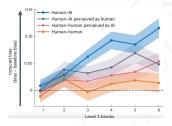
# Do Existing FAIRNESS Measures Suffice? Assessing Discrimination in Algorithmic Decision-Making

Yijun BIAN

Department of Computer Science University of Copenhagen

27 August 2025

# AI/ML is everywhere now





Social identity biases exist not only in human psychology and social behaviour, but also are present in artificial intelligence (AI) systems.<sup>1</sup>

When humans and AI interact, even minute perceptual, emotional and social biases<sup>2</sup>—originating either from AI systems or humans—leave human beliefs more biased, potentially forming a feedback loop.<sup>3</sup>

<sup>&</sup>lt;sup>1</sup>Ziad Obermeyer et al. "Dissecting racial bias in an algorithm used to manage the health of populations". In: Science 366.6464 (2019), pp. 447–453.

DOI: 10.1126/science.aax2342; Tiancheng Hu et al. "Generative language models exhibit social identity biases". In: Nat Comput Sci (2024), pp. 1–11; Richard J Chen et al. "Algorithmic fairness in artificial intelligence for medicine and healthcare". In: Nat Biomed Eng 7.6 (2023), pp. 719–742.

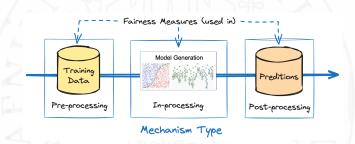
<sup>&</sup>lt;sup>2</sup>Are Skeie Hermansen et al. "Immigrant–native pay gap driven by lack of access to high-paying jobs". In: *Nature* (2025), pp. 1–7.

<sup>&</sup>lt;sup>3</sup>Moshe Glickman and Tali Sharot. "How human-AI feedback loops alter human perceptual, emotional and social judgements". In: *Nat Hum Behav* 9 (2025), pp. 345–359. DOI: 10.1038/s41562-024-02077-2; Madalina Vlasceanu and David M Amodio. "Propagation of societal gender inequality by internet search algorithms". In: *Proc Natl Acad Sci U.S.A.* 119.29 (2022), e2204529119.

# Existing work about fairness

Many sources of bias <sup>4</sup>

Mechanisms to enhance fairness<sup>5</sup>



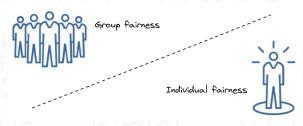
<sup>&</sup>lt;sup>4</sup>Unintential: Limited and coarse features, sample size disparity (less data by definition about minority populations), <u>skewed sample (feedback loops)</u>, <u>tainted examples</u>, features that act as proxies; <u>Intentional</u>: conscious prejudice.

<sup>&</sup>lt;sup>5</sup>Pre- and post-processing mechanisms normally function by manipulating input or output, while inprocessing mechanisms introduce fairness constraints into training procedures or algorithmic objectives.

# **Existing work about fairness**

Many sources of bias

Mechanisms to enhance fairness Types of fairness measures <sup>4,5</sup>



**Challenging**: incompatibility,<sup>6</sup> multi-attribute protection, etc.

<sup>&</sup>lt;sup>4</sup>Distributive fairness: group fairness, individual fairness, counterfactual fairness, etc.; Procedural fairness

<sup>&</sup>lt;sup>5</sup>Group fairness focuses on statistical/demographic equality among groups defined by sensitive attributes, while *individual fairness* follows a principle that "similar individuals should be evaluated or treated similarly."

<sup>&</sup>lt;sup>6</sup>Tensions between notions of fairness, between fairness and accuracy, between different methods for achieving fairness

## Three statistical non-discrimination criteria

#### Statistical non-discrimination criteria are

properties of the joint distribution of a sensitive attribute (SA, aka. protected attribute) A, target variable y, the classifier  $f(\cdot)$  or score R, sometimes including features X.

#### The three key criteria<sup>7</sup> are:

- independence
- separation
- sufficiency

Random variables (A, R) satisfy independence if  $A \perp R$  Random variables (R, A, Y) satisfy separation if  $R \perp A \mid Y$  Random variables (R, A, Y) satisfy sufficiency if  $Y \perp A \mid R$ 

These criteria are rarely satisfied all at once, except in degenerate cases.<sup>8</sup>

<sup>&</sup>lt;sup>7</sup>In essence, the latter two require the same recall or precision for each group, respectively.

<sup>&</sup>lt;sup>8</sup>Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and machine learning: Limitations and opportunities*. Cambridge, MA, USA: MIT Press, 2023; Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and machine learning*. fairmlbook.org, 2019; Reuben Binns. "Fairness in machine learning: Lessons from political philosophy". In: *FAT*. PMLR. 2018, pp. 149–159.

# #1 Independence

Random variables (A, R) satisfy independence if A  $\perp$  R



#1 statistical non-discrimination criterion

**Demographic parity** (DP),<sup>9</sup> aka. statistical parity<sup>10</sup>

$$\mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid a_1 = 0) = \mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid a_1 = 1). \tag{1}$$

<sup>&</sup>lt;sup>9</sup>Pratik Gajane and Mykola Pechenizkiy. "On formalizing fairness in prediction with machine learning". In: *FAT/ML*. 2018; Ray Jiang et al. "Wasserstein fair classification". In: *UAI*. PMLR. 2020, pp. 862–872.

<sup>10</sup> Cynthia Dwork et al. "Fairness through awareness". In: ITCS. ITCS '12. Cambridge, Massachusetts: ACM, 2012, pp. 214–226. ISBN: 9781450311151; Alexandra Chouldechova. "Fair prediction with disparate impact: A study of bias in recidivism prediction instruments". In: Big Data 5.2 (2017), pp. 153–163.

