

Bridging Explainability and Fairness in Machine Learning: From DR¹ to FairSHAP^{2,3}

Yijun Bian

Department of Computer Science
University of Copenhagen

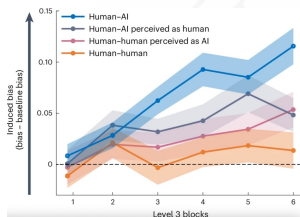
13 May 2026

¹ *Yijun Bian*^{*}, “Improving fairness with ensemble combination: Margin-dependent bounds,” in *the 9th annual ACM Conference on Fairness, Accountability, and Transparency (FAccT 2026)*, accepted.

² Lin Zhu[†], *Yijun Bian*[†], and Lei You^{*}, “FairSHAP: Preprocessing for fairness through attribution-based data augmentation,” in *the 29th International Conference on Artificial Intelligence and Statistics (AISTATS 2026)*, accepted.

³ [†]Equal contribution, ^{*}corresponding author. Joint work (Technical University of Denmark & University of Copenhagen)

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.⁴

When humans and AI interact, even minute perceptual, emotional and social biases—originating either from AI systems or humans—leave human beliefs more biased, potentially forming a feedback loop.⁵

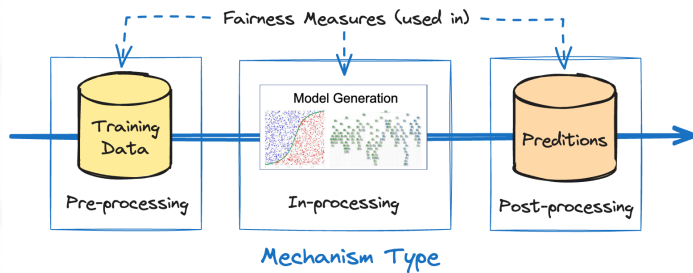
⁴Ziad Obermeyer et al. “Dissecting racial bias in an algorithm used to manage the health of populations”. In: *Science* 366.6464 (2019), pp. 447–453. DOI: 10.1126/science.aax2342; Tiancheng Hu et al. “Generative language models exhibit social identity biases”. In: *Nat Comput Sci* (2024), pp. 1–11; Richard J Chen et al. “Algorithmic fairness in artificial intelligence for medicine and healthcare”. In: *Nat Biomed Eng* 7.6 (2023), pp. 719–742.

⁵Moshe 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 bias⁶

Mechanisms to enhance fairness⁷



Types of fairness measures

⁶*Unintentional*: Limited and coarse features, sample size disparity (less data by definition about minority populations), skewed sample (feedback loops), tainted examples, features that act as proxies; *Intentional*: conscious prejudice.

⁷*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^{6,7}



Group fairness



Individual fairness

Challenging:

incompatibility,⁸ multi-attribute protection, etc.

⁶*Distributive fairness*: group fairness, individual fairness, counterfactual fairness, etc.; *Procedural fairness*

⁷*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.”

⁸Tensions between notions of fairness, between fairness and accuracy, between different methods for achieving fairness

Discriminative risk (DR)¹¹

—from an individual-level aspect

Following the principle of individual fairness (*the treatment/evaluation on one instance should not change solely due to minor changes in its sensitive attributes*), with an instance denoted by $x = (\check{x}, a)$, the fairness quality of one hypothesis⁹ $f(\cdot)$ could be evaluated by

$$\ell_{\text{bias}}(f, x) = \mathbb{I} \left(\underbrace{f(\check{x}, a)}_{\text{model prediction on the raw instance}} \neq \underbrace{f(\check{x}, \tilde{a})}_{\text{model prediction when only sensitive attribute(s) are changed}} \right) \quad (1)$$

Diagram annotations:

- the indicator function** (blue arrow) points to $\mathbb{I}(\cdot)$.
- non-sensitive attributes** (orange arrow) points to \check{x} .
- sensitive attribute(s)** (purple arrow) points to a .
- sensitive attribute(s) that are slightly perturbed (the privileged \leftrightarrow any one of the unprivileged)** (purple arrow) points to \tilde{a} .

similarly to the 0/1 loss, where \tilde{a} is a perturbed¹⁰ $a = [a_1, \dots, a_{n_a}]^T$, $a_i \in \mathcal{A}_i$, $n_a \geq 1$, and $|\mathcal{A}_i| \geq 2$. Note that Eq. (1) is evaluated on only one instance with sensitive attributes x .

⁹The hypothesis used in this equation could indicate an individual classifier or an ensemble classifier.

¹⁰A member belonging to the marginalised group(s) will be perturbed into the privileged group, while a member of the privileged group is perturbed randomly into one of the marginalised groups.

¹¹Yijun Bian. "Improving fairness with ensemble combination: Margin-dependent bounds". In: *FAccT*. Accepted. 2026.

Discriminative risk (DR)¹¹

—from an individual-level aspect
 —from a group-level aspect

To describe this characteristic of the hypothesis on multiple instances, then the **empirical discriminative risk on one dataset** S is expressed as $\hat{L}_{\text{bias}}(f, S) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{bias}}(f, \mathbf{x}_i)$, and the **true discriminative risk**⁹ of the hypothesis **over a data distribution** is $L_{\text{bias}}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\ell_{\text{bias}}(f, \mathbf{x})]$, respectively. Note that the empirical DR on S is an unbiased estimation of the true DR.

