

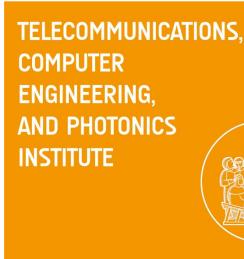
Adversarial Attacks

Deep Learning and Neural Networks: Advanced Topics

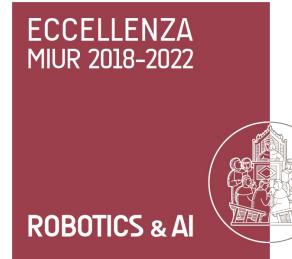
Fabio Brau

March 21, 2022

SSSA, Emerging Digital Technologies, Pisa.



Sant'Anna
School of Advanced Studies – Pisa



Sant'Anna
Scuola Universitaria Superiore Pisa



Introduction, Definitions and Taxonomy

Adversarial Examples Generation

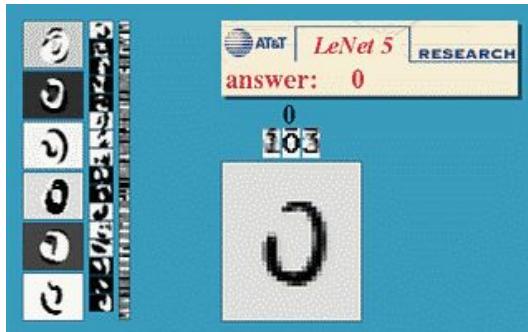
Defense Mechanisms

Conclusions



Neural Networks achieve Human Performances...

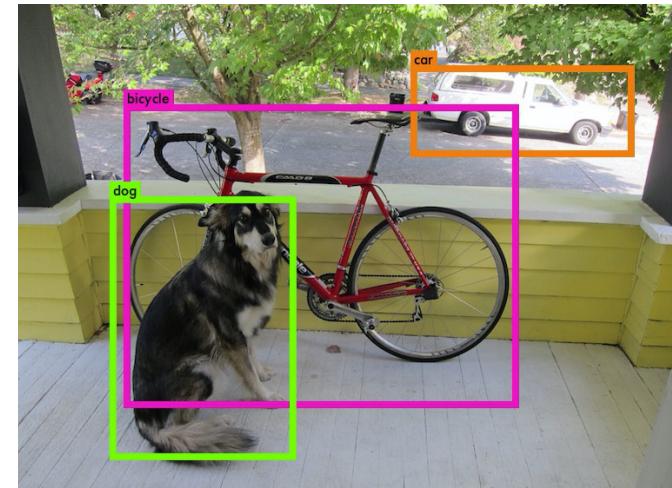
LeNet-5 (MNIST), 1998



AlexNet (ImagNet), 2013



Yolo, 2016



Neural Networks become so powerful that they reach human capabilities in the classification and object detection tasks

... but they make weird mistakes.

Evading Spam Detection (2010)

FIRECRACKER STOCK ALERT!!!
MONDAY, JULY 3, 2006
Waypoint Biomedical, Inc
Symbol: WYPH
Price: \$ 17

Sometimes these little stocks can make gains in days or weeks
accomplish, if at all! Know what we mean? You be the judge!

RECENT HEALINES:



FIRECRACKER STOCK ALCRT!!!
MONDAY, JULY 3, 2006
Waypoint Biomedical, Inc
Symbol: WYPH
Price. \$ 17

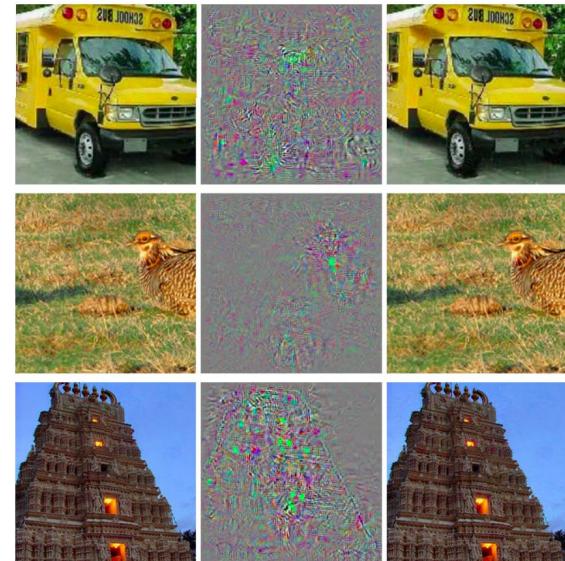
Sometimes these little stocks can make gains in days or weeks
accomplish, if at all! Know what we mean? You be the judge!

