
Draft Ethical Report Concerning Predictive Modelling in the Allegheny Babies and Families Project

Michael Veale

Draft April 11, 2019

Contents

- Introduction** **2**
- Short System Overview** **4**
- Proportionality of Targeting** **5**
- Proportionality of the Control Trial** **7**
- Data Minimisation** **9**
- Deployment and Function Creep** **11**
 - Within the County 11
 - Beyond the County 13
- Useful Transparency** **15**
- Holistic Evaluation and Maintenance** **18**
- Supporting Decision Support** **21**

1 Introduction

Public bodies and agencies have been increasingly seeking to use new forms of data analysis to provide ‘better’ public services. These reforms have come in many guises: digital service transformations such as ‘e-government 2.0’ and the creation of ‘integrated data infrastructures’ (linked administrative datasets),¹ broadly aimed at ‘improving the experience of the citizen’, ‘making government more efficient’ and ‘boosting business and the wider economy’.²

Some systems in this vein can be considered as attempting to *automate* administrative decisions—for example, streamlining the act of applying for services or documents which are relatively straightforward to provide. While there will likely always be edge-cases which require human oversight or intervention—hence rules around the world on human intervention in automated decisions³—in gen-

¹Statistics New Zealand, “Integrated Data Infrastructure” (*Government of New Zealand*, 2016) (<https://perma.cc/9RXL-SV7P>) visited on October 4, 2018.
²John Manzoni, “Big data in government: the challenges and opportunities” (*GOVUK*, February 2017) (<https://perma.cc/GF7B-5A2R>) visited on October 4, 2018.
³See eg Lilian Edwards and Michael Veale, “Enslaving the Algorithm: From a “Right to an Explanation” to a “Right to Better Decisions”?” (2018) 16(3) *IEEE Security & Privacy* 46 doi: 10/gdz29v.

11 eral these operate best in uncontentious, well-defined domains.⁴

12 Other systems, such as those which are the topic of this report, attempt to *augment* traditional
13 decision-making, such as administrative decision-making, with additional evidence which improves
14 factors such as its efficacy or consistency. When considering individual decisions over citizens (rather
15 than rule- and policy-making over phenomena of social concern), these systems usually become
16 entwined in existing processes of clinical judgement. The role of clinical judgement—assessment by
17 an individual such as a social worker, who determines a course of action—has many roles in a service
18 delivery context, such as helping individuals navigate labyrinthine processes, getting feedback on how
19 as system works to improve it, ensuring unusual, mitigating or difficult to measure circumstances are
20 substantially assessed and considered, or simply enabling face-to-face contact. Frontline workers
21 providing clinical judgement also hold important tacit knowledge about complex policy problems
22 which can be difficult to codify in rule-based systems, and the procedures they are faced with often
23 have substantial issues and grey areas which leaves them having to make calls within the grey-zones
24 of the rules, at times even acting as ‘street-level ministers’, heavily interpreting the rules in areas
25 which are not clear.⁵

26 At the same time, more ‘actuarial’ judgement, which many will be familiar with through terms such
27 as machine learning or predictive modelling, *can* bring strong benefits to a decision system. Statis-
28 tical modelling is in general more ‘accurate’ than humans at challenging prediction tasks,⁶ although
29 accuracy is by far from the only relevant measure of performance of a system more broadly. Frontline
30 individuals can exhibit undesirable biases or heuristics,⁷ both in an ‘irrational’ psychological sense
31 of mental shortcuts as well as in the broader societal senses of prejudice and discrimination. This is
32 *not* to say that issues similar to this are not present in computing systems used for decision-support,
33 as will be discussed, but that if a societal aim is to mitigate these issues and create a system which
34 is both procedurally and substantively fair to the humans faced with it, neither clinical and actuarial
35 systems should be idolised or romanticised, and opportunities for them to genuinely work together
36 and reinforce each other’s weaknesses should be sought.

37 This report forms one out of three perspectives from cross-disciplinary researchers looking at different
38 aspects of the ethics of the risk-scoring system Allegheny County plans to deploy. My perspective in
39 this third of the broader ethical review is as a researcher in technology policy surrounding the social
40 impacts of data-driven predictive systems, particularly those in the public sector. This report therefore
41 *does not* focus on issues such as the comparative efficacy of this mode of child and family protection
42 compared to other evidence-based interventions, but particularly and more narrowly considers the

⁴Compare to the notion of ‘structured problem’ in Robert Hoppe, *The Governance of Problems: Puzzling, Powering and Participation* (Policy Press 2010).

⁵Michael Lipsky, *Street-level bureaucracy: Dilemmas of the individual in public services* (Russell Sage Foundation 2010).

⁶William M Grove, David H Zald, Boyd S Lebow, Beth E Snitz, and Chad Nelson, “Clinical versus mechanical prediction: a meta-analysis” (2000) 12(1) Psychological Assessment.

⁷Amos Tversky and Daniel Kahneman, “Judgment under Uncertainty: Heuristics and Biases” (1974) 185(4157) Science 1124.

43 role of predictive modelling within the entire scheme.

44 **Short System Overview**

45 What follows is a brief lay summary of the system Allegheny County is proposing that this ethics report
46 concerns.

47 Allegheny County (pop. 1.2 million) has an integrated Department of Human Services which provides
48 an array of services aimed ensuring, among other things, the health, safety and well-being of children
49 and families in early life. A new proposed programme, provisionally named, Allegheny Babies and
50 Families (ABF), seeks to provide i) universal ‘light touch’ outreach to new mothers through a hospital-
51 based community health worker (CHW); and ii) targeted outreach for ‘priority families’ determined
52 through a validated statistical prevention model delivered in collaboration with research institutes
53 including Auckland University of Technology in New Zealand. This second service is spurred by the
54 recognition that ‘the people who most need these services are not using them’:⁸ that prior to a cri-
55 sis such as home removal (of a child), a significant proportion of families are not participating in the
56 County’s supportive services. Existing approaches have focussed on geographic outreach, seeing a
57 major barrier as transportation to one of the County’s 28 Family Support Centers, however geography
58 is a far-from-perfect proxy for maltreatment, as poverty is not a direct cause of abuse, most people in
59 poverty do not abuse or neglect children, and many abused or neglected children come from homes
60 with moderate or high levels of income. The County’s own research illustrated the socioeconomic
61 span of cases of maltreatment, and reviews of the research point to other factors for understanding
62 causes of child abuse, including untreated maternal depression or mental illness, substance abuse in
63 caregivers, inexperienced/young mothers, and intimate partner violence. Consequently, the County
64 wishes to create and validate a statistical ‘prevention model’ to target, initially, a subset of the 6,000
65 mothers who give birth annually at Magee Womens Hospital—the facility with the most births in the
66 County. This model will be used to select families with the highest need for support, and reach out
67 to them with an offer of voluntary services delivered through a range of partners and civil society
68 organisations in a proactive manner, as opposed to waiting for a crisis or for them to locate these
69 resources. 240 families each year will be offered this targeted program service. As part of a commit-
70 ment to evidence-based policy, Allegheny County wishes to rigorously measure the effectiveness of
71 this programme, and as such, plans to integrate a randomised control trial into the above approach.

⁸Internal memo from Amy Malen on Allegheny Babies and Families.

72 **Proportionality of Targeting**

73 Historical analysis undertaken by Allegheny County has indicated that the services they provide to
74 families are not being taken up by those who go on to experience a crisis such as a home removal.
75 Prevention services triggered by referrals exist, in addition to a separate actuarial screening system,
76 however approximately 65% of all child fatalities in the County, largely babies and toddlers, are never
77 referred into the Department of Human Services before the fatal incident.

78 When considering the social impacts of a targeted solution, it is important to consider the alternatives.
79 Some of these alternatives are things that the County could do today instead (or in addition to) the
80 proposed course of action. Others are more ambitious, and would likely require a more significant
81 amount of political will. All of them are important to consider however: even if the County is restricted
82 in its actions, its attitude towards the trajectory of its ambitions concerning predictive technology is
83 important, particularly given its vanguard role and consequent high profile in airing the opportunities
84 and limitations of predictive technologies from operational experience.

