

# FairTargetSim: An Interactive Simulator for Understanding and Explaining the Fairness Effects of Target Variable Definition

Dalia Gala<sup>1</sup>, Milo Phillips-Brown<sup>2,3</sup>, Naman Goel<sup>1</sup>, Carinal Prunkl<sup>4</sup>, Laura Alvarez Jubete<sup>5</sup>, medb corcoran<sup>5</sup> and Ray Eitel-Porter<sup>5</sup>

<sup>1</sup>University of Oxford

<sup>2</sup>University of Edinburgh

<sup>3</sup>Jain Family Institute

<sup>4</sup>Utrecht University

<sup>5</sup>Accenture

## Abstract

Machine learning requires defining one’s target variable for predictions or decisions, a process that can have profound implications on fairness: biases are often encoded in target variable definition itself, before any data collection or training. We present an interactive simulator, FairTargetSim (FTS), that illustrates how target variable definition impacts fairness. FTS is a valuable tool for algorithm developers, researchers, and non-technical stakeholders. FTS uses a case study of algorithmic hiring, using real-world data and user-defined target variables. FTS is open-source and available at: <http://tinyurl.com/ftsinterface>. The video accompanying this paper is here: <http://tinyurl.com/ijcaifts>.

## 1 Motivation

Machine learning requires translating real-world problems into numerical representations. Sometimes, the translation is straightforward—e.g. in predicting whether someone defaults on a loan. Other times, things are not so simple. When developing an algorithm to predict which job applicants will be good employees, for example, one must make precise the notion of a “good” employee. This is an ambiguous, subjective notion about which reasonable minds may disagree. How one translates this notion numerically—how one *defines the target variable*—can have profound implications for fairness [Passi and Barocas, 2019]. Defining “good” employee one way rather than another may result, e.g., in fewer applicants being hired from certain demographics. These issues arise in many domains. For a college admissions algorithm, one must determine who counts as a “good” student; for a search engine, one must determine what counts as a “good” search result; etc. How these notions are defined may likewise have weighty implications for fairness: which university applicants are admitted [Kizilcec and Lee, 2023]; which items appear at the top of search results [Phillips-Brown, manuscript]; etc. Target variable definition, then, is not a merely technical matter. Defining “good” employee, student, or search result is a value-laden process: it calls for close attention and transparency [Fazelpour and Danks, 2021].

But all too often, target variables are defined without transparency or attention to fairness. On one hand, technical developers may take target variable definition as a given, focusing instead on issues such as data quality, variance, accuracy of predictions, etc. On the other hand, stakeholders who are not a part of the technical process—like (hiring) managers in non-technical roles, or those working in upper management—either do not understand, or are simply unaware of, the implications of target variable definition in algorithmic settings. There is thus a pressing need for the fairness implication of target variable definition to be understood—and foregrounded—for stakeholders of all kinds.

To help meet this need, we developed an *interactive target variable simulator*, FairTargetSim (FTS): <http://tinyurl.com/ftsinterface>. FTS introduces its users to target variable definition, and reveals and explains its impact on fairness. FTS uses a case study: hiring algorithms. FTS invites the user to imagine that they are building a hiring algorithm, which mirrors a widely-used style of hiring algorithm based on psychometric tests. The user defines two target variables, using real-world psychometric test data from [Jaffe *et al.*, 2022]. With these two definitions, FTS builds two corresponding models and gives visualizations of how the models and training data differ in matters of fairness and overall performance. FTS’s interactive and visualization elements bring these issues to life—offering a more compelling and memorable illustration than one can get by, for example, reading a text.

FTS’s code is public and freely available. Its use is, then, not limited to hiring algorithms or to the dataset we use in our case-study: it can be extended to uses beyond education, and to different datasets and models.

