

Trustworthy Machine Learning

Certified Adversarial Learning

Sangdon Park

POSTECH

Motivation

- Heuristic adversarial learning often fails against powerful adversaries with the same maximum perturbation ϵ .

CIFAR10						
	Simple	Wide	Simple	Wide	Simple	Wide
Natural	92.7%	95.2%	87.4%	90.3%	79.4%	87.3%
FGSM	27.5%	32.7%	90.9%	95.1%	51.7%	56.1%
PGD	0.8%	3.5%	0.0%	0.0%	43.7%	45.8%
(a) Standard training			(b) FGSM training			(c) PGD training

- ▶ ϵ -FGSM training and ϵ -FGSM attacks: 90.9% accuracy :)
- ▶ ϵ -FGSM training and ϵ -PGD attacks: 0.0% accuracy :(

Motivation

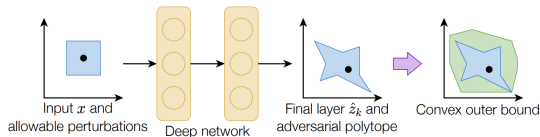
- Heuristic adversarial learning often fails against powerful adversaries with the same maximum perturbation ϵ .

CIFAR10						
	Simple	Wide	Simple	Wide	Simple	Wide
Natural	92.7%	95.2%	87.4%	90.3%	79.4%	87.3%
FGSM	27.5%	32.7%	90.9%	95.1%	51.7%	56.1%
PGD	0.8%	3.5%	0.0%	0.0%	43.7%	45.8%
(a) Standard training			(b) FGSM training		(c) PGD training	

- ▶ ϵ -FGSM training and ϵ -FGSM attacks: 90.9% accuracy :)
 - ▶ ϵ -FGSM training and ϵ -PGD attacks: 0.0% accuracy :(
- Can we learn a classifier robust to **any** small perturbations?

Certified Adversarial Learning

- Convex outer approximation [Kolter and Wong, 2017]



✓ Certified!

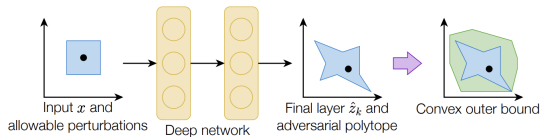
$$\max_{\|\delta\|_{\infty} \leq \varepsilon} \ell(f, x + \delta, y) \leq U(\varepsilon, f, x, y)$$

★ Essentially linear classification over overly approximated “convex polytope”-points

✗ Not scalable :(

Certified Adversarial Learning

- Convex outer approximation [Kolter and Wong, 2017]



✓ Certified!

$$\max_{\|\delta\|_{\infty} \leq \epsilon} \ell(f, x + \delta, y) \leq U(\epsilon, f, x, y)$$

★ Essentially linear classification over overly approximated “convex polytope”-points

✗ Not scalable :(

- Randomized smoothing: a post-hoc method

Certified Adversarial Robustness via Randomized Smoothing

Jeremy Cohen¹ Elan Rosenfeld¹ J. Zico Kolter^{1,2}

✓ (Probably) Certified!

✓ Scalable!

A Goodness Definition: Robustness

“Hard” Robustness

$$\forall \delta \text{ s.t. } \|\delta\|_p \leq \varepsilon, f(x + \delta) = f(x)$$

- $f : \mathcal{X} \rightarrow \mathcal{Y}$: a hard-classifier

A Goodness Definition: Robustness

“Hard” Robustness

$$\forall \delta \text{ s.t. } \|\delta\|_p \leq \varepsilon, f(x + \delta) = f(x)$$

- $f : \mathcal{X} \rightarrow \mathcal{Y}$: a hard-classifier
- The constraint on the perturbation δ can be more general.

A Goodness Definition: Robustness

“Hard” Robustness

$$\forall \delta \text{ s.t. } \|\delta\|_p \leq \varepsilon, f(x + \delta) = f(x)$$

- $f : \mathcal{X} \rightarrow \mathcal{Y}$: a hard-classifier
- The constraint on the perturbation δ can be more general.
- It does not matter whether $f(x)$ is correct.

