

The Use of Algorithms in Decision-Making: An Ethical Conundrum

Amber Francisca Ling Yi Rong



Abstract

This essay considers the ethical complexities of using algorithms in decision-making processes in the justice system in the UK, with a focus on the principle of explainability. While algorithmic decision-making tools offer improved efficiency and consistency, their implementation within justice systems raises serious accountability, transparency, and fairness issues. The paper argues that it is fundamentally challenging to establish and measure standards of explainability because of stakeholder's varying requirements, the technical obscurity of the “black box” system, and legal issues related to proprietary rights. It also considers the ethical challenge of interpretability and completeness through case studies of COMPAS and Loomis. The essay concludes that although complete explainability remains elusive, greater inclusivity in AI design, regulatory responsibility, and user autonomy can promote ethical governance and enhance public trust in algorithmic decision-making within the justice system.

Introduction

Algorithms have become an integral part of modern life, and there appears to be a growing interest in their use within the UK justice system. Lord Chancellor Shabana Mahmood has expressed support for utilising these systems to address the mounting backlog of court cases, which has reached a record-high level of 73,000 in December 2024 from 38,000 in December 2019.¹ While the promise of enhanced efficiency is appealing, studies have also highlighted the critical importance of receiving explainability from these algorithms. This matter is reflected in the European Ethical Charter, yet the charter provides little clarity on how the standard of explainability should be met.² This paper thus contends that explainability is an important consideration in algorithmic decision-making in the justice system and explores the

¹ GOV.UK, ‘Courts reform to see quicker justice for victims and keeps streets safe’ (GOV.UK, 12 December 2024) <<https://www.gov.uk/government/news/courts-reform-to-see-quicker-justice-for-victims-and-keepsstreets-safe>> 16 December 2024.

² The European Commission for the Efficiency of Justice, ‘European Ethical Charter on the use of artificial intelligence (AI) in judicial systems and their environment’ (December 2018) <<https://rm.coe.int/ethical-charteren-for-publication-4-december-2018/16808f699c>> accessed 13 December 2024.

technical, practical and ethical challenges that make defining and achieving explainability particularly difficult.

The Importance of Explainability in Algorithmic Decision-Making

Explainability, defined as a human-interpretable description of how a decision-maker arrived at a particular decision after considering a set of facts,³ has become a central focus in ethical and policy discussions worldwide. This is distinct from transparency where transparency answers “what happened” by making users aware they are dealing with AI, while explainability clarifies how the system reached its decision.⁴ Jacob Turner notes that transparency is a recurring theme across various ethical codes in different areas of expertise,

reflecting widespread concerns about the risks of opaque AI systems.⁵ However, this need for explainability is even more pronounced in the context of justice systems, as the ICO Report reveals that jurors overwhelmingly prioritise explanations in criminal justice scenarios over every other scenarios, such as healthcare.⁶ This preference stems from the value explanations provide in challenging decisions, building trust, ensuring fairness and to prove the absence of

³ Doshi-Velez and others, ‘Accountability of AI Under the Law: The Role of Explanation’ (2017) Berkman Center Research Publication, 4.

⁴ Arsen Kourianian and Mayer Brown, ‘Addressing Transparency & Explainability When Using AI Under Global Standards’ (Bloomberg Law 2024) <<https://www.mayerbrown.com-/media/files/perspectivesevents/publications/2024/01/addressing-transparency-and-explainability-when-using-ai-under-globalstandards.pdf?3Frev=8f001eca513240968f1aea81b4516757#:~:text=Global%20AI%20standards%20often%20group,decision%20was%20made%20using%20AI>> accessed 18 August 2025.

⁵ Jacob Turner, *Robot Rules: Regulating Artificial Intelligence* (Palgrave Macmillan Cham 2018)⁶
ICO, ‘Project Explain interim report’ (June 2019) <<https://ico.org.uk/media/2615039/project-explain20190603.pdf>> accessed 12 December 2024.

bias.⁶ Burrell reinforces this by stating that explanations support individual autonomy, allowing individuals to contest decisions and maintain agency of their treatment.⁷

