UniCR: Universally Approximated Certified Robustness via Randomized Smoothing

Hanbin Hong¹, Binghui Wang², and Yuan Hong¹

 ¹ University of Connecticut, Storrs CT 06269, USA hanbin.hong@uconn.edu, yuan.hong@uconn.edu
 ² Illinois Institute of Technology, Chicago IL 60616, USA bwang70@iit.edu

A Preliminary

We first briefly review the recent certified robustness scheme [4] for a general classification problem by classifying data point in \mathbb{R}^d to classes in \mathcal{Y} . Given an arbitrary base classifier f, it can be converted to a "smoothed" classifier [4] g by adding isotropic Gaussian noise to the input x:

$$g(x) = \arg\max_{c \in \mathcal{Y}} \mathbb{P}(f(x+\epsilon) = c), where \ \epsilon \sim \mathcal{N}(0, \sigma^2 I)$$
(1)

Lemma 1. (Neyman-Pearson Lemma) Let X and Y be random variables in \mathbb{R}^d with densities μ_X and μ_Y . Let $f : \mathbb{R}^d \to \{0,1\}$ be a random or deterministic function. Then:

(1) If $S = \{z \in \mathbb{R}^d : \frac{\mu_Y(z)}{\mu_X(z)} \leq t\}$ for some t > 0 and $\mathbb{P}(f(X) = 1) \geq \mathbb{P}(X \in S)$, then $\mathbb{P}(f(Y) = 1) \geq \mathbb{P}(Y \in S)$;

(2) If $S = \{z \in \mathbb{R}^d : \frac{\mu_Y(z)}{\mu_X(z)} \ge t\}$ for some t > 0 and $\mathbb{P}(f(X) = 1) \le \mathbb{P}(X \in S)$, then $\mathbb{P}(f(Y) = 1) \le \mathbb{P}(Y \in S)$.

With Lemma 1, Cohen [4] derives the certified radius when the classifier is smoothed with the Gaussian noise. As shown in Theorem 1, when the smoothed classifier's prediction probabilities satisfy Equation (2), the prediction result is guaranteed to be the most probable class c_A when the perturbation is limited within a radius R in ℓ_2 -norm.

Theorem 1. (Randomized Smoothing with Gaussian Noise [4]) Let f: $\mathbb{R}^d \to \mathcal{Y}$ be any deterministic or random function, and let $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$. Denote g as the smoothed classifier in Equation (1), and the most probable and the second probable classes as $c_A, c_B \in \mathcal{Y}$, respectively. If the lower bound of the class c_A 's prediction probability $\underline{p}_A \in [0, 1]$, and the upper bound of the class c_B 's prediction probability $\overline{p}_B \in [0, 1]$ satisfy:

$$\mathbb{P}(f(x+\epsilon) = c_A) \ge \underline{p_A} \ge \overline{p_B} \ge \max_{c \neq c_A} \mathbb{P}(f(x+\epsilon) = c)$$
(2)

Then $g(x + \delta) = c_A$ for all $||\delta||_2 \leq R$, where

$$R = \frac{\sigma}{2} \left(\Phi^{-1}(\underline{p_A}) - \Phi^{-1}(\overline{p_B}) \right) \tag{3}$$

where Φ^{-1} is the inverse of the standard Gaussian CDF.

Proof. See detailed proof in [4].

B Proofs

B.1 Proof of Theorem 1

Proof. We prove the theorem based on Neyman-Pearson Lemma (Lemma 1).

Let $x := x_0 + \epsilon$ be the random variable that follows any continuous distribution. δ be the perturbation added to the input image. $y = x_0 + \epsilon + \delta$ is the perturbed random variable. Thus, x and y are random variables with densities μ_x and μ_y . Define sets:

$$A := \{ z : \frac{\mu_y(z)}{\mu_x(z)} \le t_A \}$$
(4)

$$B := \{ z : \frac{\mu_y(z)}{\mu_x(z)} \ge t_B \}$$
 (5)

where t_A and t_B are picked to suffice:

$$\mathbb{P}(x \in A) = \underline{p_A} \tag{6}$$

$$\mathbb{P}(x \in B) = \overline{p_B} \tag{7}$$

Suppose $c_A \in \mathcal{Y}$ and $p_A, \overline{p_B} \in [0, 1]$ satisfy:

$$\mathbb{P}(f(x+\epsilon) = c_A) \ge \underline{p_A} \ge \overline{p_B} \ge \max_{c \neq c_A} \mathbb{P}(f(x+\epsilon) = c)$$
(8)