RECENT HEALINES:



Spam Detector based on OCR

Attacks on Images (2014)



Perturbed Images are no longer
properly recognized by the network

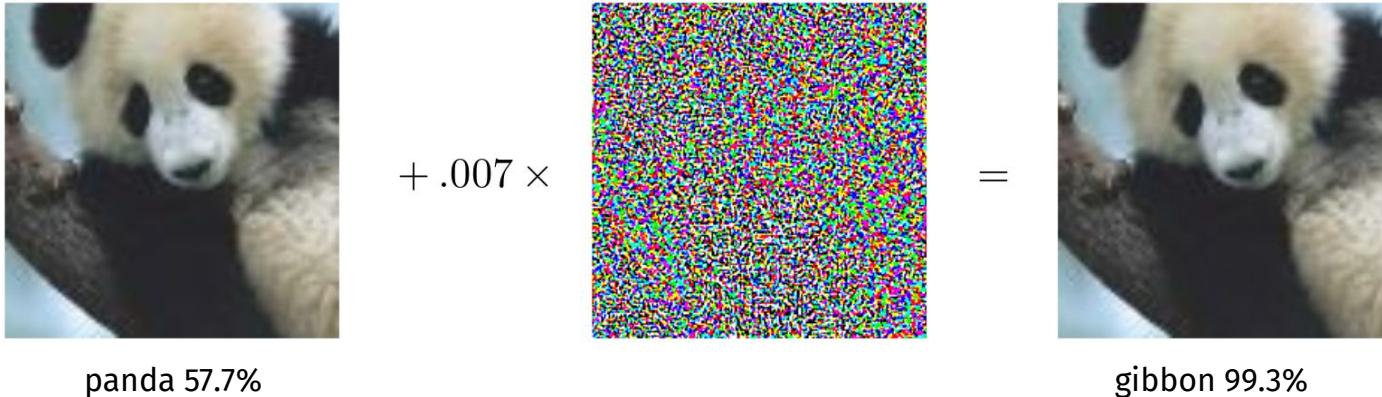
Delvi et al. "Adversarial Classification" (2004)

Battista Biggio et al. "Adversarial Pattern Classification" (2010)

Christian Szegedy et al. "Intriguing properties of neural networks" (2013)



Adversarial Examples



Adversarial examples are inputs **intentionally** designed to cause ML models to make **errors**.

They can be created so that can appear **nearly identical** to normal inputs.

Adversarial Examples

Speech-to-Text models

Can be fooled by adding an imperceptible noise to input audio



+



=



*“without the dataset the
article is useless”*

Adversarial Noise (30dB)



*“okay google browse to
evil dot com”*

Adversarial Examples become popular

WRITING

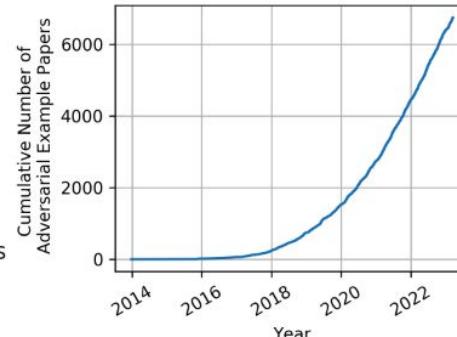
MAIN PAPERS TALKS CODE WRITING

A Complete List of All (arXiv) Adversarial Example Papers

by Nicholas Carlini 2019-06-15

It can be hard to stay up-to-date on the published papers in the field of adversarial examples, where we have seen massive growth in the number of papers written each year. I have been somewhat religiously keeping track of these papers for the last few years, and realized it may be helpful for others to release this list.

The only requirement I used for selecting papers for this list is that it is primarily a paper about adversarial examples, or extensively uses adversarial examples. Due to the sheer quantity of papers, I can't guarantee that I actually have found *all* of them.



link to the website

Why are adversarial Examples worth studying?