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

$$\text{DR}(f) = \mathbb{E}_{\mathcal{D}}[\mathbb{I}(f(\check{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \tilde{\mathbf{a}}))]$$

- Widely applicable, allowing one or more SAs, and each SA allowing binary or multiple values
- Different from existing (group/individual/counterfactual) fairness measures
- **Limitations:** small values; computational results may be affected somehow by a randomness factor¹⁰

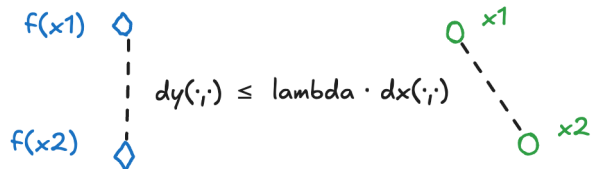
⁹The instances from S are independent identically distributed (i.i.d.) drawn from an input/feature-output/label space $\mathcal{X} \times \mathcal{Y}$ according to an unknown distribution \mathcal{D} .

¹⁰With p probability, the instance will be perturbed into another group.

¹¹Bian, see n. 11.

Differences from existing fairness measures

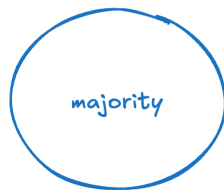
- Two distinctions from *individual fairness* measures
 - 1 relies on the choice of similarity/distance metric
 - 2 instance pairs in comparison coming from original data



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

Differences from existing fairness measures

- 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



for some metric evaluated
on different subgroups:

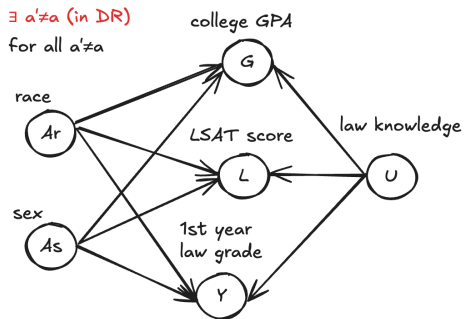


discrepancy between them?

- Four distinctions from *causal fairness*

Differences from existing fairness measures

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



Differences from existing fairness measures

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

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

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

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

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

- 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

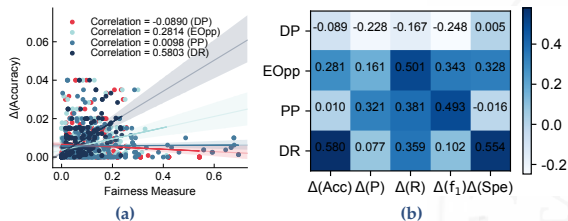


Figure 1. Comparison of the proposed discriminative risk (DR) with three group fairness measures, that is, DP, EOpp, and PP. (a) Scatter diagrams with the degree of correlation, where the x - and y -axes are different fairness measures and the variation of accuracy between the raw and disturbed data.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

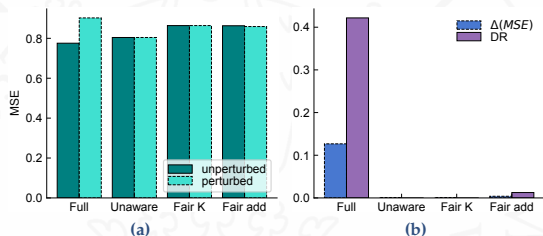


Figure 2. Example: law school success. (a) Test MSE of different models, where 'unperturbed' and 'perturbed' denote the results obtained from the original and disturbed data respectively. (b) The comparison between the change in MSE and DR, which suggests that $\text{DR} \approx 0$ when the corresponding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

Improving fairness with ensemble combination

Consider *weighted voting* as the ensemble combination.¹² For brevity, we record $\sum_{j=1}^m w_j \mathbb{I}(f_j(\mathbf{x}) = y)$ as $q_y(\mathbf{x})$, and $y^* \triangleq \mathbf{w}\mathbf{v}_\rho(\mathbf{x})$. Then we define the *top-1 voting margin* as

$$\gamma_\rho(\mathbf{x}) \triangleq q_{y^*}(\mathbf{x}) - \max_{y \neq y^*} q_y(\mathbf{x}) \in [0, 1].$$

If the ensemble exhibits discrimination risk, it means that *the leading gap in the prediction class has been reversed*, in other words, at least half of the votes for the prediction class will be altered. That is to say, the total weight of the inconsistency between the original prediction and the perturbed prediction has to hold

$$\phi_\rho(\mathbf{x}) \triangleq \mathbb{E}_\rho[\ell_{\text{bias}}(f, \mathbf{x})] = \mathbb{E}_\rho[\mathbb{I}(f(\check{\mathbf{x}}, \mathbf{a}) \neq f(\check{\mathbf{x}}, \tilde{\mathbf{a}}))] \geq \frac{1}{2} \gamma_\rho(\mathbf{x}).$$

Therefore, we have the discriminative risk of one ensemble

$$\ell_{\text{bias}}(\mathbf{w}\mathbf{v}_\rho, \mathbf{x}) \leq \mathbb{I}(\phi_\rho(\mathbf{x}) \geq \frac{1}{2} \gamma_\rho(\mathbf{x})) = \mathbb{I}(\phi_\rho(\mathbf{x}) / \gamma_\rho(\mathbf{x}) \geq \frac{1}{2}), \quad (2)$$

and can discuss some bounds concerning fairness for the weighted vote, inspired by prior work¹³.