Since $\mathbb{P}(f(x + \epsilon) = c_A) \ge \underline{p_A} = \mathbb{P}(x \in A)$ and $A = \{z : \frac{\mu_Y(z)}{\mu_X(z)} \le t_A\}$, using Neyman-Pearson Lemma (Lemma 1), we have:

$$\mathbb{P}(f(y) = c_A) \ge \mathbb{P}(y \in A) \tag{9}$$

Similarly, we have:

$$\mathbb{P}(f(y) = c_B) \le \mathbb{P}(y \in B) \tag{10}$$

To guarantee $\mathbb{P}(f(y) = c_A) \ge \mathbb{P}(f(y) = c_B)$, we need

$$\mathbb{P}(f(y) = c_A) \ge \mathbb{P}(y \in A) \ge \mathbb{P}(y \in B) \ge \mathbb{P}(f(y) = c_B)$$
(11)

In summary, to guarantee the certified robustness on class A, Equation (4), (5), (6), (7), (11) must be satisfied. The conditions can be rewritten as:

$$\mathbb{P}(\frac{\mu_y(x)}{\mu_x(x)} \le t_A) = \underline{p_A} \tag{12}$$

$$\mathbb{P}(\frac{\mu_y(x)}{\mu_x(x)} \ge t_B) = \overline{p_B}$$
(13)

UniCR 3

$$\mathbb{P}(\frac{\mu_y(y)}{\mu_x(y)} \le t_A) \ge \mathbb{P}(\frac{\mu_y(y)}{\mu_x(y)} \ge t_B)$$
(14)

where Equation (12) is from Equation (4) and Equation (6), Equation (13) is from Equation (5) and Equation (7), and Equation (14) is from Equation (11). Considering the relationship $y = x + \delta$, we can derive:

$$\mu_y(x) = \mu_x(x - \delta) \tag{15}$$

$$\mu_x(y) = \mu_x(x+\delta) \tag{16}$$

$$\mu_y(y) = \mu_x(y - \delta) = \mu_x(x) \tag{17}$$

Thus, the conditions (11), (12) and (13) can be rewritten as:

$$\mathbb{P}(\frac{\mu_x(x-\delta)}{\mu_x(x)} \le t_A) = \underline{p_A} \tag{18}$$

$$\mathbb{P}(\frac{\mu_x(x-\delta)}{\mu_x(x)} \ge t_B) = \overline{p_B}$$
(19)

$$\mathbb{P}(\frac{\mu_x(x)}{\mu_x(x+\delta)} \le t_A) \ge \mathbb{P}(\frac{\mu_x(x)}{\mu_x(x+\delta)} \ge t_B)$$
(20)

Any perturbation δ satisfying these conditions will not fool the smoothed classifier. In this case, these conditions construct a robustness area in δ space. If we want to find a ℓ_p ball within which the prediction is constant, the l_p ball should be in this robustness area. Therefore, the certified radii is the minimum $||\delta||_p$ on the boundary of this robustness area. In this case, the ℓ_p ball is exactly the maximum inscribed ball in the robustness area. Also, x can be replace by ϵ in these conditions since it is in the fraction, which means the optimization is independent to the input if given \underline{p}_A and \overline{p}_B . Therefore, the whole optimization problem is summarized as:

$$\begin{array}{ll} \underset{\delta}{\text{minimize}} & R = ||\delta||_{p} \\ \text{subject to} & \mathbb{P}(\frac{\mu_{x}(x-\delta)}{\mu_{x}(x)} \leq t_{A}) = \underline{p}_{A}, \\ & \mathbb{P}(\frac{\mu_{x}(x-\delta)}{\mu_{x}(x)} \geq t_{B}) = \overline{p}_{B}, \\ & \mathbb{P}(\frac{\mu_{x}(x)}{\mu_{x}(x+\delta)} \leq t_{A}) = \mathbb{P}(\frac{\mu_{x}(x)}{\mu_{x}(x+\delta)} \geq t_{B}) \end{array}$$

If the noise is isotropic, each dimension is independent,

$$\mu_x(x) = \prod_{i=1}^d \mu_x(x_j)$$
 (21)

Thus, conditions for the isotropic noise can be rewritten as:

$$\mathbb{P}(\prod_{j=1}^{d} \frac{\mu_x(x_j - \delta_j)}{\mu_x(x_j)} \le t_A) = \underline{p_A}$$
(22)

$$\mathbb{P}(\prod_{j=1}^{d} \frac{\mu_x(x_j - \delta_j)}{\mu_x(x_j)} \ge t_B) = \overline{p_B}$$
(23)

$$\mathbb{P}\left(\prod_{j=1}^{d} \frac{\mu_x(x_j)}{\mu_x(x_j+\delta_j)} \le t_A\right) = \mathbb{P}\left(\prod_{j=1}^{d} \frac{\mu_x(x_j)}{\mu_x(x_j+\delta_j)} \ge t_B\right)$$
(24)

Thus, this completes the proof.

