AI法案とAI制度研究会: "リスクへの対応"に関する講師の指摘を中心に



ELSI大学サミット

自己紹介

https://c-research.chuo-u.ac.jp/html/100003450_ja.html (last visited Mar. 15, 2025).



AIGO, OECD, Paris, France, in Sept. 2018.



平野 晋

(中央大学 国際情報学部 教授・学部長ニューヨーク州弁護士)

〇 中央

略歴

中央大学法学部法律学科卒 / コーネル大学(法科)大学院修了 (LL.M)・『CORNELL INT'L L.J.』編集委員 / 博士(総合政策)(中央大学) / 富士重工業㈱法規部主任、㈱NTTドコモ法務室長などを経て、中央大学総合政策学部教授、同大学国際情報学部開設準備室長など。

国際機関/政府の委員等

- 経済協力開発機構(OECD)「AI専門家会合」日本共同代表(2019年)
- 内閣府「AI制度研究会」構成員
 - 同 「人間中心のAI社会原則会議」構成員
- 総務省「AIネットワーク社会推進会議」構成員(副議長)
 - 同 「AIガバナンス検討会」構成員(座長) など

著書·論文

- 『ロボット法:ヒトとAIの共生にむけて(増補第2版)』(2024年, 弘文堂)
- 「AIに不適合なアルゴリズム回避論:機械的な人事採用選別と自動化バイアス」『情報通信政策研究』第7巻2号1頁(総務省,2024年3月)
- 『アメリカ不法行為法』(2006年, 中央大学出版部)
- 『電子商取引とサイバー法』(1999年, NTT出版) など

講演の背景



Paris, France in Oct. 2017.

内閣府「AI制度研究会」

「イノベーションの促進とリスクへの対応の両立」



首相官邸「AI戦略会 議・AI制度研究会」 令和6 年8月2日 https://www.kantei. go.jp/jp/101_kishida /actions/202408/02

ai.html (last visited

Aug. 15, 2024).



「AI法案(*)」

「イノベーションの促進とリスクへの対応の両立」

(*)「人工知能関連技術の研究開発及び活用の推進に関する法律案」

一方におきまして、今も御指摘ありましたが、AIには様々なリスクがございます。これらへの対応が必要であります。座長から御説明いただきました中間とりまとめ案に沿いまして、城内大臣を中心に、平大臣ほか関係閣僚が協力し、AIのイノベーション加速とリスク対応を両立させる新たな法案を早期に国会に提出できますよう、対応を進めていただきたいと存じます。

政府におきますAI政策の司令塔機能を強化するため、全閣僚からなります『AI戦略本部』を設置いたします。



会場のまとめを行う石窟総理1

関連動画 +

令和6年12月26日、石暖地理は、地理大臣官邸で第12回AI(AI知能)戦略会議・第6回AI制度研究会合同会議に出席しま た。



首相官邸「総理の一日」「AI戦略会議・AI制度研究会合同会議」令和6年12月26日 https://www.kantei.go.jp/jp/103/actions/2024 12/26ai.html (last visited Mar. 16, 2025).



「AI法案(*)」

「イノベーションの促進とリスクへの対応の両立」

(*)「人工知能関連技術の研究開発及び活用の推進に関する法律案」

AI制度の在り方につきましては、昨年8月から、AI戦略会議の下に設置したAI制度研究会におきまして、検討を開始し、AI関係者へのヒアリングや議論、パブリックコメントを重ね、本年2月4日に「中間取りまとめ」を決定したところでございます。この「中間取りまとめ」を踏まえ、「AI法案」の詳細についての検討を進めてまいりました。

本法案は、AIは、生産性の向上や労働力不足の解消などのメリットをもたらす一方で、偽情報・誤情報の拡散、犯罪の巧妙化などのリスクも存在することから、国民生活の向上や経済社会の発展の実現には、AIによるイノベーション促進とリスク対応を両立させることが重要との考えのもと、AI政策の司令塔機能を強化する、内閣総理大臣を本部長とする「AI戦略本部」の設置、政府が推進すべきAI政策の基本的な方針等を示す「AI基本計画」の策定、AIの適正性確保のための国際規範に即した「指針」の整備、AIの動向に関する情報収集や国民の権利利益を侵害する事案の調査などについて、規定しているものであります。





城内内閣府特命担当大臣記者会見要 旨 令和7年2月28日

https://www.cao.go.jp/minister/2411 m kiuchi/kaiken/2025 0228kaiken.html (last visited Mar. 16, 2025).



内閣府「AI制度研究会」 「イノベーションの促進とリスクへの対応の両立」

12 基本的な考え方 13 ■国際協調 (II.4.) ■イノベーション促進とリスク対応の両立 15 AIガバナンスの形成に 研究開発支援、人材育成、データや計算資 16 向けて議論をリード 源の整備などイノベーションの促進 17 ■国際整合性・相互運用 ★令とガイドライン等の適切な組合せ 18 性の確保 OECD原則、広島AIプロセス国際指針等の 19 信頼できるAI 共通的な指針等と個別の既存法令の活用 20 21 22 23 具体的な制度・施策の方向性 24 AIの研究開発・実装が最もしやすい国を目指す ■全般的な事項 (Ⅲ.1.) 速やかな法制度化が必要 26 世界のモデルになるような制度 AIの安全・安心な研究開発・活用のための戦略(基本計画)の策定 28 29 安全性の向上等 30 国による指針(広島AIプロセス準拠)の整備、事業者による協力 ■ 政府等による利用 (Ⅲ.2.) ■ 基盤サービス等における利用 (Ⅲ.3.) 各業法等による対応 等 適正なAI政府調達・利用等 1) 官房長官が議長、全閣僚が構成員となっている「統合イノベーション戦略推進会議」の下に「AI戦 36 略会議」を設置。その下に「AI制度研究会」を設置。 37 上記の政策を講じた上で、今後のリスク対応のため引き続き制度の検討を実施すべき。 38

内閣府ホームページ「中間とりまとめ(案)」 https://www8.cao.go.jp/cstp/ai/ai_kenkyu/5kai/shiryou1.pdf(last visited Mar. 14, 2025).



