)の長さは100塩基程度と短い、1千万を超える本数が容易に計測可能となったことで、RNA-seqが実現可能となった。一方RNA-seqには欠点もあり、対象種が持つ複数の遺伝子がほぼ同じ配列を有する場合、リードがどの遺伝子由来か特定できず、発現量の測定を誤る可能性がある。

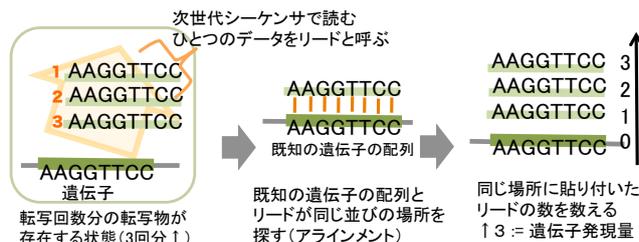


図1: RNA-seqの概要

本研究は、このRNA-seqの利点と欠点を上手に利用することで、異数倍数体における遺伝子発現の定量を行う。異数倍数体とは、異なる二種が交配し、子が両親のゲノムを持つ種で、生物の新規機能獲得プロセスの解明に有用であると考えられているが、両親のゲノム両方とも働くのか、働くとすれば同程度に働くのか、など未知の点が多い。本研究では、これらのゲノムに含まれる遺伝子の発現を同定する手法を確立する。異数倍数体の例として、本研究ではCardamine属の*C.flexuosa*を用いる。この種は*C.amara*と*C.hirsuta*を親種に持つ。親種は、同属で進化的距離が近いため、互いに数塩基のみが異なる、ほぼ同一遺伝子集合を有している(図2)。この種に対しRNA-seqを行うと、上記に述べた欠点によって、遺伝子の発現を定量化することが困難であるので、この問題点を解決する。

2 手法

RNA-seqでは、リードをアラインメントするゲノムが必要だが、*C.amara*と*C.hirsuta*のゲノムは未知のため、Cardamine属の近縁で、ゲノムや遺伝子が既知の*A.thaliana*を利用する。図2のように、*A.thaliana*は遺伝子X、*C.amara*は遺伝子X'、*C.hirsuta*は遺伝子X''、*C.flexuosa*はX'、X''を持つとすると、X、X'、X''は互いに類似した配列を持つ。子である*C.flexuosa*のX'、X''由来のリードを考えると、類似の配列を持つ*C.amara*のX'と*C.hirsuta*のX''の配列両方に対応する。RNA-seqでは、このような異なる遺伝子由来の配列が同一の遺伝子に対応することが問題となるが、今回は親種間の遺伝子配列に数塩基の違いがある事を利用していずれの親由来であるか分離する。

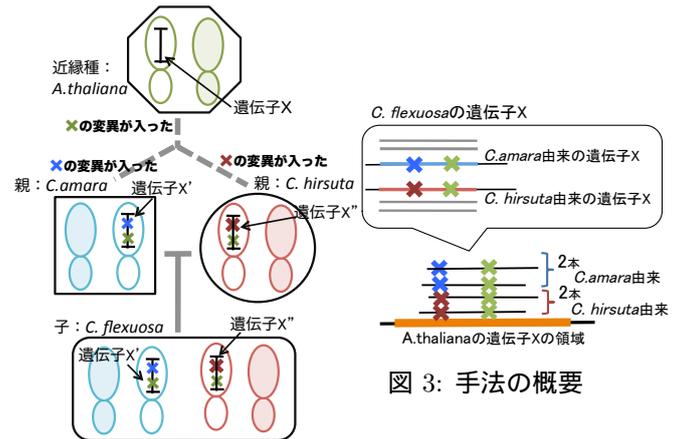


図2: 実験対象の種

2.1 変異の種類と求める方法

*A.thaliana*のゲノムと*C.flexuosa*のリードの間の変異には、親種特有の変異、祖先種から親種へ進化する過程で起きた変異(以下、それぞれ親由来の変異、祖先由来の変異)の2種類がある。図2の、*C.flexuosa*のX'、X''上のx印の、青と赤のx印が前者で、緑のx印が後者を示す。図2のX'、X''のリードを、*A.thaliana*のXにアラインメントすると図3のように表せる。x印は図2の同じ色のx印に相当し、x印の赤は*C.hirsuta*由来のX''の変異と同じで、青は*C.amara*由来、緑は祖先由来の変異とみなせる。*A.thaliana*のXと異なる部分に注目することで、*C.flexuosa*のX'、X''の変異が得られる。

