

Defenses Against Adversarial Examples

Neil Gong

Certiably robust classifier

- A classifier is (p, ε) -certifiably robust for x , if no adversarial perturbation whose L_p norm is no larger than ε exists.
- Verification
 - Given a classifier and x , verify whether the classifier is (p, ε) -certifiably robust for x
- Certification
 - Given a classifier and x , deriving p and ε

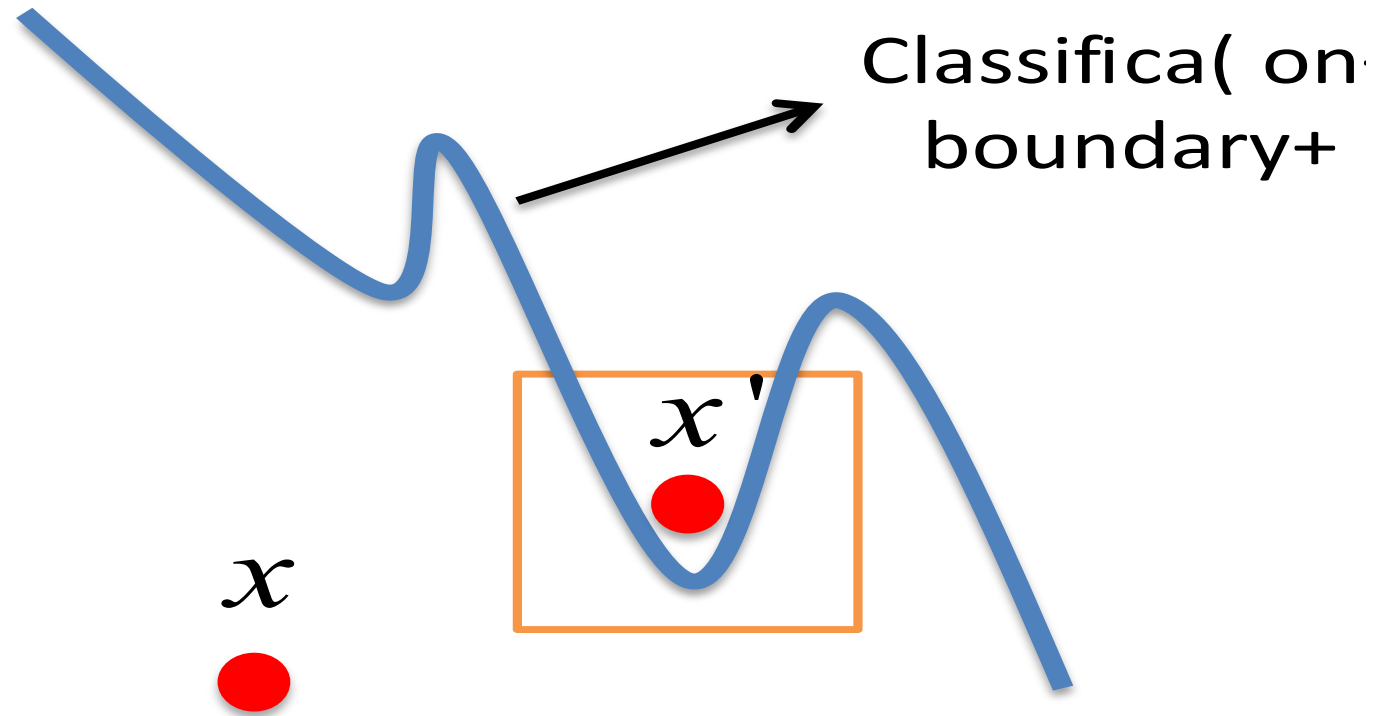
Verification via interval analysis

- Given x , $p=\infty$, ε , we propagate the intervals from the input to the output
- Limitations
 - False negatives
 - Limited to $p=\infty$
 - Not effective for certain classifiers

