

Paper: OmniFair: A Declarative System for Model-Agnostic Group Fairness in Machine Learning. [1]

Jiabin Liu

Group Reading

November 4, 2021

1 Introduction

- Group Fairness Constraints

2 OmniFair

- Declarative Specification
- Single Fairness Constraint
- Multiple Fairness Constraints

3 Experiments

Notations and group fairness constraints

- $x \in \mathbb{R}^d$, $y \in \{0, 1\}$ and $D = \{(x_i, y_i)\}_{i=1}^N$
- Goal: $h_\theta(x)$, where $h: \mathbb{R}^d \rightarrow \{0, 1\}$.

Group Fairness constraints:

- Statistical Parity (SP): $\forall g_i, g_j \in G, \Pr(h(x) = 1|g_i) \simeq \Pr(h(x) = 1|g_j)$
- False Positive Rate Parity (FPR):
 $\forall g_i, g_j \in G, \Pr(h(x) = 1|g_i, y = 0) \simeq \Pr(h(x) = 1|g_j, y = 0)$
- False Omission Rate Parity (FOR):
 $\forall g_i, g_j \in G, \Pr(y = 1|g_i, h(x) = 0) \simeq \Pr(y = 1|g_j, h(x) = 0)$
- Misclassification Rate Parity (MR):
 $\forall g_i, g_j \in G, \Pr(h(x) = y|g_i) \simeq \Pr(h(x) = y|g_j)$
- False Negative Rate Parity (FNR), False Discovery Rate Parity (FDR)

DEFINITION 1. Fairness Specification (g, f, ε) and Fairness Constraint

A fairness specification is said to be satisfied by a classifier h on D if and only if all induced fairness constraints are satisfied, i.e.,

$$\forall g_i, g_j \in g(D), |f(h, g_i) - f(h, g_j)| \leq \varepsilon$$

where $g(D)$ is the declarative grouping function and $f(h, g)$ is the declarative fairness metric function.

Problem Formulation:

- Given a dataset D , an ML algorithm \mathcal{A}
- One or multiple group fairness constraints given by one or multiple (g, f, ε) .
- Goal: h_θ that maximizes for accuracy, while satisfying given constraint(s).

Declarative Grouping Function $g(D)$

$g(D)$ is a user-defined function that partitions the input data to different groups.

Example 1

A dataset $D = \{t_1, \dots, t_{10}\}$, where t_4, t_5, t_7 and t_9 are African American and others are Caucasian. The grouping function partitions D as $g(D) = \{AfricanAmerican : [4, 5, 7, 9], Caucasian : [1, 2, 3, 6, 8, 10]\}$.

Fairness Constraint

```
def grouping(Dataset D)
  groups = {}

  # user code here with example
  groups['African-American'] = []
  groups['Caucasian'] = []
  for i in range(D.shape(0)):
    groups[D[i]['race']].append(i)

  return groups
```

Figure: Interface for group function.

Declarative Fairness Metric Function $f(h, g)$

$f(h, g)$ takes a classifier h and a group g as input, and returns $(1 + |g|)$ coefficients that specify how the metric is computed:

$$f(h, g) = \sum_{i \in g} c_i \mathbb{1}(h(x_i) = y_i) + c_0$$

	$c_i y_i = 0$	$c_i y_i = 1$	c_0
MR	$1/ g $	$1/ g $	0
SP	$-1/ g $	$1/ g $	$ \{i : i \in g, y_i = 0\} / g $
FPR	$1/ \{i : i \in g, y_i = 0\} $	0	0
FNR	0	$1/ \{i : i \in g, y_i = 1\} $	0
FOR	$1/ \{i : i \in g, h(x_i) = 0\} $	0	0
FDR	0	$1/ \{i : i \in g, h(x_i) = 1\} $	0

Figure: Coefficients for different popular group fairness metrics.

