

Wasserstein adversarial examples via projected Sinkhorn iterations

Eric Wong¹ Frank R. Schmidt² J. Zico Kolter^{1,3}

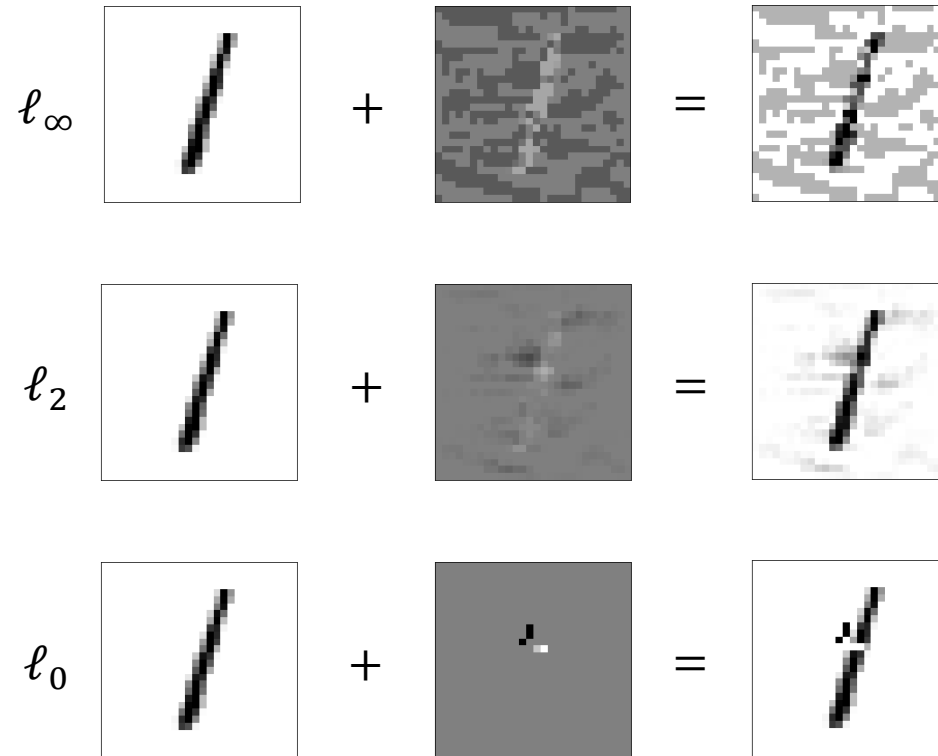
Code: https://github.com/locuslab/projected_sinkhorn/

¹School of Computer Science
Carnegie Mellon University
Pittsburgh, PA, 15213, USA

