

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

# People underestimate the errors by algorithms for credit scoring and recidivism but tolerate even fewer errors: A representative study in Germany

PREPRINT – NOT PEER REVIEWED

Felix G. Rebitschek<sup>a,b,1</sup>  
Gerd Gigerenzer<sup>a, b</sup>  
Gert G. Wagner<sup>a,b,c</sup>

<sup>a</sup>Harding Center for Risk Literacy, Faculty of Health Sciences Brandenburg, University of Potsdam,  
Potsdam

<sup>b</sup>Max Planck Institute for Human Development, Berlin

<sup>c</sup>German Socio-Economic Panel Study (SOEP), Berlin

<sup>1</sup>Felix G. Rebitschek, **Email:** [rebitschek@uni-potsdam.de](mailto:rebitschek@uni-potsdam.de)

**Author Contributions:** F.G.R, G.G. and G.G.W. designed research; F.G.R, G.G. and G.G.W. performed research; F.G.R analyzed data; and F.G.R. and G.G.W. drafted the manuscript; all authors revised the manuscript.

**Competing Interest Statement:** The authors declare no competing interest.

**Classification:** Social Sciences: Psychological and Cognitive Sciences;

**Keywords:** algorithmic decision-making, algorithm acceptance, algorithm errors

29 **Abstract**

30 This study provides the first representative analysis of error estimations and error tolerance in a  
31 Western country (Germany) with regards to algorithmic decision-making systems (ADM). We examine  
32 people's expectations about the accuracy of algorithms that predict credit default, recidivism of an  
33 offender, suitability of a job applicant, and health behavior. Also, we ask whether expectations about  
34 algorithm errors vary between these domains and how they differ from expectations about errors made  
35 by human experts. In a nationwide representative study (N=3,086) we find that most respondents  
36 underestimated the actual errors made by algorithms and are willing to tolerate even fewer errors than  
37 estimated. Error estimates and error tolerance did not differ consistently for predictions made by  
38 algorithms or human experts, but people's living conditions (e.g. unemployment, household income)  
39 affect domain-specific tolerance (job suitability, credit defaulting) of misses and false alarms. We  
40 conclude that people have unwarranted trust in the competence of ADM systems and evaluate errors  
41 in terms of potential personal consequences. Given the general public's low error tolerance, we further  
42 conclude that acceptance of ADM appears to be conditional to strict accuracy requirements.

43

44 **Significance Statement**

- 45 - The general public in Germany overestimates the accuracy of algorithmic decision-making
- 46 systems (ADM) in predicting recidivism and credit scoring.
- 47 - At the same time the public tolerates even fewer ADM errors than actually occur.
- 48 - Error tolerance reflects personal conditions and is specific to error types and domains.

49

50 **Introduction**

51 This study provides the first representative analysis of error estimations and error tolerance in  
52 a Western population (in Germany) with regards to specific algorithmic decision-making (ADM)  
53 systems. We examine how accurately algorithms are expected to perform in predicting credit  
54 defaulting, recidivism of an offender, suitability of a job applicant, and health behavior.

55 Algorithmic decision-making (ADM (1)) continues to spread into everyday life. At the same  
56 time, claims, risks (2-4), and implementations related to ADM are under debate, e.g. in criminal risk  
57 assessment (5-7) or allocation of public resources (8). Despite these controversies, it remains unclear  
58 what the general public knows about ADM potential, other than that the concept of algorithms is little  
59 understood (9). The public debate centers on laypeople's trust in and fear of algorithms (10). Yet these  
60 emotions are not independent of knowledge; they could result from misjudging ADM potential.  
61 Overestimating the accuracy of decision support could be associated with unwarranted trust in the  
62 competence of algorithms (as opposed to affective trust (11) or trust in motives (12)), while fear could  
63 arise from overestimating algorithm failures.

64 Research on algorithm aversion and appreciation (1, 13) examines both the circumstances  
65 under which people trust algorithmic advice (13, 14)—for instance because they perceive the decision  
66 problem in question to be objective or to require mechanical skills (15) or because they lack  
67 confidence in their own expertise—and the circumstances under which they are mistrustful, for  
68 instance (16, 17), in response to slow algorithm responses or to observing algorithm errors.  
69 Surprisingly, the expected level of accuracy of algorithms and the perceived competence (12) of  
70 algorithmic advice are largely neglected in research on the general population. Fifty-eight percent of  
71 Americans expect some level of human bias in ADM systems, and 47% and 49% respectively believe  
72 that resume screening of job applicants and scoring for parole are effective (10). Similarly, the Dutch  
73 population often evaluates automatic decisions in favor of AI systems (e.g. as being more useful)  
74 compared with humans (18). In the present article, we aim to reduce this gap and investigate what the  
75 general public expects regarding the accuracy of ADM in the financial, legal, occupational, and health  
76 domains. We ask people what they believe are the actual error rates made by ADM and how many  
77 errors they consider to be acceptable (1, 17) and check whether their responses meets current ADM  
78 accuracy standards (in the financial and the legal domain). To the best of our knowledge, this is the  
79 first study comparing error estimates and error tolerance of ADM.

80 In the US, people's attitudes towards algorithms has been found to be highly context-specific  
81 (10). We therefore also examine whether the degree of error tolerance is associated with personal  
82 conditions, e.g. a risk-seeking preference (19). Classifying ADM systems can balance two different  
83 types of errors, misses and false alarms, each associated with different consequences (costs).  
84 Thinking about types of decision errors and related costs can affect tolerance of errors (e.g. 'bias' in  
85 signal detection theory (20); cost-sensitive error management (21)), and this tolerance may differ in  
86 the legal (22) and in the medical domain (23). From the perspective of an unemployed person, for  
87 instance, mistakenly being overlooked for a job by an ADM (a miss) is likely more costly than being  
88 hired in spite of being unsuitable (a false positive). In our analysis, we therefore relate assessed error  
89 tolerance to critical factors such as unemployment phases. Also, we compare error tolerance with the  
90 typical error preference of exemplary stakeholders (e.g. non-recidivating offenders want to avoid a  
91 false alarm). We additionally explore factors that may influence algorithm error tolerance and  
92 underlying attitudes towards technology, such as risk preference (24), gender (9, 25), and age (26).  
93 For instance, only one third of US-Americans above 50 years of age compared with half of those 18 to  
94 29 years of age believe that algorithms can be free of human biases (10).

