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

Dinghuai Zhang*
Mila
dinghuai.zhang@mila.quebec

Mao Ye*, Chengyue Gong*
Department of Computer Science
University of Texas at Austin
{my21, cygong}@cs.utexas.edu

Zhanxing Zhu
School of Mathematical Sciences
Peking University
zhanxing.zhu@pku.edu.cn

Qiang Liu
Department of Computer Science
University of Texas at Austin
lqiang@cs.utexas.edu

Abstract

Randomized classifiers have been shown to provide a promising approach for achieving certified robustness against adversarial attacks in deep learning. However, most existing methods only leverage Gaussian smoothing noise and only work for ℓ_2 perturbation. We propose a general framework of adversarial certification with non-Gaussian noise and for more general types of attacks, from a unified functional optimization perspective. Our new framework allows us to identify a key trade-off between accuracy and robustness via designing smoothing distributions and leverage it to design new families of non-Gaussian smoothing distributions that work more efficiently for different ℓ_p settings, including ℓ_1 , ℓ_2 and ℓ_∞ attacks. Our proposed methods achieve better certification results than previous works and provide a new perspective on randomized smoothing certification.

1 Introduction

Although many robust training algorithms have been developed to overcome adversarial attacks [1, 2, 3], most heuristically developed methods can be shown to be broken by more powerful adversaries eventually (e.g., [4, 5, 6, 7]). This casts an urgent demand for developing robust classifiers with provable worst-case guarantees. One promising approach for certifiable robustness is the recent *randomized smoothing method* [8, 9, 10, 11, 12, 13, 14, 15], which constructs smoothed classifiers with certifiable robustness by introducing noise on the inputs. Compared with the other more traditional certification approaches [16, 17, 18] that exploit special structures of the neural networks (such as the properties of ReLU), the randomized smoothing approaches work more flexibly on general black-box classifiers and is shown to be more scalable and provide tighter bounds on challenging datasets such as ImageNet [19].

Most existing methods use Gaussian noise for smoothing. Although appearing to be a natural choice, one of our key observations is that the Gaussian distribution is, in fact a sub-optimal choice in high dimensional spaces even for ℓ_2 attack. We observe that there is a counter-intuitive phenomenon in high dimensional spaces [20], that almost all of the probability mass of standard Gaussian distribution concentrates around the sphere surface of a certain radius. This makes tuning the variance of Gaussian distribution an inefficient way to trade off robustness and accuracy for randomized smoothing.

*Equal contributions

Our Contributions We propose a general framework of adversarial certification using non-Gaussian smoothing noises, based on a new functional optimization perspective. Our framework unifies the methods of [9] and [14] as special cases, and is applicable to more general smoothing distributions and more types of attacks beyond ℓ_2 -norm setting. Leveraging our insight, we develop a new family of distributions for better certification results on ℓ_1 , ℓ_2 and ℓ_∞ attacks. An efficient computational approach is developed to enable our method in practice. Empirical results show that our new framework and smoothing distributions outperform existing approaches for ℓ_1 , ℓ_2 and ℓ_∞ attacking, on datasets such as CIFAR-10 and ImageNet.

2 Related Works

Certified Defenses Unlike the empirical defense methods, once a classifier can guarantee a consistent prediction for input within a local region, it is called a certified-robustness classifier. *Exact* certification methods provide the minimal perturbation condition which leads to a different classification result. This line of work focuses on deep neural networks with ReLU-like activation that makes the classifier a piece-wise linear function. This enables researchers to introduce satisfiability modulo theories [21, 22] or mix integer linear programming [23, 24]. *Sufficient* certification methods take a conservative way and bound the Lipschitz constant or other information of the network [18, 16, 25, 26]. However, these certification strategies share a drawback that they are not feasible on large-scale scenarios, *e.g.* large and deep networks and datasets.

Randomized Smoothing To mitigate this limitation of previous certifiable defenses, improving network robustness via randomness has been recently discussed [27, 28]. [8] first introduced randomization with technique in differential privacy. [12] improved their work with a bound given by Rényi divergence. In succession, [9] firstly provided a *tight* bound for *arbitrary* Gaussian smoothed classifiers based on previous theorems found by [29]. [10] combined the empirical and certification robustness, by applying adversarial training on randomized smoothed classifiers to achieve a higher certified accuracy. [11] focused on ℓ_0 norm perturbation setting, and proposed a discrete smoothing distribution which can be shown perform better than the widely used Gaussian distribution. [14] took a similar statistical testing approach with [9], utilizing Laplacian smoothing to tackle ℓ_1 certification problem. [15] extended the approach of [9] to a top-k setting. [13] extends the total variant used by [9] to f -divergences. Recent works [30, 31, 32] discuss further problems about certification methods. We also focus on a generalization of randomized smoothing, but with a different view on loosening the constraint on classifier.

Noticeably, [30] also develops analysis on ℓ_1 setting and provide a thorough theoretical analysis on many kinds of randomized distribution. We believe the [30] and ours have different contributions and were developed concurrently. [30] derives the optimal shapes of level sets for ℓ_p attacks based on the Wulff Crystal theory, while our work, based on our functional-optimization framework and accuracy-robustness decomposition (Eq.9), proposes to use distribution that is more concentrated toward the center. Besides, we also consider a novel distribution using mixed ℓ_2 and ℓ_∞ norm for ℓ_∞ adversary, which hasn't been studied before and improve the empirical results.

