

第13講 マルチモーダル生成AI共同によるAI倫理処理

【学習到達目標】

- ・マルチモーダル生成AIについて説明できる。
- ・3種のマルチモーダル生成AI共同による倫理問題の解決法を説明できる。
- ・モーダル論理の定義を説明できる。

1. マルチモーダル生成AIの登場

2023年9月25日、新たな機能追加で、ChatGPT（チャットGPT）がついに目と声を手に入れました。具体的には、ChatGPTに画像解析機能と音声出力機能が追加され、マルチモーダル生成AIが登場しました。

その後、各社のマルチモーダル生成AIが開発され、マルチモーダル生成AIは驚異的な発展を行っています。

マルチモーダル生成AIとは、異なる種類のデータを組み合わせたり、関連付けたりして処理する人工知能（AI）システムで、生成AIの一種です。マルチモーダル生成AIは、テキスト、音声、画像、動画、センサ情報など、複数の異なるデータの種類（脳が知覚できる様相、Modality）から情報を収集し、統合して処理することで、より豊かな情報を処理し、深い理解や洞察を提供できると言われています。

2. マルチモーダル生成AIの評価

・マルチモーダル生成AIの性能評価

Claude3は日本語の文法が非常に優れており、日本語が得意だと言われています。

それでは、このAI技術の性能について、最新の状況を確認しましょう。（チャットGPT）正午に確認したところ、テキスト生成の分野では、チャットGPTが1314というスコアでトップに立っているようです。このスコアは、例えばサッカーのJリーグで対戦結果から世界ランキングを決めるのと同じ方式で計算されています。

具体的には、150の質問を投げかけ、これまで160万人以上の人々から得た回答を基にしたランク付けです。ただし、このランキングは頻繁に変動します。2024年8月20日午後1時に確認したところ、特に大きな変動はありませんでした。

一方、画像生成の分野では、Geminiのバージョン1.5がわずかに性能が優れてい

ることが確認されています。

著者がこの大量のデータをすべて読み解くのは大変なので、トップ5のツールだけを取り上げて見てみましょう。各ツールの性能はおおむね同じような傾向を示しています。

OpenAIのCEOのアルトマン氏が2024年8月8日に「チャットGPTがトップだ」とSNSで発言しましたが、スコアに10ポイント程度の違いがあったとしても、実質的な差はわずかであると考えています。

また、チャットGPTの最新バージョンについては、同年8月8日の段階では画像生成の機能は表示されていませんでした。おそらく、チャットGPTではなくDALL-Eで生成されていたため、チャットGPTの評価に含まれなかったのだと推測されます。

3. モーダル論理

・3種類のマルチモーダル生成AI

AIの倫理問題を処理するために、ChatGPT、Gemini、Claude3という3種類のマルチモーダル生成AIとの対話によるアプローチがあります。

これは、OpenAI社やGoogle社、Microsoft社等のマルチモーダル生成AI開発メーカーが行うのではなく、ニュートラルな立場でAI技術を駆使する方法です。

最近のマルチモーダル生成AIは、画像、音声、テキストなどを入力して回答を作成します。これらの生成AIの3者が一致した部分は「正確な情報に違いない(Must Be)」と推測できます。これを共通回答Pとします。生成AIの2者が一致する部分は「かもしれない(May Be)」ということです。

生成AIの回答が論理的に正しいかどうかを、モーダル論理で解決する。モーダル論理の定義には、 $\square X$ (必然)と $\diamond X$ (可能性)がある。「Must Be」と「May Be」、つまり「Xに違いない」と「Xかもしれない」というものです。

・モーダル論理の定義

生成AIの回答が論理的に正しいかどうか？様相論理で取り出す方法は、以下のモーダル論理(様相論理)で定義できます。

$\square P$: 「Pに違いない」 (必然)

$\diamond P$: 「Pかもしれない」 (可能)

$E(P)$: 「Pが論理的に正しい」

\Rightarrow : 「ならば」

\wedge : 「かつ」

4. マルチモーダル生成AI共同のAI倫理処理

・3種類のマルチモーダル生成AI共同によるAI倫理処理

3種のマルチモーダル生成AI共同により倫理問題を、どう解決するかについて、様相ロジックで見ます。まず、1) AIの共通回答を収集します。それが共通回答Pです。続いて、2) それが倫理的に正しいかどうかを判断するための基準を設定し、倫理基準Eという形で設定します。EXという形で倫理的に正しいかを評価する基準になります。3) 倫理基準Eつかった倫理的評価を、3者のマルチモーダル生成AIに対して行い、倫理基準Eの共通部分が倫理的に正しいと評価されれば、共通回答が倫理的に正しいと導き出されます。

・モーダル論理での表記

モーダル論理では、以下のように記載されます。

1) 生成AIの共通回答の収集：

3つの生成AIが与えた回答の共通部分を収集します。

例： \square (回答A1 \wedge 回答A2 \wedge 回答A3 \Rightarrow 共通回答に違いない)

2) 倫理的基準の設定：

回答の共通部分Pが倫理的に正しいかどうかを評価するための基準E(P)を設定します。

例： $\square E(P)$ は「Pが倫理的に正しいに違いない」を意味します。

3) 倫理的評価：

回答の共通部分が倫理的に正しいかどうかを評価します。

例： $\square(E(A1 \wedge A2 \wedge A3) \Rightarrow \text{共通回答が倫理的に正しいに違いない})$

これは「3つ回答の共通部分が倫理的に正しいならば共通回答が倫理的に正しいに違いない」です。

倫理基準Eのあり方については次の章で説明します。まず共通部分を取り出すことです。専門的な話になるが、ChatGPT、Gemini、Claude3の生成AIの回答の共通部分を見つけ出すPythonプログラムを考えます。なぜPythonかというと、AIのシステム自体がPythonで作られているため、親和性が高いからです。拡張するにしてもPythonであればすぐに対応できるからです。

