

## 特集 2 精神科領域の個別化医療の実現可能性とバイオバンクの活用

## 2. 精神科領域の個別化医療の実現可能性とバイオバンクの活用

高山 順<sup>1,2)</sup> 田宮 元<sup>2,3)</sup>

**抄録：**精神疾患をはじめとした複雑な疾患に対し、その遺伝要因はゲノム科学的手法により、環境要因は疫学的手法により、それぞれ探求されてきた。しかし、ゲノム科学は「失われた遺伝率」問題により、また疫学は小効果のリスク因子の再現性の危機により、それぞれ限界を迎えつつある。これらの限界を乗り越えると期待されるのが、前向きゲノムコホートとよばれる研究デザインである。前向きゲノムコホートでは定義された集団に対し、ゲノム情報と環境曝露情報を収集し、血液や尿などの生体サンプルを一定の品質管理のもと保管する（バイオバンキング）。また疾患発症情報や、検査値・画像・アンケート結果などの多様で多層的な中間表現型情報を前向きに取得していく。これにより、遺伝子・環境相互作用も含めた解析が可能となり、新たなリスク因子が同定されると期待される。しかしゲノムコホート研究には特有の困難が存在する。一つは  $p \gg n$  問題であり、もう一つは多様で多層的な中間表現型情報から意味のある特徴量を抽出する問題である。これらの問題を解決すると期待されるのが統計的機械学習および深層学習技術である。本稿ではこれらの技術を適用した筆者らの研究例を紹介する。

日本生物学的精神医学会誌 32 (2) : 81-84, 2021

**Key words :** prospective genome cohort, biobank, common complex disease, machine learning, artificial intelligence

精神疾患をはじめとした遺伝要因・環境要因が複雑に寄与する有病率の高い疾患は、ありふれた複雑な疾患 (common complex diseases) とよばれる。ありふれた複雑な疾患の遺伝要因はゲノム科学的手法により、一方の環境要因は疫学的手法により、それぞれ探求されてきた。ありふれた複雑な疾患に対するゲノム科学の代表的な解析手法は、一塩基多型 (single nucleotide polymorphism : SNP) を用いたゲノムワイド関連解析 (genome-wide association study : GWAS) である。SNP-GWAS 法においては罹患者集団と対照集団を用意し、罹患者集団に有意に濃縮された遺伝的多型をリスクバリエーションとして同定する。もしくは一つの集団において量的形質を目的変数とし、各 SNP 座位のマイナーアレルコピー

数 (0, 1, 2) を説明変数として線形回帰を行い、有意な座位を同定する。しかし SNP-GWAS 法がこれまでに同定してきたリスクバリエーションは一般に、たとえ有意であっても、その効果量は非常に小さい。そのため同定されたすべてのリスクバリエーションの効果量を合算しても、古典的な家系分析で推定された遺伝要因の割合、すなわち遺伝率の大半を説明できないという「失われた遺伝率」問題が生じる<sup>4)</sup>。この「失われた遺伝率」問題は長らくゲノム科学研究者を悩ませてきた。レアバリエーションの寄与、見逃されてきた構造多型の寄与、遺伝子間相互作用や遺伝子・環境相互作用の寄与から、家系分析で得られた遺伝率推定値が誤っていた可能性などさまざまな説が提出されている<sup>2)</sup>が、決着はついていない。一方

## Personalized Medicine in Psychiatry using Biobanks

1) 東北大学未来型医療創成センター (〒980-8573 宮城県仙台市青葉区星陵町 2-1 東北メディカル・メガバンク棟 4 階) Jun Takayama : Tohoku University Advanced Research Center for Innovations in Next-Generation Medicine. Tohoku Medical Megabank Building, 4th floor, 2-1 Seiryō-machi, Aoba-ku, Sendai, Miyagi 980-8573, Japan

2) 理化学研究所革新知能統合研究センター (〒103-0027 東京都千代田区日本橋 1-4-1 日本橋一丁目三井ビルディング 15 階) Jun Takayama, Gen Tamiya : RIKEN Center for Advanced Intelligence Project. Nihonbashi 1-chome Mitsui Building, 15th floor, 1-4-1 Nihonbashi, Chuo-ku, Tokyo 103-0027, Japan