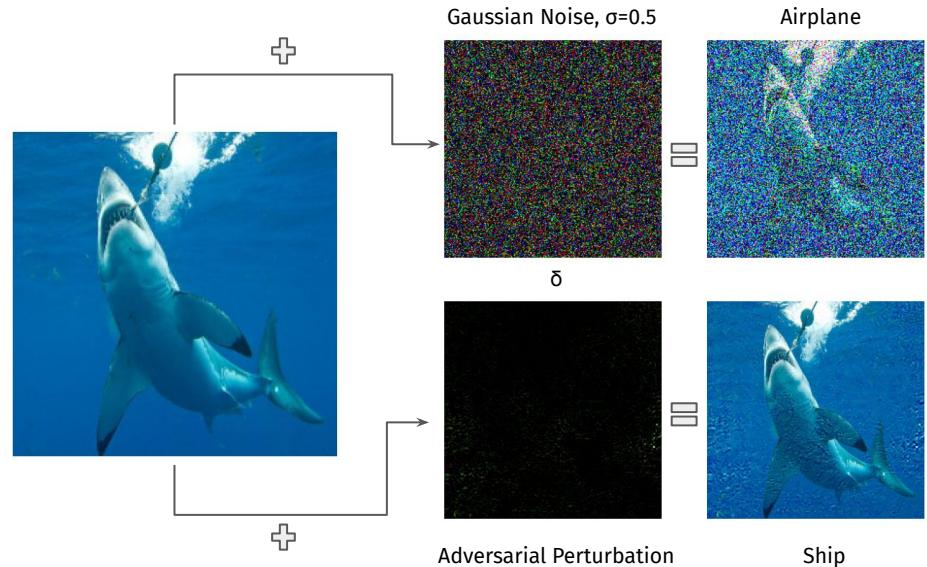
Security Threat (Known your Enemy)

Neural Networks are involved in a growing number of tasks. How can we tell safe?



DeepExplore (2021). Erroneous behaviour of DAVE-2. A tiny change in the light condition causes the model turns right.

Foundational Perspective

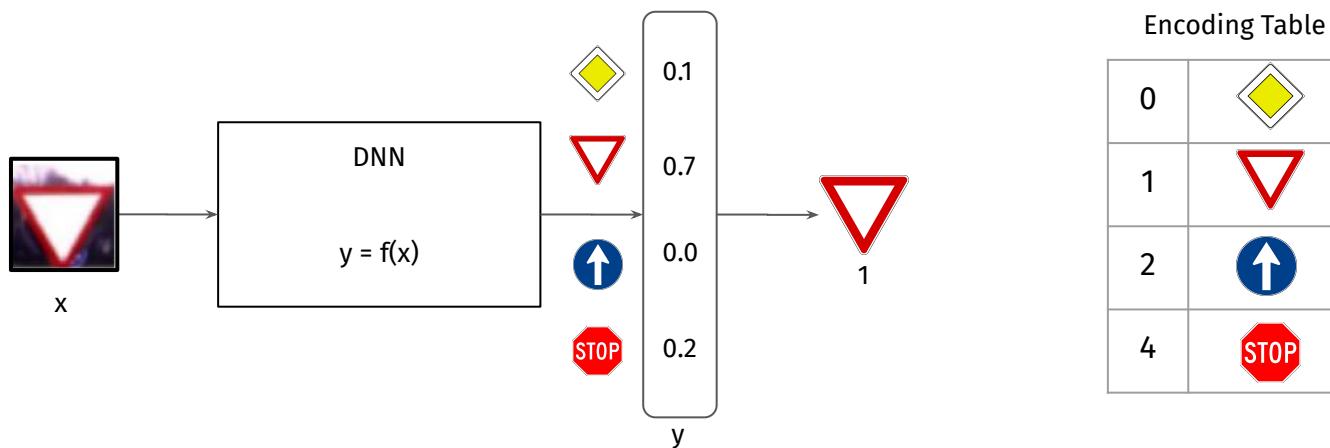
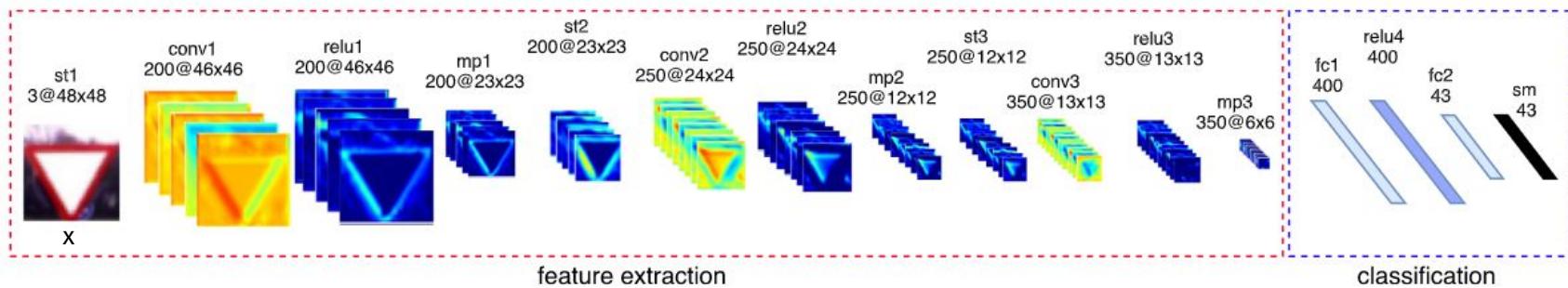


