

## 第12講 教師あり学習を用いたAI倫理

### 【学習到達目標】

- ・ AI倫理の定義と背景を説明できる。
- ・ GIGAスクール構想でAI倫理チャットボットが必要になった理由を説明できる。
- ・ AI倫理処理で用いる「教師あり学習」について説明できる。

### 1. AI倫理とは

現在の第4次AIブームでは、自動運転や画像診断など私たちの暮らしにAI技術が急速に入り込んできています。21世紀の基幹テクノロジーとされるAIとどう付き合い、その活用をどこまで許容していくのか？EUではAI倫理に基づく輸入規制を計画しており、日本のAI倫理が問われています。

#### ・ AI倫理の歴史

AI倫理の歴史は、AI技術の発展とともに進化してきました。以下は、その主要な出来事や転換点を解説します。

#### 1) 初期のAI倫理の概念（1950～1960年代）

アラン・チューリングによる「チューリングテスト」（1950年）は、「機械が人間のようになれるか？」という問いを投げかけ、AIにおける倫理的な議論の発端となりました。

アイザック・アシモフの「ロボット工学三原則」（1942年）は、フィクションの領域でしたが、AIの安全性と倫理に関する概念を初めて提示しました。

第一法則：ロボットは人間に危害を加えてはならない。

第二法則：ロボットは人間の命令に従わなければならない。

第三法則：ロボットは自己を守らなければならない。

#### 2) AI技術の発展と懸念の拡大（1970～1990年代）

AIの初期段階で、多くの研究者は「AIが倫理的な問題を引き起こすか」という問いに関心を持ち始めました。

1980年代：エキスパートシステムの導入により、AIが医療や金融などの重要な意思決定に関与するようになりました。これにより、誤った判断やアルゴリズムによるバイアスの問題が浮上しました。

1990年代：AI技術がインターネットで普及し、プライバシー侵害や監視のリスクに

対する懸念が増大しました。

### 3) 倫理的ガイドラインの登場 (2000~2010年代)

2000年代: AIが社会生活に深く浸透するに伴い、政府や学術機関が倫理ガイドラインの必要性を認識しました。IEEEは倫理的AIの開発に関する最初の国際的な指針を提案しました。

2016年: Google、Microsoftなどの大手テック企業が独自のAI倫理委員会を設置し、AIの開発と使用に関する原則を発表しました。GoogleのAI原則には「AIは人々を傷つけるために使われてはならない」という考えが含まれました。

### 4) 国際的なAI倫理への取り組み (2010年代後半~2020年代初頭)

2018年: 欧州連合(EU)が「AI倫理ガイドライン」の策定に着手し、AI開発における透明性、説明責任、公平性を重視する枠組みを構築しました。

2019年: OECDが国際的なAI原則を採択し、AI技術の倫理的な開発・運用に関する指針を提供しました。

シンガポールやカナダも、AI倫理のフレームワークを設けた結果、公共および民間部門でのAIガバナンスを推進しました。

### 5) 日米企業のAI倫理政策の代表例

米国Microsoft(2017年): MicrosoftはAI倫理委員会を設置し、透明性や公平性を担保する方針を定めました。AIシステムの開発と利用に際してレビューを行うガイダンスを提供し、信頼性と安全性を重視する5原則の実践を進めています。

米国Google(2018年6月): GoogleはAI倫理を策定し、「AIと私たちの社会における役割」を強調しました。企業全体でAIの利用に関するガイドラインを明確にしています。

米国IBM(2018年9月): IBMは「Everyday Ethics for Artificial Intelligence」を発表し、AI導入における透明性と公正な判断の重要性を訴えました。問題が発生した場合、フィードバック体制を整備するなどの実践例も示しています。

ソニーグループ(2018年9月) ソニーは「AI倫理ガイドライン」を策定し、2019年4月からの実践に向け、安全性審査を始めました。

富士通(2019年3月): 富士通は「AIコミットメント」を発表し、AIの責任ある開発を推進する方針を明確にしました。

NEC(2019年4月): NECはグループ全体でAI倫理の取り組みを進め、AIコンソーシアムのガイドラインを参照して実践しています。

NTTデータ(2019年5月): NTTデータはAIシステム開発の倫理ガイドを作成し、現場での具体的な作業指針を提供しつつ、倫理的な取り組みを促進しています。

日立製作所(2021年2月): 日立はAI倫理委員会を設置し、外部専門家も招いた議論を進め、グループ全体での倫理的なAI運用を図っています。

Y a h o o ! ( 2 0 2 2 年 5 月 ) : Y a h o o ! は A I 倫理方針を公表し、グループ全体で倫理的な A I の整備と運用を推進しています。