## 2 FairTargetSim’s audience

FTS is a valuable tool for a wide range of audiences. The first target audience is technical developers who often want to develop algorithms responsibly but have less understanding of non-algorithmic factors such as target variable definition. With FTS, they can better the behavior of their abstract algorithms under different target variable definitions. This technical audience may also have less control over non-algorithmic factors, and can use FTS to better advocate—to decision-makers with non-technical backgrounds—for re-

81 sponsible algorithmic development. This leads us to the sec-  
 82 ond target audience: non-technical stakeholders: e.g. those  
 83 who use algorithms for making decisions or those who are  
 84 impacted by the decisions. When these stakeholders better  
 85 understand the fairness implications of target variable defini-  
 86 tion, the way is paved for more responsible and accountable  
 87 use of algorithms in the real world. The third target audience  
 88 is educators. There is a pressing need for more responsible  
 89 AI education and training in universities ([Grosz *et al.*, 2018],  
 90 [Kocec *et al.*, forthcoming]), government, and the private sec-  
 91 tor [Eitel-Porter, 2021]. The ethical implications of technical  
 92 issues can be challenging to explain to learners. FTS gives  
 93 educators an accessible, hands-on way to illustrate them.

94 We emphasize that FTS illustrates not “only” the fairness  
 95 implications decisions about target variable definition. It also  
 96 illustrates, more generally, the ethical implications of deci-  
 97 sions at the intersection of technical and non-technical as-  
 98 pects of algorithmic development. While it is well understood  
 99 among theorists that such decisions are value-laden ([Fried-  
 100 man and Nissenbaum, 1996], [Johnson, forthcoming]), they  
 101 often do not wear their ethical dimensions on their sleeves.  
 102 FTS allows audiences of all kinds to see—through a concrete  
 103 example—such decisions for what they are.

### 104 3 Related Work

105 A wealth of research has established the importance of un-  
 106 derstanding and addressing the fairness implications of target  
 107 variable definition—in algorithmic systems generally ([Passi  
 108 and Jackson, 2018], [Obermeyer *et al.*, 2019], [Martin Jr.  
 109 *et al.*, 2020], [Levy *et al.*, 2021], [Barocas *et al.*, 2023])  
 110 and hiring algorithms specifically ([Băzgu and Cernea, 2019],  
 111 [Raghavan *et al.*, 2020], [Tilmes, 2022]).

112 A number of systems have been developed for  
 113 practitioners—and in some cases, non-technical  
 114 stakeholders—to understand, identify, and address algo-  
 115 rithmic bias. We list just some, and note that various of them,  
 116 like FTS, have a visualization element: [Tramèr *et al.*, 2017],  
 117 [Bellamy *et al.*, 2019], [Ribeiro *et al.*, 2018], [Cabrera  
 118 *et al.*, 2019], [Microsoft and contributors, 2019], [Saleiro  
 119 *et al.*, 2019], [Vincent and ManyOthers, 2019], [Ahn and Lin,  
 120 2020], [Wexler *et al.*, 2020], [Johnson *et al.*, 2023], [Liu  
 121 *et al.*, 2023]. FTS is an important addition to these systems  
 122 because it is, to our knowledge, the only one that addresses  
 123 target variable definition.

124 Compared to previous demonstrations at IJCAI on related  
 125 subjects (e.g. [Sokol and Flach, 2018; Juan *et al.*, 2021;  
 126 Yu *et al.*, 2019; Miguel *et al.*, 2021; Henderson *et al.*, 2021;  
 127 Baumann *et al.*, 2023]), our demonstration will focus on the  
 128 problem of fairness implications of target variable definition.

## 129 4 Overview of FairTargetSim

130 FTS’s interface works with most modern browsers; Firefox is  
 131 advised. FTS has four pages that the user visits in order.

### 132 4.1 Key Concepts Explained

133 This page introduces target variable definition to a non-  
 134 technical audience, explains how it impacts fairness, and  
 135 gives an overview of the other pages of FTS.