B.2 Binary Case for Theorem 2

Theorem 2. (Universal Certified Robustness (Binary Case)) Let $f : \mathbb{R}^d \to \mathcal{Y}$ be any deterministic or random function, and let ϵ follows any continuous distribution. Let g be defined as in (1). Suppose the most probable class $c_A \in \mathcal{Y}$ and the lower bound of the probability \underline{p}_A satisfy:

$$\mathbb{P}(f+\epsilon) = c_A \ge \underline{p_A} \ge \frac{1}{2} \tag{25}$$

Then $g(x + \delta) = c_A$ for all $||\delta||_p \leq R$, where R is given by the optimization:

$$\begin{array}{ll} \underset{\delta}{\text{minimize}} & R = ||\delta||_p\\ \text{subject to} & \mathbb{P}(\frac{\mu_x(x-\delta)}{\mu_x(x)} \leq t_A) = \underline{p}_A,\\ & \mathbb{P}(\frac{\mu_x(x)}{\mu_x(x+\delta)} \leq t_A) = \frac{1}{2} \end{array}$$

B.3 UniCR (Binary Case)

Similar to the binary case of Theorem 2, the binary case of the two-phase optimization can be easily derived:

$$\begin{split} R &= ||\lambda \delta||_{p}, where \ \delta \in \mathop{\arg\min}_{\delta} ||\lambda \delta||_{p} \\ s.t. \quad \lambda &= \mathop{\arg\min}_{\lambda} |K| \\ \mathbb{P}(\frac{\mu_{x}(x - \lambda \delta)}{\mu_{x}(x)} \leq t_{A}) = \underline{p}_{A} \\ K &= \mathbb{P}(\frac{\mu_{x}(x)}{\mu_{x}(x + \lambda \delta)} \leq t_{A}) - \frac{1}{2} \\ \underline{p}_{A} \geq \frac{1}{2} \end{split}$$

B.4 UniCR Bound

The certified radius R approximated by the two-phase optimization is tight if achieving the optimality. Under this assumption, we analysis the confidence bound for the certification. We follow [4] to compute the probabilities \underline{p}_A and \overline{p}_B using Monte Carlo method with sample number n. The confidence is $1 - \alpha_0$, where $p_A >= \alpha_0^{1/n}$. To estimate the auxiliary parameters t_A and t_B , we use Dvoretzky–Kiefer–Wolfowitz inequality [6] to bound the CDFs of the random variables $\frac{\mu_x(x-\lambda\delta)}{\mu_x(x)}$ and $\frac{\mu_x(x)}{\mu_x(x+\lambda\delta)}$, then determine the t_A and t_B using Algorithm 1.

Lemma 2. (Dvoretzky–Kiefer–Wolfowitz inequality(restate)) Let $X_1, X_2, ..., X_n$ be real-valued independent and identically distributed random variables with cumulative distribution function $F(\cdot)$, where $n \in \mathbb{N}$. Let F_n denotes the associated empirical distribution function defined by

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{X_i < =x\}}, x \in \mathbb{R}$$
(26)

The Dvoretzky-Kiefer-Wolfowitz inequality bounds the probability that the random function F_n differs from F by more than a given constant $\Delta \in \mathbb{R}^+$:

$$\mathbb{P}(\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| > \Delta) \le 2e^{-2n\Delta^2}$$
(27)

We use the Lemma 2 to estimate the CDFs in algorithm 1. In the robustness condition, we need to estimate 4 probabilities with confidence $1-2e^{-2n\Delta^2}$ as well as the p_A and p_B with confidence $1 - \alpha_0$. Therefore, the confidence that deriving the correct radius is at least $(1 - \alpha_0)^2 (1 - 2e^{-2n\Delta^2})^4$. In Figure 1, we show the confidence on a varying number of samples when $\Delta = 0.1$ and $\alpha_0 = 0.999$. As the number of samples increases to around 400 (all our experiments use more than 400 samples), our confidence is very close to Cohen's confidence [4]. Thus, the confidence is nearly 1 in all our experiments.

