# SimFair: A Unified Framework for Fairness-Aware Multi-Label Classification (AAAI 2023)

Liu et. al.

Reviewer: Jihu Lee, Jinwon Park

IDEA lab
Department of Statistics
Seoul National University

March, 21, 2024

# Table of Contents

# Multi-Label Classification

**Examples**

- An applicant may apply for multiple positions
- Undergraduates submit applications to multiple programs when applying to graduate schools
- **Target label**: Admission decision of each position

**Simple Approach**

- Decompose into multiple binary classification
- Ignores the correlations among labels
  - Applicants usually apply for positions with similar requirements of skill sets and experiences

# Fairness in Multi-Label Classification

- No existing work to define fairness directly in the context of multi-label classification
- Proposed framework $s_\gamma$-**SimFair** unifies DP and EOp
- It works even when imbalanced label distributions exist

# Table of Contents

## Notations

- Samples: $\mathcal{D} = \left\{ \left( \boldsymbol{x}^{(i)}, a^{(i)}, \boldsymbol{y}^{(i)} \right) \right\}_{i=1}^{N}$
- $\boldsymbol{x}^{(i)} \in \mathcal{X} = \mathbb{R}^M$, $a^{(i)} \in \mathcal{A} = \{1, \ldots, K\}$, $\boldsymbol{y}^{(i)} \in \mathcal{Y} = \{0,1\}^L$
- $(\boldsymbol{x}, a, \boldsymbol{y}) \sim p$
- Classifier: $h = f \circ g : \mathcal{X} \to [0,1]^L \to \mathcal{Y}$

- **Advantaged Label**: $\boldsymbol{y}_{\mathsf{adv}}$
  - Only favorable outcomes present
  - Ex) *received offers of postion A and position B* from job screening example

# DP and EOp

- **DP**: $\hat{\boldsymbol{y}} \perp a$
- **EOp**: $\hat{\boldsymbol{y}} \perp a | \boldsymbol{y}_{\text{adv}}$

### Proposition (0.1)

*For a multi-label classifier that takes the form $h = f \circ g$, where $\tilde{\boldsymbol{y}} = g(\boldsymbol{x})$ is the predicted probability and $\hat{\boldsymbol{y}} = f(\tilde{\boldsymbol{y}})$ is computed elementwisely, DP and EOp hold if for any $k \in \mathcal{A}$*

$$
\begin{aligned}
DP &: \mathbb{E}\left[\tilde{\boldsymbol{y}} | a = k\right] = \mathbb{E}\left[\tilde{\boldsymbol{y}}\right] \\
EOp &: \mathbb{E}\left[\tilde{\boldsymbol{y}} | a = k, \boldsymbol{y} = \boldsymbol{y}_{\text{adv}}\right] = \mathbb{E}\left[\tilde{\boldsymbol{y}} | \boldsymbol{y} = \boldsymbol{y}_{\text{adv}}\right]
\end{aligned}
\tag{1}
$$

**Estimation**

$$\mathbb{E}[\tilde{\boldsymbol{y}}|a = k] \approx \frac{\sum_{i=1}^{N} \tilde{\boldsymbol{y}}^{(i)} \mathbf{1}\left(a^{(i)} = k\right)}{\sum_{i=1}^{N} \mathbf{1}\left(a^{(i)} = k\right)} \tag{2}$$

$$\mathbb{E}[\tilde{\boldsymbol{y}}|a = k, \boldsymbol{y} = \boldsymbol{y}_{\mathsf{adv}}] \approx \frac{\sum_{i=1}^{N} \tilde{\boldsymbol{y}}^{(i)} \mathbf{1}\left(a^{(i)} = k\right) \mathbf{1}\left(\boldsymbol{y} = \boldsymbol{y}_{\mathsf{adv}}\right)}{\sum_{i=1}^{N} \mathbf{1}\left(a^{(i)} = k\right) \mathbf{1}\left(\boldsymbol{y} = \boldsymbol{y}_{\mathsf{adv}}\right)} \tag{3}$$

- most labels only associate with few samples $\rightarrow$ additional challenges for EOp estimation

