

# Learning Fair Classifiers in Online Stochastic Settings



Yi Sun<sup>1</sup> Ivan Ramirez<sup>1</sup> Alfredo Cuesta-Infante<sup>2</sup> Kalyan Veeramachaneni<sup>1</sup>

<sup>1</sup>Massachusetts Institute of Technology <sup>2</sup>Universidad Rey Juan Carlos

## Introduction and Motivation

### What makes policy-driven machine learning different?

- Data might not be available upfront → Online Algorithm
- Human decision makers in the loop → Meta Algorithm
- Fair and Accurate

## Online Binary Classification With Fairness

**Given:** A set of experts  $f \in \mathcal{F}$ , where  $f : (\mathcal{X}, \mathcal{Z}) \rightarrow \{0, 1\}$

**At each round:**

- An individual arrives with sensitive attributes  $z$ , and non-sensitive attributes  $x$
- Sample an expert and use its prediction
- Observe true label and update weights on experts

**Goal:**

1. Regret

$$\sum_{t=1}^T \ell(f^t(x^t, z^t), y^t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f(x^t, z^t), y^t)$$

2. Equalized Odds [2]

$$|\mathbb{E}[\hat{Y} = 1 | Y = 1, Z = A] - \mathbb{E}[\hat{Y} = 1 | Y = 1, Z = B]| \leq \epsilon$$

## Methodology

### Key Ideas

- Running separate instances of Multiplicative Weights algorithm for each group and label combination
- Randomize between instances help with fairness
- Obtain optimal selection probability between instances by optimizing regret and fairness bound

### Algorithm

Initialize  $w_{f,z,k}^1 = 1 \ \forall f, z, k$ , and  $q_{z,k}^1 = \frac{1}{2} \ \forall z, k$   
**for**  $t \leftarrow 1, \dots, T$  **do**  
    Each classifier obtains  $\hat{y}_f^t$   
    Obtain the optimal  $q^*$   
    
$$\pi^t(f|z^t) = \begin{cases} \frac{w_{f,z^t,+}^t}{\sum_{f \in \mathcal{F}} w_{f,z^t,+}^t} & \text{with probability } q_{z^t,+}^* \\ \frac{w_{f,z^t,-}^t}{\sum_{f \in \mathcal{F}} w_{f,z^t,-}^t} & \text{with probability } q_{z^t,-}^* \end{cases}$$
  
    Select classifier  $f$  according to probability  $\pi^t$  and update the regret  
    Obtain loss  $\ell_f^t = \ell(\hat{y}_f^t, y^t)$  for each classifier  $f$   
    Update weights  $w_{f,z,k}^{t+1} = w_{f,z,k}^t (1 - \eta)^{\ell_f^t 1\{z^t=z\} 1\{y^t=k\}} \ \forall f, z, k$   
**end**

**Algorithm 1:** Fairness-Aware MW algorithm

Figure 1: Fairness-aware RMW algorithm

#### Theorem 0.1: Upper Bound On Regret

Let  $\alpha_{z,-}^t = \sum_{f \in \mathcal{F}} \frac{w_{f,z,-}^t}{\sum_{f \in \mathcal{F}} w_{f,z,-}^t} \cdot \ell_{f,z,+}^t - \sum_{f \in \mathcal{F}} \frac{w_{f,z,+}^t}{\sum_{f \in \mathcal{F}} w_{f,z,+}^t} \cdot \ell_{f,z,+}^t$  and  $\alpha_{z,-}^t$  defined similarly. Thus the expected total loss of the algorithm is:

$$\mathbb{E}[L] \leq (1 + \eta)L_f + 4 \frac{\ln d}{\eta} + \alpha \quad (1)$$

where  $\alpha = \sum_{z \in \{A,B\}, y \in \{+,-\}} q_{z,y} \sum_t \alpha_{z,y}^t$

#### Theorem 0.2: Fairness Bound

In the stochastic setting, there exists  $q_{A,-}$  and  $q_{B,-}$  such that the absolute difference in FPR can be bounded as:

$$\mathbb{E}_{x,y,z} \left[ \frac{\mathbb{E}[L_{A,-}]}{C_{A,-}} - \frac{\mathbb{E}[L_{B,-}]}{C_{B,-}} \right] \leq (1 + \eta - \gamma(\eta)) \mathbb{E}_{x,y,z} \left[ \frac{L_{f^*(B,-),B,-}}{C_{B,-}} \right] + \epsilon(1 + \eta) + \left( \frac{q_{A,-} \cdot \sum_t \alpha_{A,-}^t}{p \cdot (1 - \mu_{A,+}) \cdot T} - \frac{q_{B,-} \cdot \sum_t \alpha_{B,-}^t}{(1 - p) \cdot (1 - \mu_{B,+}) \cdot T} \right) \quad (2)$$