3 Black-box Certification as Functional Optimization

3.1 Background

Adversarial Certification For simplicity, we consider binary classification of predicting binary labels $y \in \{0, 1\}$ given feature vectors $x \in \mathbb{R}^d$. The extension to multi-class cases is straightforward, and is discussed in Appendix C. We assume $f^\sharp: \mathbb{R}^d \rightarrow [0, 1]$ is a given binary classifier (\sharp means the classifier is *given*), which maps from the input space \mathbb{R}^d to either the positive class probability in interval $[0, 1]$ or binary labels in $\{0, 1\}$. In the robustness certification problem, a testing data point $x_0 \in \mathbb{R}^d$ is given, and one is asked to verify if the classifier outputs the same prediction when the input x_0 is perturbed arbitrarily in \mathcal{B} , a given neighborhood of x_0 . Specifically, let \mathcal{B} be a set of possible perturbation vectors, *e.g.*, $\mathcal{B} = \{\delta \in \mathbb{R}^d : \|\delta\|_p \leq r\}$ for ℓ_p norm with a radius r . If the classifier predicts $y = 1$ on x_0 , i.e. $f^\sharp(x_0) > 1/2$, we want to verify if $f^\sharp(x_0 + \delta) > 1/2$ still holds for any $\delta \in \mathcal{B}$. Through this paper, we consider the most common adversarial settings: ℓ_1 , ℓ_2 and ℓ_∞ attacks.

Black-box Randomized Smoothing Certification Directly certifying f^\sharp heavily relies on the smooth property of f^\sharp , which has been explored in a series of prior works [16, 8]. These methods typically depend on the special structure-property (*e.g.*, the use of ReLU units) of f^\sharp , and thus can not serve as general-purpose algorithms for any type of networks. Instead, We are interested in *black-box* verification methods that could work for *arbitrary* classifiers. One approach to enable this, as explored in recent works [9, 11], is to replace f^\sharp with a smoothed classifier by convolving it with Gaussian noise, and verify the *smoothed* classifier.

Specifically, assume π_0 is a smoothing distribution with zero mean and bounded variance, *e.g.*, $\pi_0 = \mathcal{N}(\mathbf{0}, \sigma^2)$. The randomized smoothed classifier is defined by

$$f_{\pi_0}^\sharp(\mathbf{x}_0) := \mathbb{E}_{\mathbf{z} \sim \pi_0} [f^\sharp(\mathbf{x}_0 + \mathbf{z})],$$

which returns the averaged probability of $\mathbf{x}_0 + \mathbf{z}$ under the perturbation of $\mathbf{z} \sim \pi_0$. Assume we replace the original classifier with $f_{\pi_0}^\sharp$, then the goal becomes certifying $f_{\pi_0}^\sharp$ using its inherent smoothness. Specifically, if $f_{\pi_0}^\sharp(\mathbf{x}_0) > 1/2$, we want to certify that $f_{\pi_0}^\sharp(\mathbf{x}_0 + \boldsymbol{\delta}) > 1/2$ for every $\boldsymbol{\delta} \in \mathcal{B}$, that is, we want to certify that

$$\min_{\boldsymbol{\delta} \in \mathcal{B}} f_{\pi_0}^\sharp(\mathbf{x}_0 + \boldsymbol{\delta}) = \min_{\boldsymbol{\delta} \in \mathcal{B}} \mathbb{E}_{\mathbf{z} \sim \pi_0} [f^\sharp(\mathbf{x}_0 + \mathbf{z} + \boldsymbol{\delta})] > \frac{1}{2}. \quad (1)$$

In this case, it is sufficient to obtain a *guaranteed lower bound* of $\min_{\boldsymbol{\delta} \in \mathcal{B}} f_{\pi_0}^\sharp(\mathbf{x}_0 + \boldsymbol{\delta})$ and check if it is larger than $1/2$. When π_0 is Gaussian $\mathcal{N}(\mathbf{0}, \sigma^2)$ and for ℓ_2 attack, this problem was studied in [9], which shows that a lower bound of

$$\min_{\mathbf{z} \in \mathcal{B}} \mathbb{E}_{\mathbf{z} \sim \pi_0} [f^\sharp(\mathbf{x}_0 + \mathbf{z})] \geq \Phi(\Phi^{-1}(f_{\pi_0}^\sharp(\mathbf{x}_0)) - \frac{r}{\sigma}), \quad (2)$$

where $\Phi(\cdot)$ is the cumulative density function (CDF) of standard Gaussian distribution. The proof of this result in [9] uses Neyman-Pearson lemma [29]. In the following section, we will show that this bound is a special case of the proposed functional optimization framework for robustness certification.

3.2 Constrained Adversarial Certification

We propose a **constrained adversarial certification (CAC)** framework, which yields a guaranteed lower bound for Eq.1. The main idea is simple: assume \mathcal{F} is a function class which is known to include f^\sharp , then the following optimization immediately yields a guaranteed lower bound

$$\min_{\boldsymbol{\delta} \in \mathcal{B}} f_{\pi_0}^\sharp(\mathbf{x}_0 + \boldsymbol{\delta}) \geq \min_{f \in \mathcal{F}} \min_{\boldsymbol{\delta} \in \mathcal{B}} \left\{ f_{\pi_0}(\mathbf{x}_0 + \boldsymbol{\delta}) \text{ s.t. } f_{\pi_0}(\mathbf{x}_0) = f_{\pi_0}^\sharp(\mathbf{x}_0) \right\}, \quad (3)$$

where we define $f_{\pi_0}(\mathbf{x}_0) = \mathbb{E}_{\mathbf{z} \sim \pi_0} [f(\mathbf{x}_0 + \mathbf{z})]$ for any given f . Then we need to search for the minimum value of $f_{\pi_0}(\mathbf{x}_0 + \boldsymbol{\delta})$ for all classifiers in \mathcal{F} that satisfies $f_{\pi_0}(\mathbf{x}_0) = f_{\pi_0}^\sharp(\mathbf{x}_0)$. This obviously yields a lower bound once $f^\sharp \in \mathcal{F}$. If \mathcal{F} includes only f^\sharp , then the bound is exact, but is computationally prohibitive due to the difficulty of optimizing $\boldsymbol{\delta}$. The idea is then to choose \mathcal{F} properly to incorporate rich information of f^\sharp , while allowing us to calculate the lower bound in Eq.3 computationally tractably. In this paper, we consider the set of all functions bounded in $[0, 1]$, namely

$$\mathcal{F}_{[0,1]} = \left\{ f : f(\mathbf{z}) \in [0, 1], \forall \mathbf{z} \in \mathbb{R}^d \right\}, \quad (4)$$

which guarantees to include all f^\sharp by definition.