¹²where the prediction by an ensemble of m individual classifiers, parameterised by a weight vector $\rho = [w_1, \dots, w_m]^T \in [0, 1]^m$, is given by $\mathbf{w}\mathbf{v}_\rho(\mathbf{x}) = \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{j=1}^m w_j \mathbb{I}(f_j(\mathbf{x}) = y)$, such that $\sum_{j=1}^m w_j = 1$, and w_j is the weight of individual classifier $f_j(\cdot)$.

It could be viewed as $w_j = 1/m$ for all $j \in \{1, 2, \dots, m\}$ in plurality voting and majority voting, meaning that all individual classifiers will be taken into account equally. Note that ties will be resolved arbitrarily, and both parametric and non-parametric models can serve as individual classifiers.

¹³By applying Markov's inequality, second-order Markov's inequality, and Chebyshev-Cantelli inequality (Andrés 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)

Verification of the proposed *margin-dependent* bounds

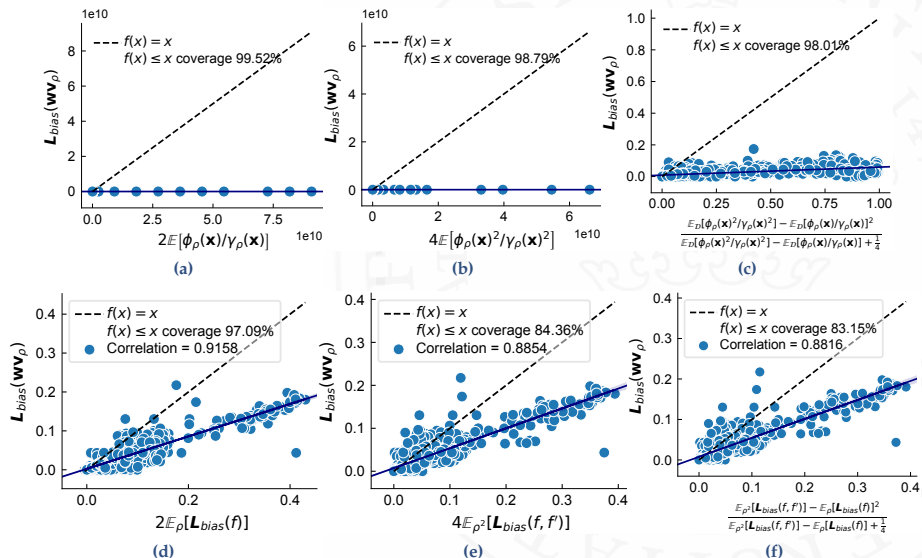


Figure 3. Verification of the proposed margin-dependent bounds, in comparison with prior work

Verification of the proposed *margin-dependent* bounds

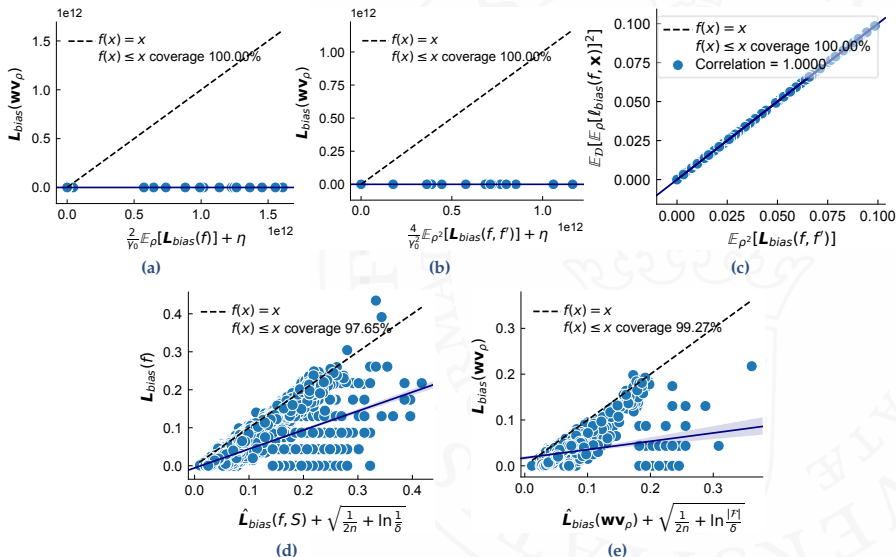
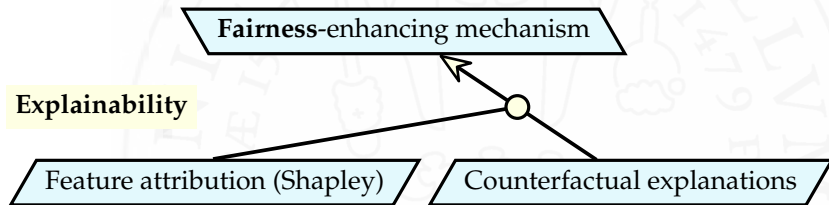


Figure 3. Verification of the proposed margin-dependent bounds' relaxation and PAC bounds