## Optimal Balance Between Regret and Fairness

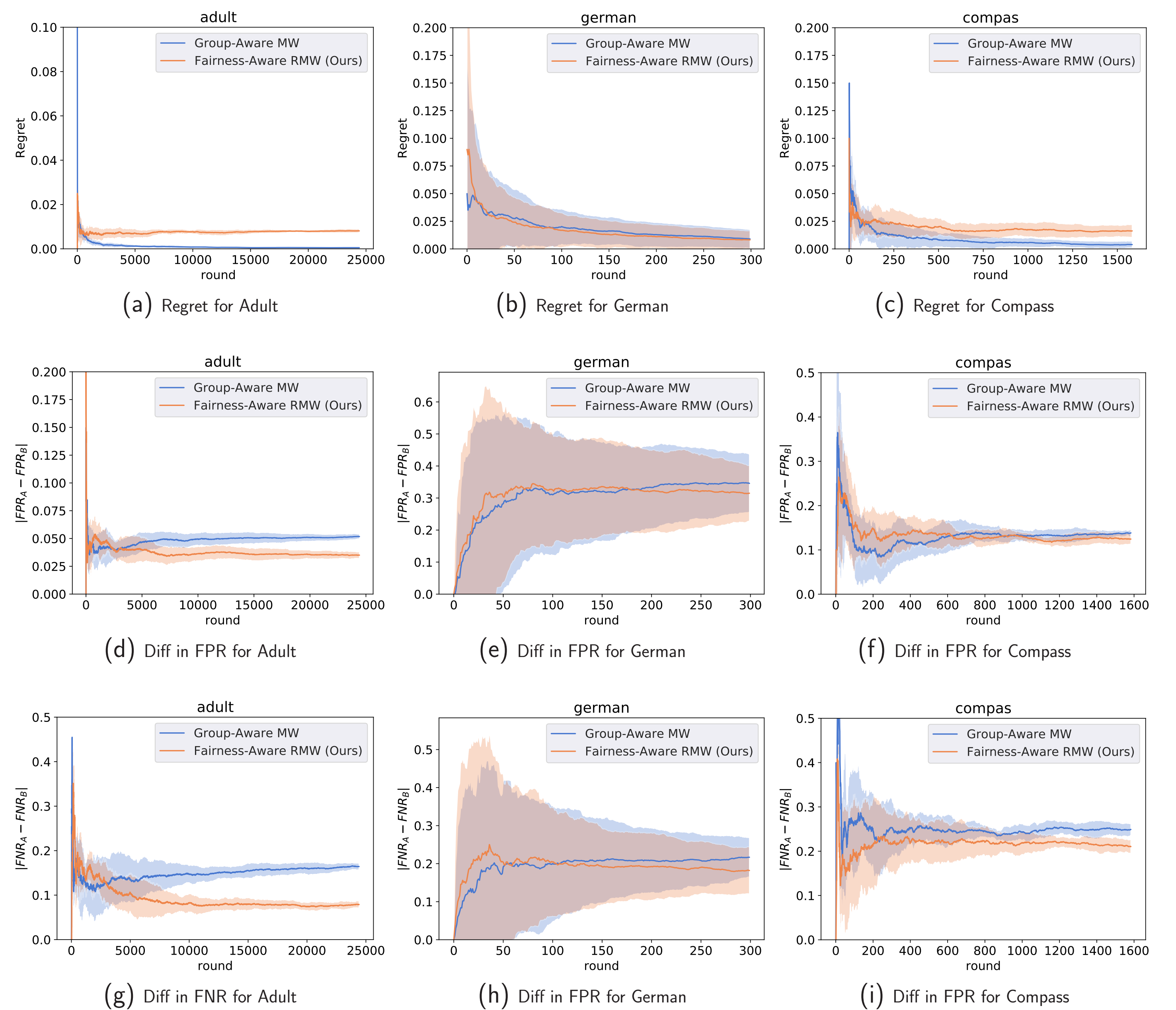
At each round, we solve the following optimization problem to minimize the fairness and regret upper bound:

$$\mathbf{q}^* = \arg \min_{\mathbf{q}} ||\lambda(\mathbf{A}\mathbf{q} - \mathbf{b})||^2 \quad (3)$$

where  $\lambda$  is a vector of balancing the importance of equalized FPR, equalized FNR and regret that can be provided on a case-by-case basis based on different potential applications.

## Experiments

The set of classifiers  $\mathcal{F}$  in our hypothesis sets are as follows: Logistic Regression (LR), Linear SVM (L SVM), RBF SVM, Decision Tree (DT), Multi-Layer Perceptron (MLP). We pre-trained each classifier for each trial by splitting the data set, with 70% for training and 30% for testing. During the simulations, the examples in the testing set arrived one by one. We compare with [1], which achieves equalized error rates by running separate instance of MW algorithm for each sensitive group.



## Conclusion

- Improvement in fairness both in terms of equalized FPR and FNR, along with a small increase in regret
- Randomization help overcome biases of experts

## References

- [1] A. Blum, S. Gunasekar, T. Lykouris, and N. Srebro. On preserving non-discrimination when combining expert advice. In *NeurIPS*, 2018.
- [2] P. E. S. N. e. a. Hardt, Moritz. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems*, pp. 3315–3323, 2016.