Small maliciously perturbations are sufficient to fool the model, why gaussian noise has not this property?

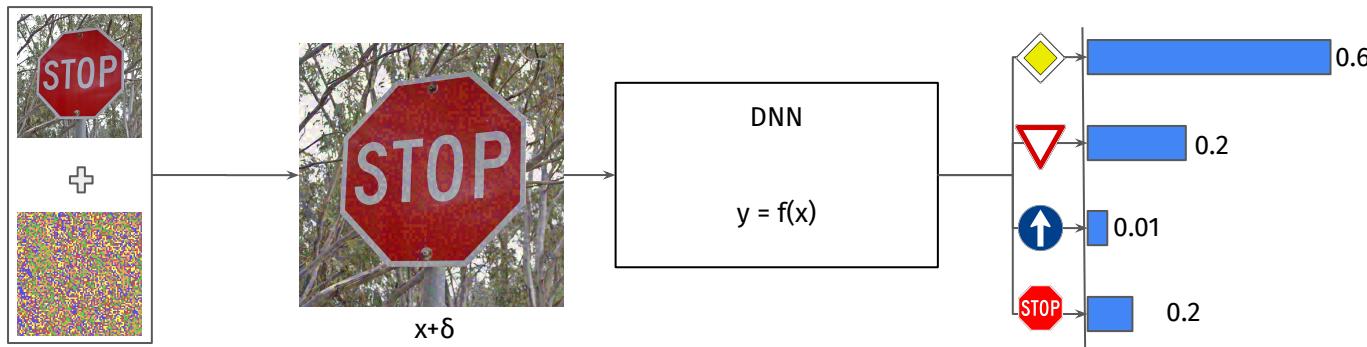
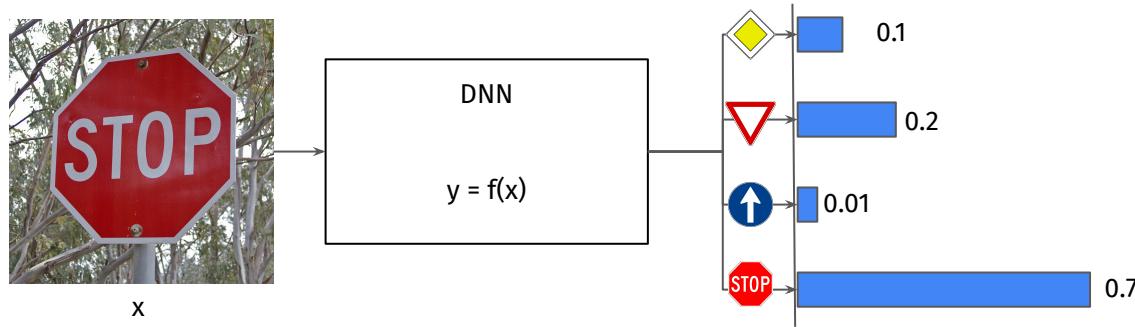


