

遺伝子発現情報の機械学習に基づく投薬効果予測

Prediction of drug therapy efficacy by support vector machine based on the personal gene expression patterns

吉田 寛輝[†] 峯田 克彦[†] 遠藤 俊徳[†]

Hiroki Yoshida Katsuhiko Mineta Toshinori Endo

1. 背景と目的

in vitro での薬剤有効率と *in vivo* での薬剤有効率とでは、その環境の違いにより大きく違い、*in vitro* では有効率が低いものでも *in vivo* では有効率が低くなる薬剤は多い。また、使用薬剤の決定も経験的判断による部分が大きく、時として無駄な投薬や重篤な副作用をもたらす可能性がある[1]。投薬前に対象患者に対してその薬剤の効果を予測し、適切な薬剤を選択することが出来れば、患者の生活の質の向上に繋がる。また、マイクロアレイによる臨床検体を対象とした解析は多くの解析方法が示されているが、未だに暗中模索の感がある[2]。本研究では、マイクロアレイを利用した遺伝子発現データからの遺伝子の絞り込み選択方法と、最適な学習判別器の構成方法を研究することにより、実際に臨床に応用できる投薬効果の予測判別方法の確立、及び、自己免疫疾患における薬効メカニズムを解明することを目的とした。また、本研究においては、精度が高いことで注目されているサポートベクターマシン(SVM)を学習判別機械とする投薬効果予測方法の確立を行った。

2. 材料と方法

匿名化した自己免疫疾患患者における、末梢血検体についてのマイクロアレイを用いた遺伝子発現量、及び、医師による各患者に対する投薬効果判定結果を用いた。患者数は73人であり、遺伝子数は789個である。そのうちの66人は学習データとして、7人は判別データとして用いた。投薬効果判定は、アメリカリウマチ学会が定めている投薬効果判定基準のACR[3]を基に判定されている。

学習データに用いた患者66人のうち、ACR70もしくはACR50と判定された患者32人を投薬効果ありの群として、ACR20もしくはACR0と判定された患者34人を投薬効果なしの群として、2群に分けた。

これらのデータをもとに *t* 検定を用いて、遺伝子の絞り込み選択を行った。多重検定の場合、偽陽性である遺伝子が採用されてしまう恐れがあるため、偽陽性と考えられる最大遺伝子数を *p*-value の低い順に除いた。

選択された遺伝子について、統計ソフト R(ver.2.7.2)で SVM のパラメータ最適化を行い、学習データを用いて学習判別の精度とサポートベクトル数を調べた。最適化は R の tune 関数を用い、精度評価は leave one out 法を用いた。用いたカーネルは Gaussian, Linear, Polynomial, Sigmoid の4つであり、用いたタイプは C-classification, Nu-classification の2つである。用いたカーネルの式は以下の通りである。

$$\text{Gaussian: } d(x,y) = \exp(-\text{gamma} * |x-y|^2)$$

Linear: マッピングは行わず、元の特徴空間で線形分離

$$\text{Polynomial: } d(x,y) = (\text{gamma} * (x * y) + \text{coef0})^{\text{degree}}$$

$$\text{Sigmoid: } d(x,y) = \tanh(\text{gamma} * (x * y) + \text{coef0})$$

続いて、学習データには用いられていない7人(投薬効果ありの群は5人、投薬効果なしの群は2人)の患者を判別データとして用い、学習データを用いて構築したモデルの、未知データに対する判別能力(汎化能力)を調べた。

