



西安电子科技大学
XIDIAN UNIVERSITY

Certified Adversarial Robustness

Mengdie Huang

April 1, 2022





Overview

1

Background

2

Differential Privacy Scheme

3

Randomized Smoothing Scheme

4

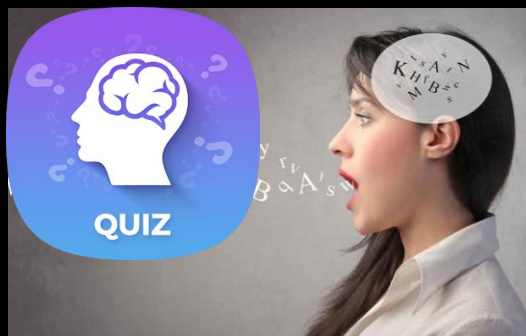
Conclusion



AI Applications



Computer Vision



Speech Recognition



Natural Language Processing



Expert system



Smart Robot

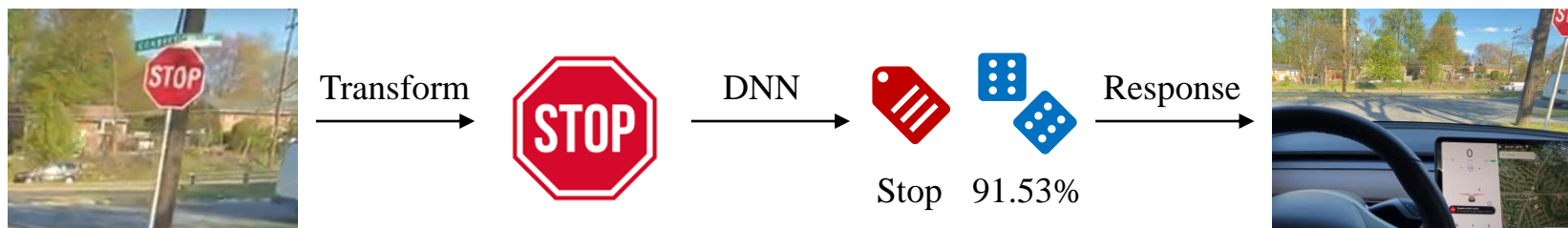


Chess Game

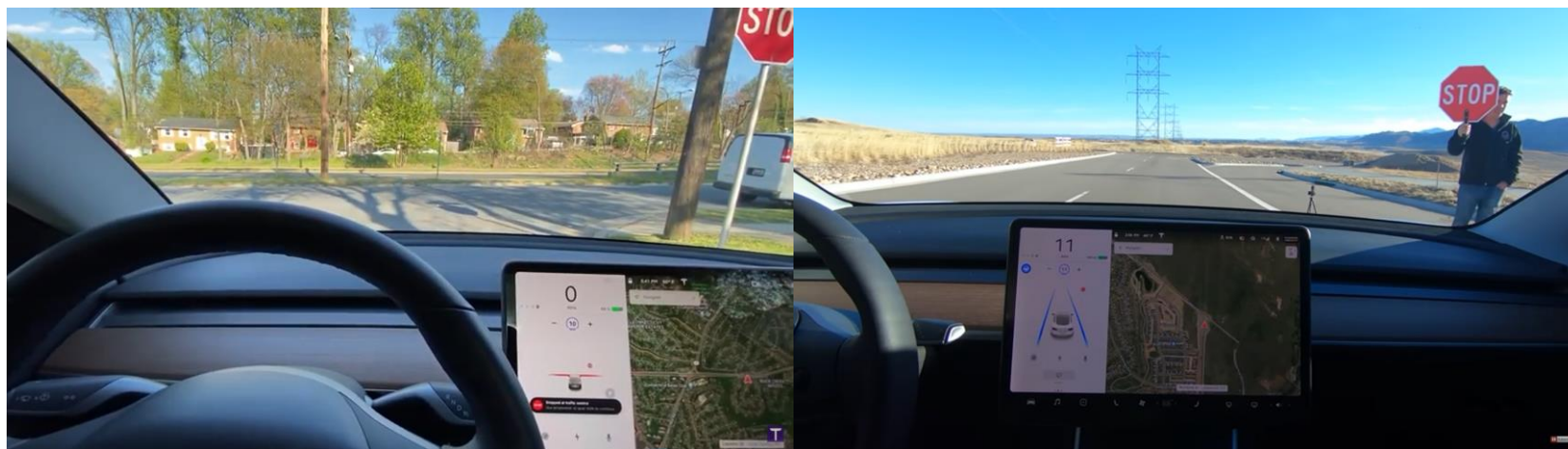


AI Error

- Traffic Light & Stop Sign Must Reads by Autopilot



- Stop Sign can be Ignored by Autopilot



Autopilot action: Stop

Autopilot action: Keep going



Attack

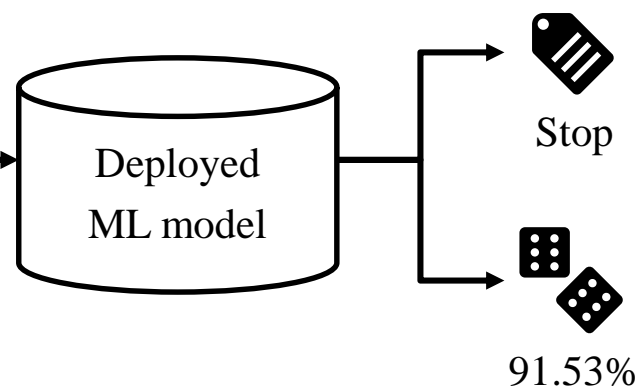
German Traffic Sign Recognition
Benchmark (GTSRB) dataset

- Adversarial attack

- Normal Prediction



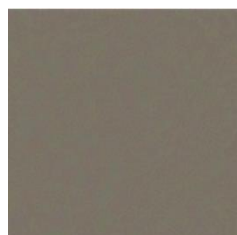
Normal Sample



- Adversarial Attack



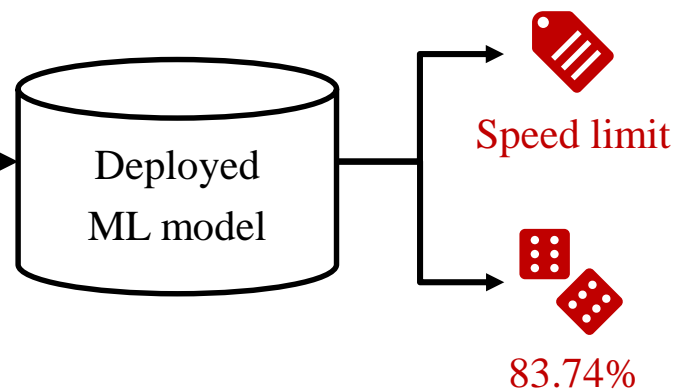
+



=



Normal Sample Perturbation Adversarial Sample

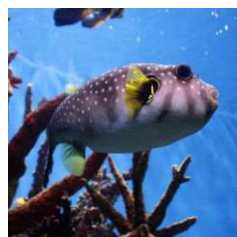




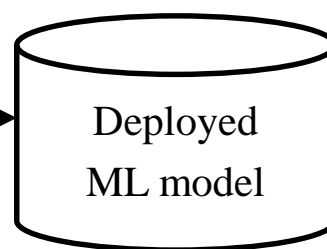
Attack

● Adversarial attack

➤ Normal Prediction

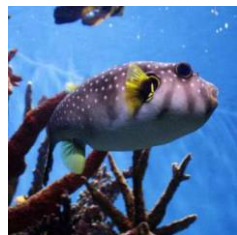


Normal Sample

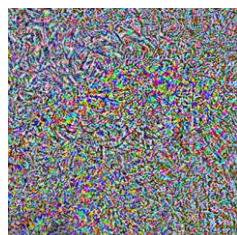


97.99%

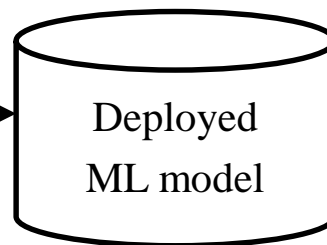
➤ Adversarial Attack



+



=



100%

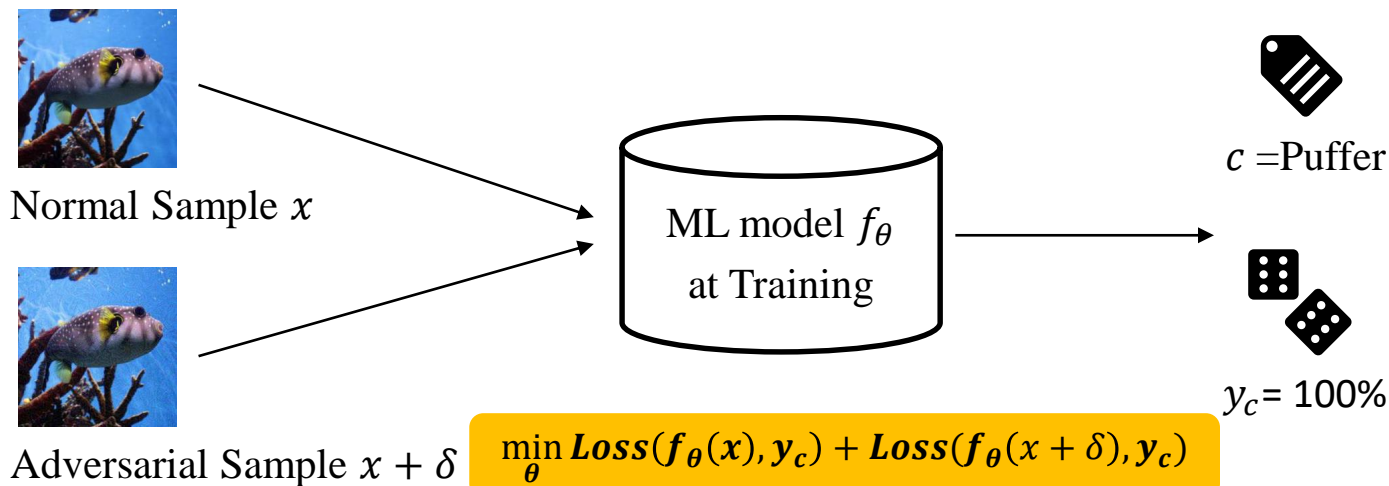
ImageNet dataset



Defense Category

● Empirical Defense

➤ Typical method: Adversarial training



➤ Drawback:

most of heuristics defenses have been shown to **fail against suitably powerful adversaries** (cat-and-mouse game).

➤ New requirement:

rigorous, theory-backed defensive approaches to stop this arms race.



Defense Category

- Empirical Defense
 - Typical method: Adversarial training

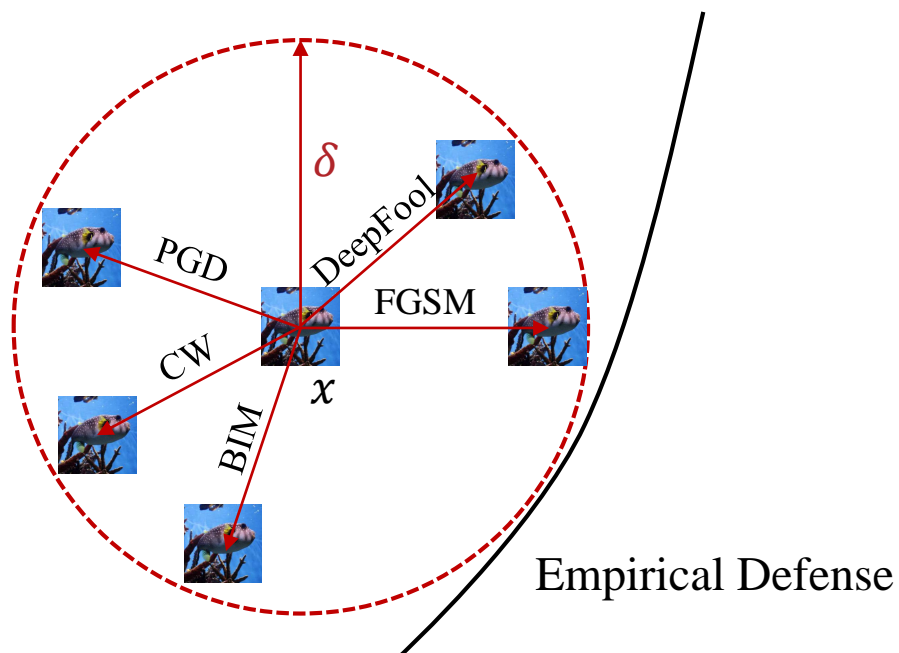


x



Defense Category

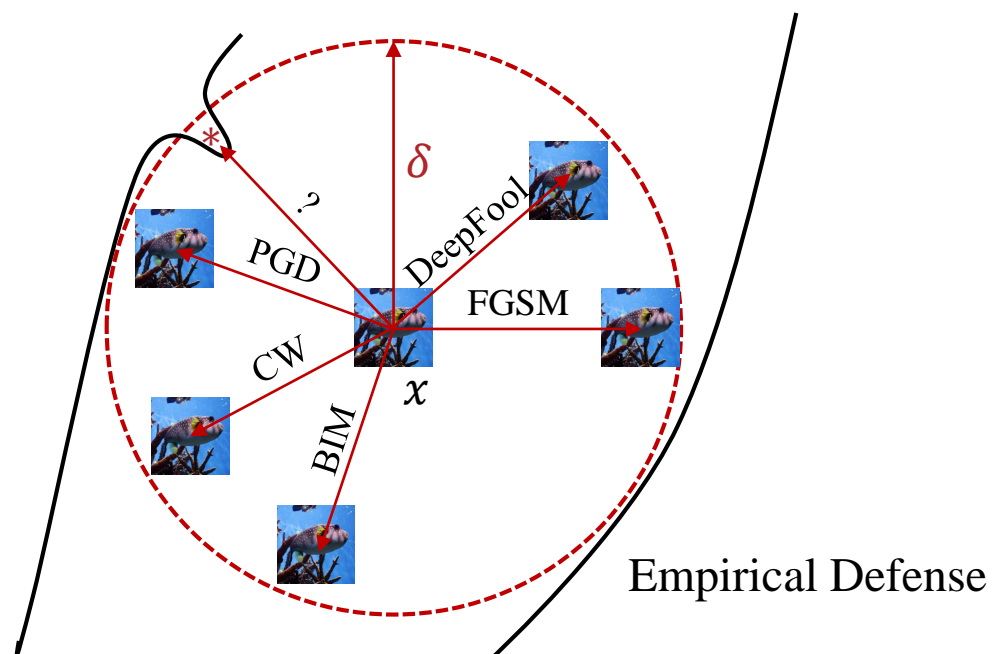
- Empirical Defense
 - Typical method: Adversarial training





Defense Category

- Empirical Defense
 - Typical method: Adversarial training



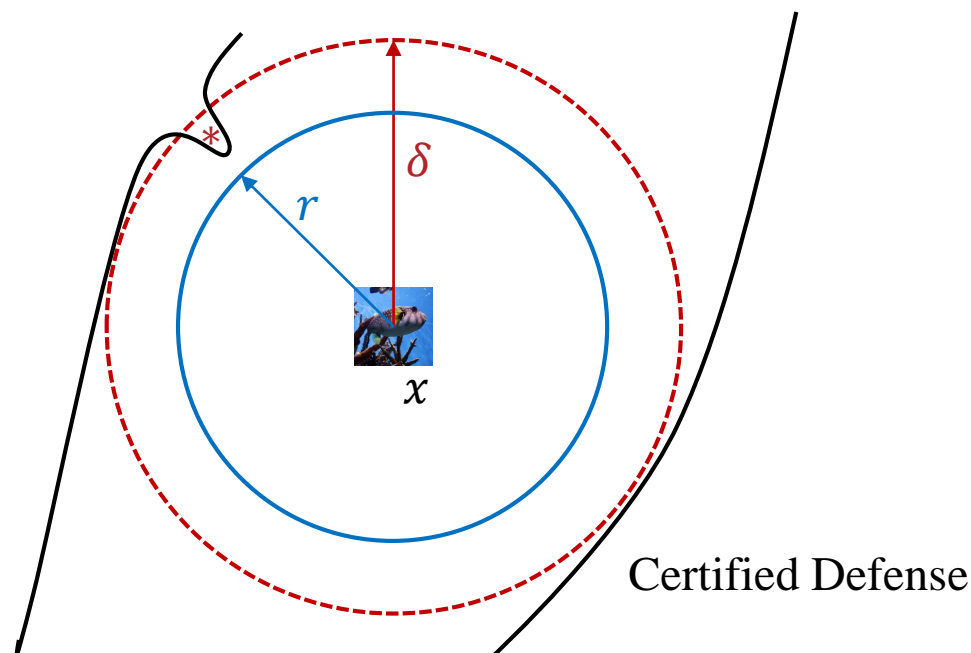


Defense Category

- Certified Defense

- Provide a certificate for adversarial robustness
- Certificate (x, f, r)

For any input x , the prediction output by the classifier on some set around x are **guaranteed** to be **constant**.