The issue of “machine bias” and inaccuracy further underscores the need for explainability in algorithmic decision-making in the justice system. This is problematic as these tools make morally consequential and socially significant decisions.⁸ One notable example is the COMPAS algorithm, which calculates the likelihood of recidivism of prisoners for their release. To illustrate this issue, ProPublica’s research revealed that COMPAS falsely flags Black defendants as future reoffenders at nearly twice the rate as White defendants (45% compared to 23%).¹⁰ Dartmouth’s finding also adds that COMPAS is no more accurate than predictions made by non-experts in the justice system.⁹ These cast significant doubt on the reliability of such systems. Richard Susskind¹⁰ emphasises that functional AI must be transparent and capable of heuristic reasoning, which means systems should be open to challenge and improvement. This reinforces why providing explanations is essential to understanding, improving and ensuring the reliability of decision-making tools, while also enabling checks and balances to prevent bias and discrimination.¹¹

⁶ *ibid.*

⁷ Jenna Burrell, ‘How the machine ‘thinks’: Understanding opacity in machine learning algorithms’ (2016) <<https://dx.doi.org/10.2139/ssrn.2660674>> accessed 12 December 2024.

⁸ Joel Walmsley, ‘Artificial intelligence and the value of transparency’ (2021) 36 *AI & Soc* 585, 588. ¹⁰ Julia Angwin and others, ‘Machine Bias’ (*ProPublica*, 23 May 2016) <<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>> accessed 3 December 2024.

⁹ Julia Dressel and Hany Farid, ‘The Accuracy, Fairness, and Limits of Predicting Recidivism’ (2018) 4(1) *Science Advances* eaao5580.

¹⁰ Richard Susskind, ‘Expert Systems in Law: A Jurisprudential Approach to Artificial Intelligence and Legal Reasoning’ (1986) 49 *MLR* 168, 173.

¹¹ The Royal Society, ‘Explainable AI: the basics’ (November 2019) <https://royalsociety.org/newsresources/projects/explainable-ai/?utm_source=report&utm_medium=print&utm_campaign=aiinterpretability> accessed 10 December 2024, 9-10.

Lastly, taking a broader view, Lisa Webley underscores that explainability is essential for upholding the core principles of natural justice.¹² These principles demand that decisions are to be made impartially, with fair notice and an opportunity to respond.¹³ Explainability would therefore allow individuals to understand the rationale behind decisions and evaluate the system's legality and robustness.¹⁴ Explainability is also especially important for legal professionals like lawyers and judges, whose actions influence not only individual clients but also public trust in the entire legal system.¹⁵ This is because lawyers uphold the rule of law and must act transparently and responsibly to maintain confidence in legal processes.¹⁶ Any breach of professional standards risks eroding public trust, triggering a chain reaction that could potentially destabilise both the legal profession and the justice system, and ultimately undermining the rule of law and societal confidence in the system.¹⁷

The Uncertainties in Defining 'Explainability' Standards

Despite these needs, there remains significant uncertainty about how we should define the standards for explainability due to the varied needs of different stakeholders. This is because different users require different forms of explanation tailored to their specific contexts.¹⁸ As Doshi posits, an explanation must provide the correct type of information in

¹² Legal Services Board Podcast, 'Ethics, Technology and Regulation' (21 May 2020) <<https://legalservicesboard.podbean.com/e/ethics-technology-and-regulation>> accessed 14 December 2024.

¹³ Committee for Privileges and Conduct, Third Report of Session 2017–19, Further report on the conduct of Lord Lester of Herne Hill, HL 252.

¹⁴ Legal Services Board Podcast, 'Ethics, Technology and Regulation' (21 May 2020) <<https://legalservicesboard.podbean.com/e/ethics-technology-and-regulation>> accessed 14 December 2024.

¹⁵ ibid.

¹⁶ ibid.

¹⁷ ibid.

¹⁸ The Royal Society, 'Explainable AI: the basics' (November 2019) <https://royalsociety.org/newsresources/projects/explainable-ai/?utm_source=report&utm_medium=print&utm_campaign=aiinterpretability/> accessed 10 December 2024.

order for it to be useful.¹⁹ A key principal for this is that they must enable the human observer to assess how much a particular input influenced the output.²⁰ In other words, explanations should present digestible information tailored to their specific audience. This shows that the explanations must be highly context-dependent, as it would only be meaningful if the AI-based decision can be explained at an individual level,²¹ highlighting the complication in defining ‘explainability’ standards. Thus, for AI developers, the critical question is not merely whether a system is explainable or whether one model is more explainable than the other.²² Instead, it is whether the system can deliver the specific type of explainability needed for a particular task or user group.²³

For example, independent auditors or system developers would require the publishing of algorithms to know the technical and transparent details about how the system functions.²⁴ However, such explanations are not suitable for laypersons, who lack the expertise to interpret technical details.²⁵ Instead, counterfactual reasonings, which detail the specific factors leading to an individual output and explore the what if scenarios, what would need to change for AI to make a different outcome, are more helpful to help them understand and contest decisions.²⁶ Regulators, by contrast, might require explanations on data processing to ensure that the system operate within the bounds of established regulations.²⁷ This highlights that explainability cannot follow a “one-size-fits-all” approach,²⁸ which contributes to the

¹⁹ Doshi-Velez and others, ‘Accountability of AI Under the Law: The Role of Explanation’ (2017) Berkman Center Research Publication, 4.