FairSHAP



Preliminaries

Shapley value

For a set of players (or elements) $\mathcal{F} = \{1, 2, \dots, n\}$ and a characteristic function $v : 2^{\mathcal{F}} \rightarrow \mathbb{R}$ that assigns a value to each coalition (or combination) $S \subseteq \mathcal{F}$, the Shapley value of player k is defined as:

$$\phi_k(v) = \sum_{S \subseteq \mathcal{F} \setminus \{k\}} \frac{|S|!(n - |S| - 1)!}{n!} (v(S \cup \{k\}) - v(S)). \quad (3)$$

Discriminative risk

With $\mathbf{x} = (\check{\mathbf{x}}, A)$, where $\check{\mathbf{x}} \in \mathcal{X}_{\mathcal{F} \setminus \{A\}}$ are the non-sensitive features and $A \in \{0, 1\}$ is the sensitive attribute ($0 = \text{unprivileged}$, $1 = \text{privileged}$), discriminative risk (DR) is defined as:

$$L_{\text{DR}}(f, \mathbf{x}) = |f(\check{\mathbf{x}}, A = 0) - f(\check{\mathbf{x}}, A = 1)| \leq \varepsilon.$$

FairSHAP

For instance i , define the coalition value on features $S \subseteq \mathcal{F}$ via an individual-fairness loss, where \mathcal{P} is the joint probability obtained by nearest neighbour matching:

$$v^{(i)}(S) = \mathbb{E}_{\tilde{\mathbf{g}} \sim \mathcal{P}(\tilde{\mathbf{g}}|\mathbf{g}_i)} \left[L_{\text{DR}}(\mathbf{g}_{i,S}; \tilde{\mathbf{g}}_{\mathcal{F} \setminus S}) \right] - \mathbb{E}_{\tilde{\mathbf{g}} \sim \mathcal{P}(\tilde{\mathbf{g}}|\mathbf{g}_i)} \left[L_{\text{DR}}(\tilde{\mathbf{g}}_{\mathcal{F}}) \right], \quad (4a)$$

$$\text{s.t. } \mathcal{P}(\mathbf{g}, \tilde{\mathbf{g}}) = \mathcal{M}_{\text{method}}(\mathcal{G}, \tilde{\mathcal{G}}). \quad (4b)$$

Framework

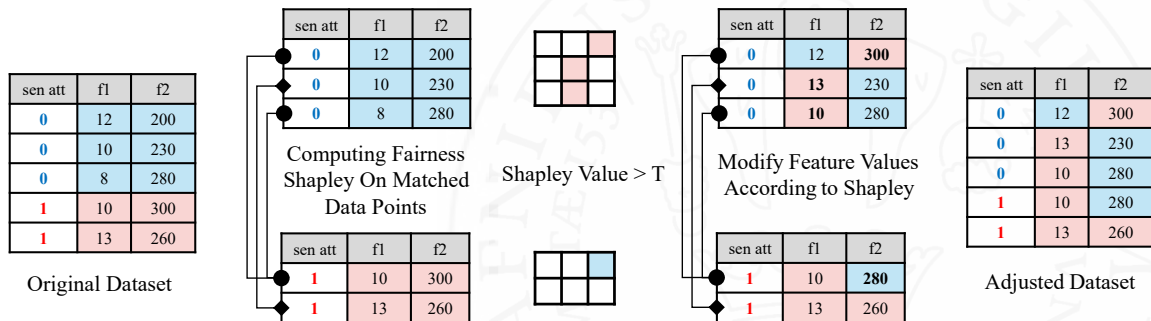
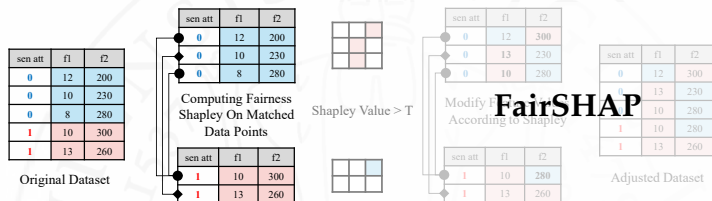


Figure 4. Overall framework of FairSHAP.¹⁴

- **Left:** Partition the original dataset by the sensitive attribute into privileged and unprivileged groups, then perform nearest-neighbour matching between two groups.
- **Middle:** Computed Shapley value matrix on the matched samples.
- **Right:** Modify feature values according to the Shapley attributions and synthesize the adjusted dataset.

¹⁴leverages Shapley value attribution to improve both individual and group fairness; *model-agnostic* and transparent; is broadly applicable to tabular data, supports various models and SHAP algorithms, and can be seamlessly integrated into existing ML pipelines

Algorithm 1



Algorithm 1. Overall framework: Enhancing fairness via matching and Shapley values

Input: Model f , dataset $\mathcal{D} \in \mathbb{R}^{(n+m) \times d}$ with sensitive attribute $A \in \{0, 1\}$, threshold T , matching method $\mathcal{M}_{\text{method}}$, where $\text{method} \in \{\text{NearestNeighbour}, \text{OptimalTransport}\}$

Output: \mathcal{D}_{new}

- 1: Split \mathcal{D} into two subgroups: $\mathcal{G} \in \mathbb{R}^{n \times d}$ (e.g., $A = 0$) and $\tilde{\mathcal{G}} \in \mathbb{R}^{m \times d}$ (e.g., $A = 1$)
- 2: $\mathcal{G}' \leftarrow \text{FairSHAP}(\text{target} = \mathcal{G}, \text{non-target} = \tilde{\mathcal{G}}, \text{model} = f, T, \mathcal{M}_{\text{method}})$ // see Algorithm 2
- 3: $\tilde{\mathcal{G}}' \leftarrow \text{FairSHAP}(\text{target} = \tilde{\mathcal{G}}, \text{non-target} = \mathcal{G}, \text{model} = f, T, \mathcal{M}_{\text{method}})$ // see Algorithm 2
- 4: $\mathcal{D}_{\text{new}} \leftarrow \text{Concat}(\mathcal{G}', \tilde{\mathcal{G}}') \in \mathbb{R}^{(n+m) \times d}$
- 5: **return** \mathcal{D}_{new}

