

On Promoting Patient Rights and Algorithmic Accountability in Automated Decision-Making Systems

Ashe Magalhaes^{a,1} and Connor Hounslow^{b,1}

^aUniversity of Edinburgh Informatics Department; ^bUniversity of Edinburgh School of Social and Political Sciences

This manuscript was compiled on May 14, 2019

1 **As automated decision-making systems proliferate into social pro-**
2 **tection systems like healthcare, it is critical that a measure of algo-**
3 **rithmic accountability be included to safeguard patient rights. This**
4 **paper examines the DeepMind-NHS Service Agreement as a case**
5 **study to explore the interpretability-explainability gap that exists in**
6 **DeepMind's recent AI research and why such a gap may be threat-**
7 **ening to patients. This paper advocates for two legislative steps to**
8 **better ensure due process: (a) disclosing if a decision is automated**
9 **and (b) to external auditing mechanisms.**

Algorithmic Accountability | Explainability | GDPR | DeepMind | NHS

1 **A**lgorithmic accountability incorporates a broader debate
2 on fairness and transparency by calling for organizations
3 to provide justification for decisions made by automated sys-
4 tems (1). As decision-making increasingly transfers from hu-
5 mans to computers, accountability is necessary to clarify con-
6 tributing factors to the decision (2) so that a clear attribution
7 of responsibility and standards exists for the subjects that
8 believe the decision is erroneous or unfair. The problem of
9 algorithmic accountability stems from the integration of black-
10 box, or opaque, deep learning classifiers within automated
11 decision-making systems (ADMS) (3). Such classifiers are
12 troubling given the discriminations and biases that permeate
13 society and their potential to be reflected and reinforced in
14 algorithms (4).

15 This is particularly relevant as ADMS are integrated within
16 social protection systems, such as the health care system. We
17 examine the Service Agreements between Google DeepMind
18 and the National Health Service (herein NHS) and note the
19 interpretability-explainability gap that exists in DeepMind's
20 recent Artificial Intelligence (AI) research, raising concerns
21 around the protections offered to patients.

22 1. The Proliferation of Automated Decision-Making 23 Systems

24 ADMS are a form of statistical risk assessment that have the
25 potential to streamline bureaucratic procedures and improve
26 services. As an example, ADMS are used in the financial sector
27 to aggregate customer data (e.g. tenure with bank, number of
28 accounts, demographic variables) in order to automate lending
29 decisions (5, 6). The rise of ADMS can be seen in their
30 increasing deployment across immigration, criminal justice,
31 and healthcare (7–9). As it relates to health and social care,
32 ADMS are promising given their ability to reduce suffering
33 through early detection of disease (10), reduced error rate in
34 diagnosis (11), and personalization in treatment (12).

35 While deep learning models offer the promise of novel pre-
36 dictive capabilities, they are notoriously unexplainable (13),

37 meaning that it is not clear why a given output was produced.
38 Although existent human bureaucratic processes could be per-
39 ceived as unexplainable, experts believe AI applied to patient
40 data must be held to a higher, more transparent, standard
41 (14).

42 2. Defining Explainability

43 In machine learning, the first step towards explainability is in-
44 terpretability, loosely defined as comprehending *what* a model
45 did or might have done as distinct from *how* it did it. The
46 difference between interpretability and explainability can be
47 better understood through the analogy of a chemistry experi-
48 ment: interpretability is the observation of a difference in color
49 or smell during a chemical reaction while explainability is the
50 understanding of the molecular interactions that produced the
51 observed output.

52 In order to develop a more robust understanding of explain-
53 ability, the United States Defense Advanced Research Projects
54 Agency developed a framework for measuring explanations.
55 An ADMS' explanation effectiveness was measured by a user's
56 satisfaction with the explanation, the ability of the user to
57 intervene in the the process at some point, a clarification of
58 the ADMS' mental model, the user trust within the system,
59 and the utility of explanation (15). This framework provides
60 an initial technical understanding of what explanation can
61 and should amount to within social protection systems.

62 3. The DeepMind-NHS Case Study

63 The DeepMind-NHS collaboration can be used as a case study
64 to support the need for stronger safeguards for patients as
65 data-subjects, given the current state of explainability in AI.
66 DeepMind is a leading AI research company and as a subsidiary
67 of Alphabet Inc., benefits from the computing resources of one
68 of the world's most valuable conglomerates. We operate under
69 the assumption that if their current (2019) body of publications
70 does not reflect an ability to integrate explainability with deep
71 learning classifiers, it is reasonable to expect that this is the
72 case for the AI field more widely.

