

テンソル分解を用いた二重鎖切断の起こりやすいヒト遺伝子の推定

Conjecturing human genes that are easy to be double strand breaks by tensor decomposition

物理学専攻 保坂 伸生

Physics major Nobuo Hosaka

1. 研究背景

二重鎖切断 (DSB) とは DNA 損傷の一つで、DNA 損傷というのは DNA を構成する塩基と呼ばれる物質が損傷することを指す。中でも DSB は損傷のレベルが高く、二重螺旋構造を持つ DNA の二本どちらにも損傷が起きている状態を指す。この DSB を含む DNA 損傷は太陽光に含まれる紫外線や呼吸によって生じる活性酸素、食物に含まれる発がん物質などによって生じ、日常的に起こる危険性をはらんでいる。

DSB は DNA 損傷の中でも修復が難しく、修復されない場合、細胞が死を選択してしまい、過剰に細胞が減ると病気になってしまうが、修復されたとしても後述する変異によって遺伝子を喪失してしまう可能性がある。DSB 以外の修復は比較的簡単で、DNA はアデニン (A)、チミン (T)、グアニン (G)、シトシン (C) の 4 つの塩基で構成されていて、中でも A と T、G と C が結合することで二重螺旋構造を作る。そのため、一本のみの損傷であれば、反対側の塩基を読み取ることで容易に修復できる。これが二本とも損傷している場合はとても難しく、修復方法には相同組換え (Homologous Recombination (HR)) と非相同末端結合 (Non-Homologous End Joining (NHEJ)) の 2 つがある。HR では、損傷した DNA 塩基配列に似た配列のものを鋳型として使い、損傷した部位を作り直す。HR は比較的精度が良いがそれでも、似た配列を使っていることから若干の違いが生じる可能性があり、また HR は常にできるわけではなく細胞周期の一部でしか行われていない。NHEJ は損傷した先端部分を両者ともに削り、平坦にしたところでそのまま結合するといったものになっている。これは損傷部位を完全に修復することをしていないことから損傷の度合いによっては塩基配列の大きな喪失となる。このような喪失を変異の一つである欠損 (Deletion) と呼ぶ。

HR、NHEJ どちらも欠損が生じた箇所に遺伝子が含まれていた場合遺伝子の塩基配列が変わることで遺伝子喪失の危険性がある。この遺伝子喪失ががんを抑制するような遺伝子で生じた場合、がんになるリスクが上昇してしまうことになる。このような点から DSB は様々な病気を引き起こす大変危険な DNA の病と言える。しかし、このような危険性に反して、DSB に関する研究はあまり芳しくない。特に機械学習を用いた解析は、多くのデータを必要とする性質上稀である。生物データは倫理、金銭、時間等様々な理由によって多くのデータを集めることが難しいためである。

2. 研究目的

本研究は、DSB と DSB を修復する経路の関係性を調べる先行研究[1]のデータに着目し、ヒトの常染色体 1-22 を平均化して解析を行っているところを常染色体毎に個別に解析することで、データとして取られている DNA についてのタンパク質と DSB の関係性を遺伝子レベルで調べて DSB の起こりやすい遺伝子の位置を推定することを目的としている。

解析手法としてはテンソル分解を用いた教師なし学習による変数選択法[2]をヒトの DNA についてのタンパク質のデータと DNA に生じた DSB を測定したデータに使い、DSB の起こりやすい遺伝子の位置を推定する。これにより、データが少ない環境下で行われる生物学に対する機械学習のモデルケースの一つという成果だけでなく、DSB が起こりやすい遺伝子のリストが得られる。DSB の修復で変異が生じやすいことがわかっているので、変異が起こりやすい遺伝子がわかったことになる。この結果から、がんなどの遺伝子変異によって生じる病気の原因である遺伝子を上述のリストから新たに見つけることができるかもしれない。また、未来の技術であるが、遺伝子に対する治療が確立されたときに変異が起こりやすい遺伝子をリストによって絞り込むことで、約 2 万個の

中から探すよりも効率的に治療が行えると考えられる。

3. 研究方法

本研究では、テンソル分解にHOSVD(Higher Order Singular Value Decomposition)を用いる。これは日本語では高次元特異値分解と訳すことができ、特異値分解が行列つまり二次元を扱っていることから三次元以上の場合の特異値分解と考えられる。行列の表記を $N \times M$ 行列 X を $X(N, M)$ と表記することになると、特異値分解は以下のように表せる(式(1))。ただし、 X は $X(N, M)$ ($N > M$)、 U は $U(N, M)$ 、 V は $V(M, M)$ 、 Σ は $\Sigma(M, M)$ で U 、 V は直交行列、 Σ は対角行列。

$$X = U\Sigma V^T \quad (1)$$

U もしくは V と Σ は、両辺に転置した X を掛けることで固有ベクトル固有値として求まる。また、求めていない V もしくは U については固有ベクトルを使い表すことができる。もし固有ベクトルの数が N もしくは M より少なければ、元の行列 X を低次元で書き表すことができる。