95 Finally, we compare error estimation and tolerance regarding algorithmic and expert advice  
96 across domains. Earlier studies indicate that in the medical domain, people prefer human experts over  
97 automated care because, among others, they believe that individual circumstances are neglected by  
98 machines (27). People were also found to trust computers less than physicians to make good  
99 recommendations (28) and, crucial for error perception, to consider algorithm-based advice as being  
100 less accurate than a clinicians' advice (29). In the domain of people analytics, human interviewers are  
101 also assumed to be more accurate and useful than algorithmic decision-aids (30). At least for those  
102 two domains, the few existing studies thus indicate that people perceive algorithms to make more  
103 errors than human experts do. On the other hand, bestseller authors and commercial companies have  
104 aggressively promoted "AI"—which in some instances has demonstrated better performance for  
105 selective tasks (e.g. image classification(31) and deceptive text detection(32))—as being superior to  
106 human experts by underscoring the accomplishments of algorithms based on big data that have, for  
107 instance, beaten the best Go players and, as in the case of IBM's supercomputer Watson, promise to  
108 revolutionize health care. From that point of view, "it would be madness not to follow their advice" (33).

109 We conducted a survey of a large population-representative sample of members of private  
110 households in Germany (by means of the Innovation Sample of the German Socio Economic Panel

111 Study, SOEP IS). In the survey, 3,086 respondents were questioned about estimated and accepted  
112 error rates in current credit scoring (introductory task) before being randomly assigned to decision  
113 scenarios either on algorithms or on experts (between-subjects). Respondents were requested to  
114 provide estimated and accepted error rates in a people analytics problem (predicting a suitable job  
115 candidate), a legal problem (predicting recidivism of an offender), and a health problem (evaluating  
116 health behavior). For instance, respondents were asked with respect to ADM and recidivism: “In the  
117 US judiciary, offenders are regularly reviewed for early release. The computer program COMPAS  
118 assesses whether a criminal will recidivate and commit another crime within the next 2 years. Now  
119 please imagine a group of 100 offenders who are actually at risk of recidivism: How many of them do  
120 you think the computer program incorrectly assesses as not being at risk of recidivism [estimated miss  
121 rate]? What would you personally find acceptable: How many of these 100 offenders at most could be  
122 wrongly assessed as not being at risk of recidivism by the computer program [accepted miss rate]?  
123 Now please imagine a group of 100 offenders who are not at risk of recidivism: How many of them do  
124 you think the computer program wrongly considers to be at risk of recidivism [estimated false alarm  
125 rate]? What would you personally find acceptable: How many of these 100 offenders at most could be  
126 wrongly estimated to be at risk of recidivism by the computer program [accepted false alarm rate]?”  
127

## 128 **Results**

### 129 ***Estimated and accepted error rates of ADM systems***

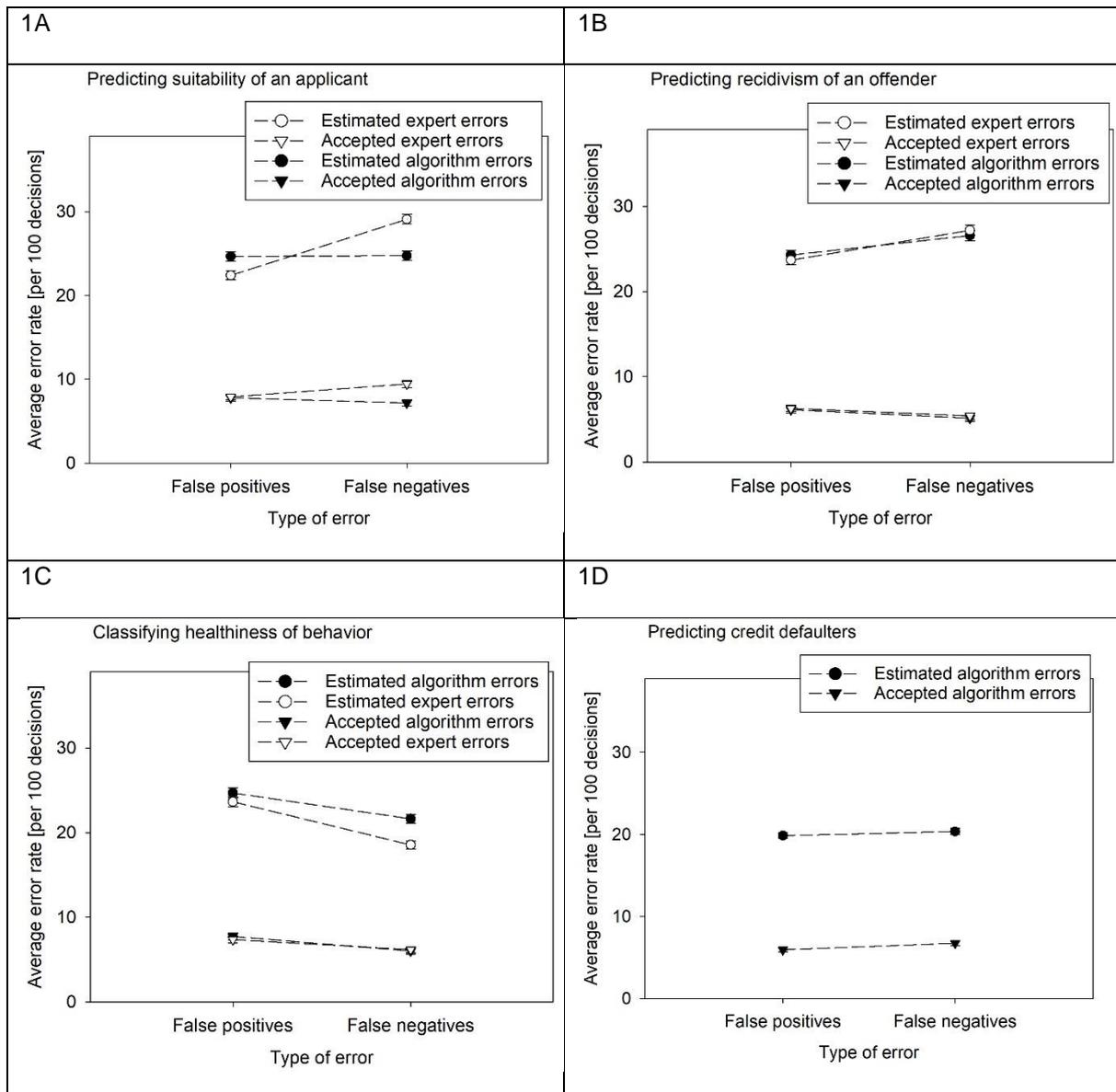
130 Across domains and types of errors, about one in every four algorithmic decisions was  
131 estimated to be wrong (Figures 1A-D), with the lowest estimates of errors for credit scoring (Figure  
132 1D). Only a few respondents believed in perfect algorithm (or expert) performance (Table S1).

