

# *Improving Fairness with Ensemble Combination: Margin-Dependent Bounds*<sup>1</sup>

Yijun Bian

Department of Computer Science  
University of Copenhagen

26 June 2026



Funded by  
the European Union



<sup>1</sup> ACM FAccT '26, Montreal, QC, Canada

# Example<sup>2</sup>

## Resume A

**Emily Walsh**

White-sounding name

Education: BA, 4 years experience

Skills: administrative, computer

Honours, volunteer work



## Resume B — identical except

**Lakisha Washington**

African-American-sounding name

Education: BA, 4 years experience

Skills: administrative, computer

Honours, volunteer work

### Motivation

If perturbing a protected attribute changes the prediction  
⇒ **discriminative risk exists**

<sup>2</sup>White-associated names received 50% more callbacks for interviews (Marianne Bertrand and Sendhil Mullainathan. "Are Emily and Greg more employable than Lakisha and Jamal? A field experiment on labor market discrimination". In: *Am Econ Rev* 94.4 [2004], pp. 991–1013)

# Discriminative risk (DR)

—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 protected/sensitive attributes*), with an instance denoted by  $\mathbf{x} = (\check{\mathbf{x}}, \mathbf{a})$ , the fairness quality of one hypothesis<sup>3</sup>  $f(\cdot)$  could be evaluated by

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

the indicator function  
 non-sensitive attributes (they may or may not include proxy attributes)  
 sensitive attribute(s)  
 sensitive attribute(s) that are slightly perturbed (the privileged  $\leftrightarrow$  any one of the unprivileged)

similarly to the 0/1 loss, where  $\mathbf{a} = [a_1, \dots, a_{n_a}]^T$ ,  $a_i \in \mathcal{A}_i$ ,  $n_a \geq 1$ ,  $|\mathcal{A}_i| \geq 2$ , and  $\tilde{\mathbf{a}}$  is a perturbed<sup>4</sup>  $\mathbf{a}$ . Note that Eq. (1) is evaluated on only one instance with sensitive attributes (SAs)  $\mathbf{x}$ .

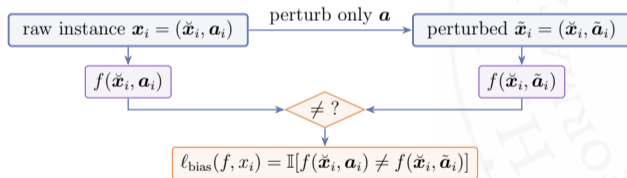
<sup>3</sup>The hypothesis used in this equation could indicate an individual classifier or an ensemble classifier.

<sup>4</sup>A 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.

# Discriminative risk (DR)

—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 (DR)** 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 DR<sup>3</sup>** 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.



$$S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

$$\underbrace{(0 \ 1 \ 0 \ 0 \ 1 \ 0 \ 1)}_{\hat{L}_{\text{bias}}(f, S) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{bias}}(f, \mathbf{x}_i)}$$

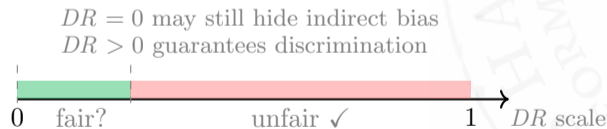
$$\boxed{L_{\text{bias}}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\ell_{\text{bias}}(f, \mathbf{x})]}$$

over population distribution

# Discriminative risk (DR)<sup>4</sup>

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

To describe this characteristic of the hypothesis on multiple instances, then the **empirical DR on one dataset**  $S$  is expressed as  $\hat{L}_{\text{bias}}(f, S) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{bias}}(f, x_i)$ , and the **true DR<sup>3</sup>** of the hypothesis over a data distribution is  $L_{\text{bias}}(f) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell_{\text{bias}}(f, x)]$ , respectively.



## Key caveat

$DR = 0$  is *necessary but not sufficient* for fairness.

When  $DR > 0$ , discrimination is *guaranteed* to exist.

<sup>3</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  $\mathcal{D}$ .

<sup>4</sup> $\hat{L}_{\text{bias}}(f, S)$  is an *unbiased estimator* of  $L_{\text{bias}}(f)$

# How *DR* relates to existing fairness measures

	Group fairness <sup>5</sup>	Individual fairness <sup>6</sup>	Causal fairness <sup>7</sup>	<i>DR</i>
Subgroup partition	required	—	—	<i>not required</i>
Similarity metric	—	required	—	<i>not required</i>
Causal graph	—	—	required	<i>not required</i>
Multi-valued SA	✗	✗	✓	✓
Original data only	Yes	Yes	No	No

*DR* follows the individual-fairness principle, aggregates like group fairness, and probes sensitivity like CFF — without requiring a specific metric, a subgroup partition, or a causal model.