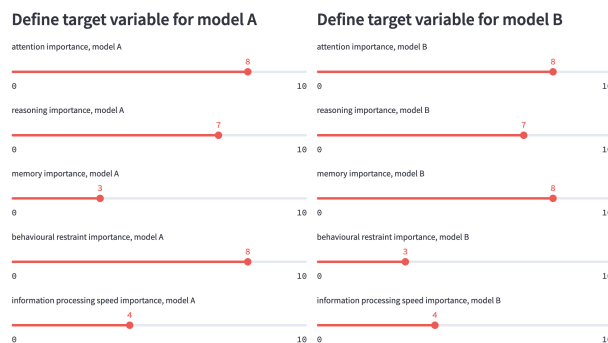


Figure 1: The user defines two target variables, using sliders representing the importance of traits of “good” employees.

### 4.2 User Defines Target Variables

This page has the user define two different target variables (Figure 1), which FTS uses to train two models, A and B.

In the real-world hiring algorithms that are based in cognitive tests, developers often define “good” employee by having an employer identify a group of current employees whom the employer deems “good” for a given role [Wilson *et al.*, 2021]. These employees then play cognitive-test games, and a model is trained to identify applicants that share cognitive traits with these employees.

FTS’s models are similar to these real-world systems in two key ways. First, like those systems, FTS uses support vector machine models to identify people who share cognitive traits with those who are identified as “good” employees. Second, FTS’s models are trained on data of real peoples’ cognitive tests; the data we use is from Jaffe *et al.*’s (2022) battery 26, which has eleven tests that we grouped into five traits: memory, information processing speed, reasoning, attention, and behavioural restraint.<sup>1</sup>

FTS’s models differ from the real-world systems in one key way: how the target variable is defined. With FTS, the user explicitly defines, using sliders depicted in Figure 1, how important the five cognitive traits are to what makes for a “good employee.” The user does this twice, creating two different target variables. Then FTS calculates the weighted average of tests scores, given the slider weightings, and assigns class label “0” to those in the bottom 85th percentile. From the top 15% subset, we randomly sample 100 “good” employees to whom we assign the class label “1” with weights ranging from 0.99 for the highest scoring candidate to 0.01 for the lowest scoring candidate, using linear distribution with the following equation for those in between:

$$f(x) = \frac{0.98}{1-n}x + \frac{0.01 - 0.99n}{1-n}$$

We assign a class label “0” to those not selected, thus intro-

<sup>1</sup>Our five categories are based on the following tests: *Memory* (forward memory span, reverse memory span, verbal list learning, delayed verbal list learning); *Information Processing Speed* (digit symbol coding, trail making part A, trail-making part B); *Reasoning* (arithmetic reasoning, grammatical reasoning); *Attention* (divided visual attention); and *Behavioral Restraint* (go/no-go).

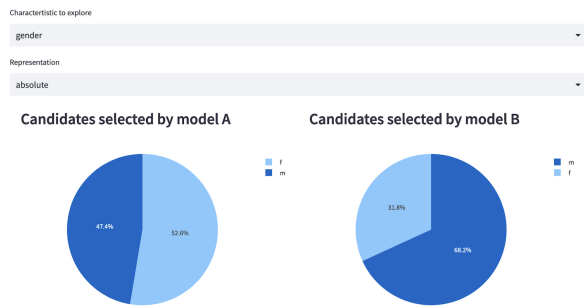


Figure 2: Charts display how the percentage of selected male and female applicants differs between models A and B.

169 ducing randomness. FTS then generates two labeled datasets  
 170 and corresponding models, each with different target variable  
 171 definitions.

172 FTS works with user-defined target variables because, first,  
 173 we do not have access to real-world target variables, and, sec-  
 174 ond, the lessons FTS offers are brought to life for the user  
 175 when she can see how her very own choices in target variable  
 176 definition can have implications for fairness. As we explain  
 177 further in Section 4.4, having user-defined target variables is  
 178 not a fundamental constraint on the idea of FTS; FTS can be  
 179 extended to use real-world labels when they are available.

### 4.3 Visualize Effects of Target Variable Definition

181 This page contains visualizations that illustrate how the user’s  
 182 two target variable definitions impact issues of fairness and  
 183 overall model performance. The visualizations are catego-  
 184 rized in to *Demographic* and *Non-demographic* sections, and  
 185 further divided into categories that (i) show features of the  
 186 models and (ii) features of the training data.