Denote by $\mathcal{L}_{\pi_0}(\mathcal{F}, \mathcal{B})$ the lower bound in Eq.3. We can rewrite it into the following minimax form using the Lagrangian function,

$$\mathcal{L}_{\pi_0}(\mathcal{F}, \mathcal{B}) = \min_{f \in \mathcal{F}} \min_{\boldsymbol{\delta} \in \mathcal{B}} \max_{\lambda \in \mathbb{R}} L(f, \boldsymbol{\delta}, \lambda) \triangleq \min_{f \in \mathcal{F}} \min_{\boldsymbol{\delta} \in \mathcal{B}} \max_{\lambda \in \mathbb{R}} \left\{ f_{\pi_0}(\mathbf{x}_0 + \boldsymbol{\delta}) - \lambda(f_{\pi_0}(\mathbf{x}_0) - f_{\pi_0}^\sharp(\mathbf{x}_0)) \right\}, \quad (5)$$

where λ is the Lagrangian multiplier. Exchanging the min and max yields the following dual form.

Theorem 1. I) (Dual Form) Denote by π_δ the distribution of $\mathbf{z} + \delta$ when $\mathbf{z} \sim \pi_0$. Assume \mathcal{F} and \mathcal{B} are compact set. We have the following lower bound of $\mathcal{L}_{\pi_0}(\mathcal{F}, \mathcal{B})$:

$$\mathcal{L}_{\pi_0}(\mathcal{F}, \mathcal{B}) \geq \max_{\lambda \geq 0} \min_{f \in \mathcal{F}} \min_{\delta \in \mathcal{B}} L(f, \delta, \lambda) = \max_{\lambda \geq 0} \left\{ \lambda f_{\pi_0}^\#(\mathbf{x}_0) - \max_{\delta \in \mathcal{B}} \mathbb{D}_{\mathcal{F}}(\lambda \pi_0 \parallel \pi_\delta) \right\}, \quad (6)$$

where we define the discrepancy term $\mathbb{D}_{\mathcal{F}}(\lambda \pi_0 \parallel \pi_\delta)$ as

$$\max_{f \in \mathcal{F}} \left\{ \lambda \mathbb{E}_{\mathbf{z} \sim \pi_0} [f(\mathbf{x}_0 + \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim \pi_\delta} [f(\mathbf{x}_0 + \mathbf{z})] \right\},$$

which measures the difference of $\lambda \pi_0$ and π_δ by seeking the maximum discrepancy of the expectation for $f \in \mathcal{F}$. As we will show later, the bound in (6) is computationally tractable with proper $(\mathcal{F}, \mathcal{B}, \pi_0)$.

II) When $\mathcal{F} = \mathcal{F}_{[0,1]} := \{f: f(x) \in [0, 1], x \in \mathbb{R}^d\}$, we have in particular

$$\mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda \pi_0 \parallel \pi_\delta) = \int (\lambda \pi_0(\mathbf{z}) - \pi_\delta(\mathbf{z}))_+ d\mathbf{z},$$

where $(t)_+ = \max(0, t)$. Furthermore, we have $0 \leq \mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda \pi_0 \parallel \pi_\delta) \leq \lambda$ for any π_0, π_δ and $\lambda > 0$. Note that $\mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda \pi_0 \parallel \pi_\delta)$ coincides with the total variation distance between π_0 and π_δ when $\lambda = 1$.

III) (Strong duality) Suppose $\mathcal{F} = \mathcal{F}_{[0,1]}$ and suppose that for any $\lambda \geq 0$, $\min_{\delta \in \mathcal{B}} \min_{f \in \mathcal{F}_{[0,1]}} L(f, \delta, \lambda) = \min_{f \in \mathcal{F}_{[0,1]}} L(f, \delta^*, \lambda)$, for some $\delta^* \in \mathcal{B}$, we have

$$\mathcal{L}_{\pi_0}(\mathcal{F}, \mathcal{B}) = \max_{\lambda \geq 0} \min_{\delta \in \mathcal{B}} \min_{f \in \mathcal{F}} L(f, \delta, \lambda).$$

Remark We will show later that the proposed methods and the cases we study satisfy the condition in part III of the theorem and thus all the lower bounds of the proposed method are tight.

Proof is deferred to Appendix A.1. Although the lower bound in Eq.6 still involves an optimization on δ and λ , both of them are much easier than the original adversarial optimization in Eq.1. With proper choices of \mathcal{F}, \mathcal{B} and π_0 , the optimization of δ can be shown to provide simple closed-form solutions by exploiting the symmetry of \mathcal{B} , and the optimization of λ is a very simple one-dimensional searching problem.

As corollaries of Theorem 1, we can exactly recover the bound derived by [14] and [9] under our functional optimization framework, different from their original Neyman-Pearson lemma approaches.

