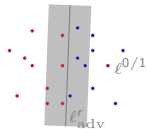


On the (Un-)Avoidability of Adversarial Examples

Sadia Chowdhury, Ruth Uerner. ICML2021 Workshop on Theoretic Foundation, Criticism, and Application Trend of Explainable AI

Phenomenon:



- Adversarial examples pose safety concerns.
- The usual formulation of the adversarial loss

$$\ell_{adv}^r(h, x, y) = 1 \left[\exists x' \in \mathcal{B}_r(x) : h(x') \neq y \right]$$

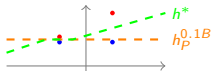
leads to inconsistencies with accuracy.



$$\text{err}_h = 0, \text{mar}_h^r = 1 \\ \Rightarrow \mathcal{L}_P^r(h) = 1$$

$$\text{err}_g = 0.5, \text{mar}_g^r = 0 \\ \Rightarrow \mathcal{L}_P^r(g) = 0.5$$

Question: Is it capturing what was intended?



⇒ We really want locally maximal robustness!

Our contributions:

- We define the notion of a **margin canonical Bayes classifier** h_P^B .

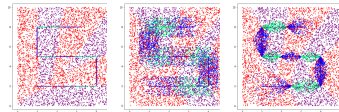


- We propose to re-define the robust loss as a **locally adaptive** requirement with respect to a margin canonical Bayes classifier:

$$\ell_{adv}^r(h, x, y) = 1 \left[h(x) \neq y \vee \mathcal{B}^{h_P^B}(x) \not\subseteq \mathcal{B}^h(x) \right],$$

where $\mathcal{B}^h(x)$ is a largest ball around x that a classifier h labels homogeneously. We also define an **empirical version** of the adaptive robust loss.

- We **introduce an adaptive data augmentation** scheme, evaluate it empirically and **prove** that it **maintains consistency of a nearest neighbor classifier**.



Conclusions:

- Robustness is not at odds with accuracy.
- It's important to analyze (and highlight) to what degree mathematical phenomena are artifacts of a framework for analysis.