

# Counterfactual Interpolation Augmentation (CIA): A Unified Approach to Enhance Fairness and Explainability of DNN

Huo Mingda

Jinan University, Guangzhou

May 11, 2023

利用因果理论中的反事实（counterfactual）  
框架来提高算法的稳定性和可解释性。

Counterfactual  
Interpolation  
Augmentation

去除敏感信息，减轻预测偏差

遵循因果关系

虚假相关性

偏置缓解技术

伪关系(spurious relation)



解释方法：

$$\text{Argmax}_E Q(E|Human, Data, Task)$$

预处理：

通过数据增强来消除偏差和提高训练集的质量

$$\text{Argmax}_{E, Model} Q(E|Model, Human, Data, Task)$$

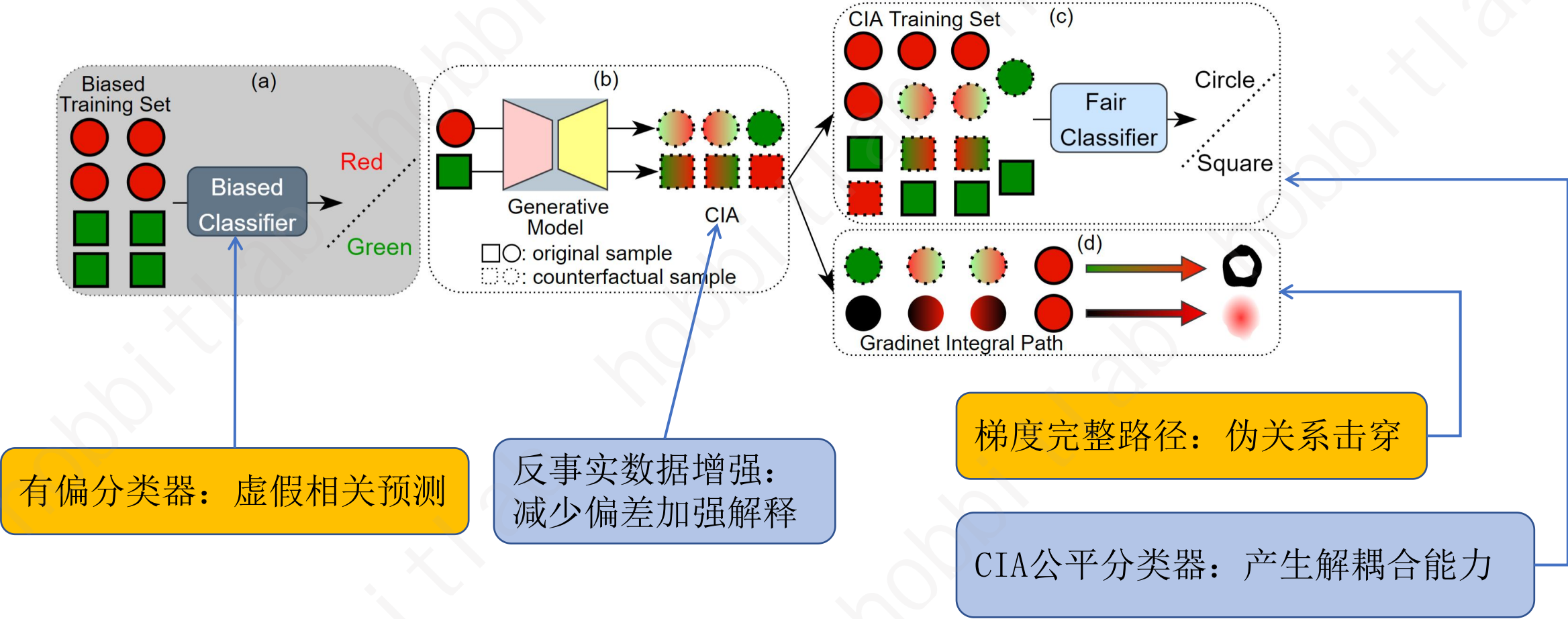
内处理：

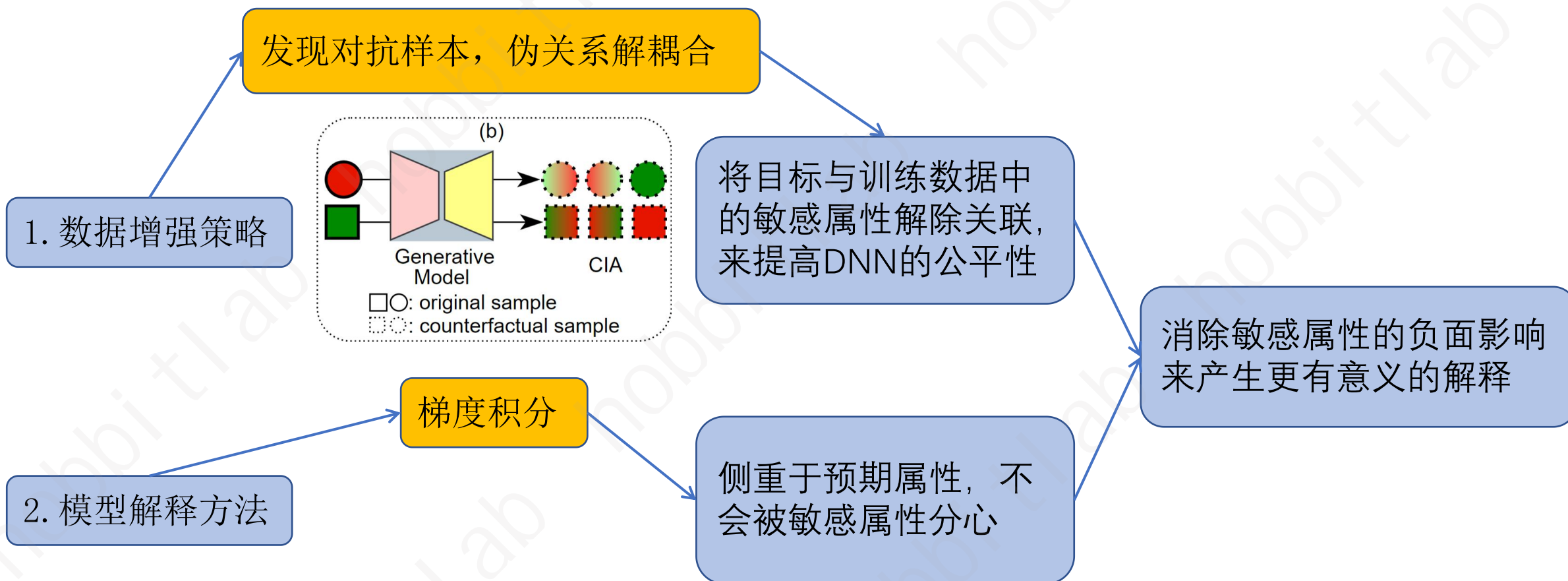
在训练过程中从学习到的特征中去除敏感信息

后处理：

根据推理时的敏感属性来校准或修改预测

# 研究内容：公平性感知预测





积分梯度

$$\text{IG}_i(x) = (x_i - x'_i) \cdot \int_0^1 \frac{\partial F(x'_i + \alpha \cdot (x_i - x'_i))}{\partial x_i} d\alpha$$