Defense Category

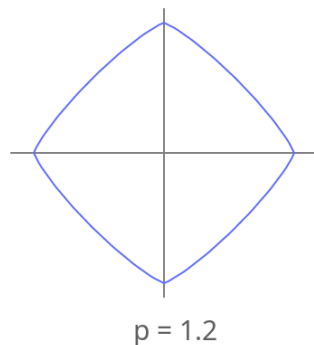
● Certified Defense

➤ Provide a certificate for adversarial robustness

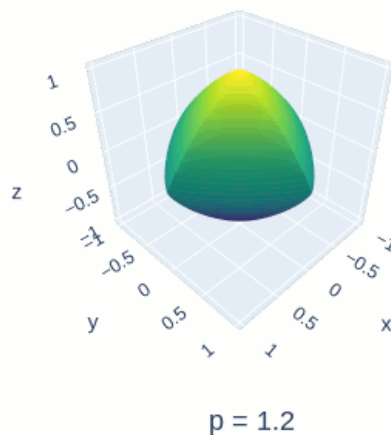
➤ Certificate (x, f, r)

For any input x , the prediction output by the classifier on some set around x are **guaranteed** to be **constant**.

➤ Some set: L_p ball with radius r



2D L_p Balls



3D L_p Balls

L_1 ball with $r = 1$

球

L_2 ball with $r = 1$

L_3 ball with $r = 1$

L_4 ball with $r = 1$

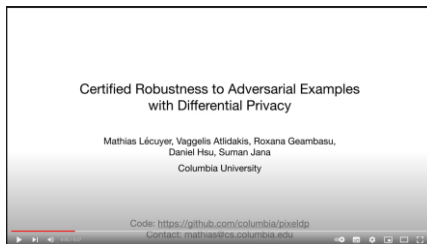
...

L_∞ ball with $r = 1$

立方体



Certified Robustness to Adversarial Examples with Differential Privacy

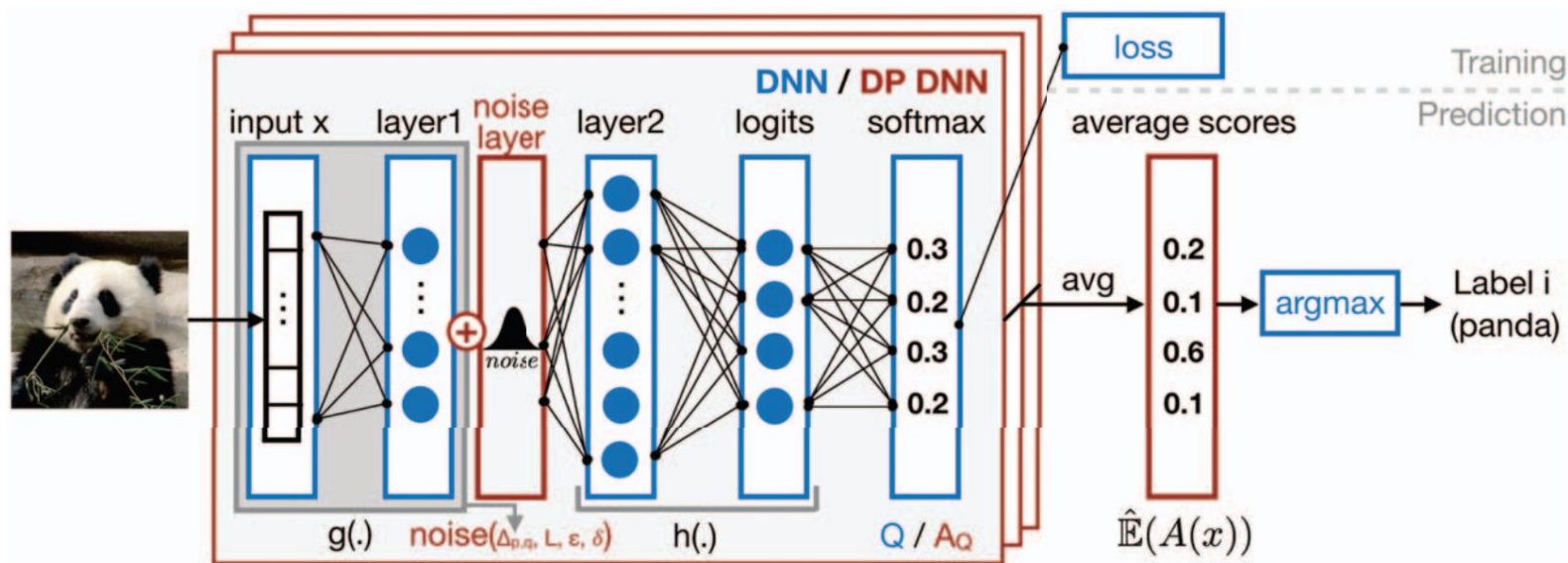


Mathias Lecuyer, ... , Suman Jana
Columbia University
IEEE S&P 2019



Overview

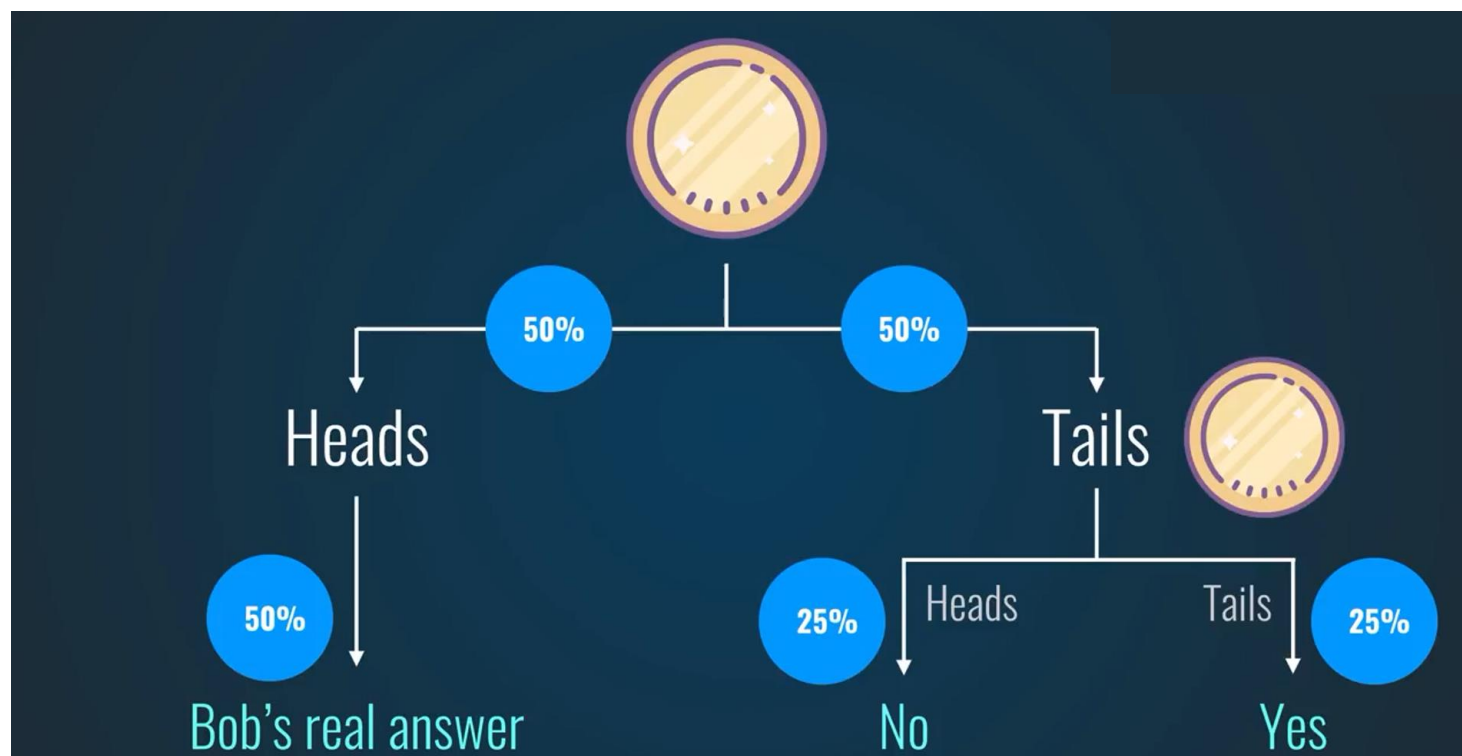
- Design a DNN classifier with differential privacy properties
 - Essence: Introduce randomness into the prediction of the classifier.
 - Way: Add a noise layer to the network.
 - Purpose: Guarantee that the output of the model is insensitive to small changes in the input.





Preliminary

- Differential Privacy
- Randomize responses

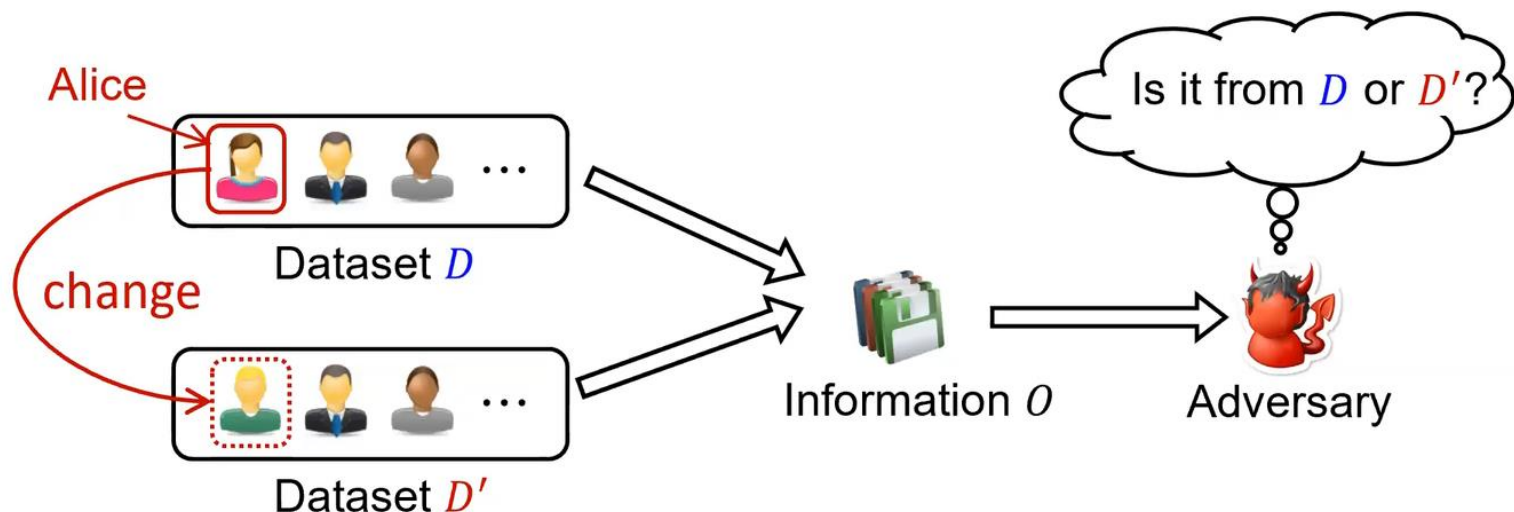




Preliminary

• ϵ -DP (Differential Privacy)

➤ Intuitive understanding



- The only difference between D and D' is Alice.
- If the attacker cannot tell whether the information O comes from D or D' , it can be considered that the privacy of Alice is protected.
- DP requires information O to be randomized before output.



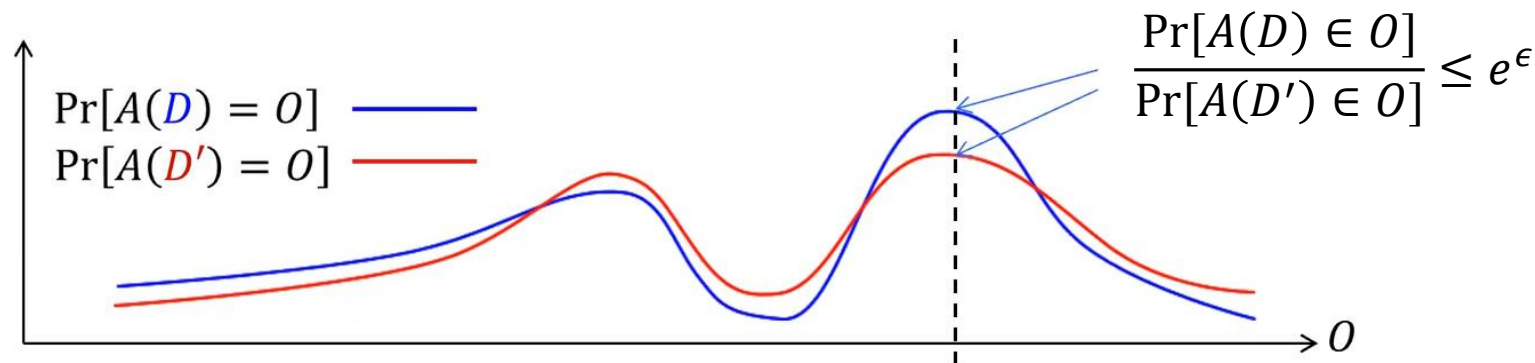
Preliminary

- ϵ -DP (Differential Privacy)

- Formalizing

Randomized algorithm A is ϵ -differentially private if for **any** $O \subseteq \text{Range}(A)$ and for **any neighboring dataset** D, D' ($\|D - D'\|_1 \leq 1$):

$$\Pr[A(D) \in O] \leq e^\epsilon \Pr[A(D') \in O]$$





Preliminary

- (ϵ, δ) -DP (Differential Privacy)

- Formalizing

Randomized algorithm A is (ϵ, δ) -differentially private if for **any** $O \subseteq \text{Range}(A)$ and for **any neighboring dataset** D, D' ($\|D - D'\|_1 \leq 1$):

$$\Pr[A(D) \in O] \leq e^\epsilon \Pr[A(D') \in O] + \delta$$

$$\frac{\Pr[A(D) \in O] - \delta}{\Pr[A(D') \in O]} \leq e^\epsilon$$

- when an event is more likely under D than under D' , δ is positive(+).
- when an event is more likely under D' than under D , δ is negative(-).
- $\|D - D'\|_1 \leq 1$ can be generalized to $\|D - D'\|_p \leq L$ by applying group privacy.



Preliminary

● Properties of DP

➤ Post-processing

If $A(x)$ satisfies (ϵ, δ) -DP, h is a x -independent mapping algorithm,
Then, the composition $h \circ A = h(A(x))$ satisfies (ϵ, δ) -DP.