図示で色が付けられているが、実データでは変異の場所と変異の塩基の種類が分かるのみなので、変異が、親由来か祖先由来かを見分ける方法が必要になる。

2.2 遺伝子発現量の算出

変異の種類に注目すると、リードの組み合わせは2種類ある。親種特有の変異(赤と青のx印)を見ると、リードの本数は赤と青で2本と2本に、祖先由来の変異(緑のx印)を見ると、緑だけで4本と0本となる。各変異が、どちらに由来するものかは分からないので、実際には各変異位置での変異の有無の比率がわかるのみである。祖先由来の変異に関わるリードの本数の比は0に近く、親由来の変異に関わるリードの本数の比は、他の親特有の変異の比に近いと考えられる。この特徴を用いて、変異が親由来か祖先由来かを見分ける。ここでは特定の変異に着目して親種の比率を求めたが1点のみの推定ではシーケンサーエラーによる影響も多く精度が低い。そこで本研究では、リード毎にいずれの親由来であるかを決定し、その本数を遺伝子全体にわたって数える事で計算の精度を向上させる。

2.3 アルゴリズム

次の4つのステップからなる。ステップ1では、アラインメントを行い、祖先種の遺伝子配列からの変異を全て抽出する。ステップ2では、抽出した変異から、親由来の変異を求める。ステップ

図3: 手法の概要

3では、各リードを親種別に分ける。ステップ4では、片親由来と判定されたリード数を数え、両親のリード数の比と遺伝子の発現量から、親ごとの遺伝子発現量を求める。以下詳細を述べる。

ステップ1：子のリードを祖先種のゲノムに、いくつかの変異を許容して全リードをアラインメントする。以降、遺伝子ごとに計算する。変異が遺伝子上の k 箇所にあるとき、変異を有する位置の集合を $M = \{M_1, M_2, \dots, M_k\}$ で表す。 $R1_i, R2_i$ を、 M_i の位置を含むリードを M_i の塩基の種類ごとに分けたリードの集合とする。ここでは変異の有無のみを考えるので、2種類が考えられる。 $R1_i$ の要素数を $|R1_i|$ と表す。

ステップ2：変異 M_i に着目した時、この変異は親由来か祖先由来の2種類の可能性がある。RNA-seqはランダムに配列を選択する手法なので、変異 M_i の状態が起こる確率 $P(M_i)$ は、ポアソン分布を用いてモデル化できる。ポアソン分布は平均 λ 、確率変数 k のとき、 $\lambda^k e^{-\lambda} / k!$ と書ける。 λ_p を親由来の平均、 λ_a を祖先由来の平均とすると、 M_i が親由来である確率は、 $P(M_i) = \lambda_p^{|R1_i|} e^{-\lambda_p} / |R1_i|!$ 、祖先由来である確率は、 $P(M_i) = \lambda_a^{|R1_i|} e^{-\lambda_a} / |R1_i|!$ と表せる。 λ_a は任意の値 t ($t \leq 1$, 実験では0.1) のとき、 $\lambda_a = t \times (|R1_i| + |R2_i|)$ で計算する。それぞれの変異は独立な事象であるとする。変異 M の尤度 $P(M)$ は $P(M_1) \times P(M_2) \times \dots \times P(M_k)$ で表せるので、この尤度が最大になるように、すべての変異を親由来か祖先由来かに分ければ良い。この時、親由来と考えられる変異セットを $M_P (\subseteq M)$ とする。

予め親種の平均 λ_p がわかっている場合には、簡単に M_P を決定することができるが、 λ_p は求めたい変数であり不明である。つまり、 $P(M)$ を最大化するには、最悪全ての考えうる M_P 及び λ_p を試さねばならない。これには膨大な時間を要する上に、現在考えているのは単一の遺伝子なので、実際にはこの計算を2万以上の遺伝子に対して行わなければならない。そこで、現実的な時間で計算できるよう、以下の工夫を行った。(1) M_P に対して、親由来を表すポアソン分布の平均 λ_p は、大体 M_P に属する変異の平均に近いと考えられるので、次の方法で計算する。 $I = \{i \mid M_i \in M_P\}$ とすると、 $k = |R1_i|$ のとき、

$$\lambda_p = (|R1_i| + |R2_i|) \times \left(\frac{\sum_{i \in I} |R1_i|}{\sum_{i \in I} |R1_i| + \sum_{i \in I} |R2_i|} \right)$$

