# TIGHT CERTIFICATES OF ADVERSARIAL ROBUSTNESS FOR RANDOMLY SMOOTHED CLASSIFIERS

## Guang-He Lee, Yang Yuan, Shiyu Chang, and Tommi S. Jaakkola
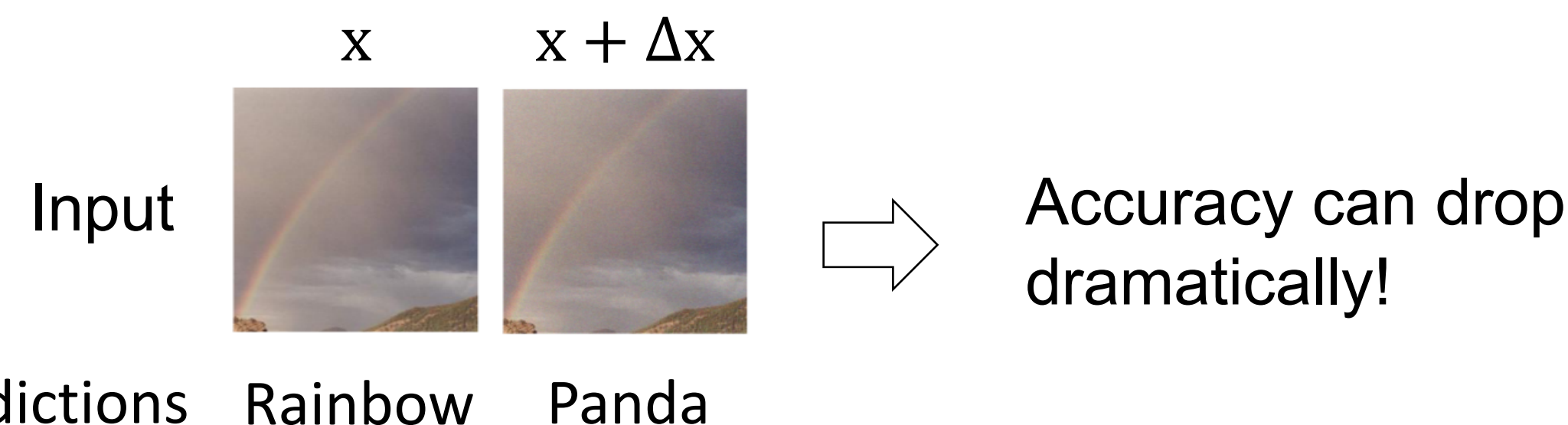
Massachusetts Institute of Technology — CSAIL — IBM

## Summary

o A general approach to deriving tight certificates of robustness for randomly smoothed classifiers.

o We focus on $\ell_0$-robustness in discrete spaces.

o We show how certificates can be tightened with additional assumptions about the classifier.

## Introduction

o Adversarial examples can be easily found on deep models

x          x + Δx

Input

Predictions    Rainbow    Panda

Accuracy can drop dramatically!

  ▪ Ideally, we want a model without adversarial example.

o If a heuristic search algorithm fails, there may still be adversarial examples.

o We need a certificate to show that no such example exists around a specified radius of the input example.

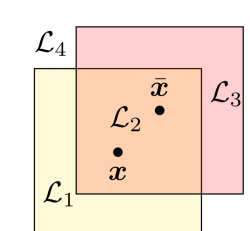o Finding certificates is particularly challenging in discrete spaces as the problem is combinatorial.

## Set-up & background

o Given an input $x \in \mathcal{X}$, a randomization scheme $\phi$ assigns a distribution $\Pr(\phi(x) = z)$ for each $z \in \mathcal{X}$.

o We use a randomly smoothed classifier $f(\phi(x))$.

  ▪ $f$ is a base classifier (e.g., a deep net / decision tree).

  ▪ $\Pr(f(\phi(x)) = y)$ is abbreviated as $p$.

o Tight certificates exist with Gaussian randomization and $\ell_2$ metric (Cohen et al., 19').