Corollary 1. With Laplacian noise $\pi_0(\cdot) = \text{Laplace}(\cdot; b)$, where $\text{Laplace}(x; b) = \frac{1}{(2b)^d} \exp(-\frac{\|x\|_1}{b})$, ℓ_1 adversarial setting $\mathcal{B} = \{\delta: \|\delta\|_1 \leq r\}$ and $\mathcal{F} = \mathcal{F}_{[0,1]}$, the lower bound in Eq.6 becomes

$$\max_{\lambda \geq 0} \left\{ \lambda f_{\pi_0}^\#(\mathbf{x}_0) - \max_{\|\delta\|_1 \leq r} \mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda \pi_0 \parallel \pi_\delta) \right\} = \begin{cases} 1 - e^{r/b}(1 - f_{\pi_0}^\#(\mathbf{x}_0)), & \text{when } f_{\pi_0}^\#(\mathbf{x}_0) \geq 1 - \frac{1}{2}e^{-r/b}, \\ \frac{1}{2}e^{-\frac{r}{b} - \log[2(1 - f_{\pi_0}^\#(\mathbf{x}_0))]}, & \text{when } f_{\pi_0}^\#(\mathbf{x}_0) < 1 - \frac{1}{2}e^{-r/b}. \end{cases} \quad (7)$$

Thus, with our previous explanation, we obtain $\mathcal{L}_{\pi_0}(\mathcal{F}, \mathcal{B}) \geq \frac{1}{2} \iff r \leq -b \log[2(1 - f_{\pi_0}^\#(\mathbf{x}_0))]$, which is exactly the ℓ_1 certification radius derived by [14]. See Appendix A.2 for proof details. For Gaussian noise setting which has been frequently adopted, we have

Corollary 2. With isotropic Gaussian noise $\pi_0 = \mathcal{N}(\mathbf{0}, \sigma^2 I_{d \times d})$, ℓ_2 attack $\mathcal{B} = \{\delta: \|\delta\|_2 \leq r\}$ and $\mathcal{F} = \mathcal{F}_{[0,1]}$, the lower bound in Eq.6 becomes

$$\max_{\lambda \geq 0} \left\{ \lambda f_{\pi_0}^\#(\mathbf{x}_0) - \max_{\|\delta\|_2 \leq r} \mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda \pi_0 \parallel \pi_\delta) \right\} = \Phi\left(\Phi^{-1}(f_{\pi_0}^\#(\mathbf{x}_0)) - \frac{r}{\sigma}\right). \quad (8)$$

Analogously, we can retrieve the main theoretical result of [9]: $\mathcal{L}_{\pi_0}(\mathcal{F}, \mathcal{B}) \geq \frac{1}{2} \iff r \leq \sigma \Phi^{-1}(f_{\pi_0}^\#(\mathbf{x}_0))$. See Appendix A.3 for proof details.

3.3 Trade-off Between Accuracy and Robustness

The lower bound in Eq.6 reflects an intuitive trade-off between the robustness and accuracy on the certification problem:

$$\max_{\lambda \geq 0} \left[\underbrace{\lambda f_{\pi_0}^\sharp(\mathbf{x}_0)}_{\text{Accuracy}} + \underbrace{\left(-\max_{\delta \in \mathcal{B}} \mathbb{D}_{\mathcal{F}}(\lambda \pi_0 \parallel \pi_\delta) \right)}_{\text{Robustness}} \right], \quad (9)$$

where the first term reflects the accuracy of the smoothed classifier (assuming the true label is $y = 1$), while the second term $-\max_{\delta \in \mathcal{B}} \mathbb{D}_{\mathcal{F}}(\lambda \pi_0 \parallel \pi_\delta)$ measures the robustness of the smoothing method, via the negative maximum discrepancy between the original smoothing distribution π_0 and perturbed distribution π_δ for $\delta \in \mathcal{B}$. The maximization of dual coefficient λ can be viewed as searching for a best balance between these two terms to achieve the largest lower bound.

More critically, different choices of smoothing distributions yields a trade-off between accuracy and robustness in Eq.9. A good choice of the smoothing distribution should ① be centripetal enough to obtain a large $f_{\pi_0}^\sharp(\mathbf{x}_0)$ and ② have large kurtosis or long tail to yield a small $\max_{\delta \in \mathcal{B}} \mathbb{D}_{\mathcal{F}}(\lambda \pi_0 \parallel \pi_\delta)$ discrepancy term. In the next section, we'll show how to design a distribution that could improve both points.

4 Improving Certification Bounds with a New Distribution Family

4.1 “Thin Shell” Phenomenon and New Distribution Family

We first identify a key problem of the usage of Laplacian and Gaussian noise in high dimensional space, due to the “thin shell” phenomenon that the probability mass of them concentrates on a sphere far away from the center points [20].

Proposition 1 ([20], Section 3.1). *Let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, I_{d \times d})$ be a d -dimensional standard Gaussian random variable. Then there exists a constant c , such that for any $\delta \in (0, 1)$, $\text{Prob}(\sqrt{d} - \sqrt{c \log(2/\delta)} \leq \|\mathbf{z}\|_2 \leq \sqrt{d} + \sqrt{c \log(2/\delta)}) \geq 1 - \delta$. See [20] for more discussion.*

This suggests that with high probability, \mathbf{z} takes values very close to the sphere of radius \sqrt{d} , within a constant distance from that sphere. There exists similar phenomenon for Laplacian distribution:

Proposition 2 (Chebyshev bound). *Let \mathbf{z} be a d -dimensional Laplacian random variable, $\mathbf{z} = (z_1, \dots, z_d)$, where $z_i \sim \text{Laplace}(1)$, $i = 1, \dots, d$. Then for any $\delta \in (0, 1)$, we have $\text{Prob}(1 - 1/\sqrt{d\delta} \leq \|\mathbf{z}\|_1/d \leq 1 + 1/\sqrt{d\delta}) \geq 1 - \delta$.*

Although choosing isotropic Laplacian and Gaussian distribution appears to be natural, this “thin shell” phenomenon makes it sub-optimal to use them for adversarial certification, because one would expect that the smoothing distribution should concentrate around the center (the original image) in order to make the smoothed classifier accurate enough in trade-off of Eq.9.

Thus it's desirable to design a distribution more *concentrated* to center. To motivate our new distribution family, it's useful to examine the density function of the distributions of the radius of spherical distributions in general.