²⁰ *ibid.*

²¹ Hans de Brujin, Martijn Warnier and Mrijn Janssen, ‘The perils and pitfalls of explainable AI: Strategies for explaining algorithmic decision-making (2022) 39 Government Information Quarterly 101666, 101669.

²² The Royal Society (n 19) 14.

²³ *ibid.*

²⁴ *ibid.*

²⁵ Jenna Burrell, ‘How the machine ‘thinks’: Understanding opacity in machine learning algorithms’ (2016) <<https://dx.doi.org/10.2139/ssrn.2660674>> accessed 12 December 2024. ²⁸ The Royal Society (n 19) 13-14.

²⁶ *ibid.* 19.

²⁷ *ibid.*

challenge in defining its standards. As such, a user-centric mindset must be adopted.²⁷ This would consider the specific context and needs of each stakeholder to ensure that explanations are both relevant and meaningful.

Uncertainty in Achieving ‘Explainability’ Standards

Technical challenges

Achieving explainability in algorithmic decision-making fundamentally requires transparency, an ability to see the inner workings of the algorithm. However, this is hindered by the technical “black box” nature of AI systems.²⁸ This opacity arises from the complexity of machine learning models, particularly those involving techniques like deep learning.²⁹ This is because these systems do not follow pre-programmed rules.³⁰ Instead, they learn relationships and patterns through a layered structure, developing their own decisional rules.³¹ These rules are typically unintelligible to humans, as they lack alignment with human concepts or symbolic reasoning. Therefore, they can even be too complicated for expert users³² and programmers³³ to understand.

A proposed solution to mitigate is the use of post-hoc explanations³⁴ or proxy models³⁵ which analyses and interpret decision-making process after decisions are made. However, even

²⁷ Margaret Hagan, *Law by Design* (2021).

²⁸ Han-Wei Liu, Ching-Fu Lin and Yu-Jie Chen, ‘Beyond *State v Loomis*: artificial intelligence, government algorithmization and accountability’ (2019) 27 IJLIT 122, 135.

²⁹ *ibid.*

³⁰ *ibid.*

³¹ *ibid.*

³² The Royal Society (n 19) 8.

³³ Han-Wei Liu, Ching-Fu Lin and Yu-Jie Chen, ‘Beyond *State v Loomis*: artificial intelligence, government algorithmization and accountability’ (2019) 27 IJLIT 122, 135.

³⁴ Hans de Bruijn, Martijn Warnier and Marijn Janssen, ‘The perils and pitfalls of explainable AI: Strategies for explaining algorithmic decision-making’ (2022) 39(2) Government Information Quarterly <<https://doi.org/10.1016/j.giq.2021.101666>> accessed 19 December 2024.

³⁵ The Royal Society (n 19) 13.

such methods are not without challenges. One example is SHAP, a widely used method for estimating feature importance in AI systems.³⁶ However, research shows that SHAP's outputs can diverge from exact computations, especially when models involve more than a handful of features.³⁷ This highlights that SHAP is better at approximating relative importance, rather than providing precise explanations. Such limitations reinforce the Royal Society's concern that post-hoc explanations may not always be accurate or faithful and could even mislead users.

³⁸ Similarly, Lipton³⁹ warned against the embracing interpretability that is optimised to satisfy subjective demands, since it can lead to explanations that are plausible but misleading. Illustrating this, Lipton drew parallels with human behaviour, where subjective rationales often conceal biases in processes like hiring and college admissions. To conclude, the technical black box remains a barrier to meaningful explainability, as alternative measures like post-hoc explanations may only provide partial transparency and risk misleading users.

Proprietary barriers

Balancing explainability with the protection of private interests presents another challenge in achieving explainability. This is also known as the legal “black box”,⁴⁰ where proprietary protections such as trade secrets and intellectual property laws restrict access to the inner workings of AI algorithms.⁴¹ This opacity often exists to maintain competitive advantage and stay ahead of adversaries.⁴² The Council of Europe in European Ethical Charter

³⁶ XuanXiang Huang and Joao Marques-Silva, ‘On the failings of Shapley values for explainability’ (2024) 171 <<https://doi.org/10.1016/j.ijar.2023.109112>> 18 August 2025.