85 **Alternative approaches to targeting using machine learning.** Machine learning and predictive
86 systems are harmful when they are seen as a panacea for deep-rooted and complex problems. Partic-
87 ularly given their popularity, they risk obscuring or impeding thoughtful policy analysis and problem
88 structuring, leading to solutions chasing problems rather than the other way around. In particular,
89 while it is avoidable, the use of vendors' off-the-shelf models and systems can be problematic, as
90 these systems can be insufficiently tailored for specific contexts, being sold in a similar form to achieve
91 economies of scale to recoup development costs despite the highly varying nature of policy issues by
92 region and sector, and can embody value-laden rules and decisions that should be subject to in-house
93 development, notice and comment given their policy-like nature.⁹

94 Allegheny County has considered, and continues to invest in, alternative approaches seeking to rec-
95 tify this challenge of non-uptake of services designed in part to reduce the risk of a crisis. These in-
96 clude placing service hubs in areas with high proportions of poverty and violence, where they can be
97 strengthened by nearby civil society organisations. These, however, suffer from the challenge that us-
98 ing location and demographic as a proxy for service need is not wholly congruent with the data, which
99 shows that demographic characteristics associated with cases of child maltreatment vary strongly
100 across the County, with very noticeable differences in median income and race. Research findings
101 instead support other factors, such as untreated maternal depression, substance abuse in caregivers,
102 young mothers and intimate partner violence, which exist throughout the county.

⁹Danielle Keats Citron, "Technological due process" (2008) 85 Washington University Law Review 1249.

Recommendation 1 *Allegheny County should continue to vigorously look for and invest in other ways to reach out to families in need of support services that do not rely on algorithmic systems, to complement the strategies that do.*

Nature of the intervention. Interventions that are coercive or that could trigger coercive effect must be considered more carefully than those that do not. For example, the opening of an investigation into child welfare, even if only partially informed by an algorithm, involves an element of wielding the power of the state against a private citizen. The services offered by Allegheny County in this instance are *voluntary* in nature, and indeed are additional to awareness raising around universal services also planned as part of the same ABF scheme. This is positive from an ethical standpoint. However, there remain concerns around future uses of this score, as well as function creep (see below, in section Deployment and Function Creep). The County should ensure that all uses that area made of this score are supportive and voluntary in nature, rather than being used to trigger investigations.

There is already an argument that such a score, despite only being linked to voluntary services, might be indirectly linked to coercive powers, insofar as the increased proximity of the family to services provided by the County gives more scope for monitoring and for referral by a competent body, were that to be deemed appropriate. However, the same can be said of all protective services however delivered, and safeguards for appropriate referrals should be in place within those services (as they are already). Furthermore, the County should monitor the impact of these services on downstream referrals and increased data collection on these individuals over time to ensure that scores do not disproportionately result in increased surveillance.

Recommendation 2 *The County should pledge that this predictive system be only used to provide voluntary supportive services, rather than to start investigations or to directly inform coercive powers.*

Recommendation 3 *The County should monitor how scoring and service targeting affects the volume of data captured on these subset of individuals and on communities, and take appropriate measures to avoid these groups being disproportionately oversurveilled.*

Targeting as the only means to access services. The proposed plan uses the risk scoring model as the gatekeeper for the targeted service(s), with different services being offered in a graduated fashion at the highest level(s) of predicted risk to the child. A problem with this system is that the predictive system is the *only* gatekeeper in this context—there is no other proposed mechanism through which to access these advanced services. As a principle, to avoid or mitigate ethical issues resulting from unexpected failure modes or biases in the system presented (eg from those who only just moved into the County with no prior data record), and to both respect and safeguard families and communities who wish to avoid data collection, the predictive system should not be the *only* means of targeting

131 families for this scheme. This is not to say that an additional stream of eligibility needs to match the
132 proposed mechanism in scope or scale: but an alternative means for screening-in families should
133 exist. This might, for example, rely on the clinical judgement of relevant trained health workers. This
134 is in line with the principle of *redundancy* in designing safety-critical systems, as well as in-line with the
135 importance of qualitative methods and tacit knowledge in understanding and grappling with complex
136 challenges.

137 Furthermore, the fact that certain services are only eligible for profiled individuals could present other
138 problems. Knowing that individuals were receiving certain services would be sufficient information
139 to reveal their minimum level of risk score on the predictive scale for receipt of that particular service.
140 It might be that this doesn't present a problem, and that the way risk is perceived is not a negative
141 one. There is a chance, however, that higher risk on the predictive scale brings some sort of stigma.
142 Such stigma might be a reason for derision in the community, could foster and perpetuate exclusion,
143 or even result in individuals refusing services for knowledge that it will effectively mean disclosing
144 sensitive data. If it became known that the reasons for the highest scores were mostly related to is-
145 sues such as health, substance abuse, or the like—highly private subjects for some—families choosing
146 to receive these services may believe that they were disclosing this information to anyone who was
147 aware that they were in receipt of them, and as a result might refuse these services.¹⁰

148 **Recommendation 4** *Allegheny County should ensure that at least one other screening-in pathway for
offering these targeted services exists. Such a pathway should have a primarily clinical element rather
than solely relying on an actuarial model, and should aim to capture families whose cases and experi-
ences might not be well quantified for procedural or substantive reasons.*

149 This report does not recommend the form that this alternative pathways should take, as it is not within
150 the expertise of the author to advise on the details of such clinical processes. However, some consider-
151 ations that align with analysis of the proposed evaluation methods below do point to some potential
152 approaches (see section Proportionality of the Control Trial)

153 **Proportionality of the Control Trial**

154 One of the more controversial aspects of this scheme is the proposal to use a randomised control
155 trial to assess the efficacy of the targeted services. The initial version of this proposal would see 240

¹⁰It should be noted that knowledge of the model might interact with the potential for stigma, and that equally transparency around models are encouraged in this report (see the section Useful Transparency). At the same time however, individuals often make 'folk theories' for systems regardless of whether they have knowledge of the innards or not, and therefore challenges such as the one hypothesised here are futile to avoid through secrecy. See e.g. Motahhare Eslami, Karrie Karahalios, Christian Sandvig, Kristen Vaccaro, Aimee Rickman, Kevin Hamilton, and Alex Kirlik, "First I 'Like' It, then I Hide It: Folk Theories of Social Feeds" in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (ACM 2016) doi: 10.1145/2858036.2858494.

156 babies per year who the model indicated would meet the threshold for targeting of service offered the
157 additional resources, in addition to the universal offering that is available to all mothers and families.
158 110 babies per year would be allocated to a control group, and would be treated the same as if they
159 had not been selected for targeting (i.e. able to access core services).

160 This controversy is not unfamiliar to the medical field, and relates to the development of effective, po-
161 tentially life-saving treatments, at the same time as assuring their efficacy in an evidence-based man-
162 ner. The political and ethical challenges of observing or abstaining from an intervention when there is
163 some reason to believe it could be effective is also not unfamiliar to the child protection field, with this
164 being a reason for the shelving of a previous project in New Zealand which some of the researchers
165 on this project had also worked on.¹¹

166 It is important to ask whether there are more proportionate ways of achieving the aim of rigorously
167 evaluating the impact of the targeted services, particularly in relation to the 110 babies a year who are
168 determined to be high-need but who are not offered targeted services. This report will now reflect on
169 some alternative options which may navigate these trade-offs with minimal impact on the rigour of a
170 study, as evidence-based interventions area also of considerable ethical merit.

171 **Option 1** Examine the opportunity of using data from hospitals elsewhere in the county which are
172 not part of this scheme, treating the analysis as a *natural experiment*. Natural experiments are gener-
173 ally found where there are reasons that can be sufficiently classified as random which separate two
174 populations that are otherwise the same. For example, postcode lotteries are a common form for
175 studying natural experiments. Given that the focus of this project is on Magee Womens Hospital, it
176 might be the case that suitable families who would be classified as high risk in a comparable way to
177 those at Magee, but are within an institutitonal context not subject at all to the pilot trial, could be lo-
178 cated. Their results could be compared after several years and any change examined. This approach
179 might not be suitable however if the County has reason to believe that the populations admitted to
180 hospitals other than Magee differ in important ways which are likely to act as confounders. It also
181 suffers from the fact that it does not integrate with the other part of the scheme: the universal service
182 offering, and therefore the comparison would be between families who are offered targeted services
183 which include the universal service offering, and families who are offered the services in the manner
184 the County currently operates at the other hospitals within its boundaries.