Proposition 3. *Assume \mathbf{z} is a symmetric random variable on \mathbb{R}^d with a probability density function (PDF) of form $\pi_0(\mathbf{z}) \propto \phi(\|\mathbf{z}\|)$, where $\phi: [0, \infty) \rightarrow [0, \infty)$ is a univariate function, then the PDF of the norm of \mathbf{z} is $p_{\|\mathbf{z}\|}(r) \propto r^{d-1} \phi(r)$. The term r^{d-1} arises due to the integration on the surface of radius r norm ball in \mathbb{R}^d . Here $\|\cdot\|$ can be any L_p norm.*

In particular, for $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_{d \times d})$, we have $\phi(r) \propto \exp(-r^2/(2\sigma^2))$ and hence $p_{\|\mathbf{z}\|_2}(r) \propto r^{d-1} \exp(-r^2/(2\sigma^2))$. We can see that the “thin shell” phenomenon is caused by the r^{d-1} term, which makes the density to be highly peaked when d is large. To alleviate the concentration phenomenon, we need to cancel out the effect of r^{d-1} , which motivates the following family of smoothing distributions:

$$\pi_0(\mathbf{z}) \propto \|\mathbf{z}\|_{n_1}^{-k} \exp\left(-\frac{\|\mathbf{z}\|_{n_2}^p}{b}\right),$$

where parameters $k, n_1, n_2, p \in \mathbb{N}$. Next we discuss how to choose suitable parameters depending on specific perturbation region.

4.2 ℓ_1 and ℓ_2 Region Certification

Based on original Laplacian and Gaussian distributions and above intuition, we propose:

$$\ell_1 : \pi_0(\mathbf{z}) \propto \|\mathbf{z}\|_1^{-k} \exp\left(-\frac{\|\mathbf{z}\|_1}{b}\right) \quad (10)$$

$$\ell_2 : \pi_0(\mathbf{z}) \propto \|\mathbf{z}\|_2^{-k} \exp\left(-\frac{\|\mathbf{z}\|_2^2}{2\sigma^2}\right) \quad (11)$$

where we introduce the $\|\mathbf{z}\|^{-k}$ term in π_0 , with k a positive parameter, to make the radius distribution more concentrated when k is large.

The radius distribution in Eq.10 and Eq.11 is controlled by two parameters: σ (or b) and k , who control the scale and shape of the distribution, respectively. The key idea is that adjusting extra parameter k allows us to control the trade-off the accuracy and robustness more precisely. As shown in Fig.1, adjusting σ moves the mean close to zero (hence ① yielding higher accuracy), but at cost of decreasing the variance quadratically (hence ② less robust). In contrast, adjusting k decreases the mean without significantly impacting the variance, thus yield a much better trade-off on accuracy and robustness.

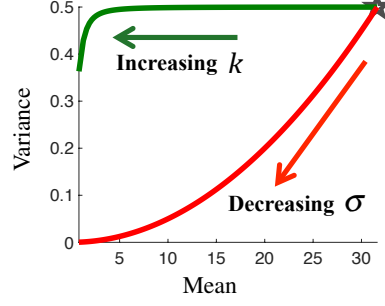


Figure 1: Starting from radius distribution in Eq.11 with $d = 100$ $\sigma = 1$ and $k = 0$ (black start), increasing k (green curve) moves the mean towards zero *without significantly reducing the variance*. Decreasing σ (red curve) can also decrease the mean, but with a cost of decreasing the variance quadratically.

Computational Method Now we no longer have the closed-form solution of the bound like Eq.7 and Eq.8. However, efficient computational methods can still be developed for calculating the bound in Eq.6 with π_0 in Eq.10 or Eq.11. The key is that the maximum of the distance term $\mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda\pi_0 \parallel \pi_\delta)$ over $\delta \in \mathcal{B}$ is always achieved on the boundary of \mathcal{B} :

Theorem 2. Consider the ℓ_1 attack with $\mathcal{B} = \{\delta : \|\delta\|_1 \leq r\}$ and smoothing distribution $\pi_0(\mathbf{z}) \propto \|\mathbf{z}\|_1^{-k} \exp\left(-\frac{\|\mathbf{z}\|_1}{b}\right)$ with $k \geq 0$ and $b > 0$, or the ℓ_2 attack with $\mathcal{B} = \{\delta : \|\delta\|_2 \leq r\}$ and smoothing distribution $\pi_0(\mathbf{z}) \propto \|\mathbf{z}\|_2^{-k} \exp\left(-\frac{\|\mathbf{z}\|_2^2}{2\sigma^2}\right)$ with $k \geq 0$ and $\sigma > 0$. Define $\delta^* = [r, 0, \dots, 0]^\top$, we have

$$\mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda\pi_0 \parallel \pi_{\delta^*}) = \max_{\delta \in \mathcal{B}} \mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda\pi_0 \parallel \pi_\delta)$$

for any positive λ .

With Theorem 2, we can compute Eq.6 with $\delta = \delta^*$. We then calculate $\mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda\pi_0 \parallel \pi_{\delta^*}) = \mathbb{E}_{\mathbf{z} \sim \pi_0} \left[\left(\lambda - \frac{\pi_{\delta^*}(\mathbf{z})}{\pi_0(\mathbf{z})} \right)_+ \right]$ using Monte Carlo approximation with i.i.d. samples $\{\mathbf{z}_i\}_{i=1}^n$ be i.i.d. samples from π_0 : $\hat{D} := \frac{1}{n} \sum_{i=1}^n (\lambda - \pi_{\delta^*}(\mathbf{z}_i)/\pi_0(\mathbf{z}_i))_+$, which is bounded in the following confidence interval $[\hat{D} - \lambda\sqrt{\log(2/\delta)/(2n)}, \hat{D} + \lambda\sqrt{\log(2/\delta)/(2n)}]$ with confidence level $1 - \delta$ for $\delta \in (0, 1)$. What's more, the optimization on $\lambda \geq 0$ is one-dimensional and can be solved numerically efficiently (see Appendix for details).

