

Increasing Confidence in Adversarial Robustness Evaluations

Javad Hezareh & Mahdi Saieedi

Security and Privacy in Machine Learning

Sharif University of Technology



Spring 2023

Increasing Confidence in Adversarial Robustness Evaluations

Roland S. Zimmermann*
University of Tübingen

Wieland Brendel
University of Tübingen

Florian Tramèr
Google

Nicholas Carlini
Google

Abstract

Hundreds of defenses have been proposed to make deep neural networks robust against minimal (adversarial) input perturbations. However, only a handful of these defenses held up their claims because correctly evaluating robustness is extremely challenging: Weak attacks often fail to find adversarial examples even if they unknowingly exist, thereby making a vulnerable network look robust.

In this paper, we propose a test to identify weak attacks, and thus weak defense evaluations. Our test slightly modifies a neural network to guarantee the existence of an adversarial example for every sample. Consequentially, any correct attack must succeed in breaking this modified network.

For eleven out of thirteen previously-published defenses, the original evaluation of the defense fails our test, while stronger attacks that break these defenses pass it. We hope that attack *unit tests* — such as ours — will be a major component in future robustness evaluations and increase confidence in an empirical field that is currently riddled with skepticism. Online version & code: zimmerrol.github.io/active-tests/

Figure: Oral at CVPR 2022 Workshop (Art of Robustness), Project [website](#).

Overview

- Can we trust a proposed defense?
- Is defense evaluation valid?

- Can we trust a proposed defense?
- Is defense evaluation valid?
- Main scheme for defense evaluation:
 - 1 Propose an attack
 - 2 Evaluate defense with this attack
 - 3 If no adversarial example is found, defense works

Outline

1 Introduction

2 Background

3 Proposed Active Test

- Classifiers with Linear Classification Readouts
- Tests for Models Leveraging Detectors

4 Evaluation

Outline

1 Introduction

2 Background

3 Proposed Active Test

- Classifiers with Linear Classification Readouts
- Tests for Models Leveraging Detectors

4 Evaluation

Theorem Analogy

Lemma 2.4 (Certificate for Empirical Mean). Let S be an (ϵ, δ) -stable set with respect to a distribution X , for some $\delta \geq \epsilon > 0$. Let T be an ϵ -corrupted version of S . Let μ_T and Σ_T be the empirical mean and covariance of T . If the largest eigenvalue of Σ_T is at most $1 + \lambda$, for some $\lambda \geq 0$, then $\|\mu_T - \mu_X\|_2 \leq O(\delta + \sqrt{\lambda\epsilon})$.

Proof of Lemma 2.4. Let $S' = S \cap T$ and $T' = T \setminus S'$. By replacing S' with a subset if necessary, we may assume that $|S'| = (1 - \epsilon)|S|$ and $|T'| = \epsilon|S|$. Let $\mu_{S'}, \mu_{T'}, \Sigma_{S'}, \Sigma_{T'}$ represent the empirical means and covariance matrices of S' and T' . A simple calculation gives that

$$\Sigma_T = (1 - \epsilon)\Sigma_{S'} + \epsilon\Sigma_{T'} + \epsilon(1 - \epsilon)(\mu_{S'} - \mu_{T'})(\mu_{S'} - \mu_{T'})^T.$$

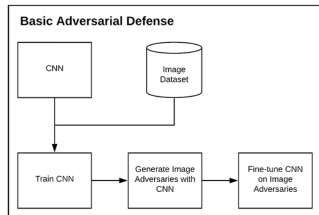
Let v be the unit vector in the direction of $\mu_{S'} - \mu_{T'}$. We have that