3) 東北大学大学院医学系研究科 (〒980-8575 宮城県仙台市青葉区星陵町 2-1) Gen Tamiya : Tohoku University Graduate School of Medicine. 2-1 Seiryō-machi, Aoba-ku, Sendai, Miyagi 980-8575, Japan

【田宮 元 E-mail : gtamiya@med.tohoku.ac.jp】



れは統計学分野における研究が想定していた  $p$  と  $n$  のアンバランスの度合いをはるかに超えるものであり、ゲノム解析に適した新しい統計的手法の開発を必要としていた。

$p \gg n$  問題を解決するための手法として、筆者らのグループでは、超高次元変数選択法である sure independence screening (SIS) に基づくゲノム解析手法を開発した。SIS においては周辺回帰の後に罰則付き回帰を行うことで、最適なモデルを求める。筆者らは SIS 法を、超並列計算が可能な GPU を用いて計算できるよう実装し、Wellcome Trust Case-Control Consortium や Stevens-Johnson 症候群のデータに適用し新たなリスクバリエーションおよび遺伝子間相互作用を現実的な時間内に同定することに成功した<sup>7, 9)</sup>。

$p \gg n$  問題はリスクバリエーションの同定だけでなく、同定したリスクバリエーションからの発症リスク予測においても問題となる。発症リスク予測においては、予測に寄与するバリエーションを最適な数だけ選ぶ必要がある。筆者らはこれを実現する手法として、smooth-threshold multivariate genetic prediction (STMGP) 法を開発した<sup>8)</sup>。STMGP 法は機械学習手法に基づき、限定的な効果を示す SNP 座位に対して適切な重み付けを行いリスク予測に組み入れるものである。シミュレーションを用いた性能評価では、STMGP 法はしばしばゲノム予測に使われる genomic best linear unbiased prediction (GBLUP) 法よりも高い性能を示した。また、実際のアルツハイマー病のデータ (Alzheimer's disease neuroimaging initiative : ADNI) や、うつ症状のゲノムリスク予測に対してもより高い予測性能を示した<sup>5, 8)</sup>。

## 2. 多様で多層的な中間表現型からの要約問題

前向きゲノムコホート研究においては多様な研究計画を実行できるように、さまざまな表現型情報を取得する。例えば身長・体重・生年月のような基礎的な情報、血液生化学検査の検査値から、脳 MRI 画像やアンケート調査に至るまで多岐にわたる情報が収集される。これら表現型のデータは分子レベル、細胞レベル、臓器レベル、個体レベルとさまざまな階層から収集され、また複数の階層にまたがるものも多い。これら個別の値を目的変数とした解析もちろん可能であるが、生体のより本質的な特徴を捉えるためには、個別の値だけではなくデータ全体をみる必要がある。特に精神科領域の疾患のような場合にはその重要性は高い。これは医師が臨床の現場

でさまざまな検査値、画像解析結果や問診結果を解釈しながら、診断仮説を組み立てることに相当する。遺伝統計学の文脈においては、潜在変数として仮定される易罹病性 (liability) 変数を、個別の表現型変数から合成し、特定の域値を超えるかどうかをもって疾患の有無を判定することに相当する。このようなことが可能となれば、前向きゲノムコホートが収集した表現型データを真に有効活用することが可能となる。しかし多様で多層的な中間表現型データを有効に活用するための手法は確立していない。

統計学の分野では複数の変数から内在する構造を見だし、より低次元の変数を合成する研究が行われてきた。代表的な手法としては、主成分分析 (principal component analysis : PCA) が挙げられる。例えば筆者らは、PCA をイネの構造に関する形質値に適用して得た主成分に対する GWAS を行い遺伝子の同定に成功している<sup>11)</sup>。しかし、PCA は線形的手法であり、より複雑な構造を持った表現型データ、特に表現型間に非線形の構造がある場合にはより発展的な手法が適しているかもしれない。オートエンコーダは、深層学習の枠組みを用いた次元削減手法であり、深層学習研究ブームのきっかけとなったものである<sup>3)</sup>。オートエンコーダでは多次元の入力値をニューラルネットワークに入力し、同数の次元で出力する。この入力層と出力層をつなぐネットワークの中間層のなかに、ごく少数のニューロンからなる層が存在する。これは PCA における主成分に相当するとみなすことができ、データに内在する構造を捉えた低次元の表現として捉えることが可能である。例えば筆者らの研究ではオートエンコーダを用いて、前立腺がんの病理画像から新しい組織学的特徴を抽出することに成功している<sup>10)</sup>。これらの成功は前向きゲノムコホートにおける中間表現型からの適切な特徴量の抽出に深層学習の枠組みが有効である可能性を示唆している。さらに今後は、画像と検査値・アンケート結果などのように大きく種類の異なるデータ (マルチモーダルデータ) からの適切な特徴量の抽出手法の開発が重要となるだろう。

