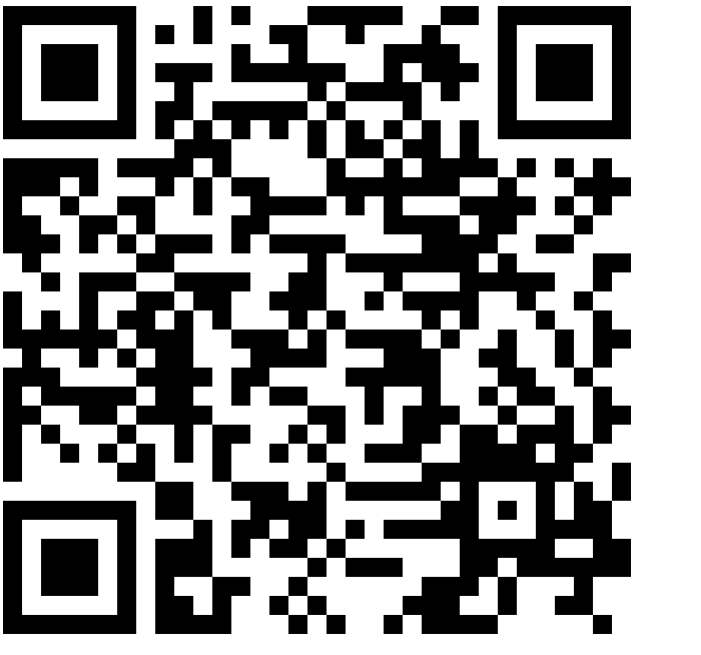


Certified defences hurt generalisation

Piersilvio De Bartolomeis¹, Jacob Clarysse¹, Fanny Yang¹, Amartya Sanyal²

¹Department of Computer Science, ETH Zürich

²ETH AI Center



ADVERSARIALLY ROBUST CLASSIFICATION

Goal: Low robust test error for a class of perturbations $\mathcal{B}_\epsilon(x)$

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{x' \in \mathcal{B}_\epsilon(x)} L(f_{\theta}(x'), y) \right]$$

Empirical defences

- Solve the inner maximisation with first-order optimisation methods
- Find a lower bound solution x^*

$$L(f_{\theta}(x^*), y) \leq \max_{x' \in \mathcal{B}_\epsilon(x)} L(f_{\theta}(x'), y)$$

Certified defences

- Solve a convex relaxation of the inner-maximisation
- Find an upper bound solution x^*

$$\max_{x' \in \mathcal{B}_\epsilon(x)} L(f_{\theta}(x'), y) \leq L(f_{\theta}(x^*), y)$$

FAIRNESS

We measure the degree of unfairness as follows:

$$\frac{\max_k \mathbf{R}^k(\theta) - \mathbf{R}(\theta)}{1 - \mathbf{R}(\theta)}$$

where

- $\mathbf{R}(\theta)$ is the error of the classifier
- $\mathbf{R}^k(\theta)$ is the error conditioned on the label k

REFERENCES

- [1] Aleksander Madry et al, Towards Deep Learning Models Resistant to Adversarial Attacks. ICLR, 2018.
- [2] Eric Wong and J. Zico Kolter, Provable Defenses against Adversarial Examples via the Convex Outer Adversarial Polytope. ICML, 2018.
- [3] Huan Zhang et al., Towards Stable and Efficient Training of Verifiably Robust Neural Networks. ICLR, 2020.