# #1 Independence

Disparate impact (i.e., "80% rule"): 11 the larger the better

$$\frac{\mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid a_1 = 0)}{\mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid a_1 = 1)} \leqslant \tau = 0.8.$$

Disparate treatment<sup>12</sup>

$$\mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid a_1) = \mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1), \ a_1 \in \{0, 1\}.$$
(3)

It is also indicated as "statistical parity (SP)" in the literature, <sup>13</sup> that is,

$$\mathbb{P}(f(\mathbf{x}, a_1) = 1 \mid \mathbf{a_1} = \mathbf{j}) = \mathbb{P}(f(\mathbf{x}, a_1) = 1), \ \forall \mathbf{j} \in \mathcal{A}_1 = \{1, 2, ..., n_{a_1}\}.$$

<sup>&</sup>lt;sup>11</sup>Michael Feldman et al. "Certifying and removing disparate impact". In: SIGKDD. 2015, pp. 259–268; Muhammad Bilal Zafar et al. "Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment". In: WWW. 2017, pp. 1171–1180.

<sup>&</sup>lt;sup>12</sup>Zafar et al., see n. 11; Christian Haas. "The price of fairness - A framework to explore trade-offs in algorithmic fairness". In: ICIS. Association for Information Systems. 2019.

<sup>&</sup>lt;sup>13</sup>Haas, see n. 12; Sam Corbett-Davies et al. "Algorithmic decision making and the cost of fairness". In: SIGKDD. New York, NY, USA: ACM, 2017, pp. 797–806. ISBN: 9781450348874; Alekh Agarwal, Miroslav Dudík, and Zhiwei Steven Wu. "Fair regression: Quantitative definitions and reduction-based algorithms". In: ICML. vol. 97. PMLR. 2019, pp. 120–129.

# #2 Separation

#### Random variables (R, A, Y) satisfy separation if $R \perp A \mid Y$

G

#2 statistical non-discrimination criterion

Equalised odds (EO)<sup>14</sup>

$$\mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid a_1 = 0, y) = \mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid a_1 = 1, y), y \in \{0, 1\}.$$
(4)

Equality of opportunity (EOpp, or EO), 15 aka. treatment equality 16

$$\mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid \mathbf{a}_1 = \mathbf{0}, y = 1) = \mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid \mathbf{a}_1 = \mathbf{1}, y = 1).$$
 (5)

<sup>&</sup>lt;sup>14</sup>Moritz Hardt, Eric Price, and Nathan Srebro. "Equality of opportunity in supervised learning". In: NIPS. vol. 29. Red Hook, NY, USA, 2016, pp. 3323–3331; Haas, see n. 12.

<sup>&</sup>lt;sup>15</sup>Hardt, Price, and Srebro, see n. 14; Gajane and Pechenizkiy, see n. 9; Haas, see n. 12.

<sup>&</sup>lt;sup>16</sup>Richard Berk et al. "Fairness in criminal justice risk assessments: The state of the art". In: Sociol Methods Res 50.1 (2021), pp. 3–44.

# #3 Sufficiency

Random variables (R, A, Y) satisfy sufficiency if  $Y \perp \!\!\! \perp A \mid R$ 



#3 statistical non-discrimination criterion

Predictive parity (PP)<sup>17</sup>

$$\mathbb{P}(y=1 \mid a_1=0, f(\check{x}, a_1)=1) = \mathbb{P}(y=1 \mid a_1=1, f(\check{x}, a_1)=1).$$
 (6)

**Calibration**<sup>18</sup> (concept). For two binary predictors  $h_1, h_0 : \mathbb{R}^{n_d+1} \mapsto [0,1], h_1$  classifies samples with  $a_1 = 1$  and  $h_0$  does samples with  $a_1 = 0$ . Any  $h_t$  ( $t \in \{1,0\}$ ) is perfectly calibrated if

$$\forall p \in [0,1], \mathbb{P}_{(\check{\mathbf{x}}, a_i = t, y)}(y = 1 \mid h_t(\check{\mathbf{x}}, a_1) = p) = p.$$

It intuitively prevents the probability scores from carrying group-specific information.

<sup>&</sup>lt;sup>17</sup>Chouldechova, see n. 10; Sahil Verma and Julia Rubin. "Fairness definitions explained". In: FairWare. 2018, pp. 1–7.

<sup>&</sup>lt;sup>18</sup>Geoff Pleiss et al. "On fairness and calibration". In: NIPS. vol. 30. 2017; Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. "Inherent trade-offs in the fair determination of risk scores". In: ITCS. 2017.

## **Individual fairness**

#### Lipschitz condition 19,20

A mapping/predictor  $h: \mathcal{X} \times \mathcal{A}_1 = \mathcal{X} \times \{0,1\} \mapsto [0,1]$  satisfies the  $\lambda$ -Lipschitz property if for any  $(\check{\mathbf{x}}, a_1), (\check{\mathbf{x}}', a_1')$ ,

$$d_{y}(\underbrace{h(\check{\mathbf{x}},a_{1})},h(\check{\mathbf{x}}',a_{1}')) \leqslant \lambda \cdot d_{x}(\underbrace{(\check{\mathbf{x}},a_{1})},(\check{\mathbf{x}}',a_{1}')), \tag{7}$$

where  $d_y$  and  $d_x$  are (task-specific) distance metrics. Note that  $\lambda$  is a positive constant.