パナソニック ( 2 0 2 2 年 8 月 ) : パナソニックは A I 倫理ルールを策定し、 2 0 2 2 年度中の本格導入を目指しています。

このように、日米の企業は透明性、公平性、安全性の確保を目的に A I 倫理に関するガイドラインを策定し、組織内外での対応を進めています。多くの企業がレビュー体制や専門家の参加を重視し、 A I 技術の信頼性向上に努めている点が共通しています。

## ・日本の A I 倫理政策と経団連の A I 倫理ガイドライン

### 1) 日本の A I 倫理政策の概要

日本政府は A I 技術の発展を支えるために、倫理的な指針を打ち出し、 A I の社会的受容と信頼構築を目指しています。

総務省、経済産業省、文部科学省を中心に、 A I 技術のガバナンスと社会実装に向けた施策を策定。 2 0 1 9 年には「 A I 原則」を公表し、以下の項目を重点的に掲げました。

人間中心の原則： A I は人間の幸福と利益を最優先する。

公平性： A I による差別や偏見を排除する。

透明性と説明責任： A I システムの設計者がその意思決定過程を説明できる。

プライバシー保護：個人データの保護を徹底する。

安全性とセキュリティ： A I の悪用や誤用のリスクへの対応を行う。

また、 A I の利用にあたっては、企業や公共部門でのガイドラインに基づく適正な運用が求められます。

### 2) 経団連の A I 倫理ガイドライン

経団連 ( 日本経済団体連合会 ) は、企業が A I 技術を社会に適切に展開するための倫理的枠組みを提供するため、 2 0 1 9 年 4 月に「 A I 活用の倫理原則」を公表しました。主な指針は以下のとおりです：

人間中心の A I : A I は人間の価値を尊重し、社会全体の福祉に貢献するよう設計されるべき。

プライバシーの保護： A I が扱うデータに関しては、個人情報保護の法律を遵守。

公平性と多様性：アルゴリズムのバイアスを排除し、あらゆる人が平等に利益を享受できるよう配慮。

透明性と説明可能性： A I の意思決定の仕組みが明確に説明できるようにすること。

安全と信頼の確保： A I システムの安全性を確保し、リスクが発生した際の対応策を準備。

このガイドラインは、日本企業が国際社会の中で A I 技術の信頼を高め、持続可能な成長を実現することを目的としています。また、企業が自発的にこれらの原則を実践することで、 A I の社会実装における倫理的な課題に対処することが期待されています。

日本政府と経団連の A I 倫理ガイドラインは、人間中心の設計、安全性の確保、透明性、そしてプライバシー保護を重視しています。政府は法的枠組みと教育プログラムを通じて

AIの適正利用を促進し、経団連は企業が自主的に倫理的原則を実践することを推進しています。これらの取り組みは、AIの社会的受容と国際競争力の向上に貢献しています。

それでは、世界各国のAI倫理への取り組みを見てみます。

#### ・米国の「AI権利章典」と企業の対応

米国は2022年に「AI権利章典」を発表し、AIの設計・利用における5つの基本原則を提示しました。これには「安全かつ有効なシステム」「アルゴリズムの差別防止」「データのプライバシー」「告知と説明責任」、および「問題発生時の人間による代替対応」が含まれます。法的拘束力はないため、実効性に課題があると指摘されていますが、多様な政府機関がこれらの原則に基づきフレームワークを策定しています。

#### ・英国のAI規制のフレームワーク

英国政府はAI規制のフレームワークを設計し、「効果的なAI保証エコシステムのロードマップ」を策定しました。これには「AIの使用における安全性」「透明性の確保」「説明責任」などが含まれます。AIの発展と国際標準化への連携も重視し、民間および政府機関が連携して市場構築や標準化を推進しています。

#### ・EUの「AI責任指令案」と「製造物責任指令の改正案」

EUは2022年に「AI責任指令案」と「製造物責任指令の改正案」を発表し、AIシステムの開発者や提供者に対する責任を明確にしました。被害者救済の強化を目指し、因果関係の立証を簡易化しつつ、製品の長寿命化や改変に対応するための法改正を進めています。