# $s_\gamma$-SimFair

EOp (Eq. (1)) is equivalent with

$$\frac{\mathbb{E}\left[\tilde{\boldsymbol{y}}\mathbf{1}\left(\boldsymbol{y}=\boldsymbol{y}_{\mathsf{adv}}\right)\right]}{\mathbb{E}\left[\mathbf{1}\left(\boldsymbol{y}=\boldsymbol{y}_{\mathsf{adv}}\right)\right]} = \frac{\mathbb{E}\left[\tilde{\boldsymbol{y}}\mathbf{1}\left(a=k\right)\mathbf{1}\left(\boldsymbol{y}=\boldsymbol{y}_{\mathsf{adv}}\right)\right]}{\mathbb{E}\left[\mathbf{1}\left(a=k\right)\mathbf{1}\left(\boldsymbol{y}=\boldsymbol{y}_{\mathsf{adv}}\right)\right]} \qquad (4)$$

### Definition 1 ($s$-SimFair)

Given a similarity function $s : \mathcal{Y} \times \mathcal{Y} \to [0,1]$, a multi-label classifier $h$ satisfiew Similarity $s$-induced Fairness if for $\forall k \in \mathcal{A}$,

$$\frac{\mathbb{E}\left[\tilde{\boldsymbol{y}}s\left(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}}\right)\right]}{\mathbb{E}\left[s\left(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}}\right)\right]} = \frac{\mathbb{E}\left[\tilde{\boldsymbol{y}}\mathbf{1}\left(a=k\right)s\left(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}}\right)\right]}{\mathbb{E}\left[\mathbf{1}\left(a=k\right)s\left(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}}\right)\right]} \qquad (5)$$

$$\mathsf{Jac}(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}}) = \frac{|\mathsf{cate}(\boldsymbol{y}) \cap \mathsf{cate}(\boldsymbol{y}_{\mathsf{adv}})|}{|\mathsf{cate}(\boldsymbol{y}) \cup \mathsf{cate}(\boldsymbol{y}_{\mathsf{adv}})|}$$

$$s_\gamma(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}}) = \exp\left(\gamma\left(\mathsf{Jac}\left(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}}\right)-1\right)\right)$$

1. DP and EOP are special cases of $s_\gamma$-SimFair
   - $s \equiv c \rightarrow$ DP
   - $s(\boldsymbol{y}, \boldsymbol{y}') = \mathbf{1}\,(\boldsymbol{y}, \boldsymbol{y}') \rightarrow$ EOp

2. $s_\gamma$-SimFair helps achieve DP and EOp
   - $\gamma$ sufficiently small $\rightarrow$ small DP violation
   - $\gamma$ sufficiently large $\rightarrow$ small EOp violation

**Violation of $s_\gamma$-SimFair**

$$
\begin{aligned}
l_{s_\gamma(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}})}(h) = \\
\sum_{k=1}^{K} \left\| \frac{\mathbb{E}\left[\tilde{\boldsymbol{y}}s_\gamma(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}})\right]}{\mathbb{E}\left[s_\gamma(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}})\right]} - \frac{\mathbb{E}\left[\tilde{\boldsymbol{y}}\mathbf{1}\left(a=k\right)s_\gamma(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}})\right]}{\mathbb{E}\left[\mathbf{1}\left(a=k\right)s_\gamma(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}})\right]} \right\|
\end{aligned}
\tag{6}
$$

**Objective**

$$
\min_{h} l_{\mathsf{mlc}}(h) + \lambda l_{s_\gamma(\boldsymbol{y},\boldsymbol{y}_{\mathsf{adv}})}(h)
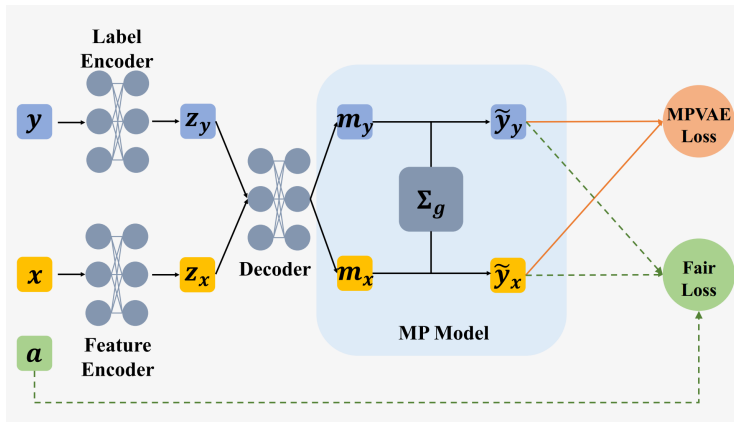\tag{7}
$$

