

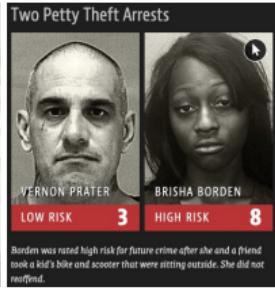
Study of FAIR ML Models from a Theoretical Perspective

Yijun BIAN

Machine Learning Section
Department of Computer Science
University of Copenhagen

5 December 2024

Examples of bias^{2,3,4}



Manuscripts under review¹

- ① Yijun Bian* and Kun Zhang. "Increasing fairness via combination with learning guarantees". In: *arXiv preprint arXiv:2301.10813v1* (2023).
- ② 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).
- ③ 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).

¹#Equal contribution; *corresponding author.

²AI 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)

³When 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)

⁴Black 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)

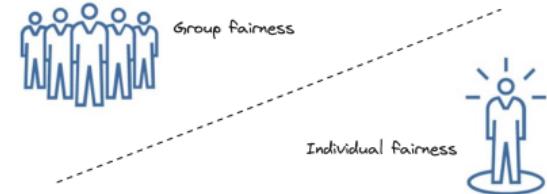
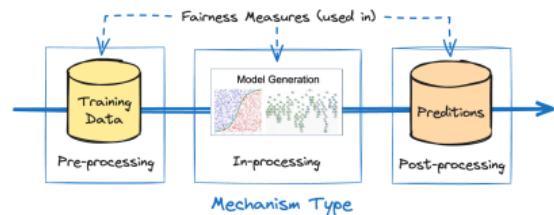
Challenging

- Fairness definitions and measures/metrics ^{a,b}
- Incompatibility among fairness measures
- Multi-attribute fairness protection
- The trade-off between fairness and accuracy
- Fairness estimation based on a finite sample
- Insufficient data

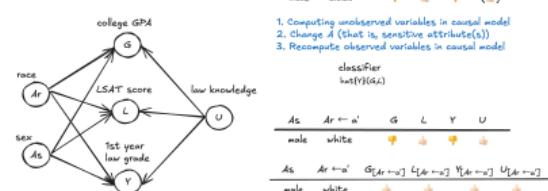
Our motivation

^a Pre- and post-processing mechanisms normally function by manipulating input or output, while inprocessing mechanisms introduce fairness constraints into training procedures or algorithmic objectives

^b 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."

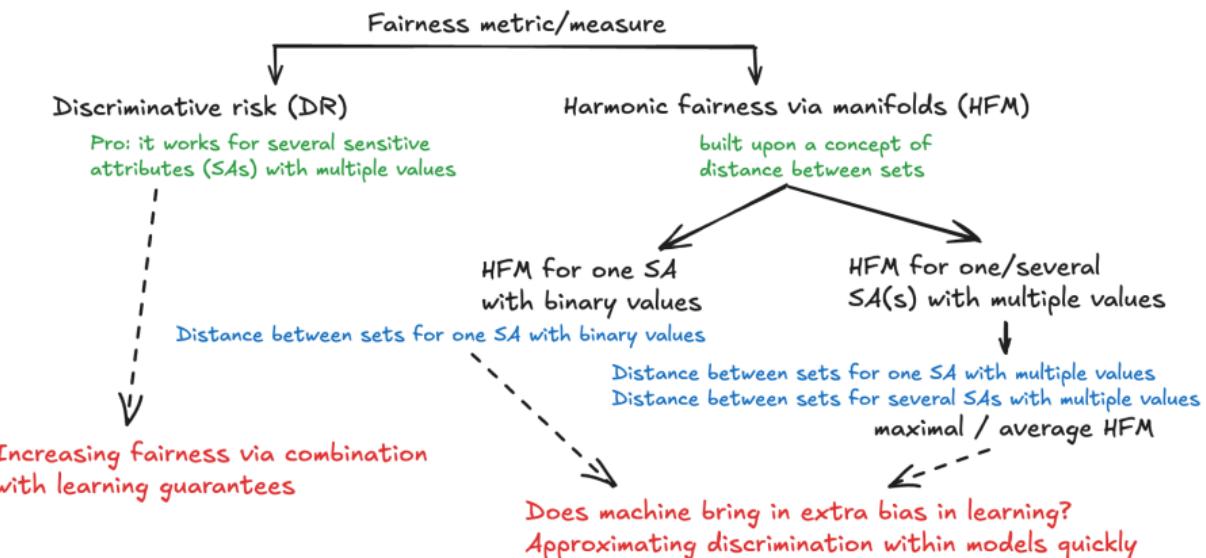
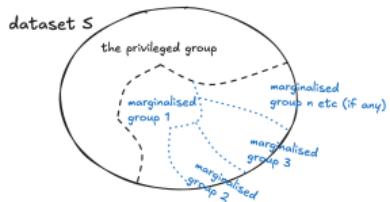


Example/Illustration:
Law school success



My research

My research in this direction, up to the present



Research question recap⁵

- 1. How to properly measure the discriminative level of a classifier from both individual and group fairness aspects?*
- 2. How to evaluate the discrimination level within classifiers when facing more than one sensitive attribute (potentially with multiple values)?*
- 3. How to make use of these proposed fairness measure/metric(s)?*

⁵Yijun Bian and Kun Zhang. "Increasing Fairness via Combination with Learning Guarantees". In: *arXiv preprint arXiv:2301.10813v1* (2023).

Discriminative risk (DR) —from an individual aspect

Following the principle of individual fairness, with an instance denoted by $x = (\check{x}, \check{a})$, the fairness quality of one hypothesis⁶ $f(\cdot)$ could be evaluated by

$$\ell_{\text{bias}}(f, x) = \mathbb{I} (f(\check{x}, \check{a}) \neq f(\check{x}, \tilde{a})) \quad (1)$$

the indicator function
 ↓
 model prediction on the raw instance
 ↓
 $\ell_{\text{bias}}(f, x) = \mathbb{I} (f(\check{x}, \check{a}) \neq f(\check{x}, \tilde{a}))$
 ↓
 non-sensitive attributes sensitive attribute(s)
 ↑
 ↓
 model prediction when only sensitive attribute(s) are changed
 ↓
 $f(\check{x}, \tilde{a})$
 ↑
 sensitive attribute(s) that are slightly disturbed

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

⁶The 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

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



and the **true discriminative risk**⁷ of the hypothesis over a data distribution is

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



respectively.

Note that the empirical DR on S is an unbiased estimation of the true DR.

⁷The 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} .

* A property of DR

$$\ell_{\text{bias}}(f, \mathbf{x}) = \mathbb{I}(f(\check{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \tilde{\mathbf{a}}))$$

$$\hat{\mathcal{L}}_{\text{bias}}(f, S) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{bias}}(f, \mathbf{x}_i)$$

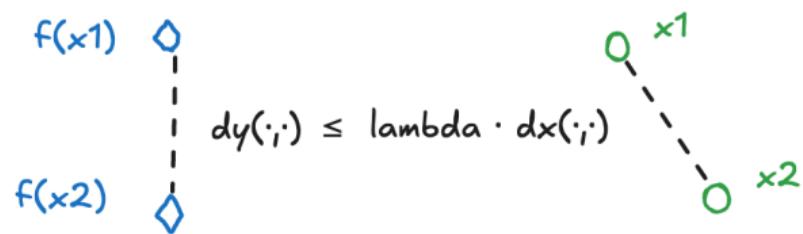
$$\mathcal{L}_{\text{bias}}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell_{\text{bias}}(f, \mathbf{x})]$$

- ① For one random variable \mathbf{X} representing instances, $\ell_{\text{bias}}(f, \mathbf{x})$ could be viewed as a new random variable obtained by using a few fixed operations on \mathbf{X} , recorded as \mathbf{Y} .
- ② For n random variables (i.e., X_1, X_2, \dots, X_n representing instances) that are independent and identically distributed (iid.), by operating them in the same way, we can get random variables Y_1, Y_2, \dots, Y_n that are iid. as well.
- ③ Then we can rewrite $\hat{\mathcal{L}}_{\text{bias}}(f, S)$ as $\frac{1}{n} \sum_{i=1}^n Y_i$ and $\mathcal{L}_{\text{bias}}(f)$ as $\mathbb{E}_{\mathbf{Y} \sim \mathcal{D}'}[\mathbf{Y}]$, where \mathcal{D}' denotes the space after operating $\mathbf{X} \sim \mathcal{D}$.
- ④ Therefore, it could be easily seen that the former is an unbiased estimation of the latter.

Our distinction

- Two distinctions from *individual fairness* measures

- 1 relies on the choice of similarity/distance metric
- 2 instance pairs in comparison coming from original data

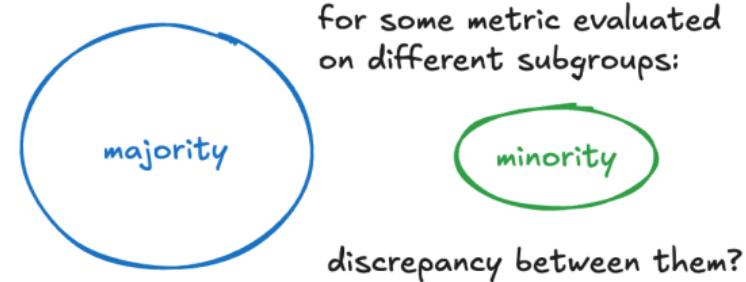


- Two distinctions from *group fairness* measures
- Five distinctions from *causal fairness*

Our distinction

- Two distinctions from *individual fairness* measures
- Two distinctions from *group fairness* measures

- ① works for only one sensitive attribute (usually bi-valued)
- ② computing separately for each subgroup, then difference

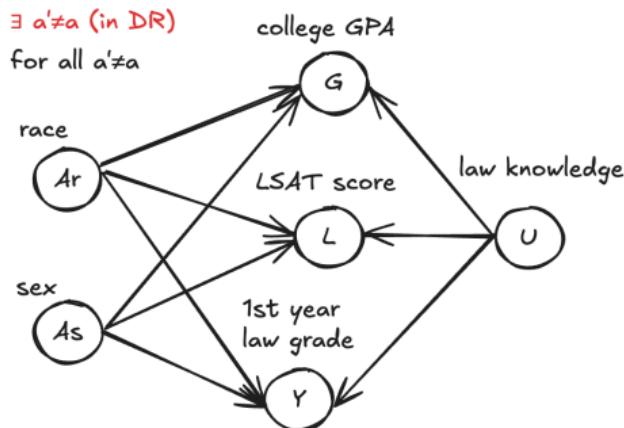


- Five distinctions from *causal fairness*

Our distinction

- Two distinctions from *individual fairness* measures
- Two distinctions from *group fairness* measures
- Five distinctions from *causal fairness*

- ① works for only one sensitive attribute (although possibly multi-valued)
- ② based on causal models/graphs, not a quantitative measure
- ③ non-sensitive attributes may vary with it in counterfactual fairness
- ④ conditions for achieving them are stronger
- ⑤ DR can be proved to be bounded vs. no such advantage