It can also be written as the probability Lipschitzness, i.e.,  $\mathbb{P}\left(\frac{d_y(h(\check{x},a_1),h(\check{x}',a_1'))}{d_x((\check{x},a_1),(\check{x}',a_1'))} \geqslant \epsilon\right) \leqslant \delta$ ; or the  $(\epsilon - \delta)$  language formulation:  $d_x((\check{x},a_1),(\check{x}',a_1')) \leqslant \epsilon \Rightarrow d_y(h(\check{x},a_1),h(\check{x}',a_1')) \leqslant \delta$ , where  $\epsilon \geqslant 0$  and  $\delta \geqslant 0$ .

In essence, individual fairness follows the principle that "similar individuals should be evaluated or treated similarly." A careful choice of distance metrics is crucial in ensuring fairness.<sup>21</sup>

<sup>19</sup> Dwork et al., see n. 10.

<sup>&</sup>lt;sup>20</sup> Additionally, a predictor satisfies individual fairness (Gajane and Pechenizkiy, see n. 9) iff:  $h(\check{x}, a_1) \approx h(\check{x}', a_1') \mid d_x((\check{x}, a_1), (\check{x}', a_1')) \approx 0$ , where  $\mathcal{X}_a \triangleq \mathcal{X} \times \mathcal{A}$  and  $d_x : \mathcal{X}_a \times \mathcal{X}_a \mapsto \mathbb{R}$  is a distance metric for individuals.

<sup>&</sup>lt;sup>21</sup>Binh Thanh Luong, Salvatore Ruggieri, and Franco Turini. "k-NN as an implementation of situation testing for discrimination discovery and prevention". In: *SIGKDD*. 2011, pp. 502–510; Laura Boeschoten et al. "Achieving fair inference using error-prone outcomes". In: *Int J Interact Multimed Artif Intell* 6.5 (2021).

## **Individual fairness**

#### General entropy indices<sup>22</sup> and the Theil index<sup>23</sup>

For a constant  $\alpha \notin \{0,1\}$ , the generalised entropy indices for a problem with n instances are defined, to quantify algorithmic unfairness, as

$$GEI^{\alpha} = \frac{1}{n\alpha(\alpha - 1)} \sum_{i=1}^{n} \left( \left( \frac{b_i}{\mu} \right)^{\alpha} - 1 \right), \tag{8}$$

where benefits  $b_i = f(\check{x}_i, a_{1i}) - y_i + 1$  and  $\mu = \sum_i b_i / n$ .

The Theil index is a special case for  $\alpha = 1$ , that is,

Theil = 
$$\frac{1}{n} \sum_{i=1}^{n} \frac{b_i}{\mu} \log \left( \frac{b_i}{\mu} \right)$$
. (9)

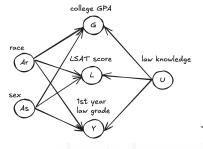
They are used additionally to group fairness measures to compare different algorithms and determine which one is considered the fairest from an individual perspective.

<sup>&</sup>lt;sup>22</sup>Till Speicher et al. "A unified approach to quantifying algorithmic unfairness: Measuring individual &group unfairness via inequality indices". In: *SIGKDD*. 2018, pp. 2239–2248.

<sup>&</sup>lt;sup>23</sup>Haas, see n. 12.

## Individual fairness

#### Example/Illustration: Law school success



As	<i>A</i> r	G	L	Y	U
male	black	9	de	4	(4)

- 1. Computing unobserved variables in causal model
- 2. Change A (that is, sensitive attribute(s))
- 3. Recompute observed variables in causal model

As 
$$Ar \leftarrow a'$$
  $G$   $L$   $Y$   $U$ 

male white  $\P$   $\P$ 

As  $Ar \leftarrow a'$   $G[Ar \leftarrow a']$   $L[Ar \leftarrow a']$   $Y[Ar \leftarrow a']$   $U[Ar \leftarrow a']$ 

Counterfactual fairness<sup>24</sup> (definition). A predictor  $\hat{Y}$  is counterfactually fair if given observations  $\mathcal{X} = \check{x}$  and A = a we have that,  $\mathbb{P}(\hat{Y}_{A \leftarrow a} = y \mid \mathcal{X} = \check{x}, A = a) = \mathbb{P}(\hat{Y}_{A \leftarrow a'} = y \mid \mathcal{X} = \check{x}, A = a)$ , for all y and  $a' \neq a$ . **Proxy discrimination** <sup>25</sup>. A predictor  $\hat{Y}$  exhibits no *individual proxy discrimination* based on a proxy P if for all p, p', we have  $\mathbb{P}(\hat{Y} \mid do(P = p)) = \mathbb{P}(\hat{Y} \mid do(P = p'))$ , where do(P = p) denotes an intervention on P.

<sup>&</sup>lt;sup>24</sup>Matt J Kusner et al. "Counterfactual fairness". In: NIPS. vol. 30. NIPS Proceedings. 2017, pp. 4069–4079.

<sup>&</sup>lt;sup>25</sup>(Niki Kilbertus et al. "Avoiding discrimination through causal reasoning". In: NIPS. vol. 30. 2017) Visually, intervening on P amounts to removing all incoming arrows of P in the graph; algebraically, it consists of replacing the structural equation of P by  $P = \mathbf{p}$ , i.e. we put point mass on the value  $\mathbf{p}$ .