テンソルは unfolding という手法でもって行列に直すことができる。HOSVD はこれを利用してテンソル分解をする。三次元テンソル X (成分は X_{ijk} ($1 \leq i \leq 2, 1 \leq j \leq 2, 1 \leq k \leq 2$)) を例にして unfolding の具体例(2)、(3)、(4)を示す。

表 1: テンソル X

k=1	j	
i	x_{111}	x_{121}
	x_{211}	x_{221}
k=2	j	
i	x_{112}	x_{122}
	x_{212}	x_{222}

$$X^{i \times (jk)} = \begin{pmatrix} x_{111} & x_{121} & x_{112} & x_{122} \\ x_{211} & x_{221} & x_{212} & x_{222} \end{pmatrix} \quad (2) \quad X^{j \times (ik)} = \begin{pmatrix} x_{111} & x_{211} & x_{112} & x_{212} \\ x_{121} & x_{221} & x_{122} & x_{222} \end{pmatrix} \quad (3) \quad X^{k \times (ij)} = \begin{pmatrix} x_{111} & x_{211} & x_{121} & x_{221} \\ x_{112} & x_{212} & x_{122} & x_{222} \end{pmatrix} \quad (4)$$

Unfolding で得られた行列に SVD をして、得られた固有ベクトルを $U^{(i)}$ 、 $U^{(j)}$ 、 $U^{(k)}$ とすると特異値分解のように固有値を掛けることで元のテンソル X に戻って欲しくなる。しかし、別々の固有ベクトルを掛けても元には戻らない。そこで、これらの行列が直交行列なので、転置行列が逆行列となり容易に逆行列が求まることから、元のテンソル X に $(U^{(i)} \times U^{(j)} \times U^{(k)})^{-1}$ をかけることで、固有ベクトルの積と掛けることで元のテンソル X に戻るようなテンソル G を求める。このテンソル G はコアテンソルと呼ばれていて、SVD によって得られた個々の U が元のテンソルに戻るための重みの役割を持つ。元のテンソルの成分を x_{ijk} とし各 SVD の固有ベクトルの数を N 、 M 、 K 、固有ベクトルの成分を $u_{l_1 i}$ 、 $u_{l_2 j}$ 、 $u_{l_3 k}$ とすると HOSVD は以下のように表せる。

$$x_{ijk} = \sum_{l_1=1}^N \sum_{l_2=1}^M \sum_{l_3=1}^K G(l_1, l_2, l_3) u_{l_1 i} u_{l_2 j} u_{l_3 k} \quad (5)$$

研究に用いるデータは [上述した論文\[1\]](#) で得られたデータを使っている。このデータはヒトの常染色体 1-22 の全塩基対について発生した DSB や付着したタンパク質を記録したものとなっている。このデータは DSB を故意に起こさせる薬品を加えたものとそうでないものの二種類ある。これらのデータからタンパク質の数 DSB の数を要素に持ち、DSB とタンパク質群、塩基対の番号、薬品の有無の 3 つの次元を持つテンソルを作った。ただし、塩基対の次元については常染色体 1-22 全てそのままでは非常に大きく解析ソフトで扱いにくいいため、この次元を分

けることで 155 個の 3 次元テンソルデータとなっている。これに HOSVD を使い、DSB が起こりやすい遺伝子を見つける。

SVD の話に戻ると、分解後は固有値と固有ベクトルのみで書き表すことができ、固有ベクトルが基底になることから固有値の絶対値が大きいかが元の行列を書き表す上で重要となる。というのも固有値の絶対値が大きいほど元の行列に戻す際、固有ベクトルの寄与が大きくなるためだ。ここで、行列をデータとして考えると、固有ベクトルはデータの特徴と考えられる。これを HOSVD に置き換えて考えると、HOSVD で得られる行列は各次元の特徴を表す固有ベクトルとなっているので、この固有ベクトルから DSB が起こりやすい遺伝子の特徴を見つければよい。

式 (5) の u_{l_1i} 、 u_{l_2j} 、 u_{l_3k} を DSB とタンパク質群、塩基対の番号、薬品の有無とすると、最終的に遺伝子が求めたいことからパラメータ l_2 にこの特徴をもたせたい。DSB を表す特徴は u_{l_1i} から i に DSB を代入し絶対値が大きいパラメータ l_1 を選ぶ。起こりやすいという特徴は、DSB が起こりやすい環境下では普段以上に DSB になるものを DSB が起こりやすいと考えて、薬品有りとなしのときで差が大きいパラメータ l_3 を選ぶ。選ばれたパラメータ l_1 、 l_3 をコアテンソル G に代入し絶対値が大きい l_2 を選ぶことでこの特徴を選ぶことができる。