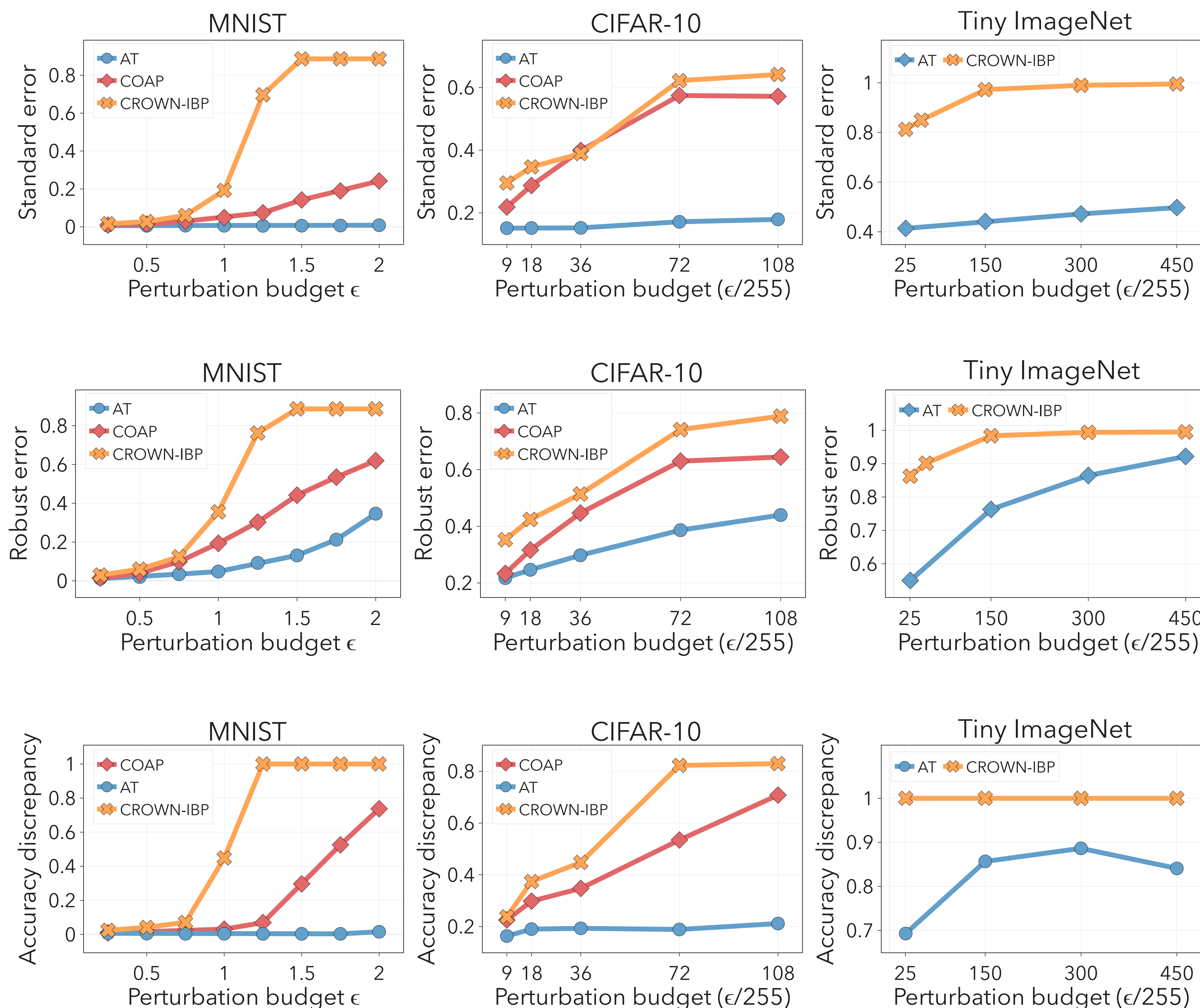
How effective are certified defences in practice?