133 Minimally more false negatives than false positives were estimated for the algorithmic  
134 prediction of recidivism ( $F(1,1314) = 15.75, P < .001, \eta_p^2 = 0.01$ ), but false negatives were considered  
135 less acceptable ( $F(1,1321) = 4.32, P = .038, \eta_p^2 < 0.01$ ). For the algorithmic evaluation of health  
136 behavior (Figure 1C), more false positives (falsely healthy) than false negatives were estimated  
137 ( $F(1,1297) = 39.17, P < .001, \eta_p^2 = 0.03$ ) and considered acceptable ( $F(1,1310) = 22.76, P < .001, \eta_p^2$   
138  $= 0.02$ ).

139 Notably, there is a consistent and large difference between the estimated and maximally  
140 accepted rates of algorithm error across all four domains. Expressed in percentage points, the  
141 differences in false positives (FP) and false negatives (FN) are for application suitability ( $\Delta_{FP} = 17, \Delta_{FN}$

142 = 17), for recidivism ( $\Delta_{FP} = 18$ ,  $\Delta_{FN} = 21$ ), for health behavior ( $\Delta_{FP} = 17$ ,  $\Delta_{FN} = 16$ ), and for credit scoring  
 143 ( $\Delta_{FP} = 14$ ,  $\Delta_{FN} = 14$ ).

144



145 Figures 1A-D: Estimated and accepted rates of false positives and false negatives related to the  
 146 assessments of job suitability (1A), recidivism (1B), and health behavior (1C) by algorithms and  
 147 experts, and to credit scoring (1D). Data are weighted for representativeness. Error bars show the  
 148 standard errors of the mean. For predicting credit defaulters, estimates of expert errors were not  
 149 elicited.

150

151

152

153 ***Estimated and actual errors of ADM systems***

154 To assess whether estimated algorithm errors correspond to actual errors, we compared  
155 respondents' estimates with actual performance metrics for recidivism and credit scoring algorithms  
156 (for the scenarios of health behavior and people analytics, reliable metrics are not publicly available).  
157 For credit scoring, most respondents' estimated error rates (false positives and false negatives) were  
158 substantially lower than those reported by the German SCHUFA, which is the major credit scorer in  
159 Germany. SCHUFA reports an accuracy of .81 of its (banking) credit score in terms of the value of the  
160 area under the receiver operating curve (AUC) (34). This is lower than the unweighted AUC of 0.90  
161 that we calculated based on median sensitivity (.90) and specificity (.90) of respondents' estimates.  
162 Figure 2A (left) shows each of the respondents' estimates for the false negative rate and false positive  
163 rate, expressed as sensitivity ( $1 - \text{false negative rate}$ ) and  $1 - \text{specificity}$  (false positive rate). Each dot  
164 corresponds to one respondent. The vast majority of dots are in the top left corner, that is, show a very  
165 high estimated sensitivity (low miss rate) and a low false alarm rate. The kinked line corresponds to  
166 the actual AUC as reported by the SCHUFA. Given that the SCHUFA does not provide data for the  
167 entire curve, we represent the value of .81 as a best case with assumed error balance. This is shown  
168 with the kinked line that is symmetric around the top-left to bottom-right diagonal. Relative to this line,  
169 the points above represent overestimation and the points below represent underestimation of the  
170 accuracy of the credit scoring algorithm.

171 The same analysis is shown for the prediction of recidivism in Figure 2B (right). The accuracy  
172 of the widely used COMPAS algorithm (General Recidivism Risk Scale) is lower than that of the credit  
173 scoring algorithms, with an AUC of .68 (5). Again, the vast majority of the respondents overestimated  
174 its accuracy. Respondents' median estimated sensitivity (.80) and specificity (.80) correspond to an  
175 unweighted AUC of .80. In other words, the performance of the algorithm does not meet the  
176 population's expectations as expressed in estimations and acceptance. We characterize this gap, and  
177 to whom it is relevant, in the following section.

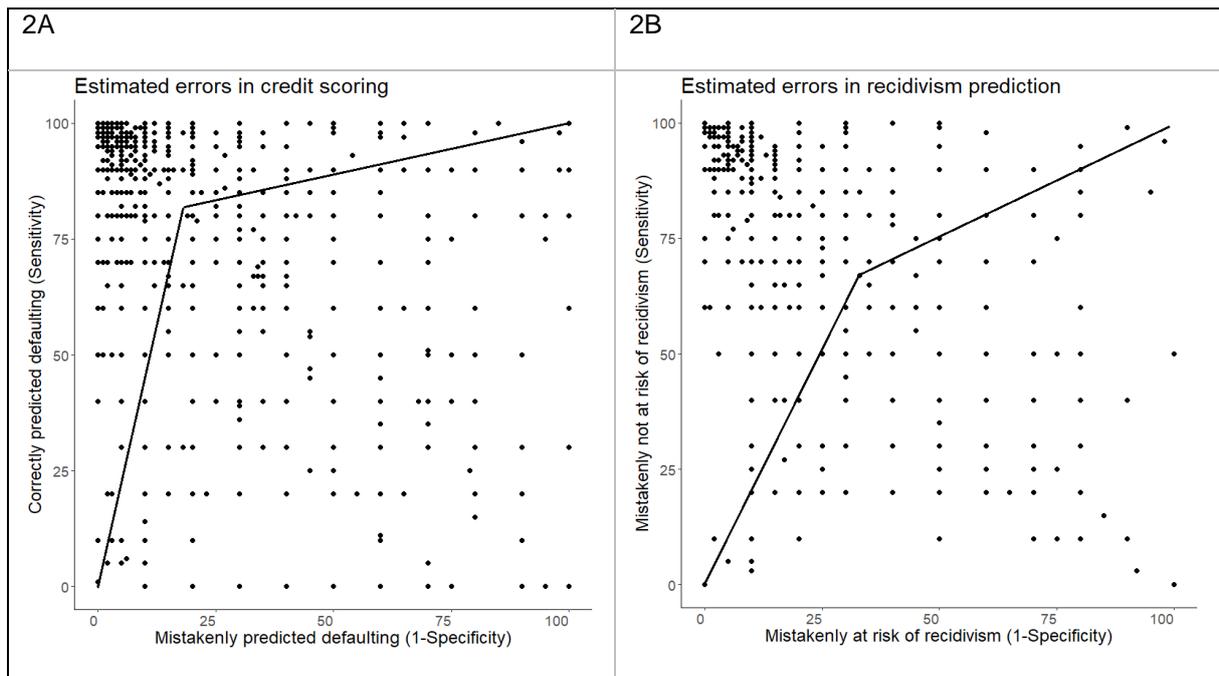
178

179

180

181

182



183 Figure 2: Respondents vastly overestimate the accuracy of algorithms for credit scoring and recidivism  
 184 prediction. The scatterplots show each respondent's estimates of the false positive rate (1 - specificity)  
 185 and false negative rate (shown as its complement, the sensitivity) in credit scoring ( $n = 2,740$ ) and in  
 186 recidivism prediction ( $n = 1,325$ ). Each point corresponds to one respondent. The kinked line shows  
 187 simplified ROC curves for current credit scoring accuracy (SCHUFA Credit Score algorithm,  
 188 AUC=0.81) and for recidivism prediction (COMPAS algorithm, AUC=0.68).