Figure 1: Framework of training MPVAE with fairness regulrization. Both probability vectors $\tilde{y}$ on two brancehs are used to construct the $s_\gamma$-SimFair regularizer

# Table of Contents

- Adult
    - Label: income, workclass, occupation
    - Sensitive: binarized age
- Credit
    - Label: default payments, education level
    - Sensitive: gender

## Settings

**Baselines**

- MPVAE with No regularizer
- MPVAE with DP regularizer
- MPVAE with EOp regularizer

**Evaluation metrics**

- micro-F1, macro-F1, example-F1

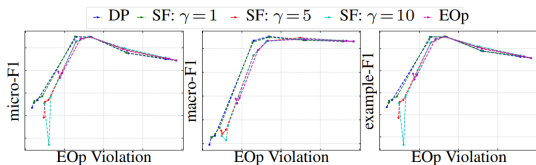| | $\boldsymbol{y}_{adv}$ | DP | $s_{0.1}$-SF | $s_{0.5}$-SF | $s_1$-SF | $s_5$-SF | $s_{10}$-SF | EOp |
|---|---|---|---|---|---|---|---|---|
| Adult | 100% | 0.11* | 0.12 | 0.12 | 0.13 | 0.18 | 0.17 | 0.17* |
| | 70% | 0.11 | 0.11 | 0.12 | 0.13 | 0.18 | 0.18 | 0.17 |
| | 30% | 0.10 | 0.10 | 0.11 | 0.12 | 0.17 | 0.18 | 0.17 |
| | 10% | 0.10 | 0.10 | 0.10 | 0.11 | 0.14 | 0.16 | 0.17 |
| | 5% | 0.10 | 0.10 | 0.10 | 0.11 | 0.15 | 0.23 | 0.27 |
| Credit | 100% | 0.03* | 0.03 | 0.03 | 0.03 | 0.03 | 0.03 | 0.03* |
| | 70% | 0.03 | 0.03 | 0.03 | 0.03 | 0.03 | 0.04 | 0.04 |
| | 30% | 0.02 | 0.02 | 0.02 | 0.02 | 0.03 | 0.03 | 0.03 |
| | 10% | 0.02 | 0.02 | 0.02 | 0.02 | 0.03 | 0.04 | 0.04 |
| | 5% | 0.02 | 0.02 | 0.02 | 0.02 | 0.03 | 0.04 | 0.05 |

Figure 2: Varying the numbers of samples in the advantaged group to different levels. Ground truth is marked with asterisk. More stable EOp estimates when EOp estimator fails.

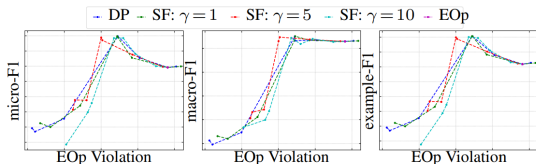|  | $|\boldsymbol{y}_{adv}|$ | Metric | Regularzier | | | | | |
|---|---|---|---|---|---|---|---|---|
|  |  |  | DP | $s_1$-SF | $s_5$-SF | $s_{10}$-SF | EOp | None |
| Adult | No.1 | DP | 0.038 | **0.031** | 0.038 | 0.043 | 0.045 | 0.111 |
|  |  | EOp | 0.051 | 0.042 | **0.030** | 0.034 | 0.035 | 0.161 |
|  | No.18 | DP | **0.038** | **0.038** | 0.043 | 0.045 | 0.094 | 0.111 |
|  |  | EOp | 0.076 | 0.072 | 0.037 | **0.027** | 0.066 | 0.095 |
| Credit | No.1 | DP | 0.018 | 0.018 | **0.017** | 0.018 | 0.018 | 0.029 |
|  |  | EOp | 0.026 | 0.026 | **0.025** | **0.025** | 0.026 | 0.038 |
|  | No.9 | DP | **0.018** | **0.018** | 0.019 | 0.019 | 0.030 | 0.030 |
|  |  | EOp | 0.202 | **0.192** | 0.193 | 0.197 | 0.241 | 0.241 |

Figure 3: DP and EOp violations of MPVAE trained with such regularizers. Best results are in bold.

(a) Credit dataset: No.1 label group



(b) Credit dataset: No.9 label group

Figure 4: EOp-accuracy tradeoffs on Credit dataset. EOp regularizer is unstable and ineffective when the advantaged group is small.