CERTIFIED DEFENCES HURT GENERALISATION AND FAIRNESS

Threat model: ℓ_2 -ball perturbations of radius ϵ

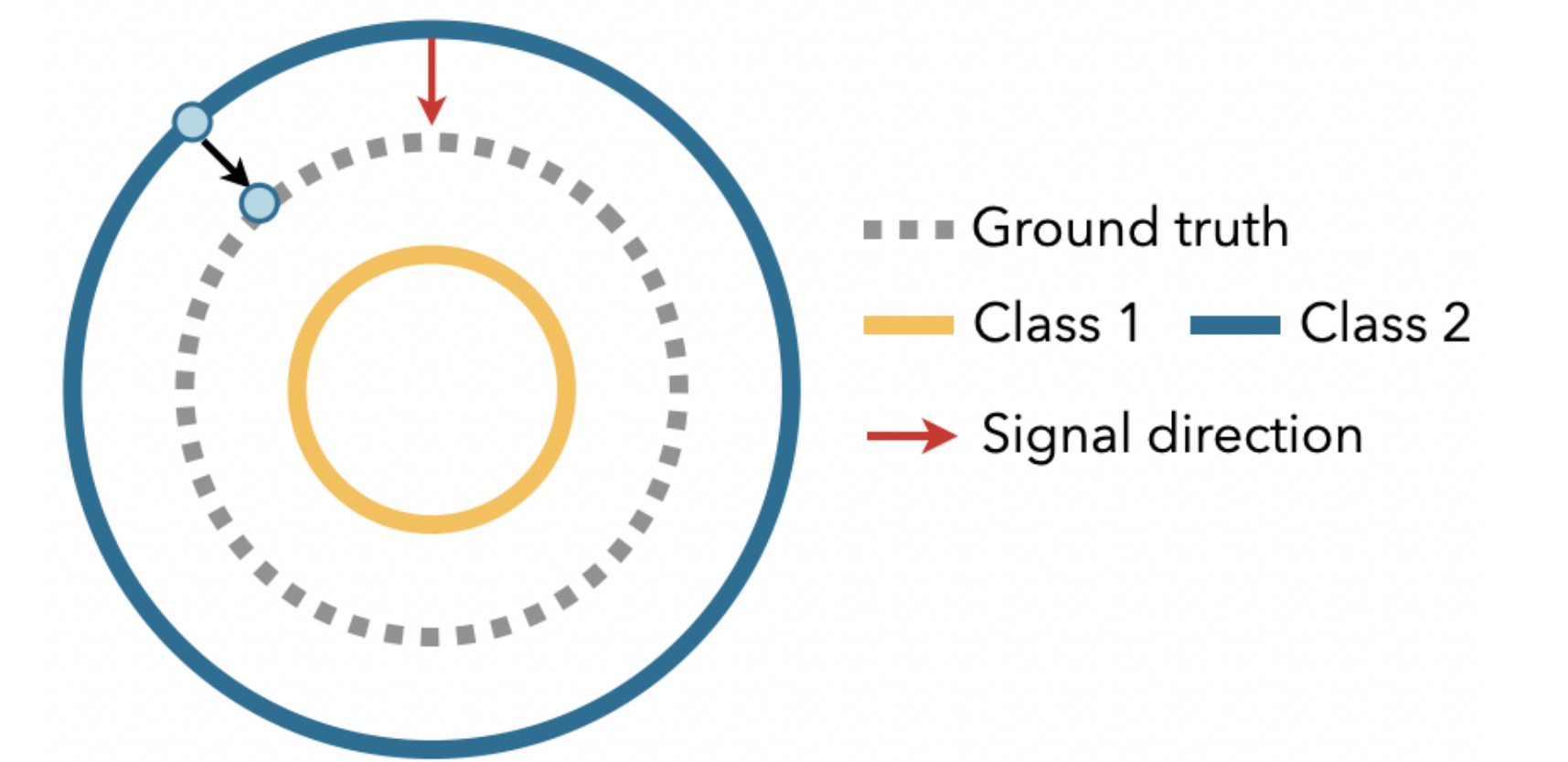
Empirical defences: adversarial training (AT) [1]

Certified defences: convex outer adversarial polytope (COAP) [2] and CROWN-IBP [3]



