



Fairness for the People, by the People: Minority Collective Action

Omri Ben-Dov¹, Samira Samadi¹, Amartya Sanyal², Alexandru Țifrea³

¹Max Planck Institute for Intelligent Systems, Tübingen AI Center

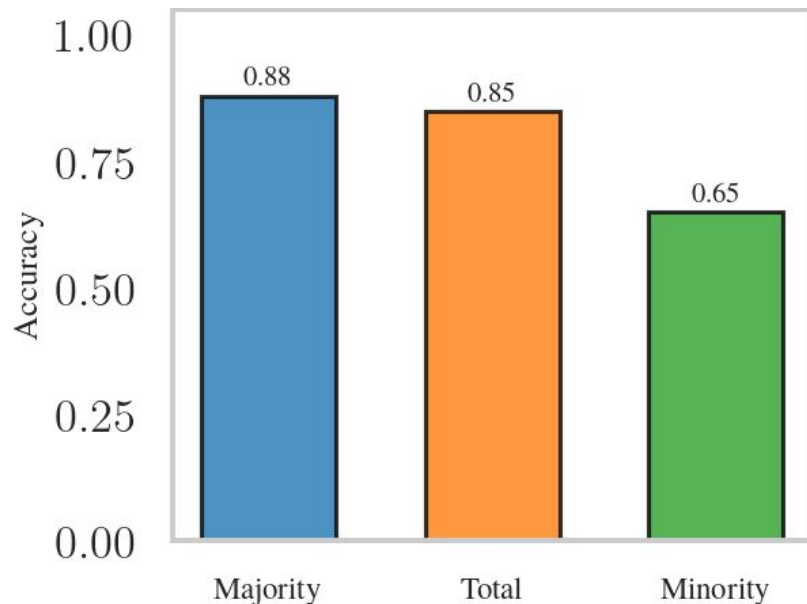
²Department of Computer Science, University of Copenhagen

³ETH Zurich

Fairness in machine learning



Features x_i , labels y_i



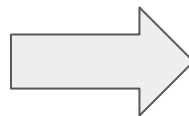
Group fairness
similar outcomes to
different groups

Solutions for fairness



Change the learning algorithm:

- Edit the data
- Modify the training objective
- Post-process the results
- ...



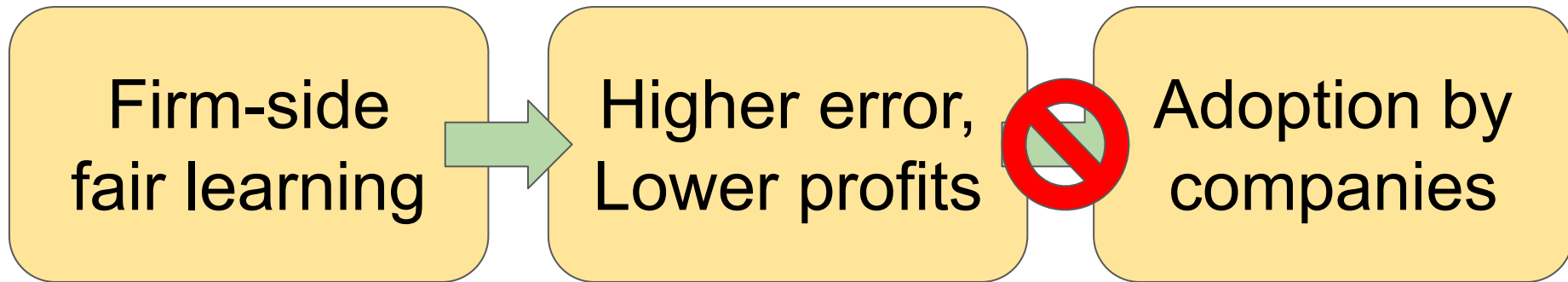
Fairness improves



Lower predictive accuracy

Menon, A.K. and Williamson, R.C., 2018, January. The cost of fairness in binary classification. In *Conference on Fairness, accountability and transparency* (pp. 107-118). PMLR.

Solutions for fairness



When a company learns from user-data,
can a minority induce fairness?

Algorithmic Collective Action

Hardt, M., Mazumdar, E., Mendler-Dünner, C. and Zrnic, T., 2023, July. Algorithmic collective action in machine learning. In *International Conference on Machine Learning* (pp. 12570-12586). PMLR.



Assume a classifier trained on data from end-users:

- One person will have a negligible effect on the classifier.
- By collaborating, people can influence the classifier.
- How many and to what goal?