A Certified Method: Randomized Smoothing

Smoothed Classifier

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

- $g : \mathcal{X} \rightarrow \mathcal{Y}$: a smoothed classifier
- σ is related to the maximum perturbation ε .
- Easier than convex outer approximation

Robustness Guarantee

Binary Classification

Theorem

Suppose that $\underline{p}_A \in (0.5, 1]$ satisfies

$$\mathbb{P} \{f(x + \delta) = c_A\} = p_A \geq \underline{p}_A \quad \text{where} \quad \delta \sim \mathcal{N}(0, \sigma^2 I).$$

Then, we have $g(x + \delta) = c_A$ if

$$\|\delta\|_2 < \sigma \Phi^{-1}(\underline{p}_A).$$

- c_A : the most probable class when f classifies $x + \varepsilon$
- p_A : the chance that f classifies $x + \delta$ by c_A
- \underline{p}_A : the lower bound of p_A
- Φ^{-1} : the inverse of the standard Gaussian CDF

Robustness Guarantee

Binary Classification

Theorem

Suppose that $\underline{p}_A \in (0.5, 1]$ satisfies

$$\mathbb{P} \{f(x + \delta) = c_A\} = p_A \geq \underline{p}_A \quad \text{where} \quad \delta \sim \mathcal{N}(0, \sigma^2 I).$$

Then, we have $g(x + \delta) = c_A$ if

$$\|\delta\|_2 < \sigma \Phi^{-1}(\underline{p}_A).$$

- c_A : the most probable class when f classifies $x + \varepsilon$
- p_A : the chance that f classifies $x + \delta$ by c_A
- \underline{p}_A : the lower bound of p_A
- Φ^{-1} : the inverse of the standard Gaussian CDF
- Here, we assume that we can compute \underline{p}_A .

Robustness Guarantee

Binary Classification

Theorem

Suppose that $\underline{p}_A \in (0.5, 1]$ satisfies

$$\mathbb{P} \{f(x + \delta) = c_A\} = p_A \geq \underline{p}_A \quad \text{where} \quad \delta \sim \mathcal{N}(0, \sigma^2 I).$$

Then, we have $g(x + \delta) = c_A$ if

$$\|\delta\|_2 < \sigma \Phi^{-1}(\underline{p}_A).$$

- c_A : the most probable class when f classifies $x + \varepsilon$
- p_A : the chance that f classifies $x + \delta$ by c_A
- \underline{p}_A : the lower bound of p_A
- Φ^{-1} : the inverse of the standard Gaussian CDF
- Here, we assume that we can compute \underline{p}_A .
- We can compute the *data-dependent* maximum perturbation to be robust!

Robustness Guarantee: A Proof Sketch (1/3)

Binary Classification

- Fix a perturbation δ .
- From the definition of g , we have

$$\begin{aligned} g(x + \delta) &:= \arg \max_c \mathbb{P} \{f(x + \varepsilon + \delta) = c\} \quad \text{where } \varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I) \\ &= \arg \max_c \mathbb{P} \{f(x + \varepsilon') = c\} \quad \text{where } \varepsilon' \sim \mathcal{N}(\delta, \sigma^2 I) \\ &\stackrel{?}{=} c_A \end{aligned} \tag{1}$$

Robustness Guarantee: A Proof Sketch (1/3)

Binary Classification

- Fix a perturbation δ .
- From the definition of g , we have

$$\begin{aligned} g(x + \delta) &:= \arg \max_c \mathbb{P} \{f(x + \varepsilon + \delta) = c\} \quad \text{where } \varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I) \\ &= \arg \max_c \mathbb{P} \{f(x + \varepsilon') = c\} \quad \text{where } \varepsilon' \sim \mathcal{N}(\delta, \sigma^2 I) \\ &\stackrel{?}{=} c_A \end{aligned} \tag{1}$$

- We wish to prove (1) for any classifier f under some condition. How?
 - ▶ f can be any classifier, which is not easy to analyze.
 - ▶ Consider a surrogate classifier that bounds the probability and is easier to analyze, e.g.,

