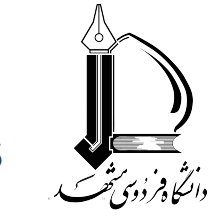


Scaling Guarantees for Nearest Counterfactual Explanations

Kiarash Mohammadi^{1,2} Amir-Hossein Karimi^{1,3}, Gilles Barthe⁴, Isabel Valera⁵
¹MPI for Intelligent Systems, Tübingen, Germany; ²Ferdowsi University of Mashhad, Mashhad, Iran;
³ETH Zürich, Zürich, Switzerland; ⁴MPI for Security and Privacy, Bochum, Germany;
⁵Department of Computer Science, Saarland University, Saarbrücken, Germany.



AAAI / ACM conference on
ARTIFICIAL INTELLIGENCE,
ETHICS, AND SOCIETY

Summary



Method	Opt. Distance	100% Coverage	Efficiency	Neural Models	Qualitative Features	Complex Constraints
Our approach	✓	✓	✓	✓	✓	✓
MACE [1]	✓	✓		✓	✓	✓
DiCE [2]		✓	✓		✓	
Efficient Search [3]	✓	✓	✓		✓	✓

Notation & Background

- Classifier: $h: \mathcal{X} \rightarrow \mathbb{R}$
- Label (e.g., loan approval): $h(\mathbf{x}) \geq 0 \Rightarrow y = +1, h(\mathbf{x}) < 0 \Rightarrow y = -1$
- Individual/factual observation: \mathbf{x}^F s.t. $h(\mathbf{x}^F) < 0$
- Two ways to formulate CFE generation:

Optimization Formulation

$$\mathbf{x}^{CFE} \in \underset{\mathbf{x} \in \mathcal{X}}{\operatorname{argmin}} \quad \operatorname{dist}(\mathbf{x}, \mathbf{x}^F) \quad (1)$$

$$\text{s.t.} \quad h(\mathbf{x}) \geq 0$$

$$\mathbf{x} \in \text{Plausible}$$

$$\mathbf{x} \in \text{Actionable}$$

Verification Formulation

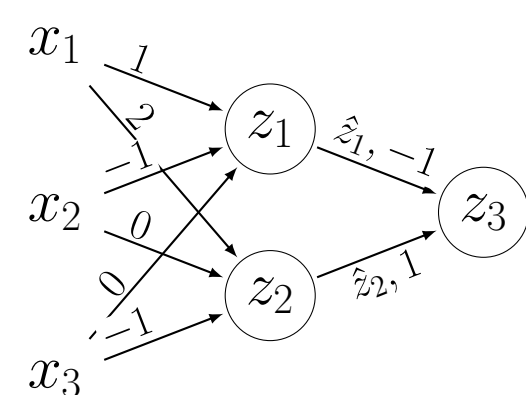
$$\exists \mathbf{x}. \operatorname{dist}(\mathbf{x}, \mathbf{x}^F) \leq \delta$$

$$h(\mathbf{x}) \geq 0$$

$$\mathbf{x} \in \text{Plausible}$$

$$\mathbf{x} \in \text{Actionable} \quad (2)$$

- Encoding a ReLU Neural Network (NN) as an SMT formula:



$$\phi_f(x) =$$

$$(z_1 = x_1 - x_2) \wedge (z_2 = 2x_1 - x_3)$$

$$\wedge ((\hat{z}_1 = z_1 \wedge z_1 \geq 0) \vee (\hat{z}_1 = 0 \wedge z_1 < 0))$$

$$\wedge ((\hat{z}_2 = z_2 \wedge z_2 \geq 0) \vee (\hat{z}_2 = 0 \wedge z_2 < 0))$$

$$\wedge (z_3 = -z_1 + \hat{z}_2)$$

MIP Encodings

For NNs, we adopt a *bounded* encoding by Tjeng and Tedrake [4], i.e., for $i \in \{1, \dots, n\}$ (\mathbf{l}_i and \mathbf{u}_i indicate the lower/upper bounds):

