



FairGBM

Gradient Boosting with Fairness Constraints

International Conference on Learning Representations
ICLR 2023

André Cruz

Catarina Belém

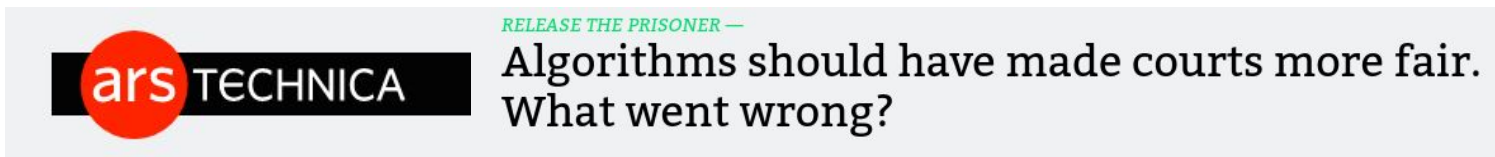
João Bravo

Pedro Saleiro

Pedro Bizarro



- ML has become ubiquitous in a multitude of high-stakes applications.
- At the same time, numerous reports have warned of AI bias in production systems.



*Biased Lending Evolves, and Blacks
Face Trouble Getting Mortgages*

The New York Times

**Amazon scraps secret AI recruiting tool that
showed bias against women**

- ML has become ubiquitous in a multitude of high-stakes applications.
- At the same time, numerous reports have warned of AI bias in production systems.

- **GBMs are still the state-of-the-art** on tabular data, often outperforming Deep Learning approaches.

- ML has become ubiquitous in a multitude of high-stakes applications.
- At the same time, numerous reports have warned of AI bias in production systems.
- **GBMs are still the state-of-the-art** on tabular data, often outperforming Deep Learning approaches.
- There are no (fairness-)constrained optimization methods tailored for GBM.

Goal: Training GBM to minimize loss under fairness constraints.

Goal: Training GBM to minimize loss under fairness constraints.

$$\theta^* = \arg \min_{\theta \in \Theta} \overbrace{l(\theta)}^{\text{predictive loss}} \quad \text{s. t.} \quad \overbrace{c_i(\theta) \leq 0}_{\text{constraints}} \quad i \in [m]$$

Goal: Training GBM to minimize loss under fairness constraints.

$$\theta^* = \arg \min_{\theta \in \Theta} l(\theta) \quad \text{s. t.} \quad c_i(\theta) \leq 0$$

predictive loss constraints

Approach

- The *standard* approach for constrained optimization is to use the method of Lagrange multipliers.

$$\mathcal{L}(\theta, \lambda) = l(\theta) + \sum_{i=1}^m \lambda_i c_i(\theta)$$

model parameters Lagrange multipliers

Goal: Training GBM to minimize loss under fairness constraints.

$$\theta^* = \arg \min_{\theta \in \Theta} \overbrace{l(\theta)}^{\text{predictive loss}} \quad \text{s. t.} \quad \overbrace{c_i(\theta) \leq 0}_{\text{constraints}} \quad i \in [m]$$

Approach

- The *standard* approach for constrained optimization is to use the method of Lagrange multipliers.

$$\mathcal{L}(\theta, \lambda) = l(\theta) + \sum_{i=1}^m \lambda_i c_i(\theta)$$


model parameters Lagrange multipliers

- Solution lies on saddle point of the Lagrangian function, $\mathcal{L}(\theta, \lambda)$.


Approach

- Finding the saddle-point of $\mathcal{L}(\theta, \lambda) = l(\theta) + \sum_{i=1}^m \lambda_i c_i(\theta)$, at $\nabla \mathcal{L}(\theta, \lambda) = \vec{0}$.

Approach

- Finding the saddle-point of $\mathcal{L}(\theta, \lambda) = l(\theta) + \sum_{i=1}^m \lambda_i c_i(\theta)$, at $\nabla \mathcal{L}(\theta, \lambda) = \vec{0}$.
- In practice, this is done via iterative gradient descent and ascent steps.


Approach

- Finding the saddle-point of $\mathcal{L}(\theta, \lambda) = l(\theta) + \sum_{i=1}^m \lambda_i c_i(\theta)$, at $\nabla \mathcal{L}(\theta, \lambda) = \vec{0}$.
- In practice, this is done via iterative gradient descent and ascent steps.


- Gradient **descent** over θ :
$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial l}{\partial \theta} + \sum_{i=1}^m \lambda_i \frac{\partial c_i}{\partial \theta}$$

- Gradient **ascent** over λ :
$$\frac{\partial \mathcal{L}}{\partial \lambda_i} = c_i(\theta)$$

Approach

- Finding the saddle-point of $\mathcal{L}(\theta, \lambda) = l(\theta) + \sum_{i=1}^m \lambda_i c_i(\theta)$, at $\nabla \mathcal{L}(\theta, \lambda) = \vec{0}$.
- In practice, this is done via iterative gradient descent and ascent steps.