B.5 Optimization Convergence

We analysis the convergence of the two-phase optimization and the certification accuracy in this section. On one hand, the optimality of the scalar optimization can be asymptotically achieved by binary search. On the other hand, it is hard to find the minimum $||\lambda\delta||_p$ in the highly-dimensional space, but some special symmetry in the direction of δ (e.g., spherical symmetry that is also found in [18, 17]), can help approximate the certified radius. The detailed algorithms are presented in Section 3.1. The defense performance of such universally approximated certified robustness against different real-world attacks is the same as certified robustness (as shown in Appendix D.2). Thus, such negligible approximation error is close to 0, but result in many significant new benefits in return.



Fig. 1. Confidence vs. number of Monte Carlo samples.

\mathbf{C} Algorithms

C.1 Computing t_A and t_B

We present the algorithm to compute the p_A and p_B in Algorithm 1.

Algorithm 1 Computing t_A and t_B

- **Input:** Lower bound of the probabilities, p_A ; upper bound of the probabilities, $\overline{p_B}$; perturbation scalar, λ ; perturbation, δ ; noise PDF, μ_x ; number of samples in the Monte Carlo method, n
- **Output:** The auxiliary parameters, t_A and t_B
- 1: Sample n noise $\epsilon \in \mathbb{R}^{n \times d}$ from a discrete version of PDF. 2: Calculate $\frac{\mu_x(x-\lambda\delta)}{\mu_x(x)}$ using these n samples of noise, μ_x , λ and δ
- 3: Estimate the CDF Φ of $\frac{\mu_x(x-\lambda\delta)}{\mu_x(x)}$ using Monte Carlo method
- 4: return $t_A = \Phi^{-1}(p_A)$ and $t_B = \Phi^{-1}(p_B)$, with inverse CDF Φ

C.2Scalar Optimization

We use the binary search to find a scale factor that minimizes |K| (the distance between δ and the robustness boundary). When K = 0, the perturbation δ is exactly on the robustness boundary. Fixing the direction of δ , we find two scalars such that K > 0 and K < 0. Specifically, we start from a scalar λ_a and compute K. If K > 0, then the scaled perturbation $\lambda_a \delta$ is within the robustness boundary, thus we enlarge the scalar to find a λ_b such that K < 0 and vice versa. After that, we iteratively compute the K using $\lambda = \frac{1}{2}(\lambda_a + \lambda_b)$: if K > 0, we let $\lambda_a = \lambda$; otherwise, we let $\lambda_b = \lambda$. We repeat this iteration until K is less than a threshold or the number of iterations is sufficiently large. The procedures are summarized in Algorithm 2.

Algorithm 2 Scalar Optimization

Input: Lower bound of the probabilities, p_A ; upper bound of the probabilities, $\overline{p_B}$; perturbation scalar, λ ; perturbation, δ ; noise PDF, μ_x ; number of samples in Monte Carlo method, n; threshold for K, K_m ; number of iterations for binary search, N **Output:** The scalar λ that minimizes |K|1: Find initial scalar λ_a and λ_b such that K > 0 and K < 02: $\lambda = (\lambda_a + \lambda_b)/2$ 3: Compute K using λ 4: while N > 0 and $|K| > K_m$ do if K > 0 then 5: $\lambda_a = \lambda$ 6: 7: else $\lambda_b = \lambda$ 8: $\lambda = (\lambda_a + \lambda_b)/2$ 9: Compute K using λ 10:N=N-111: 12: return λ

C.3 Direction Optimization

We show how to initialize the positions for different ℓ_p norms in PSO. Since some noise follows PDFs with symmetry [18, 17], we set the initial position of particles by considering this, e.g., setting the initial positions w.r.t. ℓ_p for $p \in \mathbb{R}^+$ as [0, ..., 0, a, 0, ..., 0] and the initial positions w.r.t. ℓ_{∞} as [a, a, a, ..., a], where a is a small random number. Although the search space is highly-dimensional, empirical results show that the radius given by PSO can accurately approximate the theoretical radius given by other methods, e.g., Cohen's [4] (see Figure 4 in main context). Notice that, for more complicated PDFs without symmetry (which is indeed difficult for deriving the certified radius), PSO can also approximate the certified radius with more particles and iterations.