185 Option 1 might be ethically preferable to the proposed option given that no family in the other hospi-
186 tals benefits from the ABF scheme (to my knowledge), and the infrastructure (e.g. the workers based
187 in the hospital) to offer such a scheme is not present. It would have the drawback of limiting the
188 options for families in these other hospitals to understand that their data were being processed for

¹¹See Briefing from the Children's Action Plan to Anne Tolley, Minister for Social Development of New Zealand 'Vulnerable Children Predictive Modelling: design for Testing and Trialling' (7 November 2014).

189 research purposes, however this would be less important given that their data were not used for ac-
190 tionable decision-making at any point, and therefore could be defended as a public sector research
191 use of data rather than data used for service delivery.

192 **Option 2** Allow methods for families to opt-in to a separate assessment for targeted services regard-
193 less of whether they had been screened out or not. Concerns around the basis of excluding families
194 from the heightened targeted services are discussed above in the section Proportionality of Targeting,
195 and this is closely linked to Recommendation 4. This option aims to merge a method of tackling that
196 with a method of tackling aspects of the trial issue above.

197 Given that both the universal and the targeted offerings are voluntary, the comparison group of most
198 ethical salience is arguably the group of mothers who were determined to be of high need *and* offered
199 and accepted targeted services, with the group of mothers who were determined to be of high need,
200 not offered services due to the randomised control trial, and yet (hypothetically) *would* have accepted
201 these targeted services had they been offered to her.

202 We might try to proxy this by making the assumption that the subset of mothers who would accept
203 targeted services is within the subset who would accept (aspects of) the universal services. This might
204 not be true, and it is worth looking at ways to validate this—there might be specific incentives in the
205 targeted services that would attract mothers to choose them—but it may be a reasonable assump-
206 tion to make. If this assumption is made, then the universal services become a potential venue for
207 a clinical reassessment of whether targeted services (outside of the RCT) be offered. How this could
208 be carried out is not for this report to suggest, and falls to those with different expertise. The impact
209 of the magnitude of this reassessment might also require the power calculations in the original study
210 plan to be recalculated to allow for the potential of ‘drop-outs’ from the control group. However, I
211 believe this two-stage approach might enable trade-offs around proportionality to be made in a new
212 way in relation to this project.

213 **Data Minimisation**

214 The data minimisation principle is a developed concept in law in different jurisdictions. One definition
215 of it states that data should be ‘adequate, relevant and limited to what is necessary in relation to the
216 purposes for which they are processed.’¹²

217 In the case of this programme, a risk score is being developed, and individuals above or below that
218 score can or cannot be targeted with services. At and above the targeting threshold, it may be useful to

¹²Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation) [2016] OJ L119/1 (GDPR) 5(1)(c).

219 know the exact score, but **below the threshold, there appears to be minimum value in keeping or**
220 **displaying it for operational purposes** (research purposes may differ, particularly concerning issues
221 of model calibration).

222 As a result, it is recommended here to introduce a simple feature within the model software that takes
223 the score calculated by the machine learning model deployed and looks to see if it meets or exceeds
224 a certain threshold. This threshold would, in the first instance, be the minimum score that allows eligi-
225 bility for any service, which as this author understands it is currently 15. If it does, the score is retained
226 in line with relevant standard operating procedures. If it does not, then a more generic message is
227 retained, not revealing any details of the score, and the system only records that the user was not
228 selected for targeting. The difference between these two conditions is due to the fact that targeting
229 is planned *within* the services provided: i.e. those with a score 20 will be eligible for a programme
230 of treatment that those with scores 15–19 will not. This also presents a light barrier against function
231 creep (see below section), as miminising the results with respect to a certain purpose may make it
232 more difficult to use downstream for other incompatible purposes.

233 **Recommendation 5** *Scores should be redacted on production to ‘Ineligible’ if they fall below deter-
mined thresholds for service ineligibility, rather than unnecessarily keeping them in numeric form. An
proportionate exemption to this could be determined if they were deemed necessary for validation or
research, in which case these should be kept separately and redacted in the operational system.*

234 Once service delivery begins, it is unclear what role the score has. If no further use can be established,
235 it is recommended that the score be purged entirely. In particular, this would also be useful because of
236 the perception that such as score might be used as a ‘black mark’ on somebody’s file—an argument
237 which would be easier to refute if scores were deleted rather than retained for unclear or potential
238 purposes.

239 **Recommendation 6** *At the point at which there are no plans to make use of the value of a score through
a standard operating procedure, the score should be purged.*

240 Data minimisation can also be considered in the sense of making the data coarser for any particu-
241 lar purpose. For example, if a score was desirable to retain to inform service providers that families
242 who had previously refused or dropped out of schemes were a priority to re-engage, the question
243 should be asked of what practical use retaining the difference between a 15 and a 20 would be on the
244 re-engagement intervention. If the use was limited, then the score could be replaced with a marker
245 that they were a priority to (re-)engage—a marker that could also be usefully triggered through other
246 operating procedures.

247 **Recommendation 7** *If there are justified cases for the score being retained for downstream service pro-
vision, tailoring or as flag for re-engagement, the score should be made as coarse as is possible and
compatible with those purposes.*

248 **Deployment and Function Creep**

249 One underlying anxiety concerning predictive systems in the public sector is that by virtue of being
250 created for one task, they establish an infrastructure consisting of many aspects—including data, tech-
251 nology, expertise and culture—which might expand beyond its original scope into areas its original
252 democratic and societal mandate did not permit.

253 This section explores this from two starting points. The first is *within the County*, envisaging safeguards
254 against creep in the current system by future officeholders who might want to deploy the system dif-
255 ferently or share scores from the model more widely. The second starting point is *beyond the County*,
256 considering the appetite for predictive systems in public bodies across the world more generally, and
257 the pressure on a vanguard administration such as Allegheny County to provide software and models
258 (if not data) to help them with their own challenges. Co-operation between public bodies on areas
259 of technical difficulty is naturally welcome, however ethical questions do arise if the processes of the
260 body to which the knowledge and software are being transferred to fall significantly below that of
261 Allegheny County. Significant trust and reputation issues might follow given the provenance of the
262 system.

263 **Within the County**

264 Allegheny County has, in previous work, given welcome consideration to delimiting of the role of the
265 system in operating procedures within the welfare system. A previous system implemented by the
266 County in 2016, the Allegheny Family Screening Tool (AFST),¹³ provided call-screeners with ventiles
267 (numbers 1–20 representing 5% percentiles of a probability distribution) estimating the risk that the
268 child that is the subject of the call would be removed from home were they to be screened in. These
269 inferred data were not shared beyond this stage of the process (e.g. to investigators downstream),
270 which limits the role of these data as children who were screened-in manually or who were screened
271 in on the recommendation of a high risk score are not distinguishable.

272 Some will be concerned that while that might be the policy today, it might not be robust to change
273 in the future. Similarly, those who might have lost trust in a public service more generally might not
274 trust assurances that this inferred data is deleted or not passed onto other actors in the system. From
275 a technical standpoint, there are limitations to any formal guarantees that can be provided of this
276 nature. Allegheny County could indeed delete the scores that were produced by this model, but the
277 nature of the models used by the County to date have all been deterministic in nature, meaning that
278 with a copy of the model at that time (which would be advisable to keep for purposes of auditing and

¹³See generally Alexandra Chouldechova, Diana Benavides-Prado, Oleksandr Fialko, and Rhema Vaithianathan, “A Case Study of Algorithm-Assisted Decision Making in Child Maltreatment Hotline Screening Decisions” in *Conference on Fairness, Accountability and Transparency (FAT* 2018)* (2018).

279 accountability) and a copy of the data record in question (which is likely to persist and be expanded
280 upon over time), such a score could be trivially recalculated.