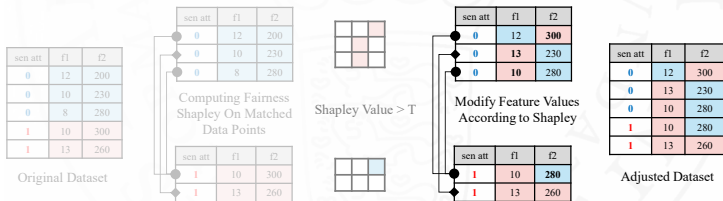
Algorithm 2

Algorithm 2. FairSHAP

Input: Target group $\mathcal{G} \in \mathbb{R}^{n \times d}$ and non-target group $\tilde{\mathcal{G}} \in \mathbb{R}^{m \times d}$, model f , threshold T , matching method $\mathcal{M}_{\text{method}}$

Output: modified dataset \mathcal{G}'

- 1: Use $\mathcal{M}_{\text{method}}(\mathcal{G}, \tilde{\mathcal{G}})$ to obtain joint probability $\mathcal{P}(\mathbf{g}, \tilde{\mathbf{g}}) \in \mathbb{R}^{n \times m}$
- 2: Use Eqs. (3) and (4) to obtain Shapley value matrix $\phi \in \mathbb{R}^{n \times d}$
- 3: Initialize reference data $\mathcal{B} \leftarrow \mathbf{0}_{n \times d}$
- 4: **for** $i = 1$ **to** n **do**
- 5: $j^* \leftarrow \arg \max_{1 \leq j \leq m} \mathcal{P}_{ij}$
- 6: $\mathcal{B}_{i,:} \leftarrow \tilde{\mathcal{G}}_{j^*,:}$
- 7: Let $\mathcal{G}' \leftarrow \mathcal{G}$ ($\mathcal{G}' \in \mathbb{R}^{n \times d}$)
- 8: **for** $i = 1$ **to** n **do**
- 9: **for** $k = 1$ **to** d **do**
- 10: **if** $\phi_{i,k} \geq T$ **then**
- 11: $\mathcal{G}'_{i,k} \leftarrow \mathcal{B}_{i,k}$
- 12: **return** \mathcal{G}'



Qualitative results

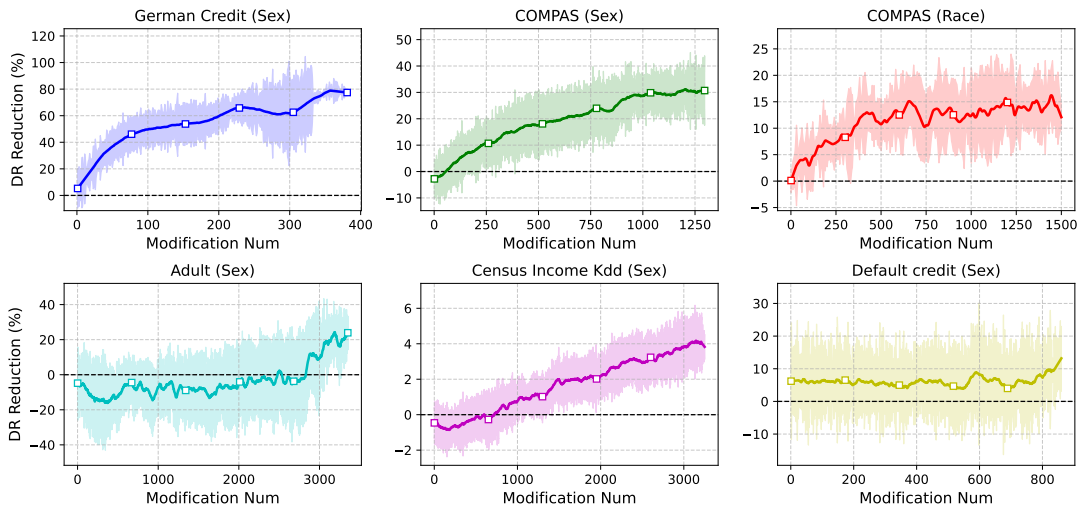


Figure 5. Percentage reduction in the DR across different datasets. The x-axis denotes the number of modifications applied (up to the maximum required under a fairness threshold $T = 0.05$), while the y-axis indicates the relative in DR, expressed as a percentage of the original value.

Qualitative results

Table 1. Compare FairSHAP with other fairness mitigation methods across different datasets.¹⁵ Here, “s.a.” denotes the sensitive attribute used in fairness evaluation. Note that TrainingAR: Training Set Adjustment Rate; TestAN: Test Set Adjustment Necessity. Data Fidelity is measured using the Wasserstein Distance to quantify the difference between the original and adjusted data distributions.