本論文に記載した筆者らの研究に関してすべて倫理的配慮を行っている。開示すべき利益相反は存在しない。

## 文 献

- 1) Collins FS (2004) The case for a US prospective cohort study of genes and environment. *Nature*, 429 (6990) : 475-477.

- 2) Eichler EE, Flint J, Gibson G, et al (2010) Missing heritability and strategies for finding the underlying causes of complex disease. *Nat Rev Genet*, 11 : 446-450.
- 3) Hinton GE and Salakhutdinov RR (2006) Reducing the dimensionality of data with neural networks. *Science*, 313 : 504-507.
- 4) Maher B (2008) Personal genomes : The case of the missing heritability. *Nature*, 456 : 18-21.
- 5) Takahashi Y, Ueki M, Tamiya G, et al (2020) Machine learning for effectively avoiding overfitting is a crucial strategy for the genetic prediction of polygenic psychiatric phenotypes. *Transl Psychiatry*, 10 : 294.
- 6) Taubes G (1995) Epidemiology faces its limits. *Science*, 269 (5221) : 164-169.
- 7) Ueki M and Tamiya G (2012) Ultrahigh-dimensional variable selection method for whole-genome gene-gene interaction analysis. *BMC Bioinform*, 13 : 72.
- 8) Ueki M and Tamiya G (2016) Smooth-threshold multivariate genetic prediction with unbiased model selection. *Genet Epidemiol*, 40 : 233-243.
- 9) Ueta M, Katsushi T, Sotozono C, et al (2012) Epistatic interaction between Toll-like receptor 3 (TLR3) and prostaglandin E receptor 3 (PTGER3) genes. *J Allergy Clin Immunol*, 129 : 1413-1416.
- 10) Yamamoto Y, Tsuzuki T, Akatsuka J, et al (2019) Automated acquisition of explainable knowledge from unannotated histopathology images. *Nat Commun*, 10 : 5642.
- 11) Yano K, Morinaka Y, Wang F, et al (2019) GWAS with principal component analysis identifies a gene comprehensively controlling rice architecture. *Proc Natl Acad Sci USA*, 116 : 21262-21267.

---

## ■ ABSTRACT

### Personalized Medicine in Psychiatry using Biobanks

Jun Takayama<sup>1,2)</sup>, Gen Tamiya<sup>2,3)</sup>

1) *Tohoku University Advanced Research Center for Innovations in Next-Generation Medicine*

2) *RIKEN Center for Advanced Intelligence Project*

3) *Tohoku University Graduate School of Medicine*

For complex diseases, including psychiatric disorders, genetic and environmental risk factors were investigated using genomic and epidemiological methods. However, genomics is faced by the so-called “missing heritability” problem, and epidemiology is limited by the reproducibility crisis for small effect risk factors. A research design called prospective genome cohort is expected to overcome these limitations. Prospective genomic cohorts collect genomic and environmental exposure information for a defined population and store biological samples such as blood and urine under a controlled quality. Besides, genome cohorts prospectively obtain disease onset information and various endophenotypes such as laboratory tests, brain imaging, and questionnaire survey. With the prospective genome cohort, new risk factors, including gene-environment interactions, will be identified. However, genome cohort studies have unique difficulties. One is the  $p \gg n$  problem, and the other is the problem in extracting meaningful features from diverse and multi-layered endophenotypic information. Here, we will present examples of statistical machine learning and deep learning techniques that are expected to address these problems.

---

(Japanese Journal of Biological Psychiatry 32 (2) : 81-84, 2021)