



南京邮电大学  
Nanjing University of Posts and Telecommunications

# Intra-Processing Methods for Debiasing Neural Networks

Yash Savani  
Abacus.AI  
San Francisco, CA 94103  
yash@abacus.ai

Colin White  
Abacus.AI  
San Francisco, CA 94103  
colin@abacus.ai

Naveen Sunda  
Govindarajulu RAIR Lab RPI  
Troy, NY 12180  
naveensundarg@gmail.com



汇报人：王亦恺（研一）



指导老师：李云

Advances in Neural Information Processing Systems 33 (NeurIPS 2020)



# 目录

CONTENTS

- 1 研究背景
- 2 基本思想
- 3 具体算法
- 4 实验结果

# 1.研究背景

机器学习在社会各领域广泛应用的同时，引发了人们对机器学习可信性的担忧，其中包括公平性。研究发现，分类器在依据性别、种族、宗教等敏感属性划分的不同群体上预测同一目标属性时，预测的准确率存在明显差异。这种差异在一些特殊的应用场景下可能被进一步放大，从而造成严重的社会问题。

例如，美国政府目前正在使用的某一面部识别系统在预测亚裔或非裔人群时，预测结果的误判率达到了欧裔人群的10倍甚至更高<sup>[1]</sup>，这可以认为是一种偏见或歧视现象。如果将该工具运用于执法系统，意味着亚裔或非裔人群更有可能被无辜逮捕。



# 1.研究背景

## 公平性算法

- **预处理方法(Pre-processing)**

通常认为，数据集的不平衡是偏见的直接来源。预处理方法从数据集出发，通过对原始数据进行去偏操作，得到相对平衡的数据集。

- **处理中方法(In-processing)**

处理中方法通过在模型训练过程中添加公平性约束项，消除源于数据的偏差。

- **后处理方法(Post-processing)**

后处理方法直接对模型的输出进行干预，从而实现公平性。因此适用于整个模型为黑盒的应用场景。



# 1.研究背景

## 方法提出

现实许多应用场景下，往往会使用一个预训练好的大模型，再根据具体应用场景对模型微调。选择何种方法合适？

A:预处理方法

B:处理中方法

C:后处理方法

需要从头训练整个模型

无法充分发挥模型的全部性能

**内处理方法(intro-processing):**  
微调模型参数以实现公平性



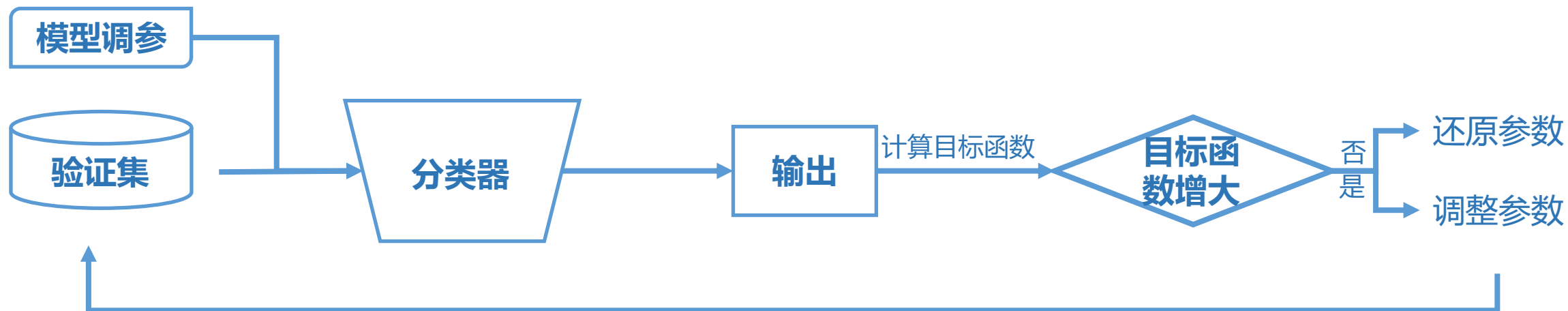


# 目录

CONTENTS

- 1 研究背景
- 2 基本思想
- 3 具体算法
- 4 实验结果

## 2.基本思想



最大化目标函数：

$$\max \phi_{\mu, \rho, \epsilon}(\mathcal{D}, \hat{\mathcal{Y}}, A) = \begin{cases} \rho & \text{if } \mu < \epsilon \\ 0 & \text{otherwise} \end{cases}$$