Dataset (s.a.)	Methods	Accuracy	DR	DP	EOpp	PP	Data Fidelity	TrainingAR	TestAN
German (sex)	Baseline	0.6650±0.0257	0.0785±0.0211	0.0512±0.0346	0.1287±0.0590	0.1341±0.0486	—	—	No
	RSA	0.6590±0.0287	—	0.0522±0.0246	0.1524±0.1007	0.2189±0.0815	—	—	Yes
	CR	<u>0.6680±0.0238</u>	0.0028±0.0029	0.0844±0.0557	0.1559±0.0609	0.0723±0.0330	0.0183±0.0211	0.9615	Yes
	DIR	0.6720±0.0337	0.0966±0.0112	0.0946±0.0373	0.1737±0.0729	0.1529±0.0634	0.0155±0.0440	0.0774	Yes
	FairSHAP	0.6630±0.0275	<u>0.0243±0.0112</u>	0.0301±0.0347	0.1126±0.0783	0.1852±0.1074	0.0049±0.0085	0.0156	No
COMPAS (sex)	Baseline	0.6698±0.0051	0.0883±0.0064	0.1548±0.0241	0.1243±0.0510	0.0492±0.0084	—	—	No
	RSA	0.6676±0.0067	—	0.1075±0.0160	0.0842±0.0475	0.0635±0.0266	—	—	Yes
	CR	0.6679±0.0045	0.0082±0.0070	0.1407±0.0248	0.1291±0.0317	0.0714±0.0517	0.0189±0.0193	0.9174	Yes
	DIR	0.6644±0.0098	0.1150±0.0091	<u>0.1155±0.0239</u>	0.0952±0.0359	0.0747±0.0370	0.0387±0.0640	0.0650	Yes
	FairSHAP	0.6609±0.0106	<u>0.0629±0.0091</u>	0.1326±0.0407	<u>0.0985±0.0603</u>	0.0452±0.0383	0.0025±0.0048	0.0113	No
COMPAS (race)	Baseline	0.6689±0.0108	0.0995±0.0076	0.1436±0.0209	0.1438±0.0233	0.0522±0.0406	—	—	No
	RSA	0.6623±0.0078	—	0.2069±0.0128	0.2024±0.0349	0.0575±0.0291	—	—	Yes
	CR	0.6611±0.0112	0.0418±0.0092	0.1502±0.0341	0.1621±0.0530	0.0592±0.0367	0.0250±0.0222	0.8920	Yes
	DIR	0.6149±0.0286	0.1185±0.0181	<u>0.1359±0.1241</u>	0.1117±0.0945	0.0399±0.0338	0.0512±0.0736	0.0701	Yes
	FairSHAP	0.6627±0.0069	<u>0.0842±0.0049</u>	0.1344±0.0332	0.1568±0.0343	<u>0.0508±0.0469</u>	0.0040±0.0055	0.0126	No

¹⁵“Kdd” is an abbreviation for the Census Income KDD dataset

Qualitative results

Table 1. Compare FairSHAP with other fairness mitigation methods (cont).