Definition and Taxonomy

The Classification Problem

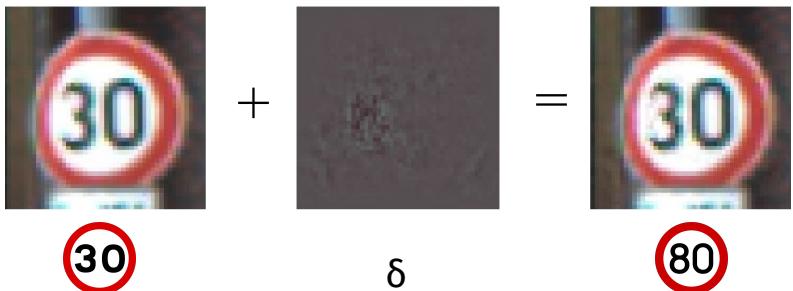


Adversarial Examples



Adversarial Examples: Digital or Real World

Digital Adversarial Perturbation



A malicious perturbation δ undetectable from human eye

Focus of this Lecture

Physical Adversarial Attacks

PATCH



POSTER



Adversarial perturbations that can be perceived from cameras, microphones, etc...

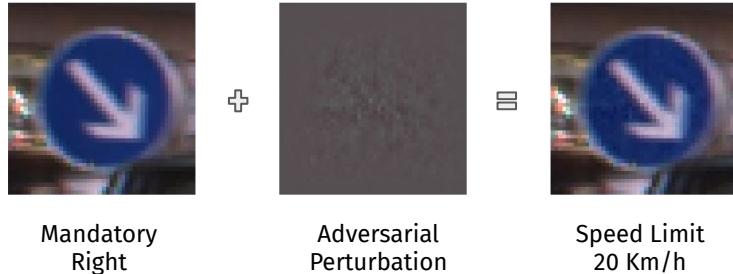
Focus of Next Lectures



Definition of Adversarial Examples

Definition (Adversarial Examples)

The result of perturbations, with a *small* $\| \cdot \|_p$ norm, that fool a classification model.



$$K(x + \delta) \neq K(x) \quad \text{and} \quad \|\delta\| \approx 0$$

Some $\| \cdot \|_p$ Norm

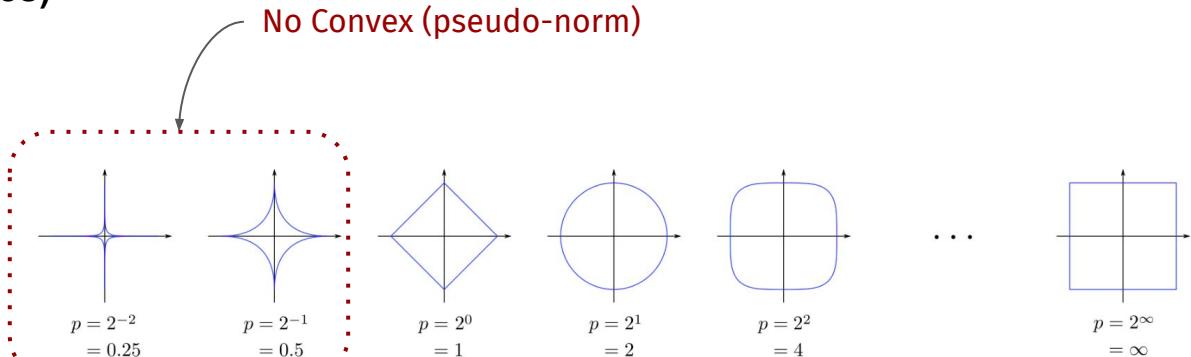
$\| \cdot \|_p$ - norms (Minkowski Distances)

$$\|\delta\|_0 = \#\{i : \delta_i \neq 0\}$$

$$\|\delta\|_p = \left(\sum_{i=1}^n |\delta_i|^p \right)^{\frac{1}{p}}$$

$$\|\delta\|_\infty = \max_i |\delta_i|$$

No Convex (pseudo-norm)