Single Fairness Constraint

- Accuracy part: $AP(\theta) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(h_{\theta}(x_i) = y_i)$
- Fairness part: $FP(\theta) = f(h_{\theta}, g_1) - f(h_{\theta}, g_2)$

$$\max_{\theta} AP(\theta) \quad (1)$$

$$\text{s.t. } |FP(\theta)| \leq \varepsilon \quad (2)$$

Lagrangian dual function:

$$\begin{aligned} h(\lambda_1, \lambda_2) &= \max_{\theta} AP(\theta) + \lambda_1(\varepsilon - FP(\theta)) + \lambda_2(\varepsilon + FP(\theta)) \\ &= \max_{\theta} AP(\theta) + (\lambda_2 - \lambda_1)FP(\theta) + (\lambda_1 + \lambda_2)\varepsilon \end{aligned}$$

* $h(\lambda_1, \lambda_2)$ provides an upper bound for eq.1 for any $\lambda_1 > 0, \lambda_2 > 0$.

$$\max_{\theta} AP(\theta) + (\lambda_2 - \lambda_1)FP(\theta) + (\lambda_1 + \lambda_2)\varepsilon \quad (3)$$

$$\max_{\theta} AP(\theta) + \lambda FP(\theta) \quad (4)$$

LEMMA 1.

Assume the primal equation is feasible and let θ^* be an optimal solution to primal problem, then

- ① for any $\lambda_1, \lambda_2 > 0$, $h(\lambda_1, \lambda_2) \geq AP(\theta^*)$; and under strong duality assumption, $\min_{\lambda_1 > 0, \lambda_2 > 0} h(\lambda_1, \lambda_2) = AP(\theta^*)$
- ② for any $\lambda_1, \lambda_2 > 0$, let $\tilde{\theta}$ be an optimal solution to eq.(3), then there exists $\lambda \in \mathbb{R}$ (i.e. $\lambda = \lambda_2 - \lambda_1$) such that $\tilde{\theta}$ also optimizes eq.(4); and under strong duality assumption, there exists λ such that θ^* optimizes eq.(4).

$$\max_{\theta} AP(\theta) + \lambda FP(\theta) \quad (5)$$

Steps:

- 1 enumerate all possible λ values;
- 2 find optimal θ for every λ
- 3 pick the one that has the maximal $AP(\theta)$.

How to solve eq. (5) for a given λ ? (step 2)

$$\begin{aligned} & \max_{\theta} AP(\theta) + \lambda FP(\theta) \\ &= \max_{\theta} \frac{1}{N} \sum_{i=1}^N \mathbb{1}_i + \lambda \left(\sum_{i \in \mathcal{G}_1} (c_i^{\mathcal{G}_1} \mathbb{1}_i + c_0^{\mathcal{G}_1}) - \sum_{i \in \mathcal{G}_2} (c_i^{\mathcal{G}_2} \mathbb{1}_i + c_0^{\mathcal{G}_2}) \right) \\ &= \max_{\theta} \frac{1}{N} \sum_{i=1}^N w_i(\lambda, h_{\theta}) \mathbb{1}(h_{\theta}(x_i) = y_i) \end{aligned}$$

How to get weight $w_i(\lambda)/w_i(\lambda, h_\theta)$?

weight	metric	$w_i y_i = 0, g_1$	$w_i y_i = 1, g_1$	$w_i y_i = 0, g_2$	$w_i y_i = 1, g_2$
$w_i(\lambda)$	MR	$1 + \lambda N / g_1 $	$1 + \lambda N / g_1 $	$1 - \lambda N / g_2 $	$1 - \lambda N / g_2 $
	SP	$1 - \lambda N / g_1 $	$1 + \lambda N / g_1 $	$1 + \lambda N / g_2 $	$1 - \lambda N / g_2 $
	FPR	$1 - \lambda N / \{i : i \in g_1, y_i = 0\} $	1	$1 + \lambda N / \{i : i \in g_2, y_i = 0\} $	1
	FNR	1	$1 - \lambda N / \{i : i \in g_1, y_i = 1\} $	1	$1 + \lambda N / \{i : i \in g_2, y_i = 1\} $
$w_i(\lambda, h_\theta)$	FOR	$1 - \lambda N / \{i : i \in g_1, h(x_i) = 0\} $	1	$1 + \lambda N / \{i : i \in g_2, h(x_i) = 0\} $	1
	FDR	1	$1 - \lambda N / \{i : i \in g_1, h(x_i) = 1\} $	1	$1 + \lambda N / \{i : i \in g_2, h(x_i) = 1\} $

Figure: Weights for different popular group fairness metrics.