次に、偽陽性を考慮して選択された遺伝子の機能を調べ、自己免疫疾患における投薬効果の有無を決定づける遺伝子の特徴を明らかにした。

最後に、薬効メカニズムを調べるために、投薬効果の見られた患者32人の投薬前後の遺伝子発現量の変化を調べた。まず、*t* 検定を用いて投薬前後で有意に変化した遺伝子を調べた。今回も多重検定であるので、FDR をコントロールする方法[4]により、偽陽性を考慮した。この方法は、新たな閾値 α を求めることで、 $p\text{-value} < \alpha$ となった時、統計的に有意とする方法である。 α を求める式は以下の通りである。

$$\alpha = (0.05 * i) / \text{全遺伝子数} \quad (i: p \text{ 値が低い方からの順位})$$

3. 結果と考察

3.1 *t* 検定の結果

t 検定の結果、 $p\text{-value} < 0.05$ で選択された遺伝子は49個であった。全部で789個の遺伝子を調べたので、今回、偽陽性と考えられる遺伝子の最大数は、 $789 * 0.05 = 39.45$ となり、*p*-value の低い順に40個の遺伝子を除いた。その結果、残った遺伝子は以下の9遺伝子である(表1)。

表1 選択された遺伝子

遺伝子名	機能
Asialoglycoprotein receptor 1*	損傷したグリコタンパクを除去(エンドサイトーシス)
Nonmetastatic protein 23*	造血幹細胞で、分化を抑制しアポトーシスを促進
Transforming growth factor, beta 1*	組織修復、炎症反応や免疫反応を終了させる調節因子
DNA fragmentation factor 40K, beta*	アポトーシスの最期の段階における、DNAの分解
Pyridoxal (pyridoxine, vitamin B6) kinase	ビタミン B6 が PLP になる際のリン酸化を促進
CX3C chemokine receptor 1*	血管内から炎症組織への白血球の遊走
t-complex-associated-testis-expressed 1-like 1	人の dynein のサブユニット
DNA fragmentation factor 45K, beta*	アポトーシスの最期の段階における、DNAの分解
ERK activator kinase	ERK を活性化させる

[†]北海道大学大学院 情報科学研究科 Graduate School of Information Science and Technology, Hokkaido University

3.2 SVMの結果

選択された9遺伝子を用い、1患者ごとに9次元データのベクトルとして、全66ベクトルが正しく判別されたかをleave one out法を用いて精度評価した。SVMの学習判別の結果、C-classificationタイプとNu-classificationタイプにおける、LinearカーネルとGaussianカーネルのそれぞれの組み合わせ(計4種類の組み合わせ)において、100%の精度を示した。また、C-classificationタイプとLinearカーネルの組み合わせによって最少のサポートベクトル数21を示した。計算量の観点から、サポートベクトル数は少ないほどよいモデルとされており、また、一般的なSVMモデルのサポートベクトル数は全ベクトル数の10%~30%である[5][6]。この結果より、C-classificationタイプとLinearカーネルの組み合わせのモデルのサポートベクトル数は、学習データの全66ベクトルの約30%であり、一般的なSVMモデルの範囲内にあることが分かった。

次に、学習データを用いて構築したモデルを用いて、未知データの判別能力(汎化能力)を調べた。結果は、7人中何人が正しく判別されたかを示している(表2)。

表2 汎化能力結果

	C-classification	Nu-classification
Gaussian	100%	100%
Linear	100%	100%
Polynomial	42.9%	85.7%
sigmoid	71.4%	N.D.

表2より、C-classificationタイプとNu-classificationタイプにおける、LinearカーネルとGaussianカーネルのそれぞれの組み合わせにおいて、100%の汎化能力を示した。よって、自己免疫疾患における投薬効果の予測に必要な十分な9遺伝子を同定できた。

3.3 9遺伝子の特徴

表1を見ると、9遺伝子中6遺伝子がアポトーシスやエンドサイトーシスや組織修復、免疫に関する遺伝子であることが分かる(*印のついた遺伝子)。よって、自己免疫疾患における投薬効果の有無を決定づける遺伝子の特徴は、アポトーシスやエンドサイトーシスによって不要な細胞を除いて事前に炎症を防ぐか、炎症した組織を修復する遺伝子が働いているかで決定していると示唆される。