281 The County may instead be able to provide some institutional and legal assurance of the path that
282 they choose to go down. In many regimes around the world, data protection law or frameworks like
283 it provide for the concept of *purpose limitation*.¹⁴ While a similar regime may not be available in Penn-
284 sylvania law, similar constructs might be achieved through contracts or declarations which bind the
285 County. For example, if such a predictive system were defined, some legally binding declaration could
286 be made delimiting the purposes of this system in advance to a sufficiently narrow scope and set of
287 actors. This agreement would then serve as a mechanism that could be used to hold future uses of
288 this model to account—at least insofar as it would have to be actively and ideally publicly removed
289 before the purposes of a score or a model could change. Organisations without such agreements have
290 been heavily accused of function creep in other contexts, such as the use of data from the education
291 sector in immigration enforcement in the United Kingdom.¹⁵

292 While it is not for this (non-US lawyer) to stipulate the legal form such an instrument would take, I can
293 illustrate some of the characteristics that any suitable form should be tested against. It should be as
294 binding as possible given its nature; it should allow for third-parties to challenge its enforcement; it
295 should not be amendable without public declaration; and its core form should be public.

296 In the absence of an overarching framework governing the use of data and prediction in public ser-
297 vices, the County should attempt to ensure that it remains proactively able to be held to account.
298 Given the intentions of the current project this does not seem to place restrictions on what the sys-
299 tem proposed intends to do, but instead acts as a guard against unwarranted function creep. In doing
300 so, the County would also set an example for other public bodies and encourage them to do the same.
301 This would be a collective benefit, as a single authority misusing scoring in a way that is seen as so-
302 cietally unacceptable is likely to tar trust across organisations seeking to use data and technology in
303 proportionate and safeguarded ways that are highly mindful of recent debates around the pitfalls of
304 these practices.

305 **Recommendation 8** *The County should examine options for creating a pledge in a binding and public
manner to ensure that the intended purposes of the system do not expand without a clear and account-
able process, and the safeguards do not diminish.*

306 This would sit alongside the existing initiatives for misuse in general that the County has in place, such
307 as regular audits of the *Client View* portal for analysing the data of individuals.

¹⁴See eg GDPR, 5(1)(b). The notion of purpose and compatible purposes also arises in the California Consumer Privacy Act, 2018 Cal. Legis. Serv. Ch. 55 (A.B. 375).

¹⁵Damien Gayle, “Schools Census Used to Enforce Immigration Laws, Minister Says” (*The Guardian*, January 13, 2019) (<https://www.theguardian.com/politics/2019/jan/13/schools-census-used-for-immigration-enforcement-minister-says>) visited on April 10, 2019.

308 Given the significant powers that the County has, both in terms of being a safeguard for extremely
309 vulnerable individuals as well as coercive powers of the state to e.g. remove children, it is noted that
310 predictive modelling is by far from the most straight-forward method to do harm, were the organ-
311 isation to find itself controlled by those with agendas that were not well-suited for addressing the
312 difficult social challenges the Department for Human Services is faced with. That being said, comput-
313 ing systems can in many cases be seen as administrative rule-making systems,¹⁶ and risk being used
314 to disguise the true intent of policies that otherwise would be easier to uncover through tools such
315 as FOIA. At the very least, they risk being *seen* as vehicles through which this could be done, even if
316 they are not, and mechanisms of binding legal recourse are one way in which to build public trust,
317 particularly in an environment where other agencies around the world may adopt less savoury uses
318 of similar systems.

319 **Beyond the County**

320 As Allegheny County is one of the regions of the United States most advanced in this form of analysis,
321 it is highly likely that other public organisations will seek to build on the experience the County has
322 gained during building and deploying these systems. This itself brings ethical challenges.

323 The responsible downstream use of computationally-advanced statistical systems has been a sub-
324 ject of contention in a number of fields. General purpose technologies, like all similar tools, can be
325 applied towards a wide variety of ends. It cannot be guaranteed that predictive tools such as those
326 developed at Allegheny County will be redeployed elsewhere with comparable safeguards or ethical
327 considerations. Not only might this harm individuals, families and communities, but the irresponsible
328 deployment of technologies elsewhere might severely erode trust in more responsible deployments
329 of similar technologies, as well as present serious reputational risk to the originators of these tools.

330 Some existing tools exist aiming to constrain the downstream use of software in society. This report
331 does not seek to mandate any particular scheme, but recommends that Allegheny County consider
332 feasible and effective safeguards within the legal regime and obligations they work under to

333 **Licensing schemes.** Licensing schemes attempt to legally constrain what can be done with code.
334 While proprietary licensing is commonplace, if not the norm, in much commercial software, licensing
335 has also been deployed to secure public value from developed code. Notable, the GNU General Public
336 License (GPL) license enables software to be freely shared under the condition that any derivative
337 products are also shared under the same license. This contrasts with permissive free software licenses,
338 such as the MIT License, which places very few restrictions on downstream reuse.

¹⁶Citron (see n. 9); Michael Veale and Irina Brass, “Administration by Algorithm? Public Management meets Public Sector Machine Learning” in Karen Yeung and Martin Lodge (eds.), *Algorithmic Regulation* (Oxford University Press 2019).

339 Allegheny County, as emphasised elsewhere in this report (see the section Useful Transparency),
340 should prioritise transparency measures aimed at empowered third parties. It should also ensure—
341 using the limited leverage it has as the developer—that the deployment of similar systems elsewhere
342 is accompanied by a similar minimal level of transparency.

343 One approach for this could be as follows. The County could release code associated with this project
344 under licensing terms which allow open reuse, but which mandate the licensing of derived code un-
345 der the same licenses as the original code was licensed under. In software licensing, this would be a
346 sibling of *copyleft* licensing. For example, such a license could state that any organisation deploying
347 such code towards a public function employ institutional methods to allow third parties to scrutinise
348 its function, or more extremely, to release all code openly. *Copyleft* regimes have been criticised for
349 extending these provisions to code which surrounds or interfaces with the code under license, insofar
350 as they are distributed as an inextricable package.

351 It should be emphasised that this approach comes with particular limitations. If these downstream
352 provisions are breached, the infringement by the party utilising the code is one of the *copyright* of the
353 County. There exist an array of specific questions that follow from that, such as the enforceability of
354 breaches in this contract by a third party.¹⁷ There may be other legal approaches which could be used
355 to this effect, and the County should consider exploring those.

356 **Institutional programmes.** Institutional programmes may also serve to support responsible dis-
357 semination of machine learning systems to other organisations. Allegheny County might wish to op-
358 erate or accomodate employees from other public bodies being seconded in to learn more about how
359 systems are deployed on the ground, and to pick up tacit knowledge. They may wish to establish
360 high-quality documentation of both the technical deployment as well as the social processes involved
361 across the organisation in relation to these systems.

362 **Co-development of future systems.** This is a more prospective arrangement that the County may
363 wish to consider. In effect, this would result in consortia of two or more distinct governments fac-
364 ing comparable challenges joining their efforts and funding together to co-develop software systems,
365 drawing on capacity that might be difficult for one alone. These consortia themselves would instigate
366 knowledge sharing agreements for both technical and social processes surrounding the deployment
367 of predictive and other systems, and further regions or governments wishing to benefit from the de-
368 veloped code and expertise would be able to join the consortium upon meeting certain determined
369 membership conditions.

¹⁷See further Andrés Guadamuz, “Viral Contracts or Unenforceable Documents? Contractual Validity of Copyleft Licenses” (2004) 26(8) European Intellectual Property Review 331 (<https://ssrn.com/abstract=569101>).

Recommendation 9 *The County should produce and make public a strategy for ensuring that the system and their expertise is disseminated in the most responsible manner possible to other interested public sector bodies. This strategy may wish to consider legal instruments such as licensing schemes, and institutional arrangements such as secondment and collaboration agreements.*

Useful Transparency

Individuals should have a right to know why they were offered targeted services upon request. These ‘algorithmic explanations’ might promote perceptions of justice in the process as a whole,¹⁸ and are commonly proposed,¹⁹ and a version of them is recommended here.

Recommendation 10 *Allegheny County should ensure that easy-to-understand explanations of why an individual was offered targeted services are provided upon request. In documentation or during outreach, the opportunity to request this information should be actively presented.*

Furthermore, input data is important as a component of understanding systems—and in some cases, may be more important than providing access to the innards,²⁰ particularly when the type of variables that would be used in this kind of system might not be well-understood by citizens. Indeed, the use of data in this way would be a salient time to highlight to individuals that they can see copies of this data, and to explain the process through which that can be done. This is usually called a subject access request in laws around the world,²¹ but in absence of it being present in a generic form in the United States, it would be recommended for Allegheny County to make particular effort to provide and publicise it in this case given the high-stakes nature of the subject matter.