➤ Expected output stability

The expected value $\mathbb{E}(A(x))$ of an (ϵ, δ) -DP algorithm A with bounded output $A(x) \in [0, b]$ is not sensitive to small changes in the input.

$$\forall \alpha \leq \text{Ball}_p(r=1), \quad \mathbb{E}(A(x)) \leq e^\epsilon \mathbb{E}(A(x + \alpha)) + b\delta$$

✓ Proof

- $\Pr[A(x) \in O] \leq e^\epsilon \Pr[(x + \alpha) \in O] + \delta$ (ϵ, δ)-DP定义
- $\int_0^b \Pr[A(x) \in O] dt \leq \int_0^b e^\epsilon \Pr[(x + \alpha) \in O] + \delta dt$ 积分
- $\int_0^b \Pr[A(x) < t] dt \leq e^\epsilon \int_0^b \Pr[(x + \alpha) > t] dt + \int_0^b \delta dt$ PDF定义
- $\mathbb{E}(A(x)) \leq e^\epsilon \mathbb{E}(A(x + \alpha)) + b\delta$ 期望定义



DP-robustness Connection

● One-to-one correspondence

Differential Privacy



Sex	Blood	...	HIV
F	B	...	Y
M	A	...	N
M	O	...	N
M	O	...	Y
F	A	...	N
M	B	...	Y

☐ Databases

$A(d)$ ☐ Randomized query algorithm

☐ Records in a database

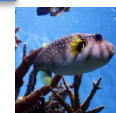
☐ Remove/alter on records

$\Delta A(d) \leq 1$ ☐ Result in a bounded change in the output

Adversarial Robustness



☐ Input images



☐ **Randomized classify algorithm**

$A(x)$

☐ Pixels in an images



☐ Increase or decrease pixels



☐ Result in a bounded change in the prediction

$A(x + \alpha) = A(x)$



Notation

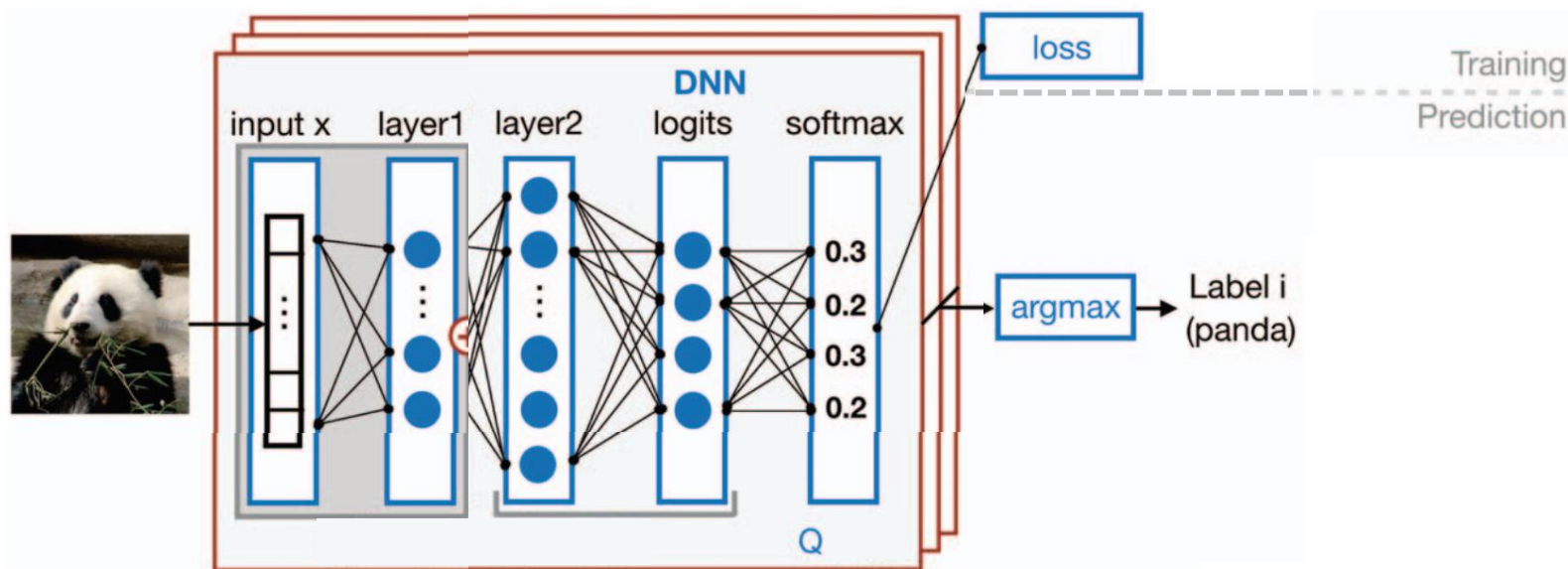
Symbol	Definition
$\mathcal{K} = \{1, \dots, K\}$	The set of all labels
$x = (x_1, \dots, x_n) \in \mathbb{R}^n$	n -dimensional input (n pixels of a image) x_i is the i th pixel in the image
$y = (y_1, \dots, y_K)$	A vector of scores $y_k(x) \in [0,1], \sum_{k=1}^K y_k(x) = 1$
$Q(x) = y = (y_1(x), \dots, y_K(x))$	Scoring function
$f(x) = \arg \max_{k \in \mathcal{K}} y_k(x)$	Prediction procedure
$\alpha = (\alpha_1, \dots, \alpha_n)$	Perturbation (or called the change in the input) α_i is the change to the i th pixel in the image
$x + \alpha$	Adversarial example
$\ \alpha\ _p = \ (\alpha_1, \dots, \alpha_n)\ _p$	p -norm of the perturbation (change) for $1 \leq p < \infty$, $\ \alpha\ _p = (\sum_{i=1}^n \alpha_i ^p)^{1/p}$ for $p = \infty$, $\ \alpha\ _p = \max \alpha_i $ for $p = 0$, $\ \alpha\ _p = \{i: \alpha_i \neq 0\} $
$B_p(r) := \{\alpha \in \mathbb{R}^n: \ \alpha\ _p \leq r\}$	p -norm ball of radius r
L	Radius of the $\alpha \in B_p(L)$ where α is such that $f(x + \alpha) \neq f(x)$



PixelDP classifier

● PixelDP DNN

- Deterministic scoring function Q : $x = (x_1, \dots, x_n) \in \mathbb{R}^n \rightarrow y = (y_1, \dots, y_K)$
 $Q(x) = (y_1(x), \dots, y_K(x))$
- The vulnerability of DNN to adversarial example $(x + \alpha)$ stems from the **unbounded sensitivity of Q** with respect to l_p changes in x .





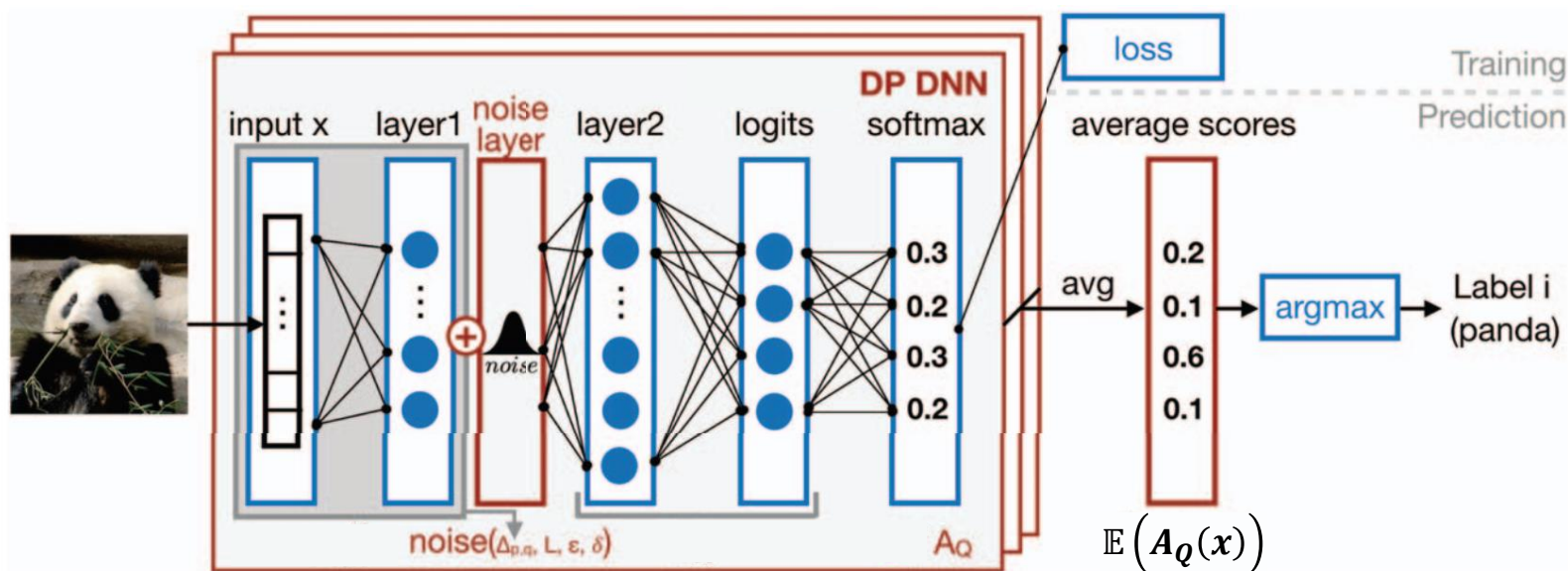
PixelDP classifier

● PixelDP DNN

- Randomized scoring function of the network that satisfy (ϵ, δ) -PixelDP:

$$A_Q(x) = (y_1(x), \dots, y_K(x))$$

- The expected output $\mathbb{E}(A_Q(x))$ of $A_Q(x)$ will have bounded sensitivity to l_p changes in x .

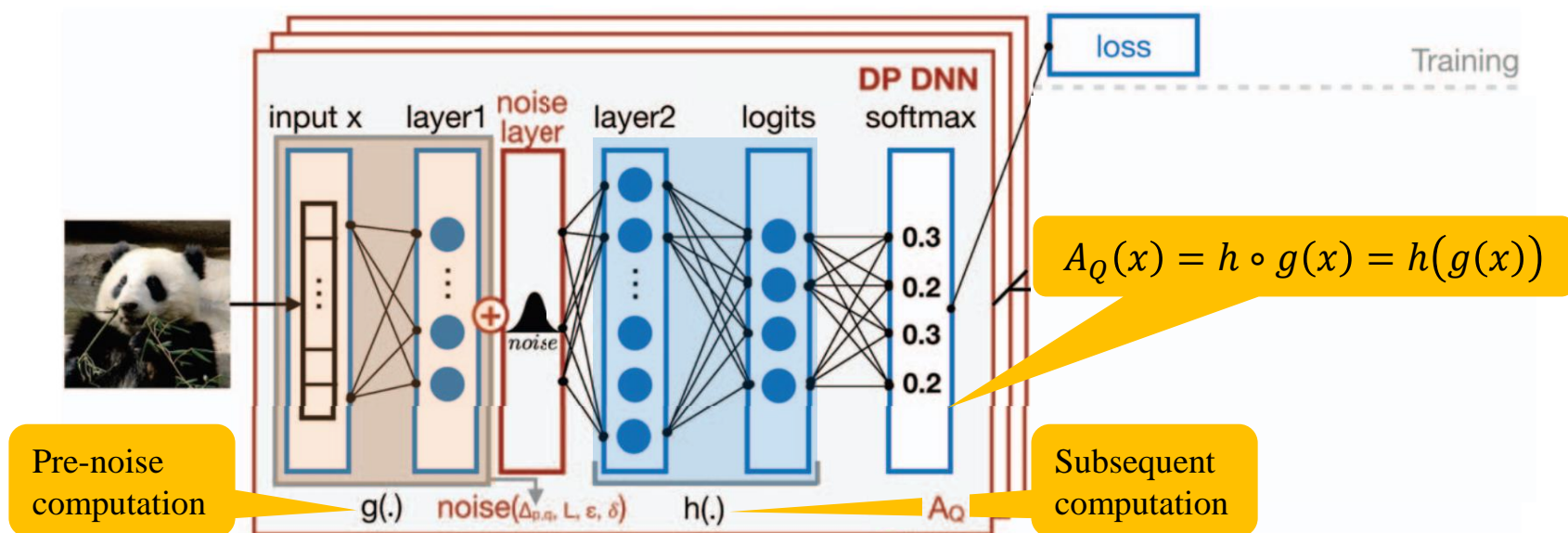




PixelDP classifier

• Training Procedure

- **Step 1** Transform g into another function \tilde{g} that has a **fixed sensitivity** $\Delta \leq 1$ to l_p changes in x .

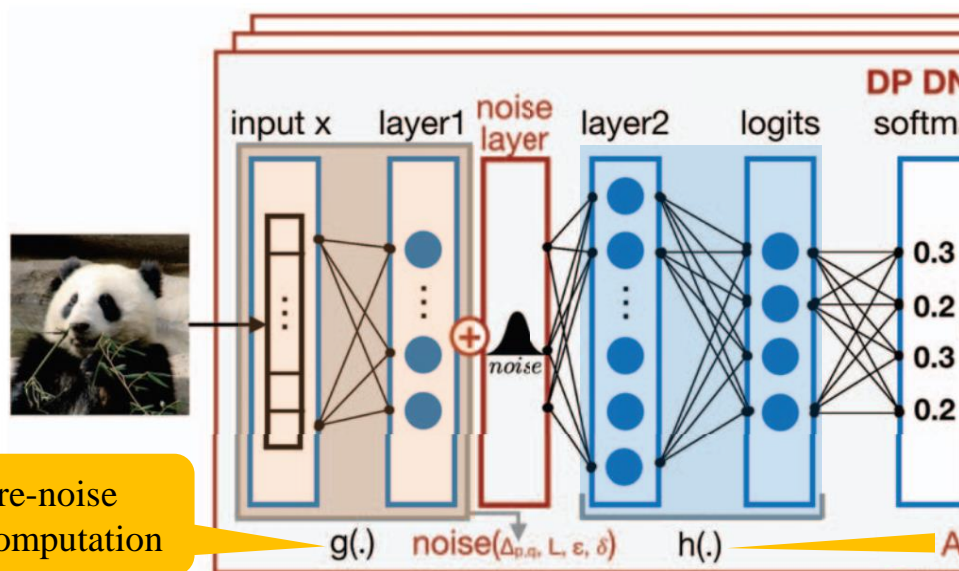




PixelDP classifier

● Training Procedure

- **Step 1** Transform g into another function \tilde{g} that has a **fixed sensitivity** $\Delta \leq 1$ to l_p changes in x .



- ✓ Sensitivity : the maximum change in output that can be produced by a change in the input.

$$\Delta_{p,q} = \max_{Sx, x': x \neq x'} \frac{\|g(x) - g(x')\|_q}{\|x - x'\|_p}$$

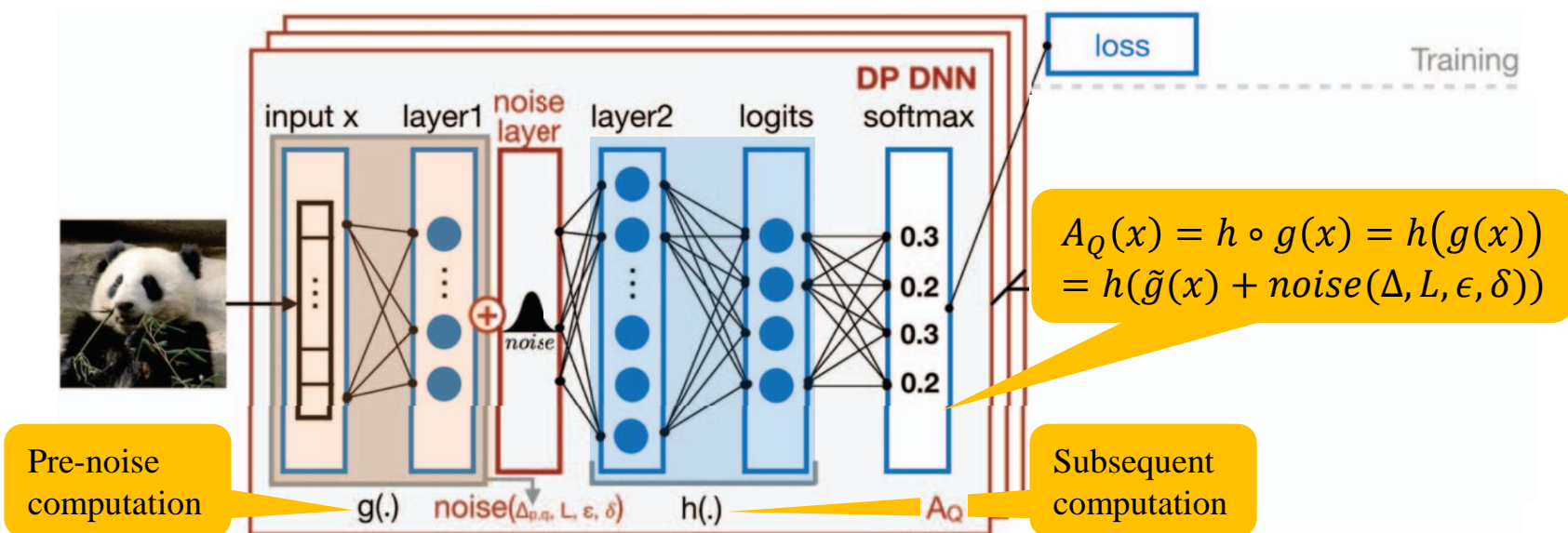
- ✓ Transform Reason:
Training procedure will enlarge the sensitivity $\Delta_{p,q}$ of the g , voiding the DP guarantees.
- ✓ Transform Purpose:
Keep g 's sensitivity $\Delta_{p,q}$ constant (eg. $\Delta \leq 1$) during training.
- ✓ Transform Ways:
 - For $\Delta_{1,1}, \Delta_{1,2}, \Delta_{\infty, \infty}$: SGD.
 - For $\Delta_{2,2}$: projection after SGD.



PixelDP classifier

● Training Procedure

- **Step 1** Transform g into another function \tilde{g} that has a **fixed sensitivity** $\Delta \leq 1$ to l_p changes in x .
- **Step 2** Add the noise layer to the output of \tilde{g} with a standard deviation scaled by Δ and L to ensure (ϵ, δ) -PixelDP for l_p changes of **size** L .





Architecture

- Sample a noise sample Z from noise distribution $\text{noise}(\Delta, L, \epsilon, \delta)$

➤ Mean: $\mu=0$

➤ Standard deviation: σ (b) is proportional to L and $\Delta_{p,q}$.