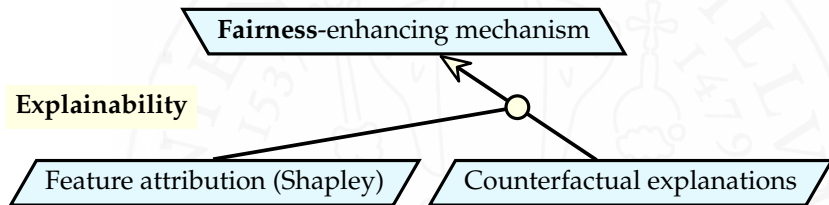
Dataset (s.a.)	Methods	Accuracy	DR	DP	EOpp	PP	Data Fidelity	TrainingAR	TestAN
Adult (sex)	Baseline	0.8722 \pm 0.0033	0.0315 \pm 0.0037	0.1805 \pm 0.0066	0.0735 \pm 0.0275	0.0275 \pm 0.0321	—	—	No
	RSA	0.8721 \pm 0.0024	—	0.1770 \pm 0.0066	<u>0.0678</u> \pm 0.0277	0.0272 \pm 0.0252	—	—	Yes
	CR	0.8706 \pm 0.0029	0.0000 \pm 0.0000	0.1824 \pm 0.0055	0.0955 \pm 0.0243	0.0278 \pm 0.0173	0.0167 \pm 0.0391	0.9887	Yes
	DIR	0.8550 \pm 0.0067	<u>0.0499</u> \pm 0.0076	<u>0.1607</u> \pm 0.0157	0.0772 \pm 0.0624	0.0360 \pm 0.0253	<u>0.0046</u> \pm 0.0417	0.0081	Yes
	FairSHAP	0.8692 \pm 0.0046	<u>0.0273</u> \pm 0.0047	0.1558 \pm 0.0130	0.0393 \pm 0.0254	0.0474 \pm 0.0319	0.0010 \pm 0.0073	0.0012	No
Adult (race)	Baseline	0.8721 \pm 0.0033	0.0398 \pm 0.0025	0.1034 \pm 0.0110	0.0808 \pm 0.0326	0.0302 \pm 0.0265	—	—	No
	RSA	0.8710 \pm 0.0044	—	0.0971 \pm 0.0043	0.0864 \pm 0.0347	0.0427 \pm 0.0142	—	—	Yes
	CR	0.8713 \pm 0.0033	0.0000 \pm 0.0000	0.1008 \pm 0.0115	0.0983 \pm 0.0389	0.0480 \pm 0.0235	0.0300 \pm 0.0450	0.9620	Yes
	DIR	0.8320 \pm 0.0173	0.0740 \pm 0.0209	0.0703 \pm 0.0515	0.0871 \pm 0.0355	0.0482 \pm 0.0730	0.0252 \pm 0.0480	0.0089	Yes
	FairSHAP	0.8720 \pm 0.0023	<u>0.0284</u> \pm 0.0017	<u>0.0851</u> \pm 0.0155	0.0287 \pm 0.0277	0.0259 \pm 0.0318	0.0030 \pm 0.0084	0.0014	No
Kdd (sex)	Baseline	0.9377 \pm 0.0026	0.0767 \pm 0.0011	0.0012 \pm 0.0008	0.0012 \pm 0.0008	0.0796 \pm 0.0054	—	—	No
	RSA	<u>0.9380</u> \pm 0.0026	—	0.0008 \pm 0.0004	0.0008 \pm 0.0004	0.0796 \pm 0.0055	—	—	Yes
	CR	0.9377 \pm 0.0025	0.0000 \pm 0.0000	0.0013 \pm 0.0004	0.0013 \pm 0.0003	0.0796 \pm 0.0054	0.0007 \pm 0.0013	0.9790	Yes
	DIR	0.9377 \pm 0.0028	0.0766 \pm 0.0014	0.0012 \pm 0.0010	0.0013 \pm 0.0011	0.0797 \pm 0.0055	0.0001 \pm 0.0007	0.0013	Yes
	FairSHAP	0.9381 \pm 0.0029	<u>0.0732</u> \pm 0.0020	<u>0.0008</u> \pm 0.0007	0.0006 \pm 0.0008	0.0794 \pm 0.0060	0.0000 \pm 0.0000	0.0003	No
DefaultCredit (sex)	Baseline	0.8141 \pm 0.0059	0.0226 \pm 0.0019	0.0335 \pm 0.0072	0.0348 \pm 0.0211	0.0269 \pm 0.0127	—	—	No
	RSA	0.8144 \pm 0.0059	—	0.0259 \pm 0.0063	<u>0.0244</u> \pm 0.0122	0.0154 \pm 0.0172	—	—	Yes
	CR	<u>0.8144</u> \pm 0.0068	0.0003 \pm 0.0002	0.0339 \pm 0.0077	0.0419 \pm 0.0203	0.0363 \pm 0.0203	0.0051 \pm 0.0064	0.9844	Yes
	DIR	0.8138 \pm 0.0051	0.0224 \pm 0.0027	0.0308 \pm 0.0091	0.0334 \pm 0.0244	<u>0.0255</u> \pm 0.0214	0.0023 \pm 0.0056	0.0741	Yes
	FairSHAP	0.8145 \pm 0.0050	<u>0.0214</u> \pm 0.0026	<u>0.0306</u> \pm 0.0054	0.0216 \pm 0.0164	0.0289 \pm 0.0075	0.0001 \pm 0.0003	0.0004	No

Qualitative results

Table 2. Compare FairSHAP with ablation studies across different datasets. The sensitive attribute is sex in all cases. Best and second-best results for fairness measures (DR, DP, EOpp, PP) are **bold**.

Dataset	Methods	Accuracy	DR	DP	EOpp	PP
German	Baseline	0.6650 \pm 0.0257	0.0785 \pm 0.0211	0.0512 \pm 0.0346	0.1287 \pm 0.0590	0.1341\pm0.0486
	Ablation study 1	0.6690\pm0.0237	0.0709 \pm 0.0239	0.0446\pm0.0106	0.0972 \pm 0.1091	0.1463 \pm 0.1133
	Ablation study 2	0.6470 \pm 0.0452	0.0640 \pm 0.0176	0.0708 \pm 0.0510	0.1619 \pm 0.0730	0.1475 \pm 0.1384
	FairSHAP	0.6630 \pm 0.0275	0.0243\pm0.0112	0.0301\pm0.0347	0.1126\pm0.0783	0.1852 \pm 0.1074
COMPAS	Baseline	0.6698 \pm 0.0051	0.0883 \pm 0.0064	0.1548 \pm 0.0241	0.1243 \pm 0.0510	0.0492 \pm 0.0084
	Ablation study 1	0.6713\pm0.0090	0.0903 \pm 0.0110	0.1519 \pm 0.0114	0.1315 \pm 0.0379	0.0490 \pm 0.0465
	Ablation study 2	0.6699 \pm 0.0076	0.0721 \pm 0.0059	0.1537 \pm 0.0177	0.1125 \pm 0.0431	0.0318\pm0.0121
	FairSHAP	0.6609 \pm 0.0106	0.0629\pm0.0091	0.1326\pm0.0407	0.0985\pm0.0603	0.0452 \pm 0.0383
Adult	Baseline	0.8722 \pm 0.0033	0.0315 \pm 0.0037	0.1805 \pm 0.0066	0.0735 \pm 0.0275	0.0275\pm0.0321
	Ablation study 1	0.8717 \pm 0.0031	0.0332 \pm 0.0020	0.1842 \pm 0.0048	0.0932 \pm 0.0186	0.0344 \pm 0.0228
	Ablation study 2	0.8722\pm0.0016	0.0278 \pm 0.0018	0.1816 \pm 0.0084	0.0761 \pm 0.0325	0.0442 \pm 0.0266
	FairSHAP	0.8692 \pm 0.0046	0.0273\pm0.0047	0.1558\pm0.0130	0.0393\pm0.0254	0.0474 \pm 0.0319
CensusIncome	Baseline	0.9377 \pm 0.0026	0.0767 \pm 0.0011	0.0012 \pm 0.0008	0.0012 \pm 0.0008	0.0796 \pm 0.0054
	Ablation study 1	0.9377 \pm 0.0027	0.0765 \pm 0.0019	0.0010 \pm 0.0006	0.0011 \pm 0.0007	0.0797 \pm 0.0055
	Ablation study 2	0.9378 \pm 0.0024	0.0744 \pm 0.0018	0.0010 \pm 0.0003	0.0010 \pm 0.0003	0.0796 \pm 0.0056
	FairSHAP	0.9381\pm0.0029	0.0732\pm0.0020	0.0008\pm0.0007	0.0006\pm0.0008	0.0794\pm0.0060
DefaultCredit	Baseline	0.8141 \pm 0.0059	0.0226 \pm 0.0019	0.0335 \pm 0.0072	0.0348 \pm 0.0211	0.0269\pm0.0127
	Ablation study 1	0.8149 \pm 0.0061	0.0221 \pm 0.0027	0.0327 \pm 0.0067	0.0431 \pm 0.0111	0.0369 \pm 0.0166
	Ablation study 2	0.8150\pm0.0051	0.0217 \pm 0.0031	0.0331 \pm 0.0070	0.0259 \pm 0.0135	0.0325 \pm 0.0137
	FairSHAP	0.8145 \pm 0.0050	0.0214\pm0.0026	0.0306\pm0.0054	0.0216\pm0.0164	0.0289 \pm 0.0075