講師が着目するリスク

- ・「不適切なJAI利用
 - ✓ 内閣府「AI制度研究会」→「AI法案(*)」
 - ✓ "ADM:" <u>Automated Decision-Making</u>、又は <u>Algorithmic Decision-Making</u>
- · AIへの過度な依存
 - ✓ 内閣府「人間中心のAI社会原則」
 - ✓ 平野『AIとヒトの共生にむけて』

(*)「人工知能関連技術の研究開発及び活用の推進に関する法律案」



「イノベーションの促進とリスクへの対応の両立」



ELSIは、ハンドルとブレーキー11

【高橋構成員】

2点確認したい。1点目は、事務局への確認である。資料1の冒頭で「ICTイン テリジェント化の影響とリスク」について議論があった。技術開発はアクセルにあた る。人文社会科学にはハンドルという側面もあるが、ブレーキでもある。ブレーキを 踏む人を確認することやそのためのルールが、正しいものになっているからこそアク セルを踏めるだろう。例えば、遺伝子工学では倫理委員会が大学毎にあるが、そうい う組織があるからこそ技術開発が進められてきた。そういった理解でよいか。

ICTインテリジェント化影響評価検討会議(AIネットワーク化検討会議「第1回議事概要」5~6頁,2016年2月2日(理研 生化学シミュレーション研究チーム チームリーダー高橋恒一構成員発言))

https://www.soumu.go.jp/main_sosiki/kenkyu/iict/index.html (last visited Aug. 15, 2024)..

(1) 構成員

須藤 修 (座長) 東京大学大学院情報学環教授

中央大学大学院総合政策研究科教授 平野 晋 (座長代理)

石井 夏生利 筑波大学図書館情報メディア系准教授

エンジンとハンドル/ブレーキ

STEM

Science,

Technology,

Engineering, and

Mathematics

<u>ELSI</u>

Ethical,

Legal, and

Social

Implications

理数工学系



人文/社会科学系

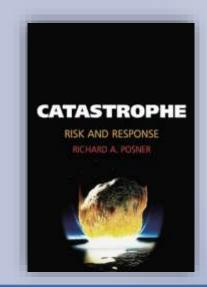
"信頼に値するAI"

"人間中心の Al"

【参考】 リチャード A. ポズナー判事の指摘 13

"[s]cientists want to advance scientific knowledge rather than to protect society from science; the policy maker's ordering of values is the reverse. Not that scientists are indifferent to public safety; but it is not their business and sometimes it is in competition with their business

RICHARD A. POSNER, CATASTROPHE: RISK AND RESPONSE 99 (2004)(emphasis added).





Conflict of the First Priorities

Policy Makers

E·L·S·I

1st Priority of Value:
Public Safety

Scientists

S-T-E-M

1st Priority of Value:

<u>Advancement of</u>

<u>Scientific Knowledge</u>

versus



AI制度研究会:「リスクへの対応」

- ・「不適切な」AI利用
 - ✓内閣府「AI制度研究会」→「AI法案」
 - ✓ "ADM:" <u>Automated Decision-Making</u>
 <u>Algorithmic Decision-Making</u>

講師が特に指摘したリスク

平野構成員 提出資料

MENORANDUM ***

Dec. 26, 2024 中央大学 国際情報学報長・教授

-

1. 個人の性格、能力、減損性、感情等をAIに評価・予測・決定 指引等させる問題。

CO (RESERVE CARROL MINE, NO FAME TRESPORT A LABOR. Ti-5. [gutomated-decision gaking) XH [glgsrithmin decision. gaking: WOMIST & d. See, e.g., Sancto Motato et al., Distriction is the Actionted Administrative State. 26 Cos. F.L. 6 Perry. 171 (2011). また、世界主要部子を禁制を1 Feartion voltage や fattert moognition; \$25020.71.0. dec. a.g.: Hooks V. Horre, Mirror for (Artificial) Intelligence: In Moose Seflection?, 41 Chep. Lan. L. 6 900'S J 47, 40 (2013) . なお AC が高様や声音等から情報。 **取力、成長性、文は単核報を提供することには何意的に検知がなく、私ればであると** 接续する例として使用になる文献としては、 Lake Stack & Javes Stotes Obpatogramic Artificial Intelligence, 30 Frames Sents. Pers Shorts & Ber. L.J. 525, 536-38 (2002). Show also Marcanata Lávesque, Atalog Privilege, 26 H.Y.J.J. Lerts, 4 Prs. Pts. 7 525, 454 (2024) (開展学習仁基づく確実完整等やの*technologies night improve over time, [] but been judgment convently outpurforms then by several orders of magnitude. "2 mm;; tages A. toges, Foliating Sections What social Psychology Cas Teach Fourth Assessment Section, 12 Nov. 3, Nov. 605, 605 (2004) ("A Large and growing body of research . . shows that there is no reliable evidence that homes can accountably and reliably Serect sectional states From Canial sepressions () (-0.000): Sandra Maybler, Confessione and Acapholes in the St Af Act and Af Alability Directives: What This Means for the Surger Online, the Online Ofster, and Septend. De Tale J.L. & Dale. 671, etc. [2020] (特殊日本ソントウエアが有機性主义

「AI制度研究会 構成 員提出資料」33/50~ 46/50頁

https://www8.cao.go.j p/cstp/ai/ai_kenkyu/5 kai/shiryou2.pdf (last visited Mar. 16, 2025).