# **Brief summary**

**Table 1:** Summary of existing fairness measures.

N	l., . 1	Meaning		Applicable situation(s) in definition			Non-binary handling <sup>4</sup>	
Name of measure	Fairness type <sup>1</sup>	quant.2	fairer	#label (n <sub>c</sub> )	#sen-att $(n_a)^3$	#values per $A_i$	$n_{a_i} > 2$	$n_a > 1$
Demographic parity (aka. statistical parity)	*, group-	yes	lower value	binary	singular	bi-valued	yes	indirectly
Disparate impact /80% rule	*, group-	yes	larger value	binary	singular	bi-valued	yes	indirectly
Disparate treatment	*, group-	poss.	lower value	binary	singular	bi- (multi- allowed)	yes	indirectly
Conditional statistical parity	*, group-	poss.	lower value	binary	singular	multi-valued	_	indirectly
Bounded group loss	*, group-	poss.	lower value	binary	singular	multi-valued		indirectly
Strategic minimax fairness	*, group-	no	_	bi-/multi-	singular	multi-valued	_	indirectly
Equalised odds	*, group-	yes	lower value	binary	singular	bi-valued	yes	indirectly
Equality of opportunity	*, group-	yes	lower value	binary	singular	bi-valued	yes	indirectly
Predictive equality	*, group-	poss.	lower value	binary	singular	multi-valued	_	indirectly
$\gamma$ -subgroup fairness	*, group-	yes	lower value	binary	singular	bi-valued	yes	indirectly
Predictive parity	*, group-	yes	lower value	binary	singular	bi-valued	yes	indirectly
Lipschitz condition	*, individual-	no	_	binary	singular	bi-valued	yes	indirectly
General entropy indices (and the Theil index)	*, individual-	yes	lower value	binary	singular	multi-valued	_	indirectly
Counterfactual fairness	*, individual-	no	_	binary	allows plural	multi- allowed	yes	indirectly
Proxy discrimination	*, individual-	no	_	binary	singular	multi- allowed	yes	indirectly
Discriminative risk <sup>26</sup>	*,5	yes	lower value	bi-/multi-	allows plural	multi-valued	_	_
Harmonic fairness via manifold <sup>27</sup>	* 5	ves	lower value	bi-/multi-	allows plural	multi-valued	_	_
Multiaccuracy	*,group-	poss.	lower value	binary	singular	multi-valued	_	indirectly
Differentially fair	*,group-	poss.		binary	allows plural	bi- (multi- allowed)	yes	indirectly
Group benefit ratio and worst-case min-max ratio	*,group-	ves	larger value	binary	allows plural	bi- (multi- allowed)	yes	indirectly
Feature-apriori fairness	procedural	yes		binary	10-A-1		yes	yes
Feature-accuracy fairness	procedural	yes	- <	binary			yes	yes
Feature-disparity fairness	procedural	yes		binary			yes	yes
FAE-based procedural fairness	procedural	yes	lower value	binary	singular	bi-valued	ves	indirectly

<sup>&</sup>lt;sup>26</sup>Yijun Bian and Kun Zhang. "Increasing fairness via combination with learning guarantees". In: arXiv preprint arXiv:2301.10813 (2023). Under review.

<sup>&</sup>lt;sup>27</sup>Yijun Bian and Yujie Luo. "Does machine bring in extra bias in learning? Approximating fairness in models promptly". In: arXiv preprint arXiv:2405.09251 (2024). Under review; Yijun Bian, Yujie Luo, and Ping Xu. "Approximating discrimination within models when faced with several non-binary sensitive attributes". In: arXiv preprint arXiv:2408.06099 (2024). Under review.

#### Statistical parity (SP)<sup>28</sup>

A predictor h satisfies *statistical parity* under a distribution over (X, A, Y) if h(x) is independent of the protected attribute a. Since  $h(x) \in [0, 1]$ , this is equivalent to

$$\mathbb{P}(f(\check{\mathbf{x}}, a_1) \geqslant z \mid \mathbf{a_1} = \mathbf{j}) = \mathbb{P}(f(\check{\mathbf{x}}, a_1) \geqslant z)$$
(10)

for all  $j \in A_1 = \{1, 2, ..., n_{a_1}\}$  and  $z \in [0, 1]$ .

**Demographic parity** (DP)  $|\mathbb{P}(f(\check{x}, a_1) = 1 \mid a_1 \neq 1) - \mathbb{P}(f(\check{x}, a_1) = 1 \mid a_1 = 1)| \le \varepsilon$ 

#1 statistical non-discrimination criterior

DP's extension and alternative<sup>29</sup> form

$$\max_{j \in \mathcal{A}_1} | \mathbb{P}(f(\check{x}, a_1) = 1 \mid a_1 = j) - \mathbb{P}(f(\check{x}, a_1) = 1) |,$$
(11a)

$$\max_{j,k \in \mathcal{A}_1, j \neq k} | \mathbb{P}(f(\check{x}, a_1) = 1 \mid a_1 = j) - \mathbb{P}(f(\check{x}, a_1) = 1 \mid a_1 = k) |.$$
 (11b)

<sup>&</sup>lt;sup>28</sup>Corbett-Davies et al., see n. 13; Agarwal, Dudík, and Wu, see n. 13.

<sup>&</sup>lt;sup>29</sup> Jiang et al., see n. 9.