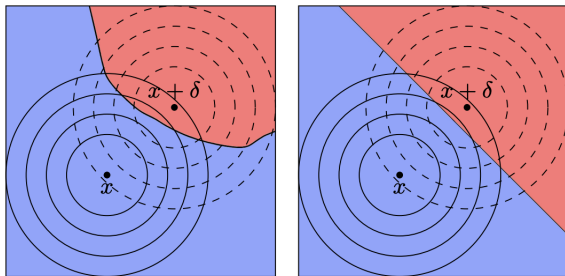
$$\mathbb{P} \{f(x + \varepsilon') = c_A\} \geq \min_{f': \mathbb{P} \{f'(x + \varepsilon) = c_A\} \geq \underline{p}_A} \mathbb{P} \{f'(x + \varepsilon') = c_A\} > \frac{1}{2} \implies g(x + \delta) = c_A.$$

Robustness Guarantee: A Proof Sketch (2/3)

Binary Classification

- Interestingly, f^* is linear (due to the Neyman-Perason lemma), where

$$f^* = \arg \min_{f': \mathbb{P}\{f'(x+\epsilon)=c_A\} \geq \underline{p}_A} \mathbb{P}\{f'(x+\epsilon')=c_A\}$$



- There could be a non-linear classifier but we can find a corresponding linear classifier with the same minimum value.

Robustness Guarantee: A Proof Sketch (3/3)

Binary Classification

- We have a closed-form solution of f^* :

$$f^*(x') := \begin{cases} c_A & \text{if } \delta^T(x' - x) \leq \sigma \|\delta\|_2 \Phi^{-1}(\underline{p}_A) \\ c_B & \text{otherwise} \end{cases}.$$

- This (non-trivially) implies the following minimum value:

$$\min_{f': \mathbb{P}\{f'(x+\varepsilon)=c_A\} \geq \underline{p}_A} \mathbb{P}\{f'(x+\varepsilon')=c_A\} = \mathbb{P}\{f^*(x+\varepsilon')=c_A\} = \Phi\left(\Phi^{-1}(\underline{p}_A) - \frac{\|\delta\|_2}{\sigma}\right)$$

- The above probability should be larger than $\frac{1}{2}$, i.e.,

$$\Phi\left(\Phi^{-1}(\underline{p}_A) - \frac{\|\delta\|_2}{\sigma}\right) > \frac{1}{2} \quad \implies \quad \|\delta\|_2 < \sigma \Phi^{-1}(\underline{p}_A).$$

Robustness Guarantee

Multi-class Classification

Theorem

Suppose that $\underline{p}_A, \overline{p}_B \in [0, 1]$ satisfy

$$\mathbb{P} \{f(x + \varepsilon) = c_A\} \geq \underline{p}_A \geq \overline{p}_B \geq \max_{c \neq c_A} \mathbb{P} \{f(x + \varepsilon) = c\}.$$

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

$$R := \frac{\sigma}{2} \left(\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B) \right).$$

- c_A : the most probable label (with probability at least \underline{p}_A)
- $c_B := \arg \max_{c \neq c_A} \mathbb{P} \{f(x + \varepsilon) = c\}$: the second-most probable label (with probability at most \overline{p}_B)

Robustness Guarantee: An Alternative

Multi-class Classification

Theorem

Suppose that we have class A and B that satisfy

$$\max_k \mathbb{P} \{f_k(x + \varepsilon)\} = p_A \geq p_B = \max_{k \neq A} \mathbb{P} \{f_k(x + \varepsilon)\}.$$

Then, we have $g(x + \delta) = A$ for all $\|\delta\|_2 < R$, where

$$R := \frac{\sigma}{2} (\Phi^{-1}(p_A) - \Phi^{-1}(p_B)).$$

- Consider a soft classifier $f_k : \mathcal{X} \rightarrow [0, 1]$ for class k
- A smoothed classifier $g_k(x) := \arg \max_k \mathbb{P}_\varepsilon \{f_k(x + \varepsilon)\}$

Robustness Guarantee: An Alternative Proof Sketch (1/2)

- Let $f_k : \mathbb{R}^n \rightarrow [0, 1]$: a soft classifier for class k
- Let $\tilde{f}_k : \mathcal{X} \rightarrow [0, 1]$: a smoothed classifier for class k , i.e.,

$$\tilde{f}_k(x) := (f_k * \mathcal{N}(0, \sigma I))(x) = \int_{\mathbb{R}^n} f_k(t) \frac{\exp\left(-\frac{1}{2\sigma^2} \|x - t\|^2\right)}{(2\pi\sigma^2)^n} dt = \mathbb{P}_\varepsilon\{f_k(x + \varepsilon)\}$$