Our distinction

- Two distinctions from *individual fairness* measures
- Two distinctions from *group fairness* measures
- Five distinctions from *causal fairness*

$$\ell_{\text{bias}}(f, \mathbf{x}) = \mathbb{I}(f(\tilde{\mathbf{x}}, \mathbf{a}) \neq f(\tilde{\mathbf{x}}, \tilde{\mathbf{a}}))$$

$$\hat{\mathcal{L}}_{\text{bias}}(f, S) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{bias}}(f, \mathbf{x}_i)$$

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

$$\begin{aligned}\mathcal{L}'_{\text{bias}}(f) = & |\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D} | \mathbf{a}=1} [\ell_{\text{bias}}(f, \mathbf{x})] \\ & - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D} | \mathbf{a}=0} [\ell_{\text{bias}}(f, \mathbf{x})]| \end{aligned}$$

- Similarities that *DR* shares with the existing fairness measures

- follows the same principle as *individual fairness* measures
- is computed over a group of instances (like one dataset or a data distribution)
- indicates the discrimination level from a statistical/demographic perspective

Validating DR, a fairness quality measure

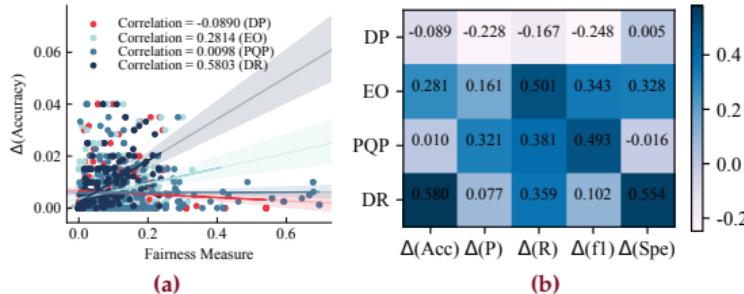


Figure 1: Comparison of the proposed discriminative risk (DR) with three group fairness measures, that is, DP, EO, and PQP. (a) Scatter diagrams with the degree of correlation, where the x - and y -axes are different fairness measures and the variation of accuracy between the raw and disturbed data. (b) Correlation among multiple criteria. Note that correlation here is calculated based on the results from all datasets.

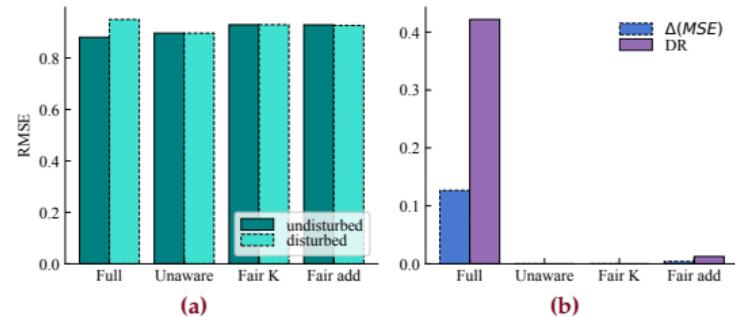


Figure 2: Example: law school success. (a) Test MSE of different models, where ‘undisturbed’ and ‘disturbed’ denote the results obtained from the original and disturbed data respectively. (b) The comparison between the change in MSE and DR , which suggests that $DR \approx 0$ when the corresponding model satisfies or nearly satisfies counterfactual fairness.

Interim summary⁸

RQ 1. *How to properly measure the discriminative level of a classifier from both individual and group fairness aspects?*

Discriminative risk (DR) is proposed, that is,

$$\ell_{\text{bias}}(f, \mathbf{x}) = \mathbb{I}(f(\check{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \check{\mathbf{a}}))$$

Widely applicable, with two reasons enlarging its applicable fields/scenarios:

- ① suitable for both binary and multi-class classification
- ② allows one or more SAs, and each SA allows binary or multiple values

Limitations

- ① The computational results of DR may be affected somehow by a randomness factor
- ② The degree of influence due to the number of values in SAs may vary, although its property remains

⁸Bian and Zhang, see n. 5.

Research question recap⁹

2'. How to evaluate the added discrimination introduced by ML models (on top of potential discrimination present in the raw data) properly?

- in the face of one sensitive attribute with binary values
- in the face of one sensitive attribute (with multiple values)
- in the face of several sensitive attributes (with multiple values)

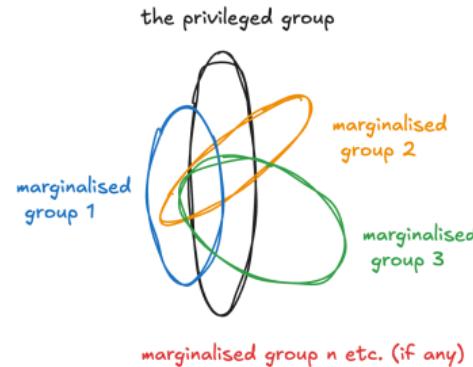
3'. How to efficiently evaluate the added discrimination introduced in the learning process?

- A quick approximation of distances between sets for the Euclidean spaces

⁹Yijun Bian and Yujie Luo. "Does Machine Bring in Extra Bias in Learning? Approximating Fairness in Models Promptly". In: *arXiv preprint arXiv:2405.09251* (2024); 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.

Harmonic Fairness via Manifolds (HFM)

If we view the *instances (with the same value of sensitive attributes)* as *data points on certain manifold(s)*, the manifold representing members from the marginalised/unprivileged group(s) is supposed to be as close as possible to that representing members from the privileged group.



To measure the fairness with respect to the sensitive attribute, we use a concept of 'distance of sets' introduced in mathematics, recorded as **D**, to *evaluate the discrepancy among groups divided by sensitive attributes*. Then a fairness measure¹⁰ is built upon the concept of distances between sets.

¹⁰indicating difference from both individual- and group- aspects

Distance between sets —for one bi-valued SA

Given a specific distance metric $\mathbf{d}(\cdot, \cdot)$ ¹¹ on the feature space, the **distance between two subsets** is defined by

$$\mathbf{D}(\mathcal{S}_1, \bar{\mathcal{S}}_1) \triangleq \max \left\{ \max_{(x,y) \in \mathcal{S}_1} \min_{(x',y') \in \bar{\mathcal{S}}_1} \mathbf{d}((\check{x}, \check{y}), (\check{x}', \check{y}')) , \right.$$

↑ to find the nearest data point in $\bar{\mathcal{S}}_1$

$$\left. \max_{(x',y') \in \bar{\mathcal{S}}_1} \min_{(x,y) \in \mathcal{S}_1} \mathbf{d}((\check{x}, \check{y}), (\check{x}', \check{y}')) \right\} ,$$

↑ works for both the true label y and the prediction \hat{y} of a trained classifier $f(\cdot)$

↑ the privileged group
↑ the marginalised/unprivileged group(s)

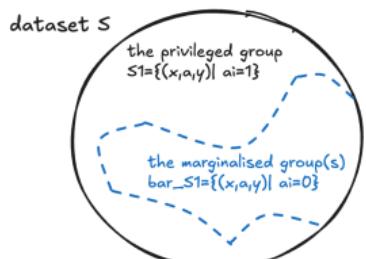
(6)

and is viewed as an approximation of the distance between the manifold of unprivileged groups and that of the privileged group.

Basic properties satisfied:

- ① For any two data sets $S_0, S_1 \in \mathcal{X} \times \mathcal{Y}$, $\mathbf{D}(S_0, S_1) = 0$ if and only if S_0 equals S_1 ; and
- ② For any sets S_0, S_1 , and S_2 , we have the triangle inequality

$$\mathbf{D}(S_0, S_2) \leq \mathbf{D}(S_0, S_1) + \mathbf{D}(S_1, S_2).$$



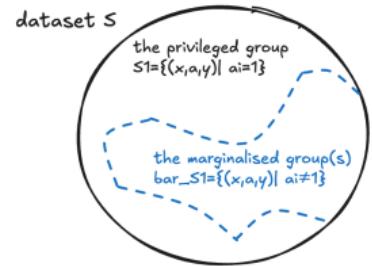
¹¹Here we use the standard Euclidean metric. In fact, any two metrics $\mathbf{d}_1, \mathbf{d}_2$ derived from norms on the Euclidean space \mathbb{R}^d are equivalent in the sense that there are positive constants c_1, c_2 such that $c_1 \mathbf{d}_1(x, y) \leq \mathbf{d}_2(x, y) \leq c_2 \mathbf{d}_1(x, y)$ for all $x, y \in \mathbb{R}^d$.

HFM¹² [Work 1 in this series]

For the distance between two subsets, and that of a trained classifier $f(\cdot)$, we have

$$\mathbf{D}(S_1, \bar{S}_1) \triangleq \max \left\{ \max_{(x,y) \in S_1} \min_{(x',y') \in \bar{S}_1} \mathbf{d}((\tilde{x},y), (\tilde{x}',y')), \max_{(x',y') \in \bar{S}_1} \min_{(x,y) \in S_1} \mathbf{d}((\tilde{x},y), (\tilde{x}',y')) \right\}, \quad (7a)$$

$$\mathbf{D}_f(S_1, \bar{S}_1) = \max \left\{ \max_{(x,y) \in S_1} \min_{(x',y') \in \bar{S}_1} \mathbf{d}((\tilde{x},\hat{y}), (\tilde{x}',\hat{y}')), \max_{(x',y') \in \bar{S}_1} \min_{(x,y) \in S_1} \mathbf{d}((\tilde{x},\hat{y}), (\tilde{x}',\hat{y}')) \right\}. \quad (7b)$$



We remark that $\mathbf{D}(S_1, \bar{S}_1)$ suggests the biases from the data and $\mathbf{D}_f(S_1, \bar{S}_1)$ suggests the biases from the algorithm. Then the following value could be used to indicate the fairness degree of this classifier, that is,

$$\mathbf{df}_{\text{prev}}(f) = \frac{\mathbf{D}_f(S_1, \bar{S}_1)}{\mathbf{D}(S_1, \bar{S}_1)} - 1. \quad (8)$$

¹²Bian and Luo, see n. 9.

Distance between sets —for one multi-valued SA

By extending Eq. (6), we introduce the following distance measures: (i) *maximal distance measure for one sensitive attribute*

$$\mathbf{D}_{\cdot,a}(S, a_i) \triangleq \max_{1 \leq j \leq n_{a_i}} \left\{ \max_{(\bar{x}, \bar{y}) \in S_j} \underbrace{\min_{(\bar{x}', \bar{y}') \in \bar{S}_j} \mathbf{d}((\bar{x}, \bar{y}), (\bar{x}', \bar{y}'))}_{\text{to find the nearest data point in } \bar{S}_j} \right\}, \quad (9)$$

and (ii) *average distance measure for one sensitive attribute*

$$\mathbf{D}_{\cdot,a}^{\text{avg}}(S, a_i) \triangleq \frac{1}{n} \sum_{j=1}^{n_{a_i}} \sum_{(\bar{x}, \bar{y}) \in S_j} \min_{(\bar{x}', \bar{y}') \in \bar{S}_j} \mathbf{d}((\bar{x}, \bar{y}), (\bar{x}', \bar{y}')) , \quad (10)$$

where $\bar{S}_j = S \setminus S_j$, and $n_{a_i} = |\mathcal{A}_i| \geq 2$ is the number of optional values for the sensitive attribute $a_i \in \mathcal{A}_i$. Notice that $\mathbf{D}_{\cdot,a}(S, a_i) = \mathbf{D}_{\cdot}(S_1, \bar{S}_1)$ when $\mathcal{A}_i = \{0, 1\}$.