$$\mathbf{z}_i = \mathbf{W}_i \hat{\mathbf{z}}_{i-1} + \mathbf{b}_i \quad (3a)$$

$$\delta_i \in \{0, 1\}^{k_i}, \quad \hat{\mathbf{z}}_i \geq 0, \quad \hat{\mathbf{z}}_i \leq \mathbf{u}_i \cdot \delta_i, \quad (3b)$$

$$\hat{\mathbf{z}}_i \geq \mathbf{z}_i, \quad \hat{\mathbf{z}}_i \leq \mathbf{z}_i - \mathbf{l}_i \cdot (1 - \delta_i)$$

As a preliminary step, a linear approximation of ReLUs from Ehlers [5] replaces (3b) to compute the lower/upper bounds:

$$\hat{\mathbf{z}}_i \geq \mathbf{z}_i, \quad \hat{\mathbf{z}}_i \geq 0, \quad \hat{z}_{i,j} \leq u_{i,j} \frac{z_{i,j} - l_{i,j}}{u_{i,j} - l_{i,j}} \quad (4)$$

These bounds are then placed in (3) to complete an *exact* MIP encoding.

CFE Generation

Within each iteration of an exponential search that gradually increases the distance interval, tight bounds of the hidden units are computed.

- MIP-SAT:** Relies on SMT solving. Removes fixed-state ReLUs from the SMT formula given their bounds. Verifies the new formula.
- MIP-EXP:** Relies on MIP solving. Uses the bounds to implement (3) and optimizes it toward and until flipping the output logit sign.
- MIP-OBJ:** No exponential search. Uses the distance MIP as the objective, with a constraint of the output logit being flipped, directly minimizes distance.