4.3 ℓ_∞ Region Certification

Going further, we consider the more difficult ℓ_∞ attack whose attacking region is $\mathcal{B}_{\ell_\infty, r} = \{\delta : \|\delta\|_\infty \leq r\}$. The commonly used Gaussian smoothing distribution, as well as our ℓ_2 -based smoothing distribution in Eq.11, is unsuitable for this region:

Proposition 4. With the smoothing distribution π_0 in Eq.11 for $k \geq 0, \sigma > 0$, and $\mathcal{F} = \mathcal{F}_{[0,1]}$ shown in Eq.4, the bound we get for certifying the ℓ_∞ attack on $\mathcal{B}_{\ell_\infty, r} = \{\delta : \|\delta\|_\infty \leq r\}$ is equivalent to that for certifying the ℓ_2 attack on $\mathcal{B}_{\ell_2, \sqrt{d}r} = \{\delta : \|\delta\|_2 \leq \sqrt{d}r\}$, that is,

$$\mathcal{L}_{\pi_0}(\mathcal{F}_{[0,1]}, \mathcal{B}_{\ell_\infty, r}) = \mathcal{L}_{\pi_0}(\mathcal{F}_{[0,1]}, \mathcal{B}_{\ell_2, \sqrt{d}r}).$$

As shown in this proposition, if we use ℓ_2 distribution in Eq.11 for certification, the bound we obtain is effectively the bound we get for verifying a ℓ_2 ball with radius \sqrt{dr} , which is too large to give meaningful results due to high dimension.

In order to address this problem, we extend our proposed distribution family with new distributions which are more suitable for ℓ_∞ certification setting:

$$\pi_0(z) \propto \|z\|_\infty^{-k} \exp\left(-\frac{\|z\|_\infty^2}{2\sigma^2}\right), \quad (12)$$

$$\pi_0(z) \propto \|z\|_\infty^{-k} \exp\left(-\frac{\|z\|_2^2}{2\sigma^2}\right). \quad (13)$$

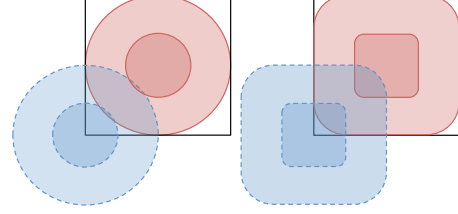


Figure 2: For ℓ_∞ attacking, compared with the distribution in Eq.11, the mixed norm distribution in Eq.13 (right) yields smaller discrepancy term (because of larger overlap areas), and hence higher robustness and better confidence bound. The distribution described in Eq.12 has the same impact.

The motivation is to allocate more probability mass along the “pointy” directions with larger ℓ_∞ norm, and hence decrease the maximum discrepancy term $\max_{\delta \in \mathcal{B}_{\ell_\infty, r}} \mathbb{D}_{\mathcal{F}}(\lambda\pi_0 \parallel \pi_\delta)$, see Fig.2.

Computational Method In order to compute the lower bound with proposed distribution, we need to establish similar theoretical results as Theorem 2, showing the optimal δ is achieved at one vertex (the pointy points) of ℓ_∞ ball.

Theorem 3. Consider the ℓ_∞ attack with $\mathcal{B}_{\ell_\infty, r} = \{\delta : \|\delta\|_\infty \leq r\}$ and the mixed norm smoothing distribution in Eq.13 with $k \geq 0$ and $\sigma > 0$. Define $\delta^* = [r, r, \dots, r]^\top$. We have for any $\lambda > 0$,

$$\mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda\pi_0 \parallel \pi_{\delta^*}) = \max_{\delta \in \mathcal{B}} \mathbb{D}_{\mathcal{F}_{[0,1]}}(\lambda\pi_0 \parallel \pi_\delta).$$

The proofs of Theorem 2 and 3 are non-trivial and deferred to Appendix. With the optimal δ^* found above, we can calculate the bound with similar Monte Carlo approximation outlined in Section 4.2.

5 Experiments

We evaluate proposed certification bound and smoothing distributions for ℓ_1 , ℓ_2 and ℓ_∞ attacks. We compare with the randomized smoothing method of [14] with Laplacian smoothing for ℓ_1 region certification. For ℓ_2 and ℓ_∞ cases, we regard the method derived by [9] with Gaussian smoothing distribution as the baseline. For fair comparisons, we use the same model architecture and pre-trained models provided by [14], [9] and [10], which are ResNet-110 for CIFAR-10 and ResNet-50 for ImageNet. We use the official code² provided by [9] for all the following experiments. For all other details and parameter settings, we refer the readers to Appendix B.2.

ℓ_1 RADIUS (CIFAR-10)	0.25	0.5	0.75	1.0	1.25	1.5	1.75	2.0	2.25
BASELINE (%)	62	49	38	30	23	19	17	14	12
OURS (%)	64	51	41	34	27	22	18	17	14

ℓ_1 RADIUS (IMAGENET)	0.5	1.0	1.5	2.0	2.5	3.0	3.5
BASELINE (%)	50	41	33	29	25	18	15
OURS (%)	51	42	36	30	26	22	16

Table 1: Certified top-1 accuracy of the best classifiers with various ℓ_1 radius.

²<https://github.com/locuslab/smoothing>. Our results are slightly different with those in original paper due to the randomness of sampling.

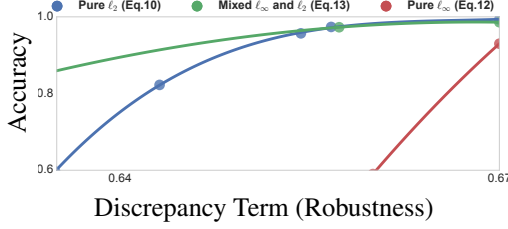


Figure 3: The Pareto frontier of accuracy and robustness (in the sense of Eq.9) of the three smoothing families in Eq.11, Eq.13, and Eq.12 for ℓ_∞ attacking, when we search for the best parameters (k, σ) for each of them. The mixed norm family Eq.13 yields the best trade-off than the other two. We assume $f^\sharp(\mathbf{x}) = \mathbb{I}(\|\mathbf{x}\|_2 \leq r)$ and dimension $d = 5$. The case when $f^\sharp(\mathbf{x}) = \mathbb{I}(\|\mathbf{x}\|_\infty \leq r)$ has similar result (not shown).

