

Adversarial Robustness is at Odds with Lazy Training

Yunjuan Wang, Enayat Ullah, Poorya Mianjy, Raman Arora
Department of Computer Science, Johns Hopkins University

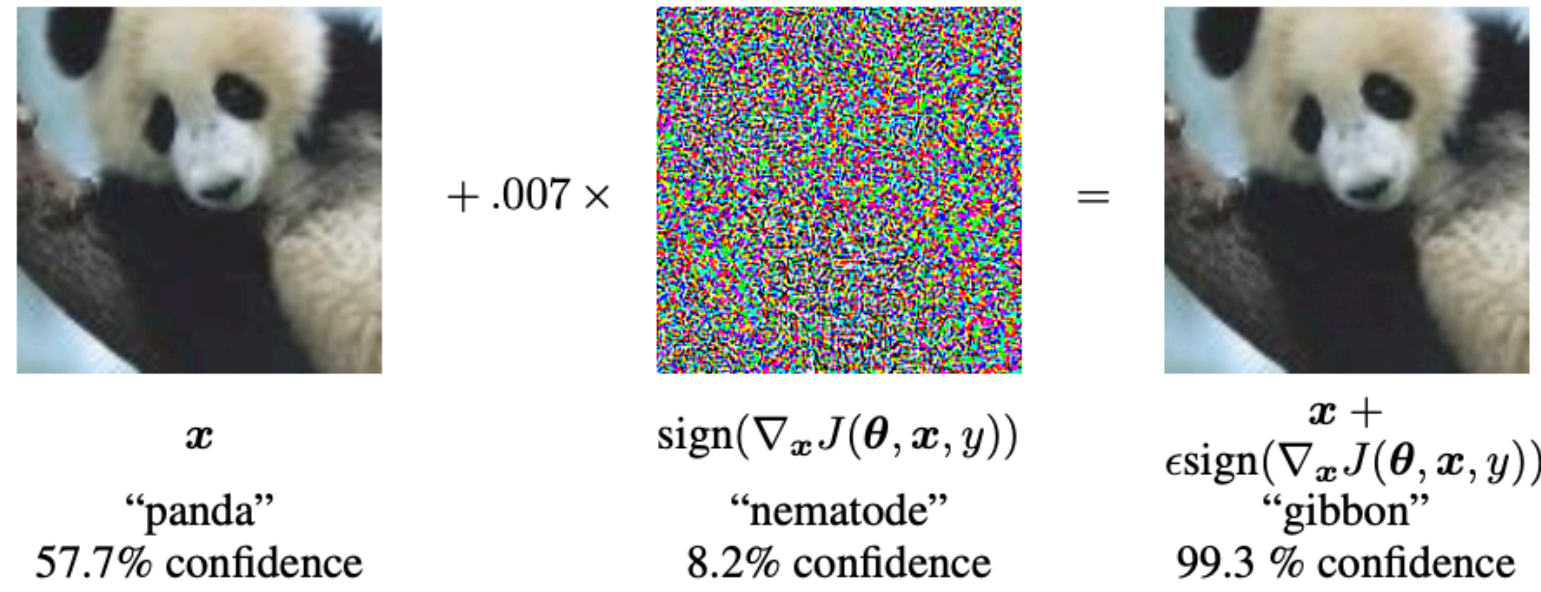


JOHNS HOPKINS UNIVERSITY



1. Background

- ML systems are fragile and susceptible to imperceptible attacks [GSS15].



2. Motivation

[SSRD19]: Propose an algorithm to generate bounded ℓ_0 -norm adversarial perturbation with guarantees for arbitrary deep networks.

[DS21]: Multi-step gradient ascent can find adversarial examples for random ReLU networks with small widths.

[BCGT21]: A single gradient step finds adversarial examples for sufficiently wide but not extremely wide randomly initialized ReLU networks.

[BBC21]: Extend the above to randomly initialized deep networks.

Question: why do trained neural networks remain vulnerable to adversarial attacks?