➤ If Gaussian :

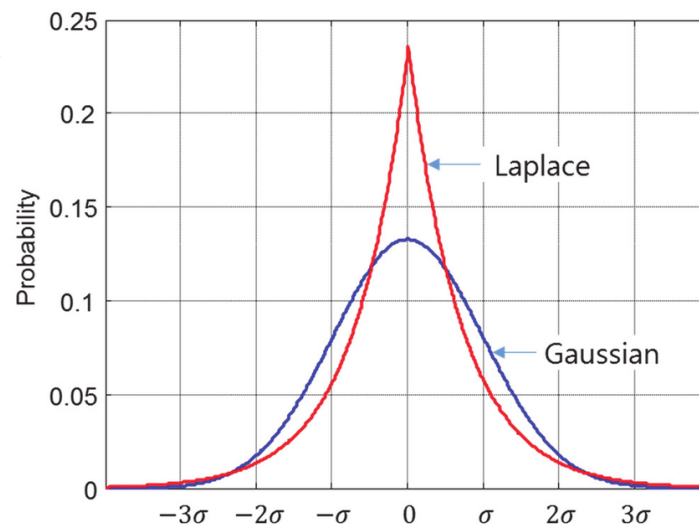
- PDF of $N(\mu, \sigma^2)$: $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$

- $Z \sim N(\mu, \sigma^2) = N(0, \sqrt{2 \ln\left(\frac{1.25}{\delta}\right)} \cdot \Delta_{p,2} \cdot \frac{L}{\epsilon})$

➤ If Laplace :

- PDF of $L(\mu, b)$: $f(x) = \frac{1}{2b} \cdot e^{-\frac{|x-\mu|}{b}}$

- $Z \sim L(\mu, b) = L(0, \sqrt{2} \cdot \Delta_{p,1} \cdot \frac{L}{\epsilon})$

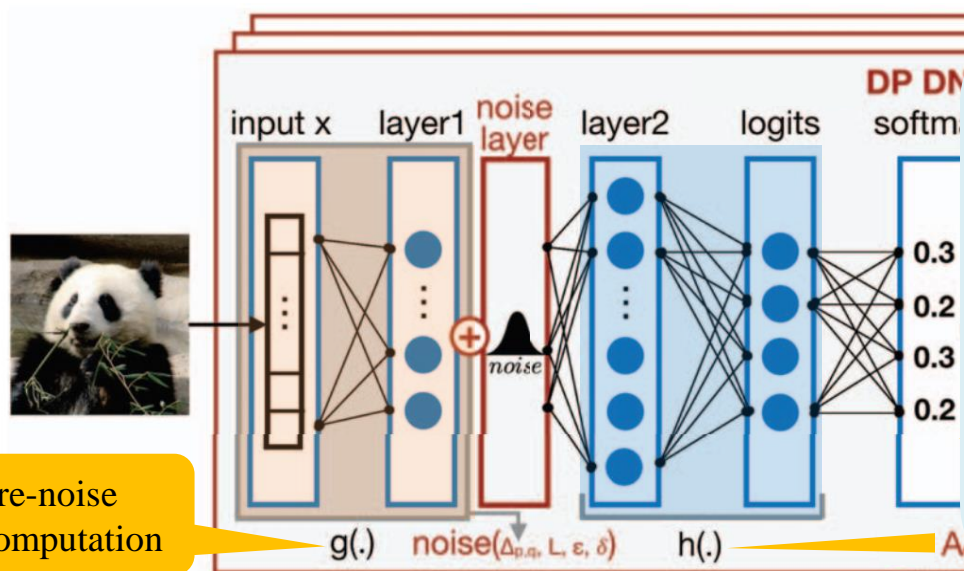




PixelDP classifier

● Training Procedure

- **Step 1** Transform g into another function \tilde{g} that has a **fixed sensitivity** $\Delta \leq 1$ to l_p changes in x .
- **Step 2** Add the noise layer to the output of \tilde{g} with a standard deviation scaled by Δ and L to ensure (ϵ, δ) -PixelDP for l_p changes of **size** L .



- ✓ Set L, ϵ, δ
- ✓ Compute fixed sensitivity Δ
- ✓ Input sample $(x, y_{\text{true}} = c_k)$
- ✓ Sample a noise sample Z from noise distribution $\text{noise}(\Delta, L, \epsilon, \delta)$
 - $\text{Gaussian}(\Delta, L, \epsilon, \delta) \rightarrow \epsilon$ -DP
 - $\text{Laplace}(\Delta, L, \epsilon, \delta) \rightarrow (\epsilon, \delta)$ -DP
- ✓ Optimization

$$\min_{\theta_1, \theta_2} \text{Loss}(h_{\theta_2}(\tilde{g}_{\theta_1}(x) + \text{noise}(\Delta, L, \epsilon, \delta)), y_{\text{true}})$$

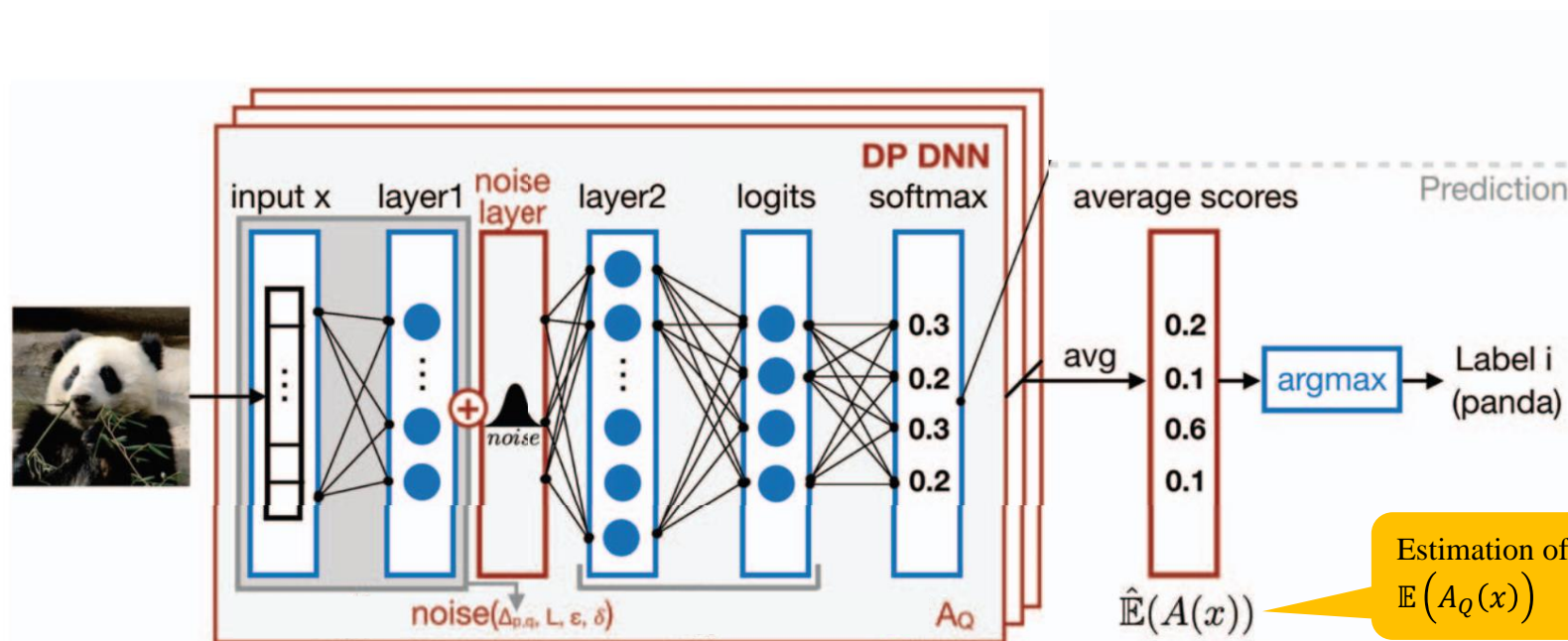


PixelDP classifier

● Prediction Procedure

- ✓ Prediction on the $A_Q(x)$ affords the robustness certification in

Proposition 1 if the prediction procedure uses $\mathbb{E}(A_Q(x))$.





PixelDP classifier

● Proposition 1

If randomized algorithm A satisfies (ϵ, δ) -PixelDP to l_p 1 in x ,

If for some $k \in K$, $\mathbb{E}(A_k(x)) > e^{2\epsilon} \max_{i:i \neq k} \mathbb{E}(A_i(x)) + (1 + e^\epsilon)\delta$

Then, the classifier is robust to any attack $\alpha \in B_p(1)$ on input x .

➤ Proof

1. $\mathbb{E}(A_k(x)) \leq e^\epsilon \mathbb{E}(A_k(x')) + \delta$ (ϵ, δ)-DP定义
2. $\mathbb{E}(A_k(x')) \geq \frac{\mathbb{E}(A_k(x)) - \delta}{e^\epsilon}$ 移项
3. $\mathbb{E}(A_k(x)) > e^{2\epsilon} \max_{i:i \neq k} \mathbb{E}(A_i(x)) + (1 + e^\epsilon)\delta$ 假设条件
4. $\mathbb{E}(A_k(x')) \geq \frac{e^{2\epsilon} \max_{i:i \neq k} \mathbb{E}(A_i(x)) + (1 + e^\epsilon)\delta - \delta}{e^\epsilon} = e^\epsilon \max_{i:i \neq k} \mathbb{E}(A_i(x)) + \delta$ 代入第二行
5. $\mathbb{E}(A_i(x')) \leq e^\epsilon \mathbb{E}(A_i(x)) + \delta, i \neq k$ (ϵ, δ)-DP定义
6. $\max_{i:i \neq k} \mathbb{E}(A_i(x')) \leq e^\epsilon \max_{i:i \neq k} \mathbb{E}(A_i(x)) + \delta, i \neq k$ 两边求最大值
7. $\mathbb{E}(A_k(x')) \geq \max_{i:i \neq k} \mathbb{E}(A_i(x'))$ 代入第四行



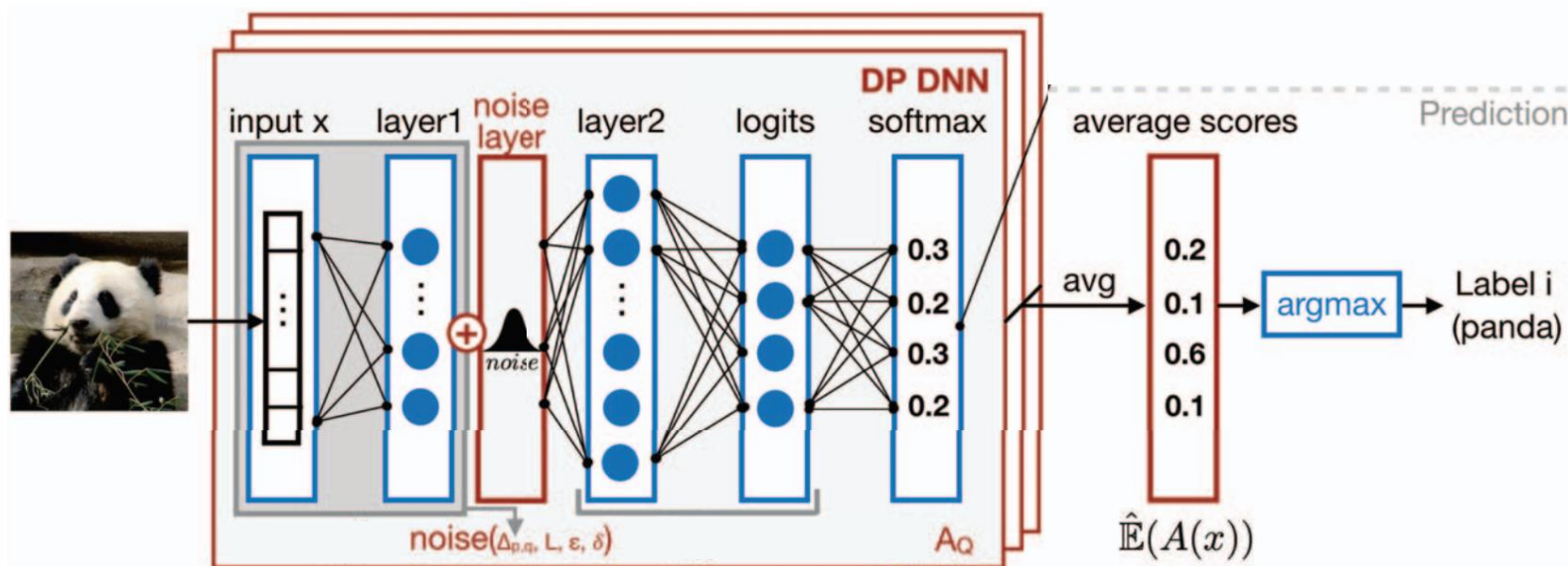
PixelDP classifier

● Prediction Procedure

- ✓ Prediction on the $A_Q(x)$ affords the robustness certification in

Proposition 1 if the prediction procedure uses $\mathbb{E}(A_Q(x))$.

- ✓ Unfortunately, $\mathbb{E}(A_Q(x))$ cannot be computed exactly.

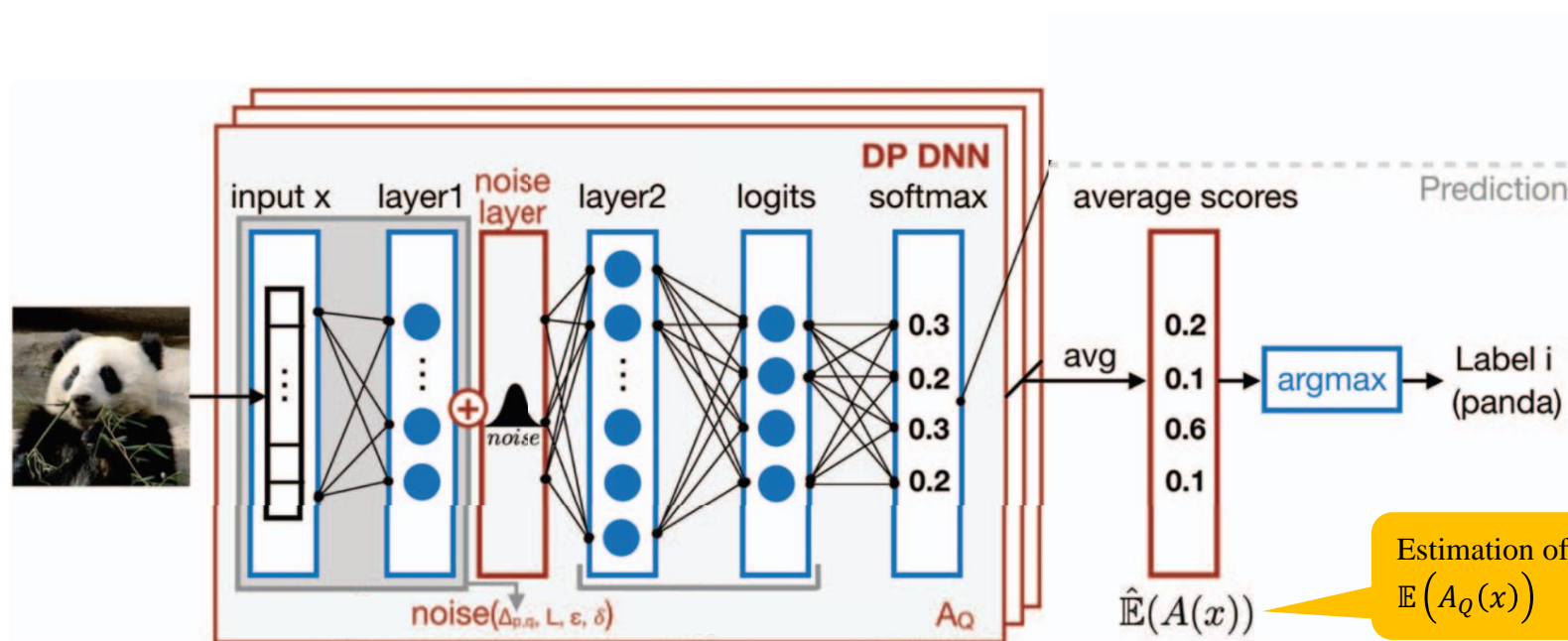




PixelDP classifier

● Prediction Procedure

- **Step 1** Use Monte Carlo methods to estimate approximate value $\hat{\mathbb{E}}(A_Q(x))$ of $\mathbb{E}(A_Q(x))$.

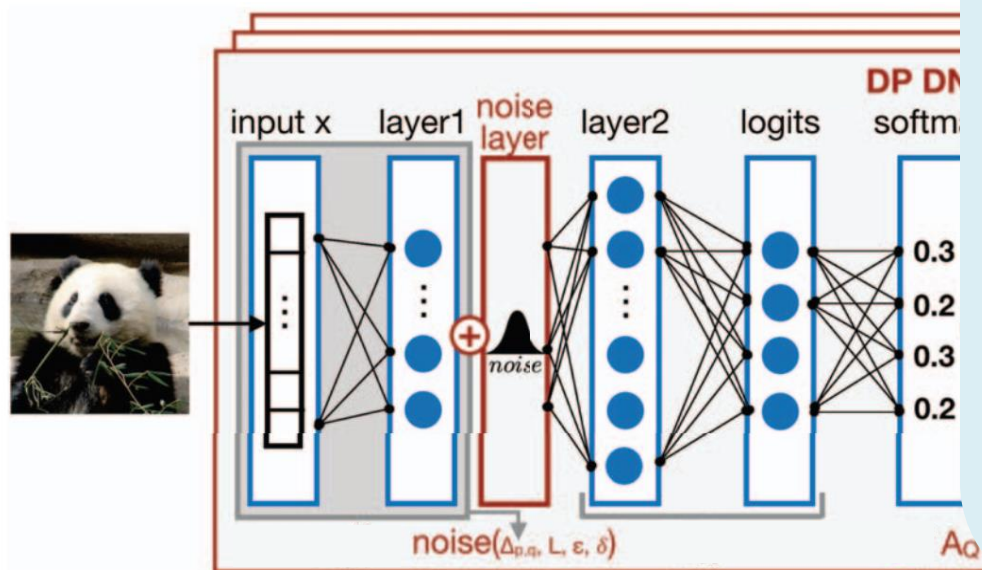




PixelDP classifier

● Prediction Procedure

- **Step 1** Use Monte Carlo methods to estimate approximate value $\hat{\mathbb{E}}(A_Q(x))$ of $\mathbb{E}(A_Q(x))$.