- Gradient **descent** over θ : $\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial l}{\partial \theta} + \sum_{i=1}^m \lambda_i \frac{\partial c_i}{\partial \theta}$ we require differentiable constraints

- Gradient **ascent** over λ : $\frac{\partial \mathcal{L}}{\partial \lambda_i} = c_i(\theta)$ but *only* for the descent step!

Approach

- Finding the saddle-point of $\mathcal{L}(\theta, \lambda) = l(\theta) + \sum_{i=1}^m \lambda_i c_i(\theta)$, at $\nabla \mathcal{L}(\theta, \lambda) = \vec{0}$.
- In practice, this is done via iterative gradient descent and ascent steps.

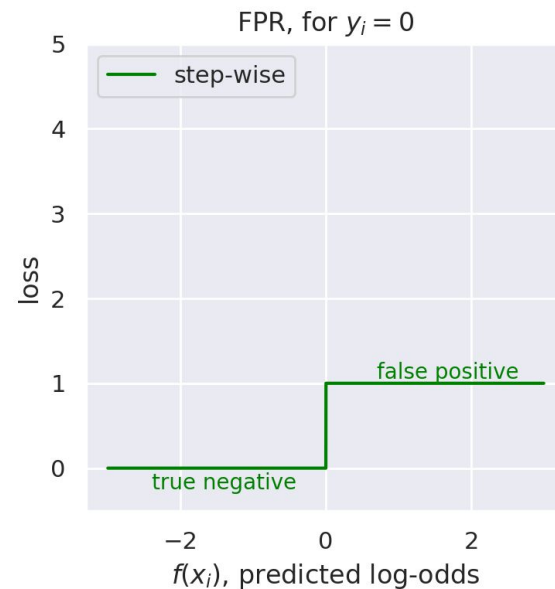
- Gradient descent over θ : $\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial l}{\partial \theta} + \sum_{i=1}^m \lambda_i \frac{\partial c_i}{\partial \theta}$ we require differentiable constraints

- Gradient ascent over λ : $\frac{\partial \mathcal{L}}{\partial \lambda_i} = c_i(\theta)$ but *only* for the descent step!

However, fairness metrics are **non-differentiable** (and also non-convex)

Issue:

Fairness metrics are non-differentiable (and non-convex).



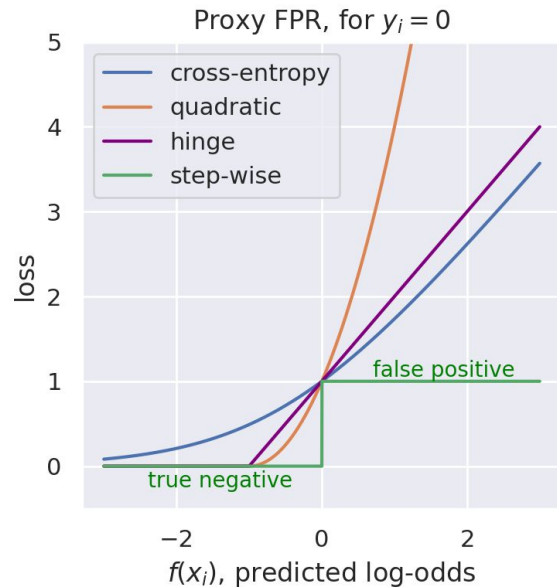
Issue:

Fairness metrics are non-differentiable (and non-convex).

Solution

- Use a differentiable upper-bound proxy: $c_i(\theta) \leq \tilde{c}_i(\theta) \leq 0$.
 - Gradient descent on the proxy-Lagrangian:

$$\frac{\partial \tilde{\mathcal{L}}}{\partial \theta} = \frac{\partial l}{\partial \theta} + \sum_{i=1}^m \lambda_i \frac{\partial \tilde{c}_i}{\partial \theta}$$



Issue:

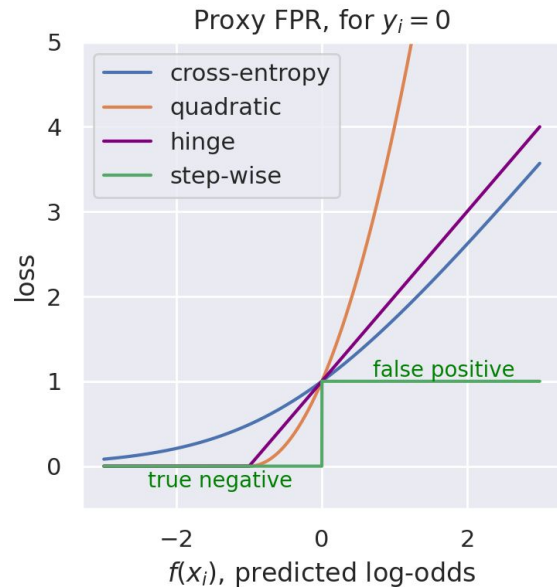
Fairness metrics are non-differentiable (and non-convex).