- 例えば、〈雇用(含、採用)分野〉や2、〈教育分野〉3。
- そもそもヒトの判断よりも、AI の方が効率的であり、かつ客観的であるから中立的であると捉える前提が問題であると指摘されている。
- しかし AI もヒトが削ったものだから、ヒトの個見から逃れられない。
- そして不公正・差別・不正確等の問題が、明らかに成って来た。
- * See, e.g., OFF. OF SCI. & TRCH. POL'Y, BLUEDRINT FOR AN AI BILL OF RIGHTS 46 (Oct. 2022) (以下のように指摘: "Automated systems with an intended use within sensitive domains, including, but not limited to, criminal justice, employment, education, and health, should additionally be tailored to the purpose, provide meaningful access for oversight, include training for any people interacting with the system, and incorporate human consideration for adverse or high-risk decisions."(emphasis added)).
- * See Stark & Butson, supra note 1, at 957-58 (コロナ福時代に使用されたカンニング防止の為に目の動きを感知するソフトの使用が学生造から猛反発された事例を例示しながら、態度・表情等から学生を評価する「人相学的 AI」の教育分野に於ける使用は、生物学的決定論や優生学や科学的人種差別を元気づけるとして批判している).
- * See, Margaret Hu, Critical Data Theory, 65 WM. & Maxr L. Rrv. 839, 862 (2024). なお、そもそもヒトの将来を予測することなどは不可能な事実を人々は常識では理解しているにも拘わらず、AI を用いた途端に人々が常識を失って、「見せかけの正確性」("veneer of accuracy")に騙されると示唆する文献例として、see Stark & Hutson, supra note 1, at 929-30.
- **Bouston Fed. of Teachers, infra note 24, 251 F.Supp.3d at 1171 (**Algorithms are human creations, and subject to error like any other human endeavor.**と法廷意見が指摘). See also Ifeoma Ajunwa, The Paradox of Automation as Anti-Bias Intervention, 41 Capaca L. Rzv. 1671 (2020) (AI はヒトが作るのだから、暴走したAI の責任をヒトが負わねばならない、あたかもフランケンシュタイン博士が造った人造人間の暴走の責任を同博士が負わねばならなかったのと同じである。と指摘): Danielle Keats Citron, Technological Due Process, 85 Wash. L. Rzv. 1249, 1253 (2008) (適正手続の保障等の政策に強いプログラマが法を無視したプログラミングを行って法を重める問題を指摘): Crystal Godfrey, Legislating Big Tech: The Effects Amazon Recognition Technology Has on Privacy Rights, 25 U.S.F. Israin. Proc. & Tech. L.J. 163, 166 (2021) (プログラマの偏見が入り込むと指摘).
- ⁶ See authorities cited in supra note 1. See also Yonathan Arbel et al., Systematic Regulation of Artificial Intelligence, 56

 日本では(も)人事採用 AI のペンダッ が活発に売り込み、かつそれを利用 する企業も少なくないように見受けるけれども、「やり過ぎ」に注意が必要 ではないか。

ARIZ. St. L.J. 545, 557 (2024); Godfrey, supra note 5, at 168(頻 認識技術—FRT—が感情を懸み取れるという主張が批判されていると指摘). 雇用 にAI を利用さると差別的結果が生じるリスクを指摘する文献としては、see, e.g., Sonderling et al., infra note 12.

* 人事採用によりを使うことが不適切であることは、日本でも、例えば平野が座長を務める「AI ネットワーク社会推進会議・AI ガバナンス検討会」の第2回会合(平成30年[2018 年] 12月10日)に於ける議者のご発表に於いても既に指摘されていた。「資料1 早稲田大学 大湾先生 御発表資料: 人事データ活用への関心とガイドライン作成に向けての議論」

https://www.soummi.go.jp/main content/000589116.pdf (last visited Sept. 12, 2024). 平野百身の文献については、平野音「AI に不適合なアルゴリズム回避論: 機械的な人事採用選別と自動化パイアス」『情報通信政策研究』第7巻2号1頁(総務者,2024年3月)

https://www.soumu.go.jp/iicp/journal/journal 07-02.html (last visited Sept. 24, 2024); 「(資料1) AIの判断に対するヒトの最終決定権の限界: Human-in-the Loopの問題」in 総務省「情報通信法学研究会 令和5年度」2023年9月6日

https://www.soumu.go.jp/main.sosiki/kenkyu/hougakuken/R05 siryou.html (last visited Sept. 24, 2024).

「やり過ぎ」についての批判としては、例えば顔の表情から感情を読み取ることは 困難であると指摘されているにも関わらず、大企業やスタートアップ企業等がこれを 売り込んでヒトの一生を左右している問題が、連邦取引委員会 (FTC) 委員等による 共著論文に於いて次のように指摘されている。この指摘が日本にも当てはまることが ないことを、筆者は祈るばかりである。

A review that analyzed more than a thousand studies on emotional expression concluded that "[e]fforts to simply 'read out' people's internal states from an analysis of their facial movements alone, without considering various aspects of context, are at best incomplete and at worst entirely lack validity, no matter how sophisticated the computational algorithms."[] Nevertheless, large companies [] --plus a host of well-funded start-ups-continue to sell questionable affect recognition technology, and it is sometimes deployed to grant or deny formative life opportunities.

A striking example of the use of affect recognition is in hiring.