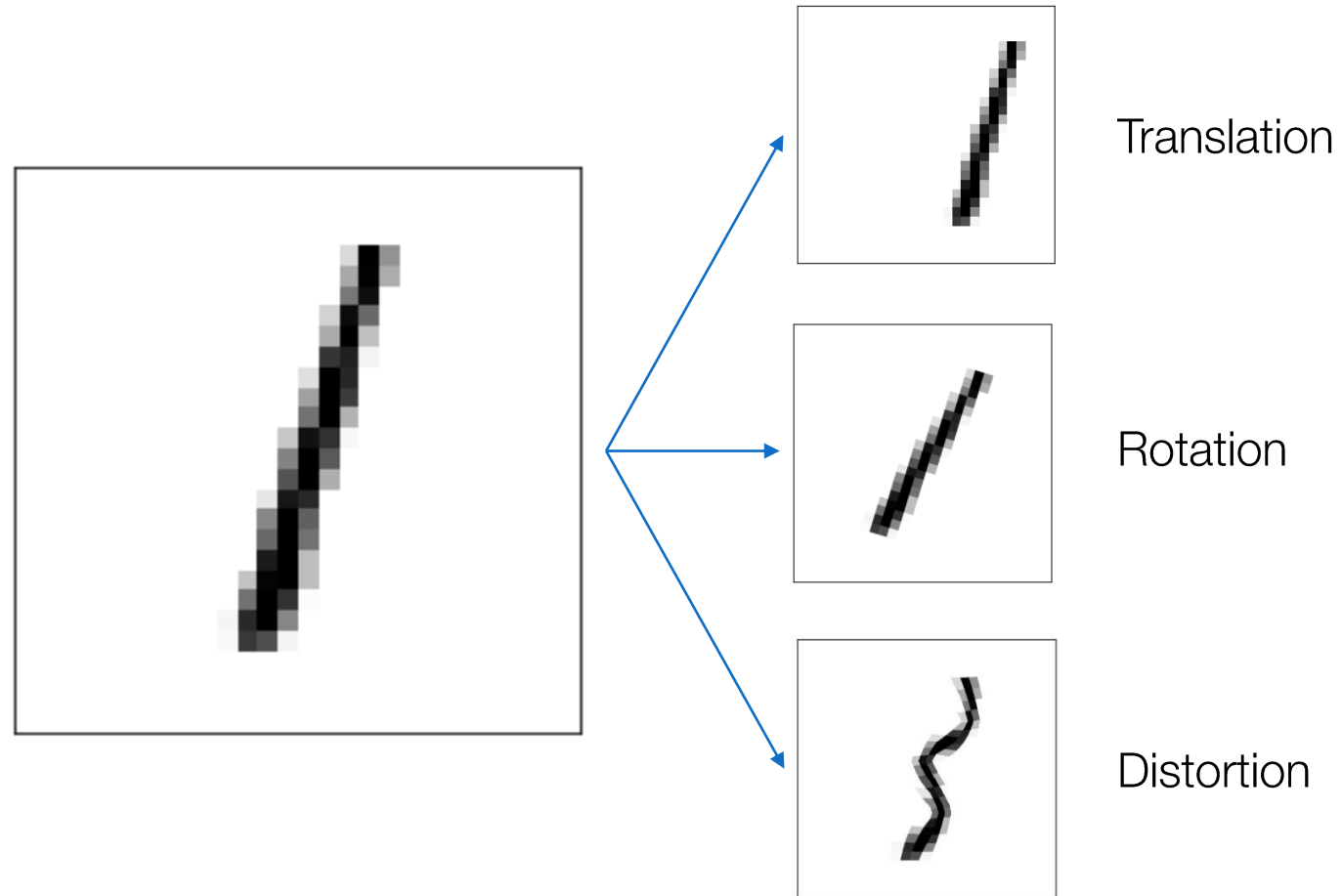
²Bosch Center for Artificial Intelligence
Pittsburgh, PA 15222, USA

³Bosch Center for Artificial Intelligence
Renningen, Germany

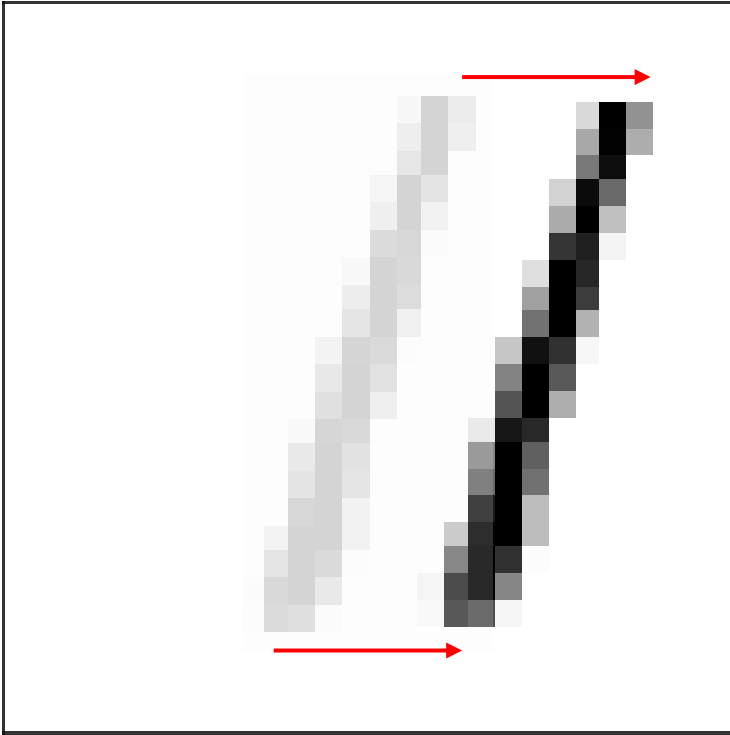
Typical threat model: norm-bounded perturbation



