



**HAL**  
open science

# Robust Neural Networks using Randomized Adversarial Training

Alexandre Araujo, Rafael Pinot, Benjamin Negrevergne, Laurent Meunier,  
Yann Chevaleyre, Florian Yger, Jamal Atif

► **To cite this version:**

Alexandre Araujo, Rafael Pinot, Benjamin Negrevergne, Laurent Meunier, Yann Chevaleyre, et al..  
Robust Neural Networks using Randomized Adversarial Training. 2019. hal-02380184v1

**HAL Id: hal-02380184**

**<https://hal.science/hal-02380184v1>**

Preprint submitted on 26 Nov 2019 (v1), last revised 6 Feb 2020 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# Robust Neural Networks using Randomized Adversarial Training

---

Alexandre Araujo<sup>1,2</sup> Rafael Pinot<sup>1,3</sup> Benjamin Negrevergne<sup>1</sup> Laurent Meunier<sup>1,4</sup>  
 Yann Chevalere<sup>1</sup> Florian Yger<sup>1</sup> Jamal Atif<sup>1</sup>

<sup>1</sup>PSL, Université Paris-Dauphine   <sup>2</sup>Wavestone   <sup>3</sup>CEA, Université Paris-Saclay   <sup>4</sup>Facebook AI Research

## Abstract

Since the discovery of adversarial examples in machine learning, researchers have designed several techniques to train neural networks that are robust against different types of attacks (most notably  $\ell_\infty$  and  $\ell_2$  based attacks). However, it has been observed that the defense mechanisms designed to protect against one type of attack often offer poor performance against the other. In this paper, we introduce *Randomized Adversarial Training* (RAT), a technique that is efficient both against  $\ell_2$  and  $\ell_\infty$  attacks. To obtain this result, we build upon adversarial training, a technique that is efficient against  $\ell_\infty$  attacks, and demonstrate that adding random noise at training and inference time further improves performance against  $\ell_2$  attacks. We then show that RAT is as efficient as adversarial training against  $\ell_\infty$  attacks while being robust against strong  $\ell_2$  attacks. Our final comparative experiments demonstrate that RAT outperforms all state-of-the-art approaches against  $\ell_2$  and  $\ell_\infty$  attacks.

## 1 Introduction

Modern neural networks achieve state of the art performances in a variety of domains such as natural language processing (Radford et al., 2018), image recognition (He et al., 2016) and speech recognition (Hinton et al., 2012). However, it has been shown that such neural networks can be vulnerable to *adversarial examples*, i.e. imperceptible variations of legitimate examples crafted to deliberately mislead a machine learning algorithm (Szegedy et al., 2014). In the literature, the magnitude of an acceptable variation is characterized using either the  $\ell_\infty$  or the  $\ell_2$  norm, and since these two metrics have very distinct

properties, the techniques used to generate adversarial examples also differ. In general,  $\ell_\infty$  adversarial examples are easy to generate, but they are also easy to detect due to the large  $\ell_2$  norm of their perturbation. In contrast generating  $\ell_2$  adversarial examples is computationally more demanding, but it is also more difficult to generate defense mechanisms that are robust against this type of attacks.

The Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2015) and other derived techniques (Madry et al., 2018) can be used to generate  $\ell_\infty$  adversarial examples that defeat off-the-shelf neural networks architectures with little computational resources. In the same work, Goodfellow et al. (2015) also proposed a defense mechanism and showed that augmenting batches with  $\ell_\infty$  adversarial examples during training could be used to significantly improve the robustness of the models against these attacks (a.k.a. Adversarial Training). However, they do not discuss the generality of Adversarial Training against other attacks.

Since the introduction of Adversarial Training, more sophisticated attacks have been proposed to craft adversarial perturbations with small  $\ell_2$  norm (Carlini & Wagner, 2017; Moosavi-Dezfooli et al., 2016) defeating existing defense mechanisms. Given the success of  $\ell_\infty$  Adversarial Training, it would be natural to consider  $\ell_2$  Adversarial Training as a candidate defense mechanism. Unfortunately for defenders, designing  $\ell_2$  adversarial training is impractical: state-of-the-art  $\ell_2$  adversarial training methods require much more computational resources and generating  $\ell_2$  adversarial examples during training is prohibitive. To cope with this limitation, several researchers have tried to design adversarial training methods based on adversarial examples generated with  $\ell_\infty$  attacks (see e.g (Madry et al., 2018)) but these methods demonstrate surprisingly little benefit in robustness against strong  $\ell_2$  attacks such as C&W (Carlini & Wagner, 2017).

In this paper, we are concerned with the problem of de-

signing a defense mechanism that is simultaneously robust against  $\ell_2$  and  $\ell_\infty$  attacks. We first give key insights to understand why  $\ell_\infty$  defense mechanisms perform poorly against  $\ell_2$  attacks: we observe that the sets of adversarial perturbations that can be generated with existing  $\ell_\infty$  or  $\ell_2$  attacks are strictly disjoint. Building on this observation we propose a new defense method called *Randomized Adversarial Training* (RAT) that combines Adversarial Training with noise injection at inference time.

Evaluating the robustness of randomized neural networks is a challenging task. Recently Athalye et al. (2018) have pointed flaws in the existing evaluation methodologies and claim that randomized networks should rather be tested against attacks in expectations. Under this new evaluation strategy, several existing methods are not as robust as they first appeared. In this paper, we follow the experimental methodology designed by Athalye et al. (2018) and evaluate RAT under attacks in expectation. Furthermore, we evaluate all our approaches against *untargeted* attacks which is the worst case scenario for defenders (an untargeted attack is successful if the attacker changes the predicted class to any other one). Our method proves to be more robust to expectation attacks than previous noise injection methods.