$$\begin{aligned} 1 + \lambda &\geq v^T \Sigma_T v = (1 - \epsilon)v^T \Sigma_{S'} v + \epsilon v^T \Sigma_{T'} v + \epsilon(1 - \epsilon)v^T (\mu_{S'} - \mu_{T'})(\mu_{S'} - \mu_{T'})^T v \\ &\geq (1 - \epsilon)(1 - \delta^2/\epsilon) + \epsilon(1 - \epsilon)\|\mu_{S'} - \mu_{T'}\|_2^2 \\ &\geq 1 - O(\delta^2/\epsilon) + (\epsilon/2)\|\mu_{S'} - \mu_{T'}\|_2^2, \end{aligned}$$

where we used the variational characterization of eigenvalues, the fact that $\Sigma_{T'}$ is positive semidefinite, and the second stability condition for S' . By rearranging, we obtain that $\|\mu_{S'} - \mu_{T'}\|_2 = O(\delta/\epsilon + \sqrt{\lambda/\epsilon})$. Therefore, we can write

$$\begin{aligned} \|\mu_T - \mu_X\|_2 &= \|(1 - \epsilon)\mu_{S'} + \epsilon\mu_{T'} - \mu_X\|_2 = \|\mu_{S'} - \mu_X + \epsilon(\mu_{T'} - \mu_{S'})\|_2 \\ &\leq \|\mu_{S'} - \mu_X\|_2 + \epsilon\|\mu_{S'} - \mu_{T'}\|_2 = O(\delta) + \epsilon \cdot O(\delta/\epsilon + \sqrt{\lambda/\epsilon}) \\ &= O(\delta + \sqrt{\lambda\epsilon}), \end{aligned}$$

where we used the first stability condition for S' and our bound on $\|\mu_{S'} - \mu_{T'}\|_2$. \square



Theorem

Proof.

we have proved that $P \neq NP$ □

- How do you refute the proof's claim?

Theorem

Proof.

we have proved that $P \neq NP$ □

- How do you refute the proof's claim?
 - Find an algorithm to solve 3-SAT in polynomial time.
 - Studying proofs line-by-line, till find some major flaw.

Defense X.

We have demonstrated that defense X improves model robustness. One can validate this claim by following the below evaluation procedure.

...



- How do you refute the authors' claim?

Defense X.

We have demonstrated that defense X improves model robustness. One can validate this claim by following the below evaluation procedure.

...



- How do you refute the authors' claim?
 - Find an adversarial attack to decrease model performance.
 - Probe defense evaluation, till find some major flaw.

- Test attack strength

Method

- Test attack strength
- Design a new task that is solvable by any sufficiently strong attack
 - Injects adversarial examples into a defense
 - Check if the attack can find them

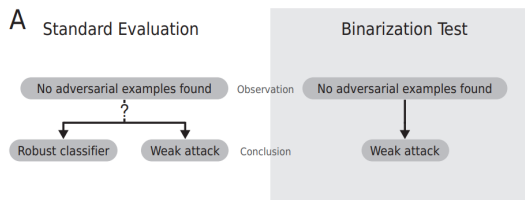


Figure: Proposed method to evaluate the attack used in defense evaluation.

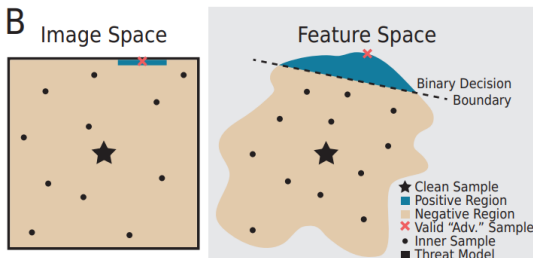


Figure: Inject adversarial examples to check whether the attack is powerful enough.

- A rejected attack doesn't necessarily mean the defense is not effective.