Rebecca Kelly Slaughter et al., Algorithm and Economic Justice: A Taxonomy of Harms and a Path Forward for the



「不適切な」AI利用 欧米に於けるADM^(*)批判と規制

(*) ADM: <u>Automated Decision-Making</u>又は <u>Algorithmic Decision-Making</u>

アメリカの学説や政策

★ 就活者の選別をAIに委ねたり依拠することは、不正確(エビデンス欠如)及び/又は差別的(不公正)なので望ましくない。

See 平野晋「AIに不適合なアルゴリズム回避論」『情報通信政策研究』7巻2号1頁 (2024年3月) https://www.jstage.jst.go.jp/article/jicp/7/2/7_1/_html/-char/ja; 平野晋「Human in the Loopの問題」 in 総務省_情報通信法学研究会, 2023 年9月6日 https://www.soumu.go.jp/main_content/000899843.pdf;「資料1早稲田大学 大湾先生 御発表資料: 人事データ活用への関心とガイドライン作成に向けての議論」

https://www.soumu.go.jp/main_content/000589116.pdf.

EUOAI ACT

- ★雇用や教育等の分野に於けるAI利用は、厳しい規制対象。
- ★職場と教育機関に於いて感情を推認するAI利用は、原則として禁止。

口不正確 口不透明 口説明無責任

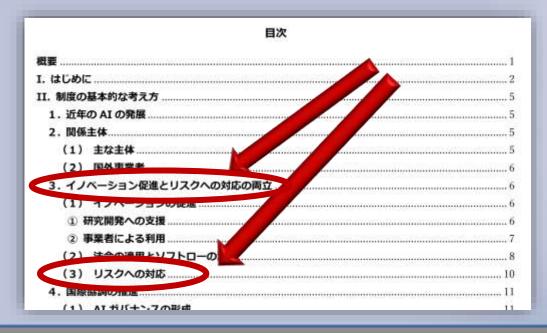
- 就活生 · 求職者に対するAI利用(濫用)例:
 - 「Jared^{ジャレッド}」という名前で高校時代に「ラクロス」 をやっていたら、仕事で成功する。
 - ゲームをやらせて仕事上の成功の資質が予測できる。
- ・実証的/科学的根拠(エヴィデンス)を欠いている。
- データマイニングで相関関係を示せても、因果関係は 証明されていない。
- peer reviewによる再現性を通じた科学的正しさが 証明されない。 See e g 平野晋 [Human in the Loop の問題 1888年 情
- ・ 営業秘密の壁。

See e. g., 平野晋「Human in the Loop の問題」総務省_情報通信法学研究会, 2023 年9月6日

https://www.soumu.go.jp/main_content/000899843.pdf (last visited Mar. 16, 2025).

内閣府「AI制度研究会中間とりまとめ(案)」

14 や用途が様々存在する中で、開発、提供、利用といった AI のライフサイクルの各場面において 15 顕在化する可能性のあるリスクとは、いかなる種類の AI モデルのどのような性質に起因するリ スクであって、誰にどのような影響を与えるものかといった要素を分析する必要がある。その前 提として、まずは AI の開発、利用等に関する実態を調査・分析し、社会全体で認識を共有した 上で必要な対応を適時適切に行うことが重要である。例えば、不適切な AI による求職者の選別 や AI による消費者の混乱に対する懸念があり、政府は実態の把握に努めるとともに必要な対策 を検討することも重要である。また、政府による実態の把握のため、各主体に協力の要請を行う



内閣府「AI制度度研究会 中間とりまとめ(案)」11頁

https://www8.cao.go.jp/cstp/ai/ ai_kenkyu/5kai/shiryou1.pdf (last visited Mar. 15, 2025).

「リスクへの対応」: AI法案

「人工知能関連技術の研究開発及び活用の推進に関する法律案」 in 内閣府「第127回 通常国会」

(調査研究等)

https://www.cao.go.jp/houan/217/index.html (last visited Mar. 15, 2025).

第十六条

国は、国内外の人工知能関連技術の研究開発及び活用の動向に関する情報の収集、不正な目的又は不適切な方法による人工知能関連技術の研究開発又は活用に伴って国民の権利利益の侵害が生じた事案の分析及びそれに基づく対策の検討その他の人工知能関連技術の研究開発及び活用の推進に資する調査及び研究を行い、その結果に基づいて、研究開発機関、活用事業者[*]その他の者に対する指導、助言、情報の提供その他の必要な措置を講ずるものとする。

(強調及び注書付加)

[*]「活用事業者」とは、「人工知能関連技術を活用した製品又はサービスの 開発又は提供をしようとする者その他の人工知能関連技術を事業活動において 活用しようとする者」をいう。(第7条) (基本理念)

第三条

• • • 0

4 人工知能関連技術の研究開発及び活用は、不正な目的又は不適切な方法で行われた場合には、犯罪への利用、個人情報の漏えい、著作権の侵害その他の国民生活の平穏及び国民の権利利益が害される事態を助長するおそれがあることに鑑み、その適正な実施を図るため、人工知能関連技術の研究開発及び活用の過程の透明性の確保その他の必要な施策が講じられなければならない。

5

(強調付加)

(国の責務)

第四条

国は、前条に定める基本理念(以下「基本理念」という。)にのっとり、人工 知能関連技術の研究開発及び活用の推進に関する施策を総合的かつ計画 的に策定し、及び実施する責務を有する。

(活用事業者の責務)

(強調付加)

第七条

人工知能関連技術を活用した製品又はサービスの開発又は提供をしようとする者その他の人工知能関連技術を事業活動において活用しようとする者(以下「活用事業者」という。)は、基本理念にのっとり、自ら積極的な人工知能関連技術の活用により事業活動の効率化及び高度化並びに新産業の創出に努めるとともに、第四条の規定に基づき国が実施する施策及び第五条の規定に基づき地方公共団体が実施する施策に協力しなければならない。