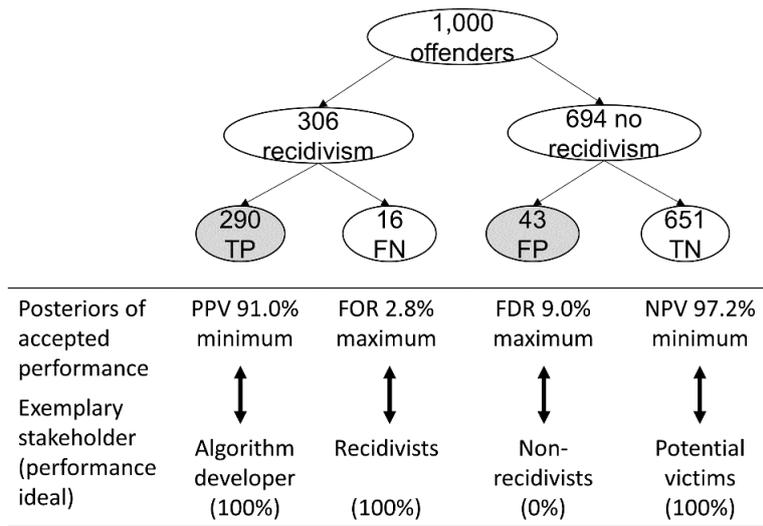
189

190 **Acceptance of different errors made by ADM systems**

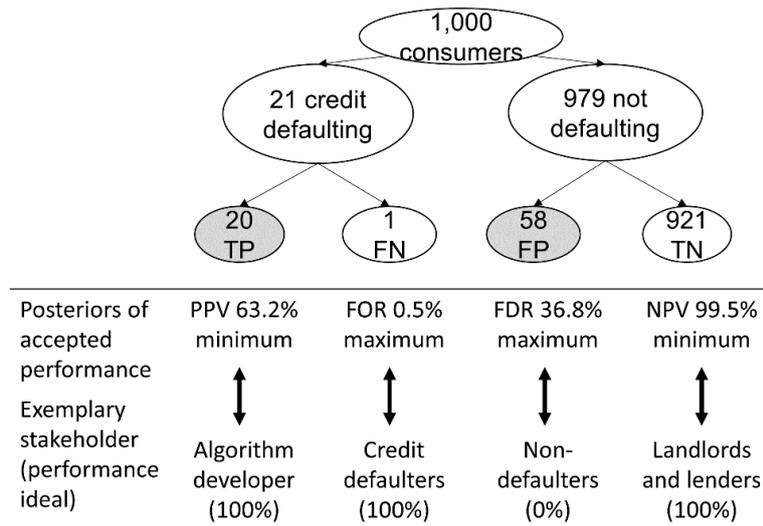
191 To interpret tolerated errors, their base rates have to be taken into consideration. What counts for an  
 192 individual decision are the posterior probabilities (how likely a certain prediction is correct or incorrect),  
 193 not the miss rates and false alarm rates. In a first step, by linking the rates of accepted algorithms  
 194 errors (weighted for representativeness) and the reported base rates of 30.6% for recidivism (within  
 195 four years, any offense, both genders; Table 2 in (35)) and of 2.1% for credit defaulting (within one  
 196 year (36)), we calculated the implied posterior probabilities. As the Figures 3A-B) show, the accepted  
 197 minimum positive predictive value (PPV) for recidivism is about 91%, whereas that of credit defaulting  
 198 is only 63%, despite the credit algorithm being more accurate than the one for recidivism.

199

3A



3B



200 Figures 3A-B: Natural frequency trees (37) of algorithm error acceptance for base rates of 30.6%  
 201 recidivists (3A) and of 2.1% credit defaulters (3B), and comparison of minimum performance  
 202 acceptance (posterior probabilities) of respondents with performance ideals of exemplary  
 203 stakeholders. Posteriors shown are the acceptable minimum PPV (positive predictive value; e.g. that a  
 204 person will commit another crime if testing positive), minimum NPV (negative predictive value: e.g.,  
 205 that a person will not commit another crime if testing negative), maximum FOR (false omission rate =  
 206 1-NPV), and maximum FDR (1 - PPV); TP (true positive), TN (true negative), FP (false positive), FN  
 207 (false negative). Note that depicted averages of individual posteriors are not equal to theoretical  
 208 posteriors based on the depicted aggregated TPs, FPs, TNs, and FNs.  
 209

210 In a second step, we compared the calculated posterior probabilities of error acceptance with  
211 ideal posteriors according to different stakeholders' goals (Figures 3A-B). For algorithm developers,  
212 the key goal is increasing PPV (precision) in the direction of 100% (38). An offender at risk of  
213 recidivism, in contrast, wants to falsely test negative and be released, that is, a test with a high FOR.  
214 In contrast to offenders at risk, those not at risk have an interest in not being falsely diagnosed as at  
215 risk of recidivism (low FDR). Finally, potential victims of defenders on bail will demand tests that  
216 maximize NPV towards 100%, that is, the probability that offenders will not commit another crime if  
217 they are tested negative and are released. The comparison showed that people's preferences for  
218 recidivism prediction overlap with performance ideals of developers, of potential victims, and of non-  
219 recidivating offenders, but not with recidivating offenders: False negatives were considered less  
220 acceptable than false positives.

221 Figure 3B provides similar examples of conflicting interests for predicting default. For instance,  
222 landlords and lenders do not want to contract with defaulting consumers; thus, their goal is to reach an  
223 NPV that approaches 100%. Customers at risk of default, in contrast, profit from tests that have low  
224 NPV, which is equivalent to high FOR. People's error tolerance did partially match performance ideals  
225 of the developers. Of note, our sample signaled substantial dissimilarity with non-defaulters because  
226 many false positives, and more of them than false negatives, were considered acceptable.

