

2026年6月23日

報道関係者各位

慶應義塾大学医学部
東京科学大学
静岡大学

統合失調症の発話に共通する“言語の3要素”を自然言語処理で同定

— 500名超・1300データセットの世界最大級の会話コーパスを基に多層自然言語処理で検証 —

慶應義塾大学岸本泰士郎教授（医学部医科学研究連携推進センター）、東京科学大学高橋英彦教授（医歯学総合研究科精神行動医科学分野）、静岡大学狩野芳伸教授（グリーン科学技術研究所／情報学部）らを中心とする研究グループは、統合失調症患者の発話を自然言語処理（NLP：注1）により多階層的に解析し、自由会話において共通して観察される中核的な言語指標として、①格助詞（注2）の低下（「誰が何をしたか」といった文法的関係を示す手がかりの弱まり）、②語の意味的類似度の上昇（語の選択が意味的に近い範囲に偏る傾向）、③副詞（注3）の使用頻度の低下、の3要素を同定しました。これら3指標を組み合わせたモデルは、独立検証においてAUC（注4）=0.87（95%CI 0.74-0.97）と高い判別性能を示しました。

統合失調症における言語障害は、形態・統語・意味・談話といった複数の階層にまたがることが知られていますが、従来研究では単一レベルの指標や単一課題に依存した評価が多く、課題を超えて再現される「核」となる特徴は明確ではありませんでした。

本研究グループは、約10年間にわたり500名を超える精神疾患を患う患者から自然会話および課題会話データを収集し、30～60分の会話記録からなる約1,300件にのぼる世界最大級のデータセットを蓄積してきました。本研究では、このうち統合失調症患者104名および健常対照者101名の半構造化面接データを用い、発話テキストを対象に、形態・統語・意味・談話にまたがる計76の言語特徴量を抽出しました。その後、因子分析を行うことで代表的特徴量を選出し、診断との関連を検証した結果、上記の3指標が最も強く統合失調症の罹患と関連することが明らかとなりました。特に、格助詞および副詞の使用頻度の低下は再現性高く確認され、統合失調症の言語特徴として頑健なものであることが示されました。

本成果は、統合失調症の言語障害を「形態統語的明示性の低下」「意味空間（注5）の狭まり」「修飾による文脈調整の低下」という3要素として整理する新たな枠組みを提示するものであり、将来的な客観的モニタリング指標の基盤となることが期待されます（※本研究は個人の診断を単独で行うものではありません）。

本研究成果は、6月23日（日本時間）に国際学術誌 *Psychological Medicine* に掲載されました。

1. 本研究のポイント

- 精神疾患の会話コーパスを蓄積・今までに500名超・1300データセットを蓄積
- 日本語の会話・物語・絵の描写の3課題の会話文をNLPで解析し、統合失調症に関連する言語特徴を多層（形態・統語・意味・談話）で整理
- 自由会話で「格助詞の減少」「副詞の減少」「語の意味的近さの増加」が顕著で、判別性能はAUC=0.87（95% CI 0.74–0.97）
- 格助詞と副詞は複数課題で再現性（クロスタスク）が示され、会話の“関係づけ”や“文脈づけ”の弱まりを反映する可能性
- 将来、統合失調症の症状の客観的モニタリング指標の基盤となることが期待される

2. 研究の背景

統合失調症の発話は、連合の弛緩や脱線、内容の貧困化など、思考形式の障害として古くから知られています。今まではある意味、エキスパートの精神科医が名人芸として見抜く/気づくことができる評価でした。ただ、近年は自然言語処理の進展により、統語的複雑性や意味的一貫性など、多様な言語指標が提案され、言語障害は多次元的であることが示唆されています。

それにもかかわらず、従来研究の多くは小規模かつ単一課題に依存しており、(1) 言語階層間の統合的变化、(2) 課題による変動、(3) 課題を超えて再現される中核的特徴が十分に整理されていませんでした。そこで本研究では、言語を階層構造として捉え、多数の特徴量を網羅的に抽出したうえで、因子分析により解釈可能な代表指標を選定し、「核となる言語マーカ」の同定を目指しました。