- ▶ The convolution of f_k and $\mathcal{N}(0, \sigma I)$, a.k.a. the Weierstrass transform of f_k

Robustness Guarantee: An Alternative Proof Sketch (1/2)

- Let $f_k : \mathbb{R}^n \rightarrow [0, 1]$: a soft classifier for class k
- Let $\tilde{f}_k : \mathcal{X} \rightarrow [0, 1]$: a smoothed classifier for class k , i.e.,

$$\tilde{f}_k(x) := (f_k * \mathcal{N}(0, \sigma I))(x) = \int_{\mathbb{R}^n} f_k(t) \frac{\exp\left(-\frac{1}{2\sigma^2} \|x - t\|^2\right)}{(2\pi\sigma^2)^n} dt = \mathbb{P}_\varepsilon\{f_k(x + \varepsilon)\}$$

- ▶ The convolution of f_k and $\mathcal{N}(0, \sigma I)$, a.k.a. the Weierstrass transform of f_k
- Let p_A is the most-probable class probability assigned by the smoothed classifier \tilde{f}_k , i.e.,

$$p_A = \tilde{f}_A(x) \quad \text{where} \quad A = \arg \max_k \tilde{f}_k(x)$$

- Let p_B is the class probability by \tilde{f}_k such that $A \neq B$ and less than p_A , i.e.,

$$p_B = \tilde{f}_B(x) \leq p_A.$$

Robustness Guarantee: An Alternative Proof Sketch (1/2)

- Let $f_k : \mathbb{R}^n \rightarrow [0, 1]$: a soft classifier for class k
- Let $\tilde{f}_k : \mathcal{X} \rightarrow [0, 1]$: a smoothed classifier for class k , i.e.,

$$\tilde{f}_k(x) := (f_k * \mathcal{N}(0, \sigma I))(x) = \int_{\mathbb{R}^n} f_k(t) \frac{\exp\left(-\frac{1}{2\sigma^2} \|x - t\|^2\right)}{(2\pi\sigma^2)^n} dt = \mathbb{P}_\varepsilon\{f_k(x + \varepsilon)\}$$

- ▶ The convolution of f_k and $\mathcal{N}(0, \sigma I)$, a.k.a. the Weierstrass transform of f_k
- Let p_A is the most-probable class probability assigned by the smoothed classifier \tilde{f}_k , i.e.,

$$p_A = \tilde{f}_A(x) \quad \text{where} \quad A = \arg \max_k \tilde{f}_k(x)$$

- Let p_B is the class probability by \tilde{f}_k such that $A \neq B$ and less than p_A , i.e.,

$$p_B = \tilde{f}_B(x) \leq p_A.$$

- Let Φ be a CDF of a Gaussian distribution, i.e.,

$$\Phi(a) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a \exp\left(-\frac{1}{2}s^2\right) ds$$

Robustness Guarantee: An Alternative Proof Sketch (2/2)

Then, we have the robustness guarantee due to the following reasons:

- For any perturbation δ and any class k , we have

$$\left| \Phi^{-1} \left(\tilde{f}_k(x) \right) - \Phi^{-1} \left(\tilde{f}_k(x + \delta) \right) \right| \leq \frac{1}{\sigma} \|\delta\|_2.$$

- ▶ $\Phi^{-1} \circ \tilde{f}$ is $\frac{1}{\sigma}$ -Lipschitz (check out the paper)
- Consider any adversarial perturbation $\bar{\delta}$ that changes the classification result, i.e.,

$$\tilde{f}_A(x + \bar{\delta}) \leq \tilde{f}_B(x + \bar{\delta}) \quad \text{for some } B \neq A$$