Outline

1 Introduction

2 Background

3 Proposed Active Test

- Classifiers with Linear Classification Readouts
- Tests for Models Leveraging Detectors

4 Evaluation

Adversarial Examples

- Imperceptible perturbations that change the decision of a deep neural network in arbitrary directions

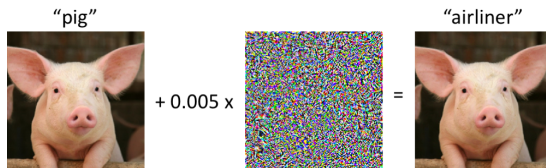


Figure: Adversarial examples.

- Add input pre-processing steps
- Introduce architectural changes
- Methods for detecting adversarial examples

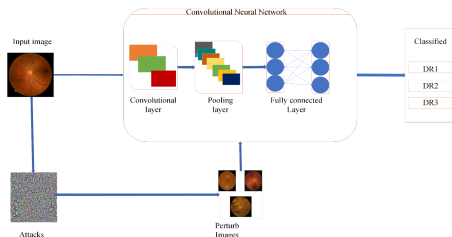


Figure: Adversarial training.

Outline

1 Introduction

2 Background

3 Proposed Active Test

- Classifiers with Linear Classification Readouts
- Tests for Models Leveraging Detectors

4 Evaluation

- Defense evaluation scheme:
 - Perform an attack on defense
 - No adversarial example found within distance $d(\mathbf{x}_c, \mathbf{x}_{adv}) \leq \epsilon$

- Defense evaluation scheme:
 - Perform an attack on defense
 - No adversarial example found within distance $d(\mathbf{x}_c, \mathbf{x}_{adv}) \leq \epsilon$
- Attack strength depends on:
 - Attack itself
 - The defense it is meant to evaluate

Outline

1 Introduction

2 Background

3 Proposed Active Test

- Classifiers with Linear Classification Readouts
- Tests for Models Leveraging Detectors

4 Evaluation

Classification Model

- f : our classifier
- f^* : model feature extractor
- $f = f^* + \text{linear classification head}$

New Task Algorithm

- New binary classification task
- Ensure having adversarial examples

Algorithm Build New Task

Require: feature extractor f^* of original classifier, test sample \mathbf{x}_c , distance ϵ , number of inner/boundary samples N_i and N_b .

```
1: function CREATEBINARYCLASSIFIER( $f^*, \mathbf{x}_c, \epsilon, N_i, N_b$ )
2:    $\mathcal{X}_i := \{\mathbf{x}_c\} \cup \{\hat{\mathbf{x}} \mid d(\mathbf{x}_c, \hat{\mathbf{x}}) \leq \alpha\epsilon \text{ and } \hat{\mathbf{x}} \neq \mathbf{x}_c\}_{1:N_i}$ 
3:    $\mathcal{X}_b := \{\hat{\mathbf{x}} \mid d(\mathbf{x}_c, \hat{\mathbf{x}}) = \epsilon\}_{1:N_b}$ 
4:    $F_i := \{f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_i\}$ 
5:    $F_b := \{f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_b\}$ 
6:    $\mathcal{D} = \{(\hat{\mathbf{x}}, 0) \mid \hat{\mathbf{x}} \in F_i\} \cup \{(\hat{\mathbf{x}}, 1) \mid \hat{\mathbf{x}} \in F_b\}$ 
7:    $g = \text{TRAINLINEAR}(\mathcal{D})$ 
8:    $h = g \circ f^*$ 
9:   return  $h$ 
```

Binarization Test

Algorithm Binarization Test for Classifiers with Linear Classification Readouts

Require: feature extractor f^* of original classifier, test samples \mathcal{X}_{test} , distance ϵ , number of inner/boundary samples N_i and N_b .

```
1: function BINARIZATIONTEST( $f^*$ ,  $\mathcal{X}_{test}$ ,  $\epsilon$ ,  $N_i$ ,  $N_b$ )
2:   attack_successful = []
3:   random_attack_successful = []
4:   for all  $\mathbf{x}_c \in \mathcal{X}_{test}$  do
5:      $h = \text{CREATBINARYCLASSIFIER}(f^*, \mathbf{x}_c, \epsilon, N_i, N_b)$ 
6:     attack_successful.append(ATTACK( $h, \mathbf{x}_c$ ))
7:     random_attack_successful.append(RANDOMATTACK( $h, \mathbf{x}_c$ ))
8:   ASR = MEAN(attack_successful)
9:   RASR = MEAN(random_attack_successful)
10:  return ASR, RASR
```