187 In the *Demographic* section, charts as in Figure 2 show  
 188 how models A and B differ in, e.g., the proportions of sel-  
 189 ected applicants across demographic groups (gender, educa-  
 190 tion level, age, and nationality—these are the demographic  
 191 groups that the Jaffe *et al.* dataset has information on). Other  
 192 charts show how the models differ across groups with respect  
 193 to “fairness metrics” ([Angwin *et al.*, 2016], [Corbett-Davies  
 194 and Goel, 2018]), such as true and false positive rates and  
 195 positive and negative predictive value.

196 The differences are stark: different target variable defini-  
 197 tions often result in major differences in the demographics  
 198 of selected applicants and in fairness metrics (see e.g. Fig-  
 199 ure 2). Visualizations in the *Demographics* section also show  
 200 how target variable definition affects models’ training data:  
 201 e.g. how positive and negative labels are distributed across  
 202 demographic groups.

203 In the *Non-demographic* section, visualizations show how  
 204 the models and training data differ in ways other than fair-  
 205 ness: e.g. how the models rank particular applicants (Figure  
 206 3), overall model confusion matrices, and accuracy metrics.

### 4.4 Further uses of FairTargetSim

208 This page gives recommendations for using FTS not just for  
 209 providing explanations and educating stakeholders, but also  
 210 for directly impacting practices in hiring and other domains.

Candidate ID	Ranking, model A	Ranking, model B	Predicted label A	Predicted label B
12	0.148166	0.987394	0	1
9	0.336614	0.978990	1	1
11	0.912242	0.960691	1	1
15	0.795276	0.953701	1	1
21	0.016506	0.952070	0	1

Figure 3: A table illustrates how individual applicants are evaluated differently by the two models.



Figure 4: Bar graphs show how choice of features of importance affects the model input scores achieved for different candidates depending on the demographic group—in this case, country of origin. For example, for model A, the median score for American candidates is approximately 0.57, while for model B, it is 0.63.

211 As noted, FTS’s code is available publicly; an organization  
 212 can extend FTS to use with their own data, models, and tar-  
 213 get variables. And, as also noted, in real-world target variable  
 214 definition, employers do not directly identify cognitive char-  
 215 acteristics of “good” employees; they identify certain current  
 216 employees as “good.” We give guidance on how to this in a  
 217 way that can promote fairness. For example, (i) consult vari-  
 218 ous managers on whom they judge “good;” these judgments  
 219 can be weighted in different ways—just as FTS weights the  
 220 cognitive tests in different ways—resulting in different target  
 221 variables. Or, (ii) use various performance metrics to evaluate  
 222 current employees (e.g., number of years to promotion, length  
 223 of tenure at a company, or role-specific metrics, such as num-  
 224 ber of sales with a sales role); these metrics can, again, be  
 225 weighted in different ways, resulting in different target vari-  
 226 ables. We also explain how to weight different judgements  
 227 and metrics in other domains: e.g. in a college admissions  
 228 algorithm or a search algorithm.

## 5 Future work

229 FTS opens up various avenues for future work, of which we  
 230 will highlight a few. One, as noted in Section 4.4, is to ap-  
 231 ply FTS to real-world hiring settings. Another, facilitated by  
 232 the fact that FTS is flexible and openly available, is to invite  
 233 the community to add more features to the simulator by, for  
 234 example using different kinds of data sets, models, or visual-  
 235 izations. Likewise, FTS could be extended to cases beyond  
 236 algorithmic hiring, such as college admissions or search en-  
 237 gines. Finally, FTS affords opportunities for human-centered  
 238 research. For example, user-studies could be run—with both  
 239 technical and non-technical stakeholders—to test how FTS  
 240 affects how they think about, develop, and use algorithms for  
 241 hiring and beyond.  
 242