227 Based on the insight that posteriors in the interest of multiple stakeholders can only be  
228 partially achieved, it is examined in the following how error tolerance in general could depend on  
229 personal consequences from errors in specific domains.

230

### 231 ***Factors of domain-specific error acceptance***

232 A linear regression ( $F_{FP}(6, 1124) = 52.79, P < .001, R_{corr}^2 = 0.22$ ) that controlled for error  
233 estimations and age-, income-, and gender-specific differences in accepting algorithms' and experts'  
234 errors showed that reporting phases of unemployment in the past ten years (std.  $\beta = 0.10$ ) went along  
235 with generally increased acceptance of falsely assessed suitability. In contrast, overlooking suitable  
236 job candidates ( $F_{FN}(6, 1125) = 49.35, P < .001, R_{corr}^2 = 0.20$ ) was less accepted by men (std.  $\beta = -$   
237  $0.09$ ) and by older people (std.  $\beta = -0.08$ ).

238 Furthermore, linear regressions ( $F_{FN}(5, 1685) = 29.93, P < .001, R_{corr}^2 = 0.08$ ) revealed that  
239 misses (false negatives) of recidivating offenders were less accepted by younger respondents (std.  
240  $\beta_{FN} = 0.08$ ), those with higher income (std.  $\beta_{FN} = -0.07$ ), and women (std.  $\beta_{FN} = -0.08$ ). Similar

241 relationships were not observed for false positives, with  $F_{FP}(5, 1669) = 40.66, P < .001, R_{corr^2} = 0.11,$   
242 and generally not in the case of the evaluation of health behavior ( $F_{FP}(8, 1625) = 42.39, P < .001,$   
243  $R_{corr^2} = 0.17; F_{FN}(8, 1635) = 66.43, P < .001, R_{corr^2} = 0.24$ ), even if controlling for body mass index,  
244 health, and insurance status.

245 With regards to errors in credit scoring, fewer false negatives were accepted ( $F(5, 1733) =$   
246  $115.03, P < .001, R_{corr^2} = .25$ ) by older respondents (std.  $\beta_{FP} = -0.05$ ) and those reporting higher  
247 household income (std.  $\beta_{FP} = -0.07$ ). The latter also accepted fewer false positives (std.  $\beta_{FP} = -0.05$ ),  
248 ( $F(5,1738) = 88.10, P < .001, R_{corr^2} = .20$ ).

249 In sum, laypeople take into account potential benefits and harms of certain types of errors,  
250 regardless of who or what caused the error. In the last step, we examined the differences in  
251 estimations and acceptance of errors made by ADM systems and experts.

252

### 253 ***Experts vs. algorithms***

254 We analyzed whether estimated and accepted errors vary for ADM and human experts (Table S2). A  
255 3x2x2 ANOVA (advisor x type of error x domain) indicated no general difference ( $F(1,2446) = 0.14, P$   
256  $= .707$ ) between expert and algorithm error estimations but that these differed depending on the  
257 domain ( $F(1,4892) = 11.79, P < .001, \eta_p^2 = 0.01$ ) and on the type of error within the domain ( $F(2,4892)$   
258  $= 32.81, P < .001, \eta_p^2 = 0.01$ ). For algorithms, (4.6 percentage points) lower rates of false negatives  
259 were estimated in the case of suitable applicants ( $F(1,2622) = 78.50, P < .001, \eta_p^2 = 0.03$ ). However,  
260 minimally higher error rates were estimated for recognizing healthy behavior, where false positives  
261 and false negatives are estimated to be 1.1 and 2.0 percentage points higher, respectively, for ADM  
262 than for humans ( $F(1,2567) = 8.87, P < .001, \eta_p^2 < 0.01$ ).

263 Differences in acceptance of errors made by algorithms or experts depended on the domain  
264 ( $F(1,4896) = 3.08, P = .046, \eta_p^2 < 0.01$ ) and on the type of error within the domain ( $F(1,4896) = 3.45, P$   
265  $= .032, \eta_p^2 = 0.01$ ) (Table S2). For suitability prediction, there was a generally lower acceptance of  
266 algorithm errors in suitability prediction ( $F(1,2631) = 4.92, P = .027, \eta_p^2 < 0.01$ ); false negatives were  
267 also considered less acceptable than false positives in the case of algorithms ( $F(1,2631) = 12.06, P =$   
268  $.001, \eta_p^2 = 0.01$ ). For recidivism prediction, in contrast, differences in acceptance are limited. False  
269 positives seem to be considered less acceptable when caused by algorithms than by experts.  
270 Mistakenly predicting false recidivism was generally less accepted ( $F(1,2780) = 4.24, P = .040, \eta_p^2 <$

271 0.01). Notably, it was generally considered less acceptable to falsely diagnose unhealthy behavior  
272 than to overlook it ( $(F(1,2575) = 37.63, P < .001, \eta_p^2 = 0.01)$ ).

273 What factors contribute to differences in evaluating the performance of an algorithm and an  
274 expert? Generally, age (std.  $\beta = -0.10, P < .001$ ) and household income (std.  $\beta = -0.08, P = .007$ ) but  
275 not gender and risk preference (Figure S1), if controlled for estimated error rates (std.  $\beta = 0.55, P <$   
276  $.001$ ), were negatively related with acceptance of ADM errors ( $F(5,933) = 89.03, P < .001, R_{\text{corr}}^2 =$   
277  $0.32$ ). However, the pattern was similar for errors made by human experts.

278 An analysis of factors that lead to ADM-specific acceptance of errors confirms that people  
279 above the age of 29 ( $(F(3,2717) = 9.01, P < .001, \eta_p^2 = 0.01)$ ) accept a higher number of human than  
280 algorithm errors (about 1 to 2 more per 100 predictions), while those younger than 29 accept fewer  
281 human errors and more algorithm errors (a difference of about 4 in 100 predictions). Similar  
282 distinctions could not be confirmed for gender ( $(F(1,2721) = 0.20, P = .659)$ ) or for monthly household  
283 net income ( $F(4,1752) = 1.82, P = .122$ ).

284 Finally, an additional analysis of the relationship between estimations and acceptance, which  
285 were strictly correlated (between  $r = .27$  and  $.43$ ), did not indicate differential developments of error  
286 acceptance with increasing perception of errors by algorithms and experts (Figure S2).

