

# Robustness and geometry of deep neural networks

**Alhussein Fawzi**

DeepMind

May 23rd 2019

The Mathematics of Deep Learning and Data Science  
University of Cambridge

# Recent advances in machine learning



GT: horse cart  
1: horse cart  
2: minibus  
3: oxcart  
4: stretcher  
5: half track



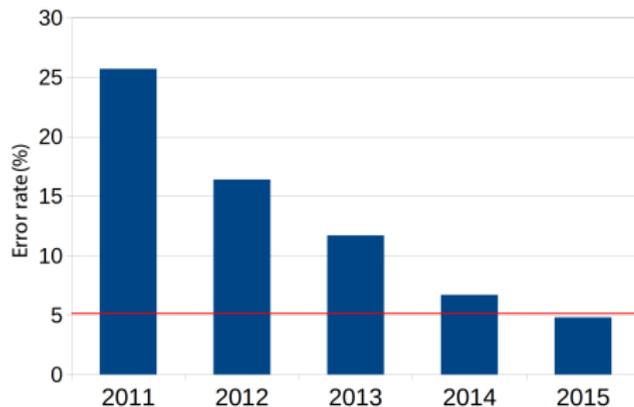
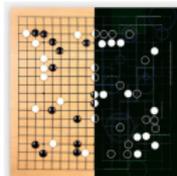
GT: birdhouse  
1: birdhouse  
2: sliding door  
3: window screen  
4: mailbox  
5: pot



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



# Robustness of classifiers to perturbations

- In real world environments, images undergo perturbations.

# Robustness of classifiers to perturbations

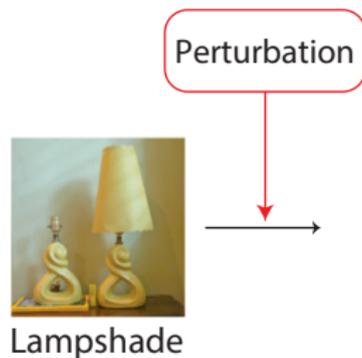
- In real world environments, images undergo perturbations.



Lampshade

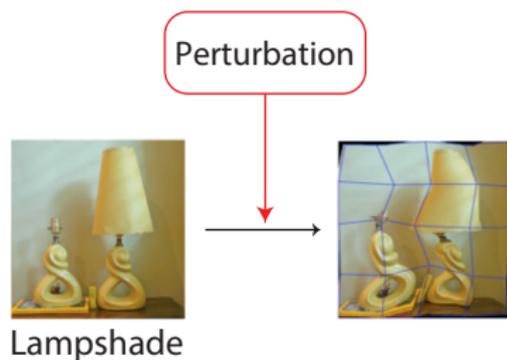
# Robustness of classifiers to perturbations

- In real world environments, images undergo perturbations.



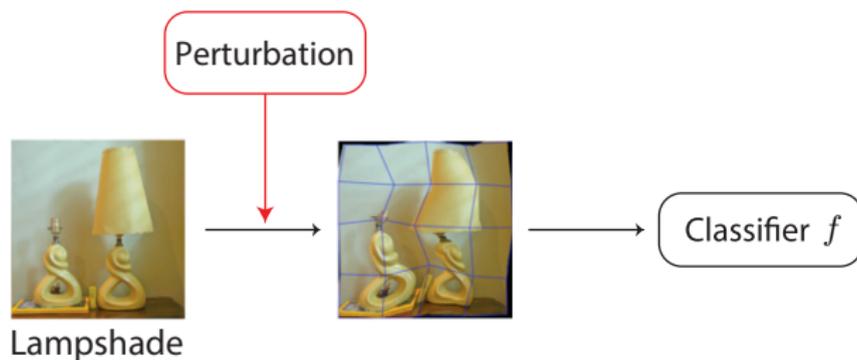
# Robustness of classifiers to perturbations

- In real world environments, images undergo perturbations.



# Robustness of classifiers to perturbations

- In real world environments, images undergo perturbations.



# Robustness of classifiers to perturbations

- In real world environments, images undergo perturbations.



# Robustness of classifiers to perturbations

- In real world environments, images undergo perturbations.



- Broad range of perturbations
  - Adversarial perturbations [Szegedy et. al. ICLR 2014], [Biggio et. al., PKDD 2013], ...
  - Random noise [Fawzi et. al., NIPS 2016], [Franceschi et. al., AISTATS 2018]
  - Structured nuisances (geometric transformations [Bruna et. al., TPAMI 2013], [Jaderberg et. al., NIPS 2015], occlusions [Sharif et. al., CCS 2016], etc...).

# Robustness of classifiers to perturbations (Cont'd)

- Safety of machine learning systems

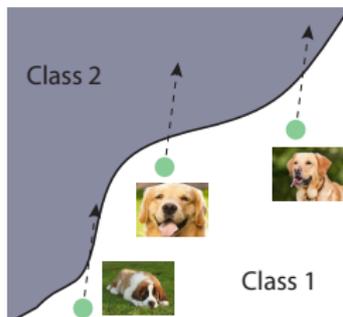


# Robustness of classifiers to perturbations (Cont'd)

- Safety of machine learning systems



- Better understanding of the geometry of state-of-the-art classifiers.



# Talk outline

- ① Fooling classifiers is easy: vulnerability to different perturbations.
- ② Improving the robustness (i.e., “defending”) is difficult.
- ③ Geometric analysis of a successful defense: adversarial training.

# Adversarial perturbations

- State-of-the-art deep neural networks have been shown to be surprisingly unstable to *adversarial* perturbations.

# Adversarial perturbations

- State-of-the-art deep neural networks have been shown to be surprisingly unstable to *adversarial* perturbations.

School bus



# Adversarial perturbations

- State-of-the-art deep neural networks have been shown to be surprisingly unstable to *adversarial* perturbations.

School bus



Ostrich



# Adversarial perturbations

- State-of-the-art deep neural networks have been shown to be surprisingly unstable to *adversarial* perturbations.

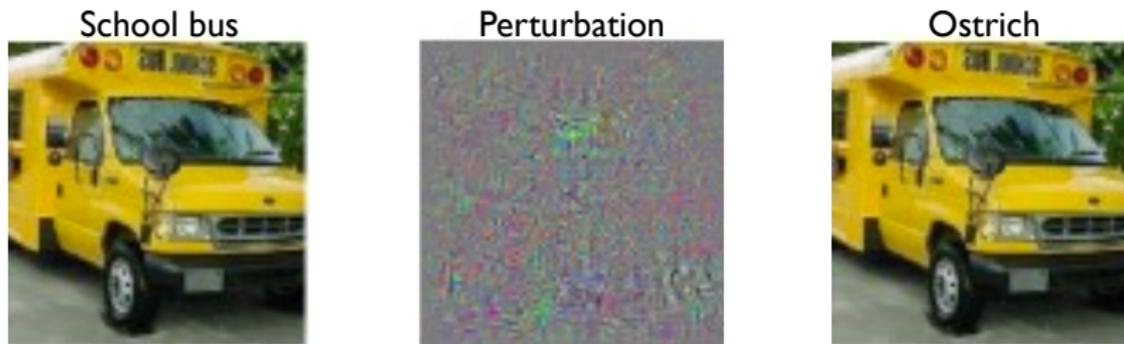


Figure from [Szegedy et. al., ICLR 2014].

# Adversarial perturbations

- State-of-the-art deep neural networks have been shown to be surprisingly unstable to *adversarial* perturbations.

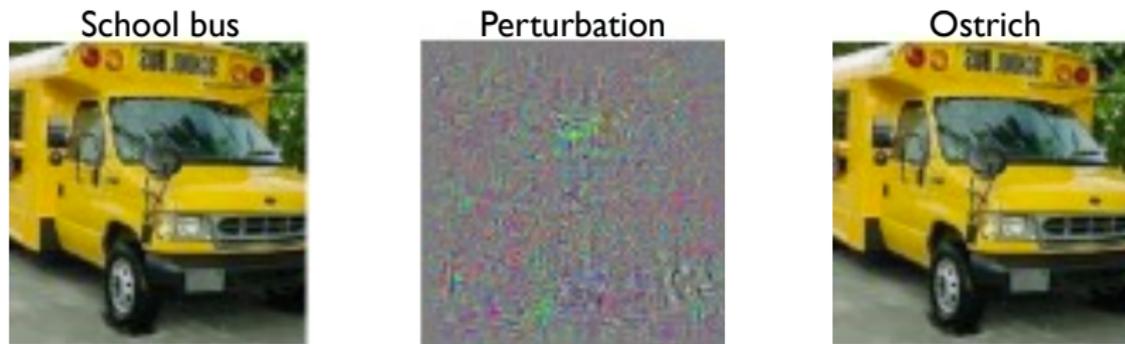
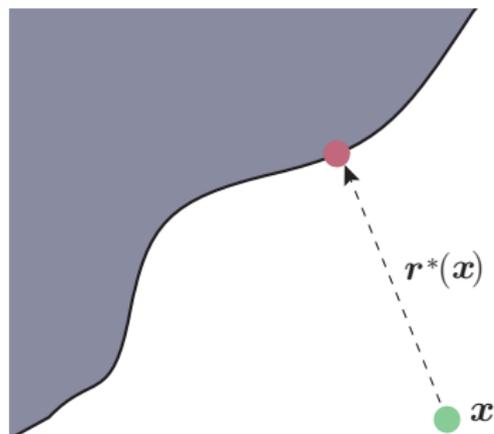


Figure from [Szegedy et. al., ICLR 2014].

- Adversarial examples are found by seeking the minimal perturbation (in the  $\ell_2$  sense) that switches the label of the classifier.

