# Latent Space Smoothing for **Individually Fair Representations**

Momchil Peychev, Anian Ruoss, Mislav Balunović, Maximilian Baader, Martin Vechev

# Individual Fairness for High-Dimensional Data

### Problem Setting

Individual fairness: treat similar individuals similarly. In the context of:

- **Data Regulator:** defines a fairness notion for the task  $\bullet$
- Data Producer: learns a fair representation of the data
- **Data Consumers:** make predictions from the representation

#### This work:

i. Design a suitable input similarity metric

ii. Enforce that similar individuals are *provably* treated similarly for high-dimensional data and real-world models.

#### Data Regulator

For a given person, all people differing only in skin tone should be classified the same.

Key challenge: high-level semantic attributes cannot be captured conveniently in the input space of the original data  $x \in \mathbb{R}^n$ .

### Key Insight: Similarity via a Generative Model

Leverage (invertable) Glow G = (E, D) and compute attribute vector  $a \in \mathbb{R}^q$  in the latent space of G. For  $x \in \mathbb{R}^n$  and  $z_G = E(x) \in \mathbb{R}^q$ :

Similarity set

- In the latent space of  $G: S(\mathbf{x}) \coloneqq \{\mathbf{z}_G + t \cdot \mathbf{a} \mid |t| \le \epsilon\} \subseteq \mathbb{R}^q$
- In the input space  $\mathbb{R}^n$  :  $S^{in}(\mathbf{x}) \coloneqq D(S(\mathbf{x})) \subseteq \mathbb{R}^n$

**Goal:** construct an end-to-end model  $M: \mathbb{R}^n \to \mathcal{Y}$  for which, given an input *x*, we can **certify** individual fairness of the form

$$\forall \mathbf{x}' \in S^{\text{in}}(\mathbf{x}) \colon M(\mathbf{x}) = M(\mathbf{x}')$$



## Latent Space Smoothing for Individually Fair **Representations (LASSI)**

#### Data Producer

**1.** Glow encoder *E* and a representation  $R: \mathbb{R}^q \to \mathbb{R}^k$  trained with: Adversarial loss: map similar individuals close together in  $\mathbb{R}^k$ 

$$\mathcal{L}_{adv}(\boldsymbol{x}) = \max_{\boldsymbol{z} \in S(\boldsymbol{x})} \|R(\boldsymbol{z}_G) - R(\boldsymbol{z})\|_2$$

**Classification loss:** ensure utility for downstream tasks

 $\mathcal{L}_{cls}(\mathbf{x}, \mathbf{y}) = \text{cross\_entropy}(C_{a})$ 

**Reconstruction loss:** for transfer learning, preserve original signal

 $\mathcal{L}_{recon}(\mathbf{x}) = \left\| \mathbf{z}_{G} - Q(R(\mathbf{z}_{G})) \right\|_{2} \text{ with } Q \colon \mathbb{R}^{k} \to \mathbb{R}^{q}$ 

R,  $C_{aux}$  and Q are trained jointly, trading off fairness, accuracy and transferability:  $L = \lambda_1 \mathcal{L}_{cls}(\mathbf{x}, \mathbf{y}) + \lambda_2 \mathcal{L}_{adv}(\mathbf{x}) + \lambda_3 \mathcal{L}_{recon}(\mathbf{x})$ .

2. Adversarial training  $\rightarrow$  no formal guarantees. Center smoothing on R,  $\hat{R}(\mathbf{z}_G)$ , provably mapping similar individuals close together:

 $\boldsymbol{r}_{cs}, \boldsymbol{d}_{cs} \leftarrow \widehat{R}(\boldsymbol{z}_{G}) \text{ and } \forall \boldsymbol{z} \in S(\boldsymbol{x}): \| \boldsymbol{z} \in S(\boldsymbol{x}) \|$ 

with probability at least  $1 - \alpha_{cs}$  [1].

Data Consumer

Randomized smoothing [2] (with confidence  $\alpha_{rs}$ ) on the downstream classifier  $C: \mathbb{R}^k \to \mathcal{Y}, \hat{C}(\boldsymbol{r}_{cs})$ , to obtain its  $\ell_2$ -robustness radius  $d_{rs}$ around  $r_{cs}$ . The classifier C is trained after training R is completed.