7. Experiments

- Binary MNIST. Networks trained using SGD in the lazy regime.
- Recall notation: m : network width; $\|\delta\|$: L2 perturbation budget; η : step size;

$$V: \text{maximal weight deviation } \max_{s \in [m]} \|\bar{w}_s - w_{s,0}\|_2 = C_0/\sqrt{m}$$

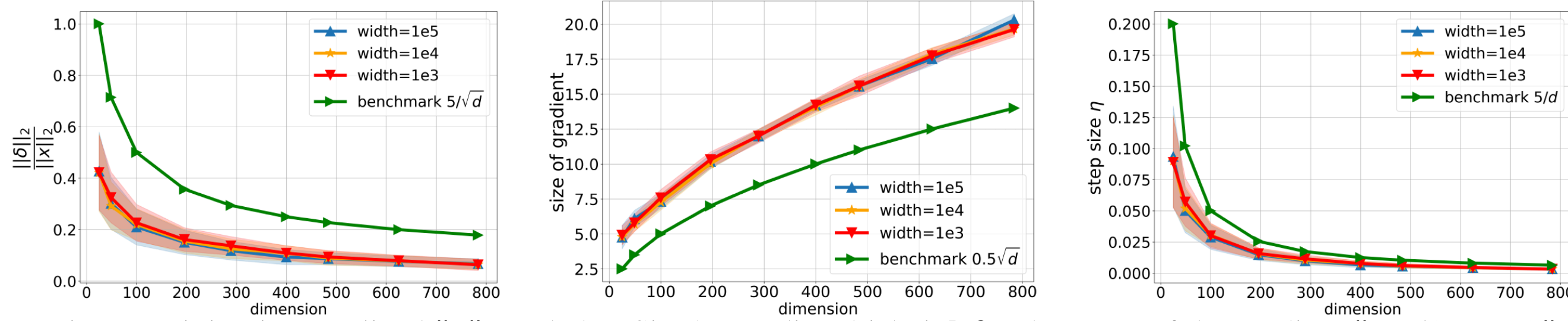


Figure: Minimal normalized $\|\delta\|$ needed to flip the predicted label (left), the norm of the gradient $\|\nabla_x f(x; a, W)\|$ for W in the lazy regime (middle), corresponding step size (right), as a function of the input dimension for a fixed value of $C_0 = 10$ and different network widths m .

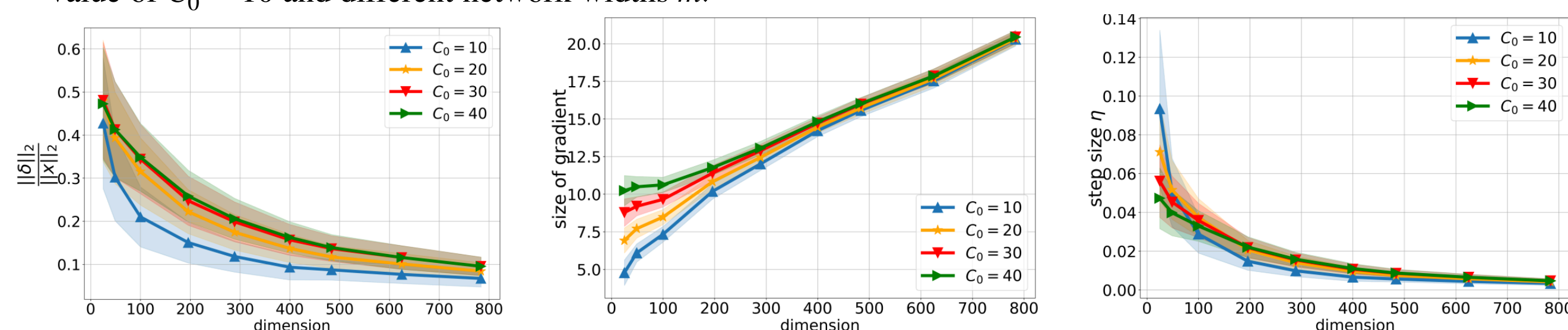


Figure: Minimal normalized $\|\delta\|$ needed to flip the predicted label (left), the norm of the gradient $\|\nabla_x f(x; a, W)\|$ for W in the lazy regime (middle), corresponding step size (right), as a function of the input dimension for a fixed network width of $m = 10^5$ and different value of C_0 .

