

# Exercise 11

## Individual Fairness

Reliable and Trustworthy Artificial Intelligence  
ETH Zurich

Recall the standard notations from the lecture:

- $\phi: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \{0, 1\}$  – binary function specifying if its inputs,  $x$  and  $x'$ , are similar;
- $f_\theta: \mathbb{R}^n \rightarrow \mathbb{R}^k$  – encoder (data producer) mapping the input  $x$  to a learned representation  $f_\theta(x)$ ;
- $h_\psi: \mathbb{R}^k \rightarrow \mathbb{R}^o$  – classifier (data consumer) mapping the learnt representation to the logits corresponding to each label  $y \in \mathcal{Y}$ ;

**Problem 1** (Similarity Sets). One of the goals for the data producer  $f_\theta$  is to map similar individuals close to each other:

$$\phi(x, x') \implies \|f_\theta(x) - f_\theta(x')\|_\infty \leq \delta, \quad (1)$$

where  $\phi$  defines input similarity and  $\delta$  is a hyperparameter. Assuming that  $x$  is given (e.g., as a sample from the dataset), let  $\mathcal{L}(\phi \implies \omega)(x, x')$  be the standard DL2 loss of Formula (1) at  $x'$  (as discussed in the lectures). In this problem, you will consider different notions of input similarity.

1. Let  $x \in \mathbb{R}^d$  represent samples from the following dataset (only the first 2 rows and a subset of the features are shown below):

Name	Gender	Age	Salary	
Alice	Female	30	\$130,000	...
Bob	Male	32	\$121,000	

Assume that categorical features are encoded in a finite range of numbers, e.g., by enumeration. Two individuals  $x$  and  $x'$  are similar if all of their attributes are equal except for:

- Their name or gender;
- Their age can differ by at most 3;
- Their salary can differ by at most \$10,000.

Define  $\phi$  formally and design a DL2 query which finds a counterexample of (1).

2. Let  $x \in \mathbb{R}^d$  be an image of a person and  $G = D \circ E$  be a generative model with an encoder  $E$  and a decoder  $D$ . Two individuals  $x$  and  $x'$  are similar if all of their attributes are the same apart from their skin color and hair color. Propose a formal definition of the similarity set  $S_\phi(x)$  of individuals similar to  $x$  and design a DL2 query which can be used to find counterexamples of (1).

*Hint: Assume access to the vectors  $\mathbf{a}_{\text{skin}}$  and  $\mathbf{a}_{\text{hair}}$  which manipulate these attributes in the latent space of the generative model.*

3. What algorithms can be used to solve the DL2 queries which you wrote in the previous subtasks?

**Problem 2** (Enforcing Individual Fairness). A typical strategy in fair representation learning is to enforce fairness by casting it as an optimization problem. In this task we will design a loss function which can be used to train the data producer  $f_\theta$ . You can assume a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  and a binary function  $\phi$  defining similarity, with respect to which the end-to-end model should be individually fair.

1. Define a loss term which enforces individual fairness and explain how you can compute or approximate it.
2. Define a loss term which makes the data producer  $f_\theta$  aware of the potential data consumer  $h_\psi$  during training and encourages high task utility.
3. Define a loss term which promotes the transferability of the learnt representations, i.e., making them useful to other downstream tasks as well.
4. Mention at least one issue which must be considered when combining all of the above loss terms together.

**Problem 3** (Certifying Individual Fairness). Prove the following lemma which formalizes the compositional individual fairness certificate:

**Lemma 1** (Individual fairness certificate). *Suppose  $M = h_\psi \circ f_\theta$  with data point  $x$  and similarity notion  $\phi$ . Furthermore, let  $z = f_\theta(x)$ ,  $S_\phi(x) = \{x' \in \mathbb{R}^d \mid \phi(x, x')\}$  and*

$$\epsilon = \max_{x' \in S_\phi(x)} \|z - f_\theta(x')\|_\infty. \text{ If}$$

$$\max_{z' \in \mathbb{B}_\infty(z, \epsilon)} h_\psi^{(y')} (z') - h_\psi^{(y)} (z') < 0 \quad (2)$$

for all labels  $y'$  different from the predicted label  $y = M(x)$ , then for all  $x' \in S_\phi(x)$  we have  $M(x) = M(x')$ .

Here,  $h_\psi^{(y)}(z)$  is the value of the logit corresponding to label  $y$  computed by the classifier  $h_\psi$  for the input  $z$ . Note that both  $\epsilon$  and the left-hand side of Eq. (2) can be computed by MILP or soundly approximated by other incomplete certification methods.

**Problem 4** (Properties of Fair Representation – from a previous exam). The encoder  $f_\theta$  and the classifier  $h_\psi$  are trained jointly to compose the end-to-end model  $M = h_\psi \circ f_\theta$ . Moreover,  $f_\theta$  is trained (via DL2 and adversarial training) to satisfy the following condition

$$\phi(x, x') \implies \|f_\theta(x) - f_\theta(x')\|_\infty \leq \delta \quad (3)$$

for a given input  $x$  and a fixed hyperparameter  $\delta$ .

1. Let  $\delta = 0$ . The model  $M$  is such that its encoder  $f_\theta$  satisfies Condition (3) for all inputs  $x \in \mathbb{R}^n$  and for all similarity functions  $\phi$ . That is,

$$\forall \phi, x, x' : \phi(x, x') \implies \|f_\theta(x) - f_\theta(x')\|_\infty \leq 0.$$

Show that  $M$  is a constant classifier, i.e., it classifies all inputs  $x \in \mathbb{R}^n$  the same.

Consider the case  $\delta > 0$ . The model  $M = h_\psi \circ f_\theta$  was trained on a given dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^m$ . Assume that the set of all valid inputs  $x$  is closed and bounded, e.g.,  $x, x' \in [0, 1]^n$ , and therefore  $L = \max_{x \in [0, 1]^n} \|f_\theta(x)\|_\infty > 0$  exists.

2. Prove that the encoder  $f'_\theta(x) := \frac{\delta}{2L} f_\theta(x)$  satisfies Condition (3) for all valid inputs  $x$  and for all similarity functions  $\phi$ . That is,

$$\forall \phi, x, x' : \phi(x, x') \implies \|f'_\theta(x) - f'_\theta(x')\|_\infty \leq \delta.$$

3. Construct a model  $M' = h'_\psi \circ f'_\theta$  with the same accuracy on  $\mathcal{D}$  as  $M$  such that  $f'_\theta$  satisfies Condition (3) for all valid inputs  $x, x'$  and for all similarity functions  $\phi$ .
4. Which of the models,  $M$  or  $M'$  (from the previous subtask), is more individually fair?

## References

- [1] David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. “Learning Adversarially Fair and Transferable Representations”. In: *Proceedings of the 35th International Conference on Machine Learning*. 2018.
- [2] Anian Ruoss, Mislav Balunovic, Marc Fischer, and Martin Vechev. “Learning Certified Individually Fair Representations”. In: *Advances in Neural Information Processing Systems 33*. 2020.