Recommendation 11 *Allegheny County should ensure that individuals can request all information that relates to them used by the County in the creation of this score, and are actively informed about their ability to do this.*

¹⁸See Reuben Binns, Max Van Kleek, Michael Veale, Ulrik Lyngs, Jun Zhao, and Nigel Shadbolt, “‘It’s Reducing a Human Being to a Percentage’; Perceptions of Justice in Algorithmic Decisions” in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI’18)* (ACM 2018) doi: 10/cvcp.

¹⁹See e.g. their presence in international law in the Modernised Council of Europe Convention 108; in French administrative law, and in EU data protection law. See Mireille Hildebrandt, “The Dawn of a Critical Transparency Right for the Profiling Era” in Jacques Bus, Malcolm Crompton, Mireille Hildebrandt, and George Metakides (eds.), *Digital Enlightenment Yearbook* (IOS Press 2012) doi: 10/cwpm; Edwards and Veale, “Enslaving the Algorithm” (see n. 3).

²⁰Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Vaughan, and Hanna Wallach, “Manipulating and Measuring Model Interpretability” [2018] arXiv:1802.07810 [cs] (<http://arxiv.org/abs/1802.07810>) visited on April 10, 2019.

²¹Jef Ausloos and Pierre Dewitte, “Shattering one-way mirrors—data subject access rights in practice” (2018) 8(1) *International Data Privacy Law* 4 doi: 10/cwcf.

385 However, algorithmic explanations aimed at the individuals subject to algorithmic decisions should
386 equally not be romanticised as a panacea or governance solution.²² While explanations can be useful
387 for a range of purposes, such as identifying errors or establishing trust, they face several limitations.
388 Explanations burden decision-subjects with understanding and determining the nature of decisions
389 that can affect them. Those who might need to challenge them most might be in vulnerable posi-
390 tions, and the least able to grapple with quite complex systems (and the bureaucracies or political
391 economies behind them). This is especially the case in the situation in question, where women being
392 approached for targeted services have just given birth, and can additionally be assumed to be in-need
393 in variety of ways due to the fact they have been targeted by this system. Furthermore, many explana-
394 tions can often be generated for each system, because of the non-linear nature of statistical systems.
395 For those individuals who were outliers in a dataset and might like an explanation most, it is not clear
396 that the explanation will make sense or be of much use.²³

397 It is therefore necessary to think of other forms of transparency that can augment individual expla-
398 nations and reduce the burden on those affected by the decisions to oversee the system at the same
399 time as experiencing it, and raising a young child. Primarily, in line with many movements around the
400 world for open data and transparency on governmental systems, the County should aim to make the
401 developed systems as open as possible. Allegheny County should consider publishing either full ver-
402 sions of its models, or parsimonious versions that estimate their core logics (eg a version of the more
403 simple SLIM²⁴ model trained). Full versions of the model may not be possible or wise if such a model
404 has sensitive variables that might leave it at risk of confidentiality attacks such as model inversion,
405 which reveal a semblance of the data used to train it.²⁵ It is possible such publication could be unwise
406 both reputationally and under laws such as HIPAA. However, many models, such as regression-based
407 approaches, are generally resistant to model inversion attacks due to being destructive of data in the
408 process of training for generalisable insights. Such systems might be published on a specific open
409 data portal, or a code repository such as *GitHub*.

410 **Recommendation 12** *Allegheny County should seek to publish version of the models it produces online, after ensuring they would not leak sensitive data.*

411 If a model cannot be released, Allegheny County should aim to release a ‘model-centric explanation’²⁶

²²See Mike Ananny and Kate Crawford, “Seeing without knowing: Limitations of the transparency ideal and its application to algorithmic accountability” [2016] *New Media & Society* doi: 10/gddxrg; Lilian Edwards and Michael Veale, “Slave to the Algorithm? Why a ‘Right to an Explanation’ is Probably Not The Remedy You Are Looking For” (2017) 16 *Duke L. & Tech. Rev.* 18 doi: 10/gdxtlj.

²³Edwards and Veale, “Slave to the Algorithm?” (see n. 22).

²⁴Berk Ustun and Cynthia Rudin, “Supersparse Linear Integer Models for Optimized Medical Scoring Systems” (2016) 102(3) *Machine Learning* 349 doi: 10/f8crhw.

²⁵Michael Veale, Reuben Binns, and Lilian Edwards, “Algorithms That Remember: Model Inversion Attacks and Data Protection Law” (2018) 376 *Phil. Trans. R. Soc. A* 20180083 doi: 10/gfc63m.

²⁶Edwards and Veale, “Slave to the Algorithm?” (see n. 22).

412 or ‘model card’²⁷ describing general features of the model. These would include summary statistics
413 on model performance on various sub-groups and in various situations, a presentation of the core
414 ‘logics’ extracted from the model (for example, random forest variable importance scores, despite
415 their caveats,²⁸ can be useful in this situation). Such information should also include process-related
416 elements such as data source, quality assurance, maintenance and re-training procedures, and the
417 user interfaces it is deployed within. Ideally, such a model would also be presented in peer-reviewed
418 work analysing it in context, as was undertaken with the AFST previously.²⁹

Recommendation 13 *Allegheny County should publish and maintain ‘model-centric’ explanations on-line: metadata about the models including at least their main inputs, logics, and optimisation targets; performance in practice, including on sensitive and salient sub-groups; and practices and processes around model maintainance and use.*

420 The most promising approach for transparency is to aim it at empowered civil society, research or jour-
421 nalistic organisations, providing access to models. If models can be totally published, then these or-
422 gansiations need only to be additionally access to the institutional and human infrastructures around
423 them in the wider decision system. If they cannot be published, then other approaches might be pos-
424 sible. Sporadic physical arrangements for e.g. computational journalists to come and, through secure
425 terminals and pre-installed software, examine data use, would be also be a possibility. Such arrange-
426 ments are commonly used in statistical agencies for examining sensitive microdata, and inspiration
427 on how that could be applied to models could be sought there—although there would be economic
428 considerations.³⁰

429 As the number of tools of this type grows within the County, it is worth the County considering what
430 other scalable methods exist for oversight. As part of this report I have heard about the general strat-
431 egy of openness towards both critics and those who want to understand more about the system by
432 observing and studying it in practice. It would be worthwhile attempting to systematise efforts for
433 internal review and oversight. Information days, either physically or online through webinars, might
434 be useful in disseminating information more broadly. Other arrangements which are of a slower na-
435 ture could be considered, such as joint-supervised PhD students with local universities, or projects
436 for groups of postgraduates. Where research projects can lead to public, peer-reviewed outputs, this
437 should be encouraged. While none of these presents a silver bullet for accountability and oversight,
438 all of them contribute to general atmosphere of openness and collaboration that set the bar for other
439 agencies to meet.

²⁷Margaret Mitchell et al., “Model Cards for Model Reporting” in *Proceedings of the 2nd ACM Conference on Fairness, Accountability and Transparency (ACM FAT* 2019)* (ACM 2019) (<https://arxiv.org/abs/1810.03993>).

²⁸Carolin Strobl, Anne-Laure Boulesteix, Achim Zeileis, and Torsten Hothorn, “Bias in random forest variable importance measures: Illustrations, sources and a solution” (2007) 8(1) *BMC Bioinformatics* 1 doi: 10.1186/1471-2105-8-25.

²⁹Chouldechova, Benavides-Prado, Fialko, and Vaithianathan (see n. 13).

³⁰See generally Felix Richie, “Access to sensitive data: Satisfying objectives rather than constraints” (2014) 30(3) *Journal of Official Statistics*.

Recommendation 14 *Allegheny County should invite researchers to examine and audit their systems and give them heightened access to software, data and front-line workers. A list of such collaborations should be published.*

Recommendation 15 *Allegheny County should seek to work with researchers to create and make publicly available peer-reviewed research on the models characteristics and their use in situ.*

Holistic Evaluation and Maintenance

Machine learning and decision support systems exist within a procedural and organisational context, which itself exists within a wider landscape of societal structures and challenges. In issues of welfare and social protection, this is especially salient.