- For $w_i(\lambda)$: any black-box ML algorithm.
- For $w_i(\lambda, h_\theta)$:
 - if $\lambda_2 - \lambda_1 \leq \delta = 0.00001$, then $h_{\theta_1|\lambda_1}(x_i) = h_{\theta_2|\lambda_2}(x_i)$.
 - So $w_i(\lambda_2, h_{\theta_2}) \simeq w_i(\lambda_2, h_{\theta_1})$.
 - Starting from $\lambda = 0$ and taking small incremental steps δ to estimate weights.

How to tune the hyperparameter λ ?

LEMMA 2. Monotonicity for Single Fairness Constraint

Consider two values λ_1 and λ_2 , where $\lambda_1 < \lambda_2$, and let θ_1 and θ_2 denote two optimal solutions given λ_1 and λ_2 . Then, the following properties hold, where $AP(\theta_1), AP(\theta_2), FP(\theta_1), FP(\theta_2)$ are evaluated on the same input training set D :

$$\begin{aligned} FP(\theta_1) &\leq FP(\theta_2) \\ AP(\theta_1) &\geq AP(\theta_2) \quad \text{where } \lambda_2 \geq \lambda_1 \geq 0 \\ AP(\theta_1) &\leq AP(\theta_2) \quad \text{where } 0 \geq \lambda_2 \geq \lambda_1 \end{aligned}$$

* ML model has the maximum accuracy when $\lambda = 0$.

Algorithm for tuning λ

Algorithm 1 Tuning Single λ

Input: Dataset D , a fairness constraint (g, f, ϵ) , an ML Algorithm \mathcal{A}

Output: A fair ML model h_θ

```
1:  $\theta_0 \leftarrow$  apply  $\mathcal{A}$  with  $w_i(0)$ 
2:  $flag \leftarrow false$  if (weights are parameterized by  $\theta$ ) else  $true$ 
3: if  $|FP(\theta_0)| \leq \epsilon$  then return  $h_{\theta_0}$ 
4: if  $FP(\theta_0) > 0$  then
5:   change the order of  $g_1$  and  $g_2$  in  $FP$ 
6:  $\lambda_l \leftarrow 0$  and  $\lambda_u \leftarrow 1$ 
7: if  $flag = true$  then
8:    $\lambda_l, \lambda_u \leftarrow$  EXPONENTIALSEARCH( $\lambda_l, \lambda_u$ )
9: else
10:   $\lambda_l, \lambda_u \leftarrow$  LINEARSEARCH( $\lambda_l, \delta$ ) // e.g.  $\delta = 0.001$ 
11: while  $\lambda_u - \lambda_l \geq \tau \delta$  //  $\tau \rightarrow 0$ ; e.g.  $\tau = 0.0001$ 
12:    $\lambda_m \leftarrow (\lambda_u + \lambda_l)/2$ 
13:   if  $flag = true$  then
14:      $\theta_m \leftarrow$  apply  $\mathcal{A}$  with  $w_i(\lambda_m)$ 
15:   else
16:      $\theta_m \leftarrow$  apply  $\mathcal{A}$  with  $w_i(\lambda_m, h_{\theta_l})$ 
17:   if  $FP(\theta_m) < -\epsilon$  then  $\lambda_l \leftarrow \lambda_m$ 
18:   else  $\lambda_u \leftarrow \lambda_m$ 
19: return  $h_{\theta_m}$ 
```

```
20:
21: function EXPONENTIALSEARCH( $\lambda_l, \lambda_u$ )
22:    $\theta_u \leftarrow$  apply  $\mathcal{A}$  with  $w_i(\lambda_u)$ 
23:   while  $FP(\theta_u) < -\epsilon$  do
24:      $\lambda_l \leftarrow \lambda_u$ 
25:      $\lambda_u \leftarrow 2 \times \lambda_u$ 
26:      $\theta_u \leftarrow$  apply  $\mathcal{A}$  with  $w_i(\lambda_u)$ 
27:   return  $\lambda_l, \lambda_u$ 
28:
29: function LINEARSEARCH( $\lambda_l, \delta$ )
30:    $\lambda_u \leftarrow \lambda_l + \delta$ 
31:    $\theta_u \leftarrow$  apply  $\mathcal{A}$  given  $\lambda_u$ 
32:   while  $FP(\theta_u) < -\epsilon$  do
33:      $\lambda_l \leftarrow \lambda_u$ 
34:      $\theta_l \leftarrow \theta_u$ 
35:      $\lambda_u \leftarrow \lambda_l + \delta$ 
36:      $\theta_u \leftarrow$  apply  $\mathcal{A}$  with  $w_i(\lambda_u, \theta_l)$ 
37:   return  $\lambda_l, \lambda_u$ 
```