**函数输入：**

$\mathcal{D}$ : 数据集

$\hat{\mathcal{Y}}$ : 预测结果

$A$ : 受保护属性

**函数输出：**

$\rho$ : 评估模型性能的函数

**其他：**

$\mu$ : 评估模型公平性的函数

$\epsilon$ : 公平性指标的阈值



## 2.基本思想

### 公平性指标

奇偶校验差(Statistical Parity Difference, SPD):

$$SPD(\mathcal{D}, \hat{Y}, A) = P_{(x_i, Y_i) \in \mathcal{D}}(\hat{Y}_i = 1 \mid a_i = 0) - P_{(x_i, Y_i) \in \mathcal{D}}(\hat{Y}_i = 1 \mid a_i = 1)$$

平等机会差异(Equal opportunity difference, EOD):

$$EOD(\mathcal{D}, \hat{Y}, A) = TPR_{A=0}(\mathcal{D}, \hat{Y}) - TPR_{A=1}(\mathcal{D}, \hat{Y})$$

平均机会差异(Average Odds Difference, AOD):

$$AOD(\mathcal{D}, \hat{Y}, A) = \frac{(FPR_{A=0}(\mathcal{D}, \hat{Y}) - FPR_{A=1}(\mathcal{D}, \hat{Y})) + (TPR_{A=0}(\mathcal{D}, \hat{Y}) - TPR_{A=1}(\mathcal{D}, \hat{Y}))}{2}$$

图像数据:  $\mu(\mathcal{D}, \hat{Y}, A) = AOD(\mathcal{D}, \hat{Y}, A)$

表格数据:  $\mu(\mathcal{D}, \hat{Y}, A) = SPD(\mathcal{D}, \hat{Y}, A)$





## 2.基本思想

### 性能评估指标

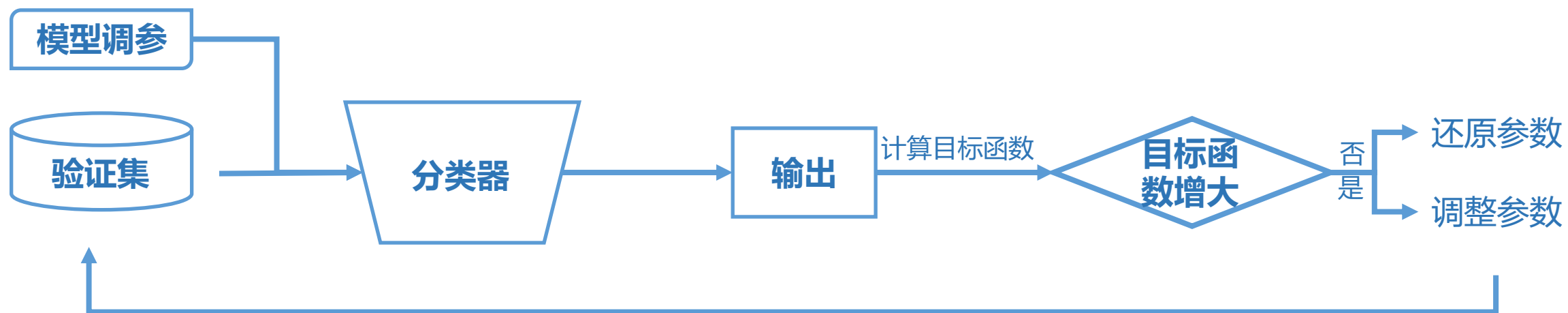
常见的性能评估指标：准确率（Accuracy）、精确率（Precision）、召回率（Recall）、AUC-ROC曲线

平衡精度(balanced accuracy):

$$\rho(y, \hat{y}) = \frac{TPR + TNR}{2}$$



## 2.基本思想



最大化目标函数：

$$\max \phi_{\mu, \rho, \epsilon}(\mathcal{D}, \hat{\mathcal{Y}}, A) = \begin{cases} \rho & \text{if } \mu < \epsilon \\ 0 & \text{otherwise} \end{cases}$$