Actuarial or machine learning systems used in high-stakes decision-making have been long accused of creating and perpetuating bias and discrimination in society.³¹ As a result, a range of methods through which different fairness properties can be assessed and statistically assured have been developed.³² It is likely useful to use such methods in this case, at least in an initial audit capacity, however the exact trade-offs will only become apparent through trying out many different types,³³ and there may be restrictions on the types of fairness issues that can be tackled based on limitations in measured or available sensitive data.³⁴ The County has prior experience of assessing deployed systems for bias and should draw upon the institutional routines and practices developed there.³⁵

³¹See e.g. Oscar H Gandy, “Engaging Rational Discrimination: Exploring Reasons for Placing Regulatory Constraints on Decision Support Systems” (2010) 12(1) *Ethics and Information Technology* 29 doi: 10/bzwqrx; Solon Barocas and Andrew D Selbst, “Big Data’s Disparate Impact” (2016) 104 *California Law Review* 671 doi: 10/gfgq9w.

³²See e.g. Faisal Kamiran and Toon Calders, “Classifying without discriminating” in *2nd International Conference on Computer, Control and Communication, Karachi, Pakistan, 17–18 February, 2009* (2009) doi: 10/dtcf5n; Andrea Romei and Salvatore Ruggieri, “Discrimination data analysis: A multi-disciplinary bibliography” in Bart Custers, Toon Calders, Bart Schermer, and Tal Zarsky (eds.), *Discrimination and Privacy in the Information Society* (Springer 2012); Sara Hajian and Josep Domingo-Ferrer, “Direct and indirect discrimination prevention methods” in Bart Custers, Toon Calders, Bart Schermer, and Tal Zarsky (eds.), *Discrimination and privacy in the information society* (Springer 2012); Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel, “Fairness through awareness” in *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference (ITCS ’12)* (2012) doi: 10/fzd3f9; Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian, “Certifying and removing disparate impact” in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2015) doi: 10/gfgrbk; Moritz Hardt, Eric Price, and Nati Srebro, “Equality of Opportunity in Supervised Learning” in DD Lee, M Sugiyama, UV Luxburg, I Guyon, and R Garnett (eds.), *Advances in Neural Information Processing Systems 29* (Curran Associates, Inc 2016); Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gummadi, “Fairness beyond Disparate Treatment & Disparate Impact: Learning Classification without Disparate Mistreatment” in *Proceedings of the 26th International Conference on World Wide Web* (International World Wide Web Conferences Steering Committee 2017) doi: 10/gfgq8r.

³³On empirical comparisons in this area, see Sorelle A Friedler, Carlos Scheidegger, Suresh Venkatasubramanian, Sonam Choudhary, Evan P Hamilton, and Derek Roth, “A Comparative Study of Fairness-Enhancing Interventions in Machine Learning” [2018] arXiv preprint (<http://arxiv.org/abs/1802.04422>).

³⁴Michael Veale and Reuben Binns, “Fairer machine learning in the real world: Mitigating discrimination without collecting sensitive data” (2017) 4(2) *Big Data & Society* doi: 10/gdfnzn.

³⁵Chouldechova, Benavides-Prado, Fialko, and Vaithianathan (see n. 13).

454 More broadly however, discrimination-aware data analysis comes with risks. Allegheny County has a
455 range of policy levers they can use in order to understand and to tackle issues of poverty and inequal-
456 ity. This report calls for all issues of discrimination and bias—which involve both the under and the
457 over targeting of certain (e.g. ethnic) groups—to be tackled in a holistic manner. While it is possible,
458 and even tempting, to try to ‘correct’ for biases in a decision-support model, or even to consider intro-
459 ducing positive discrimination,³⁶ there is a risk that this both reifies the importance of the actuarial
460 model in the broader process, falling into the trap of function creep,³⁷ as well as draws attention away
461 from more substantive and sustainable policy mechanisms to address bias and discrimination, such
462 as designing specific policy strategies to reach out or to consult with poorly represented groups. As
463 has been previously noted in relation to child protection tools, ‘fairness is a process property, not just
464 a model property’,³⁸ and this should be kept strongly in mind and at the heart of relevant interven-
465 tions.

466 Consequently, this report suggests that *while discrimination-aware data mining is useful, it must form*
467 *part of a broader strategy*. The County should resist becoming a poster child for statistical fairness, and
468 should seek instead to integrate it seamlessly as just one part of a broader strategy around how tar-
469 geted interventions—whether that be actuarially or by other means such as geographic investment—
470 interact with issues of concern such as discrimination and inequality. Many findings in actuarial fair-
471 ness indicate the difficulties present when trying to reconcile different notions of formalised fairness,
472 and that it can be mathematically impossible to have many types of fairness that might be desired at
473 once.³⁹ This does not mean however that broader strategies cannot aim at fairness, particularly when
474 the question is not a binary one (such as child removal) but a decision point with many potential in-
475 terventions that are qualitatively different from each other.

476 **Recommendation 16** *The County should begin a programme of work on discrimination and inequality
in information use and service provision more broadly which involves, but is not limited to, statistical
understandings of fairness. Where deficits in modelling are discovered, particular attention should be
paid to how different interventions might seek to mitigate these, rather than solely attempting to adjust
the model to compensate.*

477 **Recommendation 17** *The County’s work in building the strategy in Recommendation 16 should be
deeply participatory and open to comment.*

478 There are several further issues concerning modelling that are important whether personal data is
479 involved or not, and these issues are relevant to an ethical analysis of this case.

³⁶Sicco Verwer and Toon Calders, “Introducing positive discrimination in predictive models” in Bart Custers, Toon Calders, Bart Schermer, and Tal Zarsky (eds.), *Discrimination and privacy in the information society* (Springer 2013).

³⁷See the section Deployment and Function Creep.

³⁸Chouldechova, Benavides-Prado, Fialko, and Vaithianathan (see n. 13) 13.

³⁹See e.g. Alexandra Chouldechova, “Fair prediction with disparate impact: A study of bias in recidivism prediction instruments” (2017) 5(2) *Big Data* 153 doi: 10/gdcdqf.

480 Firstly, statistical models in particular suffer from a risk of *concept drift*. Concept drift or dataset shift
481 refers to changes in the conditional distributions of model inputs (e.g. DHS warehouse data) and out-
482 puts (e.g. a placement event).⁴⁰ For example, a model may have been built before some demographic,
483 economic or legal change which makes the population in question, and their social behaviour, signif-
484 icantly different to what it once was. It might also be that some individuals become more or less
485 difficult to collect data on over time, and therefore end up under- or over-represented in a modelling
486 process.

487 The consequences of concept drift can be severe. It might, for example, be that the model no longer
488 represents the social phenomena on the ground as well as it could, or potentially could misrepresent it
489 for certain demographic subgroups. In the Allegheny case, the nature and challenges of families with
490 high scores or with crisis events may change over time. There are some tools that can help understand
491 concept drift, and where technically feasible these should be considered by the modelling team.⁴¹
492 These might include refreshing the model, dropping data over time, or weighting past data less heavily
493 than recent data. Implementing these can be challenging, and they do not always cope well with
494 certain types of drift that might be sudden or cyclical in nature. As reviews of the field have indicated,
495 integrating separate, expert knowledge to detect when something has changed that the model has
496 not considered is likely to be an important component of any concept drift detection strategy.⁴²

497 Linking qualitative knowledge to understand drift in machine learning models is a challenge in public
498 sector machine learning in general.⁴³ How might Allegheny County cope with it in an ongoing man-
499 ner. Some lessons can already be taken from some of the preparatory work the County has done in
500 interviewing ‘20/20’ families to understand their needs. While instigating a regular programme of in-
501 terviews *solely* for the purpose of better understanding a predictive model might be disproportionate,
502 this report instead recommends an integrated evidence-gathering strategy across the County which
503 allows interviews and other data (such as focus groups with specialist social workers) to be fed regu-
504 larly into modelling and maintenance teams. For example, risk scores might become part of an inter-
505 view sampling strategy for more general evidence gathering across the County, enabling these inter-
506 views to play a useful and importantly a specific role in many different functions of the Department
507 of Human Services.

⁴⁰ Joaquin Quiñero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence, *Dataset shift in machine learning* (The MIT Press 2009).