Solution

- Use a differentiable upper-bound proxy: $c_i(\theta) \leq \tilde{c}_i(\theta) \leq 0$.
 - Gradient descent on the proxy-Lagrangian:

$$\frac{\partial \tilde{\mathcal{L}}}{\partial \theta} = \frac{\partial l}{\partial \theta} + \sum_{i=1}^m \lambda_i \frac{\partial \tilde{c}_i}{\partial \theta}$$

- However, we'll be doing gradient descent and ascent on **different functions**: $\tilde{\mathcal{L}}, \mathcal{L}$.



Issue

- We're doing gradient descent over $\tilde{\mathcal{L}}$, and ascent over \mathcal{L} .
 - It's no longer a zero-sum game.
- Our proxy constraints are still non-convex (although differentiable).

Issue

- We're doing gradient descent over $\tilde{\mathcal{L}}$, and ascent over \mathcal{L} .
 - It's no longer a zero-sum game.
- Our proxy constraints are still non-convex (although differentiable).

Solution

- Instead of finding a pure Nash equilibrium: a single (θ, λ) pair (a single classifier).
- The solution lies on a distribution over (θ, λ) pairs (a distribution of classifiers) [Cotter et al., 2019]
 - A **randomized classifier**.

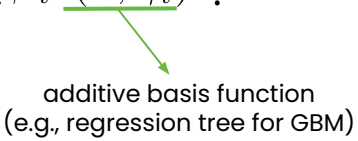
Related Work

- [Agarwal et al., 2018]: an *ensemble* of standard strong learners.
 - General method that is also compatible with GBM.
 - Increases **training time** dramatically (default: **50x increase**).

Related Work

- [Agarwal et al., 2018]: an *ensemble* of standard strong learners.
 - General method that is also compatible with GBM.
 - Increases **training time** dramatically (default: **50x increase**).

FairGBM

- Boosting fits an additive model: $f(x) = \sum_{i=1}^M \beta_i \underline{b(x, \gamma_i)}$.
 - 

additive basis function
(e.g., regression tree for GBM)
- Model iterate at iteration k : simply add the first k trees!

Related Work

- [Agarwal et al., 2018]: an *ensemble* of standard strong learners.
 - General method that is also compatible with GBM.
 - Increases **training time** dramatically (default: **50x increase**).

FairGBM

- Boosting fits an additive model: $f(x) = \sum_{i=1}^M \beta_i b(x, \gamma_i)$.

additive basis function
(e.g., regression tree for GBM)

- Model iterate at iteration k : simply add the first k trees!
- FairGBM randomized classifier carries virtually **no CPU or memory overhead**.