3. 研究の内容・成果

岸本泰士郎教授らの研究グループは、言葉に現れる精神疾患の特徴を明らかにし、新たな診断や治療につなげることを目的として、約 10 年間にわたり精神疾患患者 500 名以上から自然会話および課題会話データを収集してきました。これまでに、1 回あたり 30~60 分の会話記録からなる約 1,300 件にのぼる世界最大級のデータセットを蓄積しており（プロジェクト名：UNDERPIN）、言語データに基づく精神医学研究の基盤を構築しています。

本研究では、この UNDERPIN データの中から、統合失調症患者 104 名および健常対照者 101 名の半構造化面接データを用いました。解析対象は発話テキストであり、①自由会話、②物語を自分の言葉で再構成する課題、③絵を見て内容を説明する課題、という 3 種類の言語課題から構成されています。

研究では、各課題・各時点のデータに対して、形態・統語・意味・談話といった複数の言語階層を捉える計 76 種類の自然言語処理（NLP）特徴量を抽出しました。形態・統語情報は日本語解析器 GiNZA を用いた品詞タグ付けおよび係り受け解析から取得し、意味的特徴は Word2Vec や TF-IDF により語の意味的近さや語彙のネットワーク構造を定量化しました。さらに、SentenceBERT を用いて文と文の意味的類似度を算出し、会話全体の一貫性（談話構造）を評価しました。

これら多数の特徴量の中から本質的な要素を抽出するため、特徴量間の相関や冗長性を整理しつつ、診断情報に依存しない形で探索的因子分析を実施しました。その上で、各因子を

代表する解釈可能な指標（representative feature）を選定しました。

その結果、自由会話においては、①格助詞の使用頻度の低下、②語同士の意味的な類似性の上昇、③副詞の使用頻度の低下の3つが統合失調症と有意に関連することが明らかになりました。これらを組み合わせたモデルは、AUC=0.87（95%信頼区間 0.74-0.97）と高い識別性能を示しました。

さらに、これらの特徴が異なる課題でも一貫して現れるかを検証するため、課題横断性をpartial conjunction（PC）により評価したところ、格助詞および副詞の低下は「3課題中少なくとも2課題で有意」という基準を満たし、課題に依存しない頑健な指標であることが示されました。

以上の結果から、統合失調症における言語障害は、(1)「誰が何をしたか」といった文法的関係を示す手がかりの弱まり、(2)語の選択が意味的に近い範囲に偏る傾向、(3)副詞などによる細かなニュアンスの調整機能の低下、という相補的な3つの側面から捉えられることが示されました。

4. 社会的意義・今後の展望

本研究成果の社会的意義は、統合失調症患者にみられる言語の変化を、単一の指標や特定の課題に依存した所見としてではなく、文法・意味・会話の流れといった複数の層にまたがる現象として統合的に捉え、さらに解釈可能な少数の代表指標へ整理して示した点にあります。自由会話において、格助詞の使用低下、副詞の使用低下、語の意味的類似性の上昇という3つのコア指標が同定され、独立データにおいても高い判別性能を示したことは、臨床現場で経験的に捉えられてきた「話し言葉の変化」を、再現性のある定量指標として扱う基盤となるものです。

また、本研究では、自由会話・物語課題・絵の描写課題といった異なる場面により言語特徴の現れ方が変化することを踏まえ、課題をまたいで一貫して観察される指標（クロスタスク・マーカー）と、特定の課題でより鋭敏に現れる指標を区別して示しました。これは、自然言語処理（NLP）を用いた言語評価の標準化に向けた重要な知見といえます。

さらに本研究は、統合失調症における言語障害を、(1)文法的な関係を明示する手がかりの弱まり、(2)語の選択が意味的に近い範囲に偏る傾向、(3)副詞などによる細かなニュアンスの調整機能の低下、という3つの側面から捉える視点を提示しました。これにより、従来ひとまとめに扱われがちであった形式思考障害（FTD）を、言語産出のどの機能に偏りがあるのかという観点から、データに基づいてより細かく理解する方向性が示されます。将来的には、こうした言語プロファイルに基づく患者層別化や、研究・臨床における客観的なモニタリング指標の開発につながることを期待されます。なお、これらの指標は臨床面接に取って代わるものではなく、あくまで補助的に活用されるべきものです。