$$p = 2^{-2} \\ = 0.25$$

$$p = 2^{-1} \\ = 0.5$$

$$p = 2^0 \\ = 1$$

$$p = 2^1 \\ = 2$$

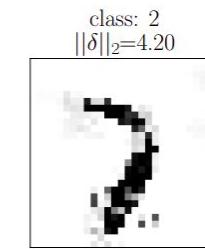
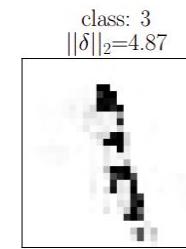
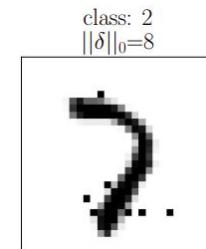
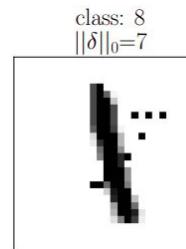
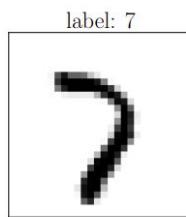
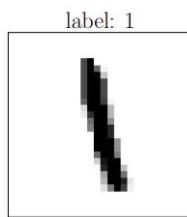
$$p = 2^2 \\ = 4$$

$$p = 2^\infty \\ = \infty$$



Different visualization of ℓ^p bounded attacks

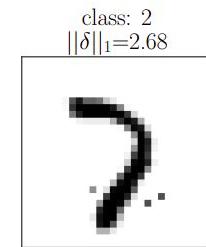
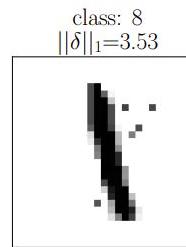
Definition in Minkowski norm covers many visually different cases



Few pixels, Unbounded values

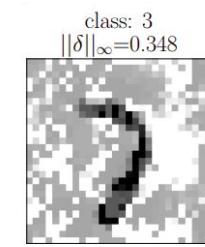
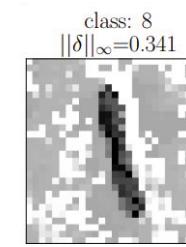
Few pixels, Bounded values

$$\|\delta\|_0 = \#\{i : \delta_i \neq 0\}$$



$$\|\delta\|_\infty = \max_i |\delta_i|$$

Few pixels, Bounded values



Many pixels, Bounded values

Targeted or Untargeted

Targeted

The attacker aims at making the model answering <target>



Un-Targeted

The attacker only aims at fooling the model, without setting the output

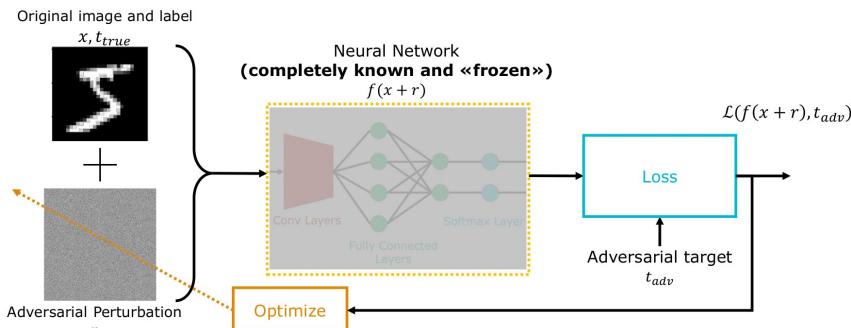


Black or White Box

White Box

The attacker has full access to:

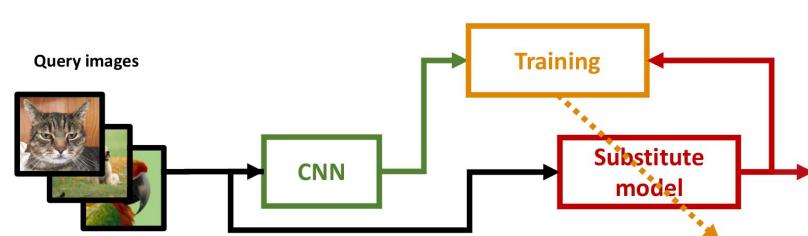
1. *Model architecture*
2. *Model weights*
3. *Data*



Black Box

The model acts as **oracle** and the attacker has full access to:

1. *Model Architecture* but not *data*
2. *Data* but not *Model Architecture*
3. Neither *Data* or *Model Architecture*



Minimal Taxonomy of Adversarial Examples

i.e.: the attacker is assumed to know both architecture and parameters of the network (**white box** setting) or anything about the model (**black box**). Intermediate settings can be considered.