このパラメータを塩基対の行列 u_{l_2j} に代入し、仮説検定をして DSB が起こりやすい塩基対を求める。帰無仮説はここで得られた u_{l_2j} が無差別に選ばれておりそれぞれの値が独立な正規分布に従うというもので、対立仮説は DSB が起こりやすいような塩基対を選ぶことができているとなる。この条件で仮説検定をし、5% 以下で帰無仮説を棄却したものを DSB が起こりやすい塩基対とする。

現状まだ塩基対を選べただけで遺伝子があるかどうかわかっていない。そこで、遺伝子は塩基対が連続していることから、仮説検定によって選ばれた塩基対の中で塩基対の番号が連続しているものをここでは遺伝子と考える。これによって DSB が起こりやすい遺伝子を解析で見つけることができる。

解析結果が実際に DSB を起こしやすい遺伝子について選んでいるか調べる。COSMIC[3] というがんに関するデータベースから欠損、挿入という DSB によって生じる変異が発生したことがある mRNA のデータを入手し、UCSC genome browser [4] というゲノムデータベースからヒトの mRNA がどの塩基対にあるかマッピングされたデータを入手する。mRNA を使う理由については、遺伝子に比べて量が多く解析データが多いに越したことはない点、タンパク質を作るのに遺伝子全てを使うわけではなく遺伝子の一部を写した mRNA が使われていることから、mRNA の変異を見るほうがどの塩基配列に変異が生じたか正確にわかる点にある。解析データを含めたこの三種類のデータからフィッシャーの正確確率検定を行い、DSB を起こしやすい遺伝子が選んでいるかどうか調べる。解析データの塩基対の範囲にある mRNA の数を a 、その範囲の mRNA の中で変異を起こしたことがあるものの総数を b 、 a の mRNA の塩基配列内に解析で得られた u_{l_2j} が含まれるような mRNA の総数を c 、 c の中で変異を起こしたことがある mRNA に限定して u_{l_2j} が含まれているものを d とする。この 4 つのパラメータからフィッシャーの正確確率検定で p 値を与える。

$$p = \frac{(a-b)!(b)!(a-c)!(c)!}{(a)!(a-b-c+d)!(c-d)!(b-d)!(d)!} \quad (6)$$

表 2: フィッシャーの正確確率検定

	解析データと一致しない	解析データと一致する	全データ
変異なし mRNA 数	$a-b-c+d$	$c-d$	$a-b$
変異あり mRNA 数	$b-d$	d	b
合計	$a-c$	c	a

4. 研究結果

塩基対について分割したため解析データは 155 個ある。この 155 個のうちフィッシャーの正確確率検定で 24 個のデータが棄却され残った。5%以下で棄却したため、155 回の施行で 24 回も 5%以下の施行が起きたこととなる。これは偶然ではありえないため、テンソル分解を用いた教師なし学習法によって、DSB が起こりやすい遺伝子を推測することができたと言える。

5. 結論

我が国において、1981 年以來日本寺の死因ランキング一位は悪性腫瘍（がん）であり、国外を見ても特に先進国では悪性腫瘍による死因は多い。そのためがんの治療法を確立することは、とても重要なことである。そんな治療法を確立するためには、その病気について詳しい知識を得ることが必要不可欠である。本研究で注目した DSB は、悪性腫瘍の原因である遺伝子変異を起こす大本の一つと言える現象でその DSB が起こりやすい遺伝子を推測することが、故意に起こされた薬品上でできたことから、実際のヒトにおいてもできると考えている。しかし、本研究では mRNA を作る遺伝子に関しての推測という遺伝子の中での制限を設けて推測を行った。そのため、更に拡張し、ノンコーディング RNA についても解析し遺伝子全てに対して推測することが今後の課題だと考えている。

6. 参考文献

[1] Clouaire T, Rocher V, Lashgari A, Arnould C, Aguirrebengoa M, Biernacka A, et al. Comprehensive mapping of histone modifications at dna double-strand breaks deciphers repair pathway chromatin signatures. *Molecular cell* 72 (2018) 250-262.

[2] Y-hTaguchi. Tensor decomposition-based and principal-component-analysis-based unsupervised feature extraction applied to the gene expression and methylation profiles in the brains of social insects with multiple castes. *BMC bioinformatics* 19 (2018) 99.

[3] <https://cancer.sanger.ac.uk/cosmic/download>

[4] <https://genome.ucsc.edu/cgi-bin/hgTables>