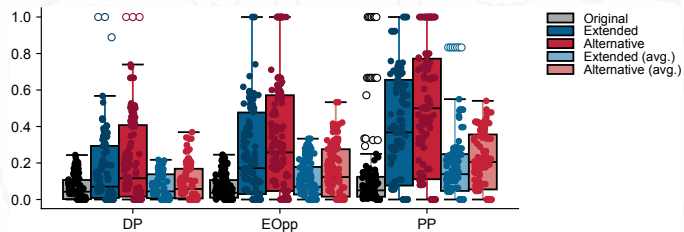
Takeaways



- Prioritising individual fairness can offer stronger and more flexible leverage than focusing only on group fairness
- Connecting explainability with fairness makes mitigation more interpretable, targeted, and practical

Additionally¹⁵

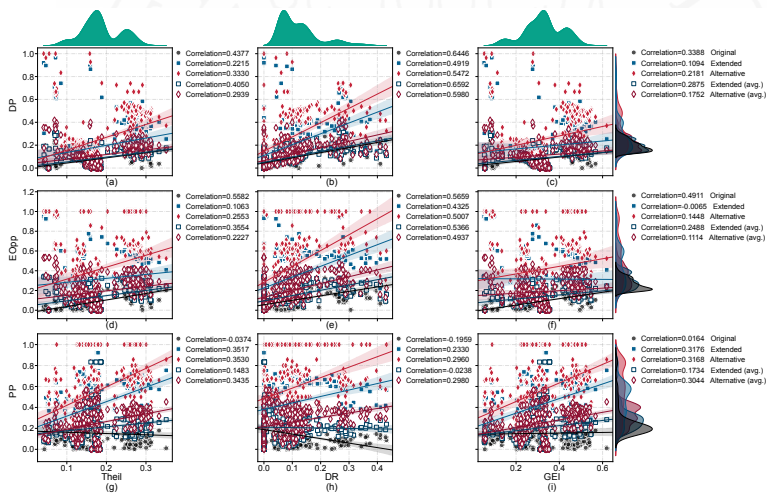
- Binarlisation systematically underestimates discrimination



¹⁵Yijun Bian et al. "Revisiting some misconceptions and limitations in algorithmic fairness". In: (2025).

Additionally¹⁵

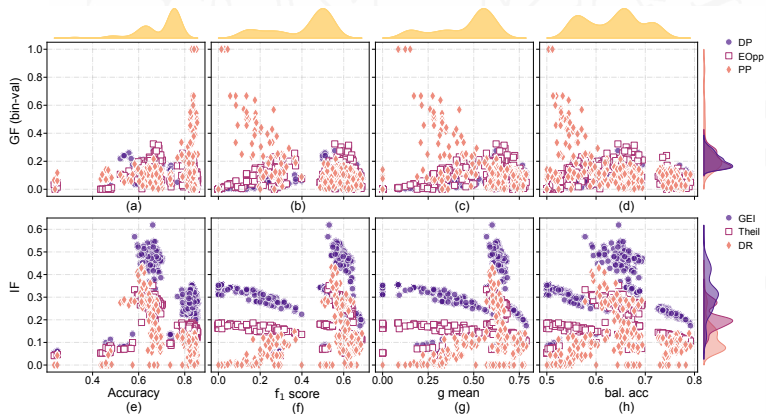
- Binarlisation systematically underestimates discrimination
- The myth of incompatibility between individual and group fairness



¹⁵Bian et al., see n. 15.

Additionally¹⁵

- Binarlisation systematically underestimates discrimination
- The myth of incompatibility between individual and group fairness
- The elusive trade-off between utility performance and fairness



¹⁵Bian et al., see n. 15.