73 In 2015, DeepMind partnered with The Royal Free London
74 NHS Foundation Trust to build a smartphone app, Streams,
75 to help clinicians manage acute kidney injury [16]. This led
76 to the transfer of an estimated 1.6 million patients' sensitive
77 medical data to DeepMind, which has since been found to be in

¹Ashe Magalhaes and Connor Hounslow contributed equally to this work.

²To whom correspondence should be addressed: ashe.magalhaes@gmail.com and chounslow11@gmail.com

78 violation of data protection legislation by the UK’s Information
79 Commission Office [17]. While DeepMind has acknowledged
80 faults in their data-handling process, work on Streams and AI
81 research on patient data continues under DeepMind Health,
82 which reports directly to Alphabet Inc. [18].

83 4. Examples of the Interpretability-Explainability Gap

84 This section examines DeepMind’s three most recent publica-
85 tions, which do not mention an application of their research to
86 healthcare. This survey is meant to illustrate the gap between
87 interpretability and explainability in deep learning models,
88 which may prove problematic if applied to patient data.

89 **A. Agents that Infer Representations From Artificial Con-**
90 **structs.** DeepMind’s research on unsupervised speech repre-
91 sentation learning includes mapping discrete representations
92 to phonemes, discrete components of speech sound [19]. While
93 this mapping adds a form of interpretability to learned repre-
94 sentations, this may be problematic for assessing agent perfor-
95 mance because phonemes are somewhat arbitrary categories
96 that humans have imposed on speech signals rather than quan-
97 titative acoustical physical waves. More broadly, this metric
98 of interpretability asks algorithms to infer representations
99 which are artificial constructs, which linguists do not agree on,
100 making full model explainability challenging.

101 **B. Agents that Infer Causal Structure.** DeepMind’s research
102 includes the first demonstration of model-free reinforcement
103 learning which generates causal reasoning, measured by an
104 agent’s ability to perform tasks dependent on causal structure
105 [20]. In healthcare, causality is likely a combination of genetic
106 factors, environmental factors, and lifestyle choices, making
107 the isolation of causal structure difficult [21]. Therefore, task
108 assessments on such tasks as applied to healthcare may be in-
109 terpretable but not entirely explainable, as with the chemistry
110 example explained in section II.

111 **C. Agents that Pose New Objectives.** Finally, open-ended
112 learning algorithms can create agents that exhibit unknown
113 or unexpected behavior, producing a population of improved
114 agents in settings such as Chess or Go [22]. This adaptive
115 approach to posing new objectives which an agent maximizes
116 may be promising for producing diverse populations that simu-
117 late human expert decisions [23] but may not be explainable or
118 could prove problematic if harmful objectives are maximized.

119 5. A Patient’s Rights in the UK

120 Given the interpretability-explainability gap that exists in
121 DeepMind’s recently published AI research, the use of ADMS
122 in healthcare challenges the existing set of standards–rights and
123 responsibilities–that defines the relationship between doctors
124 and patients within the NHS. The agreements between the
125 NHS trust and Google DeepMind operate within an existing
126 system of human rights within health and social care systems;
127 this section defines the current landscape of patient’s rights in
128 the UK.

129 As it relates to patients’ data, the primary legislative texts
130 are the European Union’s General Data Protection Regulation
131 and the Data Protection Act (2018). The GDPR establishes a
132 data subject’s ‘right to be informed’ about the logic involved
133 (Articles 13-15) and the ‘right not to be subject to automated

134 decision-making’ (Article 22). This does not amount to the
135 ‘right to explanation’ where an individual is able to explain
136 how the ADMS arrived at some conclusion, i.e. a post *hoc*
137 explanation (16).

138 The standard for a patient’s right to due process within
139 social protection systems is stated by the International Labour
140 Organization in the Convention Concerning Minimum Stan-
141 dards of Social Security, 1952 (No.101). Art 70 states that
142 ‘Every claimant shall have a right of appeal in case of refusal
143 of the benefit or complain to its quality or quantity.’

