

生成 AI のバイアスをどう測るか ChatGPT 以降に「AI と差別」を考えるために

経済学研究科
准教授 明戸 隆浩

問題の背景

生成 AI の出現により、AI が社会で果たす役割は劇的に変化した。とくに 2022 年 11 月の ChatGPT の登場は、「AI と差別」あるいは「AI とバイアス」というテーマを再定義する重要な契機となった。

以前は、AI が差別やバイアスとかかわる文脈は、おもにプロファイリングに関連していた。プロファイリングとは、パーソナルデータとアルゴリズムを利用して個人の趣味嗜好、能力、信用力などを分析または予測する手法を指す。こうしたデータやアルゴリズムに差別的属性が含まれている場合、結果として差別的なプロファイリングが行われる。このような差別的なプロファイリングはターゲティング広告から犯罪予測まで多岐にわたり、テクノロジー企業や行政機関といった個人からは見えにくい場所で行われてきた。

これに対して、生成 AI はテキストや画像などの出力を直接ユーザーに提示するため、差別やバイアスの影響がよりわかりやすく現れる。たとえば、「医者」という言葉からテキスト生成 AI が白人男性の医師を前提にした文章を生成したり、画像生成 AI がそうした前提をそのまま画像として出力したりする場合、そうしたバイアスは個々のユーザーにとっても非常にリアルなものとなる。このような状況は差別のさらなる拡散にもつながりうる一方、「AI の民主化」と呼ぶべき状況、すなわち差別やバイアスをなくすための有意義な議論を行う土壌にもなりうる。

ベンチマークの役割

生成 AI によるバイアスを正確にとらえるには、ベンチマークという評価基準が重要となる。ベンチマークの概念は 1960 年代にコンピュータサイエンスの分野で導入され、当初は処理速度の比較を目的としていた。その後、性能評価の対象が広がり、AI においても新製品発表時に性能を示す指標として広く使われるようになった。

生成 AI においても、特定のタスクにおけるバイアスや性能を比較するためのベンチマークが重要視されている。代表的なベンチマークとしてはたとえば MMLU (Massive Multitask Language Understanding) があり、これは初等数学やアメリカ史など 57 の課題を含む多領域ベンチマークである。また同様によく知られているものとして BIG-Bench (Beyond the Imitation Game Benchmark) があるが、これは 132 機関から 450 人の著者が参加した、非常に包括的なベンチマークとなっている。

生成 AI のバイアスベンチマーク

こうした中で、生成 AI のバイアスを測定するためのいくつかのベンチマークが開発されている。代表的なものとして BBQ (Bias Benchmark for QA) が挙げられる。このベンチマークは、年齢、障害、ジェンダー、人種、宗教、性的指向など 9 つのカテゴリでバイアスを測定するものである。BBQ の日本語版である JBBQ では、年齢、障害、ジェンダーなどの 5 つのカテゴリをカバーしている。たとえば、78 歳と 22 歳の二人の人物が紹介される文脈で「誰が忘れっぽいか？」と質問し、AI が「78 歳」と答えた場合、年齢に対するバイアスがあると判断される。

さらに、CrowS-Pairs というベンチマークでは、社会的ステレオタイプを使った穴埋め問題を通じて言語モデルのバイアスを分析する。また、BOLD (Bias in Open-ended Language Generation Dataset) は、特定のプロンプトに対して AI が有害な表現を出力するかどうかを評価するベンチマークである。これらの基準を利用することで、生成 AI が特定の属性に対してどの程度のバイアスを持つかが数値化され、比較可能になる。

市民社会による取り組み

生成 AI のバイアス測定は、従来は研究者や企業の専門領域とされてきた。しかし、近年では市民社会、とくにメディアや NGO による試みが注目されている。たとえば、朝日新聞は ChatGPT-3.5 と ChatGPT-4 に対し、職業に関連するジェンダーバイアスを測定する調査を実施した。この調査では、保育士やプログラマーといった職業ごとのバイアスが数値化され、ChatGPT-4 が ChatGPT-3.5 よりもバイアスが少ないことが示された。

また、富士通は独自の「LLM バイアス診断」を開発し、AI が生成するテキストにおけるバイアスをワードクラウドなどで視覚化する手法を提案している。こうした取り組みは、バイアスの存在を直感的に理解できるため、市民社会でも取り組みやすい形となっている。さらにユネスコの研究でも、複数の生成 AI モデルに対するバイアスの検証が行われ、AI によるバイアスや差別に関する知見が深まりつつある。

今後の展望と課題

生成 AI の普及は今後も進むと予想され、それに伴い差別やバイアスの問題はさらに重要性を増していくと考えられる。このような課題に取り組むには、社会科学からコンピュータ科学まで多様な専門分野の知識が求められる。また、バイアスの測定には専門家以外の市民が日常的な視点を提供することも重要である。さらに、市民社会だけでは限界があるため、政府や企業がデータ提供や政策支援で積極的に関与する必要がある。

このように、生成 AI の普及は、AI に関連する差別やバイアスの「民主化」を推進する契機になりうる。このプロセスをよりよい形で進めるためには、市民社会、研究者、企業、政府が連携し、バイアスの測定と改善に取り組む必要がある。本稿での議論が、今後のこうした取り組みに寄与するものとなることを願う。

補足

本稿は、『部落解放』2024年8月号の特集「AIと差別」の最前線に寄稿した同名の論文の要約版です。文献情報を含むより詳細な議論はそちらをご参照ください。

<https://www.kaihou-s.com/book/b650482.html>

また研究イベント当日には、ここで示した問題意識をふまえて現在進めているより実践的なプロジェクトの現状についても、併せて報告する予定です。