**Our contribution:** We introduce RAT, a Randomized Adversarial Training defense mechanism that combines adversarial training with noise injection at inference time, and demonstrate its efficiency against all state-of-the-art attacks. We then conduct experiments to highlight the difference between  $\ell_\infty$  and  $\ell_2$  attacks. Finally, we show that RAT achieves state-of-the-art robustness against  $\ell_\infty$  attacks in expectation, and outperforms existing defense mechanisms against  $\ell_2$  attacks. In other words, RAT is both more general and more robust than existing techniques. To the best of our knowledge, this is the first time that a defense mechanism obtains such good results against both types of attacks.

## 2 Related Work

Since the discovery of adversarial examples by Szegedy et al. (2014), a number of work have investigated the design of *attacks* to generate adversarial examples, and corresponding *defenses* to train models that are robust against these attacks. In their seminal work, Goodfellow et al. (2015) introduced the *Fast Gradient Sign Method* (FGSM), a technique that computes the gradient of the loss function with respect to an input  $x$  and perturbs  $x$  in the direction of the sign of the gradient to create an adversarial example. In the case of a linear classifier, their approach maximizes the change in the output for perturbations on the inputs that are bounded with the  $\ell_\infty$

norm. Then they argue that neural networks are locally linear (a property that is desirable to ease the training of deep neural nets), and demonstrate empirically that FGSM is able to generate good adversarial examples for neural networks. To protect against such attacks, Goodfellow et al. (2015) also introduced the adversarial training defense technique that we describe in the introduction (training networks using examples adversarial generated with FGSM). This defense offers satisfying protection against FGSM generated adversarial examples.

The FGSM attack can be seen as a one-step gradient method and Kurakin et al. (2017); Madry et al. (2018) have proposed variants of FGSM that perform multiple gradient descent steps. More specifically, Madry et al. (2018) introduced *Projected Gradient Descent* (PGD), an attack that can be used to generate adversarial examples defeating FGSM-based adversarial training. Naturally, Madry et al. (2018) also experimented adversarial training using adversarial examples generated with PGD. Adversarial training does offer an increased protection level against PGD adversarial examples, albeit at a higher computational cost. More recently, a number of defenses have been proposed (Song et al., 2018; Dhillon et al., 2018; Samangouei et al., 2018) and demonstrate good empirical results against some of the attacks, but do not offer a general improvement in robustness according to the thorough robustness study conducted by Athalye et al. (2018).

Several  $\ell_2$  attacks have also been developed (Moosavi-Dezfooli et al., 2016; Carlini & Wagner, 2017). In this paper, we mainly experiment with the  $\ell_2$  *C&W attack* by Carlini & Wagner (2017) which obtains the best results. This is an iterative and unbounded attack (but in practice the  $\ell_2$  distortion remains low). The attack is very efficient against all networks but requires significant computational resources (generating a single adversarial example can take minutes or even hours). Nevertheless this attack remains the most effective attack and the best existing defense mechanism so far (Madry et al. (2018)) only offers marginal protection against it.

An important idea to better protect models against strong  $\ell_2$  attacks is to add noise at inference. The presence of noise complicates the design of adversarial examples because the precise outcome of the model is not known by the attacker. The first method explicitly using this kind of technique was introduced by Xie et al. (2018) and a number of alternative noise based protection methods have been proposed (Liu et al., 2018; Rakin et al., 2018) offers some protection. However, Athalye et al. (2017) have remarked that randomization based protection mechanisms should be evaluated against attacks in expectation. Most existing noise-based protection mechanisms are weak under this evaluation methodology. Neverthe-

less, randomization remains the only viable mechanism against  $\ell_2$  attacks.

In addition, the robustness of randomized neural networks have also been studied from a theoretical perspective and these studies suggest that it is a good lead. For example, Fawzi et al. (2016) was able to compute lower and upper bounds on the minimal perturbation that fools a classifier in a randomly chosen direction. Finally Pinot et al. (2019) have demonstrated why injecting noise from exponential family at inference adds robustness, based on an in-depth analysis of the link between probability metrics (Wasserstein, Kullback Leibler, etc) and neural networks.

### 3 Preliminaries

In this section, we introduce preliminary definitions on randomized neural networks and recall the principles underlying the different attacks that we consider in this paper.

#### 3.1 Randomized neural networks

Let us first precise what we mean by a *Neural network* and a *Randomized neural network*.

**Definition 1.** An  $N$  layer neural network takes an input  $x$  and outputs a logit vector for  $x$  as follows:

$$F_\theta(x) := \phi^{(N)} \circ \dots \circ \phi^{(1)}(x)$$

$\theta$  is the parameter set of the network and  $\phi^{(i)}$  are the layers of the network parametrized by  $\theta$ . Those layers can either be convolutional, dense, or of any other form.

**Definition 2.** A randomized neural network is a neural network described as follows:

$$\tilde{F}_\theta(x) := \phi^{(N)} \circ \dots \circ (\phi^{(1)}(x) + b)$$

Where  $b$  is a multivariate Gaussian vector  $\mathcal{N}(0, \sigma I)$ .