C.4 Noise Optimization Algorithms

With the UniCR for estimating the certified radius, we can further tune the noise PDF to each input or the classifier. Specifically, let $\mu(x, \alpha)$ denote the noise PDF, where α is a set of hyper-parameters in the function, i.e., $\alpha = [\alpha_1, \alpha_2, ..., \alpha_m]$. We simply use grid-search algorithm to find the best hyper-parameters in the classifier smoothing (C-OPT). For the input noise optimization (I-OPT), we use the Hill-Climbing algorithm to find the optimal hyper-parameters in the function for each input. During the algorithm execution, hyper-parameter for the input is iteratively updated if a better solution is found in each round, until convergence. The procedures for the Hill Climbing algorithm are summarized in Algorithm 3.

Algorithm 3 I-OPT with Hill Climbing

Input: Input data, x; PDF of noise distribution, μ_x ; universally approximated certi-						
fied robustness, $\text{UniCR}(\cdot)$; initial hyper-parameters $\boldsymbol{\alpha}$; optimization range of hyper-						
parameters, $[\boldsymbol{L}, \boldsymbol{H}]$; optimization step of hyper-parameters, \boldsymbol{S}						
Dutput: The optimal hyper-parameters, $\alpha_{optimal}$						
1: Initialize the certified radius $R_0 = \text{UniCR}(x, \mu(\boldsymbol{\alpha}))$						
2: For each hyper-parameter α_i in $\boldsymbol{\alpha}$:						
3: if $L_i < \alpha_i + S_i < H_i$ then						
4: $R' = \text{UniCR}(x, \mu(\boldsymbol{\alpha} \alpha_i = \alpha_i + S_i))$						
5: if $R' > R_0$ then						
6: $\boldsymbol{\alpha}$ is updated with $\alpha_i = \alpha_i + S_i$						
7: $R_0 = R'$						
8: else if $L_i < \alpha_i - S_i < H_i$ then						
9: $R' = \text{UniCR}(x, \mu(\boldsymbol{\alpha} \alpha_i = \alpha_i - S_i))$						
10: $R_0 = R'$						
11: if $R' > R$ then						
12: $\boldsymbol{\alpha}$ is updated with $\alpha_i = \alpha_i - S_i$						
$13: R_0 = R'$						
14: else						
15: break						
16: else						
17: break						
18: return $\alpha_{optimal} = \alpha$						

D More Experimental Results

D.1 Metrics

We show the illustration of Robustness Score in Figure. 2

D.2 Defense against Real Attacks

We evaluate our UniCR's defense accuracy against a diverse set of state-ofthe-art attacks, including universal attacks [3], white-box attacks [5, 15], and black-box attacks [1, 2]. We compare UniCR with other state-of-the-art certified schemes [17, 4, 13] against ℓ_1, ℓ_2 and ℓ_{∞} perturbations. The certified radius R for each image in the test set (10,000 images in total) are computed beforehand, and the perturbation generation is constrained by $||\delta||_p = R$ for all the attack methods. We define the defense accuracy as the rate that the smoothed classifier can successfully defend against the perturbations with the ℓ_p size identical to the the certified radius:

$$acc_d = \mathbb{E}_{||\delta||_p = R}\left[\frac{\sum g(x+\delta) = c_A}{N}\right]$$
(28)

where $c_A = g(x)$, N is the total test number. In this defense study, we use 500 samples for both Monte Carlo method and testing.

Table 1 shows the defense accuracy on the smoothed classifier. The attack with "*" is re-scaled to the required norm (perturbation size R) based on their perturbation formats. UniCR universally provides a 100% defense accuracy against all the ℓ_1 , ℓ_2 and ℓ_{∞} perturbations generated by all the state-of-the-art attacks. These results validate our universally approximated certified robustness ensures the same defense performance as certified robustness in practice.



Fig. 2. An example of the Robustness Score.

Table 1. Defense against real attacks on CIFAR10 (results on MNIST & ImageNet are similar and not included due to space limit).