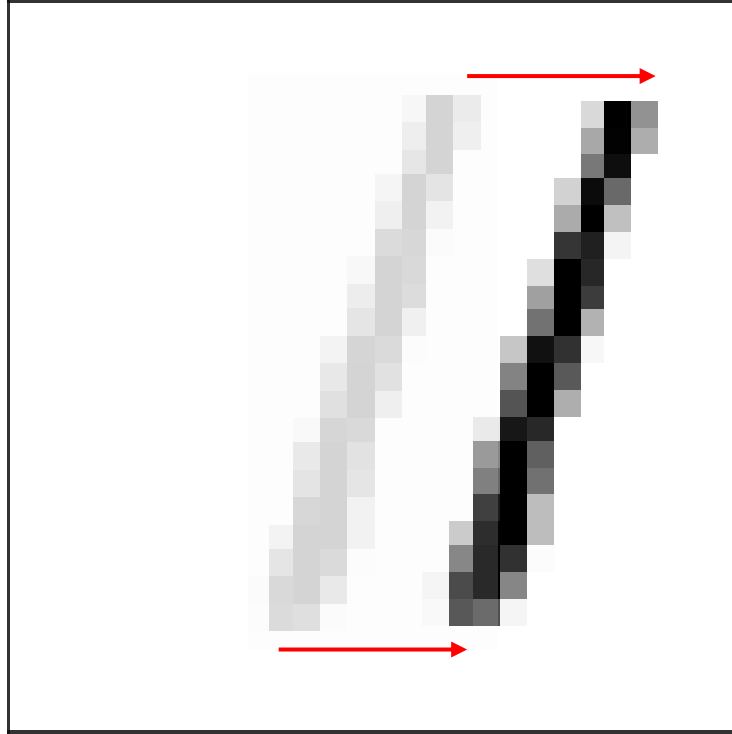
ℓ_p norms don't capture typical image transforms



These transforms move pixel mass short distances...

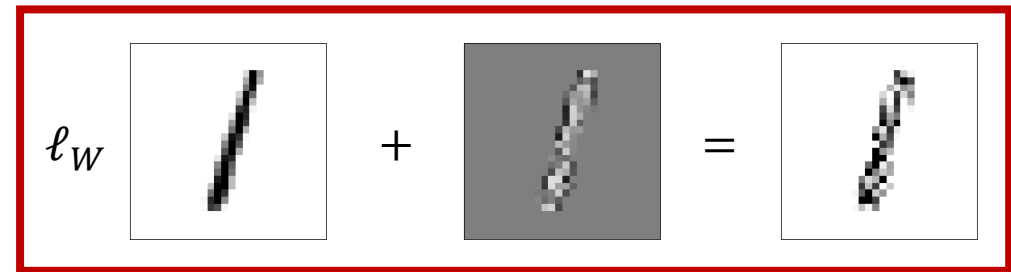
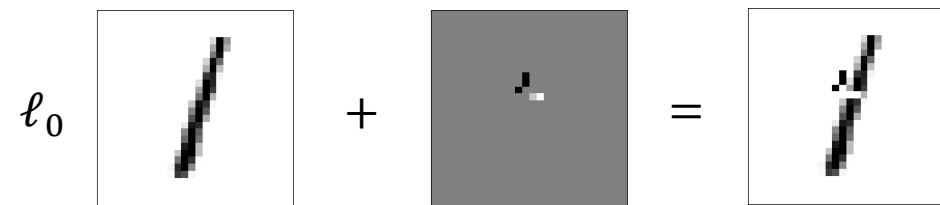
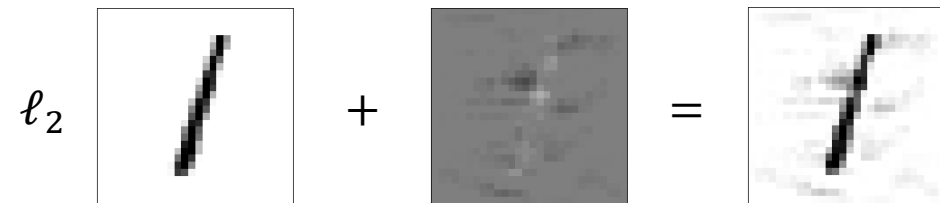
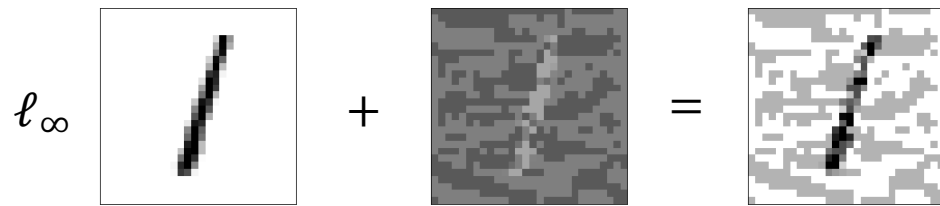


These transforms move pixel mass short distances...

