

《Certified Adversarial Robustness via Randomized Smoothing》-ICML2019

文章阅读

[ICML2019]<https://arxiv.org/pdf/1902.02918.pdf>

Designing a smoothed classifier against adversarial samples using Gaussian noise

1. 背景介绍

本文提出了“randomized smoothing”，并对其进行了理论分析。即对分类器高斯噪声处理，使得新分类器对对抗攻击足够鲁棒。

Certified Defense:

训练一个对输入样本 x 的 ℓ_2 或 ℓ_∞ 邻域内所有样本都具有鲁棒性的分类器。具体的方法分为exact method和 conservative method。

目的在于寻找一个扰动 δ ，满足 $g(x) \neq g(x + \delta)$ ，如果存在这样的扰动，则 decline to make a certification；如果找不到这样的扰动，则假设成立。然而，没有一种 exact method适用于中型复杂度的神经网络。Conservative certification 可扩展到任意大小的神经网络，但是其得到的鲁棒性 guarantee比较 loose。

given a model $h \circ f$, and a new test point (x, y) , we would like to prove $h \circ f(x') = h \circ f(x)$, for all x' in the allowed oerturbation set. That is, to provide certification about the optimality of the following equation

$$\max_{x' \in \Delta(x_{test})} \text{Loss}(h(f(x')), y_{test})$$

Randomized Smoothing:

给定一个neural network f (base classifier), 并且 $f(x) = y$. 将 f 转变成smoothed classifier g .

When queried at x , the smoothed classifier g returns whichever class the base classifier f is most likely to return when x is perturbed by isotropic Gaussian noise:

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

Certification guarantee:

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

where p_A is the prob of the top class, p_B is the prob of the runner-up class. The smoothed classifier g will return the top class, where Φ^{-1} denotes the inverse Gaussian CDF

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

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

Pseudocode for certification and prediction

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

# certify the robustness of g around x
function CERTIFY( $f, \sigma, x, n_0, n, \alpha$ )
  counts0  $\leftarrow$  SAMPLEUNDERNOISE( $f, x, n_0, \sigma$ )
   $\hat{c}_A \leftarrow$  top index in counts0
  counts  $\leftarrow$  SAMPLEUNDERNOISE( $f, x, n, \sigma$ )
   $\underline{p}_A \leftarrow$  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
```

Prediction: 创建 n 个带有高斯噪声的 x 样本, 用base classifier $f(x + \eta)$ 对其分类, 得到top two classes \hat{c}_A 和 \hat{c}_B , If $n_A \gg n_B$, then predict the top class; otherwise, abstain.

Certification: Use small #samples to identify c_A ; Use large #samples to estimate p_A ; Set $p_B = 1 - p_A$ and compute.

证明思路: <https://zhuanlan.zhihu.com/p/463037691>