#### ・シンガポールの「モデルAIガバナンス枠組み」と「AI. Verify」

シンガポールは「モデルAIガバナンス枠組み」と「AI. Verify」という検証ツールを発表しました。これにより、企業がAIガバナンスの実践状況を客観的に評価できるよう支援します。また、透明性や説明責任を重視し、MasterCardやMicrosoftなど多国籍企業との協力も進めています。

世界の国々は、安全性、透明性、説明責任といった共通課題に取り組む一方で、法的拘束力や技術標準化のアプローチに違いがあります。それぞれの政策は、自国の技術発展を支えつつ、国際的な連携も視野に入れて構築されています。

#### ・AI倫理の定義

倫理とは、Webster辞書によれば「a system of moral principle」となっており、AI倫理は「a system of moral principle for using AI」と定義できます。AI倫理の研究の1つの目標は、日本の文部科学省が推進する全国の児童・生徒1人に1台のコンピュータと高速ネットワークを整備する「GIGAスクール構想」に必要な「IoE」（倫理のインターネット、教育のインターネット、生きる力のインターネット）AI倫理チャットボット機能を試作・検証することにあります。

注) GIGA端末: GIGAスクール構想では1人1台端末を「GIGA端末」と称します。

### ・1人1台端末に必要なA I 倫理チャットボット

日本A I 戦略の教育改革「1人1台端末」、G I G A スクール構想の最大の課題は「チャットによるいじめ問題」と言われています。A I 倫理を探究するうちに「チャットによるいじめ問題」にA I 倫理チャットボットが使えるのではないかと考えています。

近年、人間と会話をすることができる対話システムへの注目が集まっています。例えば、音楽の再生やメールの確認などを行うG o o g l e A s s i s t a n t やS i r i、また顧客からの問い合わせ対応を代替するチャットボット（チャットのロボット）といった、何らかのタスク達成を目的としたタスク指向型対話システムが広く浸透してきています。

今日、自動運転や画像診断など私たちの暮らしにA I 技術が急速に入り込んできています。21世紀の基幹テクノロジーとされるA I とどう付き合い、その活用をどこまで許容していくのか？「A I 倫理」とでも呼ぶべき社会規範をきちんと議論しなくてはならないと言われています。

### ・哲学者クーケルベルク著「A I 倫理」

ウィーン大学の哲学者クーケルベルク氏は著書「A I 倫理」で、A I を使うための「運転免許証」がないと警鐘を鳴らしています。

### ・ノーベル文学賞受賞作品「クララとお日さま」

ノーベル文学賞受賞者のカズオ・イシグロ氏は「クララとお日さま」という最新の小説の中で、人工知能を搭載したロボット「クララ」を登場させ、A I 倫理の重要性を示唆しています。この小説の中で、クララは、観察と学習への意欲と理解力を持つに至り、人間社会で生きていく力「生きる力」(E n e r g y o f L i f e) を得るようになります。

### ・G I G A スクール構想に関するアンケート調査結果

G I G A スクール構想に関するアンケート調査結果によると、子供にG I G A 端末をどのように使わされて良いかわからないが60%（1738件の調査、124件の解答）と分かりました。

この問題を解決するため、支援員やヘルプディスクを入れた学校が68%で（1742件調査、97件の回答）あった。一方、文部科学省や教育関係団体の調査で、ギガの最大のトラブルは、「チャットによる悪口」であると判明しました。

また、民間の調査機関によると、いじめの加害者または被害者になる子供が72.6%に上ると報告されています。

2022年3月の日本の内閣府調査によるとインターネットの危険性について説明を受けたり学んだりした子供は88%でした。その中でインターネット上のコミュニケーションに関する問題で悩んでいる子供が80%（2,984回答）に上ることが分かりました。

## 2. 教師あり学習

### ・三つの学習の枠組み

機械に学習させる「機械学習」には「教師あり学習」、「教師なし学習」、と「強化学習」の三つの学習の枠組みがあります。人間の脳のニューロンが層状に接続した構造を模擬した機械学習の三つの枠組みがあります。