- ✓ Invoke $A_Q(x)$ n times with independent draws in the noise layer.
- ✓ For $i = 1$ to n
 - Input sample x without label
 - Sample i th noise sample from noise distribution $noise^i(\Delta, L, \epsilon, \delta)$
 - Score: $A_Q^i(x) = h(\tilde{g}(x) + noise^i(\Delta, L, \epsilon, \delta))$
 - The i th draw from the distribution of the randomized function A_Q on the k th label: $A_{Q,k}^i(x)$
 - $\hat{\mathbb{E}}(A_{Q,k}(x)) = \frac{1}{n} \sum_{i=1}^n A_{Q,k}^i(x)$



PixelDP classifier

● Proposition 1

If randomized algorithm A satisfies (ϵ, δ) -PixelDP to l_p of size 1 in x ,

If for some $k \in K$, $\mathbb{E}(A_k(x)) > e^{2\epsilon} \max_{i:i \neq k} \mathbb{E}(A_i(x)) + (1 + e^\epsilon)\delta$

Then, the classifier is robust to any attack $\alpha \in B_p(1)$ on input x .

● Proposition 2

If randomized algorithm A satisfies (ϵ, δ) -PixelDP to l_p of size L in x ,

Let $\hat{\mathbb{E}}^{ub} A_i(x)$ and $\hat{\mathbb{E}}^{lb} A_i(x)$ be the η -confidence upper and lower bound for $\mathbb{E}(A_i(x))$.

If for some $k \in K$, $\hat{\mathbb{E}}^{lb}(A_k(x)) > e^{2\epsilon} \max_{i:i \neq k} \hat{\mathbb{E}}^{ub}(A_i(x)) + (1 + e^\epsilon)\delta$

Then, the classifier is robust to any attack $\alpha \in B_p(L)$ on input x with probability $\geq \eta$.



PixelDP classifier

● Proposition 2

If for some $k \in K$, $\hat{\mathbb{E}}^{lb}(A_k(x)) > e^{2\epsilon} \max_{i:i \neq k} \hat{\mathbb{E}}^{ub}(A_i(x)) + (1 + e^\epsilon)\delta$

Then, the classifier is robust to $\alpha \in B_p(L)$ on x with probability $\geq \eta$.

➤ Proof

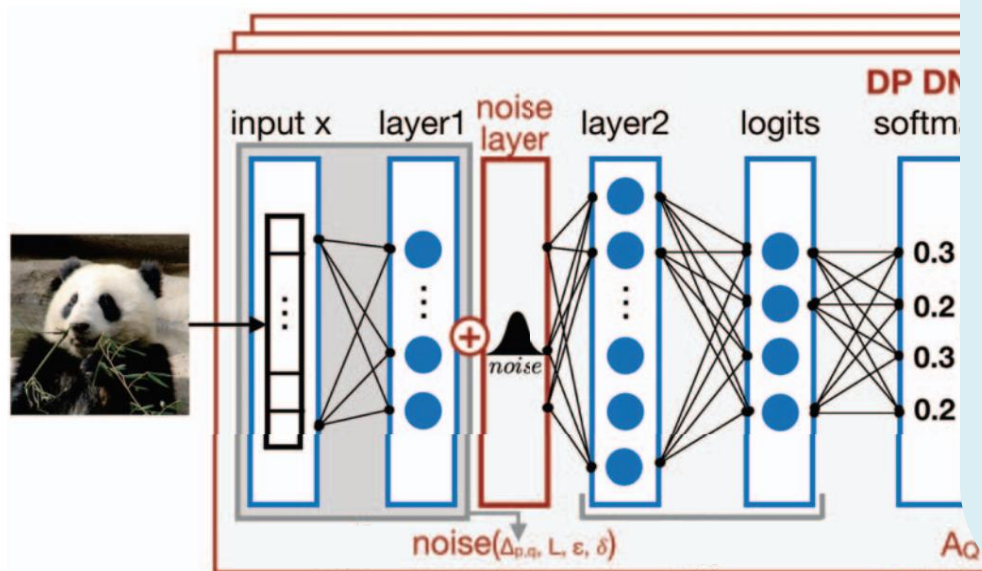
1. $\hat{\mathbb{E}}(A_k(x)) \leq e^\epsilon \hat{\mathbb{E}}(A_k(x')) + \delta$ (ϵ, δ)-DP定义
2. $\hat{\mathbb{E}}(A_k(x')) \geq \frac{\hat{\mathbb{E}}(A_k(x)) - \delta}{e^\epsilon} \geq \frac{\hat{\mathbb{E}}^{lb}(A_k(x)) - \delta}{e^\epsilon}$ 移项，取下届
3. $\hat{\mathbb{E}}^{lb}(A_k(x)) > e^{2\epsilon} \max_{i:i \neq k} \hat{\mathbb{E}}^{ub}(A_i(x)) + (1 + e^\epsilon)\delta$ 假设条件
4. $\hat{\mathbb{E}}(A_k(x')) \geq \frac{\hat{\mathbb{E}}^{lb}(A_k(x)) - \delta}{e^\epsilon} > \frac{e^{2\epsilon} \max_{i:i \neq k} \hat{\mathbb{E}}^{ub}(A_i(x)) + (1 + e^\epsilon)\delta - \delta}{e^\epsilon} =$ 代入第二行
 $e^\epsilon \max_{i:i \neq k} \hat{\mathbb{E}}^{ub}(A_i(x)) + \delta$
5. $\hat{\mathbb{E}}(A_{i:i \neq k}(x')) \leq e^\epsilon \max_{i:i \neq k} \hat{\mathbb{E}}(A_i(x)) + \delta$ (ϵ, δ)-DP定义
6. $\max_{i:i \neq k} \hat{\mathbb{E}}(A_{i:i \neq k}(x')) \leq e^\epsilon \max_{i:i \neq k} \hat{\mathbb{E}}(A_i(x)) + \delta \leq e^\epsilon \max_{i:i \neq k} \hat{\mathbb{E}}^{ub}(A_i(x)) + \delta$ 最大值
7. $\hat{\mathbb{E}}(A_k(x')) > \max_{i:i \neq k} \hat{\mathbb{E}}(A_{i:i \neq k}(x'))$ 代入第四行



PixelDP classifier

● Prediction Procedure

- **Step 1** Use Monte Carlo methods to estimate approximate value $\hat{\mathbb{E}}(A_Q(x))$ of $\mathbb{E}(A_Q(x))$.



- ✓ Compute a interval $[\hat{\mathbb{E}}^{lb}(A_Q(x)), \hat{\mathbb{E}}^{ub}(A_Q(x))]$ for $\hat{\mathbb{E}}(A_Q(x))$ holds with probability η .
- Use Hoeffding's inequality (霍夫丁不等式) to bound error in $\hat{\mathbb{E}}(A_Q(x))$.

$$\hat{\mathbb{E}}(A_Q(x)) - \sqrt{\frac{1}{2n} \ln\left(\frac{2k}{1-\eta}\right)} \leq \mathbb{E}(A_Q(x)) \leq \hat{\mathbb{E}}(A_Q(x)) + \sqrt{\frac{1}{2n} \ln\left(\frac{2k}{1-\eta}\right)}$$

$$\hat{\mathbb{E}}^{lb}(A_Q(x)) \triangleq \hat{\mathbb{E}}(A_Q(x)) - \sqrt{\frac{1}{2n} \ln\left(\frac{2k}{1-\eta}\right)}$$

$$\hat{\mathbb{E}}^{ub}(A_Q(x)) \triangleq \hat{\mathbb{E}}(A_Q(x)) + \sqrt{\frac{1}{2n} \ln\left(\frac{2k}{1-\eta}\right)}$$

- ✓ Integrate this interval into the stability bound for $\mathbb{E}(A_Q(x))$.

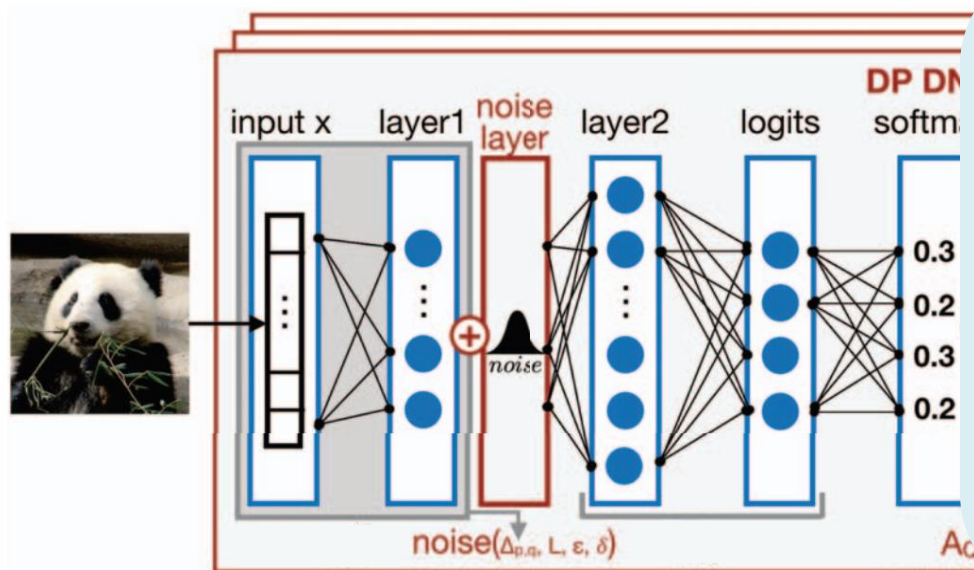
$$\mathbb{E}(A_Q(x)) \leq e^\epsilon \mathbb{E}(A_Q(x + \alpha)) + b\delta$$



PixelDP classifier

● Prediction Procedure

- **Step 1** Use Monte Carlo methods to estimate approximate value $\hat{\mathbb{E}}(A_Q(x))$ of $\mathbb{E}(A_Q(x))$.
- **Step 2** PixelDP returns a prediction for x ($\arg \max \hat{\mathbb{E}}(A_Q(x))$) and a *robustness size certificate* for that prediction.

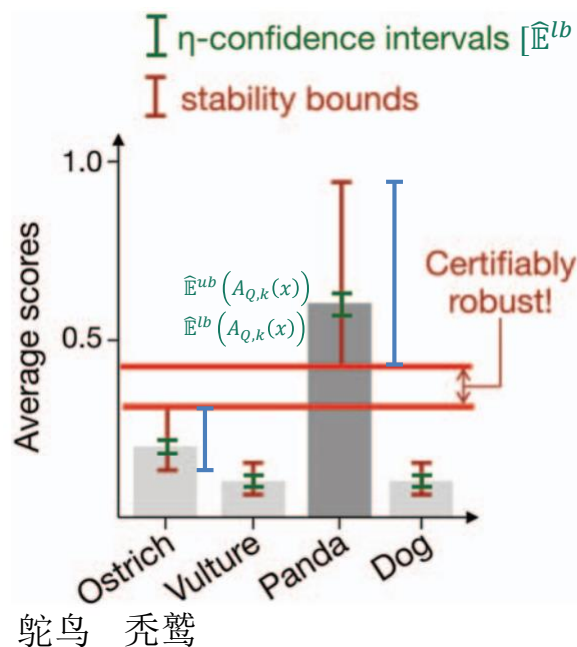


- ✓ Obtain upper and lower bounds on the change to $\hat{\mathbb{E}}(A_{Q,i}(x))$ with input change of size L with probability η .
- ✓ Compute *robustness size certificate*:
 - If,
$$\hat{\mathbb{E}}^{lb}(A_{Q,k}(x)) > e^{2\epsilon} \max_{i:i \neq k} \hat{\mathbb{E}}^{ub}(A_{Q,i}(x)) + (1 + e^\epsilon)\delta$$
 - Then, the classifier is robust to any $\alpha \in B_p(L)$ around x with probability $\geq \eta$.
 - Else, x not meet robustness check for L .



PixelDP classifier

- $\hat{\mathbb{E}}^{up}(A_{Q,i}(x))$ and $\hat{\mathbb{E}}^{lp}(A_{Q,i}(x))$
- If the **lower bound** for the label with the top average score $\hat{\mathbb{E}}^{lb}(A_{Q,k}(x))$ is strictly greater than the **upper bound** for every other label $\max_{i:i \neq k} \hat{\mathbb{E}}^{ub}(A_{Q,i}(x))$,
- Then, with probability η , the prediction for input x is robust to arbitrary attacks of l_p size L .



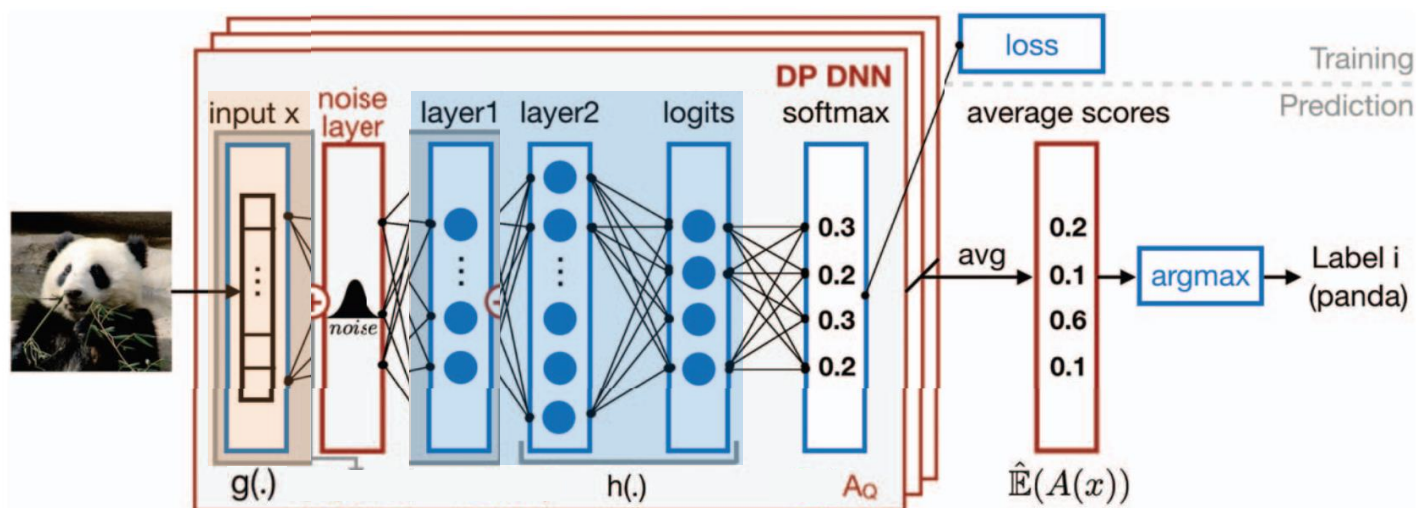


PixelDP classifier

● Noise layer place

➤ Noise in the image

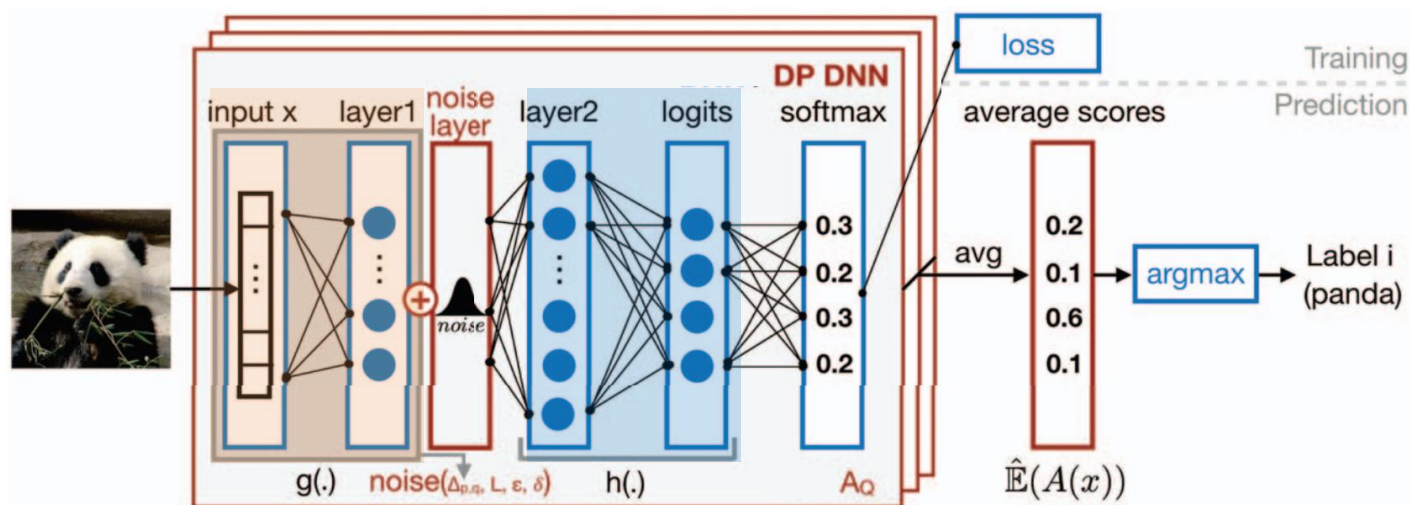
$$\begin{aligned} \bullet \Delta_{p,q} = \Delta_{1,1}^g &= \max_{x,x':x \neq x'} \frac{\|g(x) - g(x')\|_q}{\|x - x'\|_p} = \max_{x,x':x \neq x'} \frac{\|g(x) - g(x')\|_1}{\|x - x'\|_1} \\ \bullet \Delta_{p,q} = \Delta_{2,2}^g &= \max_{x,x':x \neq x'} \frac{\|g(x) - g(x')\|_q}{\|x - x'\|_p} = \max_{x,x':x \neq x'} \frac{\|g(x) - g(x')\|_2}{\|x - x'\|_2} \end{aligned}$$