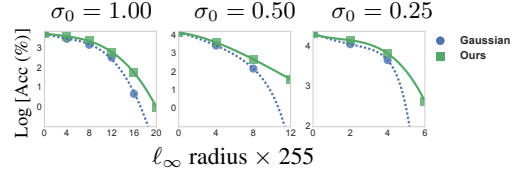


Figure 4: Results of ℓ_∞ verification on CIFAR-10, on models trained with Gaussian noise data augmentation with different variances σ_0 . Our method obtains consistently better results.

Evaluation Metrics Methods are evaluated with the certified accuracy defined in [9]. Given an input image \mathbf{x} and a perturbation region \mathcal{B} , the smoothed classifier certifies image \mathbf{x} correctly if the prediction is correct and has a guaranteed confidence lower bound larger than $1/2$ for any $\delta \in \mathcal{B}$. The certified accuracy is the percentage of images that are certified correctly. Following [10], we calculate the certified accuracy of all the classifiers in [9] or [10] for various radius, and report the best results over all of classifiers.

5.1 ℓ_1 & ℓ_2 Certification

For ℓ_1 certification, we compare our method with [14] on CIFAR-10 and ImageNet with the type 1 trained model in [14]. As shown in Table 1, our non-Laplacian centripetal distribution consistently outperforms the result of baseline for any ℓ_1 radius.

ℓ_2 RADIUS (CIFAR-10)	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

ℓ_2 RADIUS (IMAGENET)	0.5	1.0	1.5	2.0	2.5	3.0	3.5
BASELINE (%)	49	37	29	19	15	12	9
OURS (%)	50	39	31	21	17	13	10

Table 2: Certified top-1 accuracy of the best classifiers with various ℓ_2 radius.

For ℓ_2 certification, we compare our method with [9] on CIFAR-10 and ImageNet. For a fair comparison, we use the same pre-trained models as [9], which is trained with Gaussian noise on both CIFAR-10 and ImageNet dataset. Table 2 reports the certified accuracy of our method and the baseline on CIFAR-10 and ImageNet. We find that our method consistently outperforms the baseline. The readers are referred to the Appendix B.3 for detailed ablation studies.

5.2 ℓ_∞ Certification

Toy Example We first construct a simple toy example to verify the advantages of the distribution Eq.13 and Eq.12 over the ℓ_2 family in Eq.11. We set the true classifier to be $f^\sharp(\mathbf{x}) = \mathbb{I}(\|\mathbf{x}\|_2 \leq r)$ in $r = 0.65$, $d = 5$ case and plot in Fig.3 the Pareto frontier of the accuracy and robustness terms in Eq.9 for the three families of smoothing distributions, as we search for the best combinations of parameters (k, σ) . The mixed norm smoothing distribution clearly obtain the best trade-off on accuracy and robustness, and hence guarantees a tighter lower bound for certification. Fig.3 also shows that Eq.12 even performs worse than Eq.11. We further theoretically show that Eq.12 is provably not suitable for ℓ_∞ region certification in Appendix A.5.

CIFAR-10 Based on above results, we only compared the method defined by Eq.13 with [10] on CIFAR-10. The certified accuracy of our method and the baseline using Gaussian smoothing distribution and Proposition 4 are shown in Table 3. We can see that our method consistently outperforms the Gaussian baseline by a large margin. More clarification about ℓ_∞ experiments is in Appendix ??.

l_∞ 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 3: Certified top-1 accuracy of the best classifiers with various l_∞ radius on CIFAR-10.

To further confirm the advantage of our method, we plot in Fig.4 the certified accuracy of our method and Gaussian baseline using models trained with Gaussian perturbation of different variances σ_0 under different l_∞ radius. Our approach outperforms baseline consistently, especially when the l_∞ radius is large. We also experimented our method and baseline on ImageNet but did not obtain non-trivial results. This is because l_∞ verification is extremely hard with very large dimensions [32, 31]. Future work will investigate how to obtain non-trivial bounds for l_∞ attacking at ImageNet scales with smoothing classifiers.

6 Conclusion

We propose a general functional optimization based framework of adversarial certification with non-Gaussian smoothing distributions. Based on the insights from our new framework and high dimensional geometry, we propose a new family of non-Gaussian smoothing distributions, which outperform the Gaussian and Laplace smoothing for certifying ℓ_1 , ℓ_2 and ℓ_∞ attacking. Our work provides a basis for a variety of future directions, including improved methods for ℓ_p attacks, and tighter bounds based on adding additional constraints to our optimization framework.

Broader Impact

Adversarial certification via randomized smoothing could achieve *guaranteed* robust machine learning models, thus has wide application on AI security. a & b) With our empirical results, security engineers could get better performance on defending against vicious attacks; With our theoretical results, it will be easier for following researchers to derive new bounds for different kinds of smoothing methods. We don’t foresee the possibility that it could bring negative social impacts. c) Our framework is mathematically rigorous thus would never fail. d) Our method doesn’t have bias in data as we provide a general certification method for all tasks and data, and our distribution is not adaptive towards data.

Acknowledgement

Gong, Ye, Liu are supported in part by NSF CAREER 1846421. Zhu is supported in part by Beijing Nova Program (No. 202072) from Beijing Municipal Science & Technology Commission. We would like to thank Tongzheng Ren, Jiaye Teng, Yang Yuan and the reviewers for helpful suggestions that improved the paper.