# Adversarial perturbations

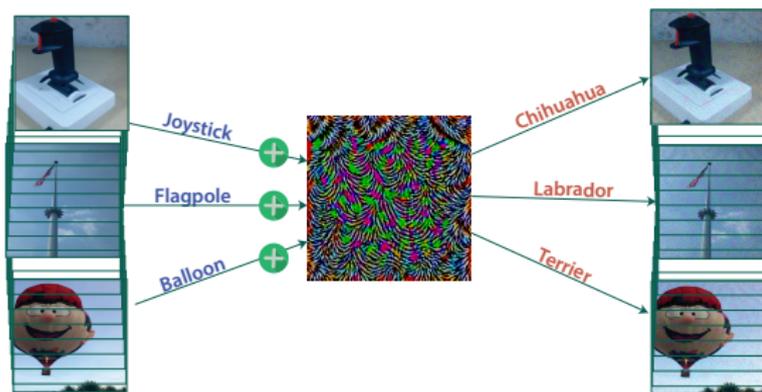
Robustness to **adversarial** noise



$$r^*(x) = \min_r \|r\|_2 \text{ subject to } f(x+r) \neq f(x).$$

# Other types of adversarial perturbations

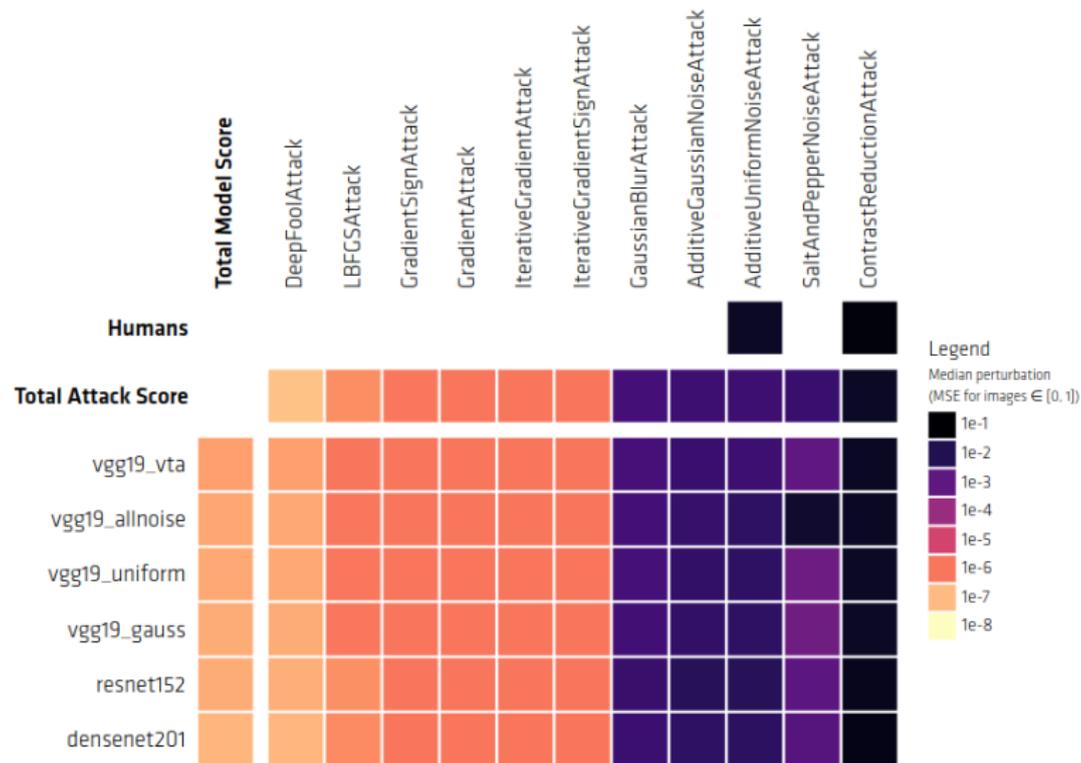
- Universal perturbations [Moosavi-Dezfooli et. al., 2017]



- Geometric transformations [Fawzi et. al., 2015, Moosavi-Dezfooli et. al., 2018, Xiao et al., 2018]



# Finding adversarial perturbations is easy...



<http://robust.vision>

# ... but designing defense mechanisms is hard!

Despite the huge number of proposed defenses, state-of-the-art classifiers are still vulnerable to small perturbations.

## Adversarial Risk and the Dangers of Evaluating Against Weak Attacks

Jonathan Uesato<sup>1</sup>, Brendan O'Donoghue<sup>1</sup>, Aaron van den Oord<sup>1</sup>, Dhruv Kohli<sup>1</sup>

### Abstract

This paper investigates recently proposed approaches for defending against adversarial examples and evaluating adversarial robustness. We motivate adversarial risk as an objective for achieving models robust to worst-case inputs. We then frame commonly used attacks and evaluation metrics as defining a tractable surrogate objective to the true adversarial risk. This suggests that models may optimize this surrogate rather than the true adversarial risk. We formalize this notion as *obscurity* as an adversary, and develop tools and heuristics for identifying obscured models and designing transparent models. We demonstrate that this is a significant problem in practice by representing gradient-free optimization techniques into adversarial attacks, which we use to decrease the accuracy of several recently proposed defenses to near zero. Our hope is that our formulations and results will help researchers to develop more powerful defenses.

### 1. Introduction

Deep learning has revolutionized the field of machine learning and led to substantial improvements in many challenging problems such as image understanding (He et al., 2016), speech recognition (Graves et al., 2013), and automatic game playing (Mnih et al., 2015). Despite these remarkable successes, we have seen some intriguing and worrying properties in the behavior of these models.

Researchers have demonstrated that certain small perturbations to the input can make neural networks perform extremely bad results (Szegedy et al., 2013; Liu and Liang, 2017). For instance, in the case of image classification, imperceptible perturbations can lead to the resulting images being misclassified to completely different object categories with high confidence (Gong et al., 2013). These so-called

<sup>1</sup>DeepMind. Correspondence to Jonathan Uesato (joesato@google.com).

Proceedings of the 35<sup>th</sup> International Conference on Machine Learning, Stockholm, Sweden, PMLR 80, 2018. Copyright 2018 by the author(s).

adversarial examples can be computed relatively easily by using optimization methods (referred to as adversarial attacks) to find perturbations that maximize the loss of the network (Goodfellow et al., 2014; Carlini and Wagner, 2016). Later work showed that this phenomena is not unique to image classification, and appears across different model architectures, as well as in other machine learning algorithms (Papernot et al., 2016, 2017).

The emergence of adversarial examples and the increasing deployment of machine learning models in real world production systems has motivated extensive research on developing models that can defend against adversarial attacks (Wilde-Furley and Goodfellow, 2016; Yuan et al., 2017). Using a variety of approaches, these works have shown that neural models can be developed that are robust against commonly used attacks (Goo et al., 2017; Song et al., 2017; Xie et al., 2017; Liu et al., 2017). However, a key question remains unanswered: *Are these models free from any adversarial examples or are they simply robust to current attack methods?*

In this paper, we formalize the intuition that in settings with the potential for catastrophic failures, minimizing expected risk may produce models with very poor worst-case performance. This motivates the study of the adversarial risk as a measure of the model's performance on worst-case inputs. However, the exact adversarial risk is computationally intractable to evaluate. In its place, we then frame commonly used attacks and adversarial evaluation metrics as defining a tractable surrogate objective to the true adversarial risk. We hypothesize that many defenses achieve robustness through obscuring, i.e., the defenses work by exploiting weaknesses of certain attacks and do not eliminate all adversarial risks.

One of the key contributions of this paper is to explicitly validate the "security by obscurity" nature of recently proposed defense methods. Specifically, we show that by using a more powerful attack (inflexible) one better able to maximize the true adversarial risk, we can dramatically reduce the performance of these defenses.

To summarize, the key contributions of this paper are:

- Formulation of adversarial attacks and defenses as optimizing surrogates of the true adversarial risk

## Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples

Anish Athalye<sup>1</sup>, Nicholas Carlini<sup>1</sup>, David Wagner<sup>2</sup>

### Abstract

We identify obfuscated gradients, a kind of gradient masking, as a phenomenon that leads to a false sense of security in defenses against adversarial examples. While defenses that cause obfuscated gradients appear to defeat iterative optimization-based attacks, we find defenses relying on this effect can be circumvented. We describe characteristic behaviors of defenses exhibiting the effect, and for each of the three types of obfuscated gradients we discover, we develop attack techniques to overcome it. In a case study, examining non-certified white-box secure defenses at ICLR 2018, we find obfuscated gradients are a common occurrence, with 7 of 9 defenses relying on obfuscated gradients. Our new attacks successfully circumvent 6 completely, and 1 partially, in the original threat model each paper considers.

### 1. Introduction

In response to the susceptibility of neural networks to adversarial examples (Szegedy et al., 2013; Biggio et al., 2013), there has been significant interest recently in constructing defenses to increase the robustness of neural networks. While progress has been made in understanding and defending against adversarial examples in the white-box setting, where the adversary has full access to the network, a complete solution has not yet been found.

As benchmarking against iterative optimization-based attacks (e.g., Karanik et al. (2016a); Madry et al. (2018); Carlini & Wagner (2017)) has become standard practice in evaluating defenses, new defenses have arisen that appear to be robust against these powerful optimization-based attacks.

We identify one common reason why many defenses provide

<sup>1</sup>Equal contribution. <sup>2</sup>Massachusetts Institute of Technology, University of California, Berkeley. Correspondence to: Anish Athalye (athalye@mit.edu), Nicholas Carlini (nc@berkeley.edu).

Proceedings of the 35<sup>th</sup> International Conference on Machine Learning, Stockholm, Sweden, PMLR 80, 2018. Copyright 2018 by the author(s).

apparent robustness against iterative optimization attacks: obfuscated gradients, a term we define as a special case of gradient masking (Papernot et al., 2017). Without a good gradient, where following the gradient does not successfully optimize the loss, iterative optimization-based methods cannot succeed. We identify three types of obfuscated gradients: *obscured gradients* are nonexistent or incorrect gradients caused either internally through non-differentiable operations or unintentionally through numerical instability; *assorted gradients* depend on test-time randomness; and *vanishing/exploding gradients* in very deep computation result in an unusable gradient.

We propose new techniques to overcome obfuscated gradients caused by these three phenomena. We address gradient obscuring with a new attack technique we call Backward Pass Differentiable Approximation, where we approximate derivatives by computing the forward pass normally and computing the backward pass using a differentiable approximation of the function. We compare gradients of randomized defenses by applying Expectation Over Transformation (Athalye et al., 2017). We solve vanishing/exploding gradients through regularization and optimize over a space where gradients do not explode/vanish.

To investigate the prevalence of obfuscated gradients and understand the applicability of these attack techniques, we use as a case study the ICLR 2018 non-certified defenses that claim white-box robustness. We find that obfuscated gradients are a common occurrence, with 7 of 9 defenses relying on this phenomenon. Applying the new attack techniques we develop, we overcome obfuscated gradients and circumvent 6 of them completely, and 1 partially, under the original threat model of each paper. Along with this, we offer an analysis of the evaluations performed in the papers. Additionally, we hope to provide researchers with a common baseline of knowledge, description of attack techniques, and common evaluation pitfalls, so that future defenses can avoid falling vulnerable to these same attack approaches.

To promote reproducible research, we release our re-implementation of each of these defenses, along with implementation of our attacks for each.

<sup>1</sup>https://github.com/nickadriano/obfuscated-gradients

arXiv:1802.05666v2 [cs.LG] 12 Jun 2018

arXiv:1802.00420v4 [cs.LG] 31 Jul 2018

# ... but designing defense mechanisms is hard!

## Despite the huge number of proposed defenses, state-of-the-art classifiers are still vulnerable to small perturbations.

### Adversarial Risk and the Dangle

Jonathan Uesato<sup>1</sup> Brendan O'Donoghue

#### Abstract

This paper investigates recently proposed approaches for defending against adversarial examples and evaluating adversarial robustness. We motivate *adversarial risk* as an objective for achieving models robust to worst-case inputs. We then focus concretely on attack and evaluation metrics as defining a tractable surrogate objective to the true adversarial risk. This suggests that models may optimize this surrogate rather than the true adversarial risk. We formalize this notion as *obscurity* as an adversary, and develop tools and heuristics for identifying obscured models and designing transparent models. We demonstrate that this is a significant problem in practice by reinterpreting gradient-free optimization techniques into adversarial attacks, which we use to decrease the accuracy of several recently proposed defenses to near zero. Our hope is that our formulations and results will help researchers to develop more powerful defenses.