144 Within the UK, the standard of due process for a doctor-
145 patient relationship is set out in the NHS Constitution, estab-
146 lished by the Health Act 2010, and the Human Rights Act 1998.
147 There is a duty of care on behalf of the NHS professionals
148 and the Trust itself (17). There is a related right ‘to be given
149 information about the test and treatment options available
150 to you, what they involve, and their risks and benefits.’ A
151 corresponding duty for doctors is communicating information
152 around the treatment options, their risks and effects. As it
153 relates to the decision of treatment, there is a duty on the side
154 of the doctors to be involved in deciding their health care. This
155 is set out in the NHS Constitution. The right to autonomy is
156 further recognized in the Human Rights Act (Art.8) and in
157 case law, where Lord Donaldson stated in *Re T (Adult) [1992]*
158 4 All ER 649 (18):

159 An adult patient who ... suffers from no mental
160 incapacity has an absolute right to choose whether to
161 consent to medical treatment ... This right of choice
162 is not limited to decisions which others might regard
163 as sensible. It exists notwithstanding the reasons for
164 making the choice are rational, irrational, unknown
165 or even non-existent.

166 6. The Potential of ADMS to Violate Patient Rights

167 The introduction of ADMS into healthcare is problematic
168 because they have the potential to violate the human rights
169 standards enumerated above. In each of the three cases below,
170 ADMS fails to meet a reasonable standard of information
171 that allows the patient to know and understand the treatment
172 options before him or her. In turn, the patient is relegated to
173 the sidelines and their autonomy is negated.

- 174 1. A patient should understand the validity of what the
175 classifier is learning. If ADMS learns contentious artificial
176 constructs, like phonemes in the case of speech, it is
177 reasonable to request experiments with varied learned
178 representations in order to interpret how the output
179 decision may change the risk from treatment. If the
180 experiments demonstrate a high variability in the output,
181 for example if the decision for the patient to take
182 a medication with painful side effects fluctuates, she
183 should have recourse for challenging the treatment option.
184
- 185 2. Similarly, the patient should have access to any structures
186 (causal or otherwise) that form the standard for AI agent
187 task assessment. If a patient wishes to contest a decision
188 around the causality of her brief smoking habit as a
189 determinant in her lung cancer, she can point to the
190 causal structure the agent learned as controversial given
191 the presence of other more relevant genetic and lifestyle

192 facts.

193

194 3. And, a patient in critical condition should have access to
195 the series of objectives that were learned in an automated
196 decision that schedules her hospital care. In the likely
197 case that this decision considers hospital staff, medical
198 resources, and the condition of existing patients, it is possible
199 the algorithm is optimizing for global utility over the
200 well-being of a singular patient, resulting in that patient
201 receiving care that is optimal for the whole ecosystem but
202 not her individually. Ultimately, this should be explained
203 to the patient so that she has the possibility of appealing
204 the automated decision or switching her care provider.
205

206 7. Safeguards for Data-Subjects

207 The standard of explainability as set out in GDPR Art 22
208 (1) does not match the technical capabilities of explanation
209 and (2) does not provide a clear and meaningful way for
210 a patient to challenge a health and social care decision
211 made by the algorithm. In order to improve the due
212 process rights, a patient's care should meet two standards
213 or rights surrounding care.

214 First, a patient should have a right to know the extent
215 to which an algorithm is being used in their healthcare
216 treatment. This right would reasonably allow a patient to
217 choose not to have an algorithm be used in their
218 healthcare process.

219 Second, the patient should also have a right to request
220 an audit of the algorithm. That is, a right to review the
221 supply chain of the algorithm. An automated decision will
222 likely not be replicable without the algorithms, data, and
223 chosen hyperparameters which comprise its 'supply chain'.
224 An ADMS' supply chain is the "training data, test data,
225 models, application program interfaces (APIs) and other
226 infrastructural components" that serve as necessary components
227 for any responsible form of auditing [24]. Without
228 an explanation of this supply chain, the data-subject or a
229 technical consultant, cannot challenge potential scientific
230 flaws in deep learning classifiers. While it may not be
231 reasonable to force a data controller to open-source all
232 three to the general public because of intellectual property
233 concerns, a third-party auditor should have access
234 in order to create some type of oversight, at least until
235 mechanisms for explainability catch up.

236 The Services agreement between Royal Free and Google
237 DeepMind provides the potential for securing due process
238 on behalf of the patients through auditing. Section 8.1 of
239 the agreement states that: '...DeepMind shall use reasonable
240 endeavours to develop and provide the Trust with a service
241 to allow the Trust to obtain an accessible audit history
242 in relation to the Data.'