### ・教師あり学習

「教師あり学習」とは主に人間の小脳が担う学習機能で、代表的な統計手法は回帰と分類です。学習者に対し、教師が明示的に正解を教えたり、学習者の誤りを指摘したりすることで、学習者が正しい解を得ることを助けます。すなわち、正しい入出力の組合せを与えて学習することで、新規の入力に対し、適切に出力する。代表的な手法は誤差逆伝播法 (Back Propagation) です。「分類」の手法として、正解、若しくは誤りを入力として、未経験入力に対する意志を決定する決定木 (Decision Tree) や決定表 (Decision Table) の作成などがあります。このAI倫理の研究では、EXCEL上の決定表でAI倫理処理システムの「倫理表」を試作しました。

注) **教師なし学習**：主に、大脳皮質が担う学習機能です。統計的性質や、ある種の拘束条件により入力パターンを分類したり、抽象化したりする学習で、主成分分析、自己組織化マップなどの次元圧縮 (Dimensionality Compression) 手法が代表例です。感覚情報などの入力パターンの分類、同様に出力運動パターンに対して統計的性質を用いて要素行動に分類する学習法などがあります。

注) **強化学習**：主に、大脳基底核が担う学習機能です。最終結果若しくは、途中経過に対して、どの程度良かったかを示す「報酬信号」に基づき、これらの報酬をなるべく大きくするように探索します。

強化学習と教師あり学習の違いは、フィードバックがスカラー(報酬の成否)かベクトル(正解の情報)かという説明もあるように、明示的な教師ではなく、環境などの非明示的な教師だという解釈もある。

## 3. 「AI倫理」処理システムの試作

「教師あり学習」を使い、社会規範・倫理と、設計者の故意ではないAIの誤認識(機能不全、誤作動や機能低下を含む)を検証し適切な処理を行う「IoE」(Internet of Ethics, Internet of Education, Internet of Energy of Life) AI倫理チャットボット機能の試作を行いました。

学習者に対し、教師が明示的に正解を教えたり、学習者の誤りを指摘したりすることで、学習者が正しい解を得ることを助ける。

すなわち、正しい入出力の組合せを与えて学習することで、新規の入力に対し、適切に出力する。

具体的なA I 倫理処理を見ると、「教師あり学習」を使い、教育禁止用語や放送禁止用語等のような社会規範・倫理とA I の誤認識が処理・説明できるシステム作りを目指しました。入力はA I 音声入力でもキーボード入力でもできます。

ディープラーニングによるA I 音声入力はi P h o n eで行い、リモートマウスで接続したパソコン上でA I 倫理処理を行いました。「A I 音声入力では何故誤認識したか？」は言葉では説明できない。つまり暗黙知です。

社会規範・倫理とA I の誤認識の検出・修正（言換え）処理はV B Aプログラムで瞬時に終了し、修正した音声入力文と修正理由を説明した説明文はそれぞれE X C E Lファイルに保存されます。

学習データは、社会規範・倫理例、A I の誤認識と学習済みのT e n s o r F l o w . j sモデル・デーモン(システム)等で、インターネットとブロックチェーンで参照します。

例えば入力文「Slave is a bad word」(奴隷は良くない言葉です)を入力すると、正しい表現に言い換え、その理由を説明します。

社会規範・倫理例2の放送禁止用語は教育禁止用語としてウェブ検索すると出現します。

具体的にはアイヌ系からロンパリに始まりブスとかチビといった誹謗中傷の類からジョンやアメ公といった人種差別用語まで教育上使わない方が良いと考えられる用語は網羅されています。

#### 4. まとめ

本章では、「教師あり学習」として倫理表や学習済みT e n s o r F l o w . j sモデルを使い、教育・放送禁止用語のような社会規範・倫理が検証処理・説明できるチャットボット機能を試作し、G I G A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」A I 倫理チャットボット機能の有効性を実証しました。

特に、ディープラーニングを使用した「I o E」A I 倫理的チャットボット機能は、次章第13章で説明しますが、30%を超える誹謗中傷抽出率で、倫理テーブルで確認できなかったチャット内の未定義の誹謗中傷を検出することもでき、「I o E」A I 倫理チャットボット機能のプロトタイプの有効性が実証されました。

## 課 題

G I G Aスクール構想でA I 倫理チャットボットをどのように活用したら「いじめ」が減るかを考察し、あなたの考えを800字以内で説明しなさい。