Defense Accuracy (%)	Gaussian*	Procedural*	[3] Auto-PGD [5]	Wasserstein [*] [15] Square* [1]	HSJ^* [2]
Teng's [13] ℓ_1 -norm R	100.00	100.00	100.00	100.00	100.00	100.00
Our ℓ_1 -norm R	100.00	100.00	100.00	100.00	100.00	100.00
Cohen's [4] ℓ_2 -norm R	100.00	100.00	100.00	100.00	100.00	100.00
Our ℓ_2 -norm R	100.00	100.00	100.00	100.00	100.00	100.00
Yang's [17] ℓ_{∞} -norm R	100.00	100.00	100.00	100.00	100.00	100.00
Our ℓ_{∞} -norm R	100.00	100.00	100.00	100.00	100.00	100.00

D.3 List of PDFs

The PDFs used in our experimental are summarized in Table 2.

Table 2. List of noise distributions.

Distribution	Probability Density Function
Gaussian	$\propto e^{- x/lpha ^2}$
Laplace	$\propto e^{- x/\alpha }$
Hyperbolic Secant	$\propto sech(x/lpha)$
General Normal	$\propto e^{- x/\alpha ^{\beta}}$
Cauthy	$\propto \frac{\alpha^2}{x^2 + \alpha^2}$
Pareto	$\propto \frac{1}{(1+ x/\alpha)^{\beta+1}}$
Laplace-Gaussian Mix.	$\propto \beta e^{- x/\alpha ^1} + (1-\beta) e^{- x/\alpha ^2}$
Exponential Mix.	$\propto e^{-\beta x/\alpha ^1 - (1-\beta) x/\alpha ^2}$

D.4 Efficiency for Radius Derivation

We show the runime of our algorithms on deriving the certified radius for the inputs with various input dimensions in Figure 3. For the common input dimen-

sions, e.g., 24×24 for MNIST, $3 \times 32 \times 32$ for CIFAR10, and $3 \times 224 \times 224$ for ImageNet, it takes less than 10 seconds for certifying an image on average. Comparing with the theoretical certified radius deriving, our method's running time is undoubtedly larger since their radius is pre-derived. However, with the significant benefits on the universality and the automatically deriving, we believe the cost of the extra running time is worthwhile and acceptable in practice.



Fig. 3. Runtime of UniCR vs. input sizes (with RTX3080 GPU).

D.5 Any p (besides $1, 2, \infty$)



Fig. 4. Radius vs. various ℓ_p pert.

Existing methods [4,13] usually focus on the certified radius in a specific norm, e.g., ℓ_1 , ℓ_2 or ℓ_{∞} norms. Some methods [18,17] provide certified robustness theories for multiple norms but specific settings are usually needed for deriving the certified radii in different norms. None of the existing methods can automatically compute the certified radius in any ℓ_p norm. In this section, we show our UniCR can automatically approximate the certified radii for various p, in which p is a real number greater than 0.

In the experiments, we set the probability $\underline{p_A} = 0.9$ and draw the lines of certified radius w.r.t. different p for p > 0. We show the results computed with different noise distributions in Figure 4. We observe that when $p \in (0, 2]$, the certified radius for different p are approximately identical. This finding also matches the theoretical results in Yang et al. [17], in which the certified radii in ℓ_1 and ℓ_2 norm are exactly the same for multiple distributions. When p > 2, we observe that the certified radius decreases as p increases.

D.6 Evaluations on Complicated PDFs

We provide a fine-grained evaluation on the complicated distributions [12], e.g., General Normal, Laplace-Gaussian Mixture, and Exponential Mixture noises with various β . It shows that the Gaussian (i.e., $\beta = 2$ for General Normal, $\beta = 0$ for Laplace-Gaussian Mixture and Exponential Mixture) is the optimal noise in these β setting. We also observe the "crash" on Laplace-based distributions when p_A is small.

D.7 Certification on Non-Smoothed Classifier

radius R		0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8
Yang's [17] vs. ℓ_1 -norm	10.6	10.4	10.4	9.8	8.8	8.2	5.4	2.2	1.0
Ours vs. ℓ_1 -norm	98.8	47.0	22.4	17.8	13.8	10.2	7.0	3.8	1.0
Cohen's [4] vs. ℓ_2 -norm	10.6	10.4	10.4	9.6	8.8	8.2	5.6	2.2	1.2
Ours vs. ℓ_2 -norm	98.8	46.0	22.4	17.6	13.8	9.8	7.0	3.8	1.2
Yang's [17] vs. ℓ_{∞} -norm (at $R/255$)	10.6	10.6	10.6	10.4	10.4	10.4	10.4	10.4	10.4
Ours vs. ℓ_{∞} -norm (at $R/255$)	98.6	92.4	69.4	61.6	53.6	46.0	37.8	27.4	24.4

Table 3. Certified accuracy on standard classifier.