## Our framework

o A tight point-wise robustness certificate for $\bar{x}$:
$$\rho_{x,\bar{x}}(p) \triangleq \min_{\bar{f} \in \mathcal{F}: \Pr(\bar{f}(\phi(x)) = y) = p} \Pr(\bar{f}(\phi(\bar{x})) = y)$$
$$\leq \Pr(f(\phi(\bar{x})) = y)$$

  ▪ It can be solved by Neyman-Pearson lemma

o A regional certificate of robustness:

  ▪ Define $\mathcal{B}_{r,q}(x) \triangleq \{\bar{x} \in \mathcal{X} : \|x - \bar{x}\|_q \leq r\}$

  ▪ $R(x, p, q) \triangleq \sup r \ s.t. \min_{\bar{x} \in \mathcal{B}_{r,q}(x)} \rho_{x,\bar{x}}(p) > 0.5$

  ▪ Implication: if $\Pr(f(\phi(x)) = y) = p$, then
$$\forall \bar{x} \in \mathcal{X} : \|x - \bar{x}\|_q < R(x, p, q),$$
$$\Pr(f(\phi(\bar{x})) = y) > 0.5$$

## A warm-up example

o A uniform randomization scheme:
$$\phi(x)_i = x_i + \epsilon_i, \epsilon_i \overset{i.i.d.}{\sim} \text{Uniform}([-\gamma, \gamma])$$

o Illustration:

o Randomization at $x$ and $\bar{x}$ divide the input space into non-overlapping regions $\mathcal{L}_1, \ldots, \mathcal{L}_4$ based on likelihood comparisons

o For any $f$ or $\bar{f}$, only the integral over a region matters; we search for $\bar{f}$ that assigns prob. [0,1] (integral value) to each region.

o Worst case $\bar{f}$ assigns high values to $\mathcal{L}_1$, low values to $\mathcal{L}_2$ and $\mathcal{L}_3$, subject to the constraint that the aggregate = $p$ across $\mathcal{L}_1$ and $\mathcal{L}_2$.

$$\Rightarrow \begin{cases} \rho_{x,\bar{x}}(p) = 0, & \text{if } 0 \leq p \leq (2\gamma)^{-d}|\mathcal{L}_1|, \\ \rho_{x,\bar{x}}(p) = p - (2\gamma)^{-d}|\mathcal{L}_1|, & \text{otherwise.} \end{cases}$$

o Regional certificate finds the worst $\bar{x} \in \mathcal{B}_{r,q}(x)$ such that $|\mathcal{L}_1|$ is maximized.

$$\Rightarrow \begin{array}{l} R(x, p, q = 1) = 2p\gamma - \gamma \\ R(x, p, q = \infty) = 2\gamma - 2\gamma(1.5 - p)^{1/d}. \end{array}$$

## A discrete distribution for $\ell_0$ robustness

o We consider the discrete space: $\mathcal{X} = \{0, \frac{1}{K}, \frac{2}{K}, \ldots, 1\}^d$.

o A discrete randomization scheme:
$$\begin{cases} \Pr(\phi(x)_i = x_i) = \alpha, \\ \Pr(\phi(x)_i = z) = (1-\alpha)/K \triangleq \beta \in (0, 1/K), & \text{if } z \in \{0, \frac{1}{K}, \frac{2}{K}, \ldots, 1\} \text{ and } z \neq x_i \end{cases}$$

o Key properties:

  > 1. for all $x, \bar{x}$ such that $\|x - \bar{x}\|_0 = r$, we have $\rho_{x,\bar{x}} = \rho_r$
  > 2. $\rho_r : [0,1] \rightarrow [0,1]$ is an increasing bijection

o Implications:

  • We can pre-compute $\rho_r^{-1}(0.5)$ (we have a $\Theta(d^3)$ algorithm).

  • If $p > \rho_r^{-1}(0.5)$, the prediction is robust within $\mathcal{B}_{r,0}(x)$.