End-to-end Fairness Certificate

Given input  $x \in \mathbb{R}^n$ , let  $z_G = E(x)$ 1.  $\boldsymbol{r}_{cs}, d_{cs} \leftarrow \hat{R}(\boldsymbol{z}_{G})$  and  $d_{rs} \leftarrow \hat{C}(\boldsymbol{r}_{cs})$ 2. If  $d_{cs} < d_{rs}$ , then **provably**  $\forall x' \in A$ with probability at least  $1 - \alpha_{cs} - \alpha_{cs}$ 

### Data Consumer





# https://github.com/eth-sri/lassi

$$aux \circ R(\mathbf{z}_G), y)$$

$$\boldsymbol{r}_{cs} - \hat{R}(\boldsymbol{z}) \big\|_2 \le d_{cs}$$

e for 
$$M = \widehat{C} \circ \widehat{R} \circ E$$

)  

$$S^{\text{in}}(\mathbf{x}): M(\mathbf{x}) = M(\mathbf{x}')$$
  
 $\alpha_{rs}$  for  $M = \hat{C} \circ \hat{R} \circ E$ .

# **Experimental Evaluation of LASSI**

Similarity set:  $a = z_{G,pos} - z_{G,neg}$  [3],  $\epsilon = 1$ 

Single & Multiple Sensitive Attributes

		Naive		DataAug		LASSI	
Task	Sensitive attribute(s)	Acc	Fair	Acc	Fair	Acc	Fair
Smiling	Pale Skin	86.3	0.6	85.7	12.2	85.9	98.0
	Young	86.3	38.2	85.9	43.0	86.3	98.8
	Blond Hair	86.3	3.4	86.6	9.4	86.4	94.7
	Heavy Makeup	86.3	0.4	85.3	13.7	85.6	91.3
	Pale Skin + Young	86.0	0.4	85.8	9.9	85.8	97.3
	Pale + Young + Blond	86.2	0.0	86.4	3.6	85.5	86.5

### Transfer Learning

Sens. attrib.:	Pale Skin		Pale +	Pale + Young		Pale+Young+Blond	
Transfer task	Acc	Fair	Acc	Fair	Acc	Fair	
Smiling	86.2	93.1	85.9	92.2	85.1	87.0	
HighCheeks	81.7	92.6	80.8	93.0	80.6	84.5	
MouthOpen	81.5	91.2	81.6	90.1	82.5	80.8	
Lipstick	88.3	94.0	85.1	90.6	86.2	81.0	
HeavyMakeup	86.5	93.0	83.7	90.0	83.3	80.4	
Wavy Hair	79.2	93.3	77.6	91.5	78.8	85.3	
Eyebrows	78.3	92.1	77.8	92.2	78.7	85.6	

The method is modular: LASSI can still achieve high certified individual fairness even when downstream tasks are not known.

#### Computing *a*

LASSI is *independent* of the actual computation of *a*. We instantiate it with four different attribute vector types in our paper [3,4,5,6].

#### **References:**

Kumar and Goldstein, Center Smoothing: Certified Robustness for N Cohen et al., Certified Adversarial Robustness via Randomized Smoo

Kingma and Dhariwal, Glow: Generative flow with invertible 1x1 cor

Denton et al., Detecting bias with generative counterfactual face at

Ramaswamy et al., Fair attribute classification through latent space

6. Li and Xu, Discover the unknown biased attribute of an image classif



Datasets: CelebA 64x64 (and FairFace – see paper), 312 test samples Naive:  $\lambda_1 = 1, \lambda_2 = \lambda_2 = 0$ , LASSI:  $\lambda_2 = 0.05$ , Transfer:  $\lambda_1 = 0, \lambda_3 = 0.1$ 

LASSI significantly increases the percentage of points for which we can certify individual fairness, without affecting the accuracy.

Networks with Structured Outputs	NeurIPS 20
othing	ICML 2019
nvolutions	NeurIPS 20
tribute augmentation	arXiv 2019
debiasing	CVPR 2021
ifier	ICCV 2021