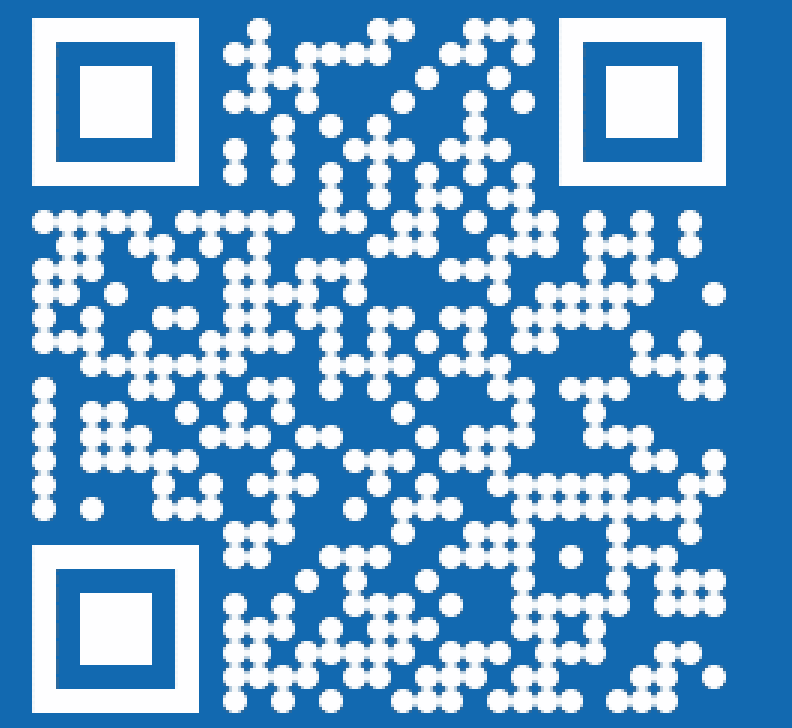


Debiasing Deep Chest X-Ray Classifiers using Intra-and Post-processing Methods

Ričards Marcinkevičs[✉], Ece Ozkan and Julia E. Vogt
 Department of Computer Science, ETH Zurich
 ✉: ricards.marcinkevics@inf.ethz.ch



1 Motivation

Given

- $X, A \in \{0,1\}, Y \in \{0,1\}$
- $\mathcal{D} = \{(x_i, y_i, a_i)\}_i = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{valid}} \cup \mathcal{D}_{\text{test}}$
- neural network $f_\theta(\cdot)$ trained on $\{(x_i, y_i)\}_i$ from $\mathcal{D}_{\text{train}}$

Goal: reduce the bias μ of $f_\theta(\cdot)$, sacrificing the performance ρ as little as possible