According to Pinot et al. (2019), it is possible to use any kind of noise injection in the network as long as it comes from an exponential family. It is possible to inject noise on any on the  $N$  layers' activation. For the sake of simplicity, we restrict our study to the injection of Gaussian noise on the first layer of the network.

#### 3.2 Adversarial attacks

As discussed in Section 1, attacks can be divided in two categories:  $\ell_2$  and  $\ell_\infty$  attacks. Finding a good attack is a bi-criteria optimization problem, in which the attacker tries to minimize the norm of the perturbation while trying to maximize the loss function. There are two different ways

of implementing this: (1) by minimizing the norm under constraint on the loss, or (2) by maximizing the loss under constraint on the norm. There is no consensus so far and the state-of-the-art attacks that we describe in the rest of this section are implemented following either (1) or (2).

**Attack 1.** (Goodfellow et al., 2015)

*Fast Gradient Method with  $\ell_p$ -norm ( $FGM_p$ ) aims to solve the following type (2) optimization problem*

$$\operatorname{argmax}_{\|r\|_p \leq \epsilon} \mathcal{L}(F_\theta(x+r), y)$$

where  $\mathcal{L}$  is the loss function of the classifier (e.g. cross entropy),  $F_\theta$  is the logit predicted by the network and  $y$  the true label of  $x$ . Assuming  $\epsilon$  to be small, the previous problem can be approximated by solving the following:

$$\operatorname{argmax}_{\|r\|_p \leq \epsilon} \nabla_x \mathcal{L}(F_\theta(x), y)^\top r$$

When  $p \neq \infty$ , this corresponds to the  $FGM_p$  attack. When considering  $p = \infty$ , the solution to this problem is  $r = \epsilon \operatorname{sign}(\nabla_x \mathcal{L})$ , which corresponds to the  $FGSM$  attack.

**Attack 2.** (Carlini & Wagner, 2017)

*To craft an  $\ell_2$  adversarial example, Carlini & Wagner (2017) solves the following type (1) optimization problem*

$$\operatorname{argmin}_r \|r\|_2 + c \times g(x+r)$$

where  $g(x+r) < 0$  if and only if the  $f(x+r) \neq y$  ( $y$  here is the true label of  $x$ ), and  $f$  the underlying classifier of the neural network  $F_\theta$ . Carlini & Wagner (2017) use a change of variable  $r = \tanh(w) - x$  to ensure that  $-1 \leq x+r \leq 1$ , a binary search to optimize the constant  $c$ , and Adam or SGD to compute an approximated solution.

C&W attack is an  $\ell_2$  attack since the objective function aims to minimize the  $\ell_2$  norm of the perturbation. C&W can also be implemented as a  $\ell_\infty$  attack but this attack is not as effective as classical  $\ell_\infty$  attacks.

**Attack 3.** (Kurakin et al., 2016; Madry et al., 2018)

*The PGD attack is a generalization of the iterative FGSM attack (Kurakin et al., 2016). The goal of the adversary is to solve the following type (2) problem*

$$\operatorname{argmax}_{\|r\|_p \leq \epsilon} \mathcal{L}(F_\theta(x+r), y)$$

*In practice, the authors propose an iterative method to compute a solution:*

$$x^{t+1} = \prod_{x \oplus r} (x^t + \alpha \operatorname{sign}(\nabla_x \mathcal{L}(F_\theta(x^t), y)))$$

Where  $x \oplus r$  is the Minkowski sum between  $\{x\}$  and  $\{r \text{ s.t. } \|r\|_p \leq \epsilon\}$ ,  $\alpha$  a gradient step size,  $\Pi_S$  is the projection operator on  $S$  and  $x^0$  is randomly chosen in  $x \oplus r$ .

As for Fast Gradient Methods, PGD can be implemented either with  $p = \infty$  or  $p = 2$ . In the rest of this paper, we denote by  $\text{PGD}^t$ , a PGD attack with  $t$  steps.

**Remark 1.** *When the network is randomized, the good practice concerning evaluation is to use the Expectation Over Transformation (Athalye et al., 2017; Carlini et al., 2019) during the attack. If  $\tilde{F}_\theta(x)$  is a randomized network with  $x$  a given input then the Expectation Over Transformation (EOT) is  $\mathbb{E}(\tilde{F}_\theta(x))$ , where the expectation is taken over the randomness of  $\tilde{F}_\theta$ . In practice, the expectation is estimated by Monte Carlos type methods. As advised by Athalye et al. (2017) we use EOT to evaluate the performances of randomized methods.*

## 4 Randomized Adversarial Training

Though seemingly simple, Adversarial Training based on PGD or FGSM is one of the only defenses that practically proved robust (Athalye et al., 2018). This method has been extensively studied and demonstrates state-of-the-art results in  $\ell_\infty$  setting. Unfortunately it gives very poor results against state-of-the-art  $\ell_2$  attacks such as C&W. In this section we first explain why Adversarial Training based on  $\ell_\infty$  attacks does protect models against  $\ell_2$  attacks and vice-versa. Then, we introduce RAT, a defense strategy that mixes  $\ell_\infty$  adversarial training with randomization in order to offer simultaneous protection against powerful  $\ell_2$  and  $\ell_\infty$  attacks.

### 4.1 No free lunch in adversarial training

To understand why  $\ell_\infty$  adversarial training does not protect against  $\ell_2$  attacks, we first need to characterize the adversarial examples generated with  $\ell_\infty$  attacks.