# 目录

CONTENTS

- 1 研究背景
- 2 基本思想
- 3 具体算法
- 4 实验结果

## 3.1 随机扰动算法

---

### Algorithm 1 Random Perturbation

---

- 1: **Input:** Trained model  $f$  with weights  $\theta$ , validation dataset  $\mathcal{D}_{\text{valid}}$ , objective  $\phi_{\mu, \rho, \epsilon}$ , parameter  $T$
  - 2: Set  $\theta^* = \theta$ ,  $\text{val}^* = -\infty$ , and  $\tau^* = 0$
  - 3: **for**  $i = 1$  to  $T$  **do**
  - 4:   Sample  $q_j \sim \mathcal{N}(1, 0.1)$  for all  $j \in \{1, 2, \dots, |\theta|\}$
  - 5:    $\theta'_j = \theta_j \cdot q_j$
  - 6:   Select threshold  $\tau \in [0, 1]$  which maximizes the objective  $\phi_{\mu, \rho, \epsilon}$  on the validation set
  - 7:   Set  $\text{val} = \phi_{\mu, \rho, \epsilon}(\mathcal{D}_{\text{valid}}, \{\mathbb{I}\{f_{\theta'}(\mathbf{x}) > \tau\} \mid (\mathbf{x}, Y) \in \mathcal{D}_{\text{valid}}\}, A)$
  - 8:   If  $\text{val} > \text{val}^*$ , set  $\text{val}^* = \text{val}$ ,  $\theta^* = \theta'$ , and  $\tau^* = \tau$ .
  - 9: **end for**
  - 10: **Output:**  $\theta^*, \tau^*$
- 

从均值为1方差为0.1的高斯分布中随机得到 $|\theta|$ 个随机数，用这些随机数扰动模型参数。



## 3.2层级优化



核心思想：迭代多棵回归树共同决策最终结果

优点：属于零阶优化器，适用于目标函数不可微的情况，因此又称为“黑盒优化器”

缺点：需要进行多次迭代，因此调参过程较为复杂，训练时间可能较长

## 3.2 层级优化

梯度增强回归树  
(GBRT, gradient-  
boosted regression  
trees)

使用优化器调整  
每一层的参数，  
而不是使用随机  
扰动

---

### Algorithm 2 Layer-wise optimization

---

- 1: **Input:** Trained model  $f = f^{(\ell)} \circ \dots \circ f^{(1)}$  with weights  $\theta_1, \dots, \theta_\ell$ , objective  $\phi_{\mu, \rho, \epsilon}$ , black-box optimizer  $\mathcal{A}$
  - 2: Set  $\theta^* = \emptyset$ ,  $\text{val}^* = -\infty$ , and  $\tau^* = 0$
  - 3: **for**  $i = 1$  to  $\ell$  **do**
  - 4:   Run optimizer  $\mathcal{A}$  to optimize weights  $\theta_i$  to  $\theta'_i$  with respect to  $\phi_{\mu, \rho, \epsilon}$ .
  - 5:   Select threshold  $\tau \in [0, 1]$  which maximizes objective  $\phi_{\mu, \rho, \epsilon}$
  - 6:   Set  $\text{val} = \phi_{\mu, \rho, \epsilon}(\mathcal{D}_{\text{valid}}, \{\mathbb{I}\{f_{\theta'}(\mathbf{x}) > \tau\} \mid (\mathbf{x}, Y) \in \mathcal{D}_{\text{valid}}\}, A)$ , where  $\theta' = (\theta_1, \dots, \theta'_i, \dots, \theta_\ell)$
  - 7:   If  $\text{val} > \text{val}^*$  set  $\text{val}^* = \text{val}$ ,  $\theta^* = \theta'$ , and  $\tau^* = \tau$ .
  - 8: **end for**
  - 9: **Output:**  $\theta^*, \tau^*$
- 



## 3.3 对抗优化

### Algorithm 3 Adversarial Fine-Tuning

```
1: Input: Trained model  $f = f^{(\ell)} \circ f'$  with weights  $\theta$ , validation dataset  $\mathcal{D}_{\text{valid}}$ , parameters  $\lambda, \epsilon, \delta, n, m, m', T$ .
2: Set  $g$  as the critic model with weights  $\psi$ .
3: for  $i = 0$  to  $n$  do
4:   for  $j = 0$  to  $m$  do
5:     Sample a minibatch  $(\mathbf{X}_k, \mathbf{Y}_k)$  with replacement from  $\mathcal{D}_{\text{valid}}$ 
6:     Evaluate the bias in the minibatch,  $\bar{\mu} \leftarrow \mu((\mathbf{X}_k, \mathbf{Y}_k), f(\mathbf{X}_k))$ .
7:     Update the critic model  $g$  by updating its stochastic gradient
       
$$\nabla_{\psi}(\bar{\mu} - (g \circ f')(\mathbf{X}_k))^2$$

8:   end for
9:   for  $j = 0$  to  $m'$  do
10:    Sample a minibatch  $(\mathbf{X}_k, \mathbf{Y}_k)$  with replacement from  $\mathcal{D}_{\text{valid}}$ 
11:    Update the original model by updating its stochastic gradient
       
$$\nabla_{\theta} [\max\{1, \lambda \cdot (|(g \circ f')(\mathbf{X}_k)| - \epsilon + \delta) + 1\} \cdot \text{BCELoss}(\mathbf{Y}_k, f(\mathbf{X}_k))]$$

12:   end for
13:   Select threshold  $\tau \in [0, 1]$  that minimizes the objective  $\phi_{\mu, \rho}$ 
14: end for
15: Output: Debaised model  $f$ , threshold  $\tau$ 
```