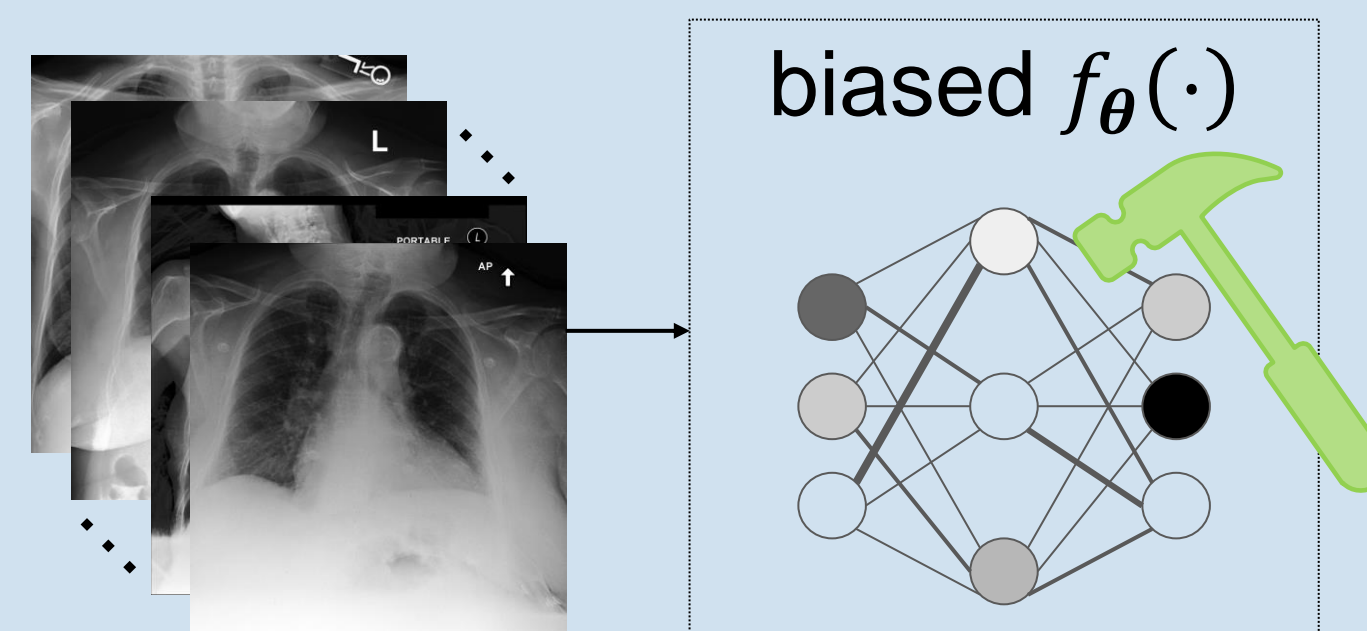
Bias Measures

$$\mu_{\text{SPD}} = \mathbb{P}_{X,A}(\hat{Y} = 1|A = 0) - \mathbb{P}_{X,A}(\hat{Y} = 1|A = 1)$$

$$\mu_{\text{EOD}} = \mathbb{P}_{X,Y,A}(\hat{Y} = 1|Y = 1, A = 0) - \mathbb{P}_{X,Y,A}(\hat{Y} = 1|Y = 1, A = 1)$$