³⁷ Ibid.

³⁸ The Royal Society (n 19) 13.

³⁹ Zachary Lipton, ‘The Mythos of Model Interpretability: In machine learning, the concept of interpretability is both important and slippery.’ (2018) 16(3) Queue 31, 42.

⁴⁰ Han-Wei Liu, Ching-Fu Lin and Yu-Jie Chen, ‘Beyond *State v Loomis*: artificial intelligence, government algorithmization and accountability’ (2019) 27 IJLIT 122, 135.

⁴¹ ibid.

⁴² Jenna Burrell, ‘How the machine ‘thinks’: Understanding opacity in machine learning algorithms’ (2016) <<https://dx.doi.org/10.2139/ssrn.2660674>> accessed 12 December 2024.

acknowledges this as a substantial barrier to achieving full transparency.⁴³ They posit that trade secret laws, in particular, limit access to the source codes of proprietary software, complicating efforts to ensure accountability and fairness.⁴⁴

This obstacle is starkly illustrated in the *Loomis*⁴⁵ case, where the defendant contests the lack of transparency in the COMPAS risk assessment system. Despite these concerns, the Court upheld the proprietary nature of COMPAS, preventing disclosure of how risk factors were weighted or how risk scores were calculated. Similarly, Northpointe, the company behind the COMPAS, also refused to disclose its calculation methods to ProPublica, which criticised them for bias and inaccuracies in predicting recidivism among Black defendants.⁴⁶ A similar concern has been raised in the UK. The legal reform organisation JUSTICE, in its review of algorithmic tools in the justice system,⁴⁷ warned that biased and opaque systems such as COMPAS and Geolitica risk becoming an unlawful and arbitrary exercise of power. These case studies underscore the practical difficulties in meeting explainability standards, as corporate secrecy obstruct access to essential information necessary for true transparency and by extension, meaningful explainability.

⁴³ The European Commission for the Efficiency of Justice, ‘European Ethical Charter on the use of artificial intelligence (AI) in judicial systems and their environment’ (December 2018) <<https://rm.coe.int/ethical-charteren-for-publication-4-december-2018/16808f699c>> accessed 13 December 2024.

⁴⁴ *ibid.*

⁴⁵ *State v Loomis* 881 NW2d 749 (Wis 2016) 754 (US).

⁴⁶ Julia Angwin and others, ‘Machine Bias’ (ProPublica, 23 May 2016) <<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>> accessed 3 December 2024.

⁴⁷ JUSTICE, ‘AI in our Justice System’ (2025) 40. <<https://files.justice.org.uk/wp-content/uploads/2025/01/29201845/AI-in-our-Justice-System-final-report.pdf>>

Ethical Dilemmas

Another obstacle in ensuring that context-based measures reliably uphold explainability standards is the ethical challenge of balancing the conflicting demands of completeness and interpretability. Interpretability seeks to make AI systems comprehensible to humans by simplifying their operations into concepts or terms that align with human knowledge.⁴⁸ Completeness, by contrast, aims to offer a full and accurate depiction of a system's functionalities.⁴⁹ As Gilpin posits, attaining both at the same time is fundamentally challenging because the most accurate explanations are difficult for people to understand; and conversely, the most interpretable descriptions often lack predictive power.⁵⁰ The most accurate explanations are often too complex to understand, while the most interpretable ones risks producing persuasive rather than transparent explanations.⁵¹ As Herman points out, this tension creates grave ethical dilemmas and poses two questions: "*When is it unethical to manipulate an explanation to better persuade users? And how do we balance concerns of transparency and ethics with our desirability for interpretability?*"⁵⁶ These questions highlight the dangers of oversimplified explanations that prioritise persuasion over credibility, which can lead to deceptive explanations that erode trust and accountability. To conclude, the balancing of completeness and interpretability presents a challenge in achieving useful explainability.

⁴⁸ Leilani Gilpin and others, 'Explaining Explanations: An Overview of Interpretability of Machine Learning' (2018) Massachusetts Institute of Technology Cambridge.

⁴⁹ ibid.

⁵⁰ ibid.

⁵¹ ibid.

⁵⁶ ibid.