## 243 References

- 244 [Ahn and Lin, 2020] Yongsu Ahn and Yu-Ru Lin. Fairsight:  
245 Visual analytics for fairness in decision making. *IEEE*  
246 *Transactions on Visualization and Computer Graphics*,  
247 26(1):1086–1095, 2020.
- 248 [Angwin *et al.*, 2016] Julia Angwin, Jeff Larson, Surya  
249 Matthu, and Lauren Kirchner. Machine bias.  
250 [https://www.propublica.org/article/machine-bias-risk-](https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing)  
251 [assessments-in-criminal-sentencing](https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing), May 23 2016.  
252 *ProPublica*.
- 253 [Bážgu and Cernea, 2019] Drago Bážgu and Mihail-  
254 Valentin Cernea. Algorithmic bias in current hiring  
255 practices: an ethical examination. *Proceedings of the*  
256 *International Management Conference*, 13(1):1068–1073,  
257 2019.
- 258 [Barocas *et al.*, 2023] Solon Barocas, Moritz Hardt, and  
259 Arvind Narayanan. *Fairness and Machine Learning: Lim-*  
260 *itations and Opportunities*. MIT Press, 2023.
- 261 [Baumann *et al.*, 2023] Joachim Baumann, Alessandro  
262 Castelnovo, Andrea Cosentini, Riccardo Crupi, Nicole  
263 Inverardi, and Daniele Regoli. Bias on demand: in-  
264 vestigating bias with a synthetic data generator. In  
265 *32nd International Joint Conference on Artificial Intelli-*  
266 *gence (IJCAI), Macao, SAR, 19-25 August 2023*, pages  
267 7110–7114. International Joint Conferences on Artificial  
268 Intelligence Organization, 2023.
- 269 [Bellamy *et al.*, 2019] R. K. E. Bellamy, K. Dey, M. Hind,  
270 S. C. Hoffman, S. Houde, K. Kannan, P. Lohia, J. Mar-  
271 tino, S. Mehta, A. Mojsilović, S. Nagar, K. Natesan Rama-  
272 murthy, J. Richards, D. Saha, P. Sattigeri, M. Singh, K. R.  
273 Varshney, and Y. Zhang. Ai fairness 360: An extensible  
274 toolkit for detecting and mitigating algorithmic bias. *IBM*  
275 *Journal of Research and Development*, 63(4/5):4:1–4:15,  
276 2019.
- 277 [Cabrera *et al.*, 2019] Angel Alexander Cabrera, Will Epper-  
278 son, Fred Hohman, Minsuk Kahng, Jamie Morgenstern,  
279 and Duen Horng Chau. Fairvis: Visual analytics for dis-  
280 covering intersectional bias in machine learning. In *2019*  
281 *IEEE Conference on Visual Analytics Science and Tech-*  
282 *nology (VAST)*. IEEE, October 2019.
- 283 [Corbett-Davies and Goel, 2018] Sam Corbett-Davies and  
284 Sharad Goel. The measure and mismeasure of fair-  
285 ness: A critical review of fair machine learning. *CoRR*,  
286 abs/1808.00023, 2018.
- 287 [Eitel-Porter, 2021] Ray Eitel-Porter. Beyond the promise:  
288 implementing ethical AI. *AI and Ethics*, 1:73–80, 2021.
- 289 [Fazelpour and Danks, 2021] Sina Fazelpour and David  
290 Danks. Algorithmic bias: Senses, sources, solutions. *Phi-*  
291 *losophy Compass*, 16(8), 2021.
- 292 [Friedman and Nissenbaum, 1996] Batya Friedman and He-  
293 len Nissenbaum. Bias in computer systems. *ACM Trans-*  
294 *actions on Information Systems*, 3(14):330–347, 1996.
- 295 [Grosz *et al.*, 2018] Barbara J. Grosz, David Gray Grant,  