Distance between sets —for several multi-valued SAs

For the general case, we introduce the generalised distance measures: (i) *maximal distance measure for sensitive attributes*

$$\mathbf{D}_{\cdot,a}(S) \triangleq \max_{1 \leq i \leq n_a} \mathbf{D}_{\cdot,a}(S, a_i) , \quad (11)$$

and (ii) *average distance measure for sensitive attributes*

$$\mathbf{D}_{\cdot,a}^{\text{avg}}(S) \triangleq \frac{1}{n_a} \sum_{i=1}^{n_a} \mathbf{D}_{\cdot,a}^{\text{avg}}(S, a_i) . \quad (12)$$

HFM¹³ [Work 2 in this series]

We remark that $\mathbf{D}_a(S)$, $\mathbf{D}_a^{\text{avg}}(S)$ reflect the biases from the data and that $\mathbf{D}_{f,a}(S)$, $\mathbf{D}_{f,a}^{\text{avg}}(S)$ reflect the extra biases from the learning algorithm. Then the following values could be used to reflect the fairness degree of this classifier, that is,

$$\mathbf{df}(f) = \log \left(\frac{\mathbf{D}_{f,a}(S)}{\mathbf{D}_a(S)} \right), \quad (13a)$$

$$\mathbf{df}^{\text{avg}}(f) = \log \left(\frac{\mathbf{D}_{f,a}^{\text{avg}}(S)}{\mathbf{D}_a^{\text{avg}}(S)} \right). \quad (13b)$$

We name the fairness degrees defined as above of one classifier by Eq. (13) as '*maximum harmonic fairness measure via manifolds (HFM)*' and '*average HFM*', respectively.

¹³Bian, Luo, and Xu, see n. 9.

Comparison between HFM^{14} and baseline fairness measures

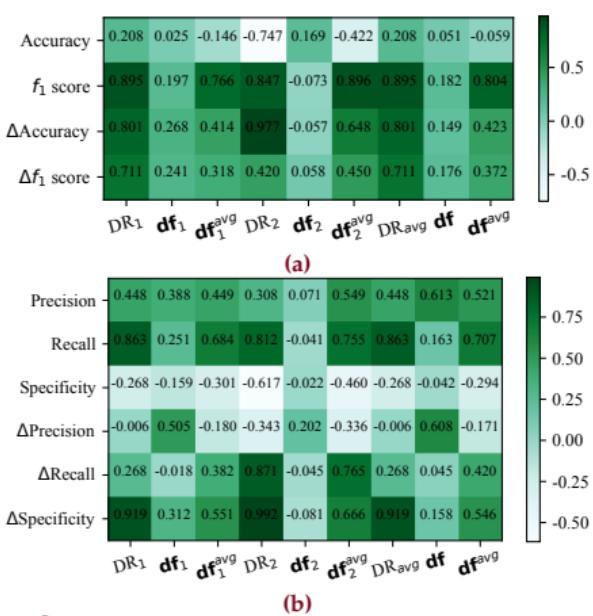


Figure 3: Correlation heatmap between normal evaluation metric and fairness measure, for all sensitive attributes within the dataset. Here we use $DR_{avg} = \frac{1}{n_a} \sum_{i=1}^{n_a} DR_i$ to reflect the bias level on the whole dataset.

¹⁴Bian, Luo, and Xu, see n. 9.

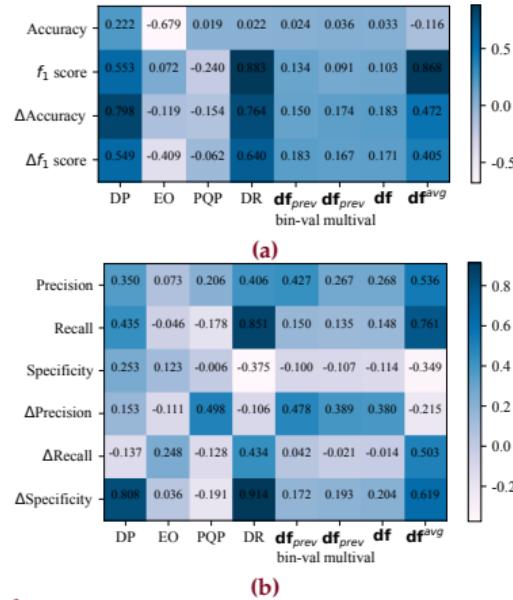


Figure 4: Correlation heatmap between normal evaluation metric and fairness, for one single sensitive attribute. Note that $df_{prev} = D_f(S_1, \bar{S}_1)/D(S_1, \bar{S}_1) - 1$ represents our previous work, and $df = \log(D_{f,a}(S, a_i)/D_a(S, a_i))$ and $df^{avg} = \log(D_{f,a}^{avg}(S, a_i)/D_a^{avg}(S, a_i))$ here represent HFM in this paper for each sensitive attribute.

*ExtendDist*¹⁵ [Work 2 in this series]



Algorithm 1 Approximation of extended distance between sets for several sensitive attributes with multiple values, aka. $\text{ExtendDist}(\{(\check{x}_i, a_i)\}_{i=1}^n, \{\check{y}_i\}_{i=1}^n; m_1, m_2)$,

Input: Dataset $S = \{(x_i, y_i)\}_{i=1}^n = \{(\check{x}_i, a_i, y_i)\}_{i=1}^n$ where $a_i = [a_{i,1}, a_{i,2}, \dots, a_{i,n_a}]^T$, prediction of S by the classifier $f(\cdot)$ that has been trained, that is, $\{\hat{y}_i\}_{i=1}^n$, and two hyper-parameters m_1 and m_2 as the designated numbers for repetition and comparison respectively

Output: Approximation of $D_{\cdot, a}(S)$ and $D_{\cdot, a}^{\text{avg}}(S)$

- 1: **for** j from 1 to n_a **do**
- 2: $d_{\max}^{(j)}, d_{\text{avg}}^{(j)} = \text{ApproxDist}(\{(\check{x}_i, a_{i,j})\}_{i=1}^n, \{\check{y}_i\}_{i=1}^n; m_1, m_2)$
- 3: **return** $\max_{1 \leq j \leq n_a} \{d_{\max}^{(j)} \mid j \in [n_a]\}$ and $\frac{1}{n_a} \sum_{j=1}^{n_a} d_{\text{avg}}^{(j)}$

¹⁵Bian, Luo, and Xu, see n. 9.

*ApproxDist*¹⁶ [Work 1&2 in this series]

Algorithm 2 (Simplified) Approximation of distance between sets, aka. $\text{ApproxDist}(\{(\check{x}_i, a_i)\}_{i=1}^n, \{\check{y}_i\}_{i=1}^n; m_1, m_2)$

Input: Dataset $S = \{(x_i, y_i)\}_{i=1}^n = \{(\check{x}_i, a_i, y_i)\}_{i=1}^n$, prediction of S by the classifier $f(\cdot)$ that has been trained, that is, $\{\hat{y}_i\}_{i=1}^n$, and two hyper-parameters m_1 and m_2 as the designated numbers for repetition and comparison respectively

Output: Approximation of distance $\mathbf{D}(S_1, \bar{S}_1)$ in Eq. (6)

- 1: **for** j from 1 to m_1 **do**
- 2: Take a random vector w from the space $\mathcal{W} = \{w = [w_0, w_1, \dots, w_{n_x}]^\top \mid \sum_{i=0}^{n_x} |w_i| = 1\} \subseteq [-1, 1]^{1+n_x}$
- 3: $d_{\max}^j = \text{AccelDist}(\{(\check{x}_i, a_i)\}_{i=1}^n, \{\check{y}_i\}_{i=1}^n, w; m_2)$
- 4: **return** $\min\{d_{\max}^j \mid j \in [m_1]\}$

Algorithm 2 Approximation of distance between sets (for one sensitive attribute with multiple values), aka. $\text{ApproxDist}(\{(\check{x}_i, a_i)\}_{i=1}^n, \{\check{y}_i\}_{i=1}^n; m_1, m_2)$

Input: Dataset $S = \{(x_i, y_i)\}_{i=1}^n = \{(\check{x}_i, a_i, y_i)\}_{i=1}^n$, prediction of S by the classifier $f(\cdot)$ that has been trained, that is, $\{\hat{y}_i\}_{i=1}^n$, and two hyper-parameters m_1 and m_2 as the designated numbers for repetition and comparison respectively

Output: Approximation of $\mathbf{D}_{\cdot, a}(S, a_i)$ and $\mathbf{D}_{\cdot, a}^{\text{avg}}(S, a_i)$

- 1: **for** j from 1 to m_1 **do**
 - 2: Take two orthogonal vectors w_0 and w_1 where each $w_k \in [-1, +1]^{1+n_x}$ ($k = \{0, 1\}$)
 - 3: **for** k from 0 to 1 **do**
 - 4: $t_{\max}^k, t_{\text{avg}}^k = \text{AccelDist}(\{(\check{x}_i, a_i)\}_{i=1}^n, \{\check{y}_i\}_{i=1}^n, w_k; m_2)$
 - 5: $d_{\max}^j = \min\{t_{\max}^k \mid k \in \{0, 1\}\} = \min\{t_{\max}^0, t_{\max}^1\}$
 - 6: $d_{\text{avg}}^j = \min\{t_{\text{avg}}^k \mid k \in \{0, 1\}\} = \min\{t_{\text{avg}}^0, t_{\text{avg}}^1\}$
 - 7: **return** $\min\{d_{\max}^j \mid j \in [m_1]\}$ and $\frac{1}{n} \min\{d_{\text{avg}}^j \mid j \in [m_1]\}$
-

¹⁶Bian and Luo, see n. 9; Bian, Luo, and Xu, see n. 9.

Distance approximation for Euclidean spaces

We observe that *the distance between similar data points tends to be closer than others after projecting them onto a general one-dimensional linear subspace* (refer to¹⁷).

To estimate the distance between data points inside $\mathcal{X} \times \mathcal{Y}$,

$$g(\mathbf{x}, \mathbf{y}; \mathbf{w}) = g(\check{\mathbf{x}}, \mathbf{a}, \mathbf{y}; \mathbf{w}) = [\mathbf{y}, \mathbf{x}_1, \dots, \mathbf{x}_{n_x}]^T \mathbf{w}, \quad (14)$$

where

- a random projection $g: \mathcal{X} \times \mathcal{Y} \mapsto \mathbb{R}$
- a non-zero random vector $\mathbf{w} = [w_0, w_1, \dots, w_{n_x}]^T$

That is to say, after sorting all the projected data points on \mathbb{R} , it is likely that *for one instance (\mathbf{x}, \mathbf{y}) in S_j , the desired instance $\operatorname{argmin}_{(\mathbf{x}', \mathbf{y}') \in \bar{S}_j} \mathbf{d}((\check{\mathbf{x}}, \mathbf{y}), (\mathbf{x}', \mathbf{y}'))$ would be somewhere near it after the projection, and vice versa*. Thus, searching for it could be **accelerated** by checking several adjacent instances rather than traversing the whole dataset.