Multiple Fairness Constraints

$$\begin{aligned} & \max_{\theta} AP(\theta) \\ & \text{s.t. } |FP_i(\theta)| \leq \varepsilon \quad \forall i \in \{1, \dots, k\} \end{aligned}$$

Lagrangian dual function is:

$$\begin{aligned} h(\Lambda_1, \Lambda_2) &= \max_{\theta} AP(\theta) + \sum_{i=1}^k \lambda_{1,i}(\varepsilon - FP_i(\theta)) + \sum_{i=1}^k \lambda_{2,i}(\varepsilon + FP_i(\theta)) \\ &= \max_{\theta} AP(\theta) + \sum_{i=1}^k (\lambda_{2,i} - \lambda_{1,i})FP_i(\theta) + (\lambda_{1,i} + \lambda_{2,i})\varepsilon \end{aligned}$$

where $\Lambda_1 = \langle \lambda_{11}, \lambda_{12}, \dots, \lambda_{1k} \rangle$ and $\Lambda_2 = \langle \lambda_{21}, \lambda_{22}, \dots, \lambda_{2k} \rangle$.

$$\begin{aligned} & \max_{\theta} AP(\theta) + \sum_{i=1}^k \lambda_i FP_i(\theta) \\ &= \max_{\theta} \frac{1}{N} \sum_{i=1}^N w_i(\Lambda, h_{\theta}) \mathbb{1}(h_{\theta}(x_i) = y_i) \end{aligned}$$

LEMMA 4. Marginal Monotonicity for Multiple Fairness Constraints

Consider two settings Λ_1 and Λ_2 , that differ only in the j^{th} dimension, namely, $\Lambda_1[j] < \Lambda_2[j]$ and $\Lambda_1[i] < \Lambda_2[i]$ for all $i \neq j$. Let θ_1 and θ_2 denote the optimal solution given Λ_1 and Λ_2 . Then

$$FP_j(\theta_1) \leq FP_j(\theta_2)$$

where $FP_j(\theta_1)$ and $FP_j(\theta_2)$ are evaluated on the same training set D .

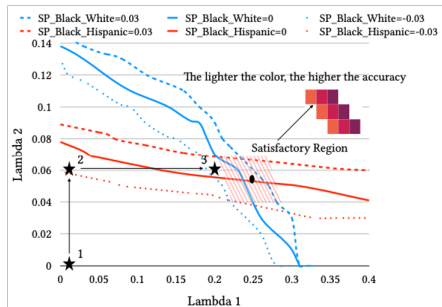
Algorithm for tuning Λ

Algorithm 2 Hill-Climbing

Input: Dataset D , a set of k fairness constraint, and an ML Algorithm \mathcal{A}

Output: A fair ML model h_θ

- 1: $\Lambda \leftarrow [0, \dots, 0]$
 - 2: $\theta \leftarrow$ apply \mathcal{A} with $w_i(0, -)$
 - 3: **while** $\exists i$, s.t. $|FP_i(\theta)| > \varepsilon_i$ **do**
 - 4: $i \leftarrow \arg \max_k |FP_k(\theta)| - \varepsilon_k$
 - 5: $\theta \leftarrow$ call Algorithm 1 to tune for the i -th fairness constraint, while fixing $\Lambda[j], \forall j \neq i$ to their current values
 - 6: **return** Λ **if** ($\Lambda \in$ intersection of all satisfactory regions) **else** "Not found after $5k$ iterations"
-



Dataset and Experimental settings

- Dataset:

Dataset	# Rows	# Attrs	Sens. Attr	Task
Adult [18]	48842	18	sex	To predict if Income > 50k
Compas [4]	11001	10	race	To predict recidivism
LSAC [45]	27477	12	race	To predict if bar exam is passed
Bank [32]	30488	20	age	To predict if marketing works

Split each dataset to 60% training, 20% validation, and 20% test.