黎曼近似求和得到线性  
路径上的积分梯度

偏导数得到网  
格灵敏度

基线直接影响归因质量

(1)最大距离基线

(2)模糊基线

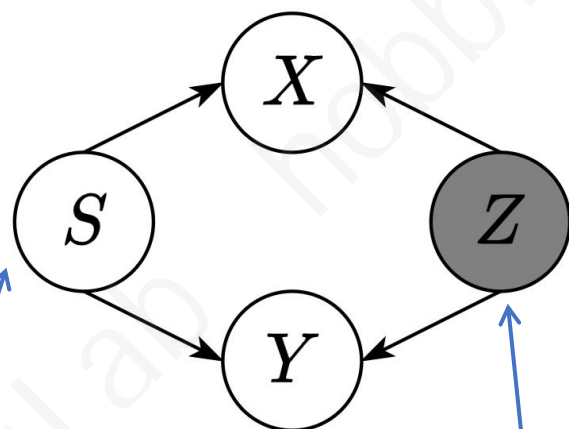
(3)高斯基线

(4)统一基线

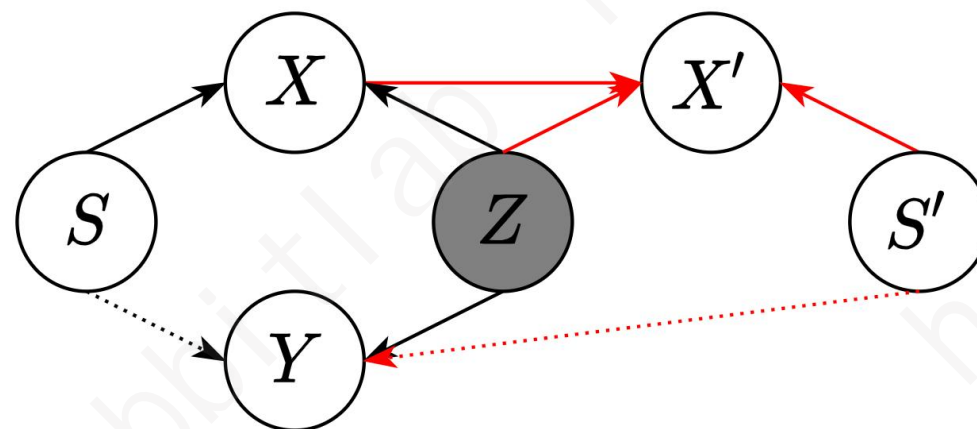
反事实插值路径  
积分梯度

可解释性

制造混淆因子 $S'$  产生反事实样例 $X'$



(a) Original causal inference



(b) Counterfactual causal inference

敏感属性

潜在因素

混淆因子 $S'$

## 反事实因果推理

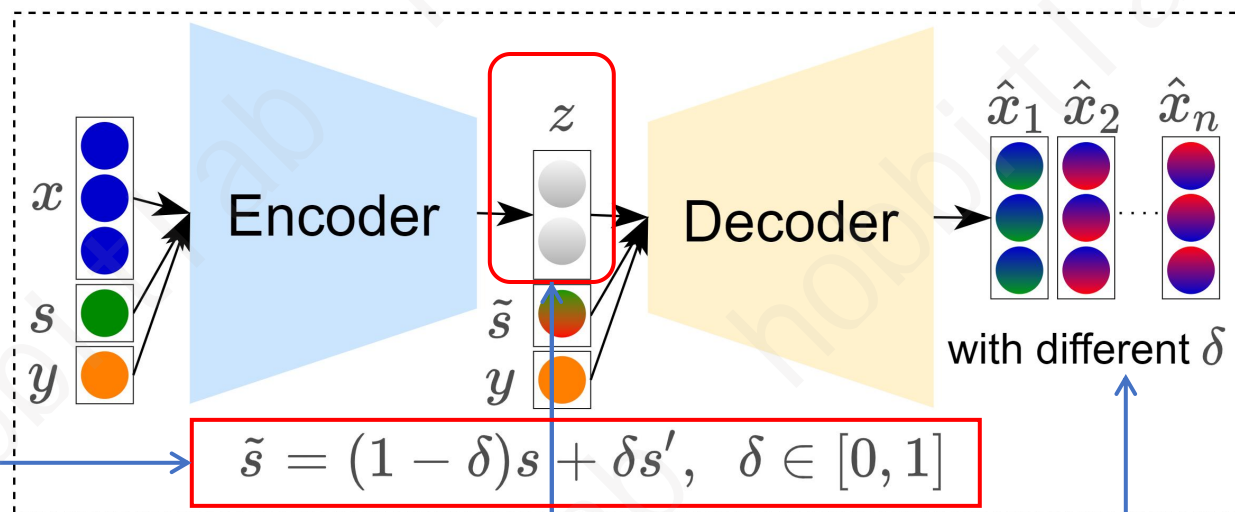
$$p(\hat{Y}_{S \leftarrow s} = y | X = x, Y = y, S = s) \\ = p(\hat{Y}_{S \leftarrow s'} = y | X = x, Y = y, S = s').$$

因果推理的视角来解决训练数据偏差问题

以事实目标 $y$ 为条件，对模型的反事实公平提出标准

因果推理的视角来解决训练数据偏差问题

从事实例子 $x$ 过渡到它的反事实例子 $x'$



学习潜在特征

条件解码：  
从特征 $Z$ 和 $Y$   
到输入空间

插值路径， $\delta$  决定插值数量



$$\text{CGI}_i(x) = (x_i - x'_i) \int_{\delta=0}^1 \frac{\partial F(\gamma(\delta))}{\partial \gamma_i(\delta)} \frac{\partial \gamma_i(\delta)}{\partial \delta} d\delta.$$