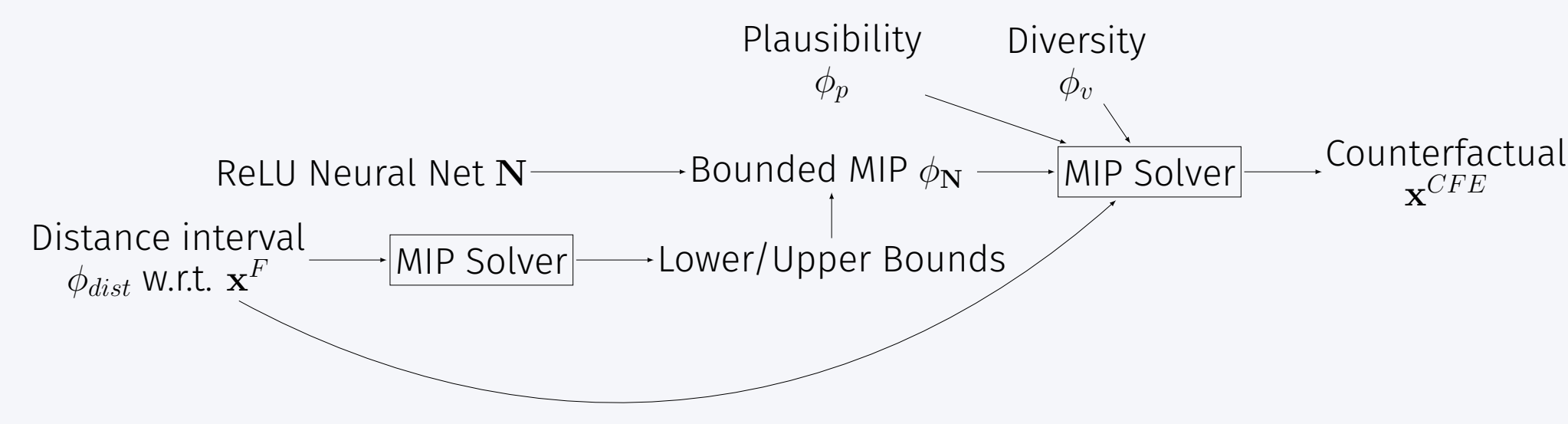
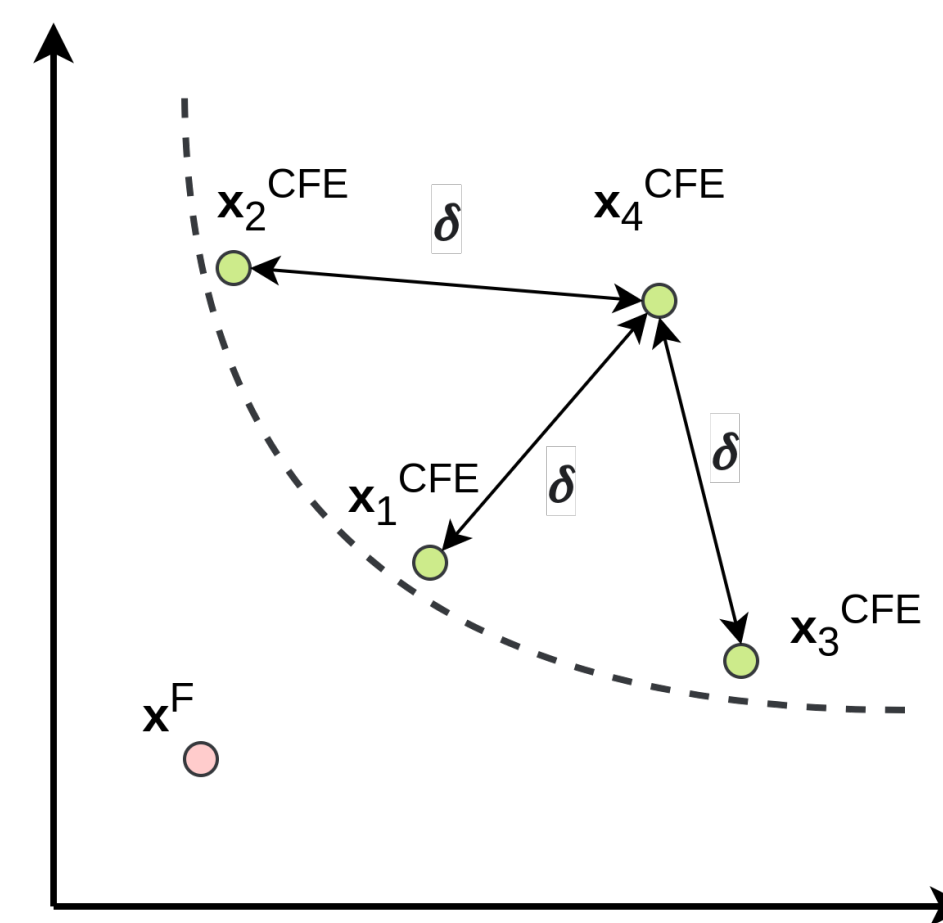


Figure 2. An iteration of the search



Different qualitative features, such as diversity, may also be encoded. For diversity, the following constraints are added for the i -th CFE to be generated:

$$\operatorname{dist}(\mathbf{x}_1^{CFE}, \mathbf{x}_i^{CFE}) \geq \delta$$

$$\dots$$

$$\operatorname{dist}(\mathbf{x}_{i-1}^{CFE}, \mathbf{x}_i^{CFE}) \geq \delta \quad (5)$$

Experiments

Setup.

We compare our approaches against MACE [1] and DiCE [2] in various settings on runtime, distance, and coverage. We employ three widely used real-world datasets from the literature: Adult ($d = 51$), COMPAS ($d = 7$), and Credit ($d = 20$). We use a two-layer ReLU-activated NN with 10 neurons for most experiments. While NNs of this scale can sufficiently discriminate between the classes of the supervised learning task, we also include experiments exploring the scalability (deepening or widening of the NN) of our approach against the opponents. Finally, we show that qualitative features, diversity in this case, can be encoded within the framework to efficiently generate sets of CFEs with guarantees

Results.