#### 1. Introduction

Deep learning has revolutionized the field of machine learning and led to substantial improvements in many challenging problems such as image understanding (He et al., 2016), speech recognition (Srivastava et al., 2015), and autonomous game playing (Mnih et al., 2015). Despite these remarkable successes, we have seen some intriguing and worrisome properties in the behavior of these models.

Researchers have demonstrated that certain exact perturbations to the input can make neural networks perform terribly bad results (Szegedy et al., 2013; Liu and Li, 2017). For instance, in the case of image classification, perceptible perturbations can lead to the resulting image being misclassified to completely different object classes with high confidence (Szegedy et al., 2013). These occur

<sup>1</sup>Corresponding. Correspondence to: Jonathan Uesato, joesato@google.com.

Proceedings of the 35th International Conference on Machine Learning, Stockholm, Sweden, PMLR 80, 2018. Copyright 2018 by the author(s).

### MagNet and “Efficient Defenses Against Adversarial Attacks” are Not Robust to Adversarial Examples

Nicholas Carlini<sup>1</sup> David Wagner<sup>2</sup>  
University of California, Berkeley

#### Abstract

MagNet and “Efficient Defenses...” were recently proposed as a defense to adversarial examples. We find that we can construct adversarial examples that defeat these defenses with only a slight increase in distortion.

#### 1. Introduction

It is an open question how to train neural networks so they will be robust to adversarial examples (Elie, 2016). Recently, these defenses have been proposed to make neural networks robust to adversarial examples:

- MagNet (H) was proposed as an approach to make neural networks robust against adversarial examples through two complementary approaches: adversarial examples near the data manifold are referred to as the data manifold set detector and rejected before classification. MagNet does not output robustness in the white-box setting, unless the authors suggest that MagNet is robust to the gray-box setting, where the adversary is aware of the distance to its class. However, the parameters of the base classifier, but not the parameters of the defense.

- An efficient defense (E) was proposed to make neural networks more robust against adversarial examples by performing cleaner data augmentation during training, and using the HMC-2 activation function. The authors do not claim perfect security, but claim this makes attacks visually identifiable.

- Adversarial Perturbation Elimination GAN (APE-GAN) (H) is similar to MagNet, only adversarial examples are projected onto the data manifold using a Generative Adversarial Network (GAN) (G). We did not set out to break this defense, but found it to be very similar to MagNet and we analyzed its in-

tegration. In this short paper, we demonstrate these defenses are not effective on the MNIST (H) and CIFAR-10 (E) datasets. We show that we are able to bypass MagNet with greater than 90% success, and the later two with 100%, with only a slight increase in distortion. We defeat MagNet by making use of the reweighted L1 (L1) property of adversarial examples: the adversary treats their concept of the defense, constructs adversarial examples on their model, and applies those adversarial examples to the defender. It turns out that these examples will also defeat the defender’s model. We defeat “Efficient Defenses Against Adversarial Attacks” and APE-GAN by showing that existing attack can defeat them with 100% success without modifications. Adversarial examples are not more visually dramatic than a standard random one.

#### 2. Background

We assume familiarity with neural networks (S), adversarial examples (H), transferability (T), generating strong attacks against adversarial examples (H) and MagNet (H). We briefly review the key details and notation.

**Notation.** Let  $F(x) = y$  be a neural network used for classification outputting a probability distribution. Call the second-to-last hidden layer before the softmax layer  $Z$ , and let  $F(Z) = softmax(Z)$ . Each output is compared to the intended probability that the object is labeled as class  $c$ . Let  $\mathcal{C}(x) = argmax_c F(x)_c$  correspond to the classification of an  $x$ . In this paper, we are concerned with neural networks used to classify images on MNIST and CIFAR-10.

**Adversarial examples (H)** are instances of  $x$  that are very close to a standard reference  $x_0$  subject to some distance metric  $\|x - x_0\|_p$  (in this paper, we have  $\mathcal{C}(x) \neq c$  for the target  $c$  chosen by the adversary). We generate adversarial examples with Carlini and Wagner’s  $L_0$ -attack algorithm (C). Specifically, we achieve

$$\max_{x'} \|x' - x_0\|_p = \epsilon(x')$$

1

- Formulation of adversarial attacks and defenses as optimizing surrogates of the true adversarial risk.

by the authors.

### Obscured Gradients Give a False Sense of Security: Commenting Defenses to Adversarial Examples

Anish Athalye<sup>1</sup> Nicholas Carlini<sup>1</sup> David Wagner<sup>2</sup>

#### Abstract

L1 gradients, a kind of gradient norm that leads to a false sense against adversarial lines that cause obscured loss. Iterative optimization-based methods cannot succeed. We identify three types of obscured gradients: *obscured gradients* are nonconvex or incorrect gradients caused either intentionally through non-differentiable operations or unintentionally through numerical instability; *stochastic gradients* depend on fast-time randomness; and *vanishing/exploding gradients* in very deep computation result in an unstable gradient.

apparent robustness against iterative optimization attacks: *obscured gradients*, a term we define as a special case of gradient masking (Shwartz et al., 2017). Without a good gradient, when following the gradient does not successfully optimize the loss, iterative optimization-based methods cannot succeed. We identify three types of obscured gradients: *obscured gradients* are nonconvex or incorrect gradients caused either intentionally through non-differentiable operations or unintentionally through numerical instability; *stochastic gradients* depend on fast-time randomness; and *vanishing/exploding gradients* in very deep computation result in an unstable gradient.

We propose new techniques to overcome obscured gradients caused by these three phenomena. We address gradient shading with a new attack technique we call Backward Post Differentiable Approximation, where we approximate derivatives by computing the forward pass normally and computing the backward pass using a differentiable approximation of the function. We compute gradients of randomized features by applying Expectation Over Transformation (Athalye et al., 2017). We solve vanishing/exploding gradients through parameterization and optimize over a space where gradients do not explode/vanish.

To investigate the prevalence of obscured gradients and circumvent of their applicability of these attack techniques, we use as a case study the ICLR 2018 non-certified defenses that claim white-box robustness. We find that obscured gradients are a common occurrence, with 7 of 9 defenses relying on this phenomenon. Applying the new attack techniques we develop, we overcome obscured gradients and circumvent of these completely, and 1 partially, under the original threat model of each paper. Along with this, we often have attacks that appear to be fully optimization-based attacks.

we why many defenses provide such attacks (Athalye et al., 2017). We solve vanishing/exploding gradients through parameterization and optimize over a space where gradients do not explode/vanish.

we why many defenses provide such attacks (Athalye et al., 2017). We solve vanishing/exploding gradients through parameterization and optimize over a space where gradients do not explode/vanish.

### Obscured Gradients Give a False Sense of Security: Commenting Defenses to Adversarial Examples

apparent robustness against iterative optimization attacks: *obscured gradients*, a term we define as a special case of gradient masking (Shwartz et al., 2017). Without a good gradient, when following the gradient does not successfully optimize the loss, iterative optimization-based methods cannot succeed. We identify three types of obscured gradients: *obscured gradients* are nonconvex or incorrect gradients caused either intentionally through non-differentiable operations or unintentionally through numerical instability; *stochastic gradients* depend on fast-time randomness; and *vanishing/exploding gradients* in very deep computation result in an unstable gradient.

We propose new techniques to overcome obscured gradients caused by these three phenomena. We address gradient shading with a new attack technique we call Backward Post Differentiable Approximation, where we approximate derivatives by computing the forward pass normally and computing the backward pass using a differentiable approximation of the function. We compute gradients of randomized features by applying Expectation Over Transformation (Athalye et al., 2017). We solve vanishing/exploding gradients through parameterization and optimize over a space where gradients do not explode/vanish.

To investigate the prevalence of obscured gradients and circumvent of their applicability of these attack techniques, we use as a case study the ICLR 2018 non-certified defenses that claim white-box robustness. We find that obscured gradients are a common occurrence, with 7 of 9 defenses relying on this phenomenon. Applying the new attack techniques we develop, we overcome obscured gradients and circumvent of these completely, and 1 partially, under the original threat model of each paper. Along with this, we often have attacks that appear to be fully optimization-based attacks.

we why many defenses provide such attacks (Athalye et al., 2017). We solve vanishing/exploding gradients through parameterization and optimize over a space where gradients do not explode/vanish.

We propose reproducible research, we release our reimplementation of each of these defenses, along with implementations of our attacks for each.<sup>1</sup>

<sup>1</sup><https://github.com/n-carlini/obscured-gradients>

University of California, Berkeley

University of California, Berkeley

# ... but designing defense mechanisms is hard!

## Despite the huge number of proposed defenses, state-of-the-art classifiers are still vulnerable to small perturbations.

### Adversarial Risk and the Dangle

Jonathan Uesato<sup>1</sup> Brendan O'Donoghue

#### Abstract

This paper investigates recently proposed approaches for defending against adversarial examples and evaluating adversarial robustness. We motivate *adversarial risk* as an objective for achieving models robust to worst-case inputs. We then frame commonly used attacks and evaluation metrics as defining a tractable surrogate objective to the true adversarial risk. This suggests that models may optimize this surrogate rather than the true adversarial risk. We formalize this notion as *obscure* as an adversary, and develop tools and heuristics for identifying obscured models and designing transparent models. We demonstrate that this is a significant problem in practice by reexpressing gradient-free optimization techniques into adversarial attacks, which we use to decrease the accuracy of several recently proposed defenses to near zero. Our hope is that our formalization and analysis will help researchers to develop more powerful defenses.

#### 1. Introduction

Deep learning has revolutionized the field of machine learning and led to substantial improvements in many challenging problems such as image understanding (He et al., 2016), speech recognition (Graves et al., 2013), and autonomous gaming (Mnih et al., 2015). Despite these remarkable successes, we have seen some intriguing and novel properties in the behaviour of these models. Researchers have demonstrated that certain small perturbations to the input can make neural networks produce entirely bad results (Szegedy et al., 2013; Liu and Li, 2017). For instance, in the case of image classification, perceptible perturbations can lead to the resulting image being misclassified to completely different object categories with high confidence (Szegedy et al., 2013; Thorelli and...

<sup>1</sup>DeepMind. Correspondence to Jonathan Uesato <uesato@google.com>.

Proceedings of the 35<sup>th</sup> International Conference on Machine Learning, Stockholm, Sweden, PMLR 80, 2018. Copyright 2018 by the author(s).