**$\ell_\infty$  adversarial examples:** let  $\epsilon$  be the maximum amount of acceptable noise measured using the  $\ell_\infty$  norm and  $r$  be the perturbation of an input example  $x \in \mathbb{R}^d$  with  $\|r\|_\infty \leq \epsilon$ . By computing the sign of the gradient of the loss, FGSM finds perturbations  $r$  in the *corners* of the  $\ell_\infty$  ball of radius  $\epsilon$ , as illustrated in Figure 1 (left). These perturbations remain visually imperceptible because they have bounded  $\ell_\infty$  norm, but they are large when measured with the  $\ell_2$  norm:  $\|r\|_2 = \epsilon \times \sqrt{d}$ .  $\text{PGD}^n$ , the other  $\ell_\infty$  attack, generalizes FGSM by performing  $n$  iterative calls to FGSM. By doing so, after a single gradient step, the perturbations  $r$  also have a norm of  $\|r\|_2 = \epsilon \times \sqrt{d}$ . When the dimensionality  $d$  is large (as it is usually the case for most deep learning applications) all the pertur-

bations are likely clustered around the corners of the  $\ell_\infty$  ball, depicted in red in Figure 1, and far away from the  $\ell_2$  ball.

Let us now consider adversarial examples generated with  $\ell_2$  attacks.

**$\ell_2$  adversarial examples:** we choose  $\epsilon'$  to be the maximum acceptable noise measured in  $\ell_2$  norm that an  $\ell_2$  attack is allowed to generate. Therefore all  $\ell_2$  attacks must lie within the  $\ell_2$  ball of radius  $\epsilon'$ , the blue area depicted in Figure 1 (left). Remark that the  $\ell_2$  ball does not include the corners of the  $\ell_\infty$  ball unless we set  $\epsilon' \geq \epsilon\sqrt{d}$ . Thus, in high dimension (large values of  $d$ ), unless  $\epsilon'$  is set to a very large value, the  $\ell_2$  ball will not include the corners of the  $\ell_\infty$  ball. As a consequence, the sets of  $\ell_\infty$  and  $\ell_2$  adversarial examples are mutually disjoint.

Because the two sets of adversarial examples are disjoint, training models with  $\ell_\infty$  adversarial examples, can only protect  $\ell_2$  adversarial examples that are inside the  $\ell_\infty$  ball. The other adversarial examples remain unprotected. We illustrate this phenomenon using Figure 1 (middle): the cross at the center is the original example  $x$  and the square around the cross represents all acceptable  $\ell_\infty$  adversarial perturbations  $\|r\|_\infty \leq \epsilon$ . By training a classifier that labels  $x$  and all the  $\ell_\infty$  adversarial examples correctly (materialized by the red line), we do not guarantee that the classifier is protected against adversarial examples that are in the cap (i.e. the  $\ell_2$  adversarial examples that are in the  $\ell_2$  ball but not in the  $\ell_\infty$  ball).

In large dimension, the red corners are very far away from the  $\ell_2$  ball. Therefore we hypothesize that in this setting, a large proportion of the  $\ell_2$  adversarial examples remain unprotected. To verify this hypothesis, we measure the proportion of adversarial examples inside of the  $\ell_2$  ball before and after  $\ell_\infty$  adversarial training. The results are presented in Figure 2 (left: without adversarial training, right: with adversarial training).

On both charts, the blue area represents the proportion of adversarial training examples that are inside the  $\ell_\infty$  ball. The red area represents the adversarial examples that are outside the  $\ell_\infty$  ball but still inside the  $\ell_2$  ball (valid  $\ell_2$  adversarial examples). Finally the brown beige area represents the adversarial examples that are beyond the  $\ell_2$  bound. The radius  $\epsilon'$  of the  $\ell_2$  ball varies along the x-axis from  $\epsilon'$  to  $\epsilon'\sqrt{d}$ .

On the left chart (without adversarial training) most  $\ell_2$  adversarial examples generated by C&W are also inside the  $\ell_\infty$  ball. On the right chart much fewer adversarial examples that remain in the  $\ell_\infty$  ball, this is the expected consequence of  $\ell_\infty$  adversarial training. However, most of these adversarial examples have shifted in the area

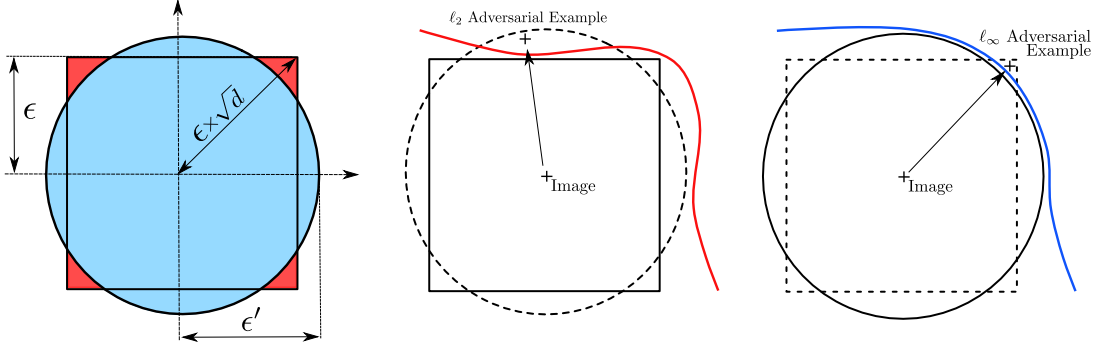


Figure 1: Left: 2D representation of the  $l_\infty$  and  $l_2$  balls of respective radius  $\epsilon$  and  $\epsilon'$ . Middle: a classifier trained with  $l_\infty$  adversarial perturbations (materialized by the red line) remains vulnerable to  $l_2$  attacks. Right: a classifier trained with  $l_2$  adversarial perturbations (materialized by the blue line) remains vulnerable to  $l_\infty$  adversarial examples (note, this case is impractical).

in the  $l_2$  ball that is not in the  $l_\infty$  ball (i.e. the cap). These examples are equally good from the  $l_2$  perspective, and the impact of  $l_\infty$  adversarial training in terms of  $l_2$  robustness is effectively null.