243 In establishing these standards, patients would have
244 clearer grounds for challenging the use of an algorithm in
245 the provision of health and social care. These standards
246 furthermore cement basic rights to health and autonomy
247 found within the International Covenant on Economic,
248 Social and Cultural Rights, and with the Human Rights
249 Act 1998.

250 8. Conclusion

251 As ADMS proliferate, it is important that we monitor
252 the technical gap that exists between interpretability and
253 explainability in order to gauge the power we give auto-
254 mated decisions in affecting citizens, especially vulnerable
255 populations such as medical patients. Without standards
256 of due process that support a patient's rights within the
257 health and social care system, we as a society risk giving
258 algorithms more decision-making power. The existing
259 duties established within the UK illustrate what general
260 human right standards must be met within health
261 and social care. However, further clarification of patient
262 rights in relation to ADMS is required in light of the
263 interpretability-explainability gap. These rights do not
264 amount to a right to explainability established in existing
265 data regulation such as the GDPR. The necessary standards
266 that should be put into place include (a) disclosing
267 if a decision is automated and (b) allowing for external
268 auditing mechanisms.

269 Notes

270 The legal component of this paper is a layperson's account
271 as both authors do not have a formal educational
272 background in law. The authors' analysis is based on
273 research into the areas which have been discussed in the
274 literature and through our own reflection on the human
275 rights issues that arise from the technical abilities.

276 **ACKNOWLEDGMENTS.** The authors would like to thank
277 Dr. Mark Sprevak and Dr. Nehal Bhuta for providing feedback
278 on the many drafts leading to submission.

- 279 1. Binns R (2017) Algorithmic accountability and public reason. *Philosophy & technology*
280 pp. 1–14.
- 281 2. Lepri B, Oliver N, Letouzé E, Pentland A, Vinck P (2018) Fair, transparent, and
282 accountable algorithmic decision-making processes. *Philosophy & Technology*
283 31(4):611–627.
- 284 3. Pasquale F (2015) *The black box society*. (Harvard University Press).
- 285 4. Chander A (2016) The racist algorithm. *Mich. L. Rev.* 115:1023.
- 286 5. Corbett-Davies S, Goel S (2018) The measure and mismeasure of fairness: A critical
287 review of fair machine learning. *arXiv preprint arXiv:1808.00023*.
- 288 6. Vieira A, Sehgal A (2018) How banks can better serve their customers through artificial
289 techniques in *Digital Marketplaces Unleashed*. (Springer), pp. 311–326.
- 290 7. Molnar P, Gill L (2019) Bots at the gate: A human rights analysis of automated decision
291 making in canada's immigration and refugee system.
- 292 8. Simmons R (2018) Big data, machine judges, and the legitimacy of the criminal justice
293 system. *UCDL Rev.* 52:1067.
- 294 9. Chouldechova A, Benavides-Prado D, Fialko O, Vaithianathan R (2018) A case study
295 of algorithm-assisted decision making in child maltreatment hotline screening decisions
296 in *Conference on Fairness, Accountability and Transparency*, pp. 134–148.
- 297 10. Amato F, et al. (2013) Artificial neural networks in medical diagnosis.
- 298 11. Esteva A, et al. (2017) Dermatologist-level classification of skin cancer with deep neural
299 networks. *Nature* 542(7639):115.
- 300 12. Dilsizian SE, Siegel EL (2014) Artificial intelligence in medicine and cardiac imaging:
301 harnessing big data and advanced computing to provide personalized medical diagnosis
302 and treatment. *Current cardiology reports* 16(1):441.
- 303 13. Montavon G, Samek W, Müller KR (2018) Methods for interpreting and understanding
304 deep neural networks. *Digital Signal Processing* 73:1–15.
- 305 14. Holzinger A, Dehmer M, Jurisica I (2014) Knowledge discovery and interactive data
306 mining in bioinformatics-state-of-the-art, future challenges and research directions.
307 *BMC bioinformatics* 15(6):11.
- 308 15. Gunning D (2017) Explainable artificial intelligence (xai). *Defense Advanced Research*
309 *Projects Agency (DARPA), nd Web*.
- 310 16. Wachter S, Mittelstadt B, Floridi L (2017) Why a right to explanation of automated
311 decision-making does not exist in the general data protection regulation. *International*
312 *Data Privacy Law* 7(2):76–99.
- 313 17. Herring J (2016) *Medical Law and Ethics*. (Oxford University Press).
- 314 18. Service NH (2013) The nhs constitution: the nhs belongs to us all.