3.4 薬効メカニズム

t検定のみ用いた時、 p -value<0.05で有意に変化した遺伝子は127個であり、投薬後に発現量が有意に低下したのは127遺伝子中わずかに3遺伝子であった。また、FDRをコントロールする方法を用いて遺伝子選択を行った結果、22遺伝子が有意に変化した。全22遺伝子において投薬後

表3 ribosomal protein以外の遺伝子

遺伝子名
T-cell cyclophilin
HLA class II DP beta 1
HLA class II DQ beta
immunoglobulin kappa
stathmin
HLA class II DM alpha
Liver-specific cytochrome c oxidase (COX VIIa-L)
Immunoglobulin rearranged gamma chain

に発現量が上昇した。22遺伝子のうち、14遺伝子はribosomal proteinであった。残りの8遺伝子は表3の通りである。

このことから、ribosomal proteinの発現量上昇により、これらのタンパク合成を促進していると考えられる。

表3の遺伝子から、IgやHLA class IIなど、免疫に関わる遺伝子が多くみられる。自己免疫疾患は自己抗体の過剰産生や、自己組織に反応するT細胞発現が原因であり、投薬後に効果が見られた患者の免疫遺伝子の発現が上昇したことは興味深い。これらから、投与された低分子化合物がハブテンとして働き、Ig産生を促し、それによって間接的にT細胞活性化を誘導したと考えられる。ここで活性化されたT細胞は、免疫反応の対象を自己組織から別の対象に移行しているのかもしれない[7]。

また、COX VIIa-L遺伝子は、年齢を重ねるにつれてプロモータ領域のメチル化が進み、遺伝子発現量が低下することで、2型糖尿病に関与していることが報告されている[8]。2型糖尿病と自己免疫疾患は、年齢を重ねるにつれて発病しやすいという表面的類似性がある。このことより、COX VIIa-L遺伝子の抑制は自己免疫疾患にも関係しているかもしれない。

4. 結論

本研究により、自己免疫疾患における投薬効果の有無を判別するのに必要十分な9遺伝子を同定することが出来た(表1)。

SVMのパラメータ最適化を行い、C-classificationタイプとlinearカーネルの組み合わせを使い、表1に示す9遺伝子を用いることで、患者負担の少ない、わずかな末梢血液細胞の遺伝子発現量検出による効果的かつ簡便な投薬効果予測を行うことが出来た。

自己免疫疾患における薬効メカニズムは、薬がハブテンとなってIg産生を促し、間接的にT細胞活性化を促進することでと考えられる。ribosomal proteinと免疫関連分子の発現上昇は、このことを示唆している。

また、2型糖尿病に関係しているCOX VIIa-L遺伝子が、自己免疫疾患に関わっている可能性があることが分かった。

参考文献

- [1]長谷川清ら; マイクロアレイデータの多変量解析による抗癌剤感受性規定因子の選択および未知検体の感受性予測; 中外製薬株式会社, 2007.
- [2]佐藤陽美; 分子ネットワークとチップデータ解析の融合~ KeyMolentの開発~; バイオテクノロジージャーナル, 2005.
- [3]Frank C. Arnett *et al*; The american rheumatism association 1987 revised criteria for the classification of rheumatoid arthritis, 1988.
- [4]Koen J.F. Verhoeven *et al*; Implementing false discovery rate control: increasing your power, 2005.
- [5]Davide Anguita Ridella *et al*; An Algorithm for Reducing the Number of Support Vectors, 2005.
- [6]Quang-Anh Train Zhang *et al*; Reduce the Number of Support Vectors by Using Clustering Techniques, 2003.
- [7]Shunichi Kumagai *et al*; Immune Responses to Hapten-Modified Self and Their Regulation in Normal Individuals and Patients with SLE, 1981.
- [8]Rönn T *et al*; Age influence DNA methylation and gene expression of COX7A1 in human skeletal muscle, 2008.