¹⁷Bian and Luo, see n. 9, Lemma 1.

AcceleDist¹⁸ [Work 1&2 in this series]

Algorithm 3 Acceleration sub-procedure in approximation, aka. *AcceleDist* ($\{(\check{x}_i, a_i)\}_{i=1}^n, \{\check{y}_i\}_{i=1}^n, w; m_2$)

Input: Data points $\{(\check{x}_i, a_i)\}_{i=1}^n$, its corresponding value $\{\check{y}_i\}_{i=1}^n$, where \check{y}_i could be its true label y_i or prediction \hat{y}_i by the classifier $f(\cdot)$, a random vector w for projection, and a hyper-parameter m_2 as the designated number for comparison

Output: Approximation of distance $\mathbf{D} \cdot (S_0, S_1)$ in Eq. (6)

Output: Approximation of $\mathbf{D}_{\cdot, a}(S, a_i)$ and $n\mathbf{D}_{\cdot, a}^{\text{avg}}(S, a_i)$

- 1: Project data points onto a one-dimensional space based on Eq. (14), in order to obtain $\{g(x_i, \check{y}_i; w)\}_{i=1}^n$
- 2: Sort original data points based on $\{g(x_i, \check{y}_i; w)\}_{i=1}^n$ as their corresponding values, in ascending order
- 3: **for** i from 1 to n **do**
- 4: Set the anchor data point (x_i, \check{y}_i) in this round
- 5: **// If** $a_i = j$ (marked for clarity), in order to approximate $\min_{(x', y') \in S_j} \mathbf{d}((\check{x}_i, \check{y}_i), (x', y'))$
- 6: Compute the distances $\mathbf{d}((\check{x}_i, \check{y}_i), \cdot)$ for at most m_2 nearby data points that meets $a \neq a_i$ and $g(\check{x}, \check{y}; w) \leq g(\check{x}_i, \check{y}_i; w)$
- 7: Find the minimum among them, recorded as d_{\min}^s
- 8: Compute the distances $\mathbf{d}((\check{x}_i, \check{y}_i), \cdot)$ for at most m_2 nearby data points that meets $a \neq a_i$ and $g(\check{x}, \check{y}; w) \geq g(\check{x}_i, \check{y}_i; w)$
- 9: Find the minimum among them, recorded as d_{\min}^r
- 10: $d_{\min}^{(i)} = \min\{d_{\min}^s, d_{\min}^r\}$
- 11: **return** $\max\{d_{\min}^{(i)} \mid i \in [n]\}$
- 12: **return** $\max\{d_{\min}^{(i)} \mid i \in [n]\}$ and $\sum_{i=1}^n d_{\min}^{(i)}$

¹⁸Bian and Luo, see n. 9; Bian, Luo, and Xu, see n. 9.

Validity of approximation¹⁹ for distances in Euclidean spaces

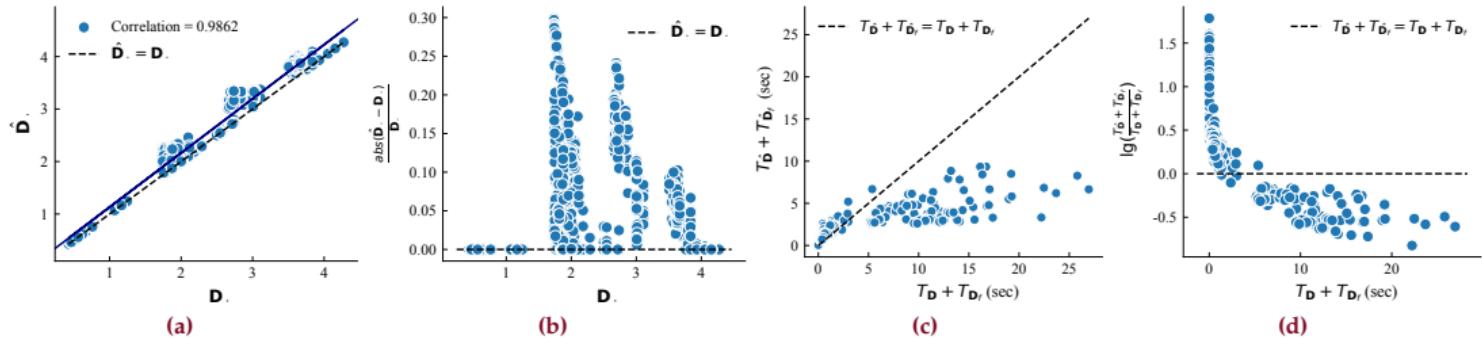


Figure 5: Comparison of approximation distances between sets with precise distances that are calculated directly by definition, evaluated on test data. (a) Scatter plot showing approximated values and precise values of distances between sets; (b) Relative difference comparison of *ApproxDist* with direct computation concerning distance values. (c-d) Comparison of time cost (second) between *ApproxDist* and direct computation based on Eq. (6).

¹⁹Bian and Luo, see n. 9.

Validity of *distance approximation*²⁰ in Euclidean spaces

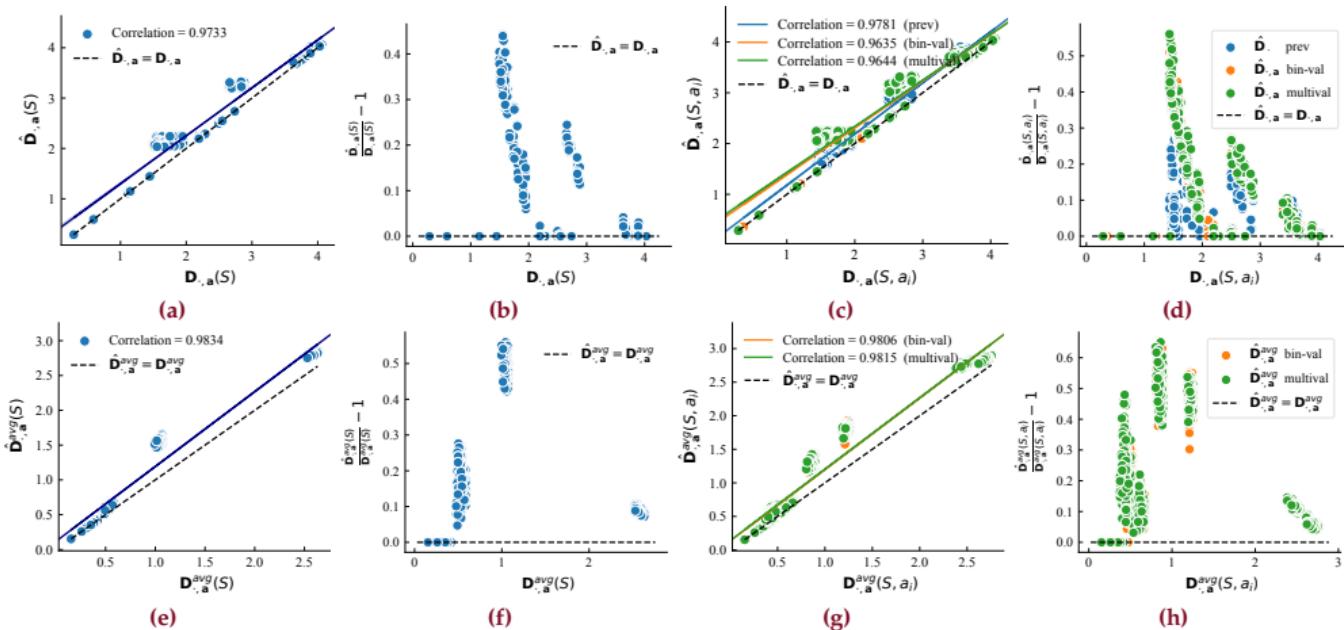


Figure 6: Comparison of approximation distances with precise distances that are calculated directly by definition, evaluated on test data. (a–b), (c–d), (e–f), and (g–h) Scatter plots for comparison between approximated and precise values of $D_{\cdot,a}(S)$, $D_{\cdot,a}(S, a_i)$, $D_{\cdot,a}^{\text{avg}}(S)$, and $D_{\cdot,a}^{\text{avg}}(S, a_i)$, respectively.

²⁰Bian, Luo, and Xu, see n. 9.

Validity of *distance approximation*²⁰ in Euclidean spaces

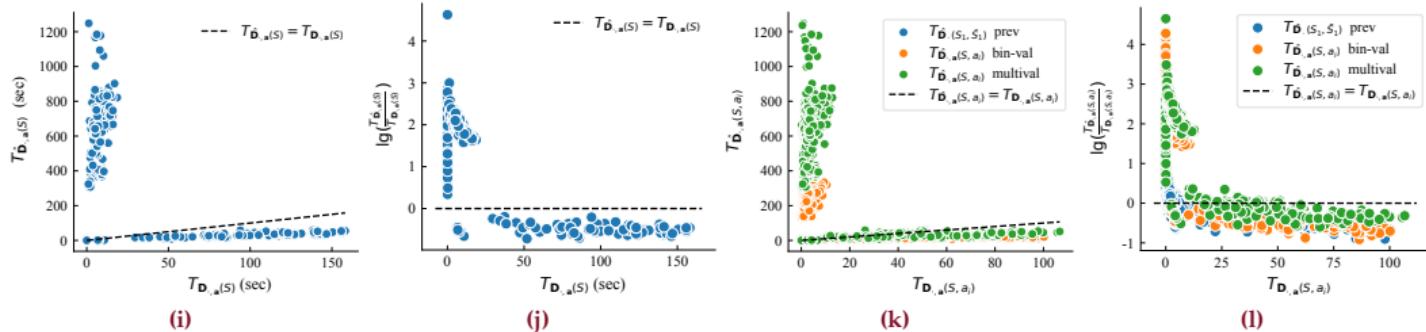


Figure 6: Comparison of approximation distances with precise distances that are calculated directly by definition, evaluated on test data. (i–j) Time cost comparison between *ExtendDist* and direct computation; (k–l) Time cost comparison between *ApproxDist* and direct computation. Note that ‘prev’ denotes approximation results obtained by the simplified Algorithm 2.

²⁰Bian, Luo, and Xu, see n. 9.

Interim summary²¹

RQ 2. How to efficiently measure the added discrimination introduced in learning by a classifier?

	Work 1	Work 2	
Distance between sets HFM	$D(S_1, \bar{S}_1)$ $\mathbf{df}_{\text{prev}}(f)$	$D_{\cdot, a}(S, a_i)(S, a_i)$ $D_{\cdot, a}(S)(S)$ $\mathbf{df}(f)$	$D_{\cdot, a}^{\text{avg}}(S, a_i)(S, a_i)$ $D_{\cdot, a}^{\text{avg}}(S)(S)$ $\mathbf{df}^{\text{avg}}(f)$
Approximation for one SA	<i>AcceleDist</i>	<i>ApproxDist</i>	<i>AcceleDist</i>
Approximation for several SA			<i>ApproxDist</i> <i>ExtendDist</i>

²¹Bian and Luo, see n. 9; Bian, Luo, and Xu, see n. 9.

