

Towards Deep Learning Models Resistant to Adversarial Attacks

Aleksander Madry, Aleksandar, Markelov, Ludwig Schmidt, Dimitris
Tsipras, Adrian Vladu

Presented by: Ethan Rathbun

Problem Statement

- Adversarial examples propose a fundamental threat to ML models
- Improving and further understanding the limitations of model robustness is important
- Defenses against a generalized adversary are most useful

The Saddle Point Problem

- They formulate a “minimax” problem over the set of possible model parameters and attack perturbations.
- The goal is to find a model that minimizes the expected loss under adversarial perturbations bounded within an l_p ball.

$$\min_{\theta} \rho(\theta), \quad \text{where} \quad \rho(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\delta \in \mathcal{S}} L(\theta, x + \delta, y) \right]$$

- This is the same framework commonly used in zero-sum games

Approximating the Inner Maximization

- Given a fixed model we can approximate a local solution to the saddle point problem by using a sufficiently strong adversarial attack.
- In this paper they essentially propose using vanilla PGD as the foundation of their optimization.

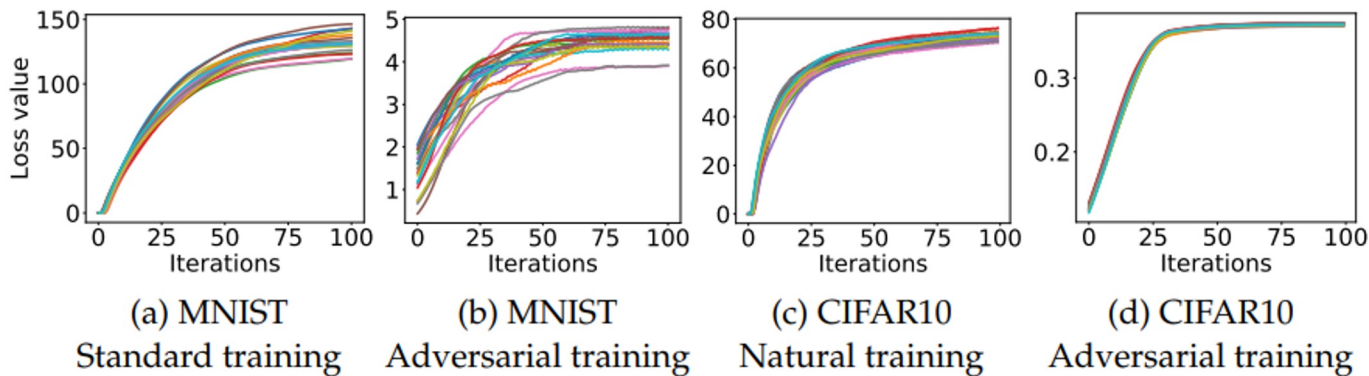
$$x^{t+1} = \Pi_{x+\mathcal{S}} (x^t + \alpha \operatorname{sgn}(\nabla_x L(\theta, x, y)))$$

Adversarial Training

- The core concept of adversarial training is to train a model against adversarial samples in order to improve robustness.
- From my reading it seems they replace the training set with adversarial examples once and then train on them, but this is unclear.

Towards Universally Robust Networks

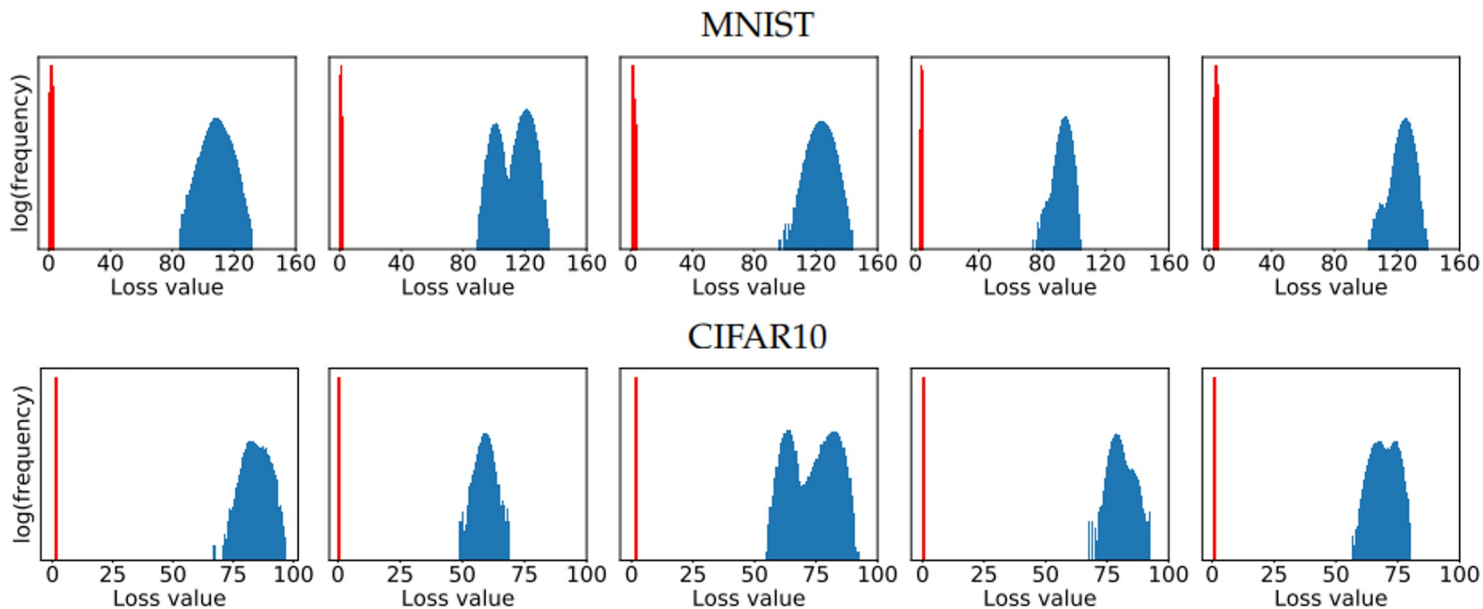
- The paper concedes that both the inner and outer maximization problems are highly non-convex, and thus exact solutions are intractable.
- In spite of this they show empirically that PGD can find good approximations to the inner maximization problem.



Here they run PGD multiple times with different starting points on the same natural example.

More Experiments

Here they attack vanilla (blue) and adversarially trained (red) models with PGD on 5 different samples. They run the experiment 10^5 times and plot the histograms of the loss.



Limitations

- Writing
 - Hinges too heavily on the idea of PGD being the “ultimate first-order attack”
 - Claims of guarantees are a little misleading
- Experiments
 - Not too many complaints here
- Method
 - The solution is approximate, and requires PGD and SGD to do a lot of heavy lifting
 - Can slow down training time significantly

Strengths

- Writing
 - The authors explore their idea fairly exhaustively, and acknowledge the paper's limitations
 - The paper is generally insightful
- Experiments
 - Lots of experiments exploring the different dimensions of their method
 - Experiments empirically verify their method, in some way, and provide further insight
- Methods
 - The method is clearly very effective and not too difficult to implement
 - Forms a good building block for further research into robustness

Discussion and Questions