### MagNet and "Efficient Defenses Against Adversarial Attacks are Not Robust to Adversarial Examples"

Nicholas Carlini<sup>1</sup> David Wagner<sup>1</sup>  
University of California, Berkeley

#### Abstract

MagNet and "Efficient Defenses..." were recently proposed as a defense to adversarial examples. We find that we can construct adversarial examples that defeat these defenses with only a slight increase in distortion.

#### 1. Introduction

It is an open question how to train neural networks so they will be robust to adversarial examples (1). Recently, these defenses have been proposed to make neural networks robust to adversarial examples:

- MagNet (8) was proposed as an approach to make neural networks robust against adversarial examples through two complementary approaches: adversarial examples near the data manifold are referred to be on the data manifold that are classified correctly, whereas adversarial examples far away from the data manifold are denoised and rejected before classification. MagNet does not reject instances in the white-box setting, unlike the authors' paper that MagNet is robust to the grey-box setting where the adversary knows the defense in its place, knows the parameters of the base classifier, but not the parameters of the defense.

- An efficient defense (12) was proposed to make neural networks more robust against adversarial attacks by pretraining Gaussian data augmentation during training, and using the IMCCT activation function. The authors do not claim complete success, but state that it makes attacks visually indistinguishable.

- Adversarial Perturbation Minimization GAN (APM-GAN) (3) is similar to MagNet, only adversarial examples are generated on-the-fly using a Generative Adversarial Network (GAN) (2). We do not see a case to bypass this defense, but we intend to try.

In this short paper, we demonstrate these defenses are not effective on the MNIST (5) and CIFAR-100 datasets. We show that we are able to hit the Net with greater than 90% success, and do this with 100% with only a slight increase in distortion.

We defeat MagNet by making use of the re-ly (11) property of adversarial examples: the basic four concepts of the defense, consistent all examples on their model, and applies to small examples to the defense. 3 items of examples will also fool the defender's model. We defeat "Efficient Defenses Against Adversarial Attacks" (APM-GAN) by showing that creating adversarial examples with 100% success without the adversarial network are not more visually than an adversarial network.

#### 2. Background

We assume familiarity with neural network terminology (see (11), specifically (14)), strong attacks against adversarial examples (1) [1]. We briefly review the key details and...

**Notation.** Let  $P(x)$  be a neural network classification output a probability distribution over the second-to-last layer (the layer before the layer  $z$ ), so that  $P(z) = \text{softmax}(P(x))$ . Each component in the probability distribution that  $x$  is labeled as class  $c$ . Let  $C(x)$  be a vector composed of the classification of  $x$  at  $P$ . As we are concerned with neural networks used against MNIST and CIFAR-100.

**Adversarial examples (11)** are instances a very close to a normal instance with respect distance metric  $\|\cdot\|$ . In this paper, we follow the classification of  $x$  from the case of the classification of  $x$  as  $\text{Instance}$ . In this paper, we follow the classification of  $x$  from the case of the classification of  $x$  as  $\text{Instance}$ . In this paper, we follow the classification of  $x$  from the case of the classification of  $x$  as  $\text{Instance}$ .

We generate adversarial examples with C-Wagner's (6) attack algorithm (1). Specifically, we minimize  $\|x' - x\|_2^2 + \lambda \|x'\|_2^2$ .

$$\min_{x'} \|x' - x\|_2^2 + \lambda \|x'\|_2^2$$

arXiv:1804.03286v1 [cs.LG] 10 Apr 2018

### On the Robustness of the CVPR 2018 White-Box Adversarial Example Defenses

Ansh Athey<sup>1</sup> Nicholas Carlini<sup>2</sup>

#### Abstract

Neural networks are known to be vulnerable to adversarial examples. In this note, we evaluate the five white-box defenses that appeared at CVPR 2018 and find that the majority, when applying existing defenses, we evaluate the accuracy of the defended models to be:

#### 1. Introduction

Training neural networks so they will be robust to adversarial examples (strongly (1), 2013) is a major challenge. Two defenses that appeared at CVPR 2018 attempt to address this problem: "On-Demand Adversarial Example with Plug Defenses" (Prabhakar et al., 2018) and "Defense against Adversarial Attacks Using High-Level Representation Guided Denoiser" (Liao et al., 2018).

In this note, we show these two defenses are not effective in the white-box threat model. We construct adversarial examples that make the classifier accuracy  $\approx 0\%$  on the ImageNet dataset (Deng et al., 2009) when bounded by a small  $L_2$  perturbation of  $4\sqrt{255}$ , a value bound that is smaller than the original papers. Our attacks can construct targeted adversarial examples with over 97% success. Our methods are a direct application of existing techniques.

#### 2. Background

We assume familiarity with neural networks, adversarial examples (Szegedy et al., 2013), generating strong attacks against adversarial examples (Miyata et al., 2016), and comparing adversarial examples for neural networks with non-differentiable layers (Doherty et al., 2016). We briefly review the key details and notation.

**Adversarial examples (Szegedy et al., 2013)** are instances of data points that are similar to a normal instance, but whose classification of  $x$  from the case of the classification of  $x$  as  $\text{Instance}$ . In this paper, we follow the classification of  $x$  from the case of the classification of  $x$  as  $\text{Instance}$ . In this paper, we follow the classification of  $x$  from the case of the classification of  $x$  as  $\text{Instance}$ .

neural network on Machine Learning, PMLR 80, 2018. Copyright 2018 by the author(s).

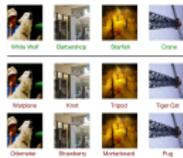


Figure 1: Original images from ImageNet validation set, its targeted adversarial examples with randomly chosen targets. In Plug Defenses (row 2) and High-Level Representation Guided Denoiser (row 3), with  $L_2$  perturbation of  $\epsilon = 4\sqrt{255}$ .

Targeted adversarial examples are instances  $x'$  whose label is equal to a adversarial target label.

We construct two defenses: Plug Defenses and High-Level Representation Guided Denoiser. We are general to the authors of these defenses for releasing their source code and pre-trained models.

**Plug Defenses (Prabhakar et al., 2018)** propose a non-differentiable preprocessing of inputs. Some grids in table 1 represent an adversarial example with a small  $L_2$  perturbation. This resulting image is often noisy, and its success accuracy, a desirable objective, is degraded.

**High-Level representation Guided Denoiser (Liao et al., 2018)** propose denoising input using a neural network before passing them to a standard classifier. This shows to be differentiable, non-complex neural network. This defense has also been evaluated by Uesato et al. (2018) and found to be ineffective.

#### 2.1. Methods

We evaluate these defenses under the white-box threat model. We construct adversarial examples with PlugDef

mechanisms of our attacks for each.

<sup>1</sup>https://github.com/ankr/white-box-adversarial-examples

- Formulation of adversarial attacks and defenses as optimizing surrogate of the true adversarial risk

# ... but designing defense mechanisms is hard!

Despite the huge number of proposed defenses, state-of-the-art classifiers are still vulnerable to small perturbations.

**Adversarial Risk and the Danger of Defenses**  
Jonathan Uesato<sup>1</sup>, Brendan O'Donoghue<sup>1</sup>

**Abstract**  
This paper investigates recently proposed defenses for defending against adversarial perturbations and evaluating adversarial robustness. We motivate adversarial risk as an objective for training models robust to worst-case perturbations and show that this metric is more robust to the true adversarial model than the true adversarial model may optimize. We show that adversarial risk is as computationally tractable as robustness to  $\ell_\infty$  perturbations and demonstrate that this metric is more robust to the true adversarial model than the true adversarial model may optimize. We show that adversarial risk is as computationally tractable as robustness to  $\ell_\infty$  perturbations and demonstrate that this metric is more robust to the true adversarial model than the true adversarial model may optimize.

**1. Introduction**  
Deep learning has revolutionized the field of machine learning and led to substantial improvements in many challenging problems such as image understanding (He et al., 2015), speech recognition (Graves et al., 2013), and autonomous driving (Mnih et al., 2015). Despite these remarkable successes, we have seen some intriguing and troubling properties in the behavior of these models. Researchers have demonstrated that certain small perturbations to the input can make neural networks produce entirely bad results (Szegedy et al., 2013; Liu and Li, 2017). For instance, in the case of image classification, perceptible perturbations can lead to the resulting image being misclassified to completely different object categories with high confidence (Szegedy et al., 2013; Thorelli et al., 2017).

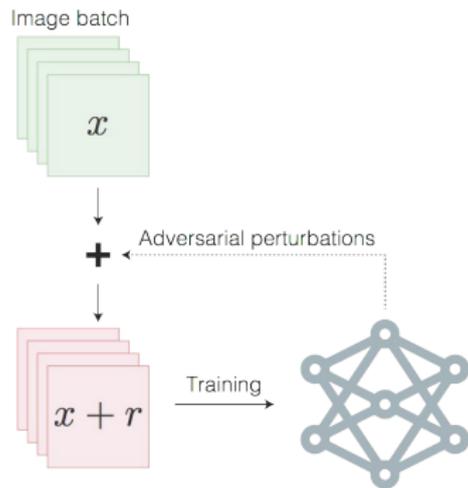
**Adversarial Examples Are Not Easily Detected: Bypassing Ten Detection Methods**  
David Winer, Nicholas Carlini, University of California, Berkeley

**Abstract**  
Adversarial examples are perturbations to legitimate inputs that cause a model to misclassify. While these perturbations are often small and imperceptible, they can be used to cause a model to misclassify with high confidence. In this paper, we propose a new method for generating adversarial examples that is able to bypass ten different detection methods. Our method is based on a novel optimization technique that allows us to generate adversarial examples that are more robust to detection than previous methods. We demonstrate that our method is able to bypass ten different detection methods, including gradient-based methods, statistical methods, and heuristic methods. Our method is able to generate adversarial examples that are more robust to detection than previous methods, and we demonstrate that our method is able to bypass ten different detection methods.

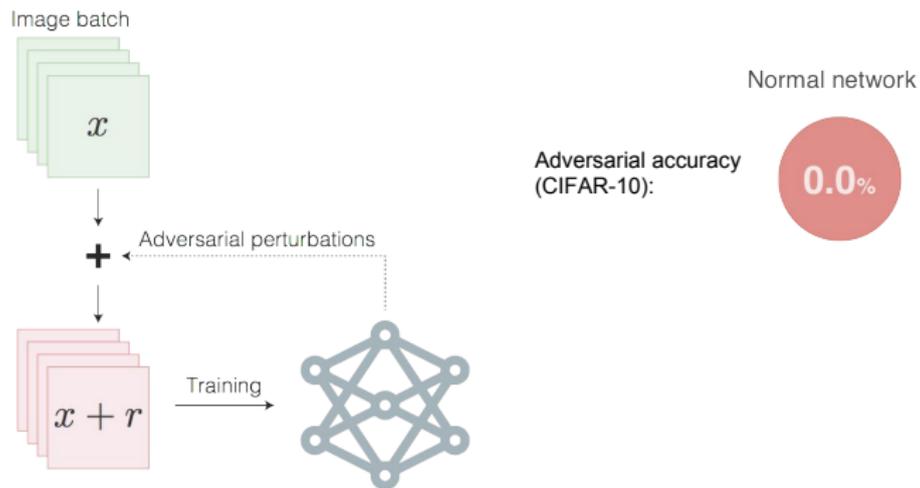
**Formulating adversarial attacks as a minimizing surrogate of the true adversary**  
David Winer, Nicholas Carlini, University of California, Berkeley

**Abstract**  
Adversarial attacks are a central problem in machine learning. In this paper, we propose a new method for generating adversarial examples that is based on a novel optimization technique. Our method is able to generate adversarial examples that are more robust to detection than previous methods, and we demonstrate that our method is able to bypass ten different detection methods.