  • $R(x, p, q)$ is simply the maximum $r$ s.t. $p > \rho_r^{-1}(0.5)$.

o Key steps for pre-computing $\rho_r^{-1}(0.5)$

  ▪ Similar to the uniform distribution, we partition the space into regions with constant likelihood ratio to simplify the problem.

    • Likelihood ratio: $\Pr(\phi(x) = z)/\Pr(\phi(\bar{x}) = z)$.

  ▪ Assigning $\bar{f}(z)$ to $y$ in ↓ likelihood ratio computes $\rho_r^{-1}(0.5)$ (Neyman-Pearson lemma. It can be done in $\Theta(d^3)$).

  ▪ A large integer algorithm is needed for high dimension setting.

## Towards tighter certification

o The certificates are tight w.r.t. measurable classifiers.

o More characterization of $f$ always improves the point-wise (and regional) certificates: if $f \in \mathcal{F}_\zeta \subset \mathcal{F}$,
$$\min_{\bar{f} \in \mathcal{F}_\zeta: \Pr(\bar{f}(\phi(x)) = y) = p} \Pr(\bar{f}(\phi(\bar{x})) = y) \geq \min_{\bar{f} \in \mathcal{F}: \Pr(\bar{f}(\phi(x)) = y) = p} \Pr(\bar{f}(\phi(\bar{x})) = y)$$

o Example 1: when $\phi(x)_i = x_i + \epsilon_i, \epsilon_i \overset{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

  ▪ If $\mathcal{X} = \{0, 1\}^d$, we can use (Cohen et al., 19) to derive $\ell_0$ certificates due to the bijection to $\ell_2$.

  ▪ If we apply denoising before feeding to model:
$$\zeta(\phi(x))_i = \mathbb{I}\{\phi(x)_i > 0.5\}, \forall i \in [d]$$

    • The resulting input is equivalent to our discrete randomization scheme.

  ▪ Our certificate is always tighter than using the one derived from the Gaussian distribution in this case.

o Example 2: when $f$ is a decision tree:

  ▪ The randomization can be expressed as a probabilistic routing scheme for each decision node.

  ▪ The exact certificate of robustness can be computed using dynamic programming over tree nodes.

## Experiment (project page: http://people.csail.mit.edu/guanghe/randomized_smoothing)

o Evaluation metrics:

  ▪ $\mu(R)$: the average certified radius in testing set.

  ▪ ACC@r: guaranteed accuracy within $\ell_0$ radius $r$.

o Binarized MNIST (CNN model).

| $\phi$ | Certificate | $\mu(R)$ | ACC@r | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | $r=1$ | $r=2$ | $r=3$ | $r=4$ | $r=5$ | $r=6$ | $r=7$ |
| Discrete | Discrete | 3.456 | 0.921 | 0.774 | 0.539 | 0.524 | 0.357 | 0.202 | 0.097 |
| Discrete | Gaussian | 1.799 | 0.830 | 0.557 | 0.272 | 0.119 | 0.021 | 0.000 | 0.000 |
| Gaussian | Gaussian | 2.378 | 0.884 | 0.701 | 0.464 | 0.252 | 0.078 | 0.000 | 0.000 |

o (Discrete) Exact ACC@1 = 0.954, ACC@2 = 0.926.

o ImageNet (ResNet50 model).

| $\phi$ and certificate | ACC@r | | | | | | |
|---|---|---|---|---|---|---|---|
| | $r=1$ | $r=2$ | $r=3$ | $r=4$ | $r=5$ | $r=6$ | $r=7$ |
| Discrete | 0.538 | 0.394 | 0.338 | 0.274 | 0.234 | 0.190 | 0.176 |
| Gaussian | 0.372 | 0.292 | 0.226 | 0.194 | 0.170 | 0.154 | 0.138 |

$\rho_r^{-1}(0.5)$ for an $\alpha$

The certified accuracy for an $\alpha$