287

## 288 Discussion

289 Our nationwide representative study in Germany delivered two major results. First, the vast  
290 majority of respondents underestimated error rates for credit scoring and recidivism prediction.  
291 Second, respondents tolerated even fewer algorithm errors than they expected to occur, only 5% to  
292 7% errors compared with their estimates of about 20% to 30% errors. This gap between tolerated error  
293 rate and both the perceived and actual accuracy of algorithms is striking. These findings indicate a  
294 lack of public trust in the *competence* of algorithms (confidence), which depends on knowledge about  
295 past performance. This is in line with the concerns of the American public about the acceptability of  
296 algorithms in recidivism prediction and resume screening of job applicants: More than 50% find them  
297 unacceptable (10). For one, they perceive those systems to be incapable of doing nuanced work  
298 across individuals.

299 The lack of trust in the competence of algorithms needs to be distinguished from a lack of trust  
300 in the *intentions* of those behind the algorithms (12), which is the common target of digitalization  
301 campaigns, such as in the EU White Paper on Artificial Intelligence (39). Our result suggests that

302 winning trust in ADM requires better accuracy of ADM systems, even beyond the performance  
303 overestimations of our respondents. Because the respondents accepted fewer errors by ADM than by  
304 human experts in some domains, the public may expect some ADM systems to be more error-free  
305 than humans in order to tolerate them. ADM was estimated to perform roughly at par with experts. The  
306 only exceptions were estimations that on the one hand algorithms detect suitable job candidates better  
307 than experts and that on the other hand experts (physicians) are perceived to make slightly fewer  
308 errors of both types in evaluating health behavior.

309         Anticipated consequences (e.g. costs) of different types of error appear to inform judgments  
310 about error tolerance: People who experienced phases of unemployment in the past ten years are  
311 more open for errors in their favor, namely false positives in predicting suitable job candidates. People  
312 with higher household income tolerate fewer errors in credit scoring. A higher tolerance of false  
313 positives than false negatives in predicting defaulting is noteworthy, given that 98% of the people in  
314 Germany are non-defaulters in the credit scoring system discussed. Given that the proportion of  
315 lenders, landlords, and entrepreneurs among the population is quite limited, future research should  
316 ask: Do people who tolerate false positives fear harms to the economy if algorithms were not that  
317 sensitive or do they simply ignore the personal consequences of being classified falsely as a  
318 defaulter?

319

### 320         *Research implications*

321 The general public's perception of the accuracy of algorithms and its general tolerance of ADM error  
322 have been a neglected issue in previous research, which has focused on issues such as people's  
323 opinions on discrimination, safety, and surveillance by algorithms, often implying the existence of high  
324 accuracy. Our representative analysis of error estimations and error tolerance of the German  
325 population shows that the accuracy of real ADM systems in credit scoring and recidivism prediction is  
326 too low from the public's perspective. Observed as well as labelled accuracy of predictive algorithms  
327 (40, 41), even if it is low (42), and also confidence information for each prediction (43) can improve  
328 trust in algorithms (although are users trusting independent of performance (44)). Contrary to studies  
329 that identify risk-seeking preferences to explain the preference of imperfect algorithms over human  
330 advice (19), our representative study across three domains did not find any indications of a  
331 relationship between risk preference and the acceptance of algorithm errors.

332 Similar to acceptance, the accuracy of predictive algorithms can often be improved by making  
333 them simpler (29). Simplicity implies transparency and understandability (e.g. explaining a prediction  
334 (42)), which are potential measures to increase trust in ADM systems (45). The design of such  
335 systems may take successful human strategies of decision making under uncertainty in the real world  
336 as a starting point (46). This is of particular importance because merely using interpretable models  
337 (47) and explaining complex models (48) does not necessarily form trust in ADM systems. A further  
338 avenue lies in equipping regulators with tools that enable lay-comprehensible algorithm assessments  
339 in terms of consequences, performance, and fairness. This calls for research that combines trust in  
340 competence, such as in the accuracy of algorithms, with trust in motivation, such as transparency.

341

#### 342 *Educational implications*

343 Currently, ADM systems are spreading through most life domains, including personalized medicine,  
344 consumer scoring, credit scoring, insurers, people analytics, learning analytics, and predictive policing.  
345 This cultural change requires citizens to acquire digital risk literacy (49), including the ability to  
346 evaluate the performance of predictive algorithms. Accordingly, schools need to teach students the  
347 requisite concepts to understand ADM systems, including input, output, consequences, features,  
348 feature weights, fairness, and performance-related properties such as false alarm rate, miss rate,  
349 positive predictive value, cross-validation, along with the skills to locate the relevant information.

350 Yet people should not only learn about the actual performance that is delivered by credit  
351 scoring and recidivism prediction. With regards to both algorithms and experts it is important to keep in  
352 mind that people are sensitive to types of error because these present different costs for them and  
353 others. Laypeople consider these potential costs when confronted with algorithmic decision support.  
354 Educational materials on decision-support systems should thus inform readers about the performance  
355 of algorithms. Facilitating this understanding, however, also requires a culture of transparency and  
356 trust created by commercial firms or imposed by regulators, where people gain easy access to  
357 relevant information about the reliability of algorithms in credit scoring, people scoring, and other  
358 domains (50).

359

#### 360 *Limitations*

361 One limitation to the present study is that the estimates for the actual accuracy of the credit scoring  
362 algorithm are self-reported by the SCHUFA. The accuracy estimates for COMPAS, in contrast, stem

363 from independent scientific studies. The SCHUFA is a commercial company, and one could argue that  
364 their estimate of their own algorithms might be inflated. We have no evidence of this possibility, but if it  
365 were the case, that would mean that the public overestimates the accuracy of the algorithm to an even  
366 greater extent than what we analyzed in Figure 2. In other words, the gap between estimated accuracy  
367 and actual accuracy would be even larger.

368 Nor can we exclude the possibility that people in Germany have a better sense of the  
369 accuracy of other ADM systems, as we could only compare respondents' estimates to the available  
370 accuracies for the domains of credit scoring and recidivism prediction. However, as long as there are  
371 no further learning possibilities or transparent information available for other ADM systems, we have  
372 no reason to assume that the overestimates are unique to recidivism and credit scoring.

373

#### 374 Materials and Methods

375 Data. The study data were collected as a part of the "Innovation Sample" of the German  
376 Socio-Economic Panel (SOEP IS) (51). For this longitudinal survey, which started in 2011, all adult  
377 and youth members of the same households are interviewed every year. The panel's  
378 representativeness relies on interviewing split-offs (younger panel members founding their own  
379 families). Our household sample with 3,086 respondents (52.1% female, M=50.4 years of age (SD =  
380 19.0)) is representative of the adult population in Germany (52).