(2) $|R1_i|/|R2_i|$ が小さいもの程、 λ_p を小さくする可能性が高いので、この比率が1に近いものから順に M_P に入れて、尤度が最大となる組み合わせを計算した。最大尤度を取る I を I' とする。

ステップ3：ステップ2までで求めた結果を元に、各リードに親ラベル A もしくは B を割り当てる。親ラベルは、 $|R1_i|/(|R1_i| + |R2_i|)$ が大きい変異 $i \in I'$ から順に割り振る。一度親ラベルが決まったリードはラベルを変更せず、次の候補へ移る。候補がなくなるまで繰り返す。

この割当を初期値とし、リード数の多い $i \in I'$ から順に、割り振ったラベル A, B を反転させ、尤度が大きくなるかを調べる。これは、ラベルが親種を取り違えていないかどうかを確認するための手順である。以上の手法で、リードごとに由来する親種を特定する。

ステップ4：RNA-seqによる遺伝子発現量はRPKM[1]と呼ばれる。ステップ3で求めた、両親それぞれのリード数の比とRPKMの積が今回求めたい両親由来別の遺伝子発現量となる。

3 実データを用いた実験の結果

次世代シーケンサ SOLiD5 で観測した *C. flexuosa* のリード約6千万本に対し、本手法を適用した。アラインメント先に *A. thaliana* (TAIR10) [4] のゲノムを選び、ゲノム上の反復配列を省くため RepeatMasker[2] を用い、アラインメントソフトには SHRiMP[3] を用いた。

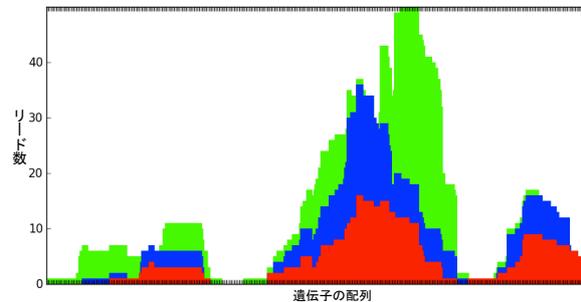


図4: 遺伝子上の親由来の変異を持つリードの分布

図4は、一つの遺伝子に対し本手法を適用した結果を示している。横軸に遺伝子上の位置、縦軸に塩基ごとのリードの本数を表す。色分けは、リードごとに判定した親種に対応し、赤と青が片親、緑は由来が不明のリードである。この遺伝子の両親種ごとの遺伝子発現量は、それぞれ5.93, 6.56だった。グラフからは、親種ごとの発現量を求めるのに使われるリードが、ほぼ遺伝子の範囲上を覆っていて全体的な変異の傾向から発現量を求められていることが分かる。

4 まとめ

本研究では、異数倍数体の *C. flexuosa* の各遺伝子について、親種を区別して遺伝子発現量を計測する手法を開発した。本研究では、対象種である異数倍数体とその親種のゲノムと遺伝子が同定されていない状態で、近縁種との遺伝子セットが類似しているが完全には一致しないという特徴を使う事で、親種別の遺伝子発現量を求める手法を開発した。

謝辞

本研究にあたり、ご指導くださいました東京工業大学情報理工学研究所 瀬々潤准教授と、貴重なデータとご助言を賜りましたチューリッヒ大学 清水健太郎教授、清水理恵研究員に深く感謝いたします。

参考文献

- [1] Mortazavi, et al. Mapping and quantifying mammalian transcriptomes by rna-seq. Nat Methods, 5(7):621-628 (2008).
- [2] Smit, AFA, Hubley, R & Green, P. RepeatMasker Open-3.0. 1996-2004. <http://www.repeatmasker.org>
- [3] Rumble, S.M. et al. SHRiMP: accurate mapping of short color-space reads. PLOS Comput. Biol. 5, e1000386 (2009).
- [4] <http://www.arabidopsis.org/>