内閣府「AI制度研究会中間とりまとめ(案)」

1		
24 (2)	政府等による利用
Z*1 1	~ /	以がせにみるが几

- 25 政府が AI を利用することにより、行政サービスや業務等の質・効率を向上させることができ
- 26 るほか、政府がユースケースやその有用性を示し、具体事例、留意点等を周知することにより、
- 27 AI の活用を促進し、国内の AI 市場の発展等にも貢献できるため、政府が率先して AI を利用し
 - 1 ていくことは重要である。ただし、国民の権利利益に重大な影響を及ぼしかねないものについて
 - 2 は、AI の出力結果を自動的に採用することのリスク 11を踏まえ、慎重に取り組むべきである。
 - 3 地方自治体についても、行政サービスの大きな部分を占めており、国民生活への影響も大きい
 - 4 ため、AI の利用を推進し、行政サービスや業務等の質・効率を向上させていくことが重要である。
 - 5 また、地方自治体における AI の利用の推進にあたっては、各自治体の先進的な取組 12を含む利
 - 6 用事例を参考にするほか、各自治体の地域課題に応じた AI の利用も重要である。

11 海外におけるリスク発現事案としては、米国における失業保険に関する事例(ミシガン州の失業保険 庁が請求者の詐欺を検出するために使用した統合データ自動化システムは、2013 年から 2 年間で 93% のエラー率を記録し、2 万人の州民を詐欺で誤って告発した。)や、教師の失職につながった事例(ヒューストン独立学区は、教師が生徒の学力成長に及ぼす影響を推定する付加価値モデル(VAM)を導入し、2007 年から 2016 年までに教師の契約更新の拒否や解雇に利用したため、多くの教師が VAM の使用を止めるための司法救済を求めた。)が存在する。

12 神戸市においては、2024年3月、一定のルール下での AI の効果的かつ安全な活用を目的として AI 条例を制定するほか、文章要約、アイデア出し、プログラミングコードの生成等 AI の活用を進めている (AI 制度研究会 (第2回) 資料3参照)。

内閣府「AI 制度度研 究会 中間 とりまとめ (案)∫17 ~18頁 https://w ww8.cao.g o.jp/cstp/a i/ai kenky u/5kai/shir you1.pdf (last visited Mar. 15, 2025).

背景:日本国政府/地方自治体も AIを利活用する方針

23	
24	
25	
26	
27	
28	
29	
30	
31	
32	
33	
34	

35

具体的な制度・施策の方向性

■全般的な事項(Ⅲ.1.)

AIの研究開発・実装が最もしやすい国を目指す

● 政府の司令塔機能の強化、戦略の策定

速やかな法制度化が必要世界のモデルになるような制度

- 全体を俯瞰する司令塔機能強化
- AIの安全・安心な研究開発・活用のための戦略(基本計画)の策定
- 安全性の向上等
 - ・国による指針(広島AIプロセス準拠)の整備、事業者による協力
 - 国による調査・情報収集、事業者・国民への指導・助言、情報提供等



- 政府等による利用 (Ⅲ.2.)
 - 適正なAI政府調達・利用等

- 基盤サービス等における利用 (Ⅲ.3.)
 - 各業法等による対応 等

R. GERROSESAGERANAS AND AND REPORT OF THE PROPERTY OF THE PARTY OF THE RECORDS 2 YOUR WINE TAKES, THE MENT OF A PROPERTY OF A PROPERTY OF **国際学校に指する中のと、「子典 30年 (12月 10日 開発的できない) に集合す。セプライチュー** シ・リスタに共和することが必要であると判断されるものではいては、内臓サイバーゼキュリア etta-Biffiendristatioti. Wetsetsationical ten. Actoro. INFORCEMENTS. -AT BRIDGE A BRIDGE ASSESSMENT OF THE STREET BURNSTANDARFECTAN, CHURCHECK, BRY A SERVERSCORE なるとは、公司をした自身関係のイドライン等の指揮や総等の、以上関するから下ライン等の開発 ODJERHOCK PRETRIE. また、このメラな**あの**記載と関するカイドラインを企業的を含えるで、AI を利用する事業を対 HODOGUSKPA-9-EZERHITEROPRICUJO BERGINA ARRESECH STREETS AND AND MAKE BUT BELLEVIEW - LOSS OF SECRETARY BY BUT BUT BELLEVIEW から各種な奇なが発展となるように発展する表である。

前掲「中間とりまとめ(案)」1,17頁.



内閣府「AI制度研究会」神戸市提出資料

AI条例が想定するリスクへの対応

神戸スてートシティ

- ▶ AI条例が想定するリスクへの対応
 - 事前に安全性を完全に担保するルールづくりは現実的ではない一方で、すべてを 事後対策に委ねることも危険である
 - 行政が積極的にAIを活用していくには、事前に十分なリスクの洗い出しを行い、 リスクが顕在化した際の対処方法を検討しておくことが重要ではないか
- ○AI条例に基づくリスクアセスメントの概要
 - ・対象 市及び受託事業者が次の業務にAIを活用するとき
 - ①市民の権利利益に影響を与える行政処分 (課税、各種給付認定など)
 - (2)基本計画等の計画策定(市の基本計画など)
 - ③その他市民生活に重大な影響を与えるおそれがあるもの
 - ・方法 AI事業者ガイドラインのワークシートを参考にしつつ、神戸市の独自項目を加えた48項目のワークシートに基づき実施
 - ・特徴 リスク軽減のために、技術的な対策以上に利用者の運用面での取り組みを重視

神戸スでートシティ

https://www8.cao.go.jp/cstp/ai/ai_k enkyu/2kai/shiryou3.pdf (emphasis added). 神戸市におけるAI活用のためのルール整備