381 Design. In accordance with an experimental between-subjects design, after responding to  
382 credit scoring items, respondents were randomly assigned to either algorithm or expert predictions  
383 (recidivism, job suitability, healthiness of behavior).

384 Measures. We used four partially known scenarios to introduce ADM systems, of which credit  
385 scoring is fully established in Germany, recidivism prediction is established outside Germany, and  
386 suitability prediction is currently being established in Germany; the evaluation of healthy behavior was  
387 derived from established health bonus programs.

388 Error rates for each algorithm and each expert were assessed with a normalized frequency  
389 format. The following is an item example for overlooking recidivism in algorithms: "In the US judiciary,  
390 offenders are regularly reviewed for early release. The computer program COMPAS assesses whether  
391 a criminal will recidivate and commit another crime within the next 2 years. Now please imagine a  
392 group of 100 offenders who are actually at risk of recidivism: How many of them do you think the  
393 computer program incorrectly assesses as not being at risk of recidivism?"; "What would you

394 personally find acceptable: At most, how many of these 100 criminals could be wrongly assessed as  
395 not being at risk of recidivism due to the computer program?" A validated self-report item about risk  
396 attitude (24) was included (11-point). Standard demographic variables in the panel that we used were  
397 gender, age, household income (per month), event of unemployment (in the last ten years), health  
398 insurance status (statutory or private), and self-rated health status (5-point).

399 Data were analyzed with regression techniques. All data were weighted on the person level to  
400 fine-tune for representativeness. All data (at the German Socio Economic Panel) and materials (at  
401 PNAS) are freely available.

402

#### 403 Acknowledgments

404 Support for this work was provided by the German Socio Economic Panel Study  
405 ([https://www.diw.de/en/diw\\_01.c.390440.en/soep\\_is.html](https://www.diw.de/en/diw_01.c.390440.en/soep_is.html)) and by a Max Planck Society Fellowship (to  
406 G.G.W.). Furthermore, F.G.R. was supported through the RisikoAtlas project ([www.risikoatlas.de](http://www.risikoatlas.de)),  
407 which was funded by the Federal Ministry of Justice and Consumer Protection (BMJV) on the basis of  
408 a resolution of the German Bundestag via the Federal Office for Agriculture and Food (BLE) within the  
409 framework of the Innovation Programme. We are grateful to Björn Meder for his advice on data  
410 presentation. We thank David Richter (SOEP) and his team as well as the team of the fieldwork  
411 organization Kantar Munich for implementing our survey questions in SOEP IS 2018/19. We thank our  
412 proofreader Rona Unrau for her critical eye.

413

414 The study was approved by the Institutional Ethics Board of the Max Planck Institute for Human  
415 Development, Berlin (Germany). It was carried out in accordance with the guidelines and regulations  
416 of the Max Planck Society for the Advancement of Science in Germany.

417

- 419 1. Burton JW, Stein M-K, & Jensen TB (2020) A systematic review of algorithm aversion in  
420 augmented decision making. *Journal of Behavioral Decision Making* n/a(n/a).
- 421 2. Russell SJ (2019) *Human Compatible: Artificial Intelligence and the Problem of Control*  
422 (Penguin).
- 423 3. Smith BC (2019) *The Promise of Artificial Intelligence: Reckoning and Judgment* (MIT Press).
- 424 4. Angwin J, Larson J, Mattu S, & Kirchner L (2016) Machine bias. *ProPublica*, May 23:2016.
- 425 5. Dressel J & Farid H (2018) The accuracy, fairness, and limits of predicting recidivism. *Science*  
426 *Advances* 4(1):eaao5580.
- 427 6. Kleinberg J, Lakkaraju H, Leskovec J, Ludwig J, & Mullainathan S (2018) Human decisions and  
428 machine predictions. *The Quarterly Journal of Economics* 133(1):237-293.
- 429 7. Stevenson MT & Doleac JL (2019) Algorithmic Risk Assessment in the Hands of Humans.  
430 *Available at SSRN*.
- 431 8. Lohninger T & Erd J (2019) SUBMISSION for the report to the UN General Assembly on digital  
432 technology, social protection and human rights. (Vienna).
- 433 9. Grzymek V & Puntschuh M (2019) What Europe Knows and Thinks About Algorithms Results  
434 of a Representative Survey. Bertelsmann Stiftung eupinions February 2019.
- 435 10. Smith A (2018) Public attitudes toward computer algorithms. in *Pew Research Center* (Pew  
436 Research Center).
- 437 11. McAllister DJ (1995) Affect-and cognition-based trust as foundations for interpersonal  
438 cooperation in organizations. *Academy of Management Journal* 38(1):24-59.
- 439 12. Twyman M, Harvey N, & Harries C (2008) Trust in motives, trust in competence: Separate  
440 factors determining the effectiveness of risk communication. *Judgment and Decision Making*  
441 3(1):111.
- 442 13. Logg JM, Minson JA, & Moore DA (2019) Algorithm appreciation: People prefer algorithmic to  
443 human judgment. *Organizational Behavior and Human Decision Processes* 151:90-103.
- 444 14. Castelo N, Bos MW, & Lehmann DR (2019) Task-dependent algorithm aversion. *Journal of*  
445 *Marketing Research* 56(5):809-825.
- 446 15. Lee MK (2018) Understanding perception of algorithmic decisions: Fairness, trust, and  
447 emotion in response to algorithmic management. *Big Data & Society*  
448 5(1):2053951718756684.
- 449 16. Efendic E, van de Calseyde P, & Evans A (2019) Slow decision speed undermines trust in  
450 algorithmic (but not human) predictions. *PrePrint*.
- 451 17. Dietvorst BJ, Simmons JP, & Massey C (2015) Algorithm aversion: People erroneously avoid  
452 algorithms after seeing them err. *Journal of Experimental Psychology: General* 144(1):114-  
453 126.
- 454 18. Araujo T, Helberger N, Kruike-meier S, & De Vreese CH (2020) In AI we trust? Perceptions  
455 about automated decision-making by artificial intelligence. *AI & Society* 1-13.
- 456 19. Jay Dietvorst B & Bharti S (2019) Risk Seeking Preferences Lead Consumers to Reject  
457 Algorithms in Uncertain Domains. in *ACR North American Advances*, eds Bagchi R, Block L,  
458 Lee L, & Duluth (Association for Consumer Research), pp 78-81.
- 459 20. Green DM & Swets JA (1966) *Signal detection theory and psychophysics* (Wiley, New York).
- 460 21. Haselton MG & Buss DM (2000) Error management theory: A new perspective on biases in  
461 cross-sex mind reading. *Journal of Personality and Social Psychology* 78(1):1-11.
- 462 22. Mitchell G & Garrett BL (2019) The impact of proficiency testing information and error  
463 aversions on the weight given to fingerprint evidence. *Behavioral Sciences & the Law*  
464 37(2):195-210.
- 465 23. Shiloh S (2010) An Experimental Investigation of the Effects of Acknowledging False Negative  
466 and False Positive Errors on Clients' Cancer Screening Intentions: The Lesser of Two Evils?  
467 *Applied Psychology: Health and Well-Being* 2(2):204-221.

