

# Increasing Fairness via Combination with Learning Guarantees<sup>1</sup>

Yijun BIAN

National University of Singapore

26 October 2023

---

<sup>1</sup>Yijun Bian et al. “Increasing Fairness via Combination with Learning Guarantees”. In: *arXiv preprint arXiv:2301.10813* (2023). Under Review.

# Overview

- 1 Background
- 2 Methodology
- 3 Discussions
- 4 Appendix

# Examples of bias



AI detectors were more likely to flag writing by international students (i.e., non-native speakers) as AI-generated<sup>2</sup>

<sup>2</sup>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.

# Examples of bias

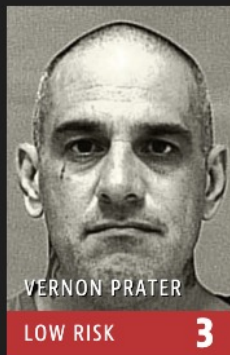


When people of color have complex medical needs, they are less likely to be referred to programmes that provide more individualised care<sup>2</sup>

<sup>2</sup>Linda Nordling. “A fairer way forward for AI in health care”. In: *Nature* 573.7775 (2019), S103–S103.

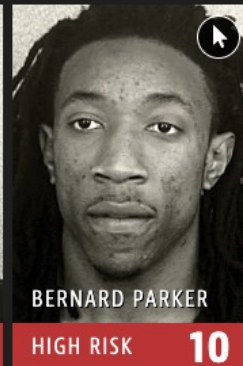
# Examples of bias

## Two Petty Theft Arrests



*Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.*

## Two Drug Possession Arrests

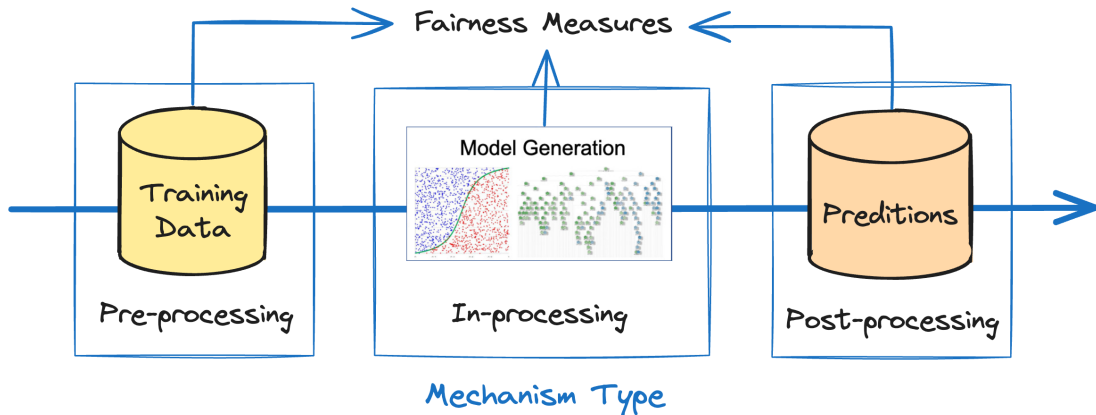


*Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.*

Black defendants were mislabelled as high risk more often than white defendants<sup>2</sup>

<sup>2</sup>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.

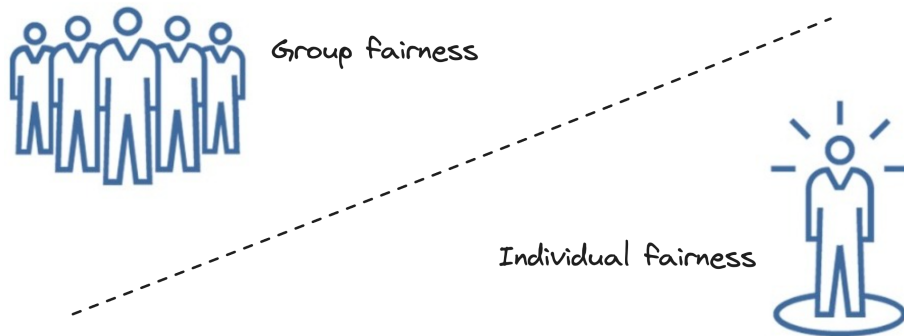
# Mechanisms to enhance fairness<sup>3</sup>



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

<sup>3</sup>Simon Caton and Christian Haas. "Fairness in machine learning: A survey". In: *ACM Comput Surv* (2020); Sorelle A Friedler et al. "A comparative study of fairness-enhancing interventions in machine learning". In: *FAT*. Atlanta, GA, USA: Association for Computing Machinery, 2019, pp. 329–338; Cynthia Dwork et al. "Decoupled classifiers for group-fair and efficient machine learning". In: *FAT*. vol. 81. PMLR, 2018, pp. 119–133.