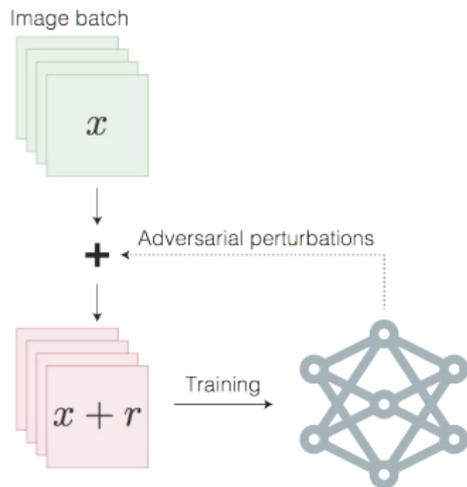
# Adversarial training



# Adversarial training



# Adversarial training



Normal network

Adversarial training

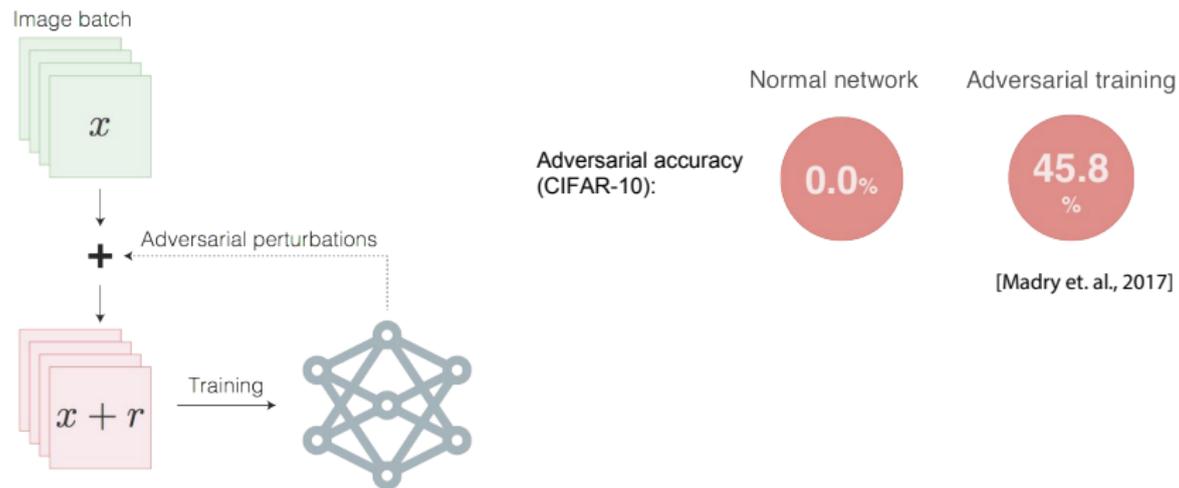
Adversarial accuracy  
(CIFAR-10):

0.0%

45.8  
%

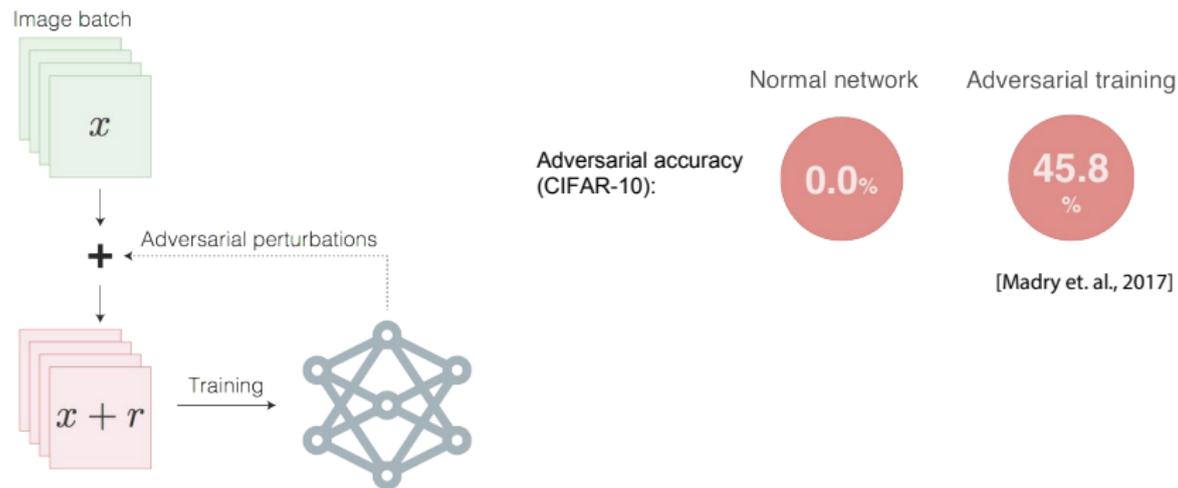
[Madry et. al., 2017]

# Adversarial training



Adversarial training leads to state-of-the-art robustness to adversarial perturbations.

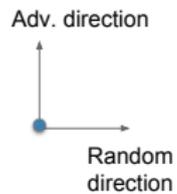
# Adversarial training



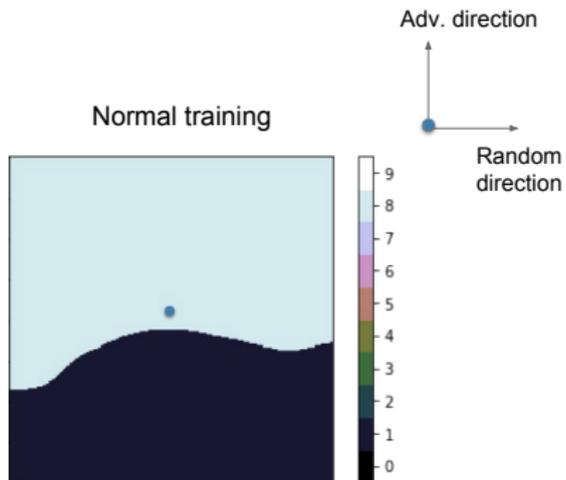
Adversarial training leads to state-of-the-art robustness to adversarial perturbations.

**But what does it actually do?**

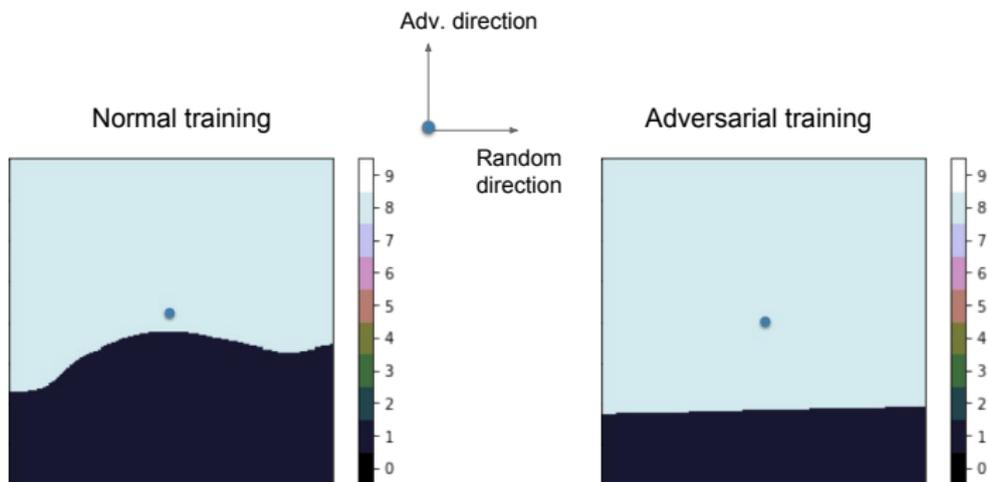
# Decision boundaries



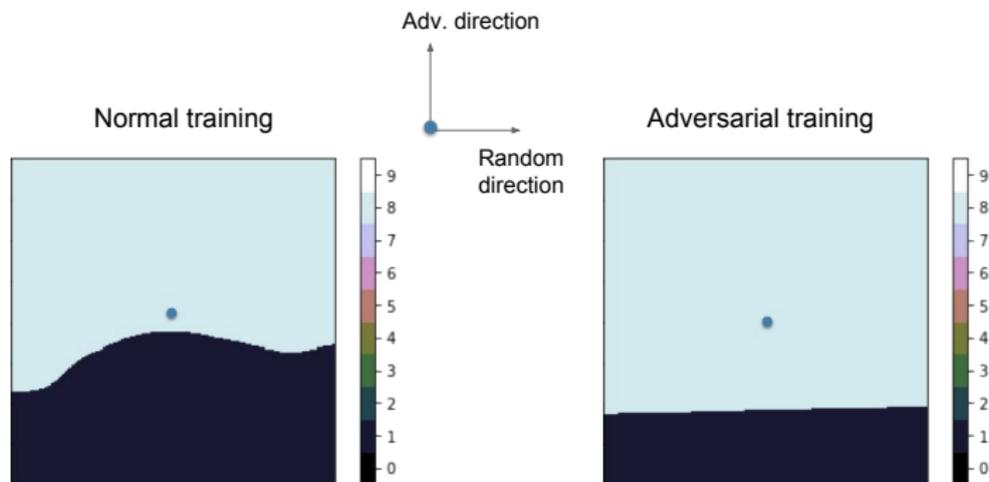
# Decision boundaries



# Decision boundaries



# Decision boundaries



After adversarial training, the decision boundaries are flatter and more regular.

