

# Improving Fairness Generalization Through a Sample-Robust Optimization Method

Julien Ferry<sup>1</sup>, Ulrich Aïvodji<sup>2</sup>, Sébastien Gambs<sup>3</sup>, Marie-José Huguet<sup>1</sup> and Mohamed Siala<sup>1</sup>

<sup>1</sup>LAAS-CNRS, Université de Toulouse, CNRS, INSA, Toulouse, France

<sup>2</sup>École de Technologie Supérieure, Montréal, Canada

<sup>3</sup>Université du Québec à Montréal, Montréal, Canada

*February 7, 2023*



- The paper is published in the **Machine Learning** journal, in July 2022
- Link to the paper [link.springer.com/article/10.1007/s10994-022-06191-y](https://link.springer.com/article/10.1007/s10994-022-06191-y)
- Preprint: <https://hal.archives-ouvertes.fr/hal-03709547>
- Open source code <https://github.com/ferryjul/FairnessSampleRobustness>

- Fairness in machine learning
- Statistical measures
- Generalisation of fairness on unseen data is one of the open challenges for trustworthy machine learning

## The COMPAS Example [Angwin et al., 2016]

- Binary classification task: Recidivism within two years
- Sensitive attribute: Ethnicity (African-American/Caucasian)
- Protected Groups:
  - ▶  $\mathcal{A}$  : African-American individuals;
  - ▶  $\mathcal{B}$  : Caucasian individuals;

## Statistical Fairness

- Principle: ensure that some measure  $\mathcal{M}$  *differs by no more than  $\epsilon$*  between several *protected groups*
- In the particular case of two protected groups ( $\mathcal{A}$ ) and ( $\mathcal{B}$ ), one need to ensure that  $|\mathcal{M}(\mathcal{A}) - \mathcal{M}(\mathcal{B})| < \epsilon$

## Supervised Fair Learning: A Bi-Objective Optimization Problem

- Notation:  $\mathcal{D}$  initial dataset,  $h$  prediction model,  $\epsilon$  unfairness tolerance
- Let  $\text{unf}(\cdot)$  be an unfairness oracle. A common formulation of the Fair Learning problem is:

$$\begin{aligned} \arg \min_{h \in \mathcal{H}} \quad & f_{obj}(h, \mathcal{D}) \\ \text{s.t.} \quad & \text{unf}(h, \mathcal{D}) \leq \epsilon \end{aligned} \tag{1}$$

where one wants to build model  $h$  minimizing objective function  $f_{obj}$  and exhibiting unfairness withing an  $\epsilon$  threshold (on training dataset  $\mathcal{D}$ )

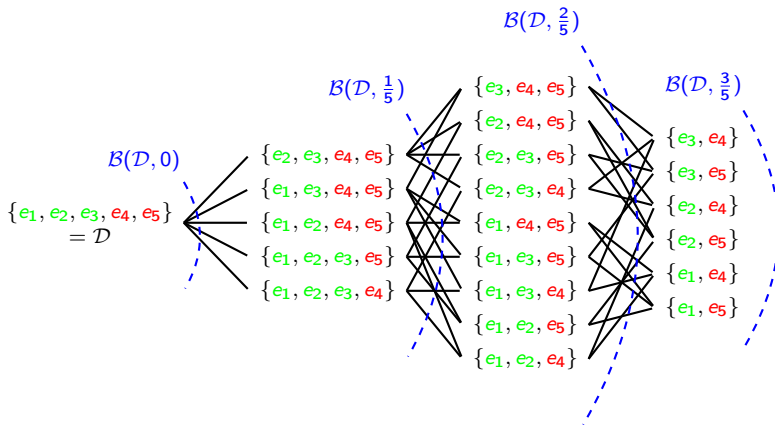
- The fairness constraint does not generalize well in practice [Cotter et al., 2018, 2019]

## Distributionally Robust Optimization (DRO)

- Instead of minimizing objective function  $f_{obj}$  for a given distribution  $\mathcal{P}$ , DRO aims at minimizing  $f_{obj}$  for a worst-case distribution among a *perturbation set* of  $\mathcal{P}$  [Sagawa et al., 2019]  $\mathcal{B}(\mathcal{P})$

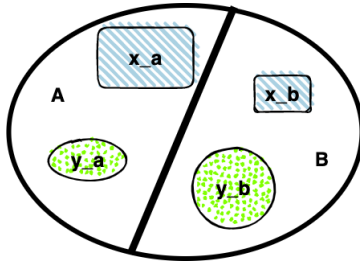
## Perturbation Set based on the Jaccard Distance

- Let  $\mathcal{D}_1$  and  $\mathcal{D}_2$  be two sample sets. The Jaccard distance between  $\mathcal{D}_1$  and  $\mathcal{D}_2$  is defined as follows:  $J_\delta(\mathcal{D}_1, \mathcal{D}_2) = 1 - \frac{|\mathcal{D}_1 \cap \mathcal{D}_2|}{|\mathcal{D}_1 \cup \mathcal{D}_2|}$



**Figure:** Example of perturbation sets for a dataset  $\mathcal{D}$  with 5 examples and two protected groups  $a$  ( $\{e_1, e_2, e_3\}$ ) and  $b$  ( $\{e_4, e_5\}$ ). Subsets that can not be used to audit a model's fairness with respect to protected groups  $a$  and  $b$  are not represented.