# Types of fairness measures



*Group fairness*<sup>4</sup> 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>4</sup>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; 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; 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.

# Our target in this work

## Research gap

- The [hard compatibility](#) among these measures means that unfair decisions may still exist even if one of them is satisfied<sup>5</sup>
- The possibility of [theoretical guarantees](#) of boosting fairness is rarely discussed in the existing fairness-aware ensemble-based methods<sup>6</sup>

---

<sup>5</sup>Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and machine learning*. fairmlbook.org, 2019; 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; Geoff Pleiss et al. “On fairness and calibration”. In: *NIPS*. vol. 30. 2017; Hardt, Price, and Srebro, see n. 4.

<sup>6</sup>Vasileios Iosifidis and Eirini Ntoutsi. “AdaFair: Cumulative fairness adaptive boosting”. In: *CIKM*. New York, NY, USA: ACM, 2019, pp. 781–790; Wenbin Zhang et al. “FARF: A fair and adaptive random forests classifier”. In: *PAKDD*. Springer. 2021, pp. 245–256; André F Cruz et al. “FairGBM: Gradient Boosting with Fairness Constraints”. In: *ICLR*. 2023.



# Our target in this work

## Research gap

- The **hard compatibility** among these measures means that unfair decisions may still exist even if one of them is satisfied<sup>5</sup>
- The possibility of **theoretical guarantees** of boosting fairness is rarely discussed in the existing fairness-aware ensemble-based methods<sup>6</sup>

## Questions that we endeavour to answer

- 1 *How to properly measure the discriminative level of a classifier from both individual and group fairness aspects?*
- 2 *Can fairness be boosted with some learning guarantee? Will COMBINATION help mitigate discrimination in multiple biased individual classifiers?*

---

<sup>5</sup>Barocas, Hardt, and Narayanan, see n. 5; Berk et al., see n. 5; Pleiss et al., see n. 5; Hardt, Price, and Srebro, see n. 4.

<sup>6</sup>Iosifidis and Ntoutsi, see n. 6; Zhang et al., see n. 6; Cruz et al., see n. 6.

# Overview

- 1 Background
- 2 Methodology**
- 3 Discussions
- 4 Appendix

# Research question recap

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

# Discriminative risk (DR) —from an individual aspect

Following the principle of individual fairness, the fairness quality of one hypothesis<sup>7</sup>  $f(\cdot)$  could be evaluated by

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

Diagram illustrating the components of the equation:

- the indicator function**: Points to  $\mathbb{I}(\cdot)$ .
- general attributes**: Points to  $\check{\mathbf{x}}$  in both  $f(\check{\mathbf{x}}, a)$  and  $f(\check{\mathbf{x}}, \tilde{a})$ .
- sensitive attribute(s)**: Points to  $a$  and  $\tilde{a}$ .
- sensitive attribute(s) that are slightly disturbed**: Points to  $\tilde{a}$ .
- model prediction on the raw instance**: Points to  $f(\check{\mathbf{x}}, a)$ .
- model prediction when only sensitive attribute(s) are changed**: Points to  $f(\check{\mathbf{x}}, \tilde{a})$ .

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

<sup>7</sup>The 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**<sup>8</sup> of the hypothesis **over a data distribution** is

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

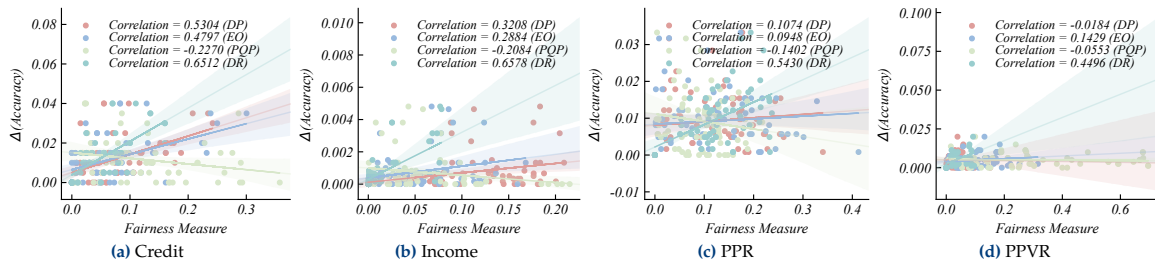
discriminative risk of  $f(\cdot)$  on one instance

respectively.

