

Black-Box Certification with Randomized Smoothing: A Functional Optimization Based Framework

Dinghuai Zhang*, Mao Ye*, Chengyue Gong*, Zhanxing Zhu,
Qiang Liu

Notation

Certification means a *guarantee* that a classifier won't change its prediction when perturbing input under some condition.

For simplicity, we consider a binary classification setting.

- | $f^J: \mathbb{R}^d \rightarrow [0, 1]$ a given binary classifier
output the probability of "positive class"
- | $f_o^J(x_0) := \mathbb{E}_z \left[f^J(x_0 + z) \right]$
randomized smoothed classifier
- | $\Phi(\cdot)$ the cdf of standard Gaussian

Cohen [2]'s result

Theorem

For a randomized smoothing classifier with Gaussian noise $z \sim \pi_0 = \mathcal{N}(0, \sigma^2 I)$, suppose $f_0^J(x_0) = p_0 \geq \frac{1}{2}$, then $f_0^J(x_0 + \delta) \geq \frac{1}{2}$ for all $\|\delta\|_2 \leq \sigma \Phi^{-1}(p_0)$.

Cohen et al. use some results from NP lemma to prove this bound.

Our Approach

The only thing we know about the given classifier is p_0 . Naturally one would think: from all classifier f with $f_0(x_0) = p_0$, which f achieves the lowest probability of $f_0(x_0 + \delta)$?

$$\begin{aligned} & \min_{f \in \mathcal{F}} \min_{\delta \in \mathcal{B}} f_0(x_0 + \delta) \\ \text{s.t. } & f_0(x_0) = f_0^J(x_0) := p_0 \end{aligned}$$

Functional Optimization

From

$$\begin{aligned} & \min_{f \in \mathcal{F}} \min_{B \in \mathcal{B}} f_0(x_0 + \cdot) \\ \text{s.t.} \quad & f_0(x_0) = f_0^J(x_0) := \rho_0 \end{aligned}$$

we write out the Lagrangian:

$$V_0(F, B) = \min_{f \in \mathcal{F}} \min_{B \in \mathcal{B}} \max_{\lambda \in \mathcal{R}} \left\{ f_0(x_0 + \cdot) - \lambda (f_0(x_0) - \rho_0) \right\}$$

Define $\pi(z) = N(\mu, \sigma^2 I)$

$$\begin{aligned}
 & V_0(F, B) \\
 &= \min_{f \in \mathcal{F}} \min_{B} \max_{R} f(x_0 + z) - \lambda(f(x_0) - \rho_0)g \\
 &= \min_{B} \min_{f \in \mathcal{F}} \max_{z \sim \pi} \mathbb{E}_z [f(x_0 + z)] + \lambda(\rho_0 - \mathbb{E}_z [f(x_0 + z)])g \\
 &= \min_{B} \max_{\rho_0} \left\{ \lambda \rho_0 + \min_{f \in \mathcal{F}} \mathbb{E}_z [f(x_0 + z)] - \lambda \mathbb{E}_z [f(x_0 + z)] \right\} \\
 &=?
 \end{aligned}$$

Specify our setting:

$$F_{[0,1]} = \left\{ f : f(z) \in [0, 1], \forall z \in \mathbb{R}^d \right\}$$

$$B = \{ f : k, k_2, r, g \}$$

Thus our bound become

$$\lambda \rho_0 \min_{f \in F_{[0,1]}} \left\{ \lambda \int f(x_0 + z) \pi_0(z) dz - \int f(x_0 + z) \pi(z) dz \right\}$$

$$= \lambda \rho_0 \min_{f \in F_{[0,1]}} \left\{ \int f(x_0 + z) (\lambda \pi_0(z) - \pi(z)) dz \right\}$$

$$= \lambda \rho_0 \int (\lambda \pi_0(z) - \pi(z))_+ dz$$

Total Variation

Bonus

For two distributions q_1, q_2 ,

$$\int (q_1(z) - q_2(z))_+ dz = TV(q_1, q_2)$$

Thus $\int (\lambda \pi_0(z) - \pi(z))_+ dz$ can also be seen as some sort of discrepancy.

Our Bound is Equivalent with Cohen's

Proposition

$$\max_{0 < k < k_2} \min_r \left\{ \lambda p_0 \int (\lambda \pi_0(z) - \pi(z))_+ dz \right\} = \Phi(\Phi^{-1}(p_0) - r/\sigma)$$

Thus

Confidence Lower Bound 0.5 , $r = \sigma \Phi^{-1}(p_0)$!!

Remark

The min max can switch order.