Besides certifying inputs with the smoothed classifier, our input noise optimization (I-OPT) can certify input with a standard classifier without degrading the classifier accuracy on clean data (on the contrary, existing works have to trade off such accuracy for certified defenses).

Specifically, since our I-OPT allows the noise for the input certification to be different from the noise used in training, a special case of the training noise is no noise ($\sigma = 0$). This means that we can certify a naturally-trained classifier



(a) General Normal vs. ℓ_1 (b) General Normal vs. ℓ_2 (c) General Normal vs. ℓ_{∞} perturbations perturbations



(d) Laplace-Gaussian Mix- (e) Laplace-Gaussian Mix- (f) Laplace-Gaussian Mixture vs. ℓ_{∞} perturbations ture vs. ℓ_{∞} perturbations ture vs. ℓ_{∞} perturbations



(g) Exponential Mixture (h) Exponential Mixture (i) Exponential Mixture vs. vs. ℓ_{∞} perturbations vs. ℓ_{∞} perturbations ℓ_{∞} perturbations

Fig. 5. p_A -R curves of General Normal, Laplace-Gaussian Mixture, and Exponential Mixture noise with a varying β .

(standard classifier). This provides an obvious benefit that the classifier can still execute normal classification on clean data with high accuracy since the standard classifier is trained without noise. Also, with I-OPT, we can tune the noise for the input to maintain the prediction accuracy. Thus, any classifier can be certifiably protected against perturbations without degrading the general performance on clean data.

To maintain the performance on standard classification, we add a condition while performing I-OPT:

$$g(x+\delta) = f(x) \tag{29}$$

We show this application on a standard ResNet110 classifier trained on CI-FAR10 (see Table 3). For the baselines, we use Gaussian noise ($\sigma = 0.35$) and its corresponding theoretical radius [17, 4] for certification. Our method uses I-OPT with General Normal noise and initializes it with the same σ . While approximating the certified radius with UniCR, we generate 4,000 samples with the Monte Carlo method on CIFAR10.

The table shows that over 98.6% of the inputs are certified by our method with a radius R > 0. This means that over 98.6% of the samples are certifiably protected while only 10.6% of inputs are certified by the baselines, which is nearly the accuracy by random guessing. This significant improvement emerges since the I-OPT could optimize the noise PDF for each input even though the classifier is not trained with noise (non-smoothed classifier). Although the certified radii are low compared to smoothly-trained classifiers, it provides a certifiable protection on perturbed data while maintaining the high accuracy for classifying clean data.

E Visual Examples of I-OPT

We present some examples of I-OPT on the ImageNet dataset against ℓ_1 , ℓ_2 and ℓ_{∞} perturbations, respectively (see Figure 6). In the first case (ℓ_1 perturbations), without executing the I-OPT, our UniCR certifies the input with a radius R = 1.24. Our I-OPT optimizes the distribution as the right-most figure shows, then the certified radius is improved to 1.48 with our UniCR. Similarly, in the rest cases, we show I-OPT can improve the certified radius significantly by optimizing the noise distribution. Especially, we improve the radius from 0.35 to 1.30 in the second case.

F Discussions

F.1 Universal Certified Robustness

It might be impractical to make a universal framework satisfy all the theoretical conditions w.r.t. all ℓ_p perturbations, especially p can be any positive real number. Thus, we admit that UniCR may not strictly satisfy certified robustness all the time due to the approximated optimization. However, extensive empirical results confirm that our derived radii highly approximate the theoretical certified radii against different ℓ_p perturbations. In addition, the defense performance against real attacks also illustrate that our method is as reliable as different theoretical certified radii. We believe that with the negligible error in practice, UniCR can be deployed as a universal framework to significantly ease the process of achieving certified robustness in different scenarios.

F.2 Certifying Perturbed Data with Randomized Smoothing

Randomized smoothing usually assumes that the input is clean and empirical defenses [11,7] are not applied, if the input data is perturbed before certification, then certification in I-OPT might be inaccurate. Indeed, the certification in traditional randomized smoothing (e.g., [4]) methods also depend on the inputs (since p_A is different for different inputs), they might be inaccurate if the

¹⁴ Hanbin Hong et al.