<sup>8</sup>The 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  $\mathcal{D}$ .

# Empirical results of *DR* in comparison with group fairness measures<sup>9,10</sup>

**Observation:** DR captures better the characteristic of the changed treatment



**Figure 1:** Comparison of the proposed DR with three group fairness measures, that is, DP, EO, and PQP. (a–d) Scatter diagrams with the degree of correlation on the credit, income, ppr, and ppvr datasets, respectively, where the x- and y-axes are different fairness measures and the variation of accuracy between the raw and disturbed data.

<sup>9</sup>They are demographic parity (DP) (Feldman et al., see n. 4; Gajane and Pechenizkiy, see n. 4), equality of opportunity (EO) (Hardt, Price, and Srebro, see n. 4), and predictive quality parity (PQP) (Chouldechova, see n. 4; Verma and Rubin, see n. 4).

<sup>10</sup>Five public datasets that we use include Ricci, Credit, Income, PPR, and PPVR, aka. Propublica-Recidivism and Propublica-Violent-Recidivism.

# Research question recap

2. *Can fairness be boosted with some learning guarantee? Will COMBINATION help mitigate discrimination in multiple biased individual classifiers?*

# Oracle bounds of fairness

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

$$\ell_{\text{fair}}(\mathbf{w}\mathbf{v}_{\rho}, x) \leq \mathbb{I}(\mathbb{E}_{\rho}[\mathbb{I}(f(\tilde{x}, a) \neq f(\tilde{x}, \tilde{a}))] \geq 0.5). \quad (4)$$

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

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

## Meaning of $\mathbf{w}\mathbf{v}_{\rho}(\cdot)$

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  **$\rho$ -weighted majority vote**  $\mathbf{w}\mathbf{v}_{\rho}(\cdot)$  predicts

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

where  $\rho$  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  $\rho$ -weighted half** of the classifiers **have made a discriminative decision** and, therefore,

$$\ell_{\text{fair}}(\mathbf{w}\mathbf{v}_{\rho}, x) \leq \mathbb{I}(\mathbb{E}_{\rho}[\mathbb{I}(f(\check{x}, a) \neq f(\check{x}, \tilde{a}))] \geq 0.5). \quad (4)$$

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

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

## Theorem 1 (First-order oracle bound)

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

discriminative risk of an individual classifier  $f$

$$\mathcal{L}_{\text{fair}}(\mathbf{w}\mathbf{v}_{\rho}) \leq 2 \mathbb{E}_{\rho}[\mathcal{L}_{\text{fair}}(f)]. \quad (5)$$

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,<sup>11</sup> by

hypothesis  $f(\cdot)$  predicts differently for similar instances

hypothesis  $f'(\cdot)$  also predicts differently for them

$$\ell_{\text{fair}}(f, f', x) = \mathbb{I}(f(\tilde{x}, a) \neq f(\tilde{x}, \tilde{a}) \wedge f'(\tilde{x}, a) \neq f'(\tilde{x}, \tilde{a})) . \quad (6)$$

tandem discriminative risk

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

$$\ell_{\text{fair}}(f, f, x) = \ell_{\text{fair}}(f, x) . \quad (7)$$

when  $f'(\cdot)$  and  $f(\cdot)$  are identical

discriminative risk of  $f(\cdot)$

<sup>11</sup>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.