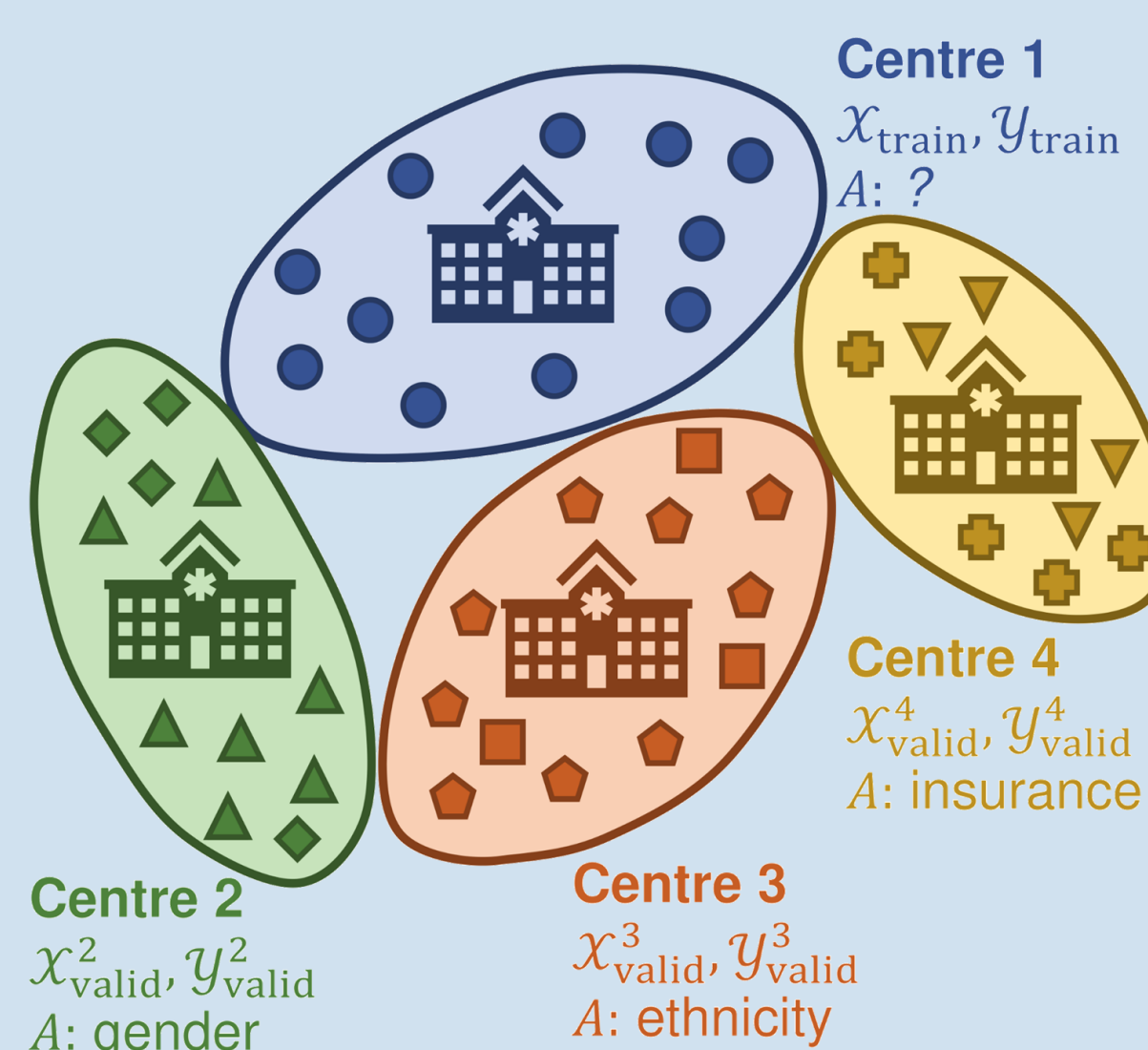
Intra-processing Setting

- debias *post hoc* on $\mathcal{D}_{\text{valid}}$
- edit the parameters θ
- A is not given at test time



Example

- train on centre 1
- deploy in centres 2, 3, 4 with local fairness constraints
- no access to $\mathcal{D}_{\text{train}}$



Contributions

- differentiable proxy functions for the SPD and EOD
- debiasing algorithms based on pruning and fine-tuning
- experiments on MIMIC-III and MIMIC-CXR

2 Debiasing Algorithms

Classification Parity Proxies

Let $\mathcal{X} = \{x_i\}_{i=1}^N, \mathcal{Y} = \{y_i\}_{i=1}^N, \mathcal{A} = \{a_i\}_{i=1}^N$. Differentiable proxies for classification disparity:

$$\tilde{\mu}_{\text{SPD}}(f_\theta(\cdot), \mathcal{X}, \mathcal{Y}, \mathcal{A}) = \frac{\sum_{i=1}^N f_\theta(x_i)(1 - a_i)}{\sum_{i=1}^N (1 - a_i)} - \frac{\sum_{i=1}^N f_\theta(x_i)a_i}{\sum_{i=1}^N a_i}$$

$$\tilde{\mu}_{\text{EOD}}(f_\theta(\cdot), \mathcal{X}, \mathcal{Y}, \mathcal{A}) = \frac{\sum_{i=1}^N f_\theta(x_i)(1 - a_i)y_i}{\sum_{i=1}^N (1 - a_i)y_i} - \frac{\sum_{i=1}^N f_\theta(x_i)a_i y_i}{\sum_{i=1}^N a_i y_i}$$