Demographic parity (DP) 
$$|\mathbb{P}(f(\check{\mathbf{x}},a_1)=1 \mid a_1 \neq 1) - \mathbb{P}(f(\check{\mathbf{x}},a_1)=1 \mid a_1 = 1)| \leqslant \varepsilon$$
 Equality of opportunity (EOpp) 
$$|\mathbb{P}(f(\check{\mathbf{x}},a_1)=1 \mid a_1 \neq 1,y=1) - \mathbb{P}(f(\check{\mathbf{x}},a_1)=1 \mid a_1 = 1,y=1)| \leqslant \varepsilon$$
 Predictive parity (PP) 
$$|\mathbb{P}(y=1 \mid a_1 \neq 1,f(\check{\mathbf{x}},a_1)=1) - \mathbb{P}(y=1 \mid a_1 = 1,f(\check{\mathbf{x}},a_1)=1)| \leqslant \varepsilon$$

Three statistical non-discrimination criteria

#### EOpp's extension and alternative form

$$\max_{j \in \mathcal{A}_1} | \mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid a_1 = j, y = 1) - \mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid y = 1) |, \tag{10a}$$

$$\max_{j,k \in \mathcal{A}_1, j \neq k} | \mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid a_1 = j, y = 1) - \mathbb{P}(f(\check{\mathbf{x}}, a_1) = 1 \mid a_1 = k, y = 1) |.$$
 (10b)

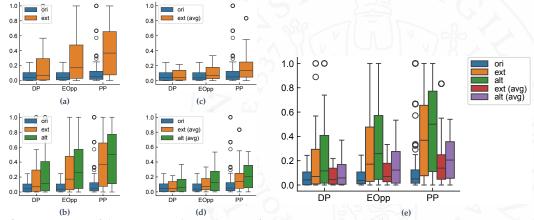
#### PP's extension and alternative form

$$\max_{j \in \mathcal{A}_1} | \mathbb{P}(y = 1 \mid a_1 = j, f(\check{x}, a_1) = 1) - \mathbb{P}(y = 1 \mid f(\check{x}, a_1) = 1) |,$$
(11a)

$$\max_{j,k \in A_1, j \neq k} | \mathbb{P}(y = 1 \mid a_1 = j, f(\check{x}, a_1) = 1) - \mathbb{P}(y = 1 \mid a_1 = k, f(\check{x}, a_1) = 1) |.$$
 (11b)

yjbian92@hotmail.com

D3A MLT workshop



**Figure 1:** Comparison of three commonly used group fairness measures and their extensions, on Income, Compas PPR, and Compas PPVR datasets.

(a–b) Comparison between binarisation and the two extension forms, analogously to Eq. (11a) and (11b); note that binarisation is equivalent to their original definitions like (1). (c–d) Comparison between binarisation and their corresponding average forms. (e) Comparison between binarisation and all four extension formulas.

dfprev

EOpp.

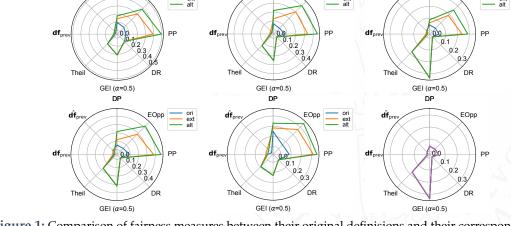
df<sub>prev</sub>

## Binarisation underestimates discrimination

EOpp

DP

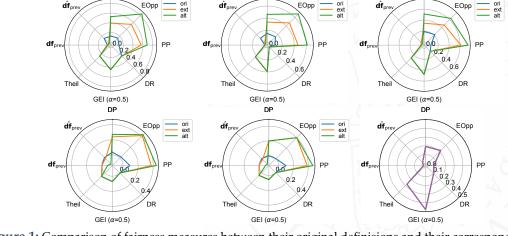
dfnres



DP

**Figure 1:** Comparison of fairness measures between their original definisions and their corresponding extension forms, on the Income dataset. (a–e) Using bagging, AdaBoost, LightGBM, AdaFair (trained using #1 sen-att), and AdaFair (trained using #2 sen-att), respectively.

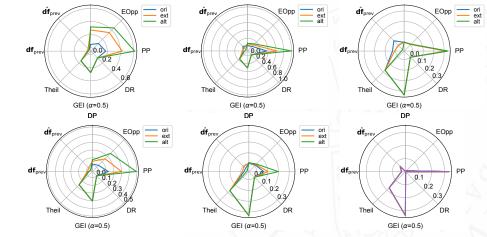
DP



DP

Figure 1: Comparison of fairness measures between their original definisions and their corresponding extension forms, on the Compas PPR dataset. (a–e) Using bagging, AdaBoost, LightGBM, AdaFair (trained using #1 sen-att), and AdaFair (trained using #2 sen-att), respectively.

DP



DP

**Figure 1:** Comparison of fairness measures between their original definisions and their corresponding extension forms, on the Compas PPVR dataset. (a–e) Using bagging, AdaBoost, LightGBM, AdaFair (trained using #1 sen-att), and AdaFair (trained using #2 sen-att), respectively.

# Traversal-based generalisation incurs computational burdens

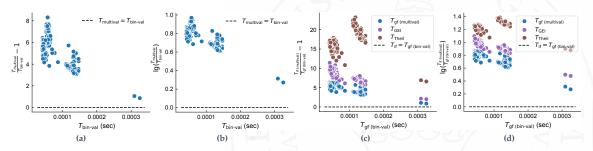
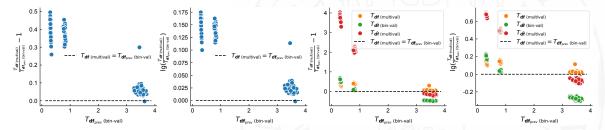


Figure 2: Time cost comparison of three commonly used group fairness measures and their extension forms, on Income, Compas PPR, and Compas PPVR datasets.