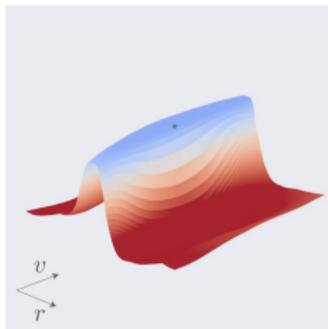
## Effect of adversarial training on loss landscape

$$\ell(x) = \text{CE}(\underbrace{f_{\theta}(x)}_{\text{Logit}}, \underbrace{y}_{\text{Label}})$$

# Effect of adversarial training on loss landscape

$$\ell(x) = \text{CE}(\underbrace{f_{\theta}(x)}_{\text{Logit}}, \underbrace{y}_{\text{Label}})$$

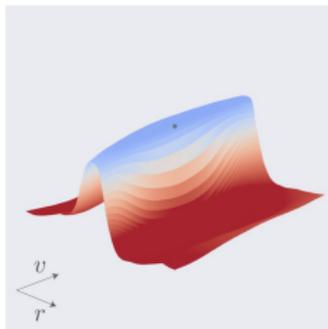
Before adv. fine-tuning



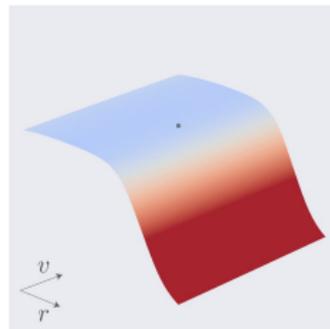
# Effect of adversarial training on loss landscape

$$\ell(x) = \text{CE}(\underbrace{f_{\theta}(x)}_{\text{Logit}}, \underbrace{y}_{\text{Label}})$$

Before adv. fine-tuning



After adv. fine-tuning



## Effect of adversarial training on loss landscape (Cont'd)

## Quantitative analysis: curvature decrease with adversarial training

We compute the Hessian matrix at a test point  $x$  with respect **to inputs**.

$$H = \left( \frac{\partial^2 \ell}{\partial x_i \partial x_j} \right)$$

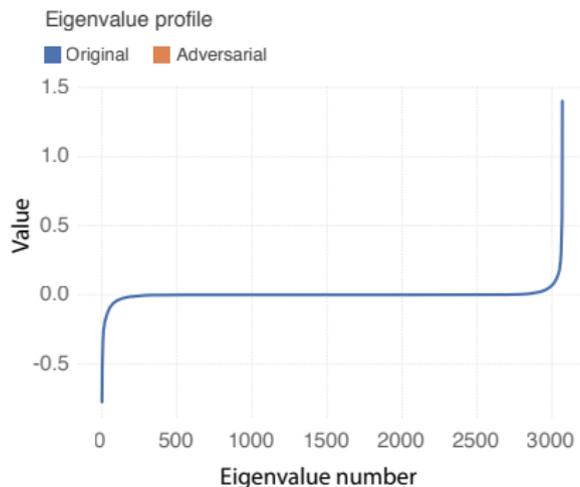
The eigenvalues of  $H$  are the curvature of  $\ell$  in the vicinity of  $x$ .

## Quantitative analysis: curvature decrease with adversarial training

We compute the Hessian matrix at a test point  $x$  with respect **to inputs**.

$$H = \left( \frac{\partial^2 \ell}{\partial x_i \partial x_j} \right)$$

The eigenvalues of  $H$  are the curvature of  $\ell$  in the vicinity of  $x$ .

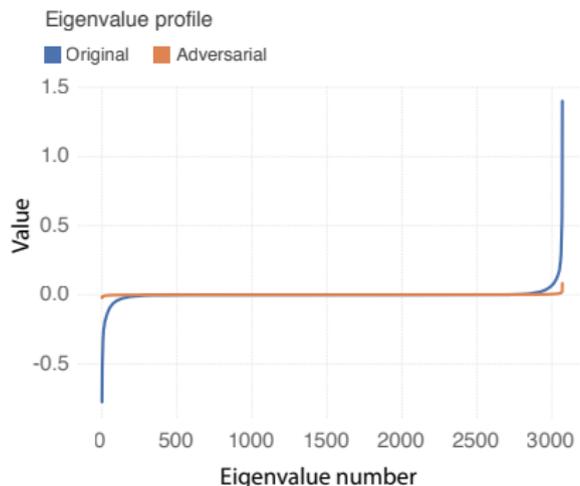


## Quantitative analysis: curvature decrease with adversarial training

We compute the Hessian matrix at a test point  $x$  with respect to **inputs**.

$$H = \left( \frac{\partial^2 \ell}{\partial x_i \partial x_j} \right)$$

The eigenvalues of  $H$  are the curvature of  $\ell$  in the vicinity of  $x$ .



## Relation between curvature and robustness

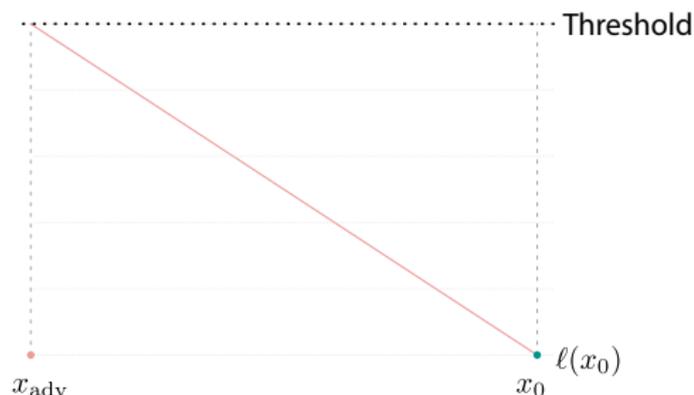
- Locally quadratic approximation of the loss function

## Relation between curvature and robustness

- Locally quadratic approximation of the loss function
- We derive upper and lower bounds on the minimal perturbation required to fool a classifier.

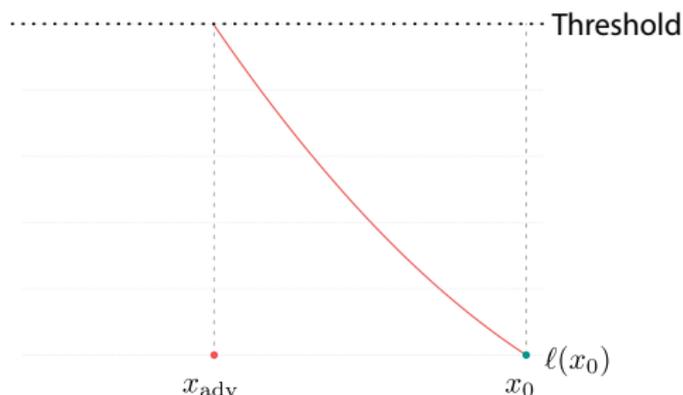
## Relation between curvature and robustness

- Locally quadratic approximation of the loss function
- We derive upper and lower bounds on the minimal perturbation required to fool a classifier.



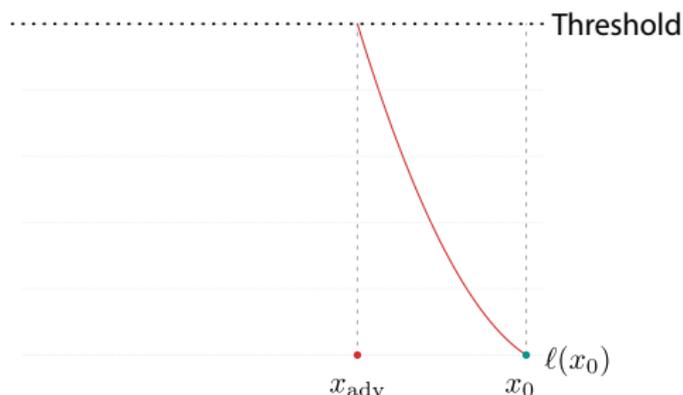
## Relation between curvature and robustness

- Locally quadratic approximation of the loss function
- We derive upper and lower bounds on the minimal perturbation required to fool a classifier.



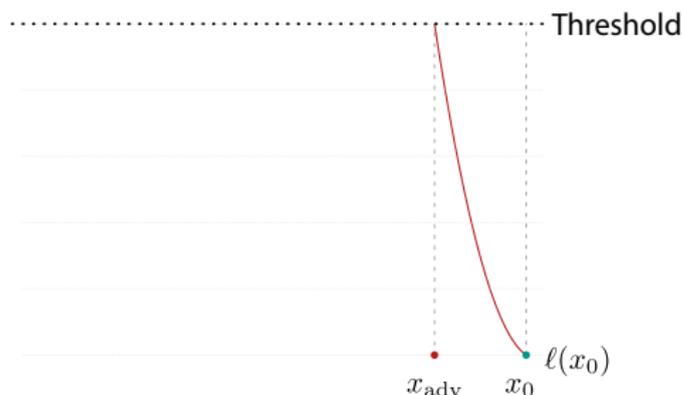
## Relation between curvature and robustness

- Locally quadratic approximation of the loss function
- We derive upper and lower bounds on the minimal perturbation required to fool a classifier.



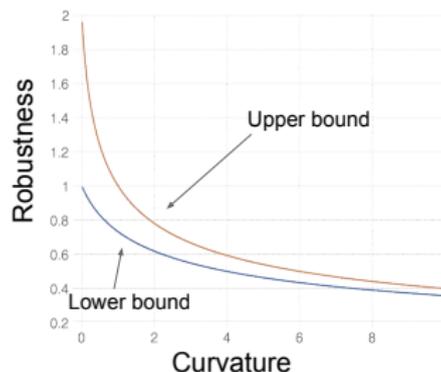
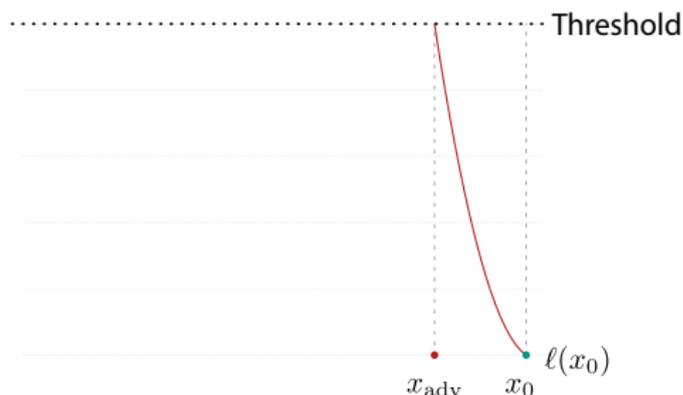
## Relation between curvature and robustness

- Locally quadratic approximation of the loss function
- We derive upper and lower bounds on the minimal perturbation required to fool a classifier.



## Relation between curvature and robustness

- Locally quadratic approximation of the loss function
- We derive upper and lower bounds on the minimal perturbation required to fool a classifier.



## How important is the curvature decrease?

Is the curvature decrease **the main effect of adversarial training** leading to improved robustness?

## How important is the curvature decrease?

Is the curvature decrease **the main effect of adversarial training** leading to improved robustness?

→ we regularize explicitly for the curvature.

## How important is the curvature decrease?

Is the curvature decrease **the main effect of adversarial training** leading to improved robustness?

→ we regularize explicitly for the curvature.

- Idea: Regularize the norm of the Hessian of the loss wrt inputs.

## How important is the curvature decrease?

Is the curvature decrease **the main effect of adversarial training** leading to improved robustness?

→ we regularize explicitly for the curvature.

- Idea: Regularize the norm of the Hessian of the loss wrt inputs.
- Use Hutchinson's estimator

$$\|H\|_F = \sqrt{\mathbb{E}_{z \sim \mathcal{N}(0, I)} \|Hz\|_2^2}$$

## How important is the curvature decrease?

Is the curvature decrease **the main effect of adversarial training** leading to improved robustness?