$$\gamma(\alpha) = x' + \alpha(x - x')$$

路径积分的方向：从反事实样本  $x'$  过渡到输入  $x$

通过CIA生成的反事实插值

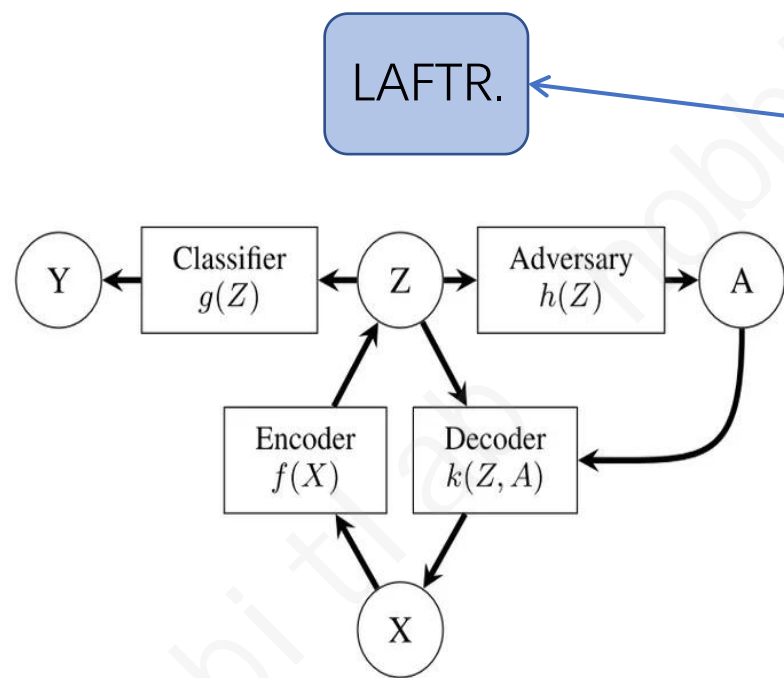


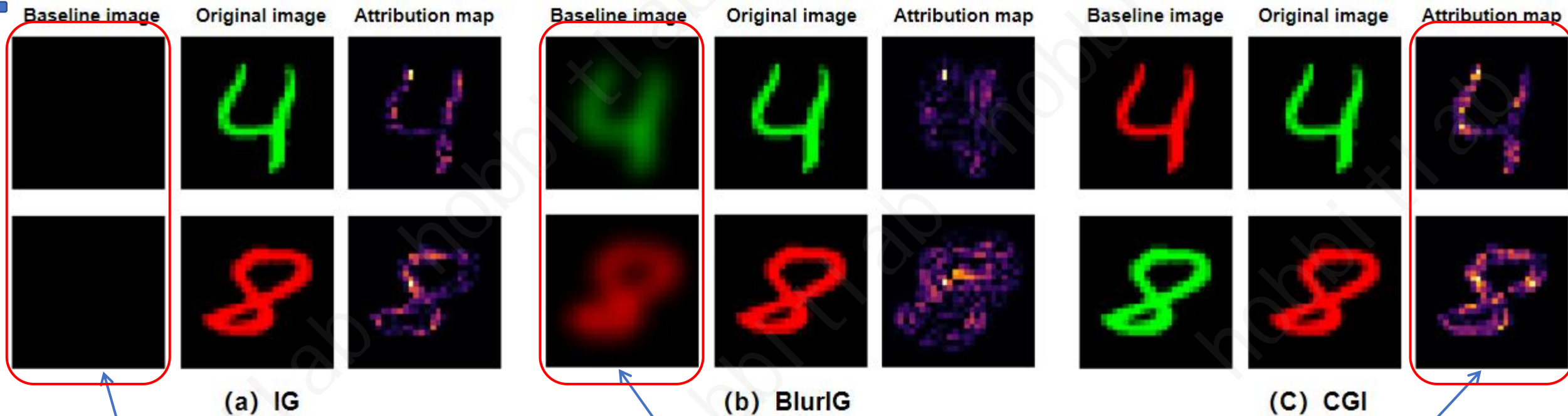
Figure 1. Model for learning adversarially fair representations. The variables are data  $X$ , latent representations  $Z$ , sensitive attributes  $A$ , and labels  $Y$ . The encoder  $f$  maps  $X$  (and possibly  $A$  - not shown) to  $Z$ , the decoder  $k$  reconstructs  $X$  from  $(Z, A)$ , the classifier  $g$  predicts  $Y$  from  $Z$ , and the adversary  $h$  predicts  $A$  from  $Z$  (and possibly  $Y$  - not shown).