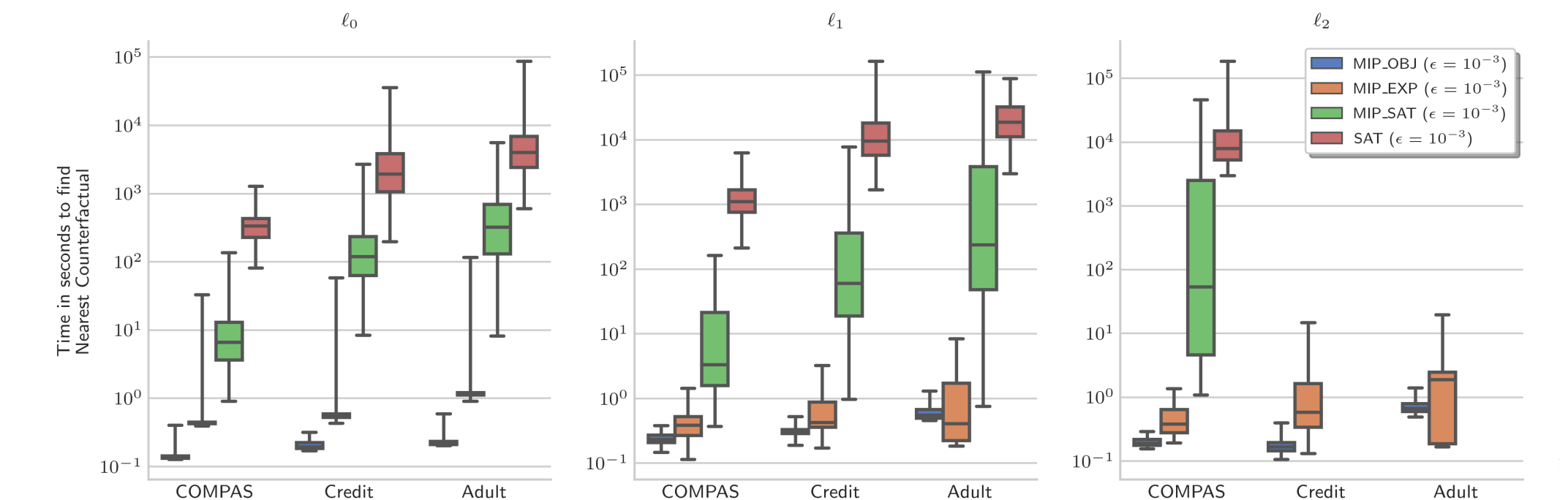


Figure 3. Full-setting runtime comparison of two-layer ReLU-activated NN with 10 neurons in each layer among our approach and MACE (SAT) [1]. Coverage is perfect by design. Each setting has been evaluated on 500 instances, however, SAT and MIP-SAT timed out on some samples. For such cases, only the samples for which all approaches have successfully finished running are included.