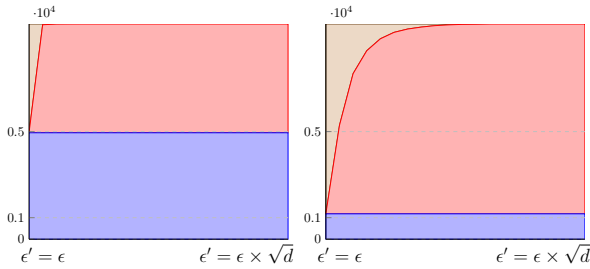


Figure 2: Comparison of the number of adversarial examples found by C&W, inside the  $l_\infty$  ball (lower, blue area), outside the  $l_\infty$  ball but inside the  $l_2$  ball (middle, red area) and outside the  $l_2$  ball (upper gray area).  $\epsilon$  is set to 0.3 and  $\epsilon'$  varies along the x-axis. Left: without adversarial training, right: with adversarial training. Most adversarial examples, have shifted from the  $l_\infty$  ball to the cap of the  $l_2$  ball, but remain at the same  $l_2$  distance from the original example.

Finally, remark that even though it is impractical,  $l_2$  adversarial training would also be inefficient against  $l_\infty$  adversarial attacks for similar reasons, as illustrated in the Figure 1 (right).

We conclude that the protection mechanisms are complementary, and that mixing  $l_\infty$  and  $l_2$  protection mechanisms is necessary to obtain a strong and general protection mechanism. Building on this insight, we introduce RAT and LeRAT, two protection mechanisms that include defense techniques against both types of attacks.

## 4.2 RAT & LeRAT

To construct a defense strategy that is both robust to  $l_2$  and  $l_\infty$  attacks, we propose Randomized Adversarial Training (RAT), described in Algorithm 1.

The main difference between AT and RAT is that the latter is designed to train randomized networks instead of deterministic networks. Recall<sup>1</sup> that each time a randomized network  $\tilde{F}_\theta$  is used to make a prediction on an example or to compute a gradient, a new noise vector  $b$  is drawn and added to the activation of the first layer of the network, thus producing a randomized prediction (respectively a randomized gradient). In RAT, this randomization is done both when computing the adversarial attack (line 5 in Algorithm 1) and when updating the weights of the network (line 7). Note that for the sake of limiting the computational cost of the training phase, EOT method isn't used to compute adversarial examples.

---

### Algorithm 1 RAT

---

- 1: Initialize  $\theta$
  - 2: **for**  $i \in [1, nb\_steps]$  **do**
  - 3:    $B_i := (x_k, y_k)_k$  the current mini-batch
  - 4:   **for**  $x_k \in B_i$  **do**
  - 5:     Compute adversarial attack  $x'_k$  of  $x_k$  against  $\tilde{F}_\theta$  using an arbitrary method
  - 6:   **end for**
  - 7:   Update  $\theta$  following the criterion:  $\frac{1}{|B_i|} \sum_k \mathcal{L}(\tilde{F}_\theta(x'_k), y_k)$
  - 8: **end for**
  - 9: **return** Network  $\tilde{F}_\theta$
- 

<sup>1</sup>see definition 2

**Learned Randomized Adversarial Training** : *Learned Randomized Adversarial Training* (LeRAT) is a natural variation of RAT where internal noise injection parameters are learned during training instead of being chosen by grid search. The rationale of this procedure is that learning the noise injection should help the network finding the good amount of noise to defend. One should note that this procedure can only be used with the adversarial loss as a classical loss would instantly shut down the noise injection.

## 5 RAT the robust

In this section, we first discuss the robustness of neural networks trained with RAT against  $\ell_\infty$  attacks. Then we discuss their performances against stronger  $\ell_2$  attacks and provide a number of empirical insights to better understand the behaviour of protection mechanisms against this type of attacks. Finally, we discuss the generality of the RAT approach. All code and datasets will be release on our website upon acceptance of this paper.

To evaluate the robustness of the protection mechanisms against  $\ell_\infty$  and  $\ell_2$  attacks we follow the rigorous experimental protocol proposed by Athalye et al. (2018), that is: we evaluate the randomized networks using EOT attacks. We also consider untargeted attacks (i.e. a misclassification from one target to any other will be considered as a failure for the defender), and using the white-box attack setting (i.e. we assume that the attacker has unlimited access to the weights of the network). Overall, this represents the most conservative setting to evaluate protection mechanisms. See Section 5.4 for a comprehensive description of our experimental setting.

