

【田中頼人特別研究】

第4回レポート

2301330039：安田直也



AI検索サービスにおける企業優遇バイアスと競争法への影響

中川 慧、平野 正徳、藤本 悠吾：大規模言語モデルを活用した金融センチメント分析における企業固有バイアスの評価、第21回テキストアナリティクス・シンポジウム、vol.124, no.173, NLC2024-15、pp.81-86（2024年9月3日）

企業固有のバイアスに特化した研究アプローチ：

多くの既存の研究は一般的なバイアスや社会的なバイアスに焦点を当てており、特定の企業や金融分野に特化したバイアス分析は限られていました。この論文は、企業名や企業属性が金融テキストの感情評価にどのように影響するかを初めて系統的に評価しており、企業特有のバイアスを評価するための新しい枠組みを提供しています。

企業の感情評価に関するバイアスと株価変動の関係の分析：

この研究では、バイアスが市場に与える潜在的影響を経済モデルで評価しています。バイアスのある感情評価が投資家の行動や株価変動に及ぼす影響を、理論的・実証的に分析するために、独自のモデルを構築しました。

例えば、ある企業に対して一貫してポジティブなバイアスが感情分析で表れると、その企業の株価が過剰に評価されるリスクがあります。この点を検証するため、感情スコアと株価変動のデータを用い、特定企業の評価バイアスが株式市場においてどのように作用するかを明示しています。これにより、LLMのバイアスが金融市場での意思決定に及ぼす影響を、経済モデルの観点からも評価している点で新しい貢献をしています。

実データに基づく日本市場の分析：

日本市場の実データを用いることで、地域特有の企業特性とバイアスの関連性を検証しています。MSCI Barraのモデルを活用し、規模やバリュエーション、成長率など20の企業特性とバイアスの相関を分析することで、LLMのバイアスが特定の市場や企業の属性に依存しているかを評価しています。

また、地域特有の金融市場の傾向や経済構造がバイアスにどう影響を与えるかを実証することにより、他の地域や業界にも適用可能な洞察を得られる可能性を示しています。これにより、金融業界におけるLLMの実用的な価値と、企業特有のバイアスの影響を地域別に評価する手法を提供している点で、新しい貢献をしています。

Kamruzzaman, M.; Nguyen, H. M.; Kim, G. L.: "Global is Good, Local is Bad?": Understanding Brand Bias in LLMs, arXiv preprint, arXiv:2406.13997, (2024)

ブランドバイアスに特化した調査:

本研究ではブランドに対するバイアスを初めて体系的に分析しています。具体的には、国や所得レベルごとに異なるブランド推奨の傾向を定量的に評価し、LLMが特定の経済圏に基づくブランドバイアスを持つことを示しました。

出身国効果※COO効果（Country-of-Origin Effect）の分析:

LLMが特定の国や地域に応じたブランド評価を行う傾向があることを示し、COO効果がどのようにLLMの応答に影響を与えるかを調査しました。これにより、国や地域ごとの文化的背景がLLMのブランド推奨にどう関わるかが明らかにされました。

高所得国と低所得国のユーザーに対するブランド推奨傾向の定量化:

本研究では、高所得国のユーザーには高級ブランドが、低所得国のユーザーには一般ブランドが推奨されやすいことが実験によって示されています。これは、LLMがデータセットから学習した社会経済的バイアスが応答に反映されることを示しており、こうした現象を数値的に評価した点で先行研究よりも詳細な分析を行っています。

公平性の向上に向けた具体的な提案:

この研究は、LLMが公平なブランド推奨を行うためのデータセットやトレーニング手法の改善を提案しており、実際の応用におけるモデルの公平性の重要性を指摘しています。

Quiñonero-Candela, J.; Wu, Y.; Hsu, B.; Jain, S.; Ramos, J.; Adams, J.; Hallman, R.; Basu, K.: Disentangling and Operationalizing AI Fairness at LinkedIn, Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency (FAccT '23), vol. 124, no. 173, pp. 81-86 (2023).

AIの公平性の明確な分離：

本論文は、AIの公平性を「等しいAI処理（Equal AI Treatment）」と「プロダクトの公平性戦略（Product Equity）」に明確に分離しています。これにより、AIの設計と製品全体の公平性戦略を個別に検討し、両者を組み合わせて総合的な公平性を実現するアプローチを提案しています。

実世界のケーススタディの提供：

LinkedInの具体的な製品での適用事例を通じて、提案する公平性フレームワークの実践的な適用方法とその効果を示しています。これにより、理論と実践の橋渡しを行い、他の企業や研究者が実際に適用可能な知見を提供しています。

正当化フレームワークの導入：

AIのバイアス緩和において、「正当化フレームワーク（justifiability framework）」を採用し、特定のグループに対する不平等が確認された場合、その原因を分析し、適切な対策を講じるプロセスを明確にしています。どうしても必要な場合のみ、特定のグループ情報を使って偏りを減らします。これにより、バイアスの原因究明と効果的な緩和策の実施が可能となります。

透明性：

透明性を持って自社の取り組みを共有しています。また、外部の専門家グループによるアルゴリズム監査も検討しており、公平性を定量的に評価するための基準を設定しようとしています。

Li: "A Survey on Fairness in Large Language Models", arXiv, 2308.10149, pp. 1-28 (2023)