- For $\bar{\delta}$, we have

$$\begin{aligned} \frac{1}{\sigma} \|\bar{\delta}\|_2 &\geq \Phi^{-1} \left(\tilde{f}_A(x) \right) - \Phi^{-1} \left(\tilde{f}_A(x + \bar{\delta}) \right) \quad \text{and} \quad \frac{1}{\sigma} \|\bar{\delta}\|_2 \geq \Phi^{-1} \left(\tilde{f}_B(x + \bar{\delta}) \right) - \Phi^{-1} \left(\tilde{f}_B(x) \right) \\ \Rightarrow \frac{2}{\sigma} \|\bar{\delta}\|_2 &\geq \left\{ \Phi^{-1} \left(\tilde{f}_A(x) \right) - \Phi^{-1} \left(\tilde{f}_B(x) \right) \right\} + \left\{ \Phi^{-1} \left(\tilde{f}_B(x + \bar{\delta}) \right) - \Phi^{-1} \left(\tilde{f}_A(x + \bar{\delta}) \right) \right\} \\ &\geq \Phi^{-1} \left(\tilde{f}_A(x) \right) - \Phi^{-1} \left(\tilde{f}_B(x) \right) = \Phi^{-1} (p_A) - \Phi^{-1} (p_B) \end{aligned}$$

Prediction

```
function  $\tilde{\text{PREDICT}}(f, \sigma, x, n, \alpha)$   
  counts  $\leftarrow \text{SAMPLEUNDERNOISE}(f, x, n, \sigma)$   
   $\hat{c}_A, \hat{c}_B \leftarrow \text{top two indices in counts}$   
   $n_A, n_B \leftarrow \text{counts}[\hat{c}_A], \text{counts}[\hat{c}_B]$   
  if  $\text{BINOMPVALUE}(n_A, n_A + n_B, 0.5) \leq \alpha$  return  $\hat{c}_A$   
  else return ABSTAIN
```

- Recall the randomized smoothing method:

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

- 1 Draw n noisy perturbations $\delta_1, \dots, \delta_n$.
 - 2 Empirically compute the most probable and the second most probable labels, *i.e.*, \hat{c}_A and \hat{c}_B .
 - 3 If \hat{c}_A is drawn from the binomial distribution with $p = 0.5$, return \hat{c}_A .
- Alternatively, you can use the (our-favorite) binomial tail bound.

Certification in Evaluation

certify the robustness of g around x

function CERTIFY($f, \sigma, x, n_0, n, \alpha$)

counts0 \leftarrow SAMPLEUNDERNOISE(f, x, n_0, σ)

$\hat{c}_A \leftarrow$ top index in counts0

counts \leftarrow SAMPLEUNDERNOISE(f, x, n, σ)

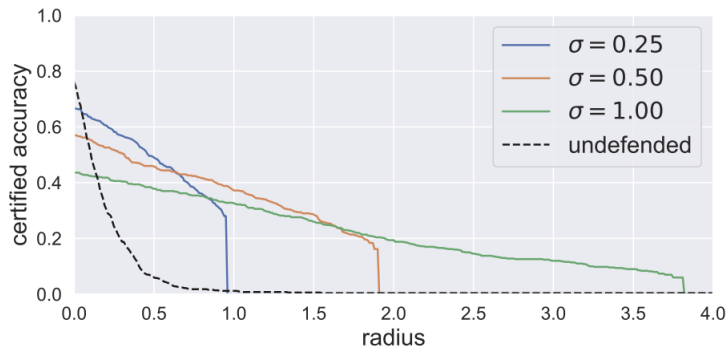
$\underline{p}_A \leftarrow$ LOWERCONFBOUND(counts[\hat{c}_A], $n, 1 - \alpha$)

if $\underline{p}_A > \frac{1}{2}$ **return** prediction \hat{c}_A and radius $\sigma \Phi^{-1}(\underline{p}_A)$

else return ABSTAIN

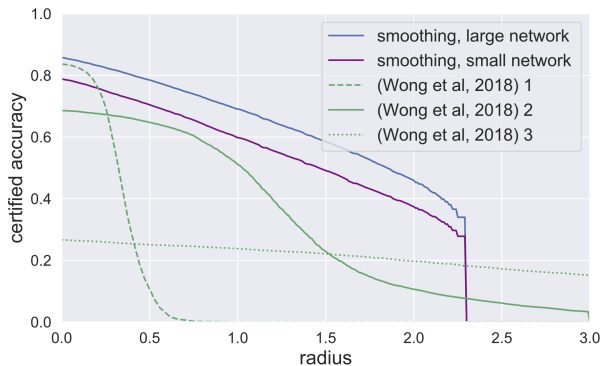
- 1 Compute \underline{p}_A via the binomial tail bound.
- 2 Compute the robust radius, i.e., $\sigma \Phi^{-1}(\underline{p}_A)$.
- 3 If (a desired radius) $\leq \sigma \Phi^{-1}(\underline{p}_A)$, then “certified”.