Why-and-how-to-use recap²²

1. How to properly measure the discriminative level of a classifier from both individual and group fairness aspects?
- 3.' Can fairness be boosted with some learning guarantee? Will COMBINATION help mitigate discrimination in multiple biased individual classifiers?
- 3." How to utilise the proposed metric to obtain better ensemble classifiers?

²²Bian and Zhang, see n. 5.

Cancellation-of-bias effect²³ in ensemble combination

- Inspired by existing work for error rates and oracle bounds
- First- and second-order oracle bounds concerning fairness
- Similarly to the *cancellation-of-errors* effect in ensemble combination

The DR of an ensemble can be bounded by a constant times the DR of the individual classifiers

²³Bian and Zhang, see n. 5.

Ensemble combination

The *weighted voting prediction* by an ensemble of m trained individual classifiers parameterised by a weight vector $\rho = [w_1, w_2, \dots, w_m]^\top \in [0, 1]^m$, such that $\sum_{j=1}^m w_j = 1$, wherein w_j is the weight of individual classifier $f_j(\cdot)$, is given by

$$\text{wv}_\rho(x) = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \sum_{j=1}^m w_j \mathbb{I}(f_j(x) = y). \quad (15)$$

where a function $f \in \mathcal{F}: \mathcal{X} \mapsto \mathcal{F}$ denotes a hypothesis in a space of hypotheses \mathcal{F} . Note that ties are resolved arbitrarily.

Ensemble combination

The **weighted voting prediction** by an ensemble of m trained individual classifiers parameterised by a weight vector $\rho = [w_1, w_2, \dots, w_m]^\top \in [0, 1]^m$, such that $\sum_{j=1}^m w_j = 1$, wherein w_j is the weight of individual classifier $f_j(\cdot)$, is given by

$$\mathbf{wv}_\rho(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{j=1}^m w_j \mathbb{I}(f_j(x) = y). \quad (15)$$

where a function $f \in \mathcal{F}: \mathcal{X} \mapsto \mathcal{F}$ denotes a hypothesis in a space of hypotheses \mathcal{F} . Note that ties are resolved arbitrarily. Ensemble classifiers predict by taking a weighted combination of predictions by hypotheses from \mathcal{F} , and the **ρ -weighted majority vote** $\mathbf{wv}_\rho(\cdot)$ predicts

$$\mathbf{wv}_\rho(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \mathbb{E}_\rho [\mathbb{I}(f(x) = y)].$$

↑
 potential ρ corresponding to an ensemble over $[0, 1]^m$

Oracle bounds of fairness

If the weighted vote makes a discriminative decision, then at least a ρ -weighted half of the classifiers have made a discriminative decision and, therefore,

$$\ell_{\text{bias}}(\mathbf{wv}_\rho, \mathbf{x}) \leq \mathbb{I}(\mathbb{E}_\rho[\mathbb{I}(f(\check{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \tilde{\mathbf{a}}))] \geq 0.5). \quad (16)$$

discriminative risk of
an ensemble $\mathbf{wv}_\rho(\cdot)$

that is, $\ell_{\text{bias}}(f, \mathbf{x})$
discriminative risk of an individual classifier $f(\cdot)$ on one instance \mathbf{x}

Meaning of $\mathbf{wv}_\rho(\cdot)$

Ensemble classifiers (via *weighted voting*)

- take a weighted combination of predictions by hypotheses, and
- predict a label that receives the largest number of votes

In other words, the ρ -weighted majority vote $\mathbf{wv}_\rho(\cdot)$ predicts

$$\mathbf{wv}_\rho(\mathbf{x}) = \operatorname{argmax}_{y \in \mathcal{Y}} \mathbb{E}_\rho[\mathbb{I}[f(\mathbf{x}) = y]],$$

where ρ corresponds to a potential ensemble over a hypothesis space.



Oracle bounds of fairness

If the weighted vote makes a discriminative decision, then at least a ρ -weighted half of the classifiers have made a discriminative decision and, therefore,

$$\ell_{\text{bias}}(\mathbf{wv}_\rho, \mathbf{x}) \leq \mathbb{I}(\mathbb{E}_\rho[\mathbb{I}(f(\check{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \tilde{\mathbf{a}}))] \geq 0.5). \quad (16)$$

discriminative risk of
an ensemble $\mathbf{wv}_\rho(\cdot)$

that is, $\ell_{\text{bias}}(f, \mathbf{x})$
discriminative risk of an individual classifier $f(\cdot)$ on one instance \mathbf{x}

Theorem 1 (First-order oracle bound)

$$\mathcal{L}_{\text{bias}}(\mathbf{wv}_\rho) \leq 2 \mathbb{E}_\rho[\mathcal{L}_{\text{bias}}(f)]. \quad (17)$$

discriminative risk of an ensemble \mathbf{wv}_ρ

discriminative risk of an individual classifier f

the worst case is controlled to a constant multiple

Tandem discriminative risk

To investigate the bound deeper, we introduce here the tandem fairness quality of two hypotheses $f(\cdot)$ and $f'(\cdot)$ on one instance (x, y) , adopting the idea of the tandem loss,²⁴ by

$$\ell_{\text{bias}}(f, f', x) = \mathbb{I}(\text{hypothesis } f(\cdot) \text{ predicts differently for similar instances} \wedge \text{hypothesis } f'(\cdot) \text{ also predicts differently for them})$$

(18)

hypothesis $f(\cdot)$ predicts differently for similar instances
 ↓
 $\ell_{\text{bias}}(f, f', x) = \mathbb{I}(f(\check{x}, a) \neq f(\check{x}, \tilde{a}) \wedge f'(\check{x}, a) \neq f'(\check{x}, \tilde{a}))$
 ↑
 tandem discriminative risk

↓
 hypothesis $f'(\cdot)$ also predicts differently for them
 ↓
 discriminative risks present in both of them

The tandem fairness quality counts a discriminative decision on the instance (x, y) if and only if **both** $f(\cdot)$ and $f'(\cdot)$ give a discriminative prediction on it. Note that in the degeneration case

$$\ell_{\text{bias}}(f, f', x) = \ell_{\text{bias}}(f, x).$$

(19)

when $f'(\cdot)$ and $f(\cdot)$ are identical
 ↑
 discriminative risk of $f(\cdot)$

²⁴Andrés R Masegosa et al. "Second order PAC-Bayesian bounds for the weighted majority vote". In: *NeurIPS*. vol. 33. Curran Associates, Inc., 2020, pp. 5263–5273.

Oracle bounds of fairness (cont.)

Then the expected tandem fairness quality is defined by $\mathcal{L}_{\text{bias}}(f, f') = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell_{\text{bias}}(f, f', x)]$.

Theorem 3 (Second-order oracle bound)

$$\mathcal{L}_{\text{bias}}(\mathbf{wv}_\rho) \leq 4 \mathbb{E}_{\rho^2} [\mathcal{L}_{\text{bias}}(f, f')] . \quad (20)$$

discriminative risk of an ensemble \mathbf{wv}_ρ ↓ tandem discriminative risk of two individuals f and f'
 the worst case is controlled to a constant multiple

Lemma 2

In multi-class classification,

$$\mathbb{E}_{\rho^2} [\mathcal{L}_{\text{bias}}(f, f')] = \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{\rho} [\ell_{\text{bias}}(f, x)]^2] . \quad (21)$$

the expected tandem discriminative risk ↓ discriminative risk of $f(\cdot)$

Oracle bounds of fairness (cont.)

Theorem 4 (C-tandem oracle bound)

If $\mathbb{E}_\rho[\mathcal{L}_{bias}(f)] < 1/2$, then

$$\mathcal{L}_{bias}(\mathbf{w}\mathbf{v}_\rho)$$

discriminative risk of an ensemble $\mathbf{w}\mathbf{v}_\rho$

the worst case is controlled, alternative bound based on Chebyshev-Cantelli inequality

$$\leq \frac{\mathbb{E}_{\rho^2}[\mathcal{L}_{bias}(f, f')] - \mathbb{E}_\rho[\mathcal{L}_{bias}(f)]^2}{\mathbb{E}_{\rho^2}[\mathcal{L}_{bias}(f, f')] - \mathbb{E}_\rho[\mathcal{L}_{bias}(f)] + \frac{1}{4}}. \quad (22)$$

All oracle bounds are expectations that can only be estimated on finite samples instead of being calculated precisely. They could be transformed into empirical bounds via PAC analysis as well to ease the difficulty of giving a theoretical guarantee of the performance on any unseen data, which we discuss in this subsection. Based on the *Hoeffding's inequality*, we can deduct generalisation bounds presented in Theorems 5 and 6.

PAC bounds for the weighted vote

Theorem 5

For any $\delta \in (0, 1)$, with probability at least $(1 - \delta)$ over a random draw of S with a size of n , for a single hypothesis $f(\cdot)$,

$$\mathcal{L}_{bias}(f) \leq \hat{\mathcal{L}}_{bias}(f, S) + \sqrt{\frac{1}{2n} \ln \frac{1}{\delta}}. \quad (23)$$

the worst case is controlled with a specific bound

Theorem 6

For any $\delta \in (0, 1)$, with probability at least $(1 - \delta)$ over a random draw of S with a size of n , for all distributions ρ on \mathcal{F} ,

$$\mathcal{L}_{bias}(\mathbf{wv}_\rho) \leq \hat{\mathcal{L}}_{bias}(\mathbf{wv}_\rho, S) + \sqrt{\frac{1}{2n} \log \frac{|\mathcal{F}|}{\delta}}. \quad (24)$$

the worst case is controlled with a specific bound

Our distinction

Despite the similar names of “first- and second-order oracle bounds” from our inspiration,²⁵ the essences of our bounds are distinct from theirs. To be specific, their work investigates **the bounds for generalisation error** and is not relevant to fairness issues, while ours focus on the theoretical support for bias mitigation. In other words, their bounds are based on the 0/1 loss

$$\ell_{\text{err}}(f, \mathbf{x}) = \mathbb{I}(f(\mathbf{x}) \neq y), \quad (25)$$

The diagram illustrates the components of the 0/1 loss function $\ell_{\text{err}}(f, \mathbf{x}) = \mathbb{I}(f(\mathbf{x}) \neq y)$. A blue bracket labeled "the loss of the classifier $f(\cdot)$ " points to the first term $\mathbb{I}(f(\mathbf{x})$. A yellow bracket labeled "label of this instance, which means it makes mistakes on the instance" points to the second term $\neq y$. A purple bracket labeled "model prediction on the raw data" points to the expression $f(\mathbf{x})$.

while ours are built on $\ell_{\text{bias}}(f, \mathbf{x})$ in Eq. (1), that is,

$$\ell_{\text{bias}}(f, \mathbf{x}) = \mathbb{I}(f(\tilde{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \tilde{\mathbf{a}})).$$

The diagram illustrates the components of the bias loss function $\ell_{\text{bias}}(f, \mathbf{x}) = \mathbb{I}(f(\tilde{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \tilde{\mathbf{a}}))$. A blue bracket labeled "the discriminative risk of $f(\cdot)$ " points to the first term $\mathbb{I}(f(\tilde{\mathbf{x}}, \mathbf{a})$. A purple bracket labeled "model prediction on the raw data" points to the expression $f(\tilde{\mathbf{x}}, \mathbf{a})$. An orange bracket labeled "model prediction when only sensitive attribute(s) are changed" points to the expression $f(\check{\mathbf{x}}, \tilde{\mathbf{a}})$.