Algorithmic Collective Action

Hardt, M., Mazumdar, E., Mendler-Dünner, C. and Zrnic, T., 2023, July. Algorithmic collective action in machine learning. In *International Conference on Machine Learning* (pp. 12570-12586). PMLR.



arXiv

Goal: Make the classifier ignores a signal g

Success:
$$S(\alpha) = \mathbb{P}_0 [h(g(x)) = h(x)]$$

Collective size Signal to "ignore" Classifier

Action:
$$x, y \rightarrow x, \operatorname{argmax}_{y' \in \{0,1\}} \mathbb{P}_0 (y' | g(x))$$

Change labels

Algorithmic Collective Action

Hardt, M., Mazumdar, E., Mendler-Dünner, C. and Zrnic, T., 2023, July. Algorithmic collective action in machine learning. In *International Conference on Machine Learning* (pp. 12570-12586). PMLR.



arXiv

Goal: Make the classifier ignores a signal g

Success:
$$S(\alpha) = \mathbb{P}_0 [h(g(x)) = h(x)]$$

Collective size Signal to “ignore” Classifier

Action: $x, y \rightarrow x, \operatorname{argmax}_{y' \in \{0,1\}} \mathbb{P}_0(y'|g(x))$
Change labels

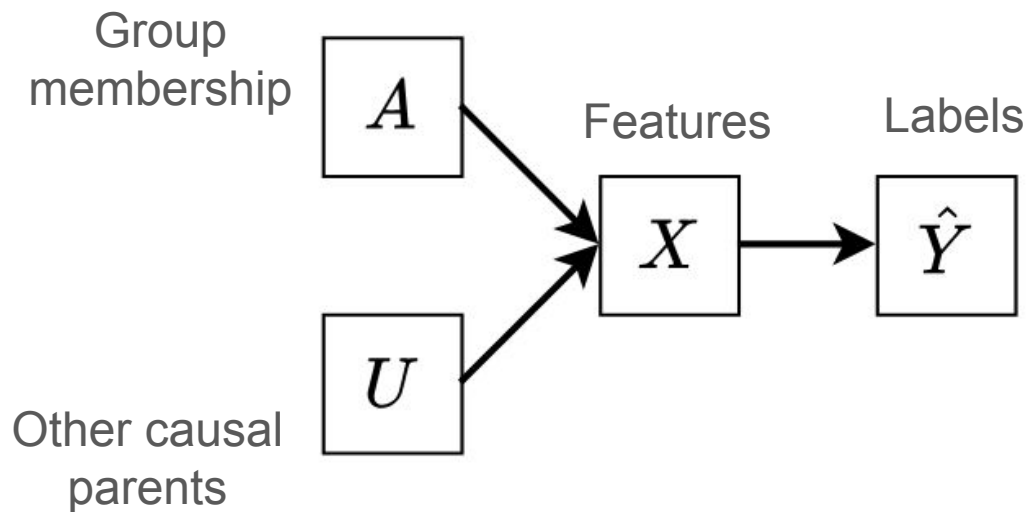
If g contains group membership,
the classifier will “ignore” group bias.



arXiv

Counterfactual Fairness

Assume a causal model



Counterfactual signal

$$g(x) = x_{A \leftarrow \text{Majority}}$$



Minority Collective Action

Cannot accurately modify all labels:

- We do not know the counterfactuals
- Only the minority may participate

$$x, y \rightarrow x, \operatorname{argmax}_{y' \in \{0,1\}} \mathbb{P}_0(y' | g(x))$$

Method:

1. Find “negative” minority members
2. Estimate the likelihood of a counterfactual “positive” label



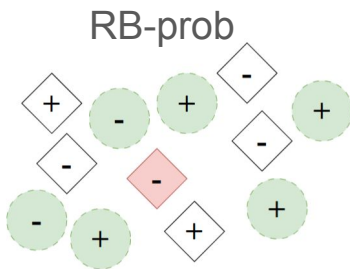
Minority Collective Action

Cannot accurately modify all labels:

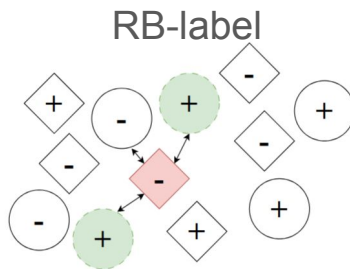
- We do not know the counterfactuals
- Only the minority may participate

Method:

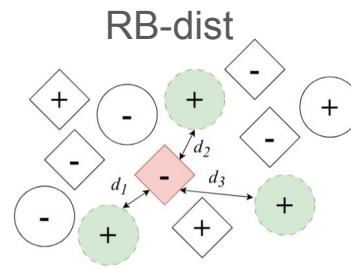
1. Find “negative” minority members
2. Estimate the likelihood of a counterfactual “positive” label