Fig. 6. Example images of applying I-OPT (based on UniCR) for smoothed classifier against different ℓ_p perturbations on the ImageNet dataset. From the left to right, the first figure shows the original image. The second and third figures show the smoothed image without I-OPT and with I-OPT, respectively. The fourth figure shows the corresponding distributions before/after I-OPT.

input data is perturbed, either. Thus, randomized smoothing focuses on certifying clean inputs rather than correcting perturbed inputs. We will study this interesting problem on certifying both clean and perturbed inputs in the future.

F.3 Can existing methods adopt noise optimization?

A question here is that if the noise optimization can improve the certified radius, can the theoretical methods provide personalized randomization for each input? The personalized randomization is actually not adaptable in the theoretical methods since they cannot automatically derive the certified radius for different noise distributions, especially for uncommon distributions, e.g., $e^{-|x/0.5|^{1.5}}$. Instead, our UniCR can automatically derive the certified radius for any distribution within the continuous parameter space.

F.4 Extensions

We evaluate our UniCR on the image classification. Indeed, our UniCR is a general method that can be directly applied to other tasks, e.g., video classification [10, 16], graph learning (e.g., node/graph classification [14] and community detection [8]), and natural language processing [9].

References

- 1. Andriushchenko, M., Croce, F., Flammarion, N., Hein, M.: Square attack: a queryefficient black-box adversarial attack via random search. In: European Conference on Computer Vision. pp. 484–501. Springer (2020)
- Chen, J., Jordan, M.I., Wainwright, M.J.: Hopskipjumpattack: A query-efficient decision-based attack. In: 2020 IEEE Symposium on Security and Privacy (2020)
- Co, K.T., Muñoz-González, L., de Maupeou, S., Lupu, E.C.: Procedural noise adversarial examples for black-box attacks on deep convolutional networks. In: ACM SIGSAC conference on computer and communications security (2019)
- 4. Cohen, J., Rosenfeld, E., Kolter, Z.: Certified adversarial robustness via randomized smoothing. In: International Conference on Machine Learning (2019)
- Croce, F., Hein, M.: Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In: International conference on machine learning. pp. 2206–2216. PMLR (2020)
- Dvoretzky, A., Kiefer, J., Wolfowitz, J.: Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator. The Annals of Mathematical Statistics pp. 642–669 (1956)
- Hong, H., Hong, Y., Kong, Y.: An eye for an eye: Defending against gradient-based attacks with gradients. arXiv preprint arXiv:2202.01117 (2022)
- Jia, J., Wang, B., Cao, X., Gong, N.Z.: Certified robustness of community detection against adversarial structural perturbation via randomized smoothing. In: Proceedings of The Web Conference 2020. pp. 2718–2724 (2020)
- 9. Jia, R., Raghunathan, A., Göksel, K., Liang, P.: Certified robustness to adversarial word substitutions. In: EMNLP/IJCNLP (2019)
- Li, S., Neupane, A., Paul, S., Song, C., Krishnamurthy, S.V., Roy-Chowdhury, A.K., Swami, A.: Stealthy adversarial perturbations against real-time video classification systems. In: NDSS. The Internet Society (2019)
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: International Conference on Learning Representations (2018)
- Mohammady, M., Xie, S., Hong, Y., Zhang, M., Wang, L., Pourzandi, M., Debbabi, M.: R2dp: A universal and automated approach to optimizing the randomization mechanisms of differential privacy for utility metrics with no known optimal distributions. In: Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security. pp. 677–696 (2020)
- Teng, J., Lee, G.H., Yuan, Y.: \$\ell_1\$ adversarial robustness certificates: a randomized smoothing approach (2020)
- 14. Wang, B., Jia, J., Cao, X., Gong, N.Z.: Certified robustness of graph neural networks against adversarial structural perturbation. In: KDD. ACM (2021)
- 15. Wong, E., Schmidt, F., Kolter, Z.: Wasserstein adversarial examples via projected sinkhorn iterations. In: International Conference on Machine Learning (2019)
- Xie, S., Wang, H., Kong, Y., Hong, Y.: Universal 3-dimensional perturbations for black-box attacks on video recognition systems. In: In Proceedings of the 43rd IEEE Symposium on Security and Privacy (Oakland'22) (2022)
- Yang, G., Duan, T., Hu, J.E., Salman, H., Razenshteyn, I., Li, J.: Randomized smoothing of all shapes and sizes. In: International Conference on Machine Learning. pp. 10693–10705. PMLR (2020)
- Zhang, D., Ye, M., Gong, C., Zhu, Z., Liu, Q.: Black-box certification with randomized smoothing: A functional optimization based framework (2020)