今後の展望として、第一に、日本語で得られたこれらの指標が、言語構造の異なる他言語においても再現されるかを検証し、普遍的な特徴と言語固有の特徴を切り分けていきます。第二に、音声情報（話す速さや間の取り方、抑揚）や表情・視線、さらには臨床・認知指標などを統合した多モーダル解析へと拡張し、異なる情報源を組み合わせても安定して検出できる特徴の同定を目指します。第三に、診断の補助にとどまらず、発症リスクの予測、治療反応や再発の予測、社会機能の回復度の予測といった臨床的な予測課題へと展開し、実装を見据えた検証を進めていきます。

5. 注意点

本研究の結果は集団平均との差に基づくもので、個人の診断を単独で行うものではありません。今後、他言語・他集団での再現や機能転帰との関連検証が必要です。

6. 特記事項

本研究は、国立研究開発法人科学技術振興機構（JST）戦略的創造研究推進事業（CREST）自然言語処理による心の病の理解：未病で精神疾患を防ぐ（JPMJCR1684）、戦略的創造研究推進事業（CREST）精神医学×メディア解析技術による心の病の定量化・早期発見と社会サービスの創出（JPMJCR19F4）、戦略的創造研究推進事業（AIP 加速研究）精神医学×メディア解析技術の展開：精神疾患への介入の挑戦（JPMJCR22U4）の支援によって行われました。

7. 論文

英文タイトル：An Integrative NLP Framework Identifies Multi-Level Linguistic Phenotypes of Schizophrenia Across Tasks

タイトル和訳：統合的自然言語処理フレームワークを用いた、課題横断的な統合失調症の多層的言語表現型の同定

掲載誌：*Psychological Medicine*

著者：中村 啓信、狩野 芳伸、杉原 玄一、竹村 亮、山口 湧声、清水 正彬、高木 俊輔、飯塚 真理、田代 紗彩、北沢 桃子、仙頭 綾子、高橋 英彦、岸本 泰士郎

DOI：<https://doi.org/10.1017/S0033291726104668>

【用語解説】

(注 1) 自然言語処理（NLP）：文章や会話をコンピュータで解析する技術

(注 2) 格助詞：「が／を／に」など、文中の役割関係を示す語

(注 3) 副詞：「とても／ゆっくり／たぶん」など、言い方の程度・時制・文脈を修飾・調整する語

(注 4) AUC：判別性能の指標（1に近いほど高い）

(注 5) 意味空間：会話の言葉同士のつながりを「言葉をベクトル化する」技術を用いてネットワークとして表したもの

※ご取材の際には、事前に下記までご一報くださいますようお願い申し上げます。

※本リリースは文部科学記者会、科学記者会、厚生労働記者会、厚生日比谷クラブ、本町記者会、各社科学部等に送信しております。

【本発表資料のお問い合わせ先】

慶應義塾大学医学部 医科学研究連携推進センター

教授 岸本 泰士郎（きしもと たいしろう）

TEL : 03-5363-3219 E-mail : tkishimoto@keio.jp

<https://mhiis.pirms.med.keio.ac.jp/>

【本リリースの配信元】

慶應義塾大学信濃町キャンパス総務課：岡見・飯塚・奈良・加納

〒160-8582 東京都新宿区信濃町 35

TEL : 03-5363-3611 FAX : 03-5363-3612 E-mail : med-koho@adst.keio.ac.jp

<https://www.keio.ac.jp/ja/med/>

東京科学大学 総務企画部 広報課

〒152-8550 東京都目黒区大岡山 2-12-1

TEL : 03-5734-2975 FAX : 03-5734-3661 E-mail : media@adm.isct.ac.jp

静岡大学 総務部 広報・基金課 広報係

〒422-8529 静岡県静岡市駿河区大谷 836

TEL : 054-238-5179 FAX : 054-238-4450 E-mail : koho_all@adb.shizuoka.ac.jp

<https://www.shizuoka.ac.jp/>