over the model iterates

Algorithm 1 FairGBM training pseudocode

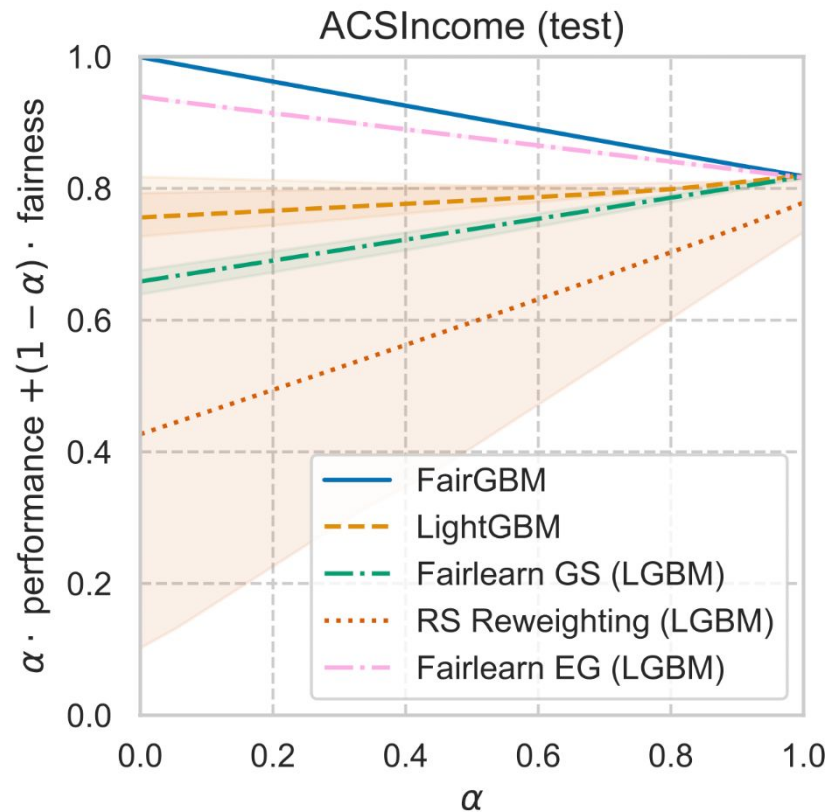
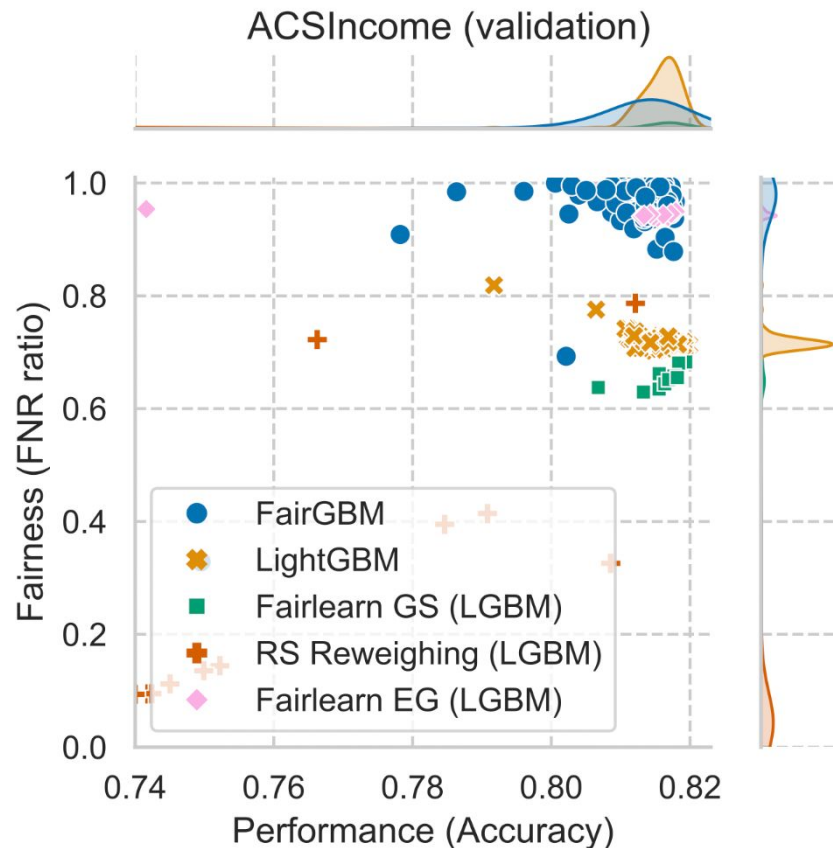
Input: $T \in \mathbb{N}$, number of boosting rounds

$\mathcal{L}, \tilde{\mathcal{L}} : \mathcal{F} \times \mathbb{R}_+^m \rightarrow \mathbb{R}$, Lagrangian and proxy-Lagrangian

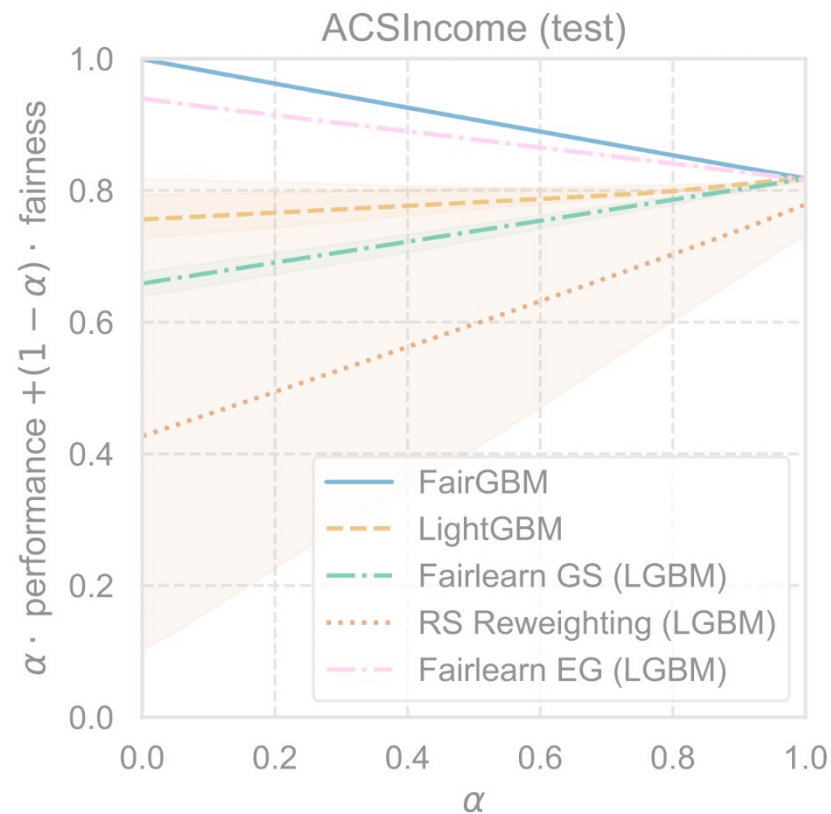
$\eta_f, \eta_\lambda \in \mathbb{R}_+$, learning rates