Besides, we have two more (PAC) generalisation bounds that they don't.

²⁵Masegosa et al., see n. ??.

Validating the oracle&PAC bounds

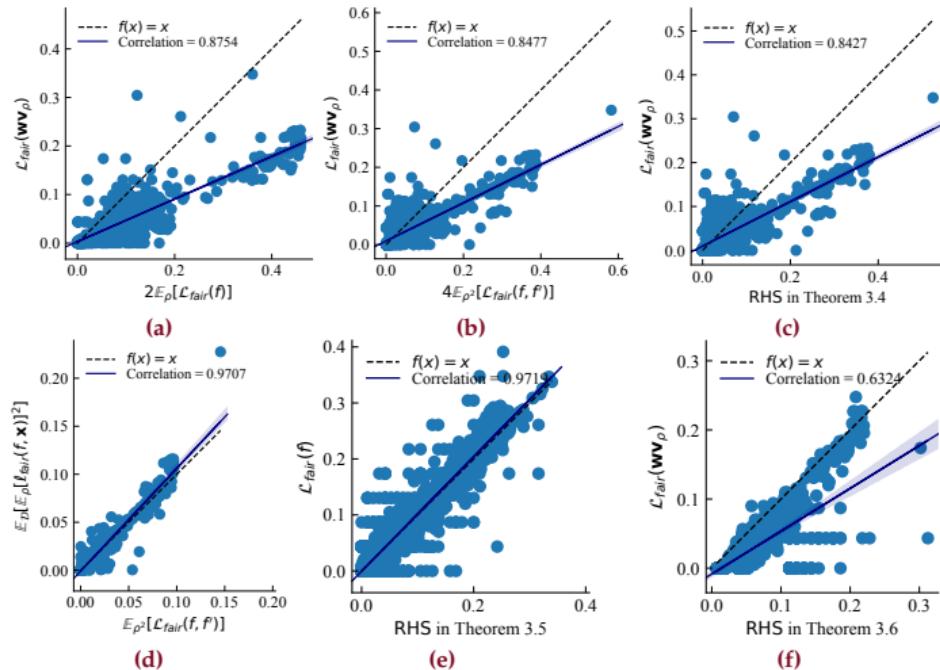


Figure 7: Correlation for oracle bounds and generalisation bounds. (a–c) Correlation between $\mathcal{L}_{\text{bias}}(\mathbf{w}\mathbf{v}_\rho)$ and oracle bounds, where $\mathcal{L}_{\text{bias}}(\mathbf{w}\mathbf{v}_\rho)$ is indicated on the vertical axis and the horizontal axes represent the right-hand sides of inequalities (17), (20), and (22), respectively. (d) The horizontal and vertical axes in (d) denote the right- and left-hand sides in (21), respectively. (e–f) Correlation between $\mathcal{L}_{\text{bias}}(\cdot)$ and generalisation bounds, where $\mathcal{L}_{\text{bias}}(\cdot)$ is indicated on the vertical axis and the right-hand sides of inequalities (23) and (24) are indicated on the horizontal axes, respectively. Note that correlation here is calculated based on the results from all datasets.

Interim summary²⁶

RQ 3. Why and how can we make use of these proposed fairness measure/metric(s)?

Ensemble combination: fairness can be boosted without being dependent on specific (hyper-)parameters, e.g.,

$$\mathcal{L}_{\text{bias}}(\mathbf{wv}_\rho) \leq 2 \mathbb{E}_\rho [\mathcal{L}_{\text{bias}}(f)]$$

cf. Theorem 1

$$\mathcal{L}_{\text{bias}}(\mathbf{wv}_\rho) \leq 4 \mathbb{E}_{\rho^2} [\mathcal{L}_{\text{bias}}(f, f')]$$

cf. Theorem 3

²⁶Bian and Zhang, see n. 5.

Why-and-how-to-use recap²⁷

- 1.** How to properly measure the discriminative level of a classifier from both individual and group fairness aspects?
- 3.'** Can fairness be boosted with some learning guarantee? Will COMBINATION help mitigate discrimination in multiple biased individual classifiers?
- 3."** How to utilise the proposed metric to obtain better ensemble classifiers?

²⁷ Bian and Zhang, see n. 5.

Concept of domination²⁸

We evaluate the accuracy quality of a hypothesis $f(\cdot)$ by the 0/1 loss $\ell_{\text{err}}(f, x)$, the empirical loss by $\hat{\mathcal{L}}_{\text{err}}(f, S) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{err}}(f, x_i)$, and the expected loss by $\mathcal{L}_{\text{err}}(f) = \mathbb{E}_{\mathcal{D}}[\ell_{\text{err}}(f, x)]$, respectively.

objective

$G = (\mathcal{L}_{\text{err}}, \mathcal{L}_{\text{bias}})$

solutions $\rho \leq \pi$

Definition 7 (Domination)

Let $\mathcal{L}_{\text{err}}(\cdot)$ and $\mathcal{L}_{\text{bias}}(\cdot)$ be two sub-objectives to be minimised, and let $G = (\mathcal{L}_{\text{err}}, \mathcal{L}_{\text{bias}})$ be the objective for a Pareto optimal solution. For two probability distributions ρ and π on \mathcal{F} that are independent of S : 1) ρ weakly dominates π if $\mathcal{L}_{\text{err}}(\mathbf{wv}_\rho) \leq \mathcal{L}_{\text{err}}(\mathbf{wv}_\pi)$ and $\mathcal{L}_{\text{bias}}(\mathbf{wv}_\rho) \leq \mathcal{L}_{\text{bias}}(\mathbf{wv}_\pi)$, denoted as \succeq_G ; 2) ρ dominates π if $\rho \succeq_G \pi$ and either $\mathcal{L}_{\text{err}}(\mathbf{wv}_\rho) < \mathcal{L}_{\text{err}}(\mathbf{wv}_\pi)$ or $\mathcal{L}_{\text{bias}}(\mathbf{wv}_\rho) < \mathcal{L}_{\text{bias}}(\mathbf{wv}_\pi)$, denoted as \succ_G .

²⁸A concept of domination (Chao Qian, Yang Yu, and Zhi-Hua Zhou. "Pareto ensemble pruning". In: AAAI. vol. 29. 1. 2015, pp. 2935–2941) is introduced to take fairness and accuracy into account simultaneously during pruning. Note that the domination relationship is used to achieve a Pareto optimal solution for this bi-objective minimisation case, where two sub-objectives are viewed as improving fairness and accuracy, respectively.

To construct fairer ensembles

Algorithm 4 Pareto Optimal Ensemble Pruning via Improving Accuracy and Fairness Concurrently (*POAF*)

Input: training set $S = \{(x_i, y_i)\}_{i=1}^n$, original ensemble $F = \{f_j(\cdot)\}_{j=1}^m$ via weighted vote, and threshold k as maximum size of the sub-ensemble after pruning

Output: Pruned sub-ensemble H ($H \subset F$ and $|H| \leq k$)

- 1: Randomly pick k individual members from F , indicated by \mathbf{r}
- 2: Initialise a candidate set for pruned sub-ensembles $\mathcal{P} = \{\mathbf{r}\}$
- 3: **for** $i = 1$ to k **do**
- 4: Randomly choose \mathbf{r} from \mathcal{P} with equal probability
- 5: Generate \mathbf{r}' by flipping each bit of \mathbf{r} with probability $1/m$
- 6: **if** $\exists \mathbf{z} \in \mathcal{P}$ such that $\mathbf{z} \succ_{\mathcal{G}} \mathbf{r}'$ **then**
- 7: **continue**
- 8: $\mathcal{P} = (\mathcal{P} \setminus \{\mathbf{z} \in \mathcal{P} \mid \mathbf{r}' \succeq_{\mathcal{G}} \mathbf{z}\}) \cup \{\mathbf{r}'\}$
- 9: Let $\mathcal{V} = \mathcal{N}_-(\mathbf{r}') \cup \mathcal{N}_+(\mathbf{r}')$
- 10: Sort \mathcal{V} by $\text{argmin}_{\mathbf{v} \in \mathcal{V}} \hat{\mathcal{L}}(\mathbf{v}, S)$ in ascending order
- 11: **for** $\mathbf{v} \in \mathcal{V}$ **do**
- 12: **if** $\exists \mathbf{z} \in \mathcal{P}$ such that $\mathbf{z} \succ_{\mathcal{G}} \mathbf{v}$ **then**
- 13: **continue**
- 14: $\mathcal{P} = (\mathcal{P} \setminus \{\mathbf{z} \in \mathcal{P} \mid \mathbf{v} \succeq_{\mathcal{G}} \mathbf{z}\}) \cup \{\mathbf{v}\}$
- 15: $H = \text{argmin}_{\mathbf{r} \in \mathcal{P}} \hat{\mathcal{L}}(\mathbf{r}, S)$



r
 r'



A solution $\text{wv}_\rho(\cdot)$ is *Pareto optimal* if there is no other solution in \mathcal{F} that dominates $\text{wv}_\rho(\cdot)$.