概要：

大規模言語モデル（LLM）における「公平性」と「バイアス」に焦点を当て、これらのモデルがトレーニングデータを通じて社会的バイアスを学習し、その偏見が下流タスクにも影響を与える点について調査しています。特に、LLMのトレーニング方法やパラメータの大きさに応じて異なるバイアスの発生パターンとその解消方法を整理・分類しています。

トレーニングパラダイムとパラメータ規模の違いに基づく分類：先行研究と異なり、本研究ではLLMのトレーニング手法（事前学習+ファインチューニング vs プロンプティング）およびパラメータの規模（中規模 vs 大規模）に基づいてLLMのフェアネス調査を分類しています。これにより、異なるトレーニングパラダイムとモデルサイズがフェアネス研究のアプローチにどのような影響を与えるかについて、より詳細な分析が可能となっています。

構造と包括性：他のフェアネス調査と異なり、本研究はLLMをトレーニングパラダイムとパラメータの規模によって2つの主要なカテゴリに分け、それぞれを独立して調査しています。これにより、より明確な構造と包括的な分類が提供され、フェアネスに関する理解が深まるとされています。

Bi, Guanqun; Shen, Lei; Xie, Yuqiang; Cao, Yanan; Zhu, Tiangang; He, Xiaodong: "A Group Fairness Lens for Large Language Models", arXiv, arXiv:2312.15478, pp.1-14 (2023)

概要：

本研究はLLMの公平性向上に向けた新しい基盤を提供し、今後もLLM設計において包括的な偏見評価と緩和メカニズムが重要であると述べています。

詳細：

問題提起：従来のLLMの偏見評価は限られた属性やグループに焦点を当てており、包括的な評価が不足している。データセット構築：新しい「GFair」データセットを構築し、性別、年齢、国籍など10の異なるバイアス次元にわたる多様なグループを含めることで、偏見評価の幅を広げました。

評価手法：「ステートメントオーガナイゼーションタスク」というオープンエンドなテキスト生成タスクを導入し、モデルがどのように偏見を持つかをより深く理解する方法を提供。

実験結果：複数のLLM（GPT-4、Llama2など）に対し、毒性バイアス、センチメントバイアス、警戒バイアスなどの指標で偏見を評価しました。結果、モデル間で偏見のレベルに差があり、特にLlama2シリーズが比較的良好な結果を示しました。

偏見緩和手法：「GF-Think」と呼ばれる新しい手法で、Chain-of-Thought（CoT）推論を用いて公平性のある回答生成を行い、モデルの偏見を緩和できることを示しました。

Fang, X.; Che, S.; Mao, M.; Zhang, H.; Zhao, M.; Zhao, X.: Bias of AI-Generated Content: An Examination of News Produced by Large Language Models, Scientific Reports, vol. 14, pp. 1234-1245 (2024)

概要：

本論文は、LLMが生成するニュースコンテンツ（AI生成コンテンツ、AIGC）における性別および人種のバイアスが依然として顕著であることを指摘しています。従来のニュースソース（ニューヨーク・タイムズやロイター）を参照し、特にジェンダーと人種に関するバイアスを測定しています。

また、RLHF人間からのフィードバックによる強化学習）の導入がバイアス低減に有効であることや、偏見のあるプロンプトに対して生成を拒否するChatGPTの防御機能が有効であることが明らかにされています。この研究は、LLMの活用においてバイアスの低減が依然として重要な課題であることを示しています

Ren, R.; Basart, S.; Khoja, A.; Gatti, A.; Phan, L.; Yin, X.; Mazeika, M.; Pan, A.; Mukobi, G.; Kim, R. H.; Fitz, S.; Hendrycks, D.: Safetywashing: Do AI Safety Benchmarks Actually Measure Safety Progress?, arXiv preprint, arXiv:2407.21792 (2024)

概要：

AI安全性ベンチマークの多くは、AIの一般的な能力（例：知識、推論力）と高い相関があり、必ずしも安全性の進展を測定していない可能性があると指摘しています。これにより「Safetywashing」（安全性の向上と見なされるが、実際には能力の向上にすぎない現象）が発生しやすくなります。本研究は、AI安全性の真の向上を評価するために能力と独立した安全性メトリクスが必要であるとし、より厳密なAI安全性研究の枠組みを提供することを提案しています。

この論文は既存の安全性ベンチマークがモデルの能力とどれほど絡み合っているかを実証的に分析した初のメタ分析であると強調されています。

Cheung, M.: A Reality check of the benefits of LLM in business, ACM, pp. 1-20 (2024)

概要：

本論文はLLMのビジネスプロセスへの実用性と制約を定量的に評価した研究であり、LLMは一定の領域で有用であるものの、バイアス、文脈理解の不足、プロンプト依存といった課題があるため、組織がLLMを導入する際には慎重な検討が必要であると結論付けています。

この実験では、LLMが新しい分野について参考になる論文をどれくらい正確に推薦できるかを調べています。具体的には、5つのサーベイ論文を選び、それぞれの論文タイトルをプロンプトとしてLLMに入力し、50件の推薦論文を出力させました。その後、それらの推薦論文が実際にサーベイ論文内で参照されているかどうかを確認しました。この実験により、LLMが新しいプロジェクトの参考文献を提供する際の精度や、テーマの人気度が推薦結果にどのように影響するかを評価しています。