批评模型：评判某一批次训练结果中的偏见

训练批评模型

训练预测模型



### 3.3 对抗训练

$$\max\{1, \lambda \cdot (|(g \circ f')(X_k)| - \epsilon + \delta) + 1\} \cdot \text{BCELoss}(\mathbf{y}_k, f(\mathbf{X}_k))$$



$$\mu(\mathcal{D}, \hat{\mathcal{Y}}, A)$$



对于一批数据，  
若符合公平性约束，则损失函数的系数为1；  
若不符合公平性约束，则提升惩罚强度。







# 目录

CONTENTS

- 1 研究背景
- 2 基本思想
- 3 具体算法
- 4 实验结果

## 4实验结果

数据集: Adult、

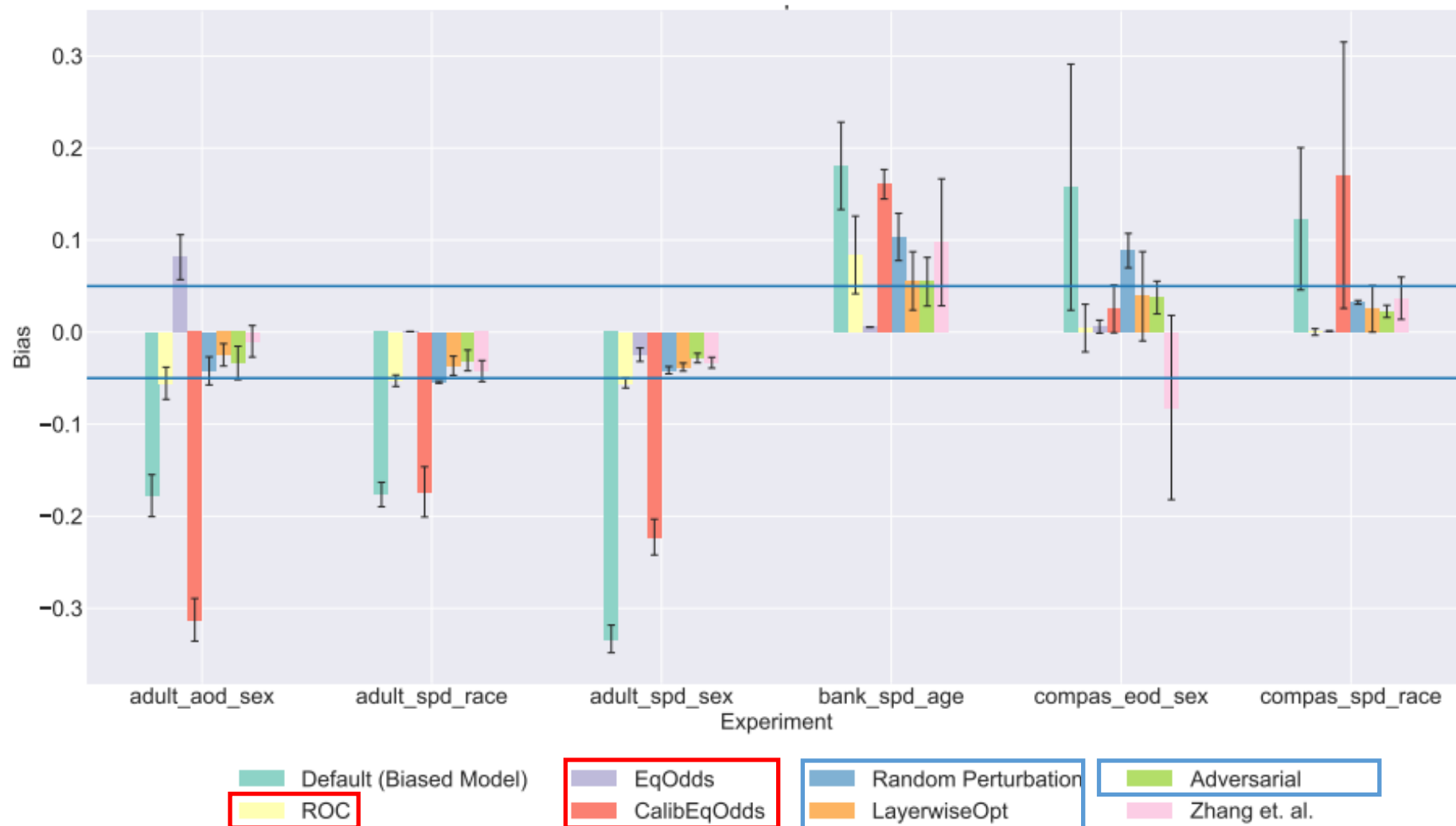
Bank、Compas

AOD:平均机会差异

SPD:奇偶校验差

EOD:平等机会差异

结论: 在公平性上,  
本文提出的内处理  
方法整体上效果**优**  
**于后处理方法**; 与  
**处理中方法效果接**  
**近。**



## 4实验结果

结论：在模型性能上，本文提出的内处理方法效果优于后处理方法；略低于处理中方法。

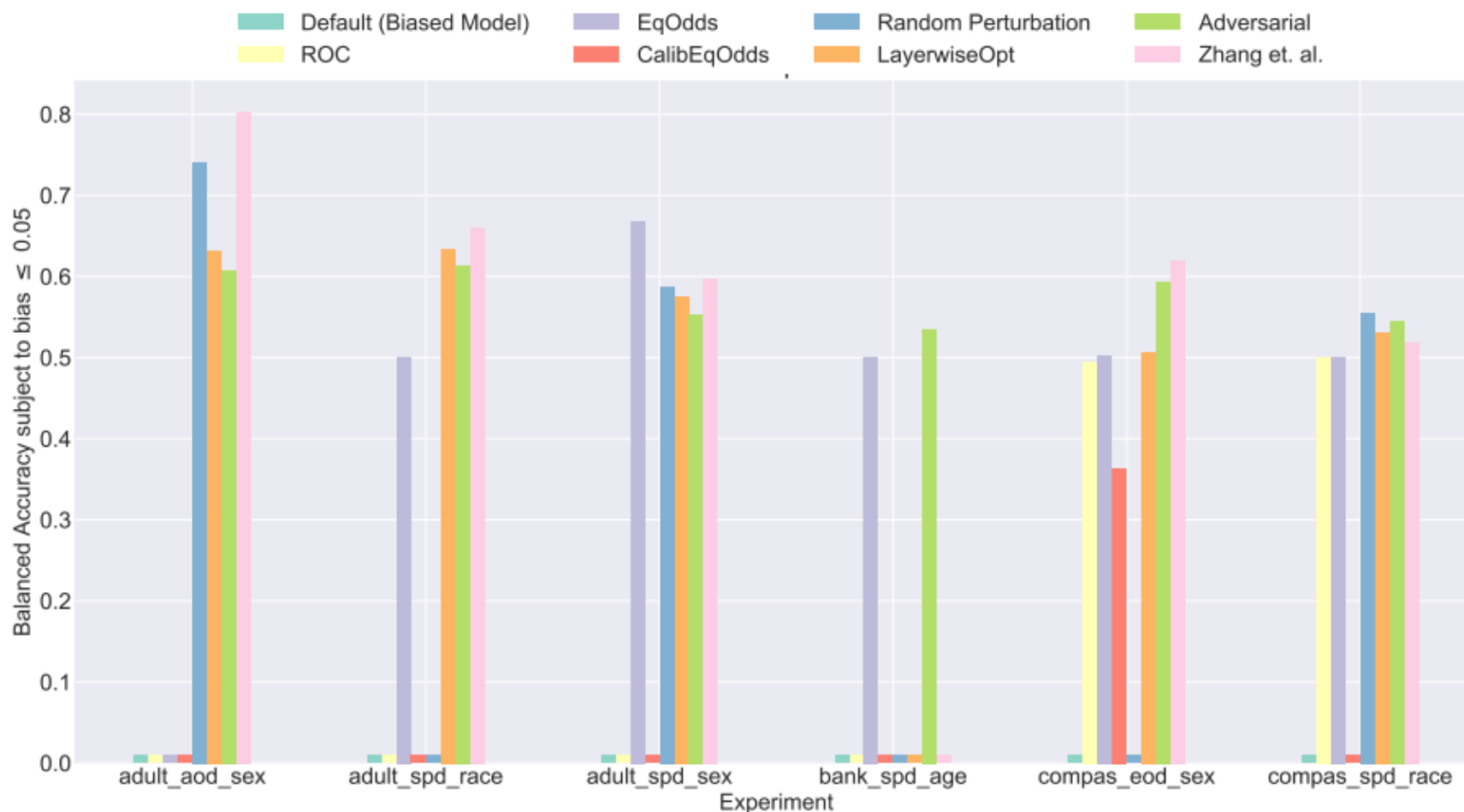