### 5.1 Robustness against $\ell_\infty$ attacks

We evaluate RAT and LeRat against  $\ell_\infty$  attacks generated with FGSM and its stronger variant: PGD (Madry et al., 2018). We also compare it against state-of-the-art Adversarial Training as well as randomized networks proposed by Pinot et al. (2019). We first experiment with a varying  $\ell_2$  bound on the perturbation. The results are presented in Figure 3.

Unsurprisingly, the accuracy lowers for all protection mechanisms when we raise the value of  $\epsilon$  from 0.05 to 0.3. (i.e. when we tolerate more perturbed adversarial examples.) On this setting, LeRAT obtains the best robustness together with state-of-the-art  $\ell_\infty$  protection mechanism, R.N. demonstrate lower performance on this setting but RAT demonstrate reasonable performances especially on strong PGD attacks.

In a second set of experiments we fix the bound  $\epsilon = 0.1$

and increase the cardinality of the sample used to compute the attack in expectation (more samples means stronger EOT attacks). The results are shown in Figure 4 (a) and (b).

As we can see plain R.N. offers a lower protection against strong EOT attacks. In contrast LeRAT is not impacted by the stronger attacks, and RAT offers a good compromise.

### 5.2 Robustness against $\ell_2$ attacks (C&W)

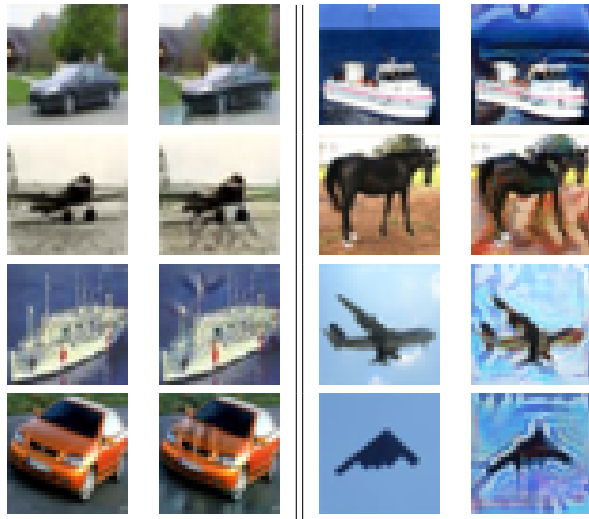


Table 1: Original images and the corresponding adversarial examples generated by C&W attack against our approach. On the left of the double line ( $\ell_2$  norm of the perturbation = 4), the perturbations easily noticeable. On the right ( $\ell_2$  norm = 10, worst adversarial examples) the perturbation can be noticed, even without comparing with the original image.

We now evaluate our methods against strong  $\ell_2$  attacks using the C&W attack. Given an input image, the C&W attack finds the minimum perturbation that will lead a classifier to a misclassification. We evaluate our approach using this setting (as in (Athalye et al., 2018)) and then discuss other relevant aspects the evaluation. The results are presented in Figure 4 (c).

First, we can see that adversarial training (A.T.) offers no protection whatsoever against  $\ell_2$  attacks. We can also see that all randomized networks outperform again A.T. trained with  $\ell_\infty$  adversarial examples. As we explained in Section 4,  $\ell_\infty$  adversarial training cannot offer proper protection against  $\ell_2$  attacks. Finally, we can see that the robustness of LeRAT is significantly lower than the other two approaches.

Remark that in this experimental the norm of the perturbation is not bounded, therefore there are no guarantees

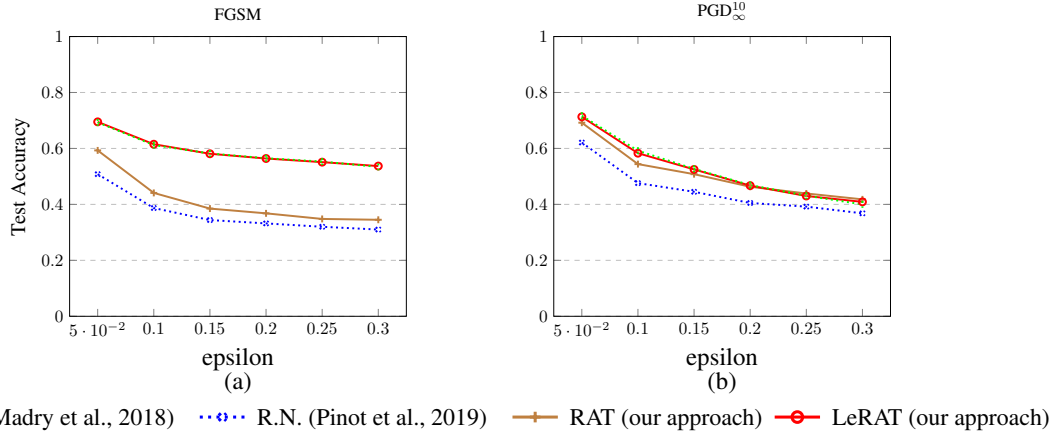


Figure 3: These figures represent the test accuracy on the CIFAR10 dataset given the bound  $\epsilon$  of the  $\ell_{\infty}$  attacks. The value range from  $5 \cdot 10^{-2}$  to 0.3 and the higher the value, the more the images is perturbed. We can see that A.T. perform equally with LeRAT while RAT outperform R.N.