296 Kate Vredenburg, Jeff Behrends, Lily Hu, Alison Sim-  
mons, and Jim Waldo. Embedded ethics: Integrating ethics  
broadly across computer science education, 2018.
- [Henderson *et al.*, 2021] Jette Henderson, Shubham Sharma,  
Alan Gee, Valeri Alexiev, Steve Draper, Carlos Marin,  
Yessel Hinojosa, Christine Draper, Michael Perng, Luis  
Aguirre, et al. Certifai: a toolkit for building trust in ai  
systems. In *Proceedings of the Twenty-Ninth International*  
*Conference on International Joint Conferences on Artifi-*  
*cial Intelligence*, pages 5249–5251, 2021.
- [Jaffe *et al.*, 2022] Paul I. Jaffe, Aaron Kaluszka, Nicole F.  
Ng, and Robert J. Schafer. A massive dataset of the neu-  
rocognitive performance test, a web-based cognitive as-  
sessment. *Scientific Data*, 9(1), 2022.
- [Johnson *et al.*, 2023] Brittany Johnson, Jesse Bartola, Rico  
Angell, Sam Witty, Stephen Giguere, and Yuriy Brun.  
Fairkit, fairkit, on the wall, who’s the fairest of them all?  
supporting fairness-related decision-making. *EURO Jour-*  
*nal on Decision Processes*, 11:100031, 2023.
- [Johnson, forthcoming] Gabbrielle M. Johnson. Are algo-  
rithms value-free? feminist theoretical virtues in machine  
learning. *Journal Moral Philosophy*, pages 1–35, forth-  
coming.
- [Juan *et al.*, 2021] Yi-Ning Juan, Yi-Shyuan Chiang, Shang-  
Chuan Liu, Ming-Feng Tsai, and Chuan-Ju Wang. Hive:  
Hierarchical information visualization for explainability.  
In *IJCAI*, pages 4988–4991, 2021.
- [Kizilcec and Lee, 2023] René F. Kizilcec and Hansol Lee.  
Algorithmic fairness in education. In Wayne Holmes and  
Kaška Porayska-Pomsta, editors, *The Ethics of Artificial*  
*Intelligence in Education*. Routledge, 2023.
- [Kopec *et al.*, forthcoming] Matthew Kopec, Meica Mag-  
nani, Vance Ricks, Roben Torosyan, John Basl, Nicholas  
Miklaucic, Felix Muzny, Ronald Sandler, Christo Wilson,  
Adam Wisniewski-Jensen, Cora Lundgren, Kevin Mills,  
and Mark Wells. The effectiveness of embedded values  
analysis modules in computer science education: An em-  
pirical study. <https://arxiv.org/abs/2208.05453>, forthcom-  
ing. forthcoming in *Nature Machine Intelligence*.
- [Levy *et al.*, 2021] Karen Levy, Kyla E. Chasalow, and Sarah  
Riley. Algorithms and decision-making in the pub-  
lic sector. *Annual Review of Law and Social Science*,  
17(1):309–334, October 2021.
- [Liu *et al.*, 2023] Jessica Liu, Huaming Chen, Jun Shen, and  
Kim-Kwang Raymond Choo. Faircompass: Operational-  
ising fairness in machine learning, 2023.
- [Martin Jr. *et al.*, 2020] Martin Martin Jr., Vinodkumar Prab-  
hakaran, Jill Kuhlberg, Andrew Smart, and William S.  
Isaac. Participatory problem formulation for fairer ma-  
chine learning through community based system dynam-  
ics, 2020.
- [Microsoft and contributors, 2019] Microsoft and contribu-  
tors. Fairlearn. <https://fairlearn.github.io/>, 2019.
- [Miguel *et al.*, 2021] Beatriz San Miguel, Aisha Naseer, and  
Hiroya Inakoshi. Putting accountability of ai systems into  
practice. In *Proceedings of the Twenty-Ninth International*

