# Preventing Fairness Gerrymandering: Auditing and Learning for Subgroup Fairness

ICML 2018

Preventing Fairness Gerrymandering: Auditing and Learning for Subgroup Fairness

# Metric of fairness



PART 1

Fairness

# Defining of Subgroup

1. group: fairness gerrymandering problem



**2.** Need of Sub group : Fairness Violation evaluation  $\rightarrow$  Auditing (subgroup searching)

### Auditing and Learning for subgroup Fairness

# (Learner-Auditor) Zero-sum game formulation



# Appendix

appendix

Experimental setting and result

Preventing Fairness Gerrymandering: Auditing and Learning for Subgroup Fairness

# **Experimental result**



Figure 1. Evolution of the error and unfairness of Learner's classifier across iterations, for varying choices of *τ*.
(a) Error "t of Learner's model vs iteration t. (b) Unfairness t of subgroup found by Auditor vs. iteration t, as measured by Definition 2.3. See text for details.

appendix

#### Experimental setting and result

Preventing Fairness Gerrymandering: Auditing and Learning for Subgroup Fairness

# **Experimental result**



Figure 2. (a) Pareto-optimal error-unfairness values, color coded by varying values of the input parameter *τ*.
(b) Aggregate Pareto frontier across all values of *τ*.
Here the *τ* values cover the same range but are sampled more densely to get a smoother frontier. See text for details.

#### Mathematical theorem

# Auditing ~ Weak Agnostic Learning reduction theorem

#### **Proving Computational equivalence**

**Theorem 3.1.** Fix any distribution  $\mathcal{P}$ , and any set of group indicators  $\mathcal{G}$ . Then for any  $\gamma, \varepsilon > 0$ , the following relationships hold:

- If there is a  $(\gamma/2, (\gamma/2 \varepsilon))$  auditing algorithm for  $\mathcal{G}$ for all D such that SP(D) = 1/2, then the class  $\mathcal{G}$  is  $(\gamma, \gamma/2 - \varepsilon)$ -weakly agnostically learnable under  $\mathcal{P}^D$ .
- If  $\mathcal{G}$  is  $(\gamma, \gamma \varepsilon)$ -weakly agnostically learnable under marginal distribution  $\mathcal{P}^D$  on (x, D(X)) for all D such that SP(D) = 1/2, then there is a  $(\gamma, (\gamma - \varepsilon)/2)$ auditing algorithm for  $\mathcal{G}$  for SP fairness under  $\mathcal{P}$ .

**Corollary 3.2.** *Fix any distribution*  $\mathcal{P}$ *, and any set of group indicators*  $\mathcal{G}$ *. The following two relationships hold:* 

- If there is a  $(\gamma/2, (\gamma/2 \varepsilon))$  auditing algorithm for  $\mathcal{G}$ for all D with FP(D) = 1/2, then  $\mathcal{G}$  is  $(\gamma, \gamma/2 - \varepsilon)$ weakly agnostically learnable under  $\mathcal{P}_{y=0}^{D}$ .
- If G is (γ, γ − ε)-weakly agnostically learnable under the conditional distribution P<sup>D</sup><sub>y=0</sub> of (X, y) conditioned on the event that D(X) = 1 for all D with FP(D) = 1/2, then there is a (γ, (γ − ε)/2) auditing algorithm for FP subgroup fairness for G under distribution P.

#### Mathematical theorem

**Theorem 3.3.** Under standard complexity-theoretic intractability assumptions, for  $\mathcal{G}$  the classes of conjunctions of boolean attributes, linear threshold functions, or boundeddegree polynomial threshold functions, there exist distributions P such that the auditing problem cannot be solved in polynomial time, for either SP or FP fairness. **Theorem 4.1.** Fix any  $\nu, \delta \in (0, 1)$ . Then given an input of n data points and accuracy parameters  $\nu, \delta$  and access to oracles  $\operatorname{CSC}(\mathcal{H})$  and  $\operatorname{CSC}(\mathcal{G})$ , there exists an algorithm runs in time  $\operatorname{poly}(1/\nu, \log(1/\delta))$ , and with probability at least  $1 - \delta$ , output a randomized classifier  $\hat{D}$  such that  $\operatorname{err}(\hat{D}, \mathcal{P}) \leq \operatorname{OPT} + \nu$ , and for any  $g \in \mathcal{G}$ , the fairness constraint violations satisfies

 $\alpha_{FP}(g, \mathcal{P}) \ \beta_{FP}(g, \hat{D}, \mathcal{P}) \leq \gamma + O(\nu).$