- Let  $\mathcal{D}'$  be the subset of  $\mathcal{D}$  where the colored sets are removed
- If  $h$  is not fair on  $\mathcal{D}'$ , then  $h$  is not robust on the perturbation set defined with the distance between  $\mathcal{D}$  and  $\mathcal{D}'$

$$\min \quad n \quad (2)$$

$$\text{s.t.} \quad n = x_a + x_b + y_a + y_b \quad (3)$$

$$\left| \frac{M_a - x_a}{N_a - x_a - y_a} - \frac{M_b - x_b}{N_b - x_b - y_b} \right| > \epsilon \quad (4)$$

$$0 \leq x_a \leq M_a$$

$$0 \leq x_b \leq M_b$$

$$0 \leq y_a \leq M_a - M_a$$

$$0 \leq y_b \leq N_b - M_b$$

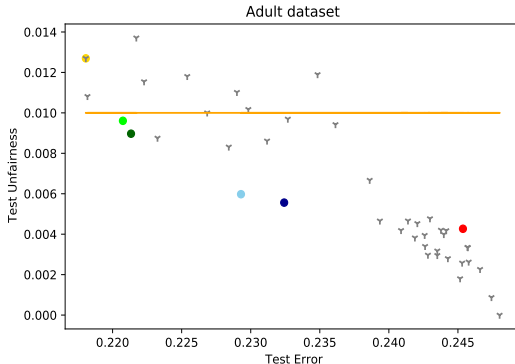
$$x_a + y_a < N_a$$

$$x_b + y_b < N_b$$

- The optimal value can be used to identified the largest perturbation set where the robustness constraint is satisfied

## Setup description

- We compare:
  - ▶ FairCORELS [Aïvodji et al., 2019]
  - ▶ TFCO: TensorFlow Constrained Optimization
  - ▶ Exact and Heuristic approaches
- Four fairness metrics:
  - ▶ Statistical Parity [Dwork et al., 2012]
  - ▶ Predictive Equality [Chouldechova, 2017]
  - ▶ Equal Opportunity [Hardt et al., 2016]
  - ▶ Equalized Odds [Hardt et al., 2016]
- Four biased datasets:
  - ▶ Adult Income dataset [Frank and Asuncion, 2010]
  - ▶ COMPAS dataset [Angwin et al., 2016]
  - ▶ Default Credit dataset [Yeh and Lien, 2009]
  - ▶ Bank Marketing dataset [Moro et al., 2014]
- The Integer Program is solved using the constraint programming solver OrTools



**Figure:** Test error and unfairness of models generated by FairCORELS using our exact and heuristic sample-robust fair methods (Statistical Parity metric,  $\epsilon = 0.01$ )

## Summary

- We address the problem of fairness generalisation via an approach based on Distributionally Robust Optimization
- We empirically show that our approach is competitive to the state-of-the-art with two learning models on many datasets in the literature using different fairness measures

## Future Work

- Can we find an efficient way to approximate the best parameters ?
- How to extend the work with other distance functions ?

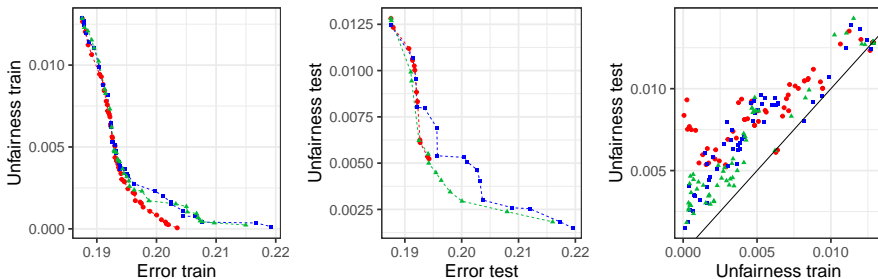
Thank you!

- Aïvodji, U., Ferry, J., Gambs, S., Huguet, M.-J., and Siala, M. (2019). Learning fair rule lists. *arXiv preprint arXiv:1909.03977*.
- Angelino, E., Larus-Stone, N., Alabi, D., Seltzer, M., and Rudin, C. (2017a). Learning certifiably optimal rule lists. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 35–44. ACM.
- Angelino, E., Larus-Stone, N., Alabi, D., Seltzer, M., and Rudin, C. (2017b). Learning certifiably optimal rule lists for categorical data. *The Journal of Machine Learning Research*, 18(1):8753–8830.
- Angwin, J., Larson, J., Mattu, S., and Kirchner, L. (2016). Machine bias. *ProPublica*, May, 23.
- Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163.
- Cotter, A., Gupta, M., Jiang, H., Srebro, N., Sridharan, K., Wang, S., Woodworth, B., and You, S. (2018). Training fairness-constrained classifiers to generalize.
- Cotter, A., Gupta, M., Jiang, H., Srebro, N., Sridharan, K., Wang, S., Woodworth, B., and You, S. (2019). Training well-generalizing classifiers for fairness metrics and other data-dependent constraints. In *International Conference on Machine Learning*, pages 1397–1405. PMLR.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. (2012). Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pages 214–226.