(a–b) Time cost comparison at different scales, and note that this is only for one 5- or 6-valued SA. Obviously, degenerating intersectional attributes ( $\mathcal{A} = \mathcal{A}_1 \times \mathcal{A}_2 = \mathbb{Z}^{n_{a_1}} \times \mathbb{Z}^{n_{a_2}}$  where  $n_{a_1}, n_{a_2} \geqslant 2$ ) into one "super" discrete SA through preprocessing is not an efficient way: It may be practical when both  $n_{a_1}$  and  $n_{a_2}$  are small enough, yet the computational cost increases exponentially as these values grow (e.g., if  $n_{a_1} = 2$  and  $n_{a_2}$  changes from 2 to 6,  $\mathcal{A}'$  transitions from  $\mathbb{Z}^4$  to  $\mathbb{Z}^{12}$ ). (c–d) Time cost comparisons, including individual fairness measures that are suitable for one multi-valued SA, indicate that individual fairness has an even heavier computational burden.

# Traversal-based generalisation incurs computational burdens



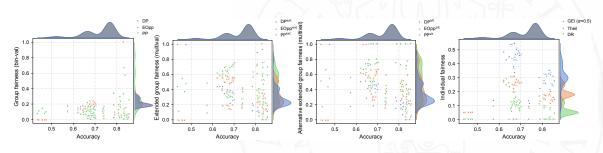
**Figure 2:** Time cost comparison of **HFM** for binary-value and multi-value cases, on Income, Compas PPR, and Compas PPVR datasets.

(a-b) Time cost comparisons of direct computation. (c-d) Comparisons including approximated results.<sup>29</sup>

<sup>&</sup>lt;sup>28</sup>Bian and Luo, see n. 27; Bian, Luo, and Xu, see n. 27.

<sup>&</sup>lt;sup>29</sup>Bian and Luo, see n. 27; Bian, Luo, and Xu, see n. 27.

# Accuracy and fairness are not strictly incompatible



**Figure 3:** Scatter plot between performance (accuracy) and fairness. Note that on the *y*-axis, the smaller the better; on the *x*-axis, the larger the better. (a) Using three commonly used group fairness measures; (b) Using their first extension forms; (c) Using their second extension forms; (d) Using individual fairness measures.

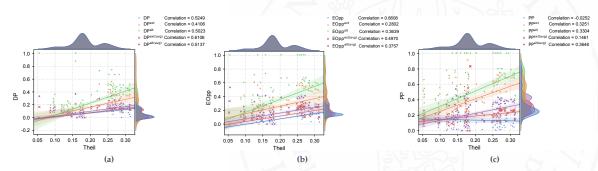


Figure 4: Relation between individual fairness and group fairness (DP, EOpp, and PP), on the Income, Compas PPR, and Compas PPVR datasets. Note that on both x- and y- axes, the smaller the better. (a–c) Using the Theil index as the individual fairness.

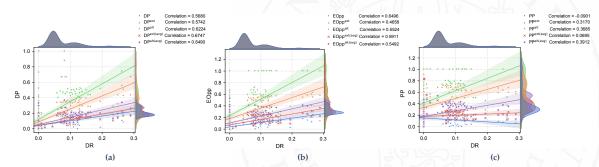
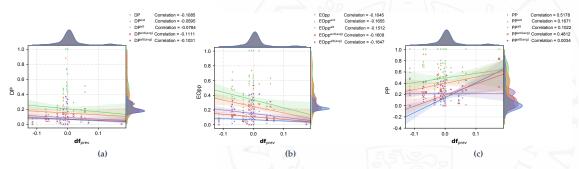


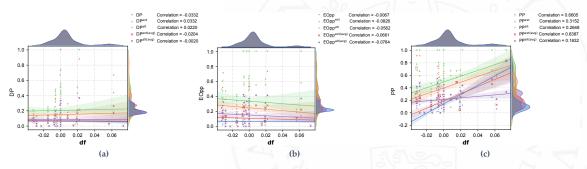
Figure 4: Relation between individual fairness and group fairness (DP, EOpp, and PP), on the Income, Compas PPR, and Compas PPVR datasets. Note that on both x- and y- axes, the smaller the better. (a–c) Using DR as the individual fairness.



**Figure 4:** Relation between individual fairness and group fairness (DP, EOpp, and PP), on the Income, Compas PPR, and Compas PPVR datasets. Note that on both x- and y- axes, the smaller the better. (a–c) Using the previous HFM $^{31}$  as the individual fairness.

<sup>&</sup>lt;sup>30</sup>Bian and Luo, see n. 27.

<sup>31</sup> Ibid.



**Figure 4:** Relation between individual fairness and group fairness (DP, EOpp, and PP), on the Income, Compas PPR, and Compas PPVR datasets. Note that on both x- and y- axes, the smaller the better. (a–c) Using the maximum HFM<sup>31</sup> as the individual fairness.

<sup>&</sup>lt;sup>30</sup>Bian, Luo, and Xu, see n. 27.

<sup>31</sup> Ibid.

Thanks! Questions?

## Fairness through unawareness<sup>32</sup>

A predictor is said to achieve *fairness through unawareness* (FTU) (or unconscious/unaware fairness) if all protected attributes  $\mathcal{A}$  are excluded from the decision-making process.

Despite its compelling simplicity, this approach has a clear shortcoming: the remaining attributes  $\mathcal{X}$  may contain discriminatory information analogous to  $\mathcal{A}$  that may not be obvious at first, acting as *proxy* attributes. As a result, discrimination cannot be guaranteed to be eliminated.

#### Discriminative risk (DR)<sup>33</sup>

$$DR(f) = \mathbb{E}[\mathbb{I}(f(\check{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \check{\mathbf{a}}))], \tag{12}$$

where  $\tilde{a}$  is a perturbed a, and  $n_a \geqslant 1$ ,  $|A_i| \geqslant 2$  ( $i \in [n_a]$ ).