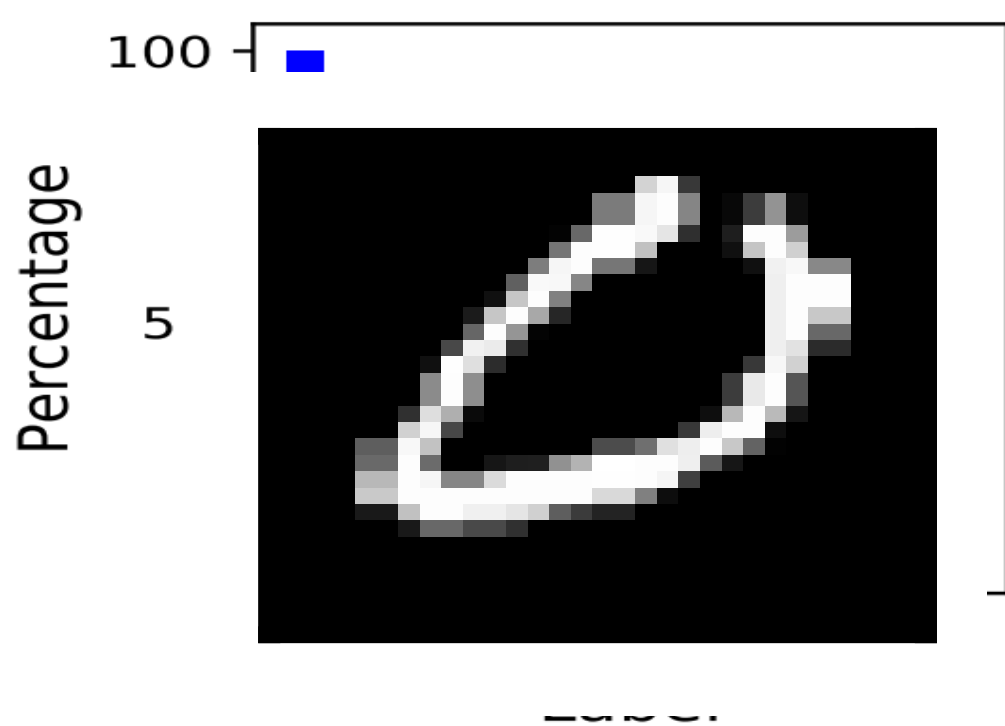
Certification via randomized smoothing

- Given a classifier and x , deriving p and ε
- Many methods have been developed
- Randomized smoothing
 - Applicable to any classifier
 - Scalable to large neural networks

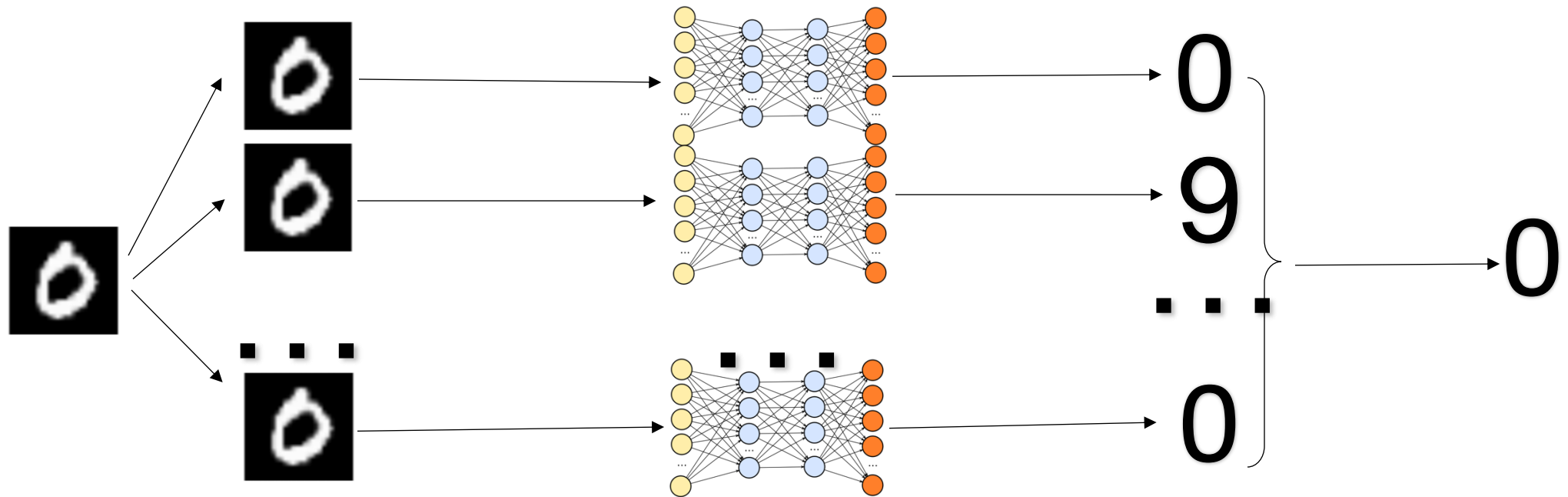
Adversarial example is close to classification boundary?



Measuring Adversarial Examples



Randomized smoothing



Formal definition of randomized smoothing

- Input
 - a classifier f
 - an example x
 - a noise distribution
- Output
 - $g(x) = \operatorname{argmax}_c \Pr(f(x + r) = c)$

Deriving (p, ε)

- Noise is isotropic Gaussian distribution
- $g(x + \delta) = c_A$ when $|\delta|_2 \leq \varepsilon$

$$\varepsilon = \frac{\sigma}{2} (\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B))$$

Certified radius

Tightness of the bound

- Given
 - No assumptions on the classifier f
 - Randomized smoothing with Gaussian noise
- The derived bound is tight

Estimating the label probabilities

- Sampling a large number of noise
- Predicting labels for the noisy examples
- Estimating label probabilities with probabilistic guarantees

Generalization to top-k

- Input
 - a classifier f
 - an example x
 - a noise distribution
- Output
 - $p_c = \Pr(f(x + r) = c)$
 - The smoothed classifier predicts k labels with the largest label probabilities
- A label is among the top- k labels if the adversarial perturbation is bounded

Training to improve certified accuracy

- Adding random noise during training
- Adding certified radius as a regularization term

$$\underbrace{\mathbf{1}_{\{g_{\theta}(x) \neq y\}}}_{0/1 \text{ Classification Error}} + \underbrace{\mathbf{1}_{\{g_{\theta}(x) = y, CR(g_{\theta}; x, y) < \epsilon\}}}_{0/1 \text{ Robustness Error}}$$

Randomized smoothing

- Strengths
 - Applicable to any classifier
 - Scalable to large classifier
- Limitations
 - Efficiency – need many predictions
 - Probabilistic guarantees