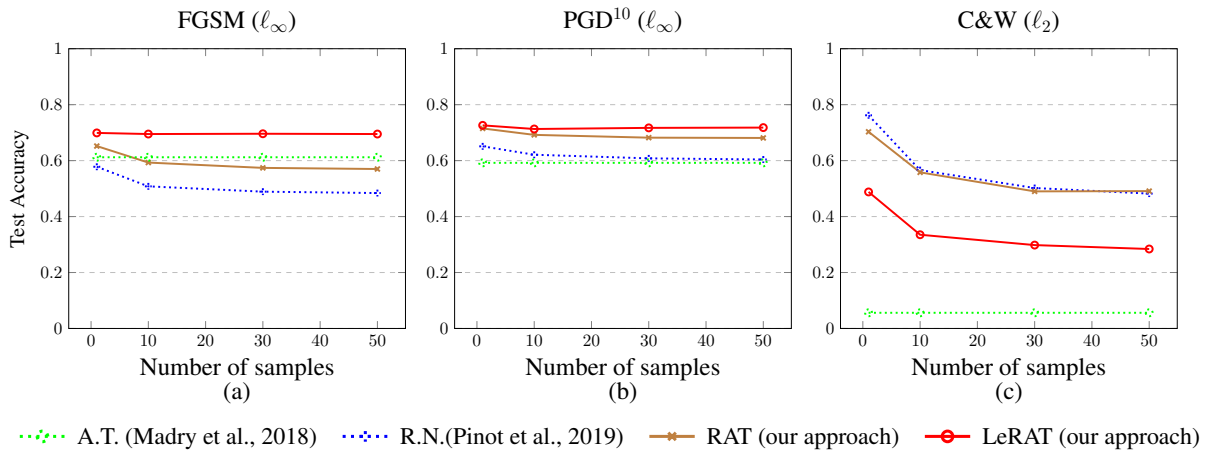


Figure 4: This figure shows the accuracy of different models against  $FGSM \ell_{\infty}$  attacks (a),  $PGD^{10} \ell_{\infty}$  attacks (b) and  $C\&W \ell_2$  attacks (c). All attacks are EOT attacks estimated using an increasing sample size. LeRAT outperforms other approaches on  $\ell_{\infty}$ . RAT and R.N. outperform LeRAT on  $\ell_2$  attacks. RAT offer the best tradeoff.

that the actual adversarial example is below the threshold of human perception. In Figure 5, we have plotted statistics about the adversarial examples generated by C&W to defeat each defense mechanism.

The statistics in Figure 5 show that the average  $\ell_2$  norm of adversarial examples required to defeat all randomized networks is surprisingly high. In Figure 1 we provide several adversarial examples of variable  $\ell_2$  norm. As we can see, adversarial examples with an  $\ell_2$  norm above 4 do defeat the defense mechanism but they are easily identified with a human eye. In other words, defeating these defense mechanisms may require perturbation that are too big to remain undetected by humans.

### 5.3 General performance against strong attacks and general discussion

We now report the performance of different approaches against stronger attacking settings. For  $\ell_{\infty}$  we use a  $PGD^{10}$  attack with a large  $\epsilon$  bound of 0.3. For  $\ell_2$ , we use C&W with 40 iterations. In both case, we use 50 Monte-Carlo iterations to estimate the best EOT perturbation with high accuracy. The results are summarized in Table 2. As we can see, RAT and LeRAT offer the best robustness against  $\ell_{\infty}$  and RAT offers the best robustness against  $\ell_2$  attacks. RAT outperforms the state-of-the-art against all existing defense mechanisms by a significant margin and offers the best compromise overall.



Model	$\ell_\infty$ Attack	$\ell_2$ Attack
Adv Training (A.T.) (Madry et al., 2018)	0.400	0.056
Randomized Network (R.N.) (Pinot et al., 2019)	0.356	0.489
Randomized Adv Training (RAT)	<b>0.409</b>	<b>0.499</b>
Learned Randomized Adv Training (LeRAT)	0.400	0.296

Table 2: Comparison of the robustness of all approaches against  $\ell_\infty$  and  $\ell_2$  attacks. Adversarial Training offers good protection against  $\ell_\infty$  attacks but performs poorly against  $\ell_2$  attacks. The  $\ell_\infty$  attack is computed with PGD<sup>10</sup> with an epsilon of 0.3, the  $\ell_2$  attack is made with C&W with 40 iterations. All attacks against randomized network are computed using a Monte Carlo with 50 iterations in order to estimate the expected value of the network. Randomize networks offer good protection against  $\ell_2$  attacks but is suboptimal against  $\ell_\infty$  attacks. Our approach (RAT) offers the best protection against all attacks.

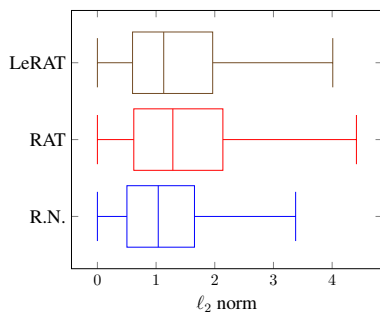


Figure 5: This chart shows statistics about the adversarial examples generated by the C&W attack to defeat different defenses. The adversarial examples generated to fight the RAT defense have higher norm (and thus are easier to identify) than for other defenses.

#### 5.4 Details of the evaluation methodology

**Dataset** For our experiments, we use the CIFAR10 dataset (Krizhevsky et al., 2009), which consist of 60 000 32x32 colour images with 10 classes. Theses dataset are split with 50 000 images for the training set and 10 000 images for evaluation.