Harmonic fairness via manifold (HFM)<sup>34</sup> (three versions, see below)

<sup>&</sup>lt;sup>32</sup>Dwork et al., see n. 10; Nina Grgić-Hlača et al. "The case for process fairness in learning: Feature selection for fair decision making". In: NIPS symposium on machine learning and the law. Vol. 1. 2. Barcelona, Spain. 2016, p. 11; Kusner et al., see n. 24; Gajane and Pechenizkiy, see n. 9.

<sup>33</sup> Bian and Zhang, see n. 26.

<sup>&</sup>lt;sup>34</sup>Bian and Luo, see n. 27; Bian, Luo, and Xu, see n. 27.

## Other distributive fairness

#### Harmonic fairness via manifold (HFM)<sup>32</sup>

Given a dataset D = (X, A, Y), it has three versions: (1) the *previous* HFM for one bi-valued SA, and (2) the *maximal (resp. average)* HFM for several multi-valued SAs.

For one bi-valued SA  $a_1 \in \mathcal{A}_1 = \{0,1\}$ , D is divided into  $D_1 = \{(\check{x}_a, y) \triangleq (\check{x}_a, a_1, y) \in D \mid a_1 = 1\}$  and  $\bar{D}_1 = D \setminus D_1$ , then given a specific distance metric  $d(\cdot, \cdot)$  (e.g., the standard Euclidean metric), the previous HFM is

$$\mathbf{df}_{\text{prev}}(f) = \frac{\mathbf{g}_f(D_1, \bar{D}_1)}{\mathbf{g}(D_1, \bar{D}_1)} - 1, \tag{12}$$

where

$$g_{\cdot}(D_1,\bar{D}_1;\boldsymbol{\mathcal{Y}}) = \max\{\max_{(\check{\mathbf{x}}_a,\boldsymbol{\mathcal{Y}})\in D_1} \min_{(\check{\mathbf{x}}'_a,\boldsymbol{\mathcal{Y}}')\in \bar{D}_1} d((\check{\mathbf{x}},\boldsymbol{\mathcal{Y}}),(\check{\mathbf{x}}',\boldsymbol{\mathcal{Y}}')) \max_{(\check{\mathbf{x}}'_a,\boldsymbol{\mathcal{Y}}')\in \bar{D}_1} \min_{(\check{\mathbf{x}}_a,\boldsymbol{\mathcal{Y}})\in D_1} d((\check{\mathbf{x}},\boldsymbol{\mathcal{Y}}),(\check{\mathbf{x}}',\boldsymbol{\mathcal{Y}}'))\},$$

and  $g_f(D_1, \bar{D}_1) = g_{\cdot}(D_1, \bar{D}_1; f(\check{\mathbf{x}}, a_1)), g(D_1, \bar{D}_1) = g_{\cdot}(D_1, \bar{D}_1; y)$  are two abbreviations for brevity.

<sup>&</sup>lt;sup>32</sup>Bian and Luo, see n. 27; Bian, Luo, and Xu, see n. 27.

## Other distributive fairness

#### Harmonic fairness via manifold (HFM) (cont.)

For one or more multi-valued SAs  $a \in A$  where  $n_a \geqslant 1$  and  $|A_i| \geqslant 2$   $(i \in [n_a])$ , the maximal (resp. average) HFM are

$$\mathbf{df} = \log \left( \frac{g_{f,a}(D)}{g_a(D)} \right) , \tag{12a}$$

$$\mathbf{df}^{\text{avg}}(f) = \log \left( \frac{\mathbf{g}_{f,a}^{\text{avg}}(D)}{\mathbf{g}_{a}^{\text{avg}}(D)} \right) , \tag{12b}$$

where

$$\mathbf{g}_{\cdot,a}(D;\ddot{y}) = \max_{1 \leq i \leq n_a} \mathbf{g}_{\cdot,a}(D,a_i;\ddot{y}), \tag{13a}$$

$$\mathbf{g}_{\cdot,a}^{\text{avg}}(D; \ddot{y}) = \frac{1}{n_a} \sum_{i=1}^{n_a} \mathbf{g}_{\cdot,a}^{\text{avg}}(D, a_i; \ddot{y}), \tag{13b}$$

$$\mathbf{g}_{\cdot,\mathbf{a}}(D,a_i;\ddot{y}) = \max_{i \in [n_{a_i}]} \{ \max_{(\breve{\mathbf{x}}_{\mathbf{a}}y) \in D_j} \min_{(\breve{\mathbf{x}}_{\mathbf{a}}'y') \in \bar{D}_j} \mathbf{d}((\breve{\mathbf{x}},\ddot{y}),(\breve{\mathbf{x}}',\ddot{y}')) \},$$

$$\mathbf{g}_{.a}^{\text{avg}}(D, a_i; \ddot{y}) = \frac{1}{n} \sum_{j \in [n_{a_i}]} \sum_{(\breve{\mathbf{x}}, y) \in D_j} \min_{(\breve{\mathbf{x}}', y') \in \bar{D}_j} \mathbf{d}((\breve{\mathbf{x}}, \ddot{y}), (\breve{\mathbf{x}}', \ddot{y}')).$$

Note that  $D_j = \{(\check{\mathbf{x}}_a, y) \in D | a_i = j\}$ ,  $\bar{D}_j = D \setminus D_j$ , and special case  $\mathbf{g}_{...a}(D, a_i; \ddot{y}) = \mathbf{g}_.(D_1, \bar{D}_1; \ddot{y})$  when  $A_i = \{0, 1\}$ .