Furthermore, the possibility of misguiding users is another ethical dilemma. Wachter⁵² notes that counterfactual explanations can reveal when an algorithmic decision was influenced by protected characteristics such as gender or race, exposing potential discrimination. However, these explanations also carry inherent limitations.⁵³ Defending his stance, Wachter explains that counterfactual reasonings merely describe dependencies between a decision and external factors, but not necessarily causal links.⁵⁴ This means that while a counterfactual explanation might show that altering a certain factor leads to a different conclusion, it does not disclose the underlying causal connection or the algorithm's internal logic. This limitation arises because counterfactual explanations only lay out specific dependencies, and not the full set of circumstances influencing the decision. Wachter illustrates this by showing that, for example, if a counterfactual reasoning demonstrates that a "Black" individual's race influenced a decision, it does not imply that the race of "White" individuals was treated as irrelevant.⁵⁵ Thus, this could lead to a superficial understanding of why a decision was made and therefore may mislead user's intuitions. In conclusion, this reveals the limits of counterfactual explanations as tools for achieving explainability as it only provides partial insights into the algorithm's behaviour.

The Way Forward

Addressing these challenges of explainability in algorithmic decision-making is a complex and ongoing issue. However, a promising temporary solution lies in fostering inclusivity in the

⁵² Sandra Wachter, Brent Mittelstadt and Chris Russel, 'Counterfactual Explanations Without Opening the Black Box: Automated Decisions and the GDPR' (2018) 31(2) Harvard Journal of Law & Technology 841, 853854.

⁵³ Sandra Wachter, Brent Mittelstadt and Chris Russel, 'Counterfactual Explanations Without Opening the Black Box: Automated Decisions and the GDPR' (2018) 31(2) Harvard Journal of Law & Technology 841, 853854.

⁵⁴ *ibid.*

⁵⁵ *ibid.*

process of designing these tools. Both Turner⁵⁶ and Al Olama⁶² highlight the importance of involving diverse voices from governments and citizens through public consultations and online questionnaires, ensuring AI's societal impact is shaped collectively and not controlled by a technocratic elite. By involving individuals in these processes, trust in these systems can be cultivated, leading to greater confidence of their use. Building on this, decision-makers should also bear the responsibility of explaining AI-influenced outcomes, which would incentivise greater scrutiny and critical evaluation of these tools.⁶³ Moreover, given the current limitations in achieving full explainability, implementing an “opt-in or optout” mechanism offers a practical interim solution which allows users the autonomy to engage with such systems.⁵⁷ This approach also aligns with the European Ethical Charter, which emphasises the importance of informing in automated decision-making processes.⁵⁸ Together, these approaches can foster public trust in algorithmic decision-making and lay the groundwork for more sustainable and ethical use.⁵⁹

Conclusion

In conclusion, the integration of algorithms into the justice system presents significant opportunities for enhanced efficiency, particularly in addressing the rising backlog cases. However, the above discussion underscores the critical importance of explainability in these

⁵⁶ Jacob Turner, *Robot Rules: Regulating Artificial Intelligence* (Palgrave Macmillan Cham 2018) 264-272. ⁶² Dom Galeon, ‘An Inside Look at the First Nation With a State Minister for Artificial Intelligence’ (*Futurism*, 12 November 2017) <<https://futurism.com/uae-minister-artificial-intelligence>> accessed 19 December 2024.

⁶³ Hans de Bruijn, Martijn Warnier and Marijn Janssen, ‘The perils and pitfalls of explainable AI: Strategies for explaining algorithmic decision-making’ (2022) 39(2) *Government Information Quarterly* <<https://doi.org/10.1016/j.giq.2021.101666>> accessed 19 December 2024.

⁵⁷ Han-Wei Liu, Ching-Fu Lin and Yu-Jie Chen, ‘Beyond *State v Loomis*: artificial intelligence, government algorithmization and accountability’ (2019) 27(2) *IJLIT* 122, 140.

⁵⁸ The European Commission for the Efficiency of Justice, ‘European Ethical Charter on the use of artificial intelligence (AI) in judicial systems and their environment’ (December 2018) <<https://rm.coe.int/ethical-charteren-for-publication-4-december-2018/16808f699c>> accessed 13 December 2024.

⁵⁹ Jacob Turner, *Robot Rules: Regulating Artificial Intelligence* (Palgrave Macmillan Cham 2018) 278.

systems as highlighted by ethical and policy discussions. The challenges of defining and achieving explainability standards are multifaceted, involving technical, proprietary and ethical dilemmas. Addressing these issues requires a nuanced approach that includes diverse stakeholder involvement and mechanisms for user engagement. By tackling these challenges, the justice system can move towards a more transparent and ethical use of algorithms, fostering public trust and upholding the principles of justice.