**Baseline** We used the Wide Resnet architecture introduced by (Zagoruyko & Komodakis, 2016) for all our experiments with 10 and 28 as widen factor and depth. Every input is normalized between -1 and +1. Our model is learned on 200 epochs with a batch size of 200 and a starting learning rate of 0.1. During the training, the learning rate decreases to 0.02, 0.004 and 0.00008 respectively after 7500, 15000 and 20 000 steps. We also use random crop and random flip as data augmentation techniques. For regularization, we use dropout and 0.002 weight decay. To train the network, we use the cross entropy loss with momentum 0.9. This network architecture achieves a 0.95 accuracy on CIFAR10.

**Randomized Adversarial Training** For adversarial training, we use the procedure described in algorithm 1 with *Projected Gradient Descent* as an attack algorithm. PGD is implemented with 1 iteration, and a global perturbation intensity of 0.03, and gradient step size of 0.02. In order to choose the variance of the injected noise, we computed a grid search that lead us to choose noise drawn from a Gaussian with mean 0 and variance 1 as a good trade off between accuracy and robustness.

**Evaluation under attack** To evaluate our approach, we compute the accuracy under attack with FGSM (Goodfellow et al., 2015), PGD (Madry et al., 2018), as  $\ell_\infty$  attacks and C&W (Carlini & Wagner, 2017), as  $\ell_2$ . For PGD, we use 10 iterations with a gradient step size of 0.02. For C&W, we use 9 iteration for the binary search and 40 iterations overall. As stated in (Athalye et al., 2017), attacks against randomized defense need to be computed with the expected value of the loss or the logits. In our experiment, we estimate the expected value of the attack with a Monte Carlo.

## 6 Conclusion

We conducted a thorough experimental analysis and provided insights to understand the behaviour of several defense mechanisms under different types of attacks. Building on these insights, we introduce two novel techniques named RAT et LeRAT that can be used to train robust networks against  $\ell_\infty$  and  $\ell_2$  attacks. Then we compare them against existing state-of-the-art techniques: adversarial training which is state-of-the-art for  $\ell_\infty$  attacks and Randomized Networks which is state-of-the-art on  $\ell_2$  attacks. Our experiments show that our approach outperform adversarial training on  $\ell_\infty$  attacks and randomized networks on  $\ell_2$  attacks, resulting in an approach that is both more general and more robust than the existing approaches.

## References

- Athalye, A., Engstrom, L., Ilyas, A., and Kwok, K. Synthesizing robust adversarial examples. *arXiv preprint arXiv:1707.07397*, 2017.
- Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 274–283, Stockholm, Sweden, 10–15 Jul 2018. PMLR.
- Carlini, N. and Wagner, D. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57. IEEE, 2017.
- Carlini, N., Athalye, A., Papernot, N., Brendel, W., Rauber, J., Tsipras, D., Goodfellow, I., and Madry, A. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705*, 2019.
- Dhillon, G. S., Azzadenesheli, K., Bernstein, J. D., Kosaiji, J., Khanna, A., Lipton, Z. C., and Anandkumar, A. Stochastic activation pruning for robust adversarial defense. In *International Conference on Learning Representations*, 2018.
- Fawzi, A., Moosavi-Dezfooli, S.-M., and Frossard, P. Robustness of classifiers: from adversarial to random noise. In *Advances in Neural Information Processing Systems*, pp. 1632–1640, 2016.
- Goodfellow, I., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- Hinton, G., Deng, L., Yu, D., Dahl, G., Mohamed, A.-r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Kingsbury, B., et al. Deep neural networks for acoustic modeling in speech recognition. *IEEE Signal processing magazine*, 29, 2012.
- Krizhevsky, A., Nair, V., and Hinton, G. Cifar-10. *Canadian Institute for Advanced Research*, 2009.
- Kurakin, A., Goodfellow, I., and Bengio, S. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016.
- Kurakin, A., Goodfellow, I. J., and Bengio, S. Adversarial machine learning at scale. In *International Conference on Learning Representations*, 2017.
- Liu, X., Cheng, M., Zhang, H., and Hsieh, C.-J. Towards robust neural networks via random self-ensemble. In *European Conference on Computer Vision*, pp. 381–397. Springer, 2018.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- Moosavi-Dezfooli, S.-M., Fawzi, A., and Frossard, P. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2574–2582, 2016.
- Pinot, R., Meunier, L., Araujo, A., Kashima, H., Yger, F., Gouy-Pailler, C., and Atif, J. Theoretical evidence for adversarial robustness through randomization: the case of the exponential family. *arXiv preprint arXiv:1902.01148*, 2019.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language models are unsupervised multitask learners. Technical report, OpenAi, 2018.
- Rakin, A. S., He, Z., and Fan, D. Parametric noise injection: Trainable randomness to improve deep neural network robustness against adversarial attack. *arXiv preprint arXiv:1811.09310*, 2018.
- Samangouei, P., Kabkab, M., and Chellappa, R. Defense-GAN: Protecting classifiers against adversarial attacks using generative models. In *International Conference on Learning Representations*, 2018.
- Song, Y., Kim, T., Nowozin, S., Ermon, S., and Kushman, N. Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. In *International Conference on Learning Representations*, 2018.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014.
- Xie, C., Wang, J., Zhang, Z., Ren, Z., and Yuille, A. Mitigating adversarial effects through randomization. In *International Conference on Learning Representations*, 2018.
- Zagoruyko, S. and Komodakis, N. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.