→ we regularize explicitly for the curvature.

- Idea: Regularize the norm of the Hessian of the loss wrt inputs.
- Use Hutchinson's estimator

$$\|H\|_F = \sqrt{\mathbb{E}_{z \sim \mathcal{N}(0, I)} \|Hz\|_2^2}$$

- In practice:
  - Compute Hessian-vector products with finite difference
  - Selective sampling on directions corresponding to high curvature

## How important is the curvature decrease?

Is the curvature decrease **the main effect of adversarial training** leading to improved robustness?

→ we regularize explicitly for the curvature.

- Idea: Regularize the norm of the Hessian of the loss wrt inputs.
- Use Hutchinson's estimator

$$\|H\|_F = \sqrt{\mathbb{E}_{z \sim \mathcal{N}(0, I)} \|Hz\|_2^2}$$

- In practice:
  - Compute Hessian-vector products with finite difference
  - Selective sampling on directions corresponding to high curvature

CURE: Regularize using  $\ell_r = \|\nabla \ell(x + hz) - \nabla \ell(x)\|$

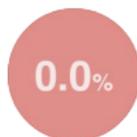
# CUREing deep networks trained on CIFAR-10

Normal network

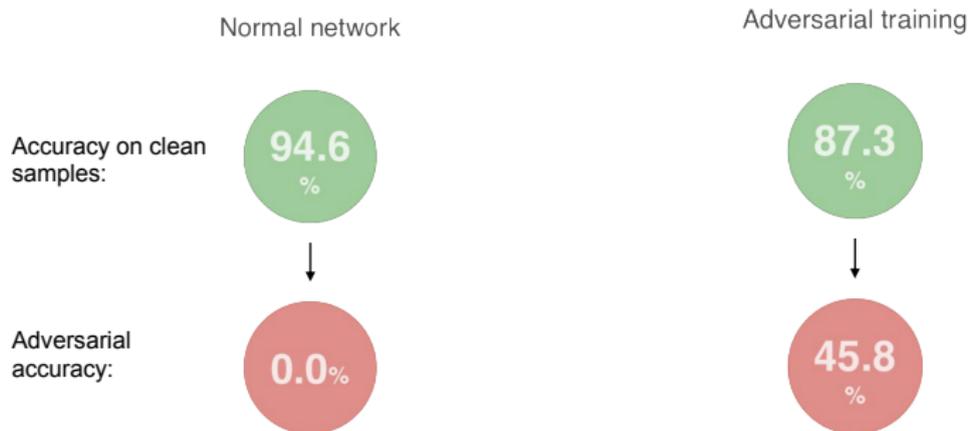
Accuracy on clean samples:



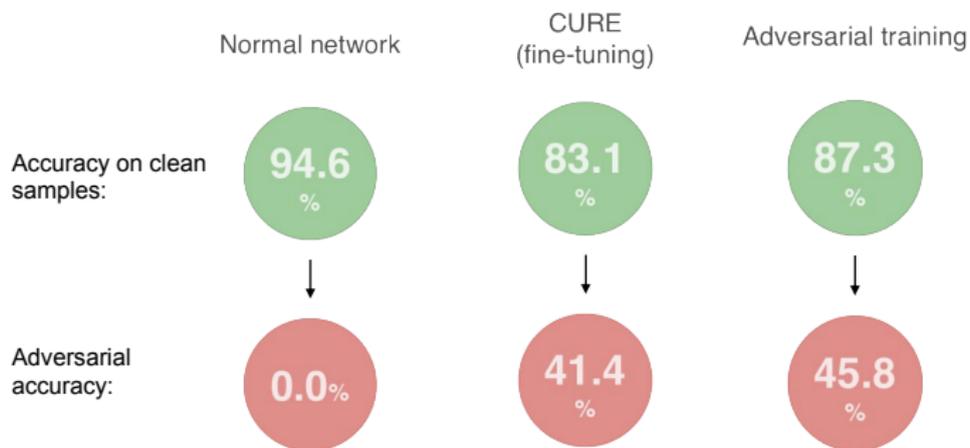
Adversarial accuracy:



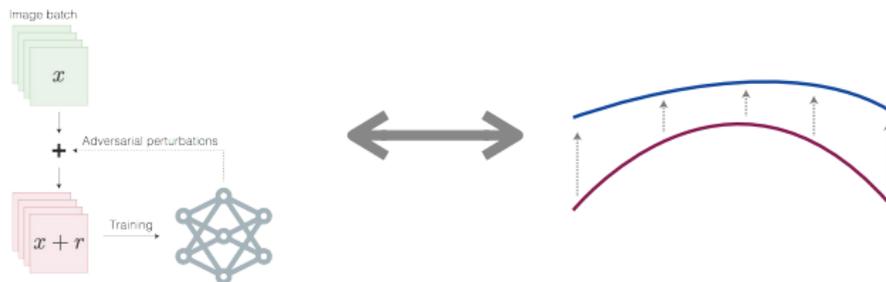
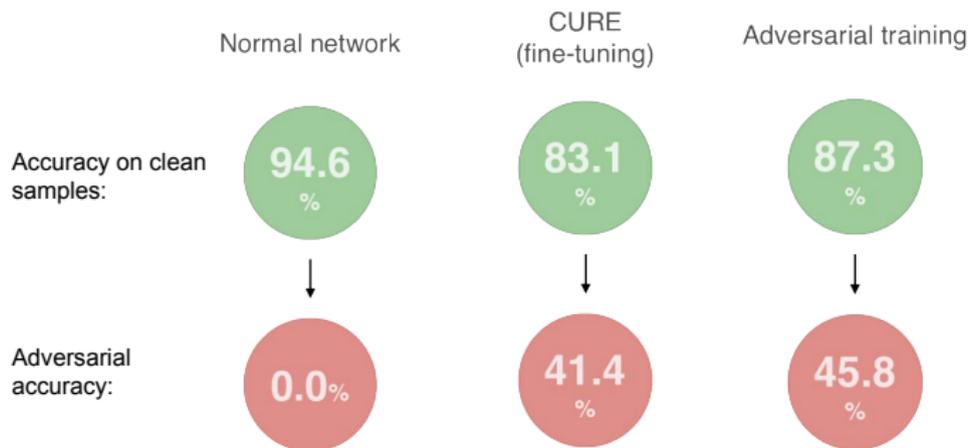
# CUREing deep networks trained on CIFAR-10



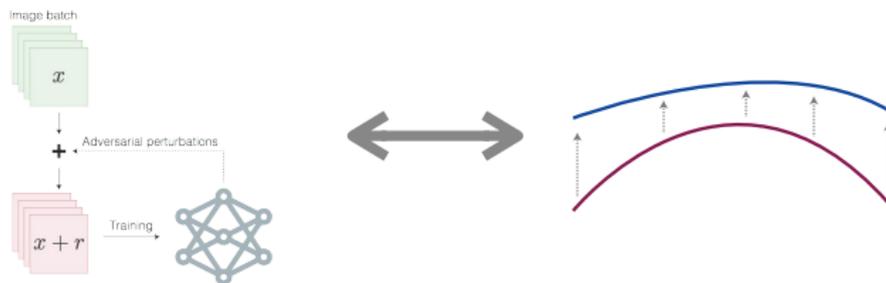
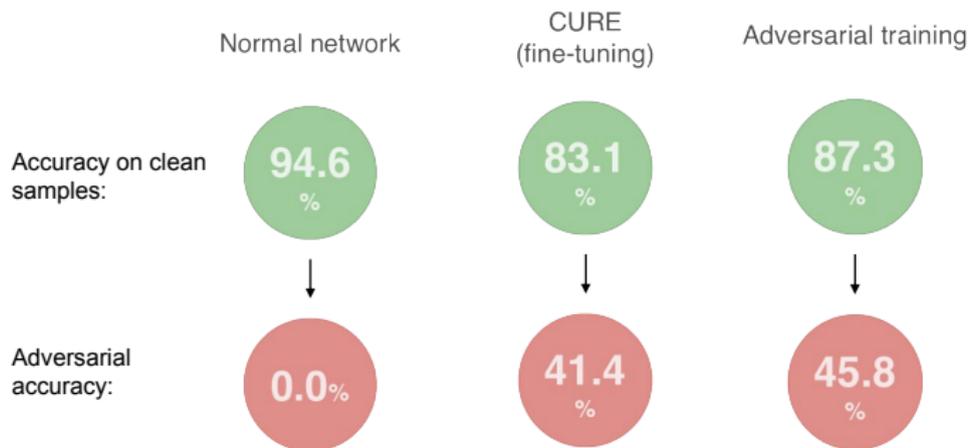
# CUREing deep networks trained on CIFAR-10



# CUREing deep networks trained on CIFAR-10



# CUREing deep networks trained on CIFAR-10



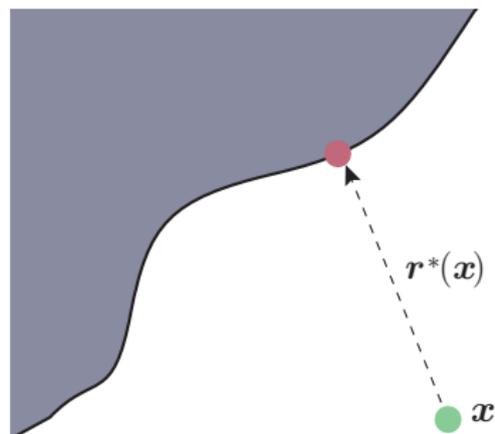
[Moosavi-Dezfooli, et. al., *Robustness via curvature regularization, and vice-versa*, CVPR 2019]

# Upper limits on adversarial robustness

- Goal: examine the existence of upper bounds on the robustness to adversarial perturbations.
- Relate to quantities that we better understand/we can better measure: *robustness to random noise*.
- Comparison to the robustness to random noise quantifies the power of an adversary having access to the model vs. no clue about the model.

# From adversarial to random noise

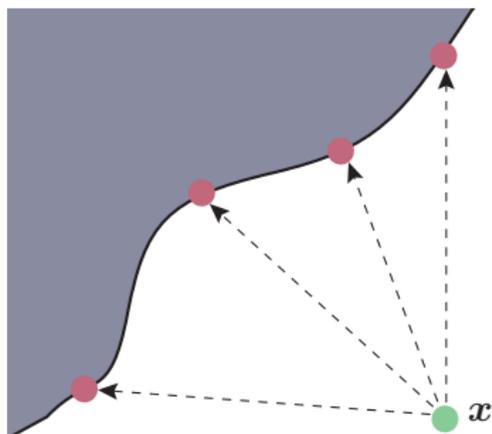
Robustness to adversarial noise



$$\min_r \|r\|_2 \text{ subject to } f(x+r) \neq f(x).$$

# From adversarial to random noise

Robustness to **random** noise



$$\min_t |t| \text{ subject to } f(x + t\mathbf{v}) \neq f(x)$$

$\mathbf{v}$  uniformly sampled from  $\mathbb{S}^{D-1}$ .