⁴¹ J Gama, Indre Žliobaitė, A Bifet, M Pechenizkiy, and A Bouchachia, “A survey on concept drift adaptation” (2013) 1(1) ACM Comput. Surv. doi: 10/gd893p.

⁴² See eg *ibid.* 30.

⁴³ See generally Michael Veale, Max Van Kleek, and Reuben Binns, “Fairness and Accountability Design Needs for Algorithmic Support in High-Stakes Public Sector Decision-Making” in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI’18)* (ACM 2018) doi: 10/ct4s 5.

Recommendation 18 *Maintenance of the model over time should have a strong qualitative dimension, seeking feedback from different sources of on-the-ground knowledge, such as front-line workers, to understand how the model performs on different groups, and how the phenomena being modelled might be changing over time.*

Supporting Decision Support

Much has been written about how decision support systems might end up determining decisions rather than just providing one of several pieces of relevant evidence. Broadly this has been characterised as an issue of *automation bias*, where individuals either over- or underrely on computerised information compared to a rational model of how much it should be relied upon.⁴⁴

In this case, automation bias presents a concern where case-workers are manually involved in reviewing a case where the score is visible. Automation bias might also work in more complex ways, such as in interaction with features such as perceived stereotypes.

Allegheny County must actively work to mitigate and pre-empt automation bias in predictive systems, as such biases may serve to undermine the aims of the tools being deployed.

In the case of the Allegheny Family Screening Tool, researchers and practitioners affiliated to the County's project analysed the level of over- or under-reliance, in terms of how often a child the system had flagged as a 'mandatory' (i.e. highly recommended) screen-in was manually overridden by supervisors. It is noted that some organisational operating procedures are already in place concerning this particular tool: for example, that only supervisors, rather than front-line screeners, are allowed to 'screen-out' scores that are in the top 15% of risk.⁴⁵ Despite this barrier, and while there is an observable drop in screen-outs over the 17–18 ventile threshold that represents this level of risk, a considerable proportion of overrides are visible.⁴⁶

Training and investment in tools beyond predictive scores In Allegheny County, the data that is used to calculate the score is available to decision-makers through their *Client View* system. This platform also contains a variety of documents, plans and assessments which decision-makers can examine. It is important that investment, training and testing in this system is maintained, in order to present decision-makers with feasible informational alternatives to the score provided. While the proposed programme involves the drawing up of lists to give to programme providers rather than decision-support in a clinical setting (as the AFST does), this recommendation is still important for

⁴⁴Linda J Skitka, Kathleen L Mosier, and Mark Burdick, "Does automation bias decision-making?" (1999) 51 *International Journal of Human-Computer Studies* 991 doi: 10/bg5rb7; Jaap J Dijkstra, "User agreement with incorrect expert system advice" (1999) 18(6) *Behaviour & Information Technology* 399 doi: 10/fsnqm9.

⁴⁵Chouldechova, Benavides-Prado, Fialko, and Vaithianathan (see n. 13) footnote 6.

⁴⁶*ibid.* 12.

534 those who are drawing up those lists or otherwise involved in monitoring and evaluating the scores.
535 As a general principle around this kind of use of predictive system, they must be deployed with con-
536 sideration for user testing, improvement, training and maintenance of existing, complementary sys-
537 tems.

538 **Workload management in the presence of risk scoring.** Where individuals are working alongside
539 systems using this score in their role, and where the score is supposed to make their jobs faster or more
540 efficient, care should be taken that there remains substantive time to oversee each individual case as
541 appropriate. The County should be mindful of the risk that where resources are constrained, there
542 is a greater tendency to want to rubber-stamp machine outputs due to limited capacity to provide
543 meaningful oversight.

544 **Organisational structures to support disagreement and monitor reliance.** Wherever risk scor-
545 ing meets organisational process, and the operating procedures permit some override, it needs to be
546 ensured that staff involved have the confidence and support to ‘disagree with the machine’. In the
547 ABF case, there do appear to be some cases where discretion is required. The DHS programme staff
548 aggregating families into groups to tier services and to provide lists to the service providers may be
549 empowered to include families who were not scored sufficiently in these lists. The community health
550 workers reaching out to mothers in the hospital may need discretion to be able to offer targeted ser-
551 vices where they see real need. Those monitoring and evaluating the system in the DHS may need
552 confidence to spot errors and to raise them with senior staff or external modellers.

553 Operating procedures may rightly discourage disagreement with the system where there is strong
554 evidence to support that judgement. In the AFST system, disagreement was discouraged through
555 defaulting those calls with a screening score of over 18 to being ‘screened-in’, with the only mechanism
556 of override available being an appeal to the supervisor on duty at the time. Despite this, a significant
557 number of overrides were recorded, which in many ways is a strong sign of critical independence and
558 discussion of the models rather than a rubber-stamping of decisions.

559 Allegheny County should establish appropriate monitoring measures for disagreement or over- or
560 under-reliance, and follow these up with qualitative studies. Such studies should be carefully de-
561 signed not to appear to blame those disagreeing with the system, but to ensure they feel that their
562 feedback is crucial to making a better system in future, and identifying the failure modes they saw
563 using the knowledge they have gained from the job. Supervisors in roles where scores are being used
564 by those they manage should receive training aimed at helping those under them feed-back on the
565 systems being used, and ensure that they encourage those under them to maintain a critical view of
566 all systems they work with. This should especially be emphasised for new starters—it is possible that
567 the levels of disagreement in the AFST system will decrease over time as the number of screeners who

568 have never known work without this system decreases. Keeping on top of these phenomena will be
 569 a continuous effort, but can be built into existing training schemes to reduce burden.

Recommendation 19 *The County should develop methods to monitor oversight, critical examination
 570 and the use of these scores by all those who have access to them or are tasked with approaching families
 for voluntary targeted services.*

Recommendation 20 *The County should ensure that investment decisions do not create undesirable
 571 over-reliance on scores, such as through under-investment in alternative sources of information, or
 through reduction of time available for each task due to efficiencies of automated systems.*

572 References

- 573 Ananny M and Crawford K, “Seeing without knowing: Limitations of the transparency ideal and its
 574 application to algorithmic accountability” [2016] *New Media & Society* doi: 10/gddxrg.
- 575 Ausloos J and Dewitte P, “Shattering one-way mirrors—data subject access rights in practice” (2018)
 576 8(1) *International Data Privacy Law* 4 doi: 10/cwcf.
- 577 Barocas S and Selbst AD, “Big Data’s Disparate Impact” (2016) 104 *California Law Review* 671 doi: 10/
 578 gfgq9w.
- 579 Binns R, Van Kleek M, Veale M, Lyngs U, Zhao J, and Shadbolt N, “‘It’s Reducing a Human Being to a Per-
 580 centage’; Perceptions of Justice in Algorithmic Decisions” in *Proceedings of the SIGCHI Conference
 581 on Human Factors in Computing Systems (CHI’18)* (ACM 2018) doi: 10/cvcp.
- 582 Chouldechova A, “Fair prediction with disparate impact: A study of bias in recidivism prediction instru-
 583 ments” (2017) 5(2) *Big Data* 153 doi: 10/gdcdqf.
- 584 Chouldechova A, Benavides-Prado D, Fialko O, and Vaithianathan R, “A Case Study of Algorithm-
 585 Assisted Decision Making in Child Maltreatment Hotline Screening Decisions” in *Conference on
 586 Fairness, Accountability and Transparency (FAT* 2018)* (2018).
- 587 Citron DK, “Technological due process” (2008) 85 *Washington University Law Review* 1249.
- 588 Dijkstra JJ, “User agreement with incorrect expert system advice” (1999) 18(6) *Behaviour & Informa-
 589 tion Technology* 399 doi: 10/fsnqm9.
- 590 Dwork C, Hardt M, Pitassi T, Reingold O, and Zemel R, “Fairness through awareness” in *Proceedings of
 591 the 3rd Innovations in Theoretical Computer Science Conference (ITCS ’12)* (2012) doi: 10/fzd3f9.
- 592 Edwards L and Veale M, “Slave to the Algorithm? Why a ‘Right to an Explanation’ is Probably Not The
 593 Remedy You Are Looking For” (2017) 16 *Duke L. & Tech. Rev.* 18 doi: 10/gdxtlj.
- 594 – “Enslaving the Algorithm: From a “Right to an Explanation” to a “Right to Better Decisions”?” (2018)
 595 16(3) *IEEE Security & Privacy* 46 doi: 10/gdz29v.