Evaluate New Task

- The efficacy of the used evaluation method:
 - Use original attack to attack h
 - ASR returns *test score*

Evaluate New Task

- The efficacy of the used evaluation method:
 - Use original attack to attack h
 - ASR returns *test score*
- Test difficulty:
 - Use randomized attack
 - #samples = #queris
 - RASR returns this metric

Outline

1 Introduction

2 Background

3 Proposed Active Test

- Classifiers with Linear Classification Readouts
- Tests for Models Leveraging Detectors

4 Evaluation

Models Leveraging Detectors

- Beside classifier f , we have a detector d
- d detects adversarial examples
- A successful attack must fool both the classifier and the detector

Models Leveraging Detectors

- Beside classifier f , we have a detector d
- d detects adversarial examples
- A successful attack must fool both the classifier and the detector
- Two test for this type:
 - Regular Test
 - Inverted Test
- A reliable evaluation must pass both tests

Algorithm Build New Task for Classifiers with Detectors

Require: feature extractor f^* of original classifier, adversarial detector d , test sample \mathbf{x}_c , distance ϵ , number of inner/boundary/reference samples $N_i/N_b/N_r$.

```
1: function CREATEBINARYCLASSIFIER( $f^*, d, \mathbf{x}_c, \epsilon, N_i, N_b, N_r$ )
2:    $\mathcal{X}_i := \{\mathbf{x}_c\} \cup \{\hat{\mathbf{x}} \mid d(\mathbf{x}_c, \hat{\mathbf{x}}) \leq \alpha\epsilon \text{ and } \hat{\mathbf{x}} \neq \mathbf{x}_c\}_{1:N_i}$ 
3:    $\mathcal{X}_b := \{\hat{\mathbf{x}} \mid d(\mathbf{x}_c, \hat{\mathbf{x}}) = \epsilon, d(\hat{\mathbf{x}}) = 1\}_{1:N_b}$ 
4:    $\mathcal{X}_r := \{\hat{\mathbf{x}} \mid d(\mathbf{x}_c, \hat{\mathbf{x}}) = \eta\epsilon, d(\hat{\mathbf{x}}) = 1\}_{1:N_r}$ 
5:    $F_i := \{f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_i\}$ 
6:    $F_b := \{f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_b\}$ 
7:    $F_r := \{f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_r\}$ 
8:    $\mathcal{D} = \{(\hat{\mathbf{x}}, 0) \mid \hat{\mathbf{x}} \in F_i\} \cup \{(\hat{\mathbf{x}}, 1) \mid \hat{\mathbf{x}} \in F_b\} \cup \{(\hat{\mathbf{x}}, 1) \mid \hat{\mathbf{x}} \in F_r\}$ 
9:    $b = \text{TRAINLINEAR}(\mathcal{D})$ 
10:  return  $b, \mathcal{X}_r$ 
```

Binarization Test For Detectors

Algorithm Binarization Test for Classifiers with Linear Classification Readouts and a Detector