3. Problem Setup

- $\mathcal{X} \subseteq \mathbb{R}^d, \mathcal{Y} = \{\pm 1\}$.

A two-layer width m ReLU net (a, W) , $f(x; a, W) := \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \sigma(w_s^T x)$.

- Attack model: ℓ_2 attack with perturbation budget R .

4. Lazy Regime

The dominant model for (non-robust) deep learning [JGH18, JT19, ADHL19].

Initialization: 1) $a_s \sim \text{unif}(\{-1, +1\})$, fixed; 2) $w_{s,0} \sim \mathcal{N}(0, I_d), \forall s \in [m]$.

Why lazy regime?

- Provable generalization:** there exists $\bar{W} : \|\bar{w}_s - w_{s,0}\|_2 = \mathcal{O}(1/\sqrt{m}), \forall s \in [m]$, such that the non-robust generalization error is small.

$\forall s \in [m]$, such that the non-robust generalization error is small.

- Computational Tractability:** Such \bar{W} can be found by efficient first-order methods such as Stochastic Gradient Descent (SGD).

Definition: The lazy regime is the set of all networks parameterized by (a, W) , such that

$$W \in \mathcal{B}_{2,\infty}(W_0, C_0/\sqrt{m}) = \left\{ W : \|\bar{w}_s - w_{s,0}\|_2 \leq C_0/\sqrt{m}, \forall s \in [m] \right\}.$$

Question: Are lazy regime networks susceptible to adversarial attacks?

5. Main Result

A single gradient step suffices to attack lazy regime networks.

Theorem: With probability at least $1 - \gamma$, for all $W \in \mathcal{B}_{2,\infty}(W_0, C_0/\sqrt{m})$

$$\text{sign}(f(x; a, W)) \neq \text{sign}(f(x + \delta; a, W))$$

where $\delta = \eta \nabla_x f(x; a, W)$ with $|\eta| = \mathcal{O}(1/d)$,

$$\max \{d^{2.4}, \mathcal{O} \log(1/\gamma)\} \leq m \leq \mathcal{O}(\exp(d^{0.24}))$$

Remark: Imperceptible perturbation $\|\delta\| = \mathcal{O}(1/\sqrt{d})$.

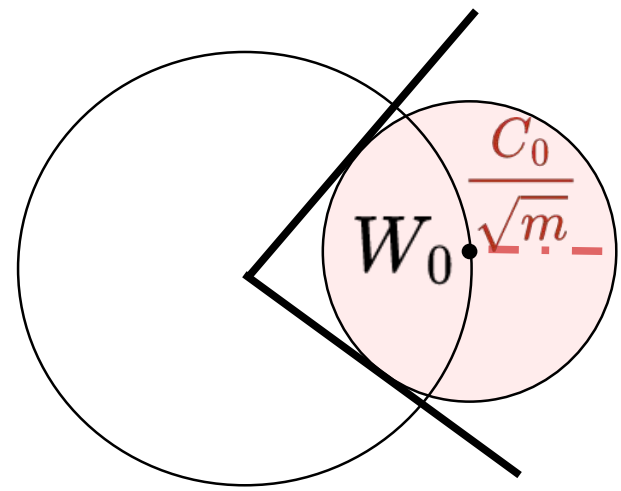
Corollary: W.h.p., for all $W \in \mathcal{B}_{2,\infty}(W_0, C_0/\sqrt{m})$ the robust error is greater than 0.9.

6. Beyond Lazy Regime

- Scaling the network weights with a positive factor for binary classification does not change the prediction.