<sup>5</sup>e.g. DP, EO, EOpp, and PP. (Empirically, *DR* has 0.58 Pearson correlation with accuracy change under perturbation, versus near-zero for DP and PP.)

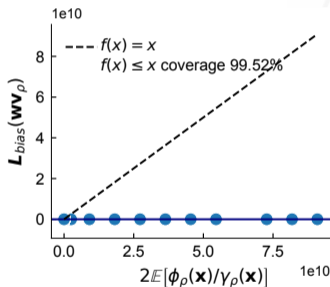
<sup>6</sup>e.g. Lipschitz condition

<sup>7</sup>e.g. CFF, proxy discrimination. (*DR* reflects counterfactual fairness (CFF) in a quantifiable way, without requiring a causal graph.)

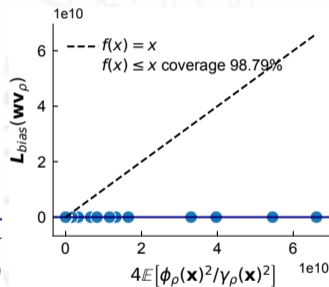
# Oracle bounds: can ensemble combination reduce DR?

1st order:  $L_{\text{bias}}(\mathbf{wv}_\rho) \leq 2 \mathbb{E}_{\mathcal{D}}[\phi_\rho(x)/\gamma_\rho(x)]$

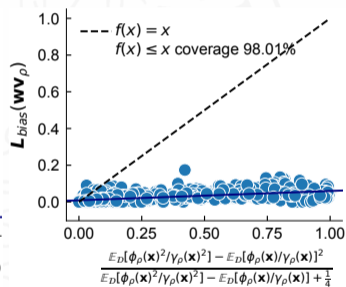
2nd order:  $L_{\text{bias}}(\mathbf{wv}_\rho) \leq 4 \mathbb{E}_{\mathcal{D}}[\phi_\rho(x)^2/\gamma_\rho(x)^2]$



(a)



(b)



(c)

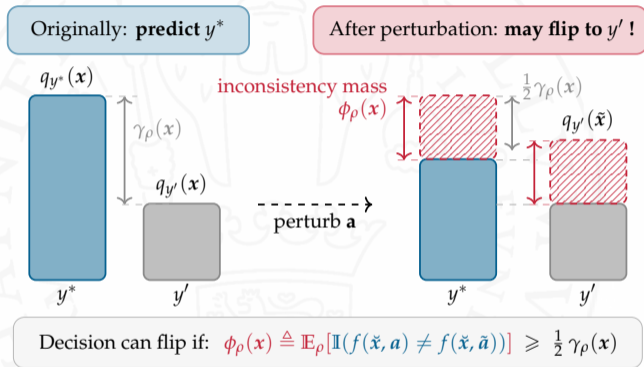
Bounds hold >98%, 2nd tighter, C-tandem tighter still  
Points below the identity line = bound holds

# What the bounds tell us

Ensembles of  $m$  classifier with weights  $\rho$   
Top-1 voting margin

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

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



DR of the ensemble  $\leq 2 \times$  (individual DR / voting margin)

Larger margin  $\rightarrow$  tighter bound  $\rightarrow$  fairer ensemble

A *cancellation-of-biases* effect exists — but it depends on margins, unlike cancellation of errors

# Summary<sup>8</sup>

## What we proposed

- *DR* captures both individual and group fairness via SA perturbation (no metric, no causal graph, no subgroup partition needed), and applies naturally to multi-valued & multiple SAs
- *Oracle bounds* indicate that a cancellation-of-biases effect exists in combination, margin-dependent
- *POAF*: ensemble pruning via Pareto optimality, improving fairness with acceptable accuracy loss

## What this means

- Fairness *can* be improved via ensembling with theoretical guarantees, not just empirical hope
- Larger voting margin  $\Rightarrow$  tighter fairness bounds: boosting performance helps fairness too
- *POAF* saves effort on repetitive hyperparameter tuning

## Limitations & future work

- *DR* involves randomness from SA perturbation
- *POAF* has higher time cost than simpler baselines
- Acceleration of *POAF* and stability of *DR* worth exploring

---

<sup>8</sup>Code: <https://github.com/eustomaqua/FairML>