2024.8.23

神戸市企画調整局デジタル戦略部

口不正確 口不透明 口説明無責任

- アルゴリズムやソフトウエアに依拠した不利益処分等が厳しく批判されたアメリカの先例に学ぼう
 - 1. Houston Fed. of Teachers v. Houston
 Independent, 251 F. Supp. 1168 (S.D. Tex.
 2017). アルゴリズムの評価に従って教員を解雇した際の説明無責任と不透明性が違憲とされた事例。
 - 2. <u>Cahoo v. SAS Analytics Inc.</u>, 912 F.3d 887 (6th Cir. 2019). 失業保険受給申請の詐欺を自動的に決定等するシステムの不正確性等が問題になった事例。
- ・ 自動的/アルゴリズム的意思決定(ADM)の問題/行政処分・不利益処分の問題

☑ 不透明 ☑ 説明無責任 [☑不正確]

Houston Fed. of Teachers, Local 2415 et al. <u>対</u>

<u>Houston Independent School District事件</u>

https://www.housto npublicmedia.org/a rticles/news/2017/1 0/10/241724/feder al-lawsuit-settledbetween-houstonsteacher-union-andhisd/ (last visited Mar. 16, 2025).





Houston Fed.事件の概要

- ヒューストン市の教員の評価を、EVAASと呼ばれるアルゴリズムに算出させ、そこで最低評価になった教員が解雇や契約非継続となった。教員組合と解雇等された教員達(π)がEVAAS使用の恒久的な差止等を求めて、ヒューストン独立教区(△)に対する訴えを提起。
- EVAASは訴外ベンダのSAS社が供給していて、そのソースコード や統計的手法等の情報は営業秘密であるとして、△にもπにも開 示せず、いわば「ブラックボックス」であった。しかし、ソースコード等 の情報が開示されなければ、評価結果を再現できず、その誤りの 有無をπが検証できないことが問題になった。
- 裁判所は、<u>再現できなければ、手続的な適正手続保障違反に当た</u> <u>る</u>と指摘。
 - → SH: 求職者を選別するAIに対する批判と一致!
 - → SH: peer reviewに服して正確性が再現できなければ、怪しいという考え方。

☑ 不透明 ☑ 説明無責任 [☑不正確]

会」2023年9月6日 9 8 3.pdf @ 総務坐

//www.soumu.go.Jp/

Human in the Loopの問題」 4頁 main_content/

営業秘密 v. 透明性

	()
	HOUSTON FEDERATION OF TEACHERS, LOCAL 2415,
	et al., Plaintiffs,
	8
H	DUSTON INDEPENDENT SCHOOL
	DISTRICT, Defendant.
	CIVIL ACTION H-14-1189
	United States District Court,
	S.D. Texas, Houston Division.
	Signed 05/04/2017

VALUE-ADDED RATING	EVAAS @TGI	RELATIONSHIP TO EXPECTED AVERAGE GROWTH	
Well above	Equal to or greater than 2	Students on average substantially exceeded expected average growth	
Above	Equal to or greater than 1 but less than 2	t Students on average exceeded average growth	
No detectable difference	Equal to or greater than -1 but less than 1		
Below	Equal to or greater than -2 but less than -1	Students on average fell short of average growth	
Well below	Less than -2	Students on average fell substantially short of expected average growth	

[W]ithout access to [the vender's] proprietary information—the . . . computer source codes, decision rules, and assumptions-EVAAS scores will remain a mysterious 'black box,' imperative to challenge.

HISD teachers have no meaningful way to ensure correct calculation of their EVAAS scores, and as a result are unfairly subject to mistaken deprivation of constitutionally protected property interests in their jobs.

Houston Fed'n of Tchrs., Loc. 2415 v. Houston Indep. Sch. Dist., 251 F. Supp. 3d 1168, 1179 (S.D. Tex. 2017) (emphasis added).

[夕不正確性]

"Algorithms are <u>human</u> creations, and <u>subject to</u> error like any other human endeavor."

Houston Fed. of Teachers v.
Houston Independent School
District, 251 F.Supp.3d 1168, 1177
(S.D.Tex. 2017) (emphasis added).

☑ 不正確性 ☑ 説明無責任

<u>Cahoo et al.</u> <u>対</u> SAS Analytics Inc., et al.事件



Cahoo事件の概要

- ミシガン州失業保険庁(UIA: <u>Unemployment Insurance Agency</u>)の採用した失業保険受給者の詐欺を自動的に判断するMiDAS (<u>Michigan Integrated Automated System</u>)が、大量(93%)の誤判断・偽陽性(faulse positive)を下した (robo-adjudicated)。
- その被害者であるの原告達(π)が、MiDASを監督・使用していた同省職員達被告(Δ)に対して提訴。
- 訴状では、△が合衆国憲法上の適正手続保障に違反していて、 〈資格に基づく免責〉に値しない旨の主張が十分記載されているとして、原審の判断を支持した控訴裁判所による抗告審が、 本件。

国際規範との整合性



人工知能関連技術の研究開発及び活用の推進に関する 法律案(AI法案)

	目的	国民生活の向上、国民経済の発展		
	基本理念	経済社会及び安全保障上重要 → 研究開発力の保持、国際競争力の向上基礎研究から活用まで総合的・計画的に推進 適正な研究開発・活用のため透明性の確保等 国際協力において主導的役割		
	AI戦略本部	本部長:内閣総理大臣 構成員:全閣僚 関係行政機関等に対して必要な協力を求める		
	AI基本計画	研究開発・活用の推進のために政府が実施すべき施策の基本的な方針等		
法案の 概要	基本的施策	研究開発の推進、施設等の整備・共用の促進 人材確保 教育振興 国際的な規範等定人の参画 適正性のための国際規範に即した指針の整備 情報収集、権利利益で受害する事業の分析・対策検討、調査 事業者・国民への指導・助言・情報提供		
	責務	国、地方公共団体、研究開発機関、事業者、国民の責務 関係者間の連携強化 事業者は国等の施策に協力しなければならない		
	附則	見直し規定(必要な場合は所要の措置)		

内閣府「第127回国 会」「人工知能関連技 術の研究開発及び活 用の推進に関する法 律案」「概要」

https://www.cao.go.jp/ houan/217/index.html (last visited Mar. 14, 2025).