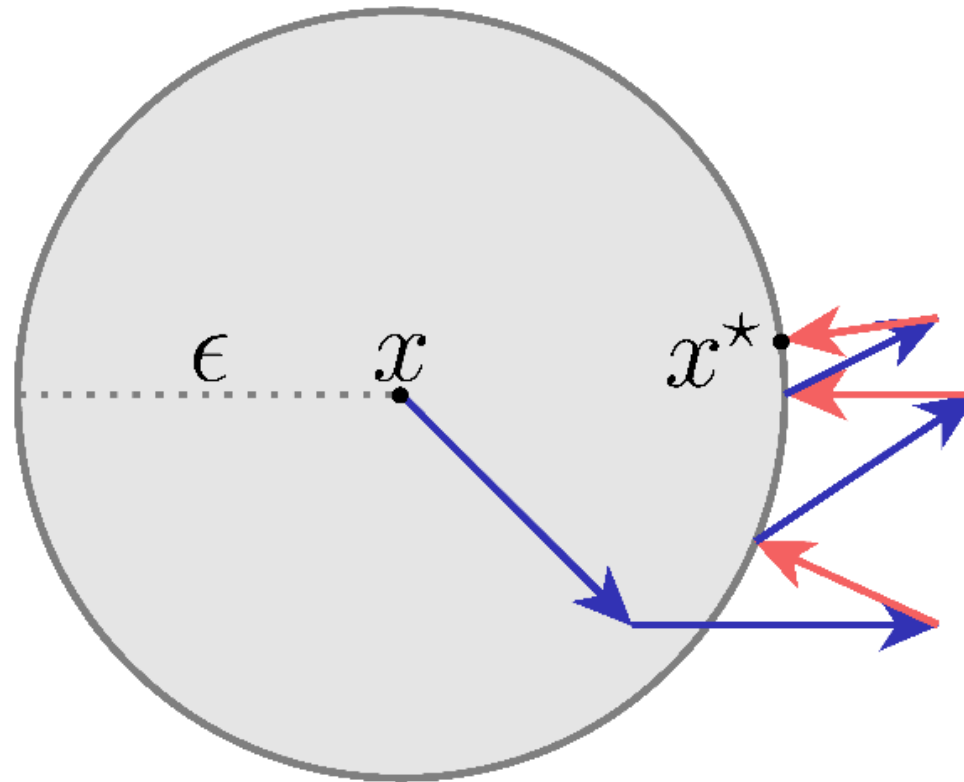


and the Wasserstein metric measures “moving mass”

We propose Wasserstein balls as a threat model



The strongest known method for generating adversarial examples is projected gradient descent



How to project onto the Wasserstein ball?

How to project onto the Wasserstein ball?

$$\underset{z \in \mathbb{R}^n, \Pi \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \|w - z\|_2^2$$

$$\text{subject to} \quad \begin{aligned} \Pi \mathbf{1} &= x \\ \Pi^T \mathbf{1} &= z \\ \langle \Pi, C \rangle &\leq \epsilon \end{aligned}$$

How to project onto the Wasserstein ball?

$$\begin{array}{ll} \text{minimize} & \boxed{\|w - z\|_2^2} \leftarrow \text{Closest point} \\ \text{subject to} & \boxed{\begin{array}{l} \Pi \mathbf{1} = x \\ \Pi^T \mathbf{1} = z \\ \langle \Pi, C \rangle \leq \epsilon \end{array}} \leftarrow \text{Wasserstein ball constraints} \end{array}$$