## Linear classifiers

Theorem (Fawzi et. al., NIPS '16, Franceschi et. al., AISTATS '18)

*For affine classifiers, we have*

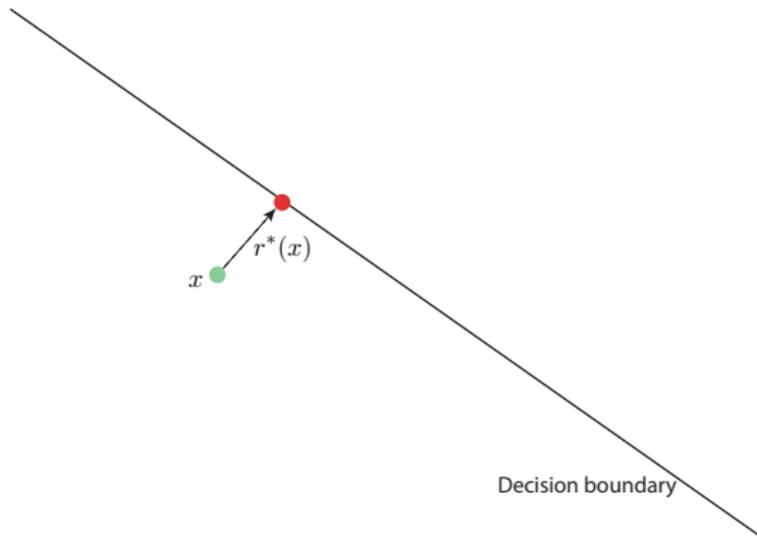
$$\|r^*\|_2 = \Theta \left( \frac{1}{\sqrt{D}} \|r_{rand}^*\|_2 \right),$$

*with high probability (over the choice of random perturbation).*

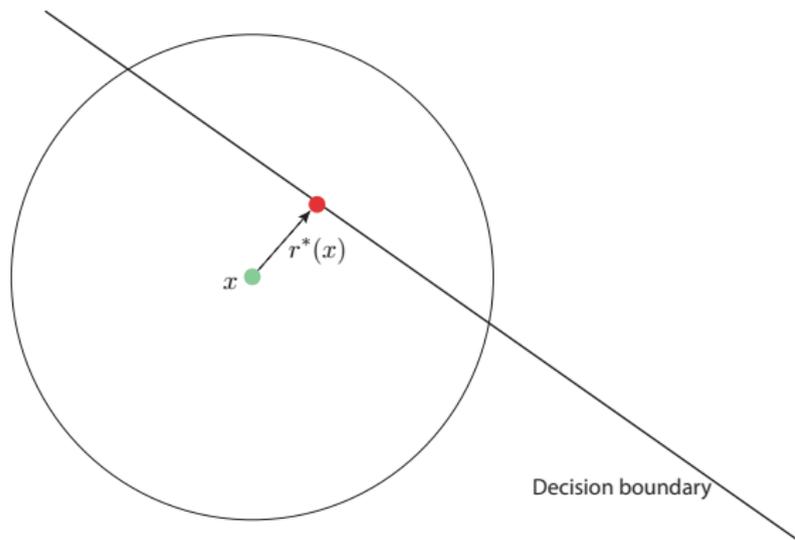
In high dimensions, very large gap between both robustness measures.

When data is bounded, we therefore typically get  $\|r^*\|_2 = O\left(\frac{1}{\sqrt{D}}\right) \rightarrow$   
Achieving robustness in high dimensions is difficult!

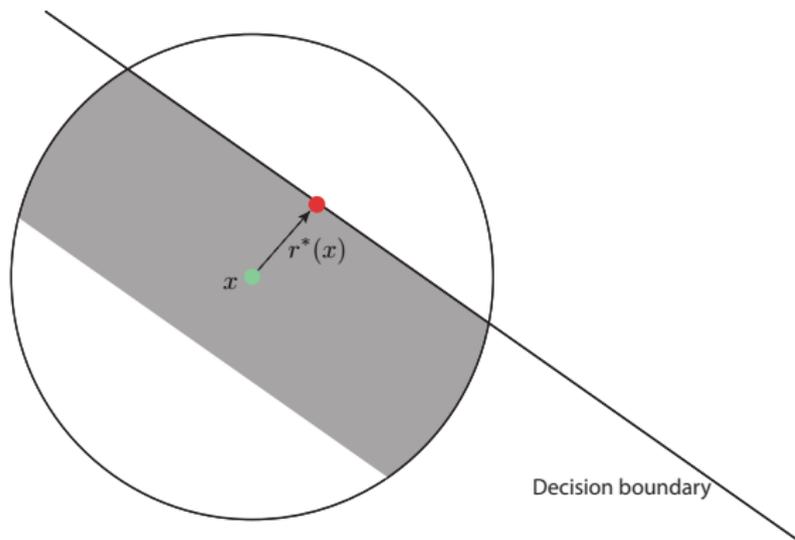
# Intuition



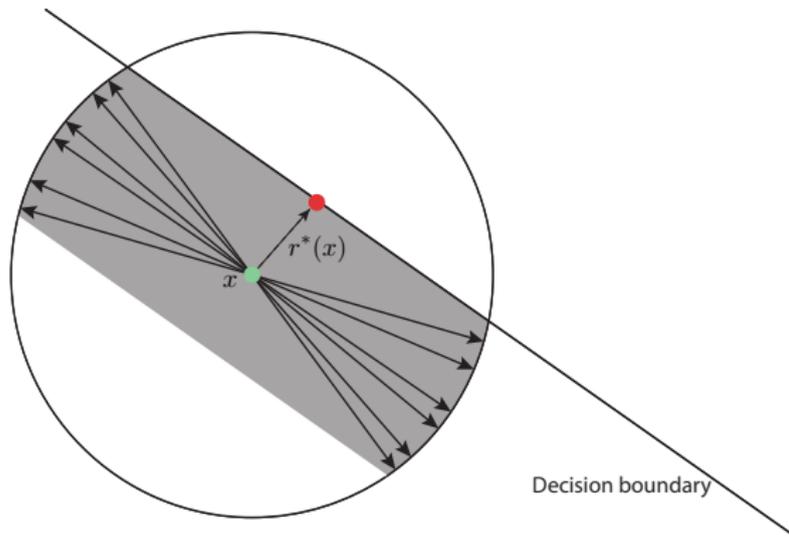
# Intuition



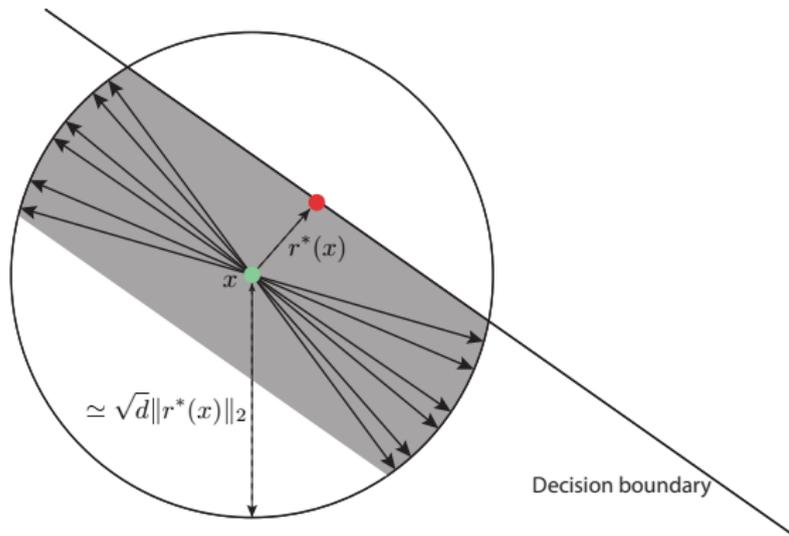
# Intuition



# Intuition

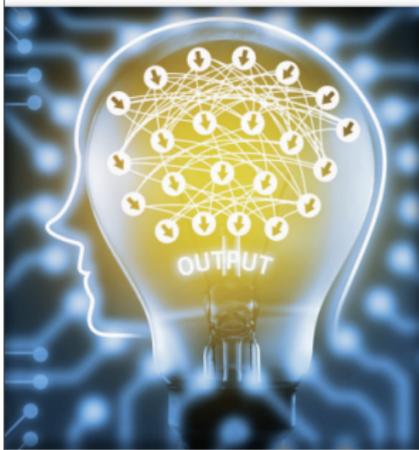


# Intuition



## The Robustness of Deep Networks

*A geometrical perspective*



Deep neural networks have recently shown impressive classification performance on a diverse set of visual tasks. When deployed in real-world (noise-prone) environments, it is equally important that these classifiers satisfy robustness guarantees: small perturbations applied to the samples should not yield significant loss to the performance of the predictor. The goal of this article is to discuss the robustness of deep networks to a diverse set of perturbations that may affect the samples in practice, including adversarial perturbations, random noise, and geometric transformations. This article further discusses the recent works that build on the robustness analysis to provide geometric insights on the classifier's decision surface, which help in developing a better understanding of deep networks. Finally, we present recent solutions that attempt to increase the robustness of deep networks. We hope this review article will contribute to shed ding light on the open research challenges in the robustness of deep networks and stir interest in the analysis of their fundamental properties.

### Introduction

With the dramatic increase of digital data and the development of new computing architectures, deep learning has been developing rapidly as a predominant framework for data representation that can contribute in solving very diverse tasks. Despite this success, several fundamental properties of deep neural networks are still not understood and have been the subject of intense analysis in recent years. In particular, the robustness of deep networks to various forms of perturbations has

# Conclusions

- Very simple strategies to fool classifiers: adversarial, geometric, universal perturbations, ...
- Difficulty of defending against adversarial perturbations; most “defenses” make the classifier robust against specific directions.
- Adversarial training is essentially reducing the curvature of the loss function, leading to an increased robustness.

## References

- Szegedy et. al., *Intriguing properties of neural networks*, ICLR 2014
- Fawzi et. al, *Manitest: Are classifiers really invariant?*, BMVC 2015
- Moosavi-Dezfooli et. al., *Deepfool: a simple and accurate approach to fool deep neural networks*, CVPR 2016
- **Fawzi et. al., Robustness of classifiers: from adversarial to random noise, NIPS 2016**
- Moosavi-Dezfooli et. al., *Universal adversarial perturbations*, CVPR 2017
- Uesato et al., *Adversarial risk and the dangers of evaluating against weak attacks*, arXiv 2018
- Fawzi et. al., *Adversarial vulnerability of any classifier*, NeurIPS 2018
- **Moosavi-Dezfooli et. al., Robustness via curvature regularization, and vice-versa, CVPR 2019**
- **Survey:** [Fawzi et. al., *Robustness of deep networks: a geometric perspective*, IEEE Signal Processing Magazine 2017]