Require: feature extractor f^* of original classifier, adversarial detector d , test samples \mathcal{X}_{test} , distance ϵ , number of inner/boundary/reference samples $N_i/N_b/N_r$.

```
1: function BINARIZATIONTEST( $f^*, d, \mathcal{X}_{test}, \epsilon, N_i, N_b, N_r, \eta$ )
2:   attack_successful = []
3:   random_attack_successful = []
4:   for all  $\mathbf{x}_c \in \mathcal{X}_{test}$  do
5:      $b, \mathcal{X}_r = \text{CREATBINARYCLASSIFIER}(f^*, d, \mathbf{x}_c, \epsilon, N_i, N_b, N_r)$ 
6:     attack_successful.append(ATTACK( $b, d, \mathbf{x}_c, \mathcal{X}_r$ ))
7:     random_attack_successful.append(RANDOMATTACK( $b, d, \mathbf{x}_c, \mathcal{X}_r$ ))
8:   ASR = MEAN(attack_successful)
9:   RASR = MEAN(random_attack_successful)
10:  return ASR, RASR
```

- Why do we need the inverted test?

Algorithm Inverted Test

```
1: function INVERTEDBINARIZATIONTEST( $f^*, d, \mathcal{X}_{test}, \epsilon, N_i, N_b, N_r, \epsilon, \eta$ )  
2:   return BINARIZATIONTEST( $f^*, \neg d, \mathcal{X}_{test}, \epsilon, N_i, N_b, N_r, \epsilon, \eta$ )
```

Outline

1 Introduction

2 Background

3 Proposed Active Test

- Classifiers with Linear Classification Readouts
- Tests for Models Leveraging Detectors

4 Evaluation

Evaluations

- Defenses without Detectors
- Defenses with Detectors

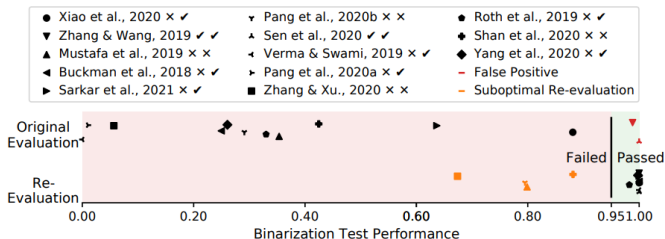


Figure: Binarization Test result for 13 Defenses.

Evaluation

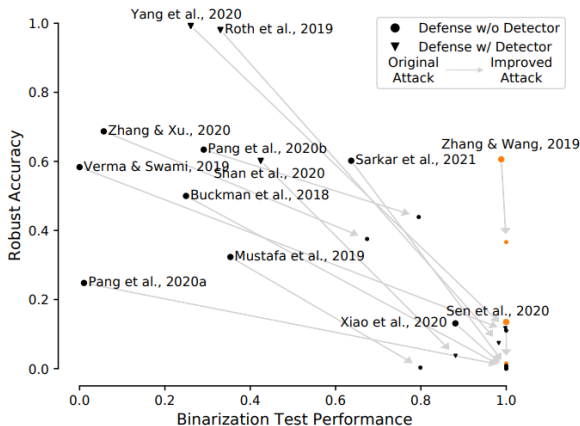


Figure: Robust accuracy as a function of the test performance.

Hardness of Test

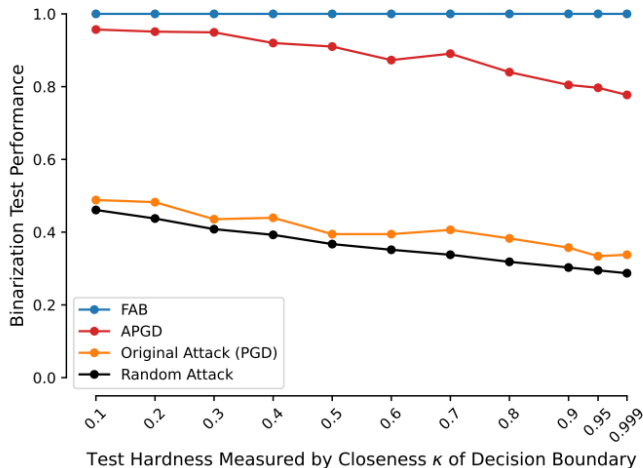


Figure: Hyperparameters influence the test's hardness

Thanks

Any questions?