Figure 1: Results for tabular datasets over 5 runs with different seeds. mean Bias with std error bars (Top) and the median of the objective function in Equation 1 (Bottom).



## 4实验结果

数据集: CelebA

Table 1: Results on the CelebA datasets for a pretrained ResNet with three initial random seeds. Results are the balanced accuracy scores after fine-tuning. The crossed out scores are those that did not have biases lower than 0.05.

	Default	ROC	EqOdds	CalibEqOdds	Random	LayerwiseOpt	Adversarial
1	<del>0.533</del>	0.533	0.983	0.519	<del>0.567</del>	<del>0.530</del>	0.914
2	<del>0.523</del>	0.521	0.983	0.487	0.529	0.508	0.917
3	0.535	0.533	0.982	0.514	<del>0.591</del>	0.529	0.905

结论：在实现公平性的模型中（本文的标准是公平性指标小于阈值0.05），Baseline模型、后处理方法得到的模型、随机扰动算法和层级优化方法得到的模型在性能上相似，通过对抗优化微调得到的模型在性能上优于其他方法得到的模型。在更复杂的模型与数据上，更复杂的方法更加有效。



## 4实验结果

数据集：CelebA

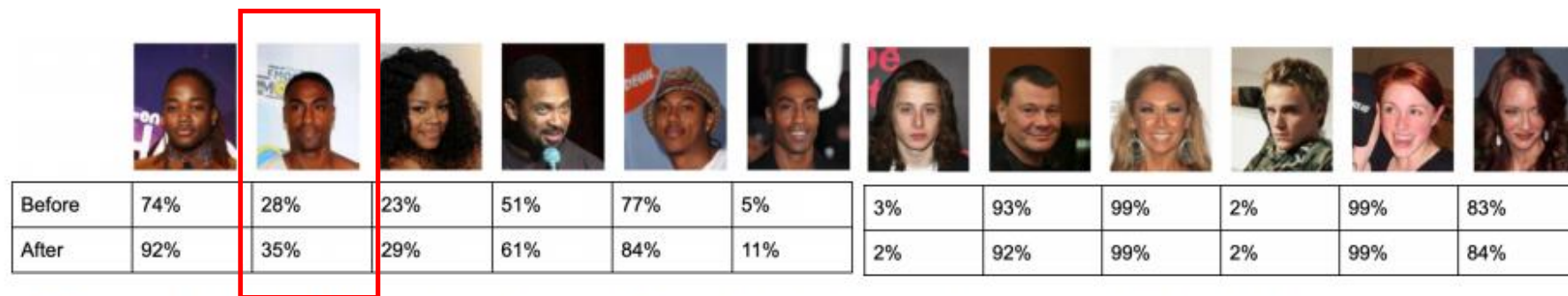


Figure 2: Probability of smiling on the CelebA dataset, before and after debiasing w.r.t. race.

结论：Baseline模型在预测目标是否微笑时，预测概率会受其种族因素影响而不平衡。经过对抗优化微调得到的模型有效提升了模型的公平性。

# 总结

- 提出了一种介于处理中方法与后处理方法之间的黑盒模型公平性方法——“内处理方法”，实验证明该方法能够有效弥补后处理方法导致模型性能下降的缺点，并具有现实应用价值。
- 对内处理方法的性质进行研究，发现公平性问题与模型的初始条件密切相关。
- 通过调整超参数，能够将一些处理中方法转化为内处理方法。





南京邮电大学  
Nanjing University of Posts and Telecommunications

敬请各位老师批评指正

---