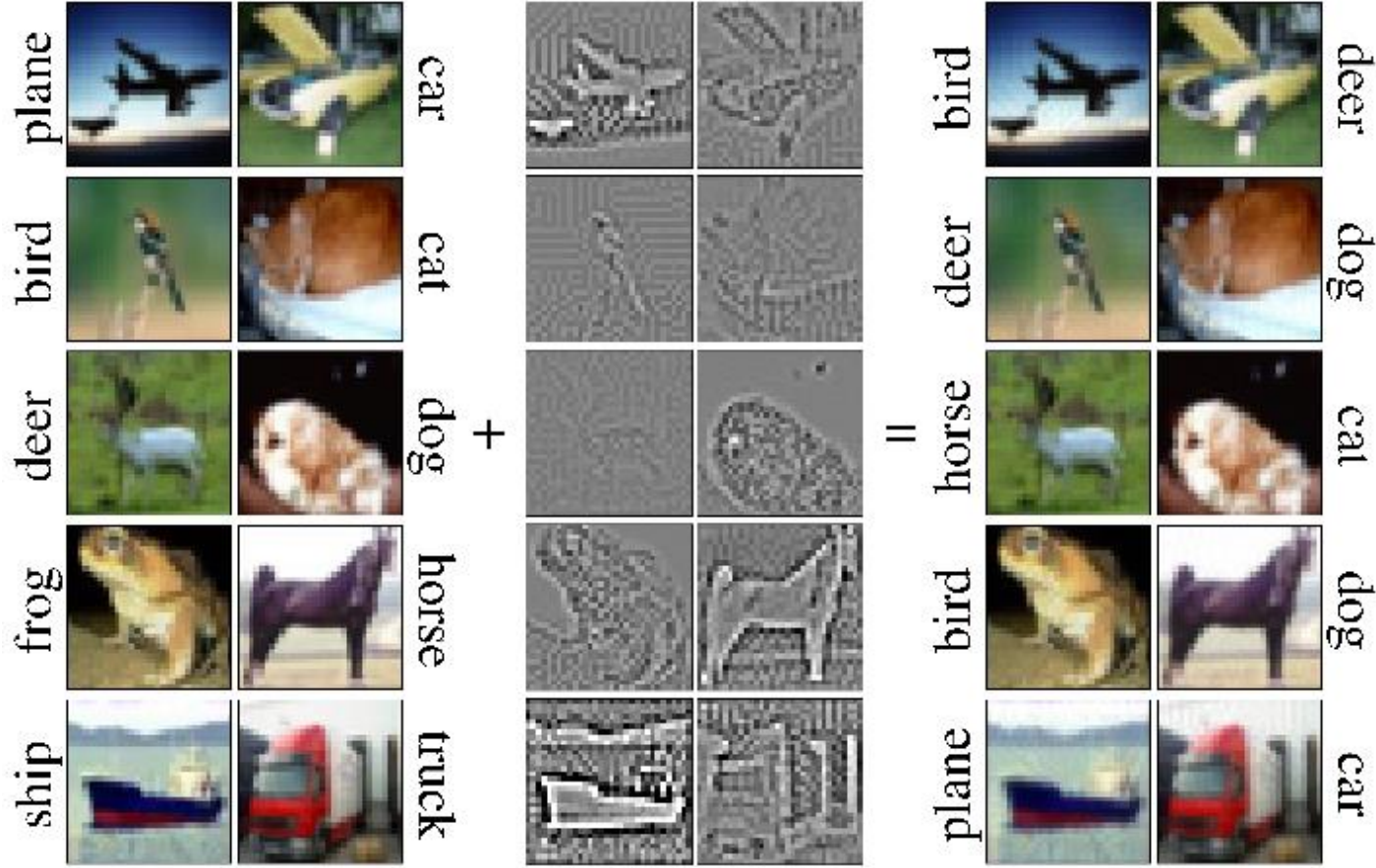
Quadratic program, quadratic number of variables \rightarrow Costly!

Projected Sinkhorn Iteration: a fast (approximate) projection algorithm onto the Wasserstein ball

$$\begin{aligned} & \underset{z \in \mathbb{R}^n, \Pi \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \|w - z\|_2^2 + \frac{1}{\lambda} \sum_{ij} \Pi_{ij} \log(\Pi_{ij}) \\ & \text{subject to} \quad \Pi \mathbf{1} = x \\ & \quad \quad \quad \Pi^T \mathbf{1} = z \\ & \quad \quad \quad \langle \Pi, C \rangle \leq \epsilon \end{aligned}$$

- Entropy regularization term
- Local transport plans
- Block coordinate descent on the dual problem