Comparison between *POAF* and fairness-aware ensemble-based methods²⁹

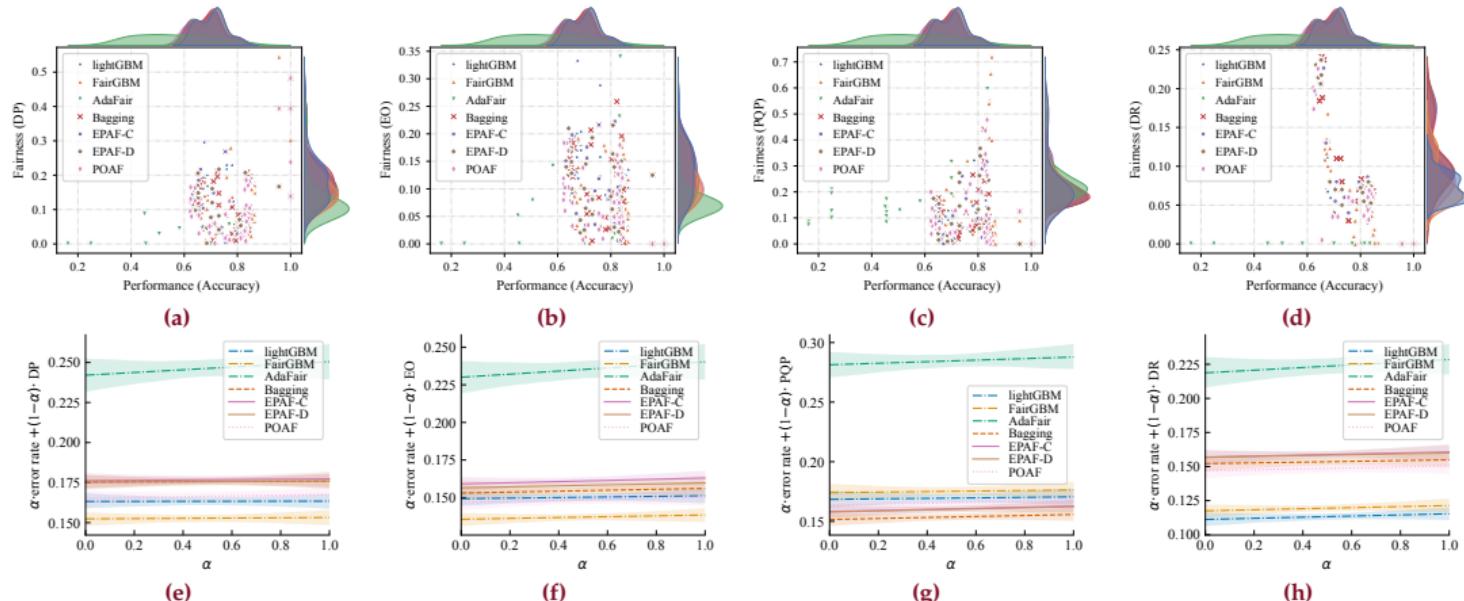


Figure 8: Comparison.. (a-d) Scatter plots showing fairness and accuracy of each algorithm, evaluated on the test data. (e-h) Plots of best test-set fairness-accuracy trade-offs per algorithm, where fairness is DP, EO, PQP, and DR, respectively. Lines show the mean value, and shades show 95% confidence intervals; The smaller the better.

²⁹ André F Cruz et al. "FairGBM: Gradient Boosting with Fairness Constraints". In: ICLR. 2023.

Comparison of *POAF* with *EPAF-C* and *EPAF-D* (& ensemble pruning methods)

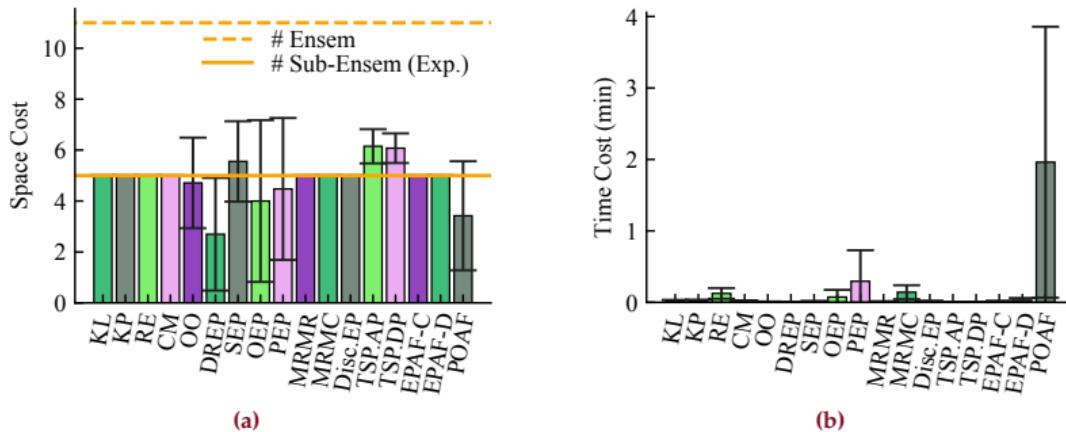


Figure 9: Comparison of the state-of-the-art pruning methods with *POAF*, where the dashed and solid lines indicate the size of the original ensemble and the expected size of the pruned sub-ensemble, respectively.
(a–b) Comparison over the space and time cost, respectively.

Alternative: Bi-objective

An *alternative* way to compare two hypotheses by considering both fairness and accuracy is to define an objective function based on an *adaptive weighted sum* method, i.e.,

$$\mathcal{L}(f, f') = \lambda \frac{\mathcal{L}_{\text{err}}(f) + \mathcal{L}_{\text{err}}(f')}{2} + (1 - \lambda) \mathcal{L}_{\text{bias}}(f, f'), \quad (26)$$

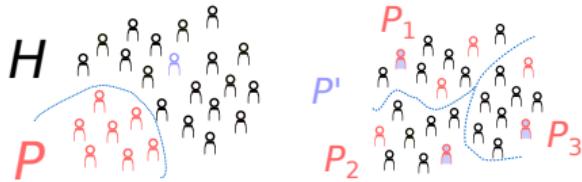
wherein $\lambda \in (0, 1)$ represents a regularisation factor introduced to balance fairness and accuracy and to indicate their relative importance as well. Note that this is a symmetrical function, that is, $\mathcal{L}(f, f') = \mathcal{L}(f', f)$, and that in the degeneration case $\mathcal{L}(f, f) = \lambda \mathcal{L}_{\text{err}}(f) + (1 - \lambda) \mathcal{L}_{\text{bias}}(f)$ when two hypotheses $f(\cdot)$ and $f'(\cdot)$ are identical. Then we can obtain the objective for the weighted vote as

$$\mathcal{L}(\mathbf{wv}_\rho) = \lambda \mathbb{E}_\rho [\mathcal{L}_{\text{err}}(f)] + (1 - \lambda) \mathbb{E}_{\rho^2} [\mathcal{L}_{\text{bias}}(f, f')], \quad (27)$$

which is expected to be minimised to improve fairness and accuracy concurrently.

Alternative: EPAF

Ensemble pruning via increasing accuracy and fairness concurrently (*EPAF*)



Algorithm 5 Centralised Version of Ensemble Pruning via Improving Accuracy and Fairness Concurrently (*EPAF-C*)

Input: training set $S = \{(x_i, y_i)\}_{i=1}^n$, original ensemble $F = \{f_j(\cdot)\}_{j=1}^m$, and threshold k as maximum size after pruning

Output: Pruned sub-ensemble H ($H \subset F$ and $|H| \leq k$)

1: $H \leftarrow$ an arbitrary individual member $f_i \in F$

2: **for** $i = 2$ **to** k **do**

3: $f^* \leftarrow \operatorname{argmin}_{f_i \in F \setminus H} \sum_{f_j \in H} \hat{\mathcal{L}}(f_i, f_j, S)$

4: Move f^* from F to H

Algorithm 6 Distributed Version of Ensemble Pruning via Improving Accuracy and Fairness Concurrently (*EPAF-D*)

Input: training set $S = \{(x_i, y_i)\}_{i=1}^n$, original ensemble $F = \{f_j(\cdot)\}_{j=1}^m$, threshold k as maximum size after pruning, and number of machines n_m

Output: Pruned sub-ensemble H ($H \subset F$ and $|H| \leq k$)

1: Partition F randomly into n_m groups as equally as possible, i.e., F_1, \dots, F_{n_m}

2: **for** $i = 1$ **to** n_m **do**

3: $H_i \leftarrow \text{EPAF-C}(F_i, k)$

4: $H' \leftarrow \text{EPAF-C}(\bigcup_{i=1}^{n_m} H_i, k)$

5: $H \leftarrow \operatorname{argmin}_{T \in \{H_1, \dots, H_{n_m}, H'\}} \hat{\mathcal{L}}(T, S)$

Thanks! Questions?

Commonly used group fairness measures

There are three commonly-used group fairness measures, that is, *demographic parity (DP)*,³⁰ *equality of opportunity (EO)*,³¹ and *predictive quality parity (PQP)*³².

These three commonly used group fairness measures of one classifier $f(\cdot)$ are evaluated as

$$\text{DP}(f) = |\mathbb{P}_{\mathcal{D}}[f(x)=1 | a=1] - \mathbb{P}_{\mathcal{D}}[f(x)=1 | a=0]|, \quad (28a)$$

$$\text{EO}(f) = |\mathbb{P}_{\mathcal{D}}[f(x)=1 | a=1, y=1] - \mathbb{P}_{\mathcal{D}}[f(x)=1 | a=0, y=1]|, \quad (28b)$$

$$\text{PQP}(f) = |\mathbb{P}_{\mathcal{D}}[y=1 | a=1, f(x)=1] - \mathbb{P}_{\mathcal{D}}[y=1 | a=0, f(x)=1]|, \quad (28c)$$

respectively, where $x = (\check{x}, a)$, y , and $f(x)$ are respectively features, the true label, and the prediction of this classifier for one instance. Note that $a = 1$ and 0 respectively mean that the instance x belongs to the privileged group and marginalised groups.

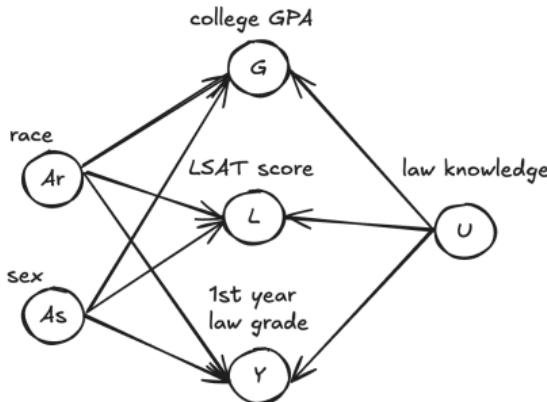
³⁰ Michael Feldman et al. "Certifying and removing disparate impact". In: *SIGKDD*. Sydney, NSW, Australia: Association for Computing Machinery, 2015, pp. 259–268; Pratik Gajane and Mykola Pechenizkiy. "On formalizing fairness in prediction with machine learning". In: *FAT/ML*. 2018.

³¹ Moritz Hardt, Eric Price, and Nathan Srebro. "Equality of opportunity in supervised learning". In: *NIPS*. vol. 29. Barcelona, Spain: Curran Associates Inc., 2016, pp. 3323–3331.

³² Alexandra Chouldechova. "Fair prediction with disparate impact: A study of bias in recidivism prediction instruments". In: *Big Data* 5.2 (2017), pp. 153–163; Sahil Verma and Julia Rubin. "Fairness definitions explained". In: *FairWare*. IEEE. 2018, pp. 1–7.

Fairness based on causality

Example/Illustration:
Law school success



As	Ar	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
 $\hat{y} = \{Y\}(G, L)$

As	Ar $\leftarrow a'$	G	L	Y	U
male	white	👎	👍	👎	👍

As	Ar $\leftarrow a'$	G $_{[Ar \leftarrow a']}$	L $_{[Ar \leftarrow a']}$	Y $_{[Ar \leftarrow a']}$	U $_{[Ar \leftarrow a']}$
male	white	👍	👍	👍	👍

Counterfactual fairness³³ definition: A predictor \hat{Y} is *counterfactually fair* if given observations $\mathcal{X}=x$ and $A=a$ we have that, $\mathbb{P}(\hat{Y}_{A \leftarrow a} = y | \mathcal{X}=x, A=a) = \mathbb{P}(\hat{Y}_{A \leftarrow a'} = y | \mathcal{X}=x, A=a)$, for all y and $a' \neq a$.