- ML algorithms:

- Logistic Regression (LR)
- Random Forest (RF)
- XGBoost (XGB)
- Neural Networks (NN)

Experiments

Algorithm	COMPAS					Adult					LSAC					Bank				
	LR	RF	XGB	NN	CMA-ES	LR	RF	XGB	NN	CMA-ES	LR	RF	XGB	NN	CMA-ES	LR	RF	XGB	NN	CMA-ES
OmniFair	-1.2%	-0.8%	-0.7%	-1.2%	NA(2)*	-2.1%	-1.9%	-1.7%	-1.7%	NA(2)*	-0.3%	-0.3%	-0.4%	-0.1%	NA(2)*	-0.1%	-0.3%	-0.2%	-0.1%	NA(2)*
Kamiran et al. [28]	-2.5%	-1.3%	-1.2%	-1.5%	NA(2)*	-2.7%	-2.3%	-1.8%	-1.9%	NA(2)*	-0.4%	-5.6%	-2.2%	-0.4%	NA(2)*	+0.1%	-1.2%	-0.2%	-0.3%	NA(2)*
Calmon et al. [11]	-1.8%	-0.5%	-0.3%	-0.9%	NA(2)*	-3.7%	-3.1%	-3.0%	-2.4%	NA(2)	NA(1)	NA(1)	NA(1)	NA(1)	NA(2)*	NA(1)	NA(1)	NA(1)	NA(1)	NA(2)*
Zafar et al. [47]	-0.9%	NA(2)	NA(2)	NA(2)	NA(2)	-2.2%	NA(2)	NA(2)	NA(2)	NA(2)	-0.2%	NA(2)	NA(2)	NA(2)	NA(2)	-0.1%	NA(2)	NA(2)	NA(2)	NA(2)
Celis et al. [12]	NA(1)	NA(2)	NA(2)	NA(2)	NA(2)	NA(1)	NA(2)	NA(2)	NA(2)	NA(2)	NA(1)	NA(2)	NA(2)	NA(2)	NA(1)	NA(2)	NA(2)	NA(2)	NA(2)	NA(2)
Agarwal et al. [3]	-2.4%	-1.2%	-2.0%	-1.8%	NA(2)*	-2.8%	-2.2%	-2.0%	-2.0%	NA(2)*	-0.6%	-5.8%	-0.2%	-0.5%	NA(2)*	-0.1%	-0.3%	-0.0%	-0.1%	NA(2)*
Thomas et al. [43]	NA(2)	NA(2)	NA(2)	NA(2)	-1.1%	NA(2)	NA(2)	NA(2)	NA(2)	-1.7%	NA(2)	NA(2)	NA(2)	NA(2)	-0.4%	NA(2)	NA(2)	NA(2)	NA(2)	-0.1%

Figure: Accuracy drop compared with no fairness constraints when $\epsilon = 0.03$ under SP.

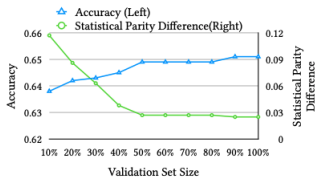


Figure: Ablation study for the validation size on the COMPAS set.

Experiments

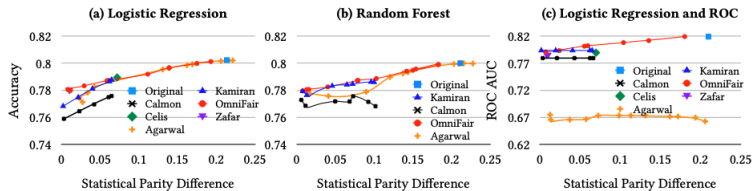


Figure: The trade-off between fairness metric and accuracy on Adult dataset.

ϵ	Accuracy	SP	FNR
Baseline	63.9%	0.233	0.180
0.01	N/A	N/A	N/A
0.02	N/A	N/A	N/A
0.03	63.6%	0.03	0.044
0.04	63.4%	0.016	0.035
0.05	63.3%	0.028	0.007
0.06	62.7%	0.057	0.032

Figure: Enforcing SP and FNR on COMPAS dataset.

- [1] Hantian Zhang et al. “OmniFair: A Declarative System for Model-Agnostic Group Fairness in Machine Learning”. In: *Proceedings of the 2021 International Conference on Management of Data*. 2021, pp. 2076–2088.