CIFAR10 Wasserstein adversarial examples



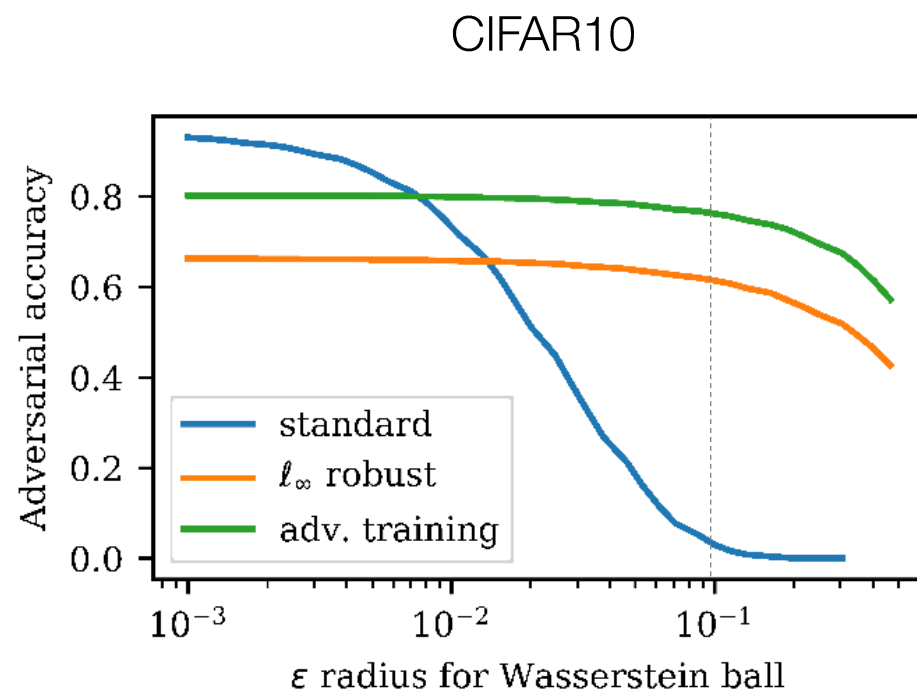
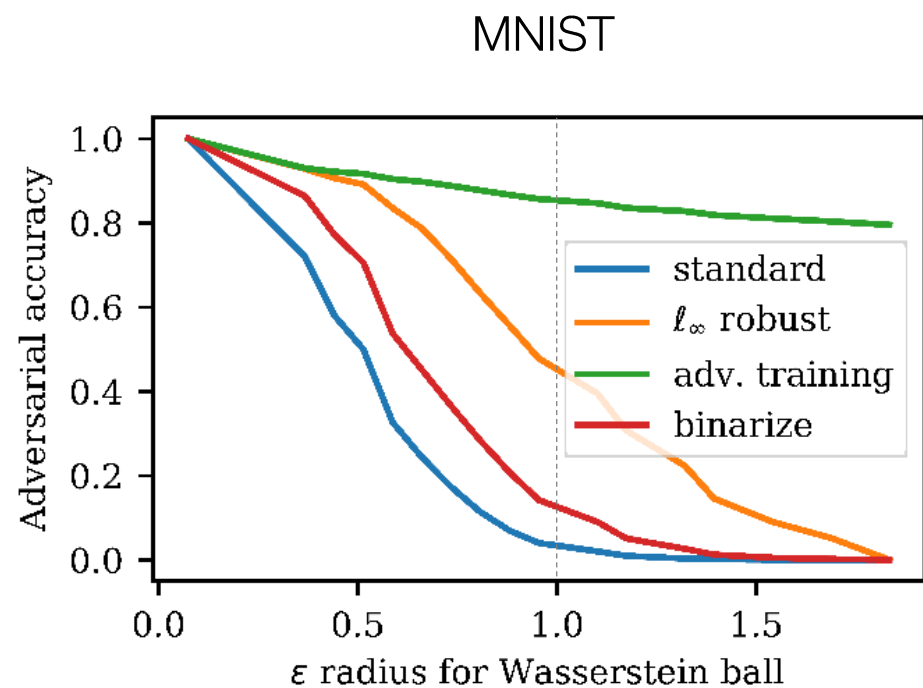
We can also adversarially train robust networks

Adversarial training

PGD adversary

Projection algorithm
onto Wasserstein balls

MNIST and CIFAR10 robustness curves



Wasserstein adversarial examples via projected Sinkhorn iterations

Poster #67 in the Pacific Ballroom

Code: https://github.com/locuslab/projected_sinkhorn

