

PART 7 3 世界の潮流

AIの倫理ガイドラインと透明性の原則①

国際社会経済研究所 (NECグループ) 主幹研究員



小泉 雄介

められており、代表的な取り組みとして17年のアシロマAI原則がある。また、経済協力開発機構(OECD)は19年5月にAIに関するAIの要件④信頼できるAIを実現するAIの要件⑤信頼できるAIを実現するための技術的/非技術的手段⑤信頼できるAIのアクセスメントリスの五つから成るフレームワークを提示している。

人工知能(AI)技術は、個人にも企業や社会にとっても多大なベネフィット(便益)をもたらすことが期待されているが、同時に個人や社会に対して新たな害悪を及ぼすリスクについても懸念されている。

人間中心のAI

AI開発や利用の諸原則については、既に国内外の政府機関などがガイドライン類を発している。わが国では総務省のAIネットワーク社会推進会議が

6月にAI活用ガイドライン案を、内閣官房の「人間中心のAI社会原則会議」が3月市で6月上旬に開催された主要20カ国・地域(G20)貿易・デジタル経済相会合では、OECD原則と同じ内容のAI原則を含む閣僚声明が採択された。

7つの要件

②は「人間の白律の尊重の原則」「害悪防止の原則」「公平性の原則」

欧州委員会AI倫理ガイドラインの7要件

Table with 2 columns: 要件 (Requirements) and サブカテゴリー (Sub-categories). It lists 7 requirements such as 'Human autonomy and supervision', 'Technical robustness and security', etc.

要件への適合を評価するためのパイロット版のリストである。

AIの透明性

本稿のテーマであるAIの透明性に関するものとして、「理解可能性の原則」では、「なぜシステムが特定のアウトプットや判断を行ったのか、またどのインプットの組み合わせがアウトプットに寄与したかについての説明は、常に可能なわけではない(いわゆるブラックボックス問題)。それらのケースでは、他の理解可能性

害悪回避へ倫理原則

「理解可能性の原則」の四つから成る。これらの倫理原則が挙げられている。さらには設計段階からの倫理、説明できるAI(XAI)など、非技術的手段には法規制、行動規範、標準化、認証などが含まれている。また⑤は、上記7