Notation for ℓ_1 Setting

$$B = \{z : \|z\|_1 \leq r\}$$

$$\pi_0(z) \propto \exp\left(-\frac{\|z\|_1}{\sigma}\right)$$

$$\pi(z) \propto \exp\left(-\frac{\|z\|_1}{\sigma}\right)$$

Optimization with Laplacian Smoothing

Still, we have

$$\min_{2B} \max_0 \left\{ \lambda \rho_0 \int (\lambda \pi_0(z) - \pi(z))_+ dz \right\}$$

Proposition

the bound = $\frac{1}{2} \exp(-\log[2(1 - \rho_0)] - \frac{r}{\sigma})$

Thus lower bound = $0.5 \cdot \exp(-\log[2(1 - \rho_0)] - \frac{r}{\sigma})$

This is the same as [1].

Outline

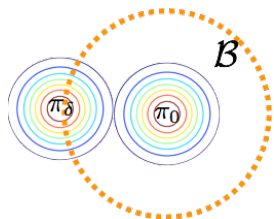
Framework: Constrained Adversarial Certification

Filling the Soap Bubbles

Motivation

Bound Decomposition

$$\max_0 \left[\underbrace{\lambda \rho_0}_{\text{Accuracy}} \underbrace{\max_{2B} \int (\lambda \pi_0(z) - \pi(z))_+ dz}_{\text{Robustness}} \right]$$



New distribution can improve:

- | More "center-massed" distribution can boost the accuracy term
- | A heavy tail can boost the robustness term

Gaussian Issues

For high dimensional Gaussian distribution, the samples will concentrate on a "soap bubble":

$$\|z\|_2^2 = d \frac{\sum_i z_i^2}{d} \approx d \sigma^2$$

Samples far from center will cause a small "accuracy" term!

Filling the Soap Bubbles

We propose a new Centripetal distribution family:

$$\pi_{\mathbf{0}}(z) \propto \|z\|_2^k \exp\left(-\frac{\|z\|_2^2}{2\sigma^2}\right)$$

The distribution of its radius is

$$p_{\|z\|_2}(r) \propto r^{d-k-1} \exp\left(-\frac{r^2}{2\sigma^2}\right)$$

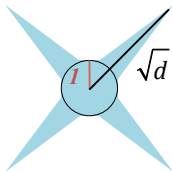
l_2 RADIUS	0.25	0.5	0.75	1.0	1.25	1.5	1.75	2.0	2.25
baseline (%)	60	43	34	23	17	14	12	10	8
OURS (%)	61	46	37	25	19	16	14	11	9

Table: Certified top-1 accuracy with various l_2 radius on CIFAR-10.

l_2 RADIUS	0.5	1.0	1.5	1.0	2.0	2.5	3.0
baseline (%)	49	37	29	19	15	12	9
OURS (%)	50	39	31	21	17	13	10

Table: Certified top-1 accuracy with various l_2 radius on ImageNet.

Cracking ℓ_1

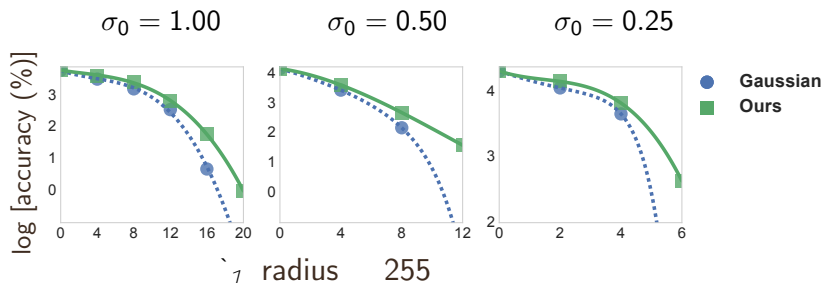


Infinite norm setting is very challenging, we propose the following centripetal distribution:

$$\pi_{\mathbf{0}}(z) \propto \|z\|_1^k \exp\left(-\frac{\|z\|_2^2}{2\sigma^2}\right)$$



l_1 RADIUS	2/255	4/255	6/255	8/255	10/255	12/255
baseline (%)	58	42	31	25	18	13
OURS (%)	60	47	38	32	23	17

Table: Certified top-1 accuracy with various l_∞ radius on CIFAR-10.



Thank you for listening!

References I

-  ANONYMOUS, *\$ell_1\$ adversarial robustness certificates: a randomized smoothing approach*, in Submitted to International Conference on Learning Representations, 2020.
under review.
-  J. M. COHEN, E. ROSENFELD, AND J. Z. KOLTER, *Certified adversarial robustness via randomized smoothing*, arXiv preprint arXiv:1902.02918, (2019).