AI法案:

国際規範に即したガイドライン整備を国の義務化

(適正性の確保)

第十三条

国は、人工知能関連技術の研究開発及び活用の<u>適正な</u> 実施を図るため、<u>国際的な規範の趣旨に即した指針の整</u>

備その他の必要な施策を講ずるものとする。

(強調付加)

内閣府「AI制度研究会」の冒頭提出講師資料





透明性と 説明可能性

- 確実に人々が、
- AIシステムと関わっていることを知ることができ、かつ
- その結果に異議を申し立てることができるようにする為の、
- ・ 透明性と責任ある開示に関する原則

OECD AI原則 1.3

https://oecd.ai/en/ai-principles (last visited July 27, 2024) (拙訳・強調付加).

内閣府「AI 戦略会議 (第11回)• AI制度研 究会(第1 回)※合同 開催1 https://w ww8.cao.g o.jp/cstp/a i/ai senrya ku/11kai/1 1kai.html (last visited Mar. 15, 2025).

「OECD AI原則」(原則1.3; 2024 rev.)



Transparency and explainability (Principle 1.3)







This principle is about transparency and responsible disclosure around Al systems to ensure that people understand when they are engaging with them and can challenge outcomes.



Al Actors should commit to transparency and responsible disclosure regarding Al systems. To this end, they should provide meaningful information, appropriate to the context, and consistent with the state of art:

- > to foster a general understanding of Al systems, including their capabilities and limitations,
- > to make stakeholders aware of their interactions with AI systems, including in the workplace,
- > where feasible and useful, to provide plain and easy-to-understand information on the sources of data/input, factors, processes and/or logic that led to the prediction, content, recommendation or decision, to enable those affected by an Al system to understand the output, and,
- > to provide information that enable those adversely affected by an AI system to challenge its output.



Rationale for this principle

The term transparency carries multiple meanings. In the context of this Principle, the focus is first on disclosing when AI is being used (in a prediction, recommendation or decision, or that the user is interacting directly with an Alpowered agent, such as a chatbot). Disclosure should be made with proportion to the importance of the interaction. The growing ubiquity of Al applications may influence the desirability, effectiveness or feasibility of disclosure in some cases.

Transparency further means enabling people to understand how an Al system is developed, trained, operates, and deployed in the relevant application domain, so that consumers, for example, can make more informed choices. Transparency also refers to the ability to provide meaningful information and clarity about what information is provided and why. Thus transparency does not in general extend to the disclosure of the source or other proprietary code or sharing of proprietary datasets, all of which may be too technically complex to be feasible or useful to understanding an outcome. Source code and datasets may also be subject to intellectual property, including trade secrets.

Explainability means enabling people affected by the outcome of an AI system to understand how it was arrived at. This entails providing easy-to-understand information to people affected by an AI system's outcome that can enable those adversely affected to challenge the outcome, notably - to the extent practicable the factors and logic that led to an outcome. Notwithstanding, explainability can be achieved in different wave depending on the context (such as the significance

Therefore, when Al actors provide an explanation of an outcome, they may consider providing - in clear and simple terms, and as appropriate to the context the main factors in a decision, the determinant factors, the data, logic or algorithm behind the specific outcome, or explaining why similar-looking rircumstances. generated a different outcome. This should be done in a way that allows individuals to understand and challenge the outcome while respecting personal data protection obligations, if relevant.

AI諸原則 / AIガイドライン: 日本が国際標準を主導

7

G20 2019.6.

・デジタル経済大臣会合



OECD 2018.9. ~ 2019.5.

- 理事会勧告 2019.5
- AI専門家会合~2019.2



内閣府 2018.4. ~ 2019.3.

• 人間中心のAI社会原則[検討]会議



総務省 2016.1. ~ 現在

- •AIネットワーク社会推進会議 2016.10.~現在
 - ・OECDに於けるカンフェレンスを総務省と共催/参加。^{♀艸○ECD}
 - AIネットワーク社会推進フォーラムに於ける国際シンポを総務省主催/参加。
 - カーネギー平和財団のカンフェレンスを日本大使館後援/総務省参加。
 - OFCD技術予測フォーラムに於ける発表。◎》○ECD
 - ✓ 高松 香川G7情報通信大臣会合(高市早苗総務大臣)
- ・AIネットワーク化検討会議 ~2016.6.

Alにおける国際的な議論への貢献

G 7情報通信大臣会合(高松、2016年4月)

・高市総務大臣(当時)からの提案※:

"G7各国が中心となり、OECD等国際機関の協力も得て、AIネットワーク化が社会・経済に与える影響や、AI開発原則の策定等AIネットワーク化をめぐる社会的・経済的・倫理的・法的課題に関し、産学民官の関係ステークホルダーの参画を得て、国際的な議論を進める"

⇒ 参加各国からの賛同を得る

Ministers' Meeting

E Friday, April 29 - 30, 2016

2-1. Suspect. Takamatos, Espaya, 260-6019, Ispan

Kagawa International Conference Hall, Tak

※提案に先立ち、叩き台として、8項目のAI開発原則を配付。

Proposal of Discussion toward Formulation of Al R&D Guideline

Referring OECD guidelines governing privacy, security, and so on, it is necessary to begin discussions and considerations toward formulating an international guideline consisting of principles governing R&D of Al to be networked ("Al R&D Guideline") as framework taken into account of in R&D of Al to be networked.