WHY? INTUITION ON SYNTHETIC DATA

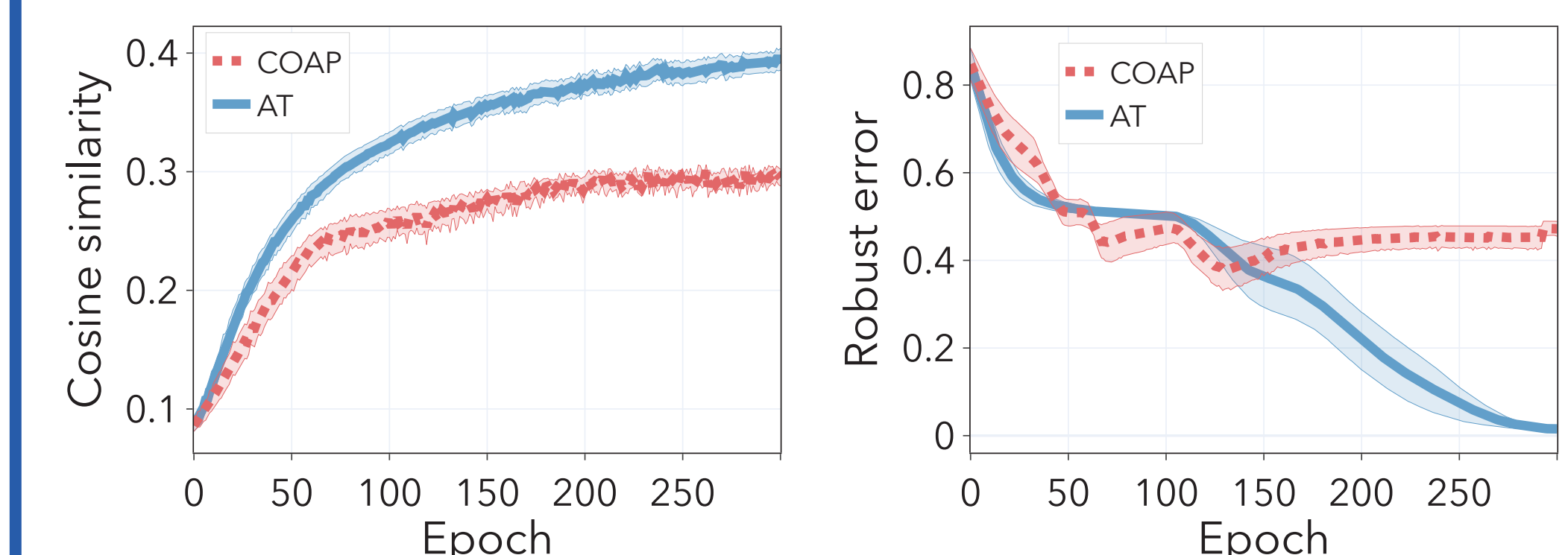
A more controlled setting: concentric spheres



Certified defences hurt generalisation when

- perturbation magnitude is large and
- perturbation aligns with the signal direction

Large ϵ



THEORETICAL RESULT

Setting

- $y \sim \{-1, 1\}$, $x_1 = \gamma \text{sgn}(y)$, $x_{2:d} \sim \mathcal{N}(0, \sigma^2 I_{d-1})$
- $\mathcal{B}_\epsilon(x) = \{z_1 = x + e_1 \beta \mid |\beta| \leq \epsilon\}$
- $f_\theta(x) = a \text{ReLU}(\theta^\top x) + b$
- θ_1 is the only trainable parameter

Theorem (informal)

For $\frac{2}{3}\gamma < \epsilon < \gamma$ and d large enough, after one step of gradient descent with respect to the COAP and AT objectives, we have: $\mathbf{R}_\epsilon(\theta^{\text{COAP}}) > \mathbf{R}_\epsilon(\theta^{\text{AT}})$