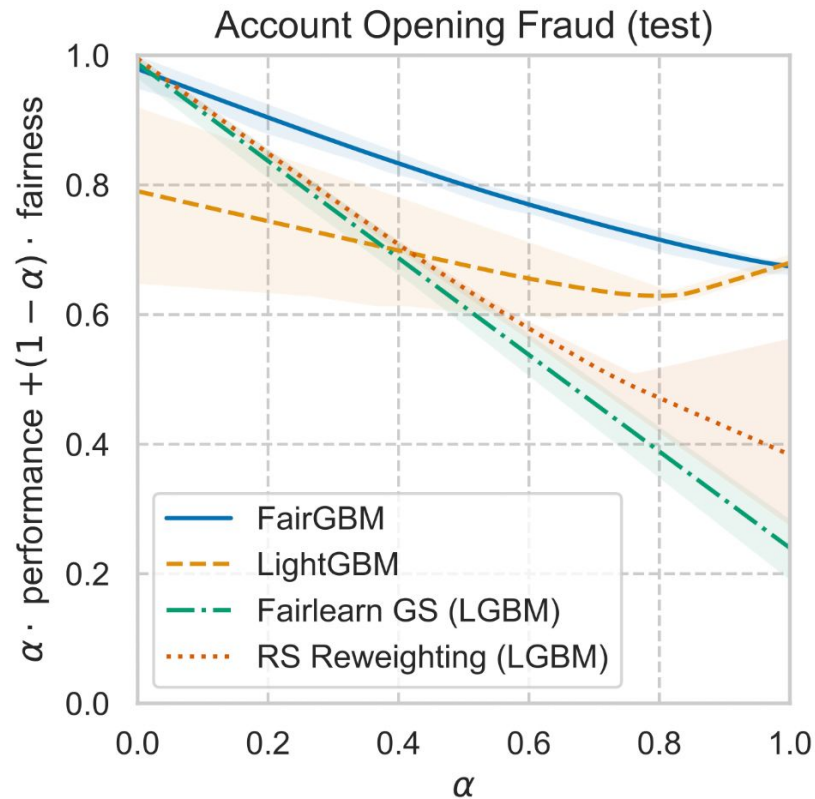
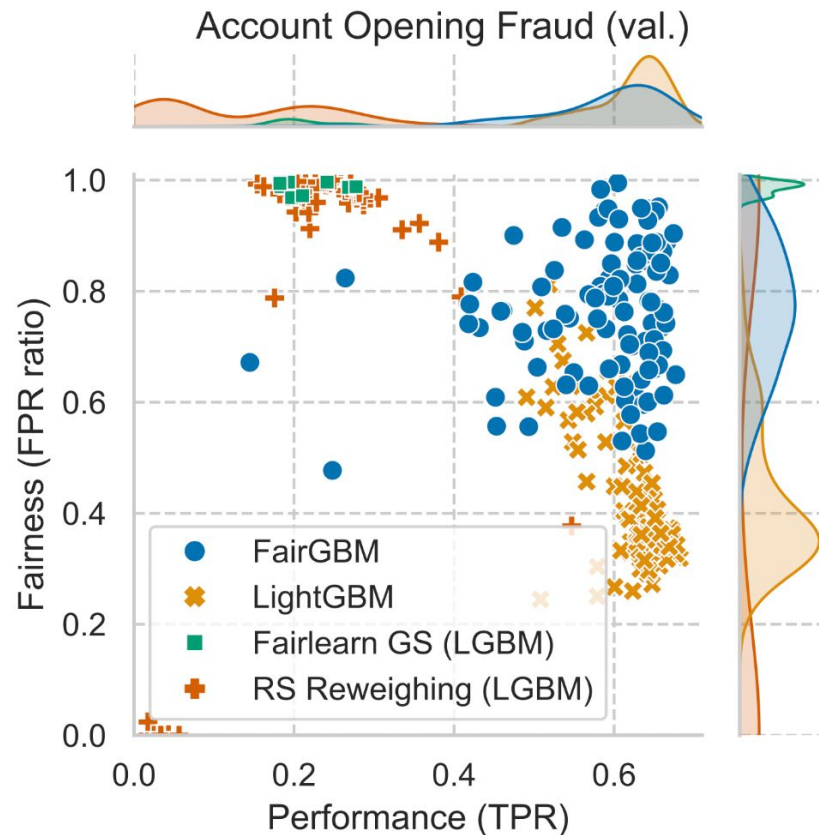
- 1: Let $h_0 = \arg \min_{\gamma \in \mathbb{R}} \tilde{\mathcal{L}}(\gamma, 0)$ ▷ Initial constant “guess”
 - 2: Initialize $f \leftarrow h_0$
 - 3: Initialize $\lambda \leftarrow 0$
 - 4: **for** $t \in \{1, \dots, T\}$ **do**
 - 5: Let $g_i = \frac{\partial \tilde{\mathcal{L}}(f(x_i), \lambda)}{\partial f(x_i)}$ ▷ Gradient of proxy-Lagrangian w.r.t. model
 - 6: Let $\Delta = \frac{\partial \mathcal{L}(f(x_i), \lambda)}{\partial \lambda}$ ▷ Gradient of Lagrangian w.r.t. multipliers
 - 7: Let $h_t = \arg \min_{h_t \in \mathcal{H}} \sum_{i=1}^N (-g_i - h_t(x_i))^2$ ▷ Fit base learner
 - 8: Update $f \leftarrow f + \eta_f h_t$ ▷ Gradient descent
 - 9: Update $\lambda \leftarrow (\lambda + \eta_\lambda \Delta)_+$ ▷ Projected gradient ascent
 - 10: **return** h_0, \dots, h_T
-