PixelDP classifier

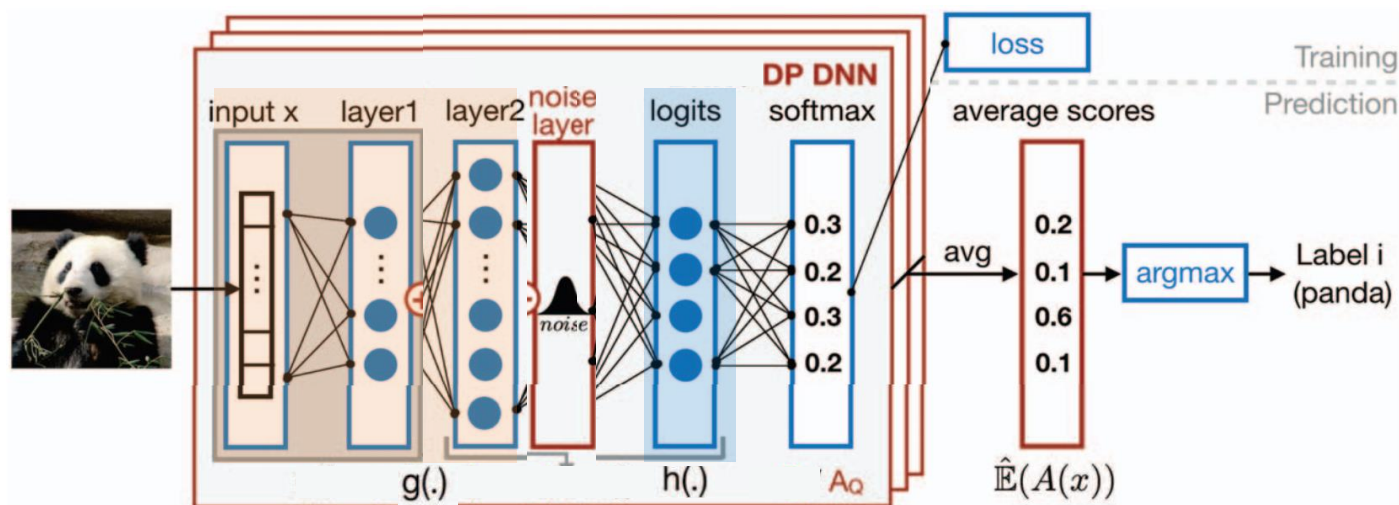
- Noise layer place
 - Noise in the image
 - Noise after first hidden layer
- $g(x) = f_1(x)$
- $\Delta_{p,q} = \Delta_{p,q}^g = \Delta_{p,q}^{f_1}$





PixelDP classifier

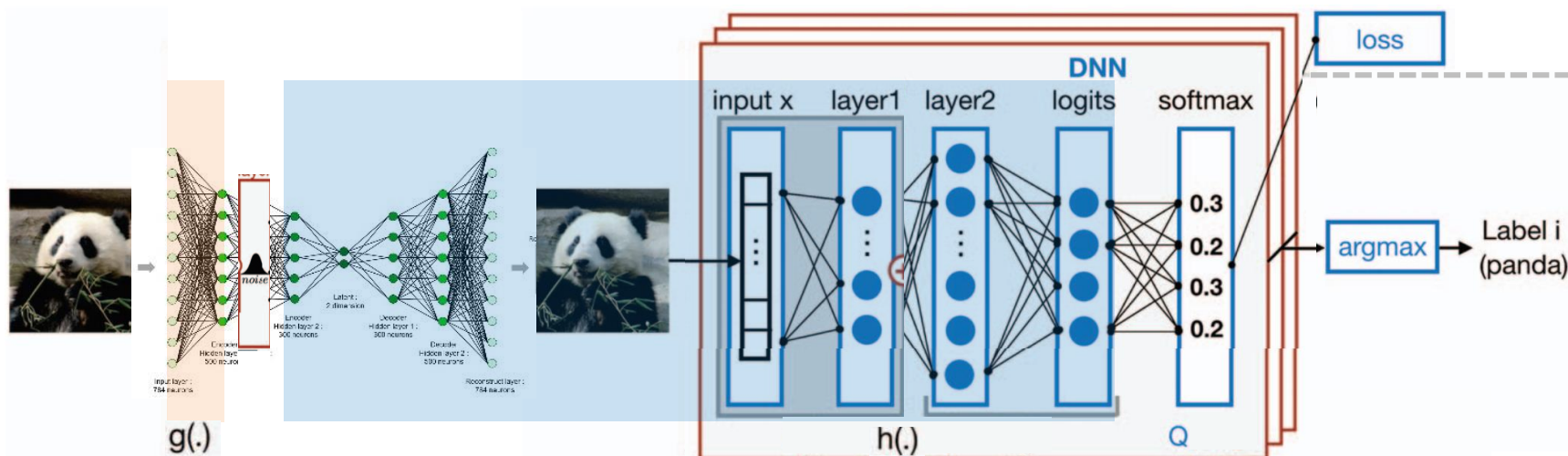
- Noise layer place
 - Noise in the image
 - Noise after first hidden layer
 - Noise after deeper hidden layer
- $g(x) = f_1(f_2(x))$
- $\Delta_{p,q} = \Delta_{p,q}^g = \Delta_{p,q}^{(f_1 \circ f_2)}$





PixelDP classifier

- Noise layer place
 - Noise in the image
 - Noise after first hidden layer
 - Noise after deeper hidden layer
 - Noise in Auto-encoder
- Auto-encoders are smaller than DNN, much faster to train. (ImageNet)





Contribution

- Establish a connection between adversarial robustness and differential privacy.
- Develop the first certified defense that scales to large networks (Google's Inception network) and datasets (ImageNet).
- Develop the certified defense that applies broadly to arbitrary model types.
- Datasets:
 - ImageNet,
 - CIFAR10, CIFAR100, SVHN
 - MNIST
- Networks:
 - Inception
 - ResNet



Authors



Mathias Lecuyer
PhD
Columbia University



Vaggelis Atlidakis
PhD
Columbia University



Roxana Geambasu
Associate Professor
Columbia University



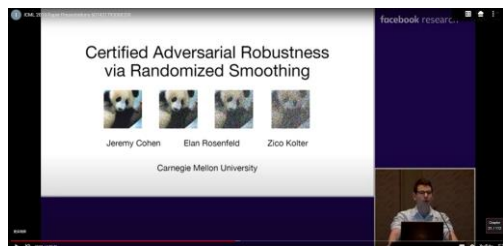
Daniel Hsu
associate professor
Columbia University



Suman Jana
associate professor
Columbia University



Certified Adversarial Robustness via Randomized Smoothing

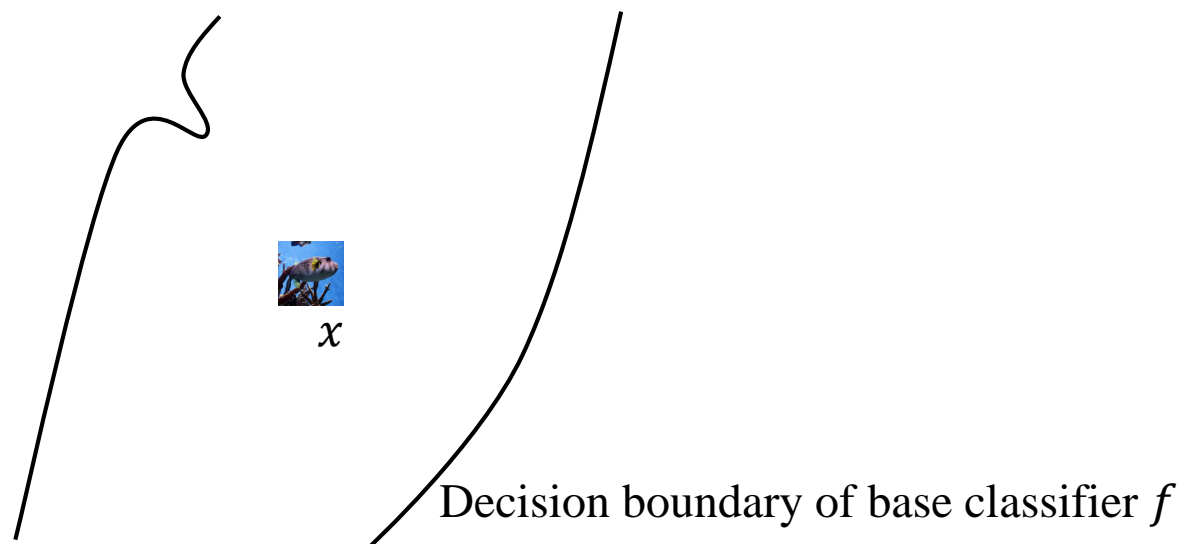


Jeremy Cohen, Elan Rosenfeld, [J. Zico Kolter](#)
Carnegie Mellon University
ICML 2019



Overview

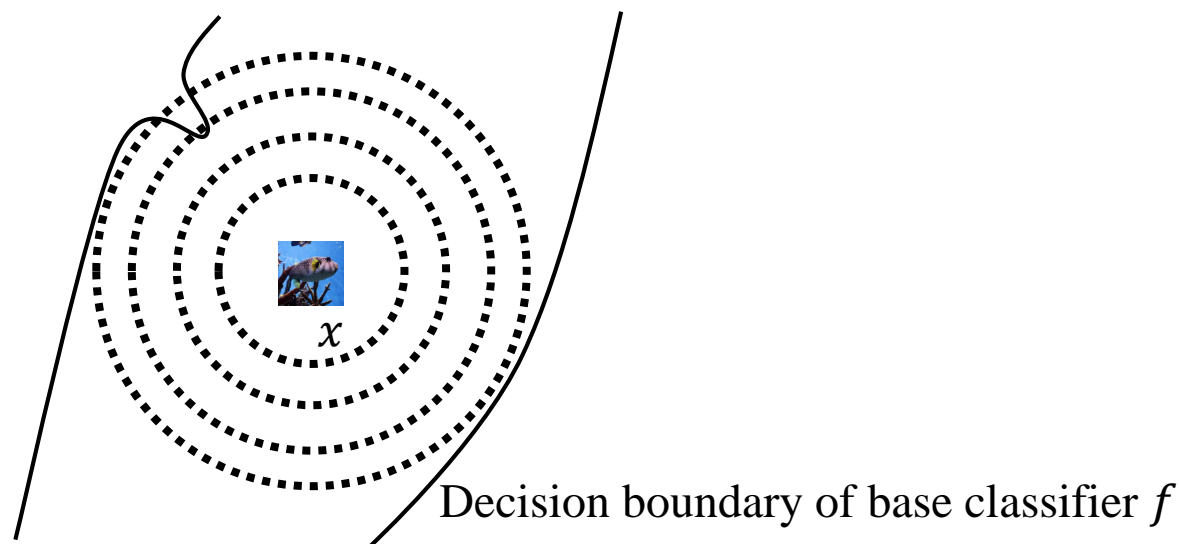
- Applying Gaussian noise and taking majority class label.
- Essence: Smoothing the decision boundary of the base classifier f .
- Way: Sampling Gaussian noise to perturb x multiple times, then vote on the labels most frequently given by the base classifier f .
- Purpose: Output the majority class label as the prediction of x .





Overview

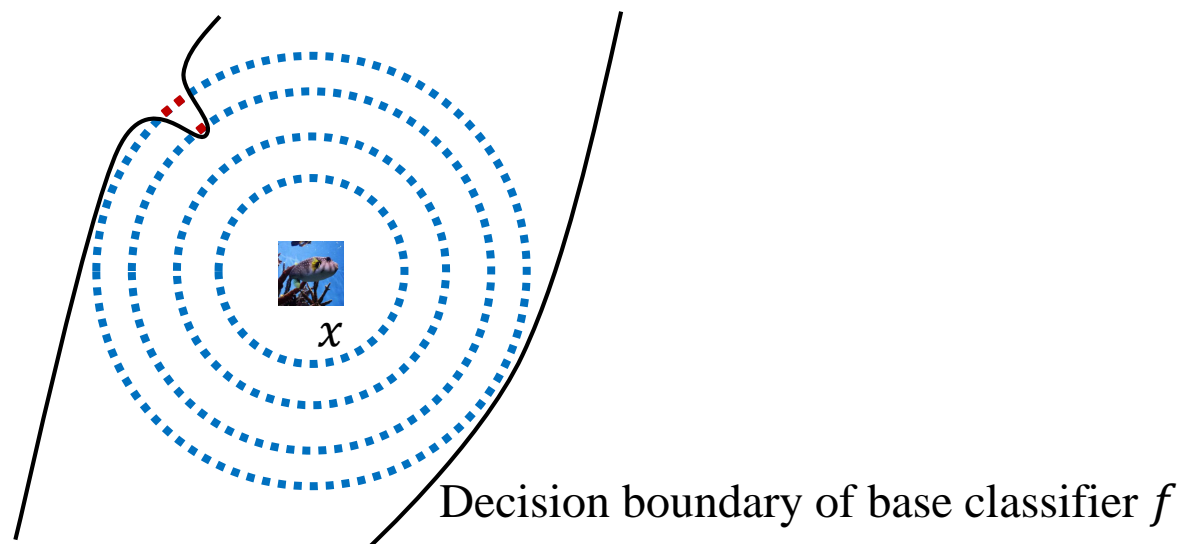
- Applying Gaussian noise and taking majority class label.
- Essence: Smoothing the decision boundary of the base classifier f .
- Way: Sampling Gaussian noise to perturb x multiple times, then vote on the labels most frequently given by the base classifier f .
- Purpose: Output the majority class label as the prediction of x .





Overview

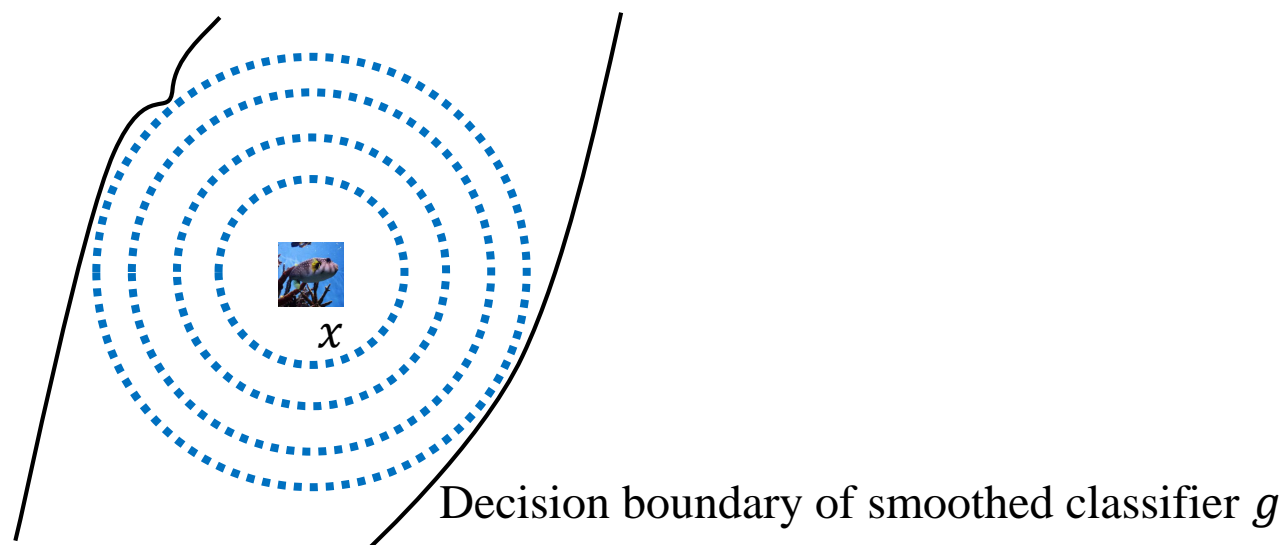
- Applying Gaussian noise and taking majority class label.
- Essence: Smoothing the decision boundary of the base classifier f .
- Way: Sampling Gaussian noise to perturb x multiple times, then vote on the labels most frequently given by the base classifier f .
- Purpose: Output the majority class label as the prediction of x .