- Frank, A. and Asuncion, A. (2010). UCI machine learning repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California. *School of information and computer science*, 213:2–2.
- Hardt, M., Price, E., and Srebro, N. (2016). Equality of opportunity in supervised learning. *arXiv preprint arXiv:1610.02413*.
- Moro, S., Cortez, P., and Rita, P. (2014). A data-driven approach to predict the success of bank telemarketing. *Decision Support Systems*, 62:22–31.
- Rivest, R. L. (1987). Learning decision lists. *Machine learning*, 2(3):229–246.
- Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P. (2019). Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*.
- Yeh, I.-C. and Lien, C.-h. (2009). The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. *Expert Systems with Applications*, 36(2):2473–2480.

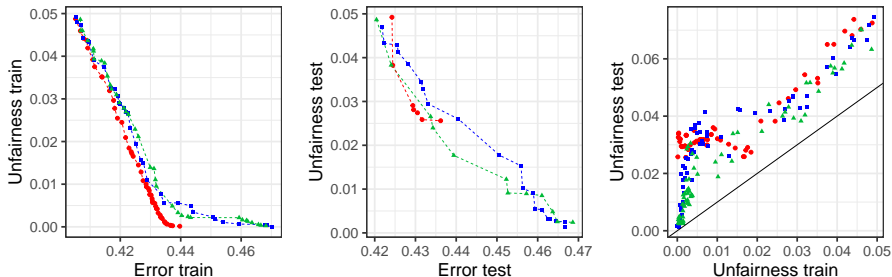


Strategy ■ 10 masks ▲ 30 masks ● no mask



**Figure:** Results obtained on the Default Credit dataset, for the Predictive Equality metric

Strategy ■ 10 masks ▲ 30 masks ● no mask



**Figure:** Results obtained on the COMPAS dataset, for the Statistical Parity metric

## Rule Lists: Definition

*Rule lists* [Rivest, 1987] are classifiers formed by an ordered list of *if-then* rules with antecedents in the *if* clauses and predictions in the *then* clauses. More precisely, a rule list  $r = (\{p_{k,k \in \{1..K\}}\}, \{q_{k,k \in \{1..K\}}\}, q_0)$  consists of  $K$  distinct association rules  $p_k \rightarrow q_k$ , in which  $p_k$  is the antecedent of the association rule and  $q_k$  its associated consequent, followed by a default prediction  $q_0$ .

**A possible rule list for the example dataset of slide ?? (with 100% accuracy)**

---

```
if [Education: Dropout] then [low]
else if [Gender: Male AND Age > 45] then [high]
else [low]
```

---

## FairCORELS Problem Formulation

- Based on the CORELS algorithm [Angelino et al., 2017a,b]
- FairCORELS [Aïvodji et al., 2019] returns rule list  $r^*$  that is a solution to the following problem:

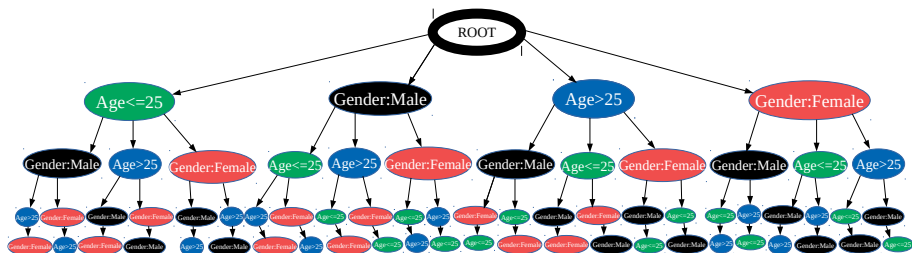
$$\begin{aligned} \arg \min_{r \in \mathcal{R}} \quad & \text{misc}(h, \mathcal{D}) + \lambda \cdot K_r \\ \text{s.t.} \quad & \text{unf}(h, \mathcal{D}) \leq \epsilon \end{aligned}$$

where:

- ▶  $\mathcal{R}$  is the space of rule lists
- ▶  $\mathcal{D}$  denotes the training dataset
- ▶  $K_r$  is the length of rule list  $r$
- ▶  $\lambda$  is a regularization parameter balancing sparsity and accuracy
- ▶  $\text{misc}(\cdot)$  is the misclassification error and  $\text{unf}(\cdot)$  measures unfairness

## FairCORELS search space

- FairCORELS represents the search space of rule lists as a prefix tree (trie)
- FairCORELS leverages several bounds to efficiently explore this search space (including CORELS' original bounds)



**Figure:** Example prefix tree with 4 attributes