

Developing Fair ML Models from Theoretical Aspects

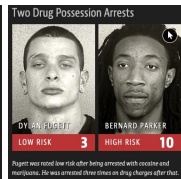
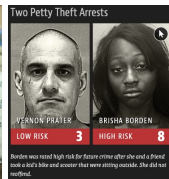
Yijun BIAN

Machine Learning Section
Department of Computer Science
University of Copenhagen

August 2024

My research

- ❶ **Yijun Bian^{*}**, Kun Zhang, Anqi Qiu, and Nanguang Chen. "Increasing fairness via combination with learning guarantees". In: arXiv preprint arXiv:2301.10813 (2023). In Revision.
- ❷ **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). In Revision.
- ❸ **Yijun Bian^{#*}**, Yujie Luo^{#*}, and Ping Xu. "Approximating Discrimination within Models When Faced With Several Non-Binary Sensitive Attributes". (2024). Under Review.



Examples ^{1,2,3}

¹AI detectors were more likely to flag writing by international students (i.e., non-native speakers) as AI-generated (Weixin Liang et al. "GPT detectors are biased against non-native English writers". In: *ICLR 2023 Workshop on Trustworthy and Reliable Large-Scale Machine Learning Models*. 2023)

²When people of color have complex medical needs, they are less likely to be referred to programmes that provide more individualised care (Linda Nordling. "A fairer way forward for AI in health care". In: *Nature* 573.7775 [2019], S103–S103)

³Black defendants were mislabelled as high risk more often than white defendants (Lorenzo Belenguer. "AI bias: exploring discriminatory algorithmic decision-making models and the application of possible machine-centric solutions adapted from the pharmaceutical industry". In: *AI and Ethics* 2.4 [2022], pp. 771–787)

Challenging

- Fairness estimation based on a finite sample
- Insufficient data
- Fairness measures/metrics ⁴
- The trade-off between fairness and accuracy

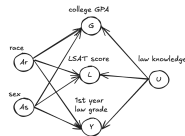


Group fairness



Individual fairness

Example/Illustration:
Law school success



A_s	A_r	G	L	Y	U
male	black	👉	👉	👉	(👉)

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

classifier
 $\text{test}(Y(G, A))$

A_s	A_r	A'	G	L	Y	U
male	white		👉	👉	👉	👉

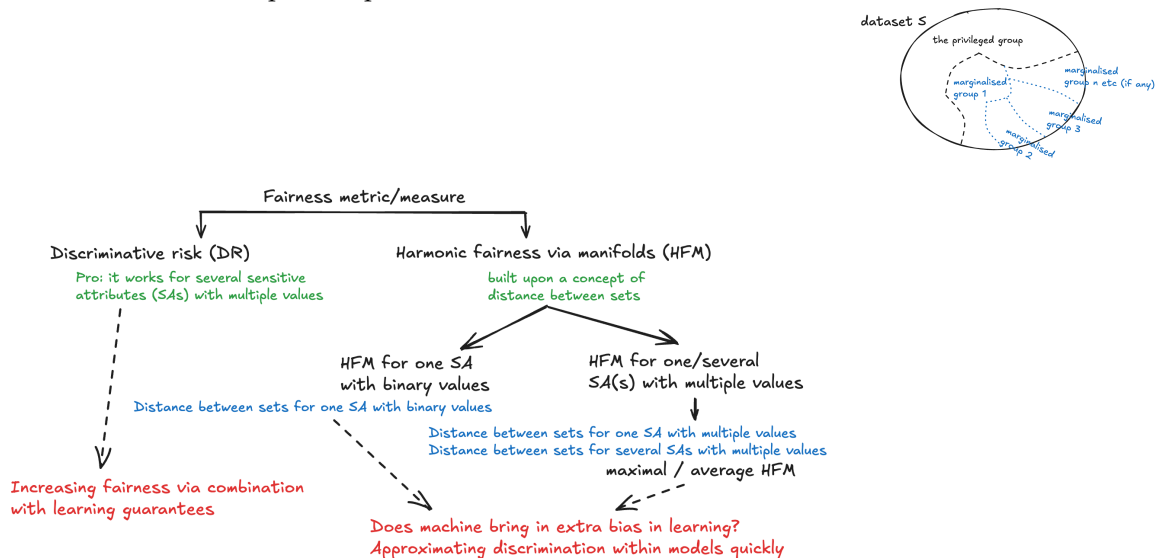
A_s	A_r	A'	$G_{(A_r \leftarrow A')}$	$L_{(A_r \leftarrow A')}$	$Y_{(A_r \leftarrow A')}$	$U_{(A_r \leftarrow A')}$
male	white		👉	👉	👉	👉

Our motivation

⁴Group 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.”

Appendix

Our current research up to the present



Discriminative risk (DR) —from an individual aspect

Following the principle of individual fairness, the fairness quality of one hypothesis⁵ $f(\cdot)$ could be evaluated by

$$\begin{aligned} \ell_{\text{fair}}(f, \mathbf{x}) &= \mathbb{I}(f(\check{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \tilde{\mathbf{a}})) \\ &= \mathbb{I}(f(\check{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \tilde{\mathbf{a}})), \end{aligned} \quad (1)$$

similarly to the 0/1 loss. Note that Eq. (1) is evaluated on only one instance with sensitive attributes \mathbf{x} .

⁵The hypothesis used in this equation could indicate an individual classifier or an ensemble classifier.

Discriminative risk (DR) —from a group aspect

To describe this characteristic of the hypothesis on multiple instances (aka. from a group level), then the **empirical discriminative risk on one dataset** S is expressed as

$$\hat{\mathcal{L}}_{\text{fair}}(f, S) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{fair}}(f, x_i), \quad (2)$$

discriminative risk of $f(\cdot)$ on one instance

and the **true discriminative risk**⁶ of the hypothesis **over a data distribution** is

$$\mathcal{L}_{\text{fair}}(f) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell_{\text{fair}}(f, x)], \quad (3)$$

discriminative risk of $f(\cdot)$ on one instance

respectively.

⁶The instances from S are independent identically distributed (i.i.d.) drawn from an input/feature-output/label space $\mathcal{X} \times \mathcal{Y}$ according to an unknown distribution \mathcal{D} .