Results: ImageNet



- Classifier: ResNet-50
- undefended: a classifier with heuristic adversarial training (using ℓ_2 adversarial attacks)
- perturbation: $\|\delta\|_2 \leq (\text{radius}) = (\text{maximum perturbation size})$

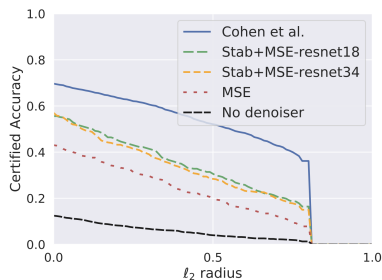
Results: Comparison



- (maybe) on MNIST
- Baseline: deterministic robustness guarantee
- randomized smoothing: high-probability guarantee

Limitation of Randomized Smoothing

- Randomized smoothing requires retraining (e.g., Gaussian data augmentation).



- ▶ Cohen et al.: Randomized smoothing with retraining
- ▶ No denoiser: Randomized smoothing without retraining
- How to avoid retraining?

Denoised Smoothing: A Provable Defense for Pretrained Classifiers

Hadi Salman
hasalman@microsoft.com
Microsoft Research

Mingjie Sun
mingjies@cs.cmu.edu
CMU

Greg Yang
gragyang@microsoft.com
Microsoft Research

Ashish Kapoor
akapoor@microsoft.com
Microsoft Research

J. Zico Kolter
zkolter@cs.cmu.edu
CMU

- A classifier randomized smoothing needs to be robust to Gaussian noise for better certified robustness.
- How about using denoised smoothing and then use the randomized smoothing?

Denoised Smoothing

Randomized Smoothing:

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

- Applicable for any classifier f

Denoised Smoothing:

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

- $\mathcal{D} : \mathcal{X} \rightarrow \mathcal{X}$: a denoiser that hopefully removes δ .
- Consider a NEW classifier $f \circ \mathcal{D}$ and then enjoy randomized smoothing.
- Retraining f is not required.

How to Train a Denoiser?

MSE objective:

$$L_{\text{MSE}} := \mathbb{E}_{x,y,\delta} \|\mathcal{D}(x + \delta) - x\|_2^2$$

How to Train a Denoiser?

MSE objective:

$$L_{\text{MSE}} := \mathbb{E}_{x,y,\delta} \|\mathcal{D}(x + \delta) - x\|_2^2$$

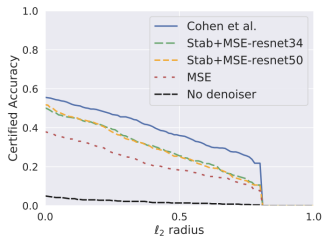
✗ Does not consider the accuracy of a classifier.

Stability objective:

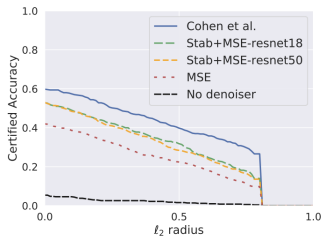
$$L_{\text{Stab}} := \mathbb{E}_{x,y,\delta} \ell_{CE}(F(\mathcal{D}(x + \delta)), f(x)) \quad \text{where} \quad \delta \sim \mathcal{N}(0, \sigma^2 I)$$

- $f : \mathcal{X} \rightarrow \mathcal{Y}$: a hard classifier
- $F : \mathcal{X} \rightarrow [0, 1]^{|\mathcal{Y}|}$: a soft classifier
- ✓ Find a denoiser \mathcal{D} that does not change predictions by the classifier f .