- Datasets
 - ACS folktables datasets [Ding et al., 2021]
 - Five public datasets based on US census data;
 - Includes modern-day version of UCI Adult dataset;
 - 1M to 3M rows;
 - Account Opening Fraud dataset
 - In-house real-world data stream of account opening fraud on a major European bank;
 - 500K rows;
- Literature baselines:
 - Fairlearn Exponentiated Gradient Reduction [Agarwal et al., 2018]
 - Fairlearn Grid Search [Agarwal et al., 2018]
 - Fairlearn Random Search
 - Unconstrained LightGBM [Ke et al., 2017]

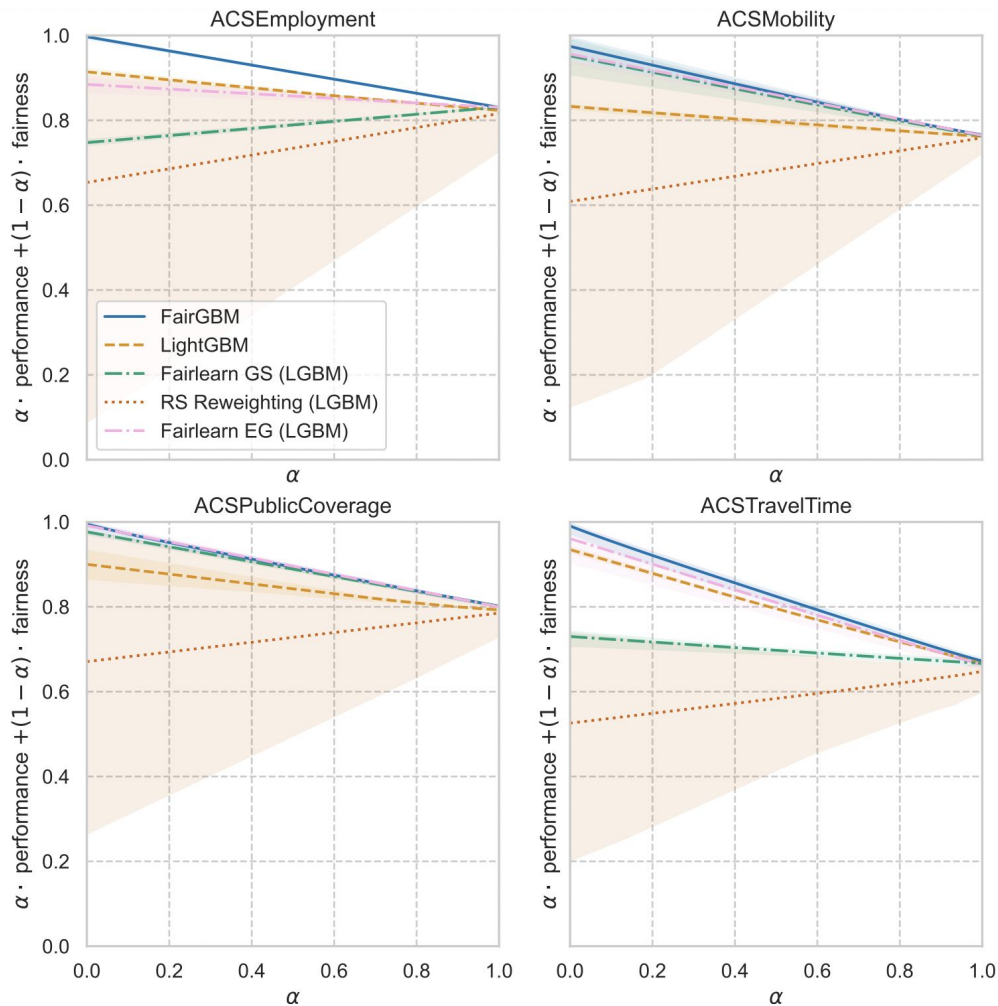


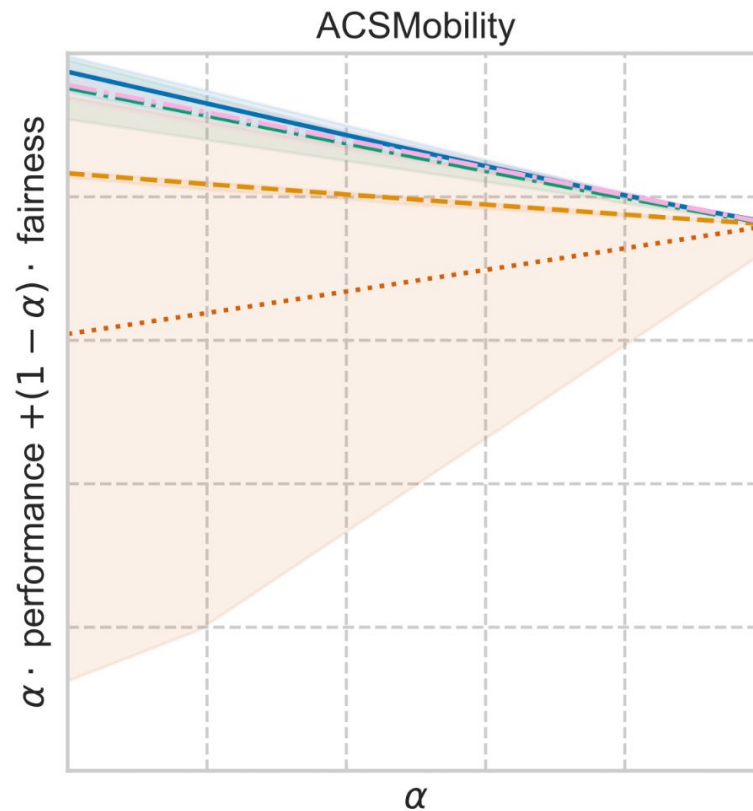
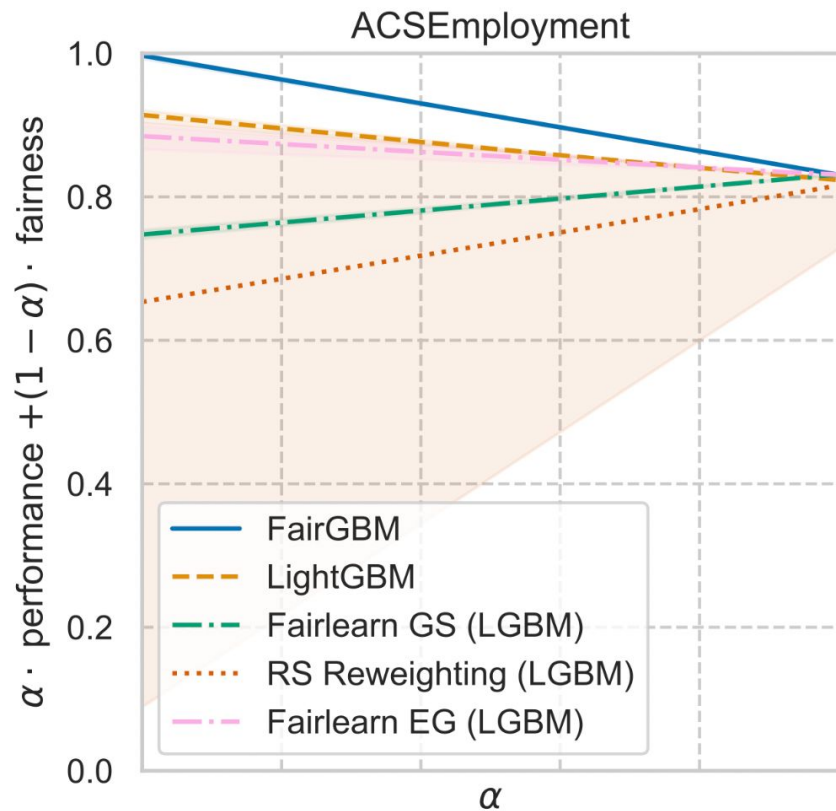
Algorithm	Run-time	
	Total (h)	Relative
<i>ACSIncome</i>		
FairGBM	9.9	x2.1
LightGBM	4.6	<i>baseline</i>
GS	43.8	x9.6
RS	37.1	x8.1
EG	99.4	x21.7