Method	Training Acc	Test Acc
Vanilla	79.48	18.08
LAFTR	74.14	75.22
Prior Training	74.62	75.46
CIA-10	79.64	78.16
CIA-20	79.95	78.23
CIA-30	79.97	78.69

错误依赖颜色和数字的虚假相关性

CIA: 样本的反事实插值越多准确度越高

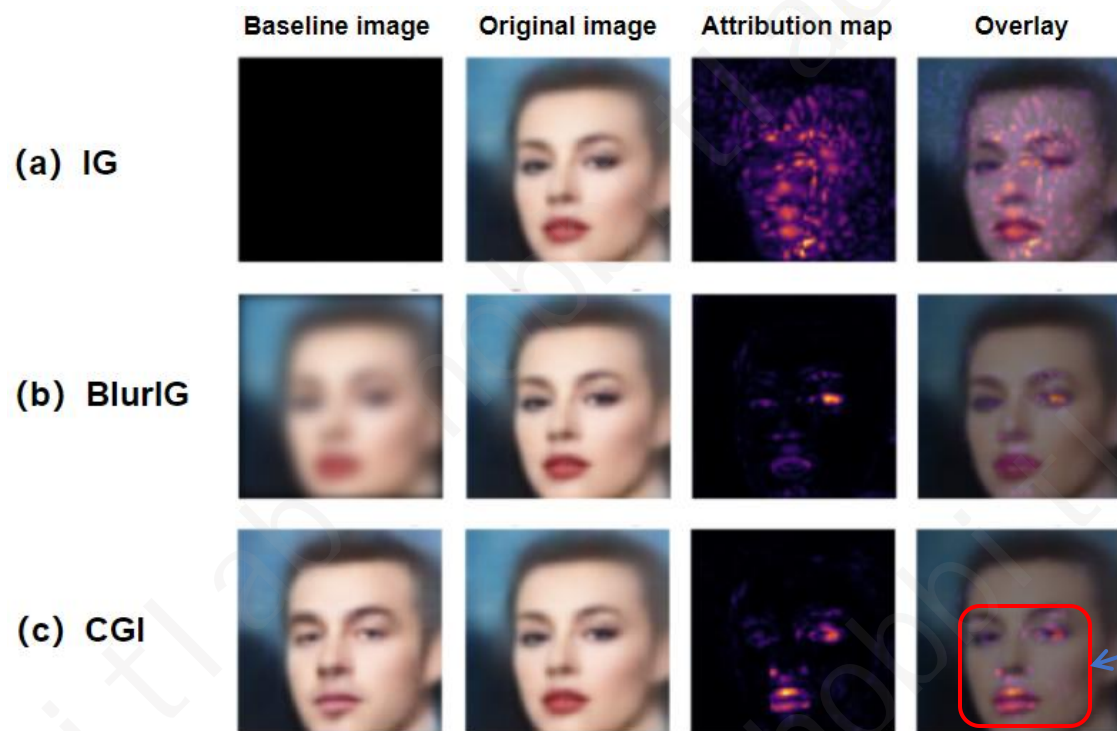
Vanilla: 训练过程中无法学习数字形状



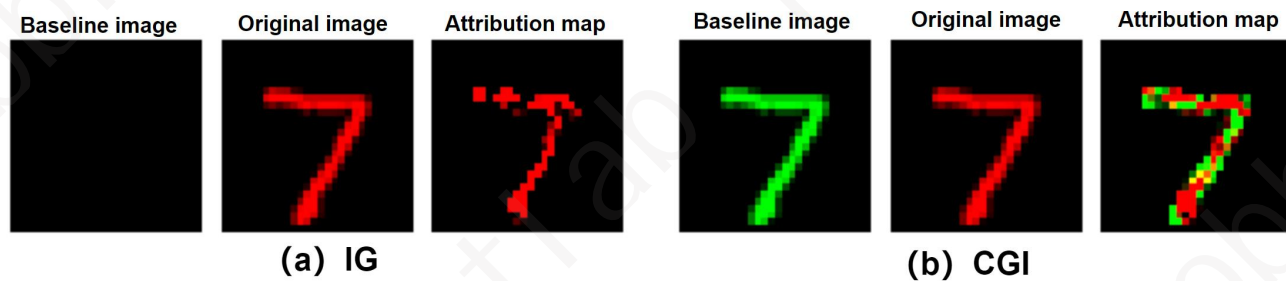
IG积分梯度算法使用黑色图像作为梯度积分的基线

模糊IG:BlurIG通过连续模糊原始输入来消除路径积分

反事实梯度积分算法CGI具备清晰的归因热图



归因属性更准确



可解释公平性

**Thanks!**