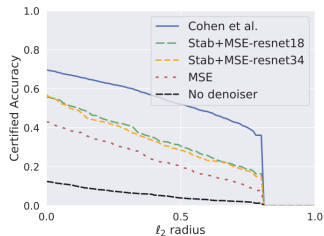
Results



(a) ResNet-18



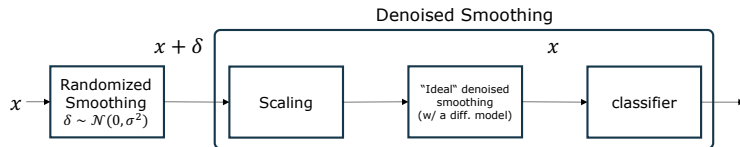
(b) ResNet-34



(c) ResNet-50

- The denoised smoothing without retraining is quite similar to the randomized smoothing with retraining.
- But not outperform the retraining one. How can we train a better denoiser?

Diffusion Models as Denoisers



Assumptions:

- A diffusion model assumes the following noise model:

$$x_t := \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t} \cdot \mathcal{N}(0, \mathbf{I}),$$

where x_0 is an initial example, t is a timestep, and α_t is any noise scheduler (monotonically decreasing in t).

- ▶ Under this noise model, an ideal denoiser finds x_0 from x_t .

Method:

- Find t^* for randomized smoothing that fits to the noise model for a diffusion model, *i.e.*,

$$\text{find } t \text{ subj. to } x_0 + \mathcal{N}(0, \sigma^2 \mathbf{I}) \approx \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t} \cdot \mathcal{N}(0, \mathbf{I})$$

Conclusion

- Randomized smoothing provides a simple defense mechanism.
- Denoised smoothing does not require to retrain a classifier (but still requires training the denoiser).
- Recently, the denoised smoothing was improved via denoising diffusion models [Carlini et al., 2023].

Method	Off-the-shelf	Extra data	Certified Accuracy at ϵ (%)				
			0.5	1.0	1.5	2.0	3.0
PixelDP (Lecuyer et al., 2019)	○	✗	(33.0) 16.0	-	-		
RS (Cohen et al., 2019)	○	✗	(67.0) 49.0	(57.0) 37.0	(57.0) 29.0	(44.0) 19.0	(44.0) 12.0
SmoothAdv (Salman et al., 2019)	○	✗	(65.0) 56.0	(54.0) 43.0	(54.0) 37.0	(40.0) 27.0	(40.0) 20.0
Consistency (Jeong & Shin, 2020)	○	✗	(55.0) 50.0	(55.0) 44.0	(55.0) 34.0	(41.0) 24.0	(41.0) 17.0
MACER (Zhai et al., 2020)	○	✗	(68.0) 57.0	(64.0) 43.0	(64.0) 31.0	(48.0) 25.0	(48.0) 14.0
Boosting (Horváth et al., 2022a)	○	✗	(65.6) 57.0	(57.0) 44.6	(57.0) 38.4	(44.6) 28.6	(38.6) 21.2
DRT (Yang et al., 2021)	○	✗	(52.2) 46.8	(55.2) 44.4	(49.8) 39.8	(49.8) 30.4	(49.8) 23.4
SmoothMix (Jeong et al., 2021)	○	✗	(55.0) 50.0	(55.0) 43.0	(55.0) 38.0	(40.0) 26.0	(40.0) 20.0
ACES (Horváth et al., 2022b)	●	✗	(63.8) 54.0	(57.2) 42.2	(55.6) 35.6	(39.8) 25.6	(44.0) 19.8
Denoised (Salman et al., 2020)	●	✗	(60.0) 33.0	(38.0) 14.0	(38.0) 6.0	-	-
Lee (Lee, 2021)	●	✗	41.0	24.0	11.0	-	-
Ours	●	✓	(82.8) 71.1	(77.1) 54.3	(77.1) 38.1	(60.0) 29.5	(60.0) 13.1

Reference I

- N. Carlini, F. Tramer, K. D. Dvijotham, L. Rice, M. Sun, and J. Z. Kolter. (certified!!) adversarial robustness for free!, 2023.
- J. Z. Kolter and E. Wong. Provable defenses against adversarial examples via the convex outer adversarial polytope. *arXiv preprint arXiv:1711.00851*, 2017.