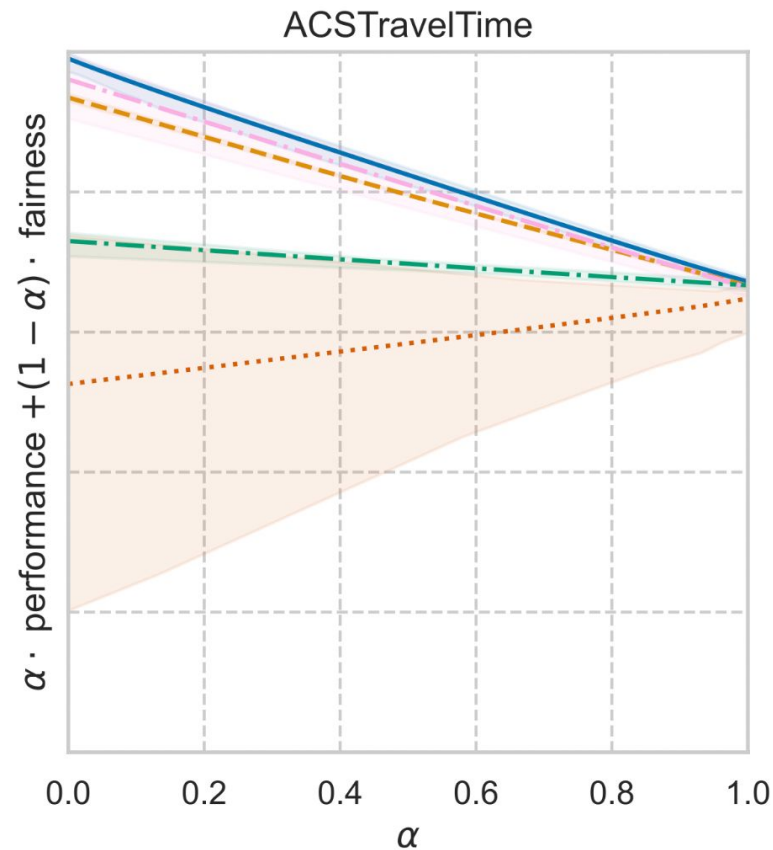
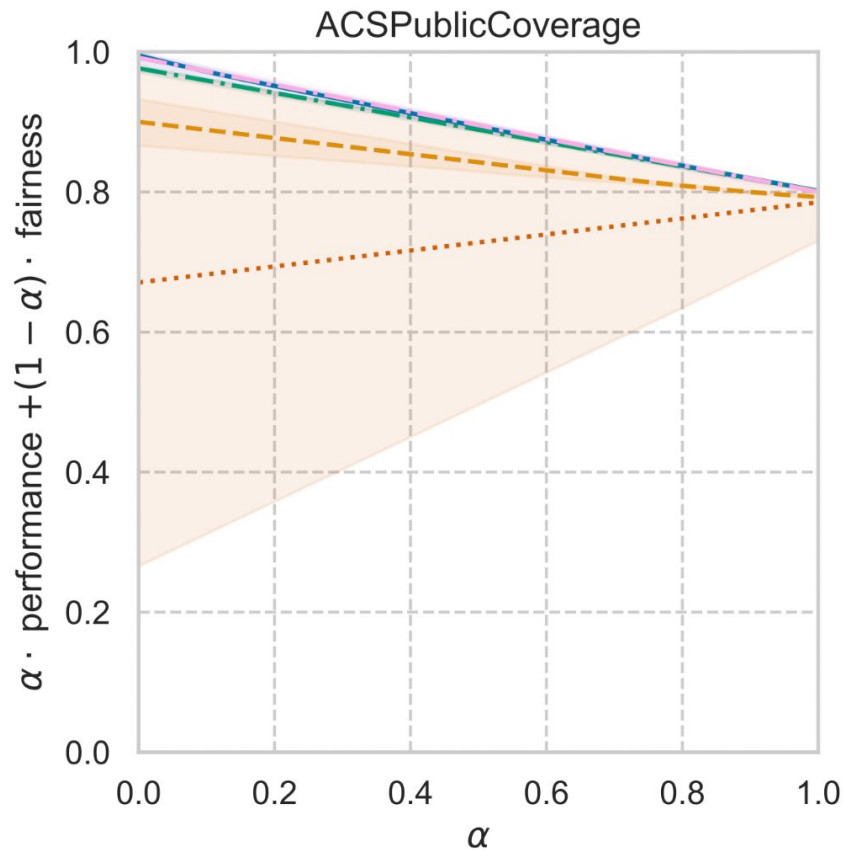


Other Datasets





Results on ACS datasets



- Open-source package available at <https://github.com/feedzai/fairgbm/>

- Try it out for yourself (on linux)



```
from fairgbm import FairGBMClassifier

# Instantiate
fairgbm_clf = FairGBMClassifier(
    constraint_type="FNR", # constraint on equal group-wise TPR (equal opportunity)
    n_estimators=200,     # core parameters from vanilla LightGBM
    random_state=42,     # ...
)

# Train using features (X), labels (Y), and sensitive attributes (S)
fairgbm_clf.fit(X, Y, constraint_group=S)
# NOTE: labels (Y) and sensitive attributes (S) must be in numeric format

# Predict
Y_test_pred = fairgbm_clf.predict_proba(X_test)[:, -1] # Compute continuous class probabilities (re
# Y_test_pred = fairgbm_clf.predict(X_test)             # Or compute discrete class predictions
```

Thank You

Check out the paper at <https://arxiv.org/pdf/2209.07850>

Check out the code at <https://github.com/feedzai/fairgbm>