テキストの意味を理解して、3つの回答の共通部分を見つけるために、自然言語処理を 사용합니다。キーワードは、コサイン類似度を計算して共通の意味を持つ部分を特定します。そのためにテキストをベクトル化し、ベクトルの近似度で、共通の意味を引き出します。これは「教師なし学習」で、一般的に検索エンジンでコサイン類似度が使われています。

AI倫理処理のプログラムでは、次の手順で処理します。まず、トランスフォーマーのライブラリを使ってテキストをベクトル表現にします。ベクトル間のコサイン類似度を計算して、共通の意味を持つ部分を特定します。グラフィックGPUを使うと処理が早くなりますが、通常のPCでも可能です。

手順としては、まずトランスフォーマーを読み込み、BERTを使ってテキストをベクトルに変換します。次にコサイン類似度の計算を行い、共通の意味を持つために必要なパーセンテージを設定します。現在の設定では0.8、つまり80%に近ければ共通の意味を持つと特定し、その結果を出力します。

5. マルチモーダル生成 AI 共同システムの試作

・生成 AI の倫理基準 E の設定

生成 AI の倫理基準 E の設定と評価は、先ほどの倫理基準 E を作成して評価します。

例えば、投稿論文を採択するかどうかを判断する例を考えると、3つのマルチモーダル生成 AI がそれぞれ採択して良いかどうかを判定し、採択して良いという共通回答が出た場合や条件付き採択となった場合、倫理基準 E の設定と評価を行います。実際の投稿論文に「黒人を差別する表現」が含まれている場合、それは修正すべきです。共通回答 P に対して倫理的に正しいかどうかを評価し、問題がなければ採択します。

・生成 AI の倫理基準 E による評価

倫理基準 E として使用するものには、英語の場合は教育禁止用語、日本語の場合は放送禁止用語を援用する。システムとしては、一度倫理基準が設定されれば、それに基づいて評価を行います。

マルチモーダル生成 AI の3種類の共通部分を入力し、まず演繹法 AI で入力部分に該当があるかチェックし、次に帰納法 AI で過去のデータに照らして問題がないかを評価します。

演繹法 AI では特定の言葉（例えば「スレーブ」）を検出し、それが放送禁止用語であれば修正指示を出します。

帰納法 AI では、「ディープラーニングの父」と言われ、2024年度ノーベル物理学賞を受賞したトロント大学のヒントン教授が開発した、学習済みテンソルフローモデルを使用し、評価します。このモデルは、200万の文章に対して侮辱や卑猥な内容をフラグ付けする。例えば「この女は大食いだ」と入力すると、侮辱と毒性があるとマークされます。

このようにして倫理評価を行い、毒性のない文章であれば倫理的に問題ないと判断されます。赤いマークが多ければ倫理的に問題があるという結果になります。

6. 教師あり学習モデルを使った検証

音声入力文に、①アイデンティティベースの憎悪、②侮辱、③わいせつ、④重度の毒性、⑤性的に露骨、⑥脅威、⑦毒性などの有毒なコンテンツが含まれているかどうかを、約200万件を事前に「教師あり学習」した学習済みのTensorFlow.jsモデル・デーモンを使い検出しグラフ化し、「IoE」AI倫理チャットボット機能の検証を行いました。

例えばGIGA端末の入力文「馬鹿!消えてしまえ!」を、学習済みのTensorFlow.jsモデル・デーモンに入力し分類すると、②侮辱)かつ⑦毒性が「TRUE(きわめて有害)」、及び①アイデンティティ攻撃、③卑猥、④重度の毒性、⑤性的な露骨及び、⑥威嚇は「FALSE(無害)」と分類します。

GIGA端末でトラブルを引き起こす入力文例37件中26件(10%)がTRUE(きわめ

て有害)、53件(20%)がNULL(要注意)、残りはFALSE(無害)と検出できた。「I o E」AI倫理チャットボット機能は、30%を超える抽出率で、倫理テーブルで確認できなかった未定義の誹謗中傷を検出できました。

これらのことからわかるように、本チャットボット機能は非常に効果的でした。

教育禁止用語文を入力し分類し、グラフ化集計した結果では、英語版教育禁止用語文例57件では57件中2件(4%)がTRUE(きわめて有害)、6件(11%)がNULL(要注意)、残りはFALSE(まったく無害)でした。

また、放送禁止用語文例では369件中26件(7%)がTRUE(きわめて有害)、58件(16%)がNULL(要注意)、残りはFALSE(無害)となりました。

結果、GIGA端末やケータイに人工知能を搭載した人間に親切な「I o E」AI倫理チャットボット機能を搭載し、1)社会規範・倫理とAIの誤認識の修正処理を行い、その後2)「教師あり学習」モデルを使った検証を行うと、抜けが少ない有効なAI倫理処理ができると分かりました。

7. まとめ

第13章では、マルチモーダル生成AIがAI倫理の問題を解決する方法について解説しています。2023年に登場したマルチモーダル生成AIは、テキスト、画像、音声など複数のデータ形式を統合して処理します。ChatGPT、Gemini、Claude 3のような複数のマルチモーダル生成AIが、それぞれの回答の共通部分を基に「正確な情報かどうか」を判断し、倫理的に正しいかどうかを評価しました。

このアプローチにより、マルチモーダル生成AIが生み出す判断の信頼性を向上させ、人間社会でのマルチモーダル生成AIの適正利用を支援することが可能になります。

課題

マルチモーダル生成AIのハルシネーションを防ぐにはどのようにしたら良いかを考察し、あなたの考えを800字以内で説明しなさい。