# Oracle bounds of fairness (cont.)

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

## Theorem 3 (Second-order oracle bound)

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

*tandem discriminative risk of two individuals  $f$  and  $f'$*

$$\mathcal{L}_{\text{fair}}(\mathbf{w}\mathbf{v}_\rho) \leq 4 \mathbb{E}_{\rho^2}[\mathcal{L}_{\text{fair}}(f, f')] . \quad (8)$$

*the worst case is controlled to a constant multiple*

### Lemma 2

In multi-class classification,

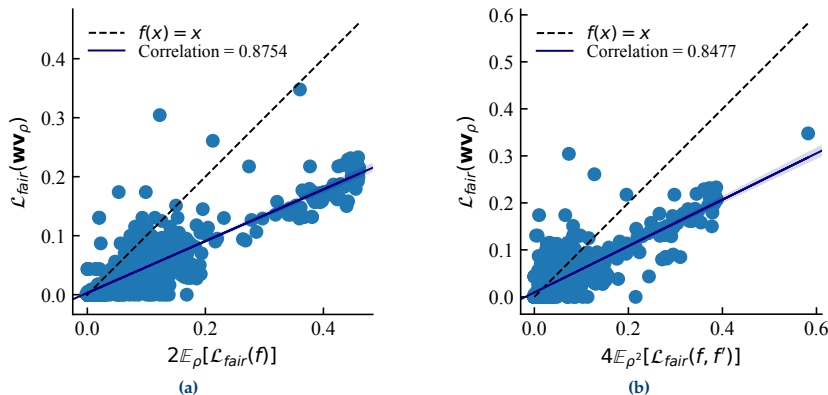
*discriminative risk of  $f(\cdot)$*

$$\mathbb{E}_{\rho^2}[\mathcal{L}_{\text{fair}}(f, f')] = \mathbb{E}_{\mathcal{D}}[\mathbb{E}_{\rho}[\ell_{\text{fair}}(f, x)]^2] . \quad (9)$$

*the expected tandem discriminative risk*

# Empirical results of oracle bounds

**Observation:** The discriminative risk (DR) of an ensemble is indeed smaller than the bounds presented in Theorems 1 and 3 in most cases, indicating that these inequalities are reliable



**Figure 2:** Correlation for oracle bounds. (a–b) Correlation between  $\mathcal{L}_{\text{fair}}(\mathbf{w}\mathbf{v}_\rho)$  and oracle bounds, where  $\mathcal{L}_{\text{fair}}(\mathbf{w}\mathbf{v}_\rho)$  is indicated on the vertical axis and the horizontal axes represent the right-hand sides of inequalities (5), and (8), respectively.

# Overview

- 1 Background
- 2 Methodology
- 3 Discussions**
- 4 Appendix

# Summary<sup>12</sup>

*RQ 2. Can fairness be boosted with some learning guarantee? Will COMBINATION help mitigate discrimination in multiple biased individual classifiers?*

*Ensemble combination: fairness can be boosted without being dependent on specific (hyper-)parameters*

$$\mathcal{L}_{\text{fair}}(\mathbf{w}\mathbf{v}_{\rho}) \leq 2 \mathbb{E}_{\rho} [ \mathcal{L}_{\text{fair}}(f) ]$$

cf. Theorem 1

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

cf. Theorem 3

<sup>12</sup>**P.S.** Please refer to our paper for full methodology and empirical results

# Summary<sup>12</sup>

*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,

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

DR is widely applicable, with two reasons enlarging its applicable fields/scenarios:

- ① suitable for both binary and multi-class classification
- ② allows one or multiple sensitive attributes, and each sensitive attribute allows binary and multiple values

<sup>12</sup>**P.S.** Please refer to our paper for full methodology and empirical results

# Future work

## 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 sensitive attributes may vary, although its property remains

## PROS

1. *Discriminative risk (DR)* is widely applicable
2. *Ensemble combination*: fairness can be boosted without *being dependent on specific (hyper-)parameters*



*Thanks! Questions?*

# Overview

- 1 Background
- 2 Methodology
- 3 Discussions
- 4 Appendix**

# Sources of bias<sup>13</sup>



Biases from data

Data collected from

- biassed device measurements
- erroneous reports
- historically biassed human decisions
- or other reasons

---

<sup>13</sup>Verma and Rubin, see n. 4.

# Sources of bias<sup>13</sup>



Biases from data



Biases from algorithms<sup>a</sup>

<sup>a</sup>e.g., caused by proxy attributes for sensitive attributes or tendentious algorithmic objectives

<sup>13</sup>Verma and Rubin, see n. 4.