- 596 Eslami M, Karahalios K, Sandvig C, Vaccaro K, Rickman A, Hamilton K, and Kirlik A, “First I “Like” It, then
597 I Hide It: Folk Theories of Social Feeds” in *Proceedings of the SIGCHI Conference on Human Factors
598 in Computing Systems* (ACM 2016) doi: 10.1145/2858036.2858494.
- 599 Feldman M, Friedler SA, Moeller J, Scheidegger C, and Venkatasubramanian S, “Certifying and remov-
600 ing disparate impact” in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge
601 Discovery and Data Mining* (2015) doi: 10/gfgrbk.
- 602 Friedler SA, Scheidegger C, Venkatasubramanian S, Choudhary S, Hamilton EP, and Roth D, “A
603 Comparative Study of Fairness-Enhancing Interventions in Machine Learning” [2018] arXiv preprint
604 (<http://arxiv.org/abs/1802.04422>).
- 605 Gama J, Žliobaitė I, Bifet A, Pechenizkiy M, and Bouchachia A, “A survey on concept drift adaptation”
606 (2013) 1(1) ACM Comput. Surv. doi: 10/gd893p.
- 607 Gandy OH, “Engaging Rational Discrimination: Exploring Reasons for Placing Regulatory Constraints
608 on Decision Support Systems” (2010) 12(1) Ethics and Information Technology 29 doi: 10/bzwqrx.
- 609 Gayle D, “Schools Census Used to Enforce Immigration Laws, Minister Says” (*The Guardian*, January 13,
610 2019) (<https://www.theguardian.com/politics/2019/jan/13/schools-census-used-for-immigration-enforcement-minister-says>) visited on April 10, 2019.
- 612 Grove WM, Zald DH, Lebow BS, Snitz BE, and Nelson C, “Clinical versus mechanical prediction: a meta-
613 analysis” (2000) 12(1) Psychological Assessment.
- 614 Guadamuz A, “Viral Contracts or Unenforceable Documents? Contractual Validity of Copyleft Licenses”
615 (2004) 26(8) European Intellectual Property Review 331 (<https://ssrn.com/abstract=569101>).
- 616 Hajian S and Domingo-Ferrer J, “Direct and indirect discrimination prevention methods” in B Custers,
617 T Calders, B Schermer, and T Zarsky (eds.), *Discrimination and privacy in the information society*
618 (Springer 2012).
- 619 Hardt M, Price E, and Srebro N, “Equality of Opportunity in Supervised Learning” in DD Lee, M
620 Sugiyama, UV Luxburg, I Guyon, and R Garnett (eds.), *Advances in Neural Information Processing
621 Systems* 29 (Curran Associates, Inc 2016).
- 622 Hildebrandt M, “The Dawn of a Critical Transparency Right for the Profiling Era” in J Bus, M Crompton,
623 M Hildebrandt, and G Metakides (eds.), *Digital Enlightenment Yearbook* (IOS Press 2012) doi: 10/
624 cwpm.
- 625 Hoppe R, *The Governance of Problems: Puzzling, Powering and Participation* (Policy Press 2010).
- 626 Kamiran F and Calders T, “Classifying without discriminating” in *2nd International Conference on Com-
627 puter, Control and Communication, Karachi, Pakistan, 17–18 February, 2009* (2009) doi: 10/dtcfsn.
- 628 Lipsky M, *Street-level bureaucracy: Dilemmas of the individual in public services* (Russell Sage Founda-
629 tion 2010).
- 630 Manzoni J, “Big data in government: the challenges and opportunities” (*GOVUK*, February 2017) (<https://perma.cc/GF7B-5A2R>) visited on October 4, 2018.
631

- 632 Mitchell M et al., “Model Cards for Model Reporting” in *Proceedings of the 2nd ACM Conference on Fair-*
633 *ness, Accountability and Transparency (ACM FAT* 2019)* (ACM 2019) ([https://arxiv.org/abs/1810.](https://arxiv.org/abs/1810.03993)
634 [03993](https://arxiv.org/abs/1810.03993)).
- 635 Poursabzi-Sangdeh F, Goldstein DG, Hofman JM, Vaughan JW, and Wallach H, “Manipulating and Mea-
636 *suring Model Interpretability*” [2018] arXiv:1802.07810 [cs] (<http://arxiv.org/abs/1802.07810>)
637 visited on April 10, 2019.
- 638 Quiñonero-Candela J, Sugiyama M, Schwaighofer A, and Lawrence ND, *Dataset shift in machine learn-*
639 *ing* (The MIT Press 2009).
- 640 Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protec-
641 *tion of natural persons with regard to the processing of personal data and on the free movement of*
642 *such data, and repealing Directive 95/46/EC (General Data Protection Regulation)* [2016] OJ L119/1.
- 643 Richie F, “Access to sensitive data: Satisfying objectives rather than constraints” (2014) 30(3) *Journal*
644 *of Official Statistics*.
- 645 Romei A and Ruggieri S, “Discrimination data analysis: A multi-disciplinary bibliography” in B Custers,
646 T Calders, B Schermer, and T Zarsky (eds.), *Discrimination and Privacy in the Information Society*
647 (Springer 2012).
- 648 Skitka LJ, Mosier KL, and Burdick M, “Does automation bias decision-making?” (1999) 51 *International*
649 *Journal of Human-Computer Studies* 991 DOI: 10/bg5rb7.
- 650 Statistics New Zealand, “Integrated Data Infrastructure” (*Government of New Zealand*, 2016) ([https:](https://perma.cc/9RXL-SV7P)
651 [//perma.cc/9RXL-SV7P](https://perma.cc/9RXL-SV7P)) visited on October 4, 2018.
- 652 Strobl C, Boulesteix A.-L, Zeileis A, and Hothorn T, “Bias in random forest variable importance mea-
653 *sures: Illustrations, sources and a solution*” (2007) 8(1) *BMC Bioinformatics* 1 DOI: 10.1186/1471-
654 [2105-8-25](https://doi.org/10.1186/1471-2105-8-25).
- 655 Tversky A and Kahneman D, “Judgment under Uncertainty: Heuristics and Biases” (1974) 185(4157)
656 *Science* 1124.
- 657 Ustun B and Rudin C, “Supersparse Linear Integer Models for Optimized Medical Scoring Systems”
658 (2016) 102(3) *Machine Learning* 349 DOI: 10/f8crhw.
- 659 Veale M and Binns R, “Fairer machine learning in the real world: Mitigating discrimination without
660 *collecting sensitive data*” (2017) 4(2) *Big Data & Society* DOI: 10/gdcfnz.
- 661 Veale M, Binns R, and Edwards L, “Algorithms That Remember: Model Inversion Attacks and Data Pro-
662 *tection Law*” (2018) 376 *Phil. Trans. R. Soc. A* 20180083 DOI: 10/gfc63m.
- 663 Veale M and Brass I, “Administration by Algorithm? Public Management meets Public Sector Machine
664 *Learning*” in K Yeung and M Lodge (eds.), *Algorithmic Regulation* (Oxford University Press 2019).
- 665 Veale M, Van Kleek M, and Binns R, “Fairness and Accountability Design Needs for Algorithmic Support
666 *in High-Stakes Public Sector Decision-Making*” in *Proceedings of the SIGCHI Conference on Human*
667 *Factors in Computing Systems (CHI’18)* (ACM 2018) DOI: 10/ct4s.

- 668 Verwer S and Calders T, “Introducing positive discrimination in predictive models” in B Custers,
669 T Calders, B Schermer, and T Zarsky (eds.), *Discrimination and privacy in the information society*
670 (Springer 2013).
- 671 Zafar MB, Valera I, Gomez Rodriguez M, and Gummadi KP, “Fairness beyond Disparate Treatment & Dis-
672 parate Impact: Learning Classification without Disparate Mistreatment” in *Proceedings of the 26th*
673 *International Conference on World Wide Web* (International World Wide Web Conferences Steering
674 Committee 2017) DOI: 10/gfgq8r.