- 352 *Conference on International Joint Conferences on Artificial Intelligence*, pages 5276–5278, 2021.
- 353
- 354 [Obermeyer *et al.*, 2019] Ziad Obermeyer, Brian Powers,  
355 Christine Vogeli, and Sendhil Mullainathan. Dissecting  
356 racial bias in an algorithm used to manage the health of  
357 populations. *Science*, 366(6464):447–453, 2019.
- 358 [Passi and Barocas, 2019] Samir Passi and Solon Barocas.  
359 Problem formulation and fairness. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 39–48, 2019.
- 360
- 361
- 362 [Passi and Jackson, 2018] Samir Passi and Steven J. Jackson.  
363 Trust in data science: Collaboration, translation, and ac-  
364 countability in corporate data science projects. *Proc. ACM Hum.-Comput. Interact.*, 2(CSCW), nov 2018.
- 365
- 366 [Phillips-Brown, manuscript] Milo Phillips-Brown. Al-  
367 gorithmic neutrality. <https://arxiv.org/abs/2303.05103>,  
368 manuscript.
- 369 [Raghavan *et al.*, 2020] Manish Raghavan, Solon Barocas,  
370 Jon Kleinberg, and Karen Levy. Mitigating bias in algo-  
371 rithmic hiring: evaluating claims and practices. In *Pro-  
372 ceedings of the 2020 Conference on Fairness, Account-  
373 ability, and Transparency*. ACM, January 2020.
- 374 [Ribeiro *et al.*, 2018] Marco Tulio Ribeiro, Sameer Singh,  
375 and Carlos Guestrin. Anchors: High-precision model-  
376 agnostic explanations. *Proceedings of the AAAI Confer-  
377 ence on Artificial Intelligence*, 32(1), Apr. 2018.
- 378 [Saleiro *et al.*, 2019] Pedro Saleiro, Benedict Kuester, Loren  
379 Hinkson, Jesse London, Abby Stevens, Ari Anisfeld,  
380 Kit T. Rodolfa, and Rayid Ghani. Aequitas: A bias and  
381 fairness audit toolkit, 2019.
- 382 [Sokol and Flach, 2018] Kacper Sokol and Peter A Flach.  
383 Glass-box: Explaining ai decisions with counterfactual  
384 statements through conversation with a voice-enabled vir-  
385 tual assistant. In *IJCAI*, pages 5868–5870, 2018.
- 386 [Tilmes, 2022] Nicolas Tilmes. Disability, fairness, and al-  
387 gorithmic bias in ai recruitment. *Ethics and Information  
388 Technology*, 21(24), 2022.
- 389 [Tramèr *et al.*, 2017] Florian Tramèr, Vaggelis Atlidakis,  
390 Roxana Geambasu, Daniel Hsu, Jean-Pierre Hubaux,  
391 Mathias Humbert, Ari Juels, and Huang Lin. Fairtest: Dis-  
392 covering unwarranted associations in data-driven applica-  
393 tions. In *2017 IEEE European Symposium on Security and  
394 Privacy (EuroS&P)*, pages 401–416, 2017.
- 395 [Vincent and ManyOthers, 2019] Matthijs Vincent and  
396 ManyOthers. Ccikit-fairness. <https://github.com/koaning/scikit-fairness>, 2019.
- 397
- 398 [Wexler *et al.*, 2020] James Wexler, Mahima Pushkarna,  
399 Tolga Bolukbasi, Martin Wattenberg, Fernanda Viégas,  
400 and Jimbo Wilson. The what-if tool: Interactive probing  
401 of machine learning models. *IEEE Transactions on Visu-  
402 alization and Computer Graphics*, 26(1):56–65, 2020.
- 403 [Wilson *et al.*, 2021] Christo Wilson, Avijit Ghosh, Shan  
404 Jiang, Alan Mislove, Lewis Baker, Janelle Szary, Kelly  
Trindel, and Frida Polli. Building and auditing fair algo- 405  
rithms: A case study in candidate screening. In *Proceed- 406  
ings of the 2021 ACM Conference on Fairness, Account- 407  
ability, and Transparency*, FAccT ’21, page 666–677, 408  
2021. 409
- [Yu *et al.*, 2019] Han Yu, Yang Liu, Xiguang Wei, Chuyu 410  
Zheng, Tianjian Chen, Qiang Yang, and Xiong Peng. Fair 411  
and explainable dynamic engagement of crowd workers. 412  
In *IJCAI*, pages 6575–6577, 2019. 413