- 468 24. Frey R, Pedroni A, Mata R, Rieskamp J, & Hertwig R (2017) Risk preference shares the  
469 psychometric structure of major psychological traits. *Science advances* 3(10):e1701381.
- 470 25. Pierson E (2017) Demographics and discussion influence views on algorithmic fairness. *arXiv*  
471 *preprint arXiv:1712.09124*.
- 472 26. Mossberger K, Tolbert CJ, & Stansbury M (2003) *Virtual inequality: Beyond the digital divide*  
473 (Georgetown University Press).
- 474 27. Longoni C, Bonezzi A, & Morewedge CK (2019) Resistance to Medical Artificial Intelligence.  
475 *Journal of Consumer Research* 46(4):629-650.
- 476 28. Promberger M & Baron J (2006) Do patients trust computers? *Journal of Behavioral Decision*  
477 *Making* 19(5):455-468.
- 478 29. Eastwood J, Snook B, & Luther K (2012) What people want from their professionals: Attitudes  
479 toward decision-making strategies. *Journal of Behavioral Decision Making* 25(5):458-468.
- 480 30. Diab DL, Pui S-Y, Yankelevich M, & Highhouse S (2011) Lay Perceptions of Selection Decision  
481 Aids in US and Non-US Samples. *International Journal of Selection and Assessment* 19(2):209-  
482 216.
- 483 31. He K, Zhang X, Ren S, & Sun J (2015) Delving deep into rectifiers: Surpassing human-level  
484 performance on imagenet classification. *Proceedings of the IEEE international conference on*  
485 *computer vision*, pp 1026-1034.
- 486 32. Ott M, Choi Y, Cardie C, & Hancock JT (2011) Finding deceptive opinion spam by any stretch  
487 of the imagination. *Proceedings of the 49th annual meeting of the association for*  
488 *computational linguistics: Human language technologies-volume 1*, (Association for  
489 Computational Linguistics), pp 309-319.
- 490 33. Harari YN (2016) *Homo Deus: A brief history of tomorrow* (Random House).
- 491 34. SCHUFA (2019) Zuverlässiger Score. Sichere Bank. - Der Schufa Score für Banken 3.0.
- 492 35. Brennan T, Dieterich W, & Ehret B (2009) Evaluating the predictive validity of the COMPAS  
493 risk and needs assessment system. *Criminal Justice and Behavior* 36(1):21-40.
- 494 36. SCHUFA (2019) Kredit Kompass 2019.
- 495 37. Gigerenzer G & Hoffrage U (1995) How to improve Bayesian reasoning without instruction:  
496 frequency formats. *Psychological Review* 102(4):684-704.
- 497 38. Demartini G & Mizzaro S (2006) A classification of IR effectiveness metrics. *European*  
498 *Conference on Information Retrieval*, (Springer), pp 488-491.
- 499 39. Ciucci M & Gouardères F (2020) The White Paper on Artificial Intelligence.
- 500 40. Yin M, Wortman Vaughan J, & Wallach H (2019) Understanding the effect of accuracy on  
501 trust in machine learning models. *Proceedings of the 2019 CHI Conference on Human Factors*  
502 *in Computing Systems*, pp 1-12.
- 503 41. Yu K, Berkovsky S, Taib R, Zhou J, & Chen F (2019) Do I trust my machine teammate? an  
504 investigation from perception to decision. *Proceedings of the 24th International Conference*  
505 *on Intelligent User Interfaces*, pp 460-468.
- 506 42. Lai V & Tan C (2019) On human predictions with explanations and predictions of machine  
507 learning models: A case study on deception detection. *Proceedings of the Conference on*  
508 *Fairness, Accountability, and Transparency*, pp 29-38.
- 509 43. Zhang Y, Liao QV, & Bellamy RK (2020) Effect of Confidence and Explanation on Accuracy and  
510 Trust Calibration in AI-Assisted Decision Making. *arXiv preprint arXiv:2001.02114*.
- 511 44. Springer A, Hollis V, & Whittaker S (2017) Dice in the black box: User experiences with an  
512 inscrutable algorithm. *2017 AAAI Spring Symposium Series*.
- 513 45. Ribeiro MT, Singh S, & Guestrin C (2016) "Why should I trust you?" Explaining the predictions  
514 of any classifier. *Proceedings of the 22nd ACM SIGKDD international conference on*  
515 *knowledge discovery and data mining*, pp 1135-1144.
- 516 46. Hafenbrädl S, Waeger D, Marewski JN, & Gigerenzer G (2016) Applied decision making with  
517 fast-and-frugal heuristics. *Journal of Applied Research in Memory and Cognition* 5(2):215-  
518 231.

- 519 47. Poursabzi-Sangdeh F, Goldstein DG, Hofman JM, Vaughan JW, & Wallach H (2018)  
520 Manipulating and measuring model interpretability. *arXiv preprint arXiv:1802.07810*.  
521 48. Cheng H-F, *et al.* (2019) Explaining decision-making algorithms through UI: Strategies to help  
522 non-expert stakeholders. *Proceedings of the 2019 CHI Conference on Human Factors in*  
523 *Computing Systems*, pp 1-12.
- 524 49. Gigerenzer G (2015) *Risk savvy: How to make good decisions* (Penguin).
- 525 50. O'Neill O (2018) Linking Trust to Trustworthiness. *International Journal of Philosophical*  
526 *Studies* 26(2):293-300.
- 527 51. Richter D & Schupp J (2015) The SOEP Innovation Sample (SOEP IS). *Schmollers Jahrbuch:*  
528 *Journal of Applied Social Science Studies/Zeitschrift für Wirtschafts-und Sozialwissenschaften*  
529 135(3):389-400.
- 530 52. Rohrer JM, Egloff B, & Schmukle SC (2015) Examining the effects of birth order on  
531 personality. *Proceedings of the National Academy of Sciences* 112(46):14224-14229.
- 532