# 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

$$\text{wv}_\rho(x) = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \sum_{j=1}^m w_j \mathbb{I}(f_j(x) = y) . \quad (10)$$

Diagram annotations:

- weight corresponding to  $f_j(\cdot)$*  (purple arrow pointing to  $w_j$ )
- individual classifier* (purple arrow pointing to  $f_j(x)$ )
- ensemble combination via weighted vote* (orange arrow pointing to  $\text{wv}_\rho(x)$ )
- optional choice of labels* (orange arrow pointing to  $y$ )

where a function  $f \in \mathcal{F}: \mathcal{X} \mapsto \mathcal{Y}$  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

$$\mathbf{wv}_\rho(x) = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \sum_{j=1}^m \overbrace{w_j}^{\text{weight corresponding to } f_j(\cdot)} \overbrace{\mathbb{I}(f_j(x) = y)}^{\text{individual classifier}} . \quad (10)$$

*ensemble combination via weighted vote* *optional choice of labels*

where a function  $f \in \mathcal{F}: \mathcal{X} \mapsto \mathcal{Y}$  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  *$\rho$ -weighted majority vote*  $\mathbf{wv}_\rho(\cdot)$  predicts

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

*potential  $\rho$  corresponding to an ensemble over  $[0, 1]^m$*

# Oracle bounds regarding fairness for weighted vote

## Theorem 4 (C-tandem oracle bound)

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

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

$$\mathcal{L}_{fair}(\mathbf{w}\mathbf{v}_\rho) \leq \frac{\mathbb{E}_{\rho^2}[\mathcal{L}_{fair}(f, f')] - \mathbb{E}_\rho[\mathcal{L}_{fair}(f)]^2}{\mathbb{E}_{\rho^2}[\mathcal{L}_{fair}(f, f')] - \mathbb{E}_\rho[\mathcal{L}_{fair}(f)] + \frac{1}{4}}. \quad (11)$$

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

*tandem discriminative risk* *discriminative risk*

# Oracle bounds regarding fairness for weighted vote

## Theorem 4 (C-tandem oracle bound)

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

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

$$\mathcal{L}_{fair}(\mathbf{w}\mathbf{v}_\rho) \leq \frac{\mathbb{E}_{\rho^2}[\mathcal{L}_{fair}(f, f')] - \mathbb{E}_\rho[\mathcal{L}_{fair}(f)]^2}{\mathbb{E}_{\rho^2}[\mathcal{L}_{fair}(f, f')] - \mathbb{E}_\rho[\mathcal{L}_{fair}(f)] + \frac{1}{4}}. \quad (11)$$

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

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-Bayesian 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 Hoeffding's inequality, we can deduct generalisation bounds presented in Theorems 5 and 6.



# PAC-Bayesian 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  $f(\cdot)$ ,

$$\underbrace{\mathcal{L}_{\text{fair}}(f)}_{\text{discriminative risk of a hypothesis}} \leq \underbrace{\hat{\mathcal{L}}_{\text{fair}}(f, S)}_{\text{empirical discriminative risk of this hypothesis}} + \underbrace{\sqrt{\frac{1}{2n} \ln \frac{1}{\delta}}}_{\text{the worst case is controlled with a specific bound}}. \quad (12)$$

## 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  $\rho$  on  $\mathcal{F}$ ,

$$\underbrace{\mathcal{L}_{\text{fair}}(\mathbf{w}\mathbf{v}_{\rho})}_{\text{discriminative risk of an ensemble}} \leq \underbrace{\hat{\mathcal{L}}_{\text{fair}}(\mathbf{w}\mathbf{v}_{\rho}, S)}_{\text{empirical discriminative risk of this ensemble}} + \underbrace{\sqrt{\frac{1}{2n} \log \frac{|\mathcal{F}|}{\delta}}}_{\text{the worst case is controlled with a specific bound}}. \quad (13)$$

# Our distinction

Despite the similar names of “first- and second-order oracle bounds” from our inspiration,<sup>14</sup> 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

$$\ell_{\text{err}}(f, x) = \mathbb{I}(f(x) \neq y), \quad (14)$$

the loss of the classifier  $f(\cdot)$

label of this instance, which means it makes mistakes on the instance

model prediction on the raw data

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

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

model prediction on the raw data

the discriminative risk of  $f(\cdot)$

model prediction when only sensitive attribute(s) are changed

<sup>14</sup>Masegosa et al., see n. 11.