References

- [1] Harini Kannan, Alexey Kurakin, and Ian Goodfellow. Adversarial logit pairing. *arXiv preprint arXiv:1803.06373*, 2018.
- [2] Haichao Zhang and Jianyu Wang. Defense against adversarial attacks using feature scattering-based adversarial training, 2019.
- [3] Runtian Zhai, Tianle Cai, Di He, Chen Dan, Kun He, John Hopcroft, and Liwei Wang. Adversarially robust generalization just requires more unlabeled data. *arXiv preprint arXiv:1906.00555*, 2019.

- [4] Anish Athalye, Nicholas Carlini, and David A. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. pages 274–283, 2018.
- [5] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. 2018.
- [6] Dinghui Zhang, Tianyuan Zhang, Yiping Lu, Zhanxing Zhu, and Bin Dong. You only propagate once: Accelerating adversarial training via maximal principle. *Advances in neural information processing systems (NeurIPS)*, 2019.
- [7] Dilin Wang, Chengyue Gong, and Qiang Liu. Improving neural language modeling via adversarial training. pages 6555–6565, 2019.
- [8] Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified robustness to adversarial examples with differential privacy. *arXiv preprint arXiv:1802.03471*, 2018.
- [9] Jeremy M Cohen, Elan Rosenfeld, and J Zico Kolter. Certified adversarial robustness via randomized smoothing. *arXiv preprint arXiv:1902.02918*, 2019.
- [10] Hadi Salman, Greg Yang, Jerry Li, Pengchuan Zhang, Huan Zhang, Ilya Razenshteyn, and Sebastien Bubeck. Provably robust deep learning via adversarially trained smoothed classifiers. *arXiv preprint arXiv:1906.04584*, 2019.
- [11] Guang-He Lee, Yang Yuan, Shiyu Chang, and Tommi S Jaakkola. A stratified approach to robustness for randomly smoothed classifiers. *Advances in neural information processing systems (NeurIPS)*, 2019.
- [12] Bai Li, Changyou Chen, Wenlin Wang, and Lawrence Carin. Second-order adversarial attack and certifiable robustness. *Advances in neural information processing systems (NeurIPS)*, 2019.
- [13] Krishnamurthy (Dj) Dvijotham, Jamie Hayes, Borja Balle, Zico Kolter, Chongli Qin, Andras Gyorgy, Kai Xiao, Sven Gowal, and Pushmeet Kohli. A framework for robustness certification of smoothed classifiers using f-divergences. In *International Conference on Learning Representations*, 2020.
- [14] Jiaye Teng, Guang-He Lee, and Yang Yuan. ℓ_1 adversarial robustness certificates: a randomized smoothing approach, 2020.
- [15] Jinyuan Jia, Xiaoyu Cao, Binghui Wang, and Neil Zhenqiang Gong. Certified robustness for top-k predictions against adversarial perturbations via randomized smoothing. In *International Conference on Learning Representations*, 2020.
- [16] Eric Wong and J Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. *arXiv preprint arXiv:1711.00851*, 2017.
- [17] Krishnamurthy Dvijotham, Robert Stanforth, Sven Gowal, Timothy A. Mann, and Pushmeet Kohli. A dual approach to scalable verification of deep networks. In *Proceedings of the Thirty-Fourth Conference on Uncertainty in Artificial Intelligence, UAI 2018, Monterey, California, USA, August 6-10, 2018*, pages 550–559, 2018.
- [18] Matt Jordan, Justin Lewis, and Alexandros G Dimakis. Provable certificates for adversarial examples: Fitting a ball in the union of polytopes. *arXiv preprint arXiv:1903.08778*, 2019.
- [19] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [20] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press, 2018.
- [21] Nicholas Carlini, Guy Katz, Clark Barrett, and David L. Dill. Provably minimally-distorted adversarial examples, 2017.

- [22] Ruediger Ehlers. Formal verification of piece-wise linear feed-forward neural networks. In *International Symposium on Automated Technology for Verification and Analysis*, pages 269–286. Springer, 2017.
- [23] Chih-Hong Cheng, Georg Nührenberg, and Harald Ruess. Maximum resilience of artificial neural networks. In *Automated Technology for Verification and Analysis - 15th International Symposium, ATVA 2017, Pune, India, October 3-6, 2017, Proceedings*, pages 251–268, 2017.
- [24] Souradeep Dutta, Susmit Jha, Sriram Sankaranarayanan, and Ashish Tiwari. Output range analysis for deep feedforward neural networks, 2018.
- [25] Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Semidefinite relaxations for certifying robustness to adversarial examples, 2018.
- [26] Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. Efficient neural network robustness certification with general activation functions. In *Advances in neural information processing systems*, pages 4939–4948, 2018.
- [27] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. Mitigating adversarial effects through randomization, 2018.
- [28] Xuanqing Liu, Minhao Cheng, Huan Zhang, and Cho-Jui Hsieh. Towards robust neural networks via random self-ensemble. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 369–385, 2018.
- [29] Wenbo V Li and James Kuelbs. Some shift inequalities for gaussian measures. In *High dimensional probability*, pages 233–243. Springer, 1998.
- [30] Greg Yang, Tony Duan, Edward Hu, Hadi Salman, Ilya Razenshteyn, and Jerry Li. Randomized smoothing of all shapes and sizes. *arXiv preprint arXiv:2002.08118*, 2020.
- [31] Avrim Blum, Travis Dick, Naren Manoj, and Hongyang Zhang. Random smoothing might be unable to certify linf robustness for high-dimensional images. *arXiv preprint arXiv:2002.03517*, 2020.
- [32] Aounon Kumar, Alexander Levine, Tom Goldstein, and Soheil Feizi. Curse of dimensionality on randomized smoothing for certifiable robustness. *arXiv preprint arXiv:2002.03239*, 2020.
- [33] Stephen Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge University Press, New York, NY, USA, 2004.