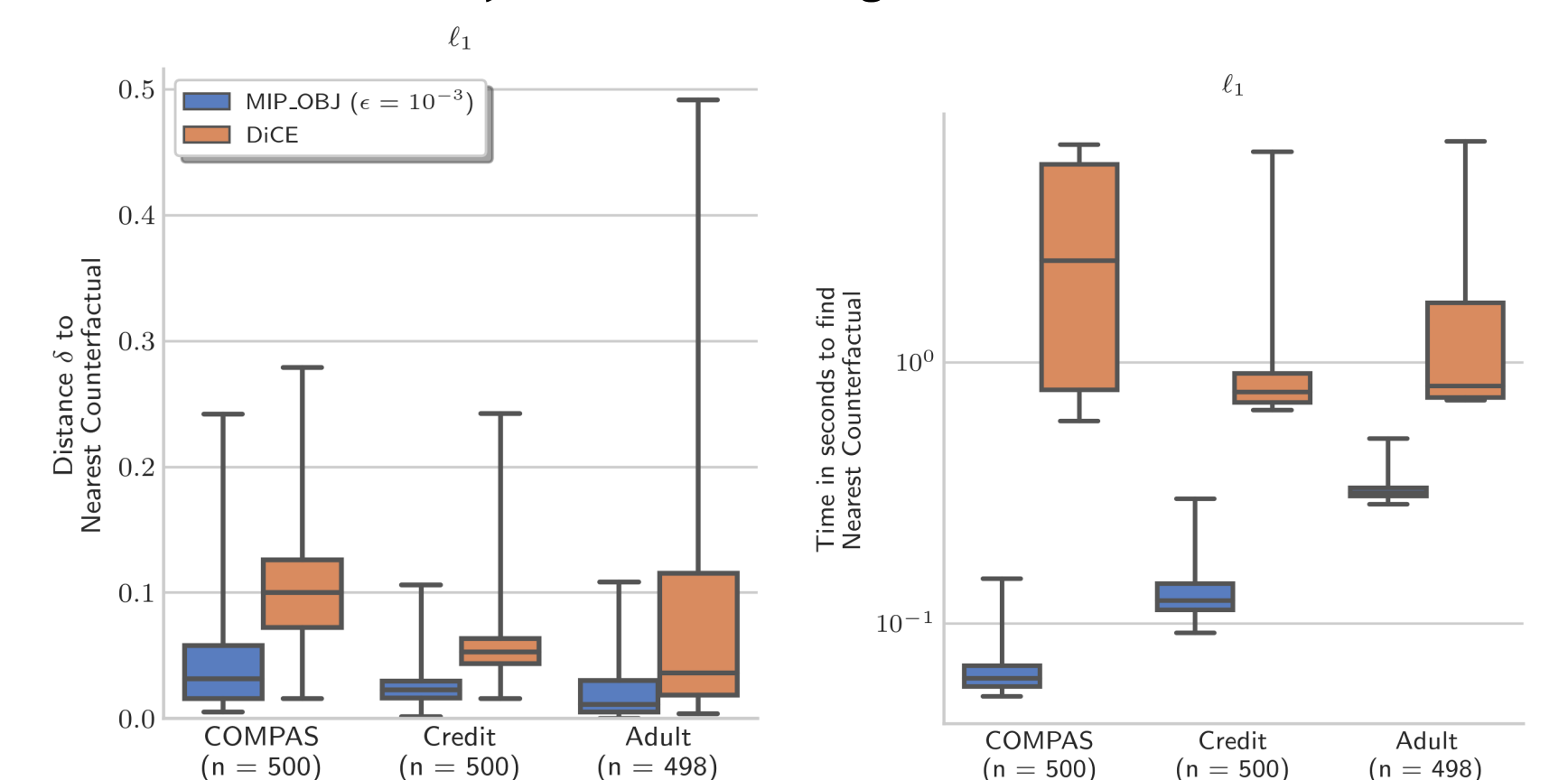


Figure 4. Distance and time comparison against DiCE as a gradient-based optimization approach. The NN model is same as above. MIP-OBJ coverage is perfect by design and DiCE coverage is also perfect but slightly dips for Adult dataset (99.6%).

References

- Amir-Hossein Karimi, Gilles Barthe, Borja Balle, and Isabel Valera. Model-agnostic counterfactual explanations for consequential decisions. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108, 2020.
- Ramaravind K. Mothilal, Amit Sharma, and Chenhao Tan. Explaining machine learning classifiers through diverse counterfactual explanations. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 2020.
- Chris Russell. Efficient search for diverse coherent explanations. In *Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT* '19*, page 20–28, New York, NY, USA, 2019. Association for Computing Machinery.
- Vincent Tjeng and Russ Tedrake. Verifying neural networks with mixed integer programming. *CoRR*, abs/1711.07356, 2017.
- Rüdiger Ehlers. Formal verification of piece-wise linear feed-forward neural networks. *CoRR*, abs/1705.01320, 2017.