Overview

- Applying Gaussian noise and taking majority class label.
- Essence: Smoothing the decision boundary of the base classifier f .
- Way: Sampling Gaussian noise to perturb x multiple times, then vote on the labels most frequently given by the base classifier f .
- Purpose: Output the majority class label as the prediction of x .





Overview

- Evaluating the smoothed classifier at an input x .
- Smoothed classifier g : a virtual classifier with smoother boundary than f
- Here, $g(x) = C_A$.

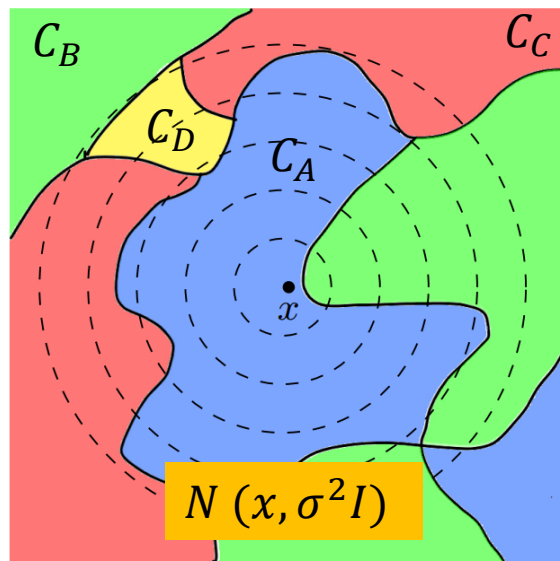


Figure 1

Left:

- Colors: decision regions of f .
- Dotted lines: level sets of $N(x, \sigma^2 I)$.



Problem

- Randomized smoothing-based heuristic defense

- Method

- Insert noise layer before each convolution layer. (Liu. ECCV 2018)
- Ensemble information in a region centered at x to predict. (Cao. ACSAC2017)

- Problem: did not prove any guarantees.

- Randomized smoothing-based certified defense

- Methods

- Smooth the classifier with noise and use DP inequalities to prove robustness guarantee. (Lecuyer. IEEE S&P 2019)
- Use Renyi divergence to prove a stronger guarantee. (Li. arXiv2018)

- Problem: existing robustness guarantees are **loose**.

Liu, Xuanqing. Towards robust neural networks via random self-ensemble. ECCV 2018.

Cao, Xiaoyu. Mitigating evasion attacks to deep neural networks via region-based classification. ACSAC2017.

Lecuyer, Mathias. Certified robustness to adversarial examples with differential privacy. IEEE S&P 2019.

Li, Bai, et al. "Second-order adversarial attack and certifiable robustness." arXiv2018.



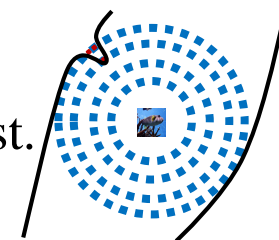
Problem

- Advantage of randomized smoothing

- Makes no assumptions about classifier architecture.
- Permits to use arbitrarily large neural networks. (others do not support)
- The only certified defense been shown feasible on ImageNet task. (before this work)

- Challenge of randomized smoothing

- Not possible to exactly compute the probabilities with which f classifies $\mathcal{N}(x, \sigma^2 I)$ as each class.
- Not possible to exactly evaluate the smoothed classifier g
- Not possible to exactly compute the radius in which g is robust.



- Solution

- Use Monte Carlo method to evaluate prediction of the classifier around x .



Notation

Symbol	Definition
\mathcal{Y}	Classes set
g	Random smoothed classifier
$f: \mathbb{R}^d \rightarrow \mathcal{Y}$	Base classifier
$x \in \mathbb{R}^d$	Input space
$\varepsilon \sim N(0, \sigma^2 I)$	Isotropic Gaussian noise
$C_A \in \mathcal{Y}$	most probable class returned by $f(x + \varepsilon)$
$p_A \in [0, 1]$	$\mathbb{P}(f(x + \varepsilon) = C_A)$.
C_B	“runner-up” class returned by $f(x + \varepsilon)$
$p_B \in [0, 1]$	$\mathbb{P}(f(x + \varepsilon) = C_B)$.
$\underline{p}_A \in [0, 1]$	lower bound of p_A
$\overline{p}_B \in [0, 1]$	upper bound of p_B
$\delta = (\delta_1, \delta_2, \dots, \delta_n)$	n -dimensional perturbation
$\ \delta\ _2$	2 norm of the vector δ , $\sqrt{(\delta_1^2 + \delta_2^2 + \dots + \delta_n^2)}$



Algorithm

- Smoothed Classifier g

- Definition:

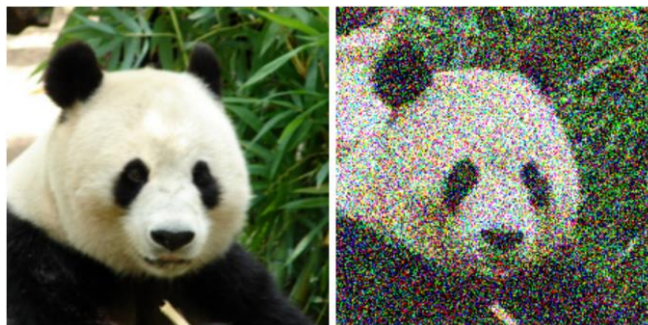
Smoothed classifier g returns whichever class the **base classifier f** is most likely to return when x is perturbed by isotropic Gaussian noise:

$$g(x) = \arg \max_{c \in \mathcal{Y}} \mathbb{P}(f(x + \varepsilon) = c) \quad (1)$$

Where $\varepsilon \sim N(0, \sigma^2)$

- Interpretation:

Randomized smoothing in high dimension is that these large **random perturbations ε** drown out small **adversarial perturbations δ** .



random Gaussian corruptions of x
($\sigma = 0.5$)



Algorithm

- Smoothed Classifier g

- Definition:

Smoothed classifier g returns whichever class the **base classifier f** is most likely to return when x is perturbed by isotropic Gaussian noise:

$$g(x) = \arg \max_{c \in \mathcal{Y}} \mathbb{P}(f(x + \varepsilon) = c) \quad (1)$$

Where $\varepsilon \sim N(0, \sigma^2)$

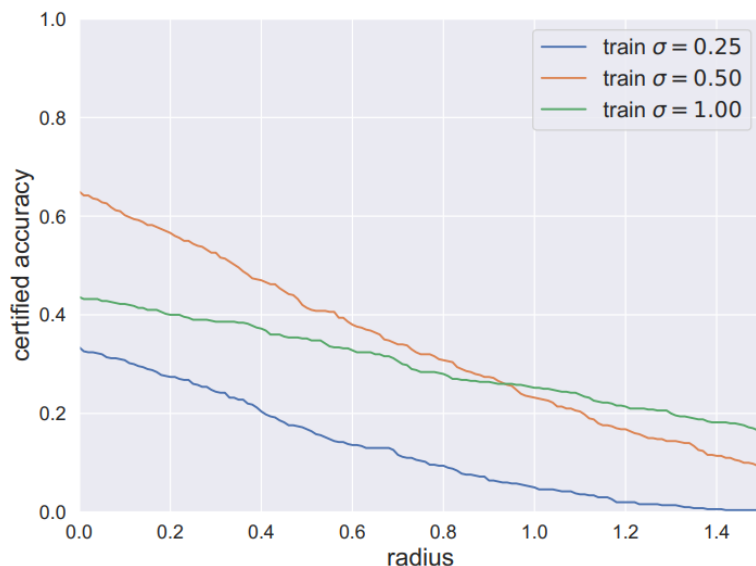
- Base classifier f

- In order for g to classify the labeled example (x, c) correctly and robustly, f needs to consistently classify $N(x, \sigma^2 I)$ as c at training.
- Train f with Gaussian data augmentation at variance σ^2 , with training noise level $\sigma_{train} \geq \sigma_{pred}$ prediction noise level .

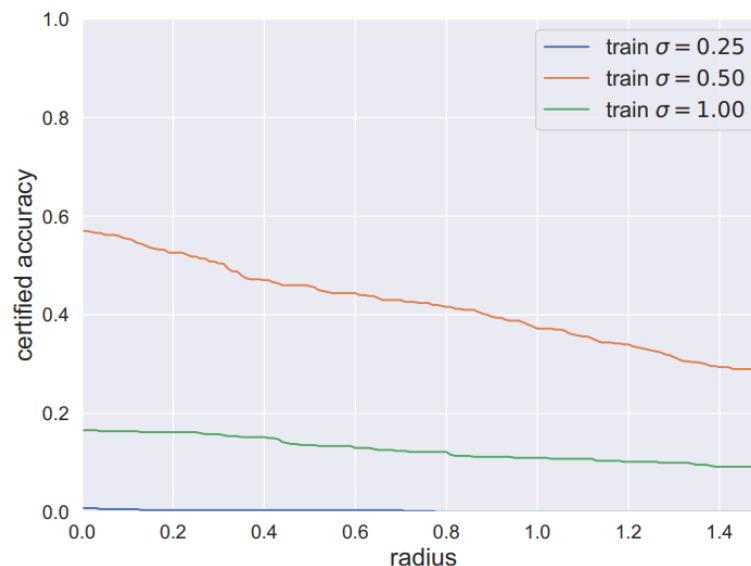


Algorithm

- How much noise to use when training the base classifier ?
- Holding prediction noise level fixed at $\sigma_{pred} = 0.5$.
- If f was trained with a different noise level ($\sigma_{train} \neq 0.5$), g has lower certified accuracy. (blue and green)
- It is better to train with $\sigma_{train} > \sigma_{pred}$ than to train with $\sigma_{train} < \sigma_{pred}$.



(a) CIFAR-10



(b) ImageNet

Figure 15



Algorithm

● PREDICT(f, σ, x, n, α)

evaluate g at x

function PREDICT(f, σ, x, n, α)

- 1 counts \leftarrow SAMPLEUNDERNOISE(f, x, n, σ)
 $\hat{c}_A, \hat{c}_B \leftarrow$ top two indices in counts
 $n_A, n_B \leftarrow$ counts[\hat{c}_A], counts[\hat{c}_B]
- 2 **if** BINOMPVALUE($n_A, n_A + n_B, 0.5$) $\leq \alpha$ **return** \hat{c}_A
else return ABSTAIN

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

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^2)
- 3 $\underline{p}_A \leftarrow$ 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

Independent algorithms for the two tasks:

(A) Evaluating $g(x)$

(B) Evaluating and certifying $g(x)$



Algorithm

- **PREDICT** (f, σ, x, n, α)

evaluate g at x

function PREDICT(f, σ, x, n, α)

- 1 counts \leftarrow SAMPLEUNDERNOISE(f, x, n, σ)
 $\hat{c}_A, \hat{c}_B \leftarrow$ top two indices in counts
 $n_A, n_B \leftarrow$ counts[\hat{c}_A], counts[\hat{c}_B]
- 2 **if** BINOMPVALUE($n_A, n_A + n_B, 0.5$) $\leq \alpha$ **return** \hat{c}_A
else return ABSTAIN

➤ Requirement:

Only need to identify the class c_A with maximal weight in $f(x + \varepsilon)$



Algorithm

● PREDICT (f, σ, x, n, α)

evaluate g at x

function PREDICT(f, σ, x, n, α)

- 1 counts \leftarrow SAMPLEUNDERNOISE(f, x, n, σ)
 $\hat{c}_A, \hat{c}_B \leftarrow$ top two indices in counts
 $n_A, n_B \leftarrow$ counts[\hat{c}_A], counts[\hat{c}_B]
- 2 **if** BINOMPVALUE($n_A, n_A + n_B, 0.5$) $\leq \alpha$ **return**
else return ABSTAIN

SampleUnderNoise (f, x, n, σ)

- Draw n samples of noise:
 $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n \sim N(0; \sigma^2 I)$.
- Run noisy images through the base classifier f to obtain the predictions:
 $f(x + \varepsilon_1), f(x + \varepsilon_2), \dots, f(x + \varepsilon_n)$

Return the *counts* for each class, where the *count* for class c is defined as:

$$\sum_{i=1}^n 1[\text{if } f(x + \varepsilon_i) = c]$$
$$n_A = \sum_{i=1}^n 1[\text{if } f(x + \varepsilon_i) = c_A]$$
$$n_B = \sum_{i=1}^n 1[\text{if } f(x + \varepsilon_i) = c_B]$$



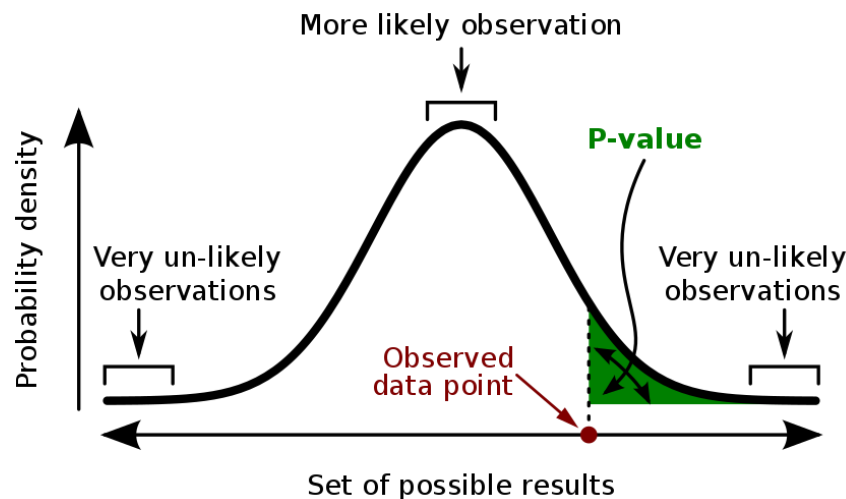
Algorithm

• PREDICT (f, σ, x, n, α)

evaluate g at x

function PREDICT(f, σ, x, n, α)

- 1 counts \leftarrow SAMPLEUNDERNOISE(f, x, n, σ)
 $\hat{c}_A, \hat{c}_B \leftarrow$ top two indices in counts
 $n_A, n_B \leftarrow$ counts[\hat{c}_A], counts[\hat{c}_B]
- 2 **if** BINOMPVALUE($n_A, n_A + n_B, 0.5$) $\leq \alpha$ **return** \hat{c}_A
else return ABSTAIN



BinomPValue ($n_A, n_A + n_B, p = 0.5$)

- Return the **p-value** of the **two-sided hypothesis test** that $n_A \sim \text{Binomial}(n_A + n_B, p = 0.5)$.
- **P-value:** probability of obtaining very unlikely observations when assuming null hypothesis is correct.
- Two-sided hypothesis test:
 $P\text{-value} = 2 * P\text{-value}_{\text{right}}$
- When α is small, abstains frequently, but rarely returns wrong class.
- When α is large, seldom abstains, but often return the wrong class.



Algorithm

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

➤ Requirement:

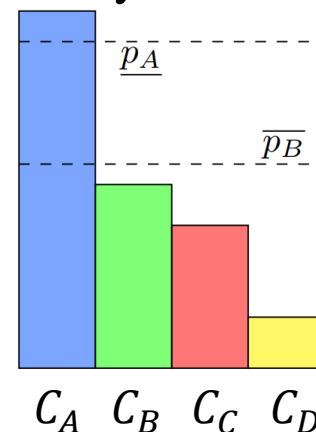
- ✓ Identify the class c_A with maximal weight in $f(x + \varepsilon)$
- ✓ Estimate a lower bound \underline{p}_A on $p_A := P(f(x + \varepsilon) = c_A)$
- ✓ Estimate an upper bound \overline{p}_B on $p_B := \max_{c \neq c_A} P(f(x + \varepsilon) = c)$

➤ Problem

Statistically speaking, estimating \underline{p}_A and \overline{p}_B while simultaneously identifying the top class c_A is a little bit tricky.