Proposed Principles in "AI R&D Guideline"

1. Principle of Transparency

Ensuring the abilities to explain and verify the behaviors of the AI network system

2. Principle of User Assistance

Giving consideration so that the AI network system can assist users and appropriately provide users with opportunities to take choices

3. Principle of Controllability

Ensuring controllability of the AI network system by humans

4. Principle of Security

Ensuring the robustness and dependability of the AI network system

5. Principle of Safety

Giving consideration so that the AI network system will not cause danger to the lives/bodies of users and third parties

6. Principle of Privacy

Giving consideration so that the AI network system will not infringe the privacy of users and third parties

7. Principle of Ethics

Respecting human dignity and individuals' autonomy in conducting research and development of AI to be networked

8. Principle of Accountability

Accomplishing accountability to related stakeholders such as users by researchers/developers of AI to be networked

講師が近年着目するリスク (cont'd)

- [「不適切な」AI利用]
- ・AIへの過度な依存
 - ✓ 内閣府「人間中心のAI社会原則」
 - ✓ 平野『AIとヒトの共生にむけて』

「人間中心のAI社会原則」

...。AI が活用される社会において、人々が AI に過度に依存したり、AI を悪用して人の意思決定を操作したりすることのないよう、我々は、リテラシー教育や適正な利用の促進などのための適切な仕組みを導入することが望ましい。

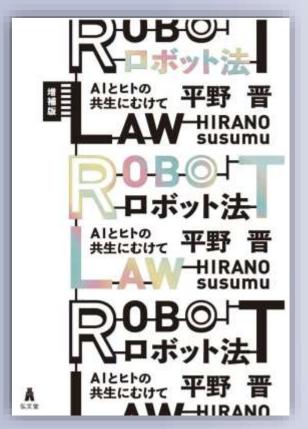
(強調付加)

内閣府「人間中心のAI社会原則」8頁(2019年3月29日)内の「1)人間中心の原則」 https://www8.cao.go.jp/cstp/ai/aigensoku.pdf (last visited Mar. 15, 2025).

「AIとヒトの共生に向けて」



平野晋 『ロボット法: AIとヒトの 共生にむけて(増補第2版)』(弘 文堂, <u>2024年</u>5月)



同『ロボット法: AIとヒトの 共生にむけて』(弘文堂, 2019年10月)



同『ロボット法: Alとヒト の共生にむけて』(弘文 堂, <u>2017年</u>11月)



背景:「AI^{アイ}ちゃんがそう言ったから」 が通用しない場合もある(?!)

AIの方がヒトよりも正しいという誤解(偏見)

- * AIは客観的だから。
- * AIは数値に基づいているから。
- * AIは偏見に左右されないから。

人間中心の/信頼に値するAl

- ・ 人間疎外のAI
- ・ 信頼に値しないAI

- · 人間中心のAI
- ・信頼に値するAI

- 制御不能
- 予見不能
- 不透明
- 説明無責任
- 不正確
- 一不公正

but not limited to, including

but not limited to,

予見可能 透明

説明責任

制御可能

- 正確
- 公正
- 適合



不適合

- 文脈が読めない。大局を理解しない。
 - · 紛争地域に於いてAK-47を持った老人は戦闘員か?
 - 紛争地域に於けるデュアル・ユースなピックアップ・トラックは非戦闘用か?(*1)
- データ化されていない情報は使えない。(*2)
- ・ 裁量が必要な判断には使えない(次葉)。
 - 「制限速度60km/h」⇔「安全運転をする義務」
 - プロ・テニスに於ける「Hawk Eye」利用 (*3)

SH: 50-50を達成した大谷翔平選手の次の打者に対して 15秒以内に投球しなかったマイケル・バウマン投手に 対して、ピッチ・クロックを宣言しなかった審判の判断。

^(*1) Charles P. IV Trumbull, <u>Autonomous Weapons</u>: <u>How Existing Law Can Regulate Future Weapons</u>, 34 EMORY INT'L L. REV. 533 (2020).

^(*2) 大湾先生 御発表資料, 前掲, at page 17.

^(*3) Ehan Lowens, Note, <u>Accuracy Is Not Enough: The Task Mismatch Explanation of Algorithm Aversion and Its Policy Implications</u>, 39 HARV. J. L. & TECH. 259 (2020).



☑ 不適合

平野「Human in the Loopの問題」39頁

採用AIに於けるHITLの欠点に対する対策方針

・ 更にAIが不得手な判断の場合とは、事前の rulesが適している場合ではなく、 standardsの当てはめによる個別的判断 が適している場合。

Green, supra, at 12-13; See also Margot E. Kaminski, Binary Governance: Lessons from the GDPR'S Approach to Algorithmic Accountability, 92 S. Cal. L. Rev. 1529, 1542-43 (2019).

- 〈法的安定性〉v. 〈個別具体的妥当性〉
- common law v. equity
- そもそも採用活動は、後者に該当する のではないか。

Ben Green, The Flaws of Policies Requiring Human Oversight of Government Algorithms, 45 COMPUT. L. SEC. REV. 1, 12-13 (2022).



法的安定性 v. 具体的妥当性



See リチャード・A・ポズナー著/平野晋監訳 『法と文学(上)』 206~07頁(2011年).



Thank you !!! ;-)

INFORMATION TECHNOLOGY & LAW ICHIGAYA TAMACHI LINK