Proxy discrimination³⁴ definition: A predictor \hat{Y} exhibits no *individual proxy discrimination* if given observations $\mathcal{X}=x$ and $A=a$ we have that, $\mathbb{P}(\hat{Y}=y | \text{do}(A=a), \mathcal{X}=x) = \mathbb{P}(\hat{Y}=y | \text{do}(A=a'), \mathcal{X}=x)$, for all y and $a' \neq a$.

³³ Matt J Kusner et al. "Counterfactual fairness". In: NIPS. vol. 30. Curran Associates, Inc., 2017, pp. 4069–4079.

³⁴ Niki Kilbertus et al. "Avoiding discrimination through causal reasoning". In: NIPS. vol. 30. 2017.

Comparison between HFM³⁵ and baseline fairness measures

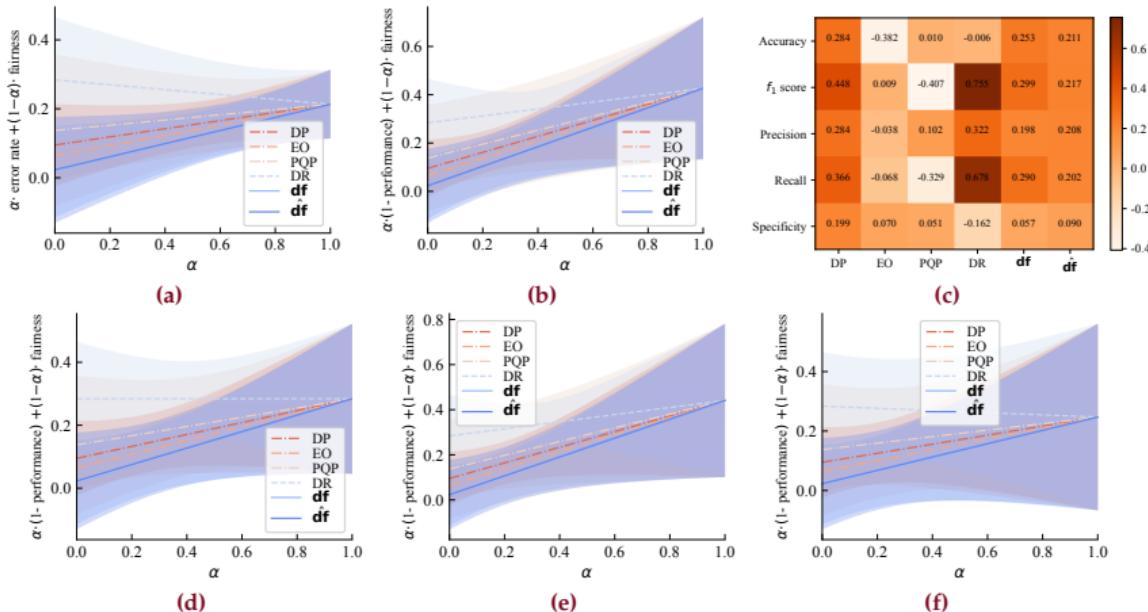


Figure 10: Comparison of baseline fairness measures and the proposed HFM, evaluated on test data. (a) Plot of trade-offs between fairness and error rates per fairness measure; (b) and (d-f) Plots of trade-offs between fairness and $(1 - \text{performance})$ per fairness measures, where performance here are f_1 score, precision, recall/sensitivity, and specificity, respectively. (c) Correlation heatmap between evaluation metric and fairness.

³⁵Bian and Luo, see n. 9; Cruz et al., see n. 28.

Comparison between HFM³⁶ and baseline fairness measures

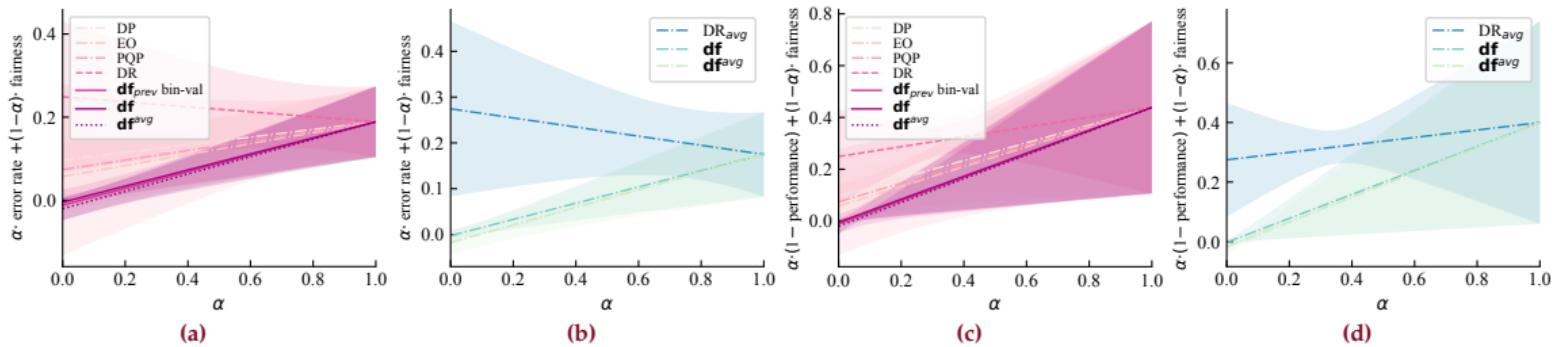


Figure 11: Plots of best test-set fairness-performance trade-offs per fairness metric (the smaller the better).
 (a) Plot of fairness-accuracy trade-off for one single sensitive attribute; (b) Plot of fairness-accuracy trade-off for all sensitive attributes; (c-d) Plots of fairness-f₁ score trade-off for one sensitive attribute and for all sensitive attributes, respectively. Note that the notations in (a) and (c) refer to those in Figure 4, and that in (b) and (d) refer to those in Figure 3.

³⁶Bian, Luo, and Xu, see n. 9; Cruz et al., see n. 28.

Comparison of the state-of-the-art pruning method with *POAF*

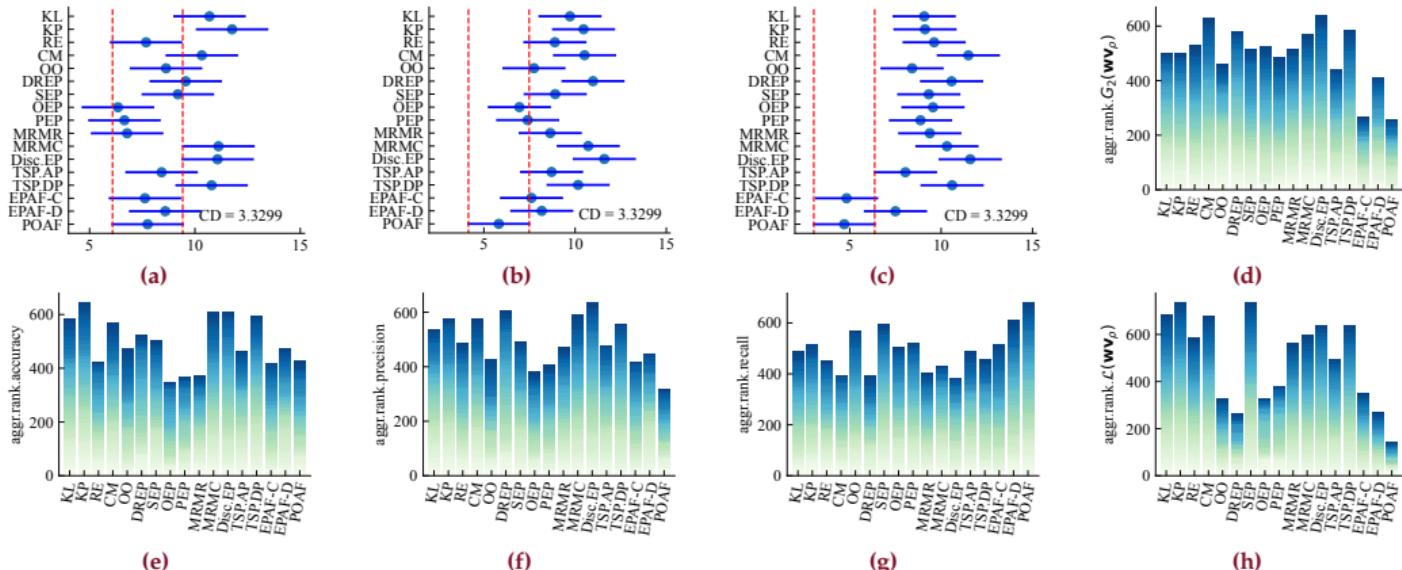


Figure 12: Comparison of the state-of-the-art pruning method with *POAF*, using bagging to conduct homogeneous ensembles. (a–c) Friedman test chart on the test accuracy, precision, and $\mathcal{L}_{\text{bias}}(\mathbf{wv}_\rho)$, respectively, of which each rejects the null hypothesis at the significance level of 5%; (d–h) The aggregated rank for each pruning method over the $\mathcal{L}_{\text{bias}}(\mathbf{wv}_\rho)$, test accuracy, precision, recall, and $\mathcal{L}(\mathbf{wv}_\rho)$, respectively.

Comparison between *EPAF-C* and *EPAF-D*

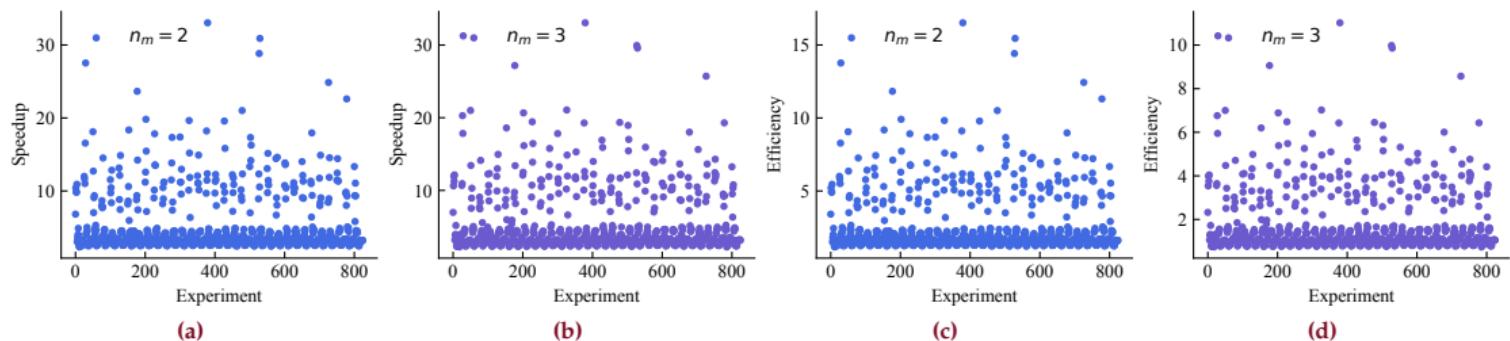


Figure 13: Comparison of speedup and efficiency between *EPAF-C* and *EPAF-D*. (a–b) Speedup with two and three machines, respectively; (c–d) Efficiency with two and three machines, respectively.