➤ Solution

Two-step procedure





Algorithm

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

➤ Requirement:

- ✓ Identify the class c_A with maximal weight in $f(x + \varepsilon)$
- ✓ Estimate a lower bound \underline{p}_A on $p_A := P(f(x + \varepsilon) = c_A)$
- ✓ Estimate an upper bound \overline{p}_B on $p_B := \max_{c \neq c_A} P(f(x + \varepsilon) = c)$

➤ Procedure:

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^2)

③ $\underline{p}_A \leftarrow$ 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



Algorithm

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

➤ Requirement:

- ✓ Identify the class c_A with maximal weight in $f(x + \varepsilon)$
- ✓ Estimate a lower bound \underline{p}_A on $p_A := P(f(x + \varepsilon) = c_A)$
- ✓ Estimate an upper bound \overline{p}_B on $p_B := \max_{c \neq c_A}$

➤ Procedure:

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^2)

③ $\underline{p}_A \leftarrow$ 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

➤ Use a small number n_0 of samples from $f(x + \varepsilon)$ to take a guess at c_A , because $f(x + \varepsilon)$ tends to put most of its weight on the top class,



Algorithm

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

➤ Requirement:

✓ Identify the class c_A with maximal weight in f

✓ Estimate a lower bound \underline{p}_A on $p_A := P(f(x) = c_A)$

✓ Estimate an upper bound \overline{p}_B on $p_B := \max_{c \neq c_A} P(f(x) = c)$

➤ Procedure:

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^2)

③ $\underline{p}_A \leftarrow$ 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

➤ Use a larger number n of samples to estimate \underline{p}_A and \overline{p}_B

✓ Use LowerConfBound($k, n, 1 - \alpha$) to return \underline{p}_A of $[\underline{p}_A, \overline{p}_A]$ that holds with probability at least $1 - \alpha$ over the sampling of $k \sim \text{Binomial}(n, p_A)$.

■ Clopper-Person confidence interval (二项式比例置信区间): 根据一系列伯努利成功(c_A)-失败(c_B)实验的结果计算出的成功(c_A)概率(p_A)的置信区间 $[\underline{p}_A, \overline{p}_A]$

➤ Take $\underline{p}_B = 1 - \underline{p}_A$.



Algorithm

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

➤ Requirement:

✓ Identify the class c_A with maximal weight in \mathcal{C}

✓ Estimate a lower bound \underline{p}_A on $p_A := P(f(x) = c_A)$

✓ Estimate an upper bound \overline{p}_B on $p_B := \max_{c \neq c_A} P(f(x) = c)$

➤ Procedure:

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^2)

③ $\underline{p}_A \leftarrow$ 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

➤ If $\underline{p}_A > 0.5$ (即 $\overline{p}_B < 0.5$)

• Compute robustness guarantee

$$\begin{aligned} R &= \frac{\sigma}{2} \left(\phi^{-1}(\underline{p}_A) - \phi^{-1}(\overline{p}_B) \right) \\ &= \frac{\sigma}{2} \left(\phi^{-1}(\underline{p}_A) - \phi^{-1}(1 - \underline{p}_A) \right) \\ &= \frac{\sigma}{2} \cdot 2\phi^{-1}(\underline{p}_A) = \sigma \phi^{-1}(\underline{p}_A) \end{aligned}$$

• Return \hat{c}_A and R



Radius

- **Theorem 1.** Let $f : \mathbb{R}^d \rightarrow \mathcal{Y}$ be any deterministic or random function, and let $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$. Let g be defined as in (1). Suppose $c_A \in \mathcal{Y}$ and $\underline{p}_A, \overline{p}_B \in [0, 1]$ satisfy:

$$\mathbb{P}(f(x + \varepsilon) = c_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{c \neq c_A} \mathbb{P}(f(x + \varepsilon) = c) \quad (2)$$

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

$$R = \frac{\sigma}{2} (\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B)) \quad (3)$$

- 证明

- 要证

$$g(x + \delta) = \arg \max_{c \in \mathcal{Y}} \mathbb{P}(f(x + \varepsilon + \delta) = c) = c_A$$

- 就要证

$$\mathbb{P}(f(x + \varepsilon + \delta) = c_A) > \max_{c \neq c_A} \mathbb{P}(f(x + \varepsilon + \delta) = c)$$

- 对任意 $c_B \neq c_A$, 证 $\mathbb{P}(f(x + \varepsilon + \delta) = c_A) > \mathbb{P}(f(x + \varepsilon + \delta) = c_B)$



Radius

● Illustration of the proof of Theorem 1.

➤ 当且仅当 $\mathbb{P}(f(Y \in A)) > \mathbb{P}(f(Y \in B))$ 成立, 有 $\mathbb{P}(f(x + \varepsilon + \delta) = c_A) > \mathbb{P}(f(x + \varepsilon + \delta) = c_B)$

- 定义: $A := \{z: \delta^T(z - x) \leq \sigma \|\delta\| \phi^{-1}(\underline{p}_A)\}$ $B := \{z: \delta^T(z - x) \leq \sigma \|\delta\| \phi^{-1}(\overline{p}_B)\}$
- $\mathbb{P}(f(Y \in A)) = \phi(\phi^{-1}(\underline{p}_A) - \frac{\|\delta\|}{\sigma})$ $\mathbb{P}(f(Y \in B)) = \phi(\phi^{-1}(\overline{p}_B) + \frac{\|\delta\|}{\sigma})$

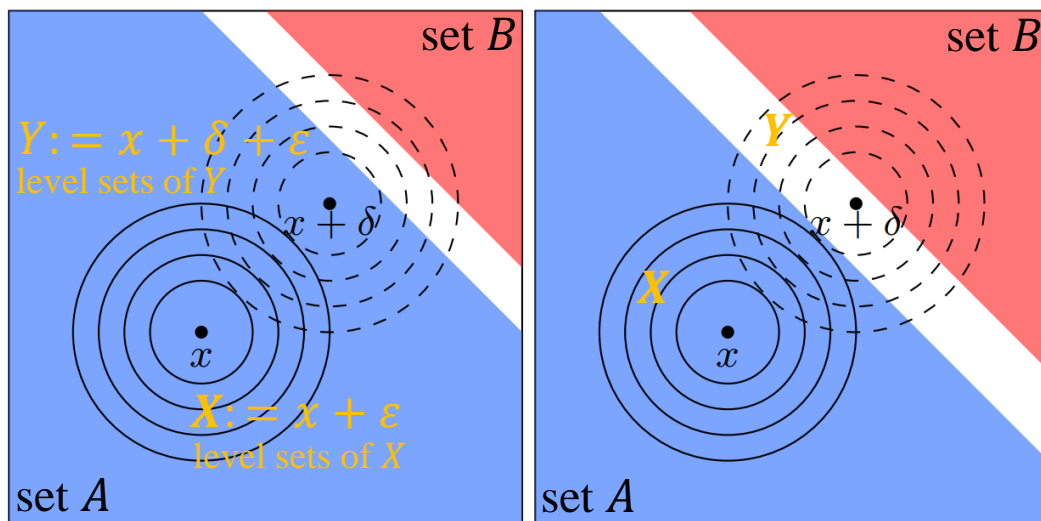


Figure 9

- The figure on the **left** depicts a situation where $\mathbb{P}(Y \in A) > \mathbb{P}(Y \in B)$, and hence $g(x + \delta) = g(Y)$ may equal c_A .
- The figure on the **right** depicts a situation where where $\mathbb{P}(Y \in A) < \mathbb{P}(Y \in B)$ and hence $g(x + \delta) = g(Y) \neq c_A$.



Evaluation

- Metric of normal classification

- Standard test set accuracy:

$$STD.ACC = \frac{N_{g-Correctly}}{N_{g-Total}}$$

- Metric of certified classification

- Certified test set accuracy at radius r :

$$CERT.ACC = \frac{N_{g-Correctly-Robust}}{N_{g-Total}}$$

- g classifies correctly (without abstaining) and certifies robust with a radius $R \geq r$.
- r is similar to Threshold size T defined in Lecuyer. IEEE S&P 2019.



Evaluation

- Experiments with randomized smoothing on ImageNet with $\sigma = 0.25$.
- Left: certified accuracies obtained using CERTIFY VS. those obtained using the guarantees derived in prior works (SP2019).
- Middle: certified accuracy if the **number of samples n** used by CERTIFY had been larger or smaller.
- Right: certified accuracy as the **failure probability α** of CERTIFY is varied.

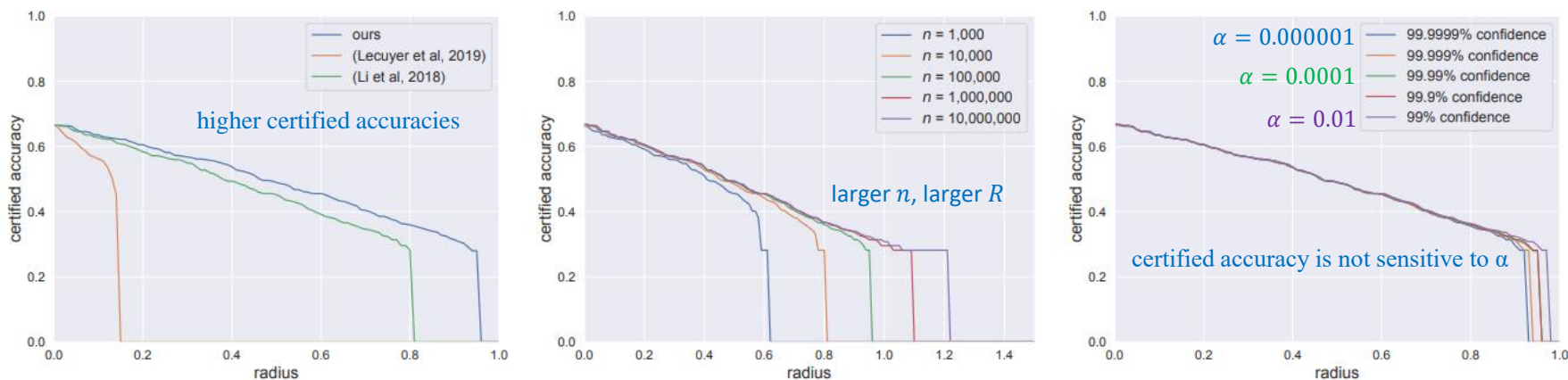


Figure 8



Contribution

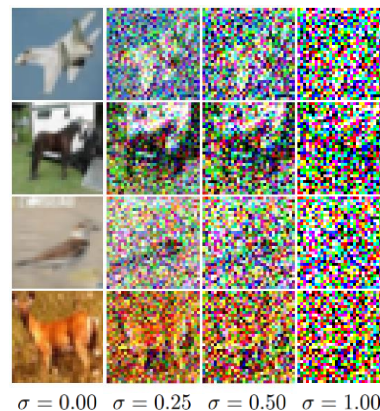
- Prove the first tight robustness guarantee for randomized smoothing.
- Analysis reveals that smoothing with Gaussian noise naturally induces certifiable robustness under the l_2 norm.
- Suspect that other noise distributions might induce robustness to other perturbation sets such as general l_p norm balls.
- Enables the use of large networks on large scale datasets and does not constrain the architecture of the classifier.

➤ Datasets:

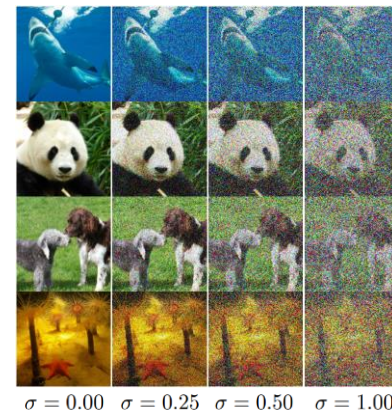
- ImageNet
- CIFAR10

➤ Networks:

- ResNet



CIFAR10



ImageNet



Authors



Jeremy Cohen
PhD at CMU



Elan Rosenfeld
PhD at CMU



J. Zico Kolter
Associate
Professor at CMU



Overview

1

Background

2

Differential Privacy Scheme

3

Randomized Smoothing Scheme

4

Conclusion

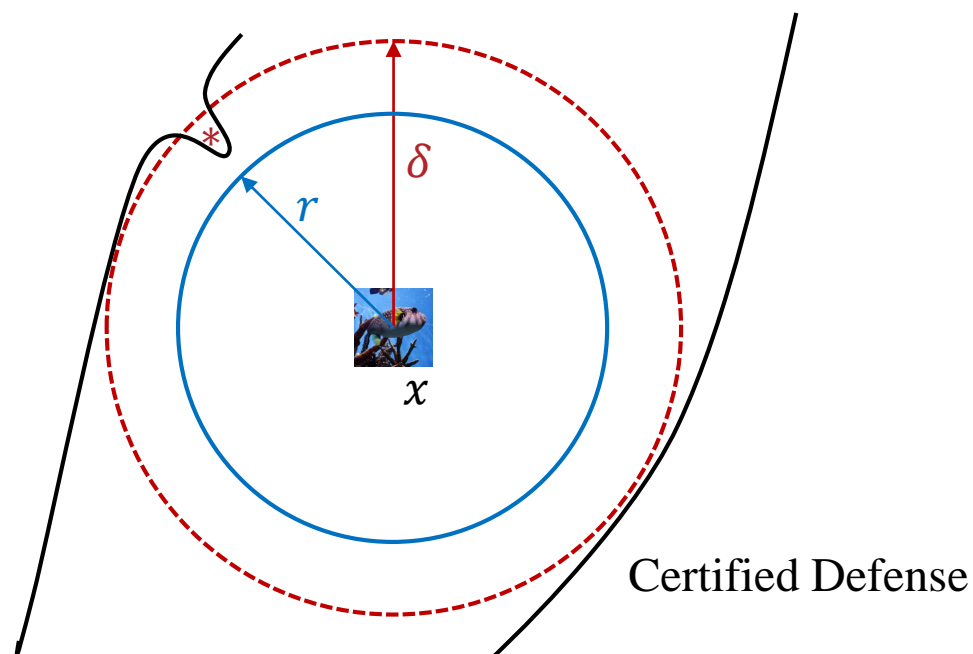


Defense Category

● Certified Defense

- Provide a certificate (r) for adversarial robustness

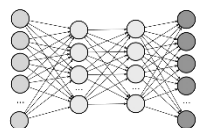
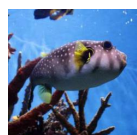
For any input x , the prediction output by the classifier f on samples in l_p ball centered at x are guaranteed to be constant.





Certified Defense Category

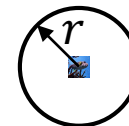
Prediction



Puffer



97.99



certified radii r

Certification

➤ Exact Certified Defense

(x, f, r)



Certification
Procedure



Does not exist δ within $||\delta|| \leq r$ for which $f(x) \neq f(x + \delta)$.



Exists δ within $||\delta|| \leq r$ for which $f(x) \neq f(x + \delta)$.

➤ Conservative Certified Defense

(x, f, r)



Certification
Procedure



Does not exist δ within $||\delta|| \leq r$ for which $f(x) \neq f(x + \delta)$.

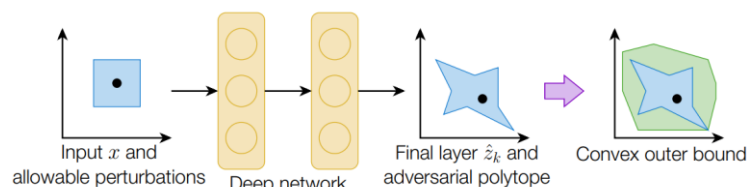


Abstain.



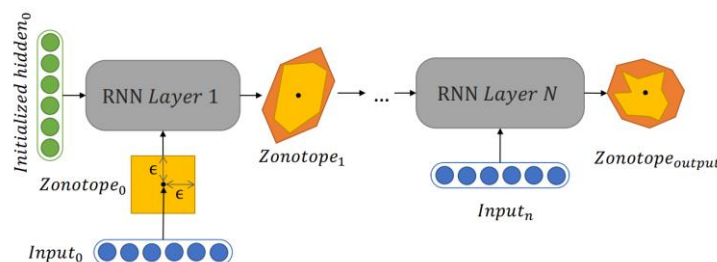
Future

● CNN Certified Defense



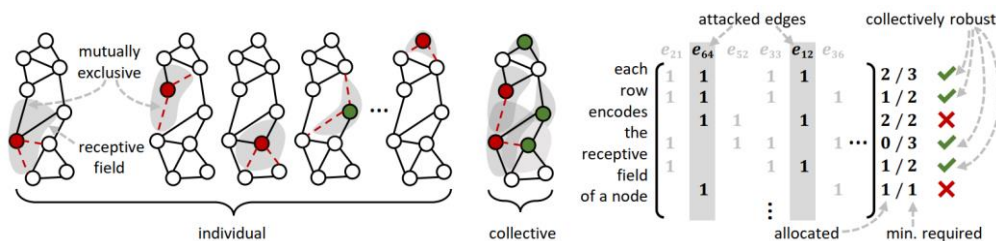
Zico Kolter.
Provable defenses against
adversarial examples via the
convex outer adversarial
polytope.
ICML2018.

● RNN Certified Defense



Du Tianyu.
Cert-RNN: Towards
Certifying the Robustness of
Recurrent Neural Networks.
CCS2021.

● GNN Certified Defense



Schuchardt.
Collective robustness
certificates: Exploiting
interdependence in graph
neural networks.
ICLR2021.



西安电子科技大学
XIDIAN UNIVERSITY

Thank You

Mengdie Huang
AI Security
Lab Ruiyun
2022-04-01