$$s_i = f(x_i)$$



$$s_i = \sum_{j \in K_i} \mathbf{1}\{y_j = 1\}$$



$$s_i = -\frac{1}{k} \sum_{j \in K_i} \|x_i - x_j\|_2$$

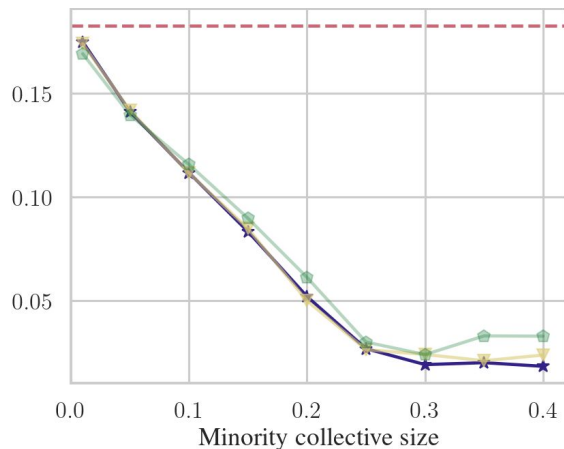
3. Top candidates flip their labels



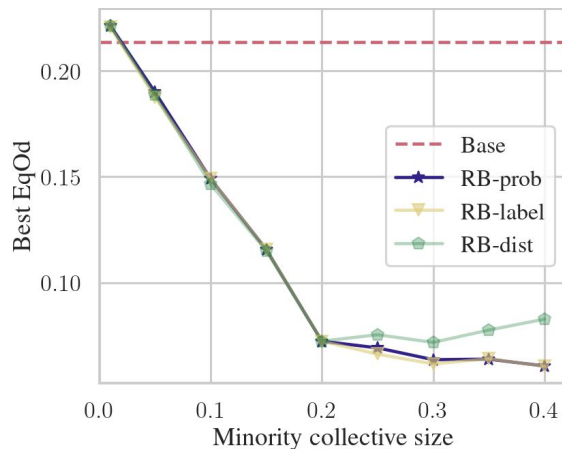
arXiv

Minority Collective Size

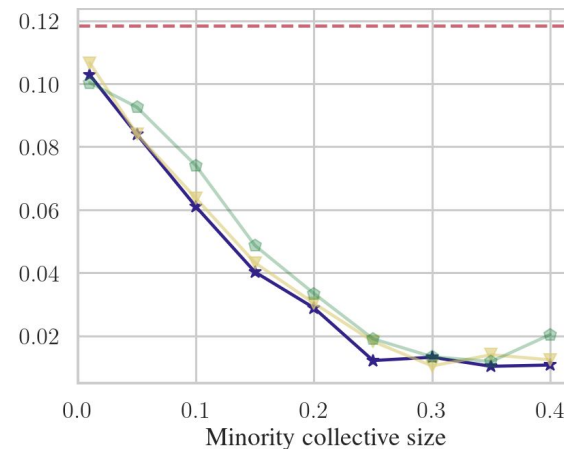
ACS-Income



COMPAS



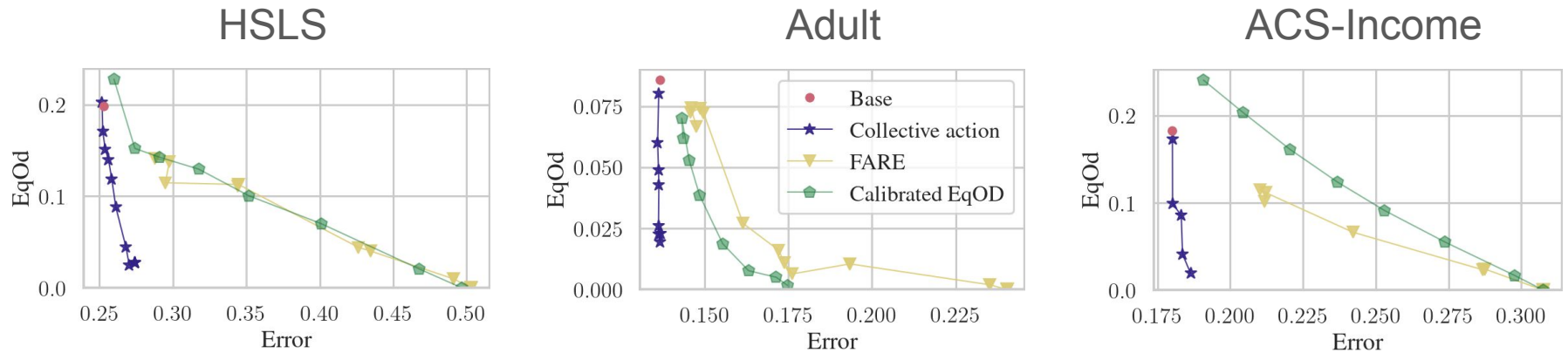
Waterbirds



20–30% of the minority attains the least fairness violation.



Compare With Firm-Side Methods



Unlike firm-side FARE and calibrated equalized odds, a minority cannot get perfect fairness, but adds smaller error.



Conclusion

1. Collective action framework for fairness.
2. Fairness improves linearly with collective size until around 30%.
3. Generally, a minority cannot achieve perfect fairness.
4. Smaller error than firm-side fairness methods.

More details and experiments in the preprint <https://arxiv.org/abs/2508.15374>