- Our results can be extended to the following cone

$$C = \left\{ W : \exists r > 0 : \|rW - W_0\|_{2,\infty} \leq C_0/\sqrt{m} \right\}$$



Observations:

- Experimental results closely match the theoretical bound for different network widths and weight deviations:

- $\|\delta\| = \mathcal{O}(1/\sqrt{d})$; 2. $|\eta| = \mathcal{O}(1/d)$

Main Takeaway: Networks within the lazy training regime are vulnerable to adversarial attacks.

[SSRD19]: Shamir, Adi, et al. "A simple explanation for the existence of adversarial examples with small hamming distance." *arXiv preprint arXiv:1901.10861* (2019).

[DS21]: Daniely, Amit, and Hadas Schacham. "Most ReLU Networks Suffer from ℓ_2 Adversarial Perturbations." *NeurIPS* (2020).

[BCGT21]: Bubeck, Sébastien, et al. "A single gradient step finds adversarial examples on random two-layers neural networks." *NeurIPS* (2021).

[BBC21]: Bartlett, Peter, Sébastien Bubeck, and Yeshwanth Cherapanamjeri. "Adversarial examples in multi-layer random relu networks." *NeurIPS* (2021).

[JT19]: Ji, Ziwei, and Matus Telgarsky. "Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow relu networks." *In ICLR* (2020).

- Binary MNIST. Networks trained using adversarial training in the lazy regime.

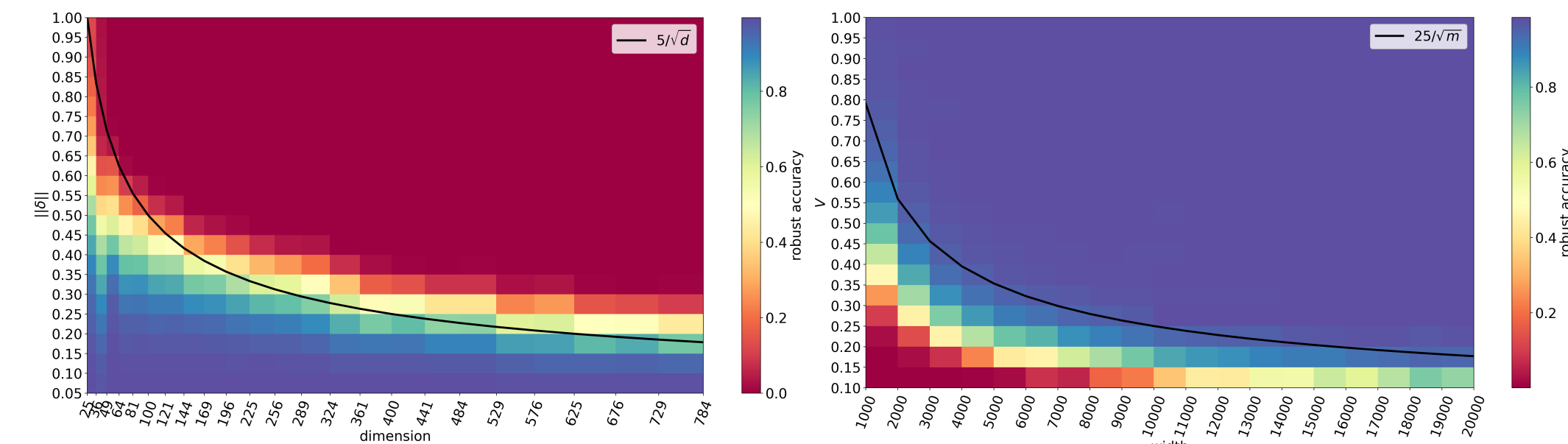


Figure: Robust test accuracy as a function of perturbation size $\|\delta\|$ and dimension d given a fixed network width and fixed value of the maximal deviation of weight vectors (left); Robust test accuracy as a function of width m and maximal deviation V given a fixed dimension and fixed perturbation size (right).

Observations:

- a sharp drop in robust accuracy about the $\mathcal{O}(1/\sqrt{d})$ threshold for the perturbation budget $\|\delta\|$ as predicted by the theorem.
- a phase transition in the robust test accuracy at the value of V around $\mathcal{O}(1/\sqrt{m})$ as required by the theorem.