## Procedural fairness

Here we use D to denote the set of all instances/members (or queried users of society), and S the set of all possible features that might be used in the decision-making process (in other words,  $|S| \le n_d + n_a$ ). Given a set of features S', let  $f_{S'}$  denote the classifier that uses those features S'.

**Feature-apriori fairness**. For a given feature  $s \in S$ , let  $D_s \subseteq D$  denote the set of all members that consider the feature s fair to use *without a priori knowledge* of how its usage affects outcomes. Then

$$PF_{apr}(f_{S'}) \triangleq \frac{|\bigcap_{s_i \in S'} D_{s_i}|}{|D|}. \tag{14}$$

Feature-accuracy fairness. (see below) Feature-disparity fairness. (see below)

These three measures<sup>32</sup> accommodate scenarios with multiple SAs, each potentially having multiple values. Despite this advantage, they rely heavily on features and on a set of members/users who perceive these features as fair, which may still introduce hidden discrimination or human prejudice. Moreover, their computation is complex and time-consuming, as user judgments may evolve with learning, limiting their practical applicability.

<sup>&</sup>lt;sup>32</sup>Grgić-Hlača et al., "The case for process fairness in learning: Feature selection for fair decision making", see n. 32; Nina Grgić-Hlača et al. "Beyond distributive fairness in algorithmic decision making: Feature selection for procedurally fair learning". In: AAAI, vol. 32. 1. 2018.

#### Procedural fairness

**Feature-accuracy fairness**. Let  $D_s^{\text{acc}} \subseteq D$  denote the set of all members that consider the feature s fair to be use if it increases the accuracy of the classifier. Note that typically  $D_s \subseteq D_s^{\text{acc}}$  is expected, though this need not always hold exactly (due to either noise in estimating member preferences, or some members attaching some sort of negative connotation to the notion of accuracy). Then

$$PF_{acc}(f_{S'}) \triangleq \frac{|\bigcap_{s_i \in S'} Cond(D_{s_i}, D_{s_i}^{acc})|}{|D|}, \tag{14}$$

where

$$\operatorname{Cond}(D_{s_i}, D_{s_i}^{\operatorname{acc}}) = \begin{cases} D_{s_i}, & \text{if } \operatorname{acc}(f_{S'}) \leqslant \operatorname{acc}(f_{S' \setminus \{s_i\}}); \\ D_{s_i} \cup D_{s_i}^{\operatorname{acc}} = D_{s_i}^{\operatorname{acc}}, & \text{otherwise}. \end{cases}$$

**Feature-disparity fairness.** (cf. Appendices) Let  $D_s^{\text{disp}} \subseteq D$  denote the set of all members that consider the feature s fair to use *even if it increases a measure of disparity* (i.e., disparate impact or disparate mistreatment) of the classifier. Typically  $D_s^{\text{disp}} \subseteq D_s$  is expected, though this need not always hold strictly due to estimation error or other reasons. Let  $\text{disp}(f_{S'})$  denote the disparity it induces, and then

$$PF_{disp}(f_{S'}) \triangleq \frac{|\bigcap_{s_i \in S'} Cond(D_{s_i}, D_{s_i}^{aisp})|}{|D|},$$
(15)

where

$$\operatorname{Cond}(D_{s_i}, D_{s_i}^{\operatorname{disp}}) = \begin{cases} D_{s_i}^{\operatorname{disp}}, & \text{if } \operatorname{disp}(f_{S'}) > \operatorname{disp}(f_{S'\setminus \{s_i\}}); \\ D_{s_i} \cup D_{s_i}^{\operatorname{disp}} = D_{s_i}, & \text{otherwise}. \end{cases}$$

## Procedural fairness

Additionally,<sup>32</sup> proposed an FAE-based (feature attribution explanation) metric to assess group procedural fairness, which depends on the specific FAE techniques employed.

**FAE-based group procedural fairness**. A given dataset D is divided into two subsets by the values of a single SA, that is,  $D_1 = \{(\mathbf{x}_i, a_{1i}, y_i) \in D \mid a_{1i} = 1\}$  and  $D_2 = \{(\mathbf{x}_i, a_{1i}, y_i) \in D \mid a_{1i} = 0\}$ . A local FAE function  $g(\cdot)$  takes a model  $f(\cdot)$  and an explained data point  $(\mathbf{x}_i, a_{1i})$  as inputs and returns explanations (i.e., feature importance scores)  $\mathbf{e}_i = g(f, \mathbf{x}_i, a_{1i}) \in \mathbb{R}^{n_d+1}$ , where its j-th component  $e_{ij}$  is the importance score of the feature  $x_{ij}$  for the model's prediction  $f(\mathbf{x}_i, a_{1i})$ . For a distance measure  $d_e(\cdot, \cdot)$  between two sets of FAE explanation results  $E_1$  and  $E_2$ , then

GPF<sub>FAE</sub> = 
$$d_e(E_1, E_2)$$
;  
 $E_1 = \{ \mathbf{e}_i \mid \mathbf{e}_i = g(f, \mathbf{x}_i, a_{1i}), (\mathbf{x}_i, a_{1i}) \in D_1' \},$   
 $E_2 = \{ \mathbf{e}_j \mid \mathbf{e}_j = g(f, \mathbf{x}_j, a_{1j}), (\mathbf{x}_j, a_{1j}) \in D_2' \},$ 
(14)

where  $D'_1$  and  $D'_2$  are sets of n data points from  $D_1$  and  $D_2$ , respectively, generated by [its Algorithm 1].

<sup>&</sup>lt;sup>32</sup>Ziming Wang, Changwu Huang, and Xin Yao. "Procedural fairness in machine learning". In: arXiv preprint arXiv:2404.01877 (2024).