Observation: $\tilde{\mu}_{\text{SPD}}(f_\theta(\cdot), \mathcal{X}, \mathcal{Y}, \mathcal{A}) \propto \widehat{\text{Cov}}(A, f_\theta(X))$ and $\tilde{\mu}_{\text{EOD}}(f_\theta(\cdot), \mathcal{X}, \mathcal{Y}, \mathcal{A}) \propto \widehat{\text{Cov}}(A, f_\theta(X)|Y = 1)$

Intuition: minimise/maximise a differentiable proxy directly

Pruning for Debiasing

i. for layer $1 \leq l \leq L$, evaluate influence

$$S_{l,j} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \tilde{\mu}(f_\theta(\cdot), \mathcal{X}, \mathcal{Y}, \mathcal{A})}{\partial z_j^l(x_i)}$$

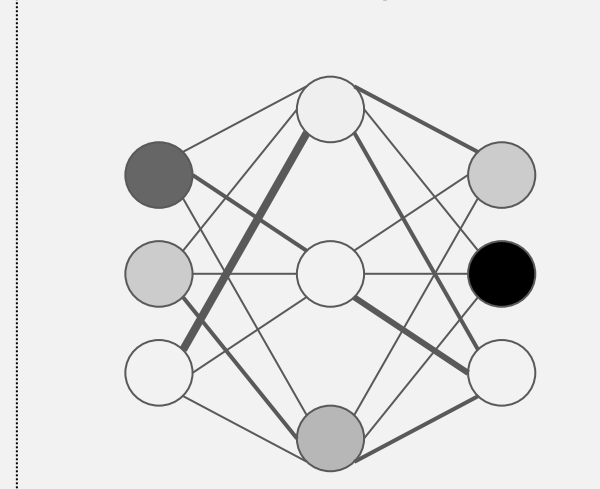
ii. prune the most influential units

iii. evaluate bias and performance on $\mathcal{D}_{\text{valid}}$

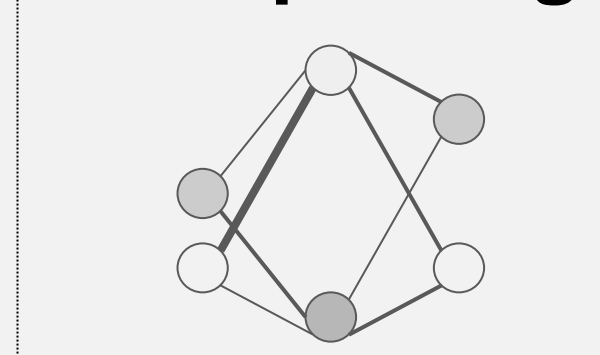
iv. repeat steps i-iii

Return the model with the minimal bias and performance $> \varrho$

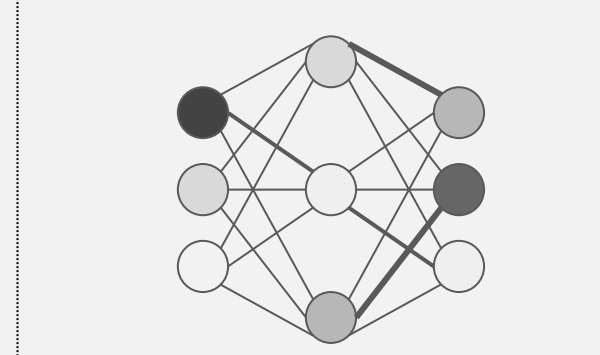
biased $f_\theta(\cdot)$



after pruning



after bias GD/A

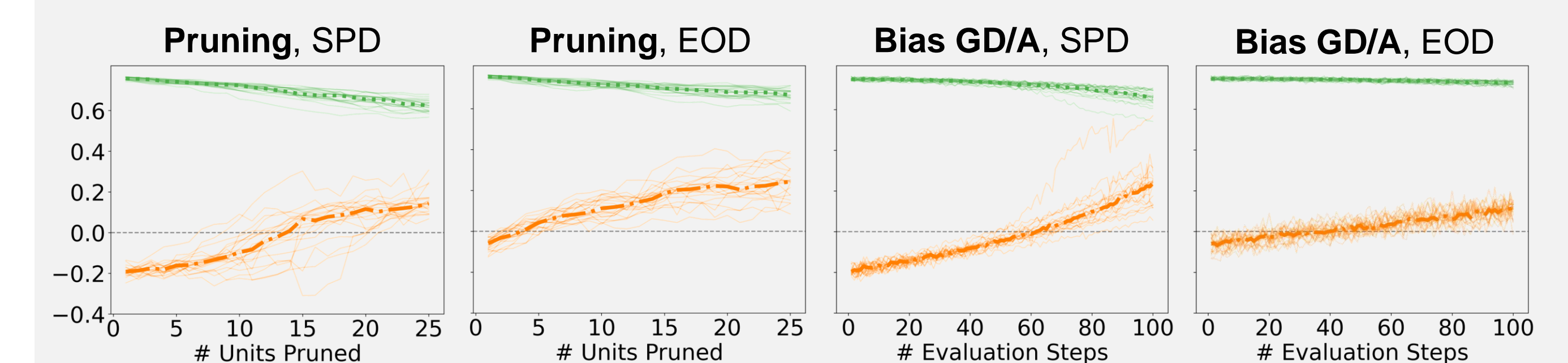


Bias Gradient Descent/Ascent (GD/A)

- fine-tune $f_\theta(\cdot)$ on $\mathcal{D}_{\text{valid}}$
 - few epochs, small LR
- mini-batch gradient descent/ascent
- stop early: performance $> \varrho$

3 Results

Bias and performance when debiasing on MIMIC-III:



Results for the VGG-16 and ResNet-18 trained on MIMIC-CXR:

Enlarged cardiomeastinum, Sex, VGG-16

Method	EOD	BA
STANDARD	-0.05±0.02	0.77±0.01
RANDOM	-0.03±0.03	0.75±0.01
ROC	-0.05±0.02	0.75±0.03
Eq. Odds	0.01±0.03	0.75±0.01
ADV. INTRA	-0.04±0.03	0.73±0.01
PRUNING	0.00±0.02	0.76±0.02
BIAS GD/A	-0.01±0.04	0.76±0.01

Pneumonia, Ethnicity, VGG-16

Method	EOD	BA
STANDARD	-0.14±0.04	0.73±0.02
RANDOM	-0.11±0.06	0.71±0.02
ROC	-0.07±0.06	0.65±0.06
Eq. Odds	0.00±0.06	0.70±0.01
ADV. INTRA	-0.13±0.05	0.70±0.02
PRUNING	-0.09±0.05	0.71±0.03
BIAS GD/A	-0.08±0.06	0.71±0.02

Enlarged cardiomeastinum, Sex, ResNet-18

Method	EOD	BA
STANDARD	-0.05±0.04	0.76±0.01
RANDOM	0.00±0.03	0.73±0.02
ROC	-0.05±0.03	0.74±0.04
Eq. Odds	0.01±0.03	0.74±0.01
ADV. INTRA	-0.04±0.04	0.73±0.02
PRUNING	-0.01±0.03	0.74±0.02
BIAS GD/A	0.00±0.03	0.76±0.01

Pneumonia, Ethnicity, ResNet-18

Method	EOD	BA
STANDARD	-0.14±0.05	0.73±0.02
RANDOM	-0.06±0.06	0.65±0.04
ROC	-0.07±0.04	0.65±0.05
Eq. Odds	-0.01±0.06	0.70±0.01
ADV. INTRA	-0.14±0.03	0.71±0.02
PRUNING	-0.11±0.05	0.70±0.02
BIAS GD/A	-0.11±0.05	0.73±0.02

4 Outlook

- Other diseases and protected attributes in MIMIC-CXR
- Other architectures, e.g. SqueezeNet, DenseNet
- Beyond SPD and EOD
- Multiple and multicategorical protected attributes and labels
- Beyond gradient-based criteria for pruning