Knowledge of the adversary

Targeted/untargeted

i.e.: the attacker aims at forcing the classifier to choose any wrong class (**untargeted**) vs one specific class (**targeted**)

Type of norm

i.e.: different norms (**2-norm**, **1-norm**, ∞ -norm) can be used to express the magnitude of the adversarial perturbation.

In the following slides no specific norm is considered, any of these will work (with different effects)



Generation of Adversarial Examples

Generation of Adversarial Examples

Methods based on the Loss Function

Un-Targeted

Increase the Loss Function in the correct class close to the input x

$$\max_{\delta \in \mathbb{R}^n} \mathcal{L}(f(x + \delta), l_{true})$$

$$\text{s.t. } \|\delta\|_p \leq \varepsilon$$

$$0 \leq x + \delta \leq 1$$

$$Q * (x + \delta) \in \{0, \dots, Q\}$$

Box Constraint

Integer Constraint,
 $Q=255$

Targeted

Decrease the Loss Function in the target class close to the input x

$$\min_{\delta \in \mathbb{R}^n} \mathcal{L}(f(x + \delta), l_{target})$$

$$\text{s.t. } \|\delta\|_p \leq \varepsilon$$

$$0 \leq x + \delta \leq 1$$

$$Q * (x + \delta) \in \{0, \dots, Q\}$$

Box and Integer Constraints ensure that the Adv.Ex is a representable image.

Generation of Adversarial Examples

Methods based on computing the
Minimal Adversarial Perturbation

Un-Targeted

Searching for the closest sample to x
that is classified in a different class

$$\min_{\delta \in \mathbb{R}^n} \|\delta\|_p$$

$$\text{s.t. } \mathcal{K}(x + \delta) \neq l_{true}$$

$$0 \leq x + \delta \leq 1$$

$$Q * (x + \delta) \in \{0, \dots, Q\}$$

Box Constraint

Integer Constraint,
 $Q=255$

Targeted

Searching for the closest sample to x
that is classified in the target class

$$\min_{\delta \in \mathbb{R}^n} \|\delta\|_p$$

$$\text{s.t. } \mathcal{K}(x + \delta) = l_{target}$$

$$0 \leq x + \delta \leq 1$$

$$Q * (x + \delta) \in \{0, \dots, Q\}$$

Box and Integer Constraints ensure that the Adv.Ex is a representable image.

Fast Gradient Sign Method

Leveraging the Loss Function

The loss function is used to search for the **worst** direction that fools the model and **not** to train the model $f : \mathbb{R}^n \rightarrow \mathbb{R}^c$

$$\delta^* = \boxed{\varepsilon} \cdot \text{sign} (\boxed{\nabla_x \mathcal{L}(f(x), l_{\text{true}})})$$

↓

Desired Magnitude
in ℓ^∞ norm

The gradient provides the
direction of the largest growth

Observation.

The adversarial perturbation has bounded uniform norm

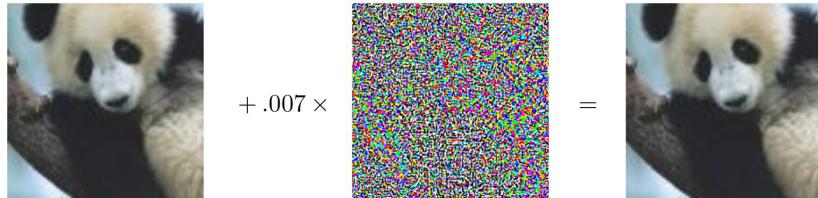
$$x^* = x + \delta^*, \quad \|\delta^*\|_\infty = \varepsilon$$

Fast Gradient Sign Method

FGSM

$$\delta^* = \varepsilon \cdot \text{sign}(\nabla_x \mathcal{L}(f(x), l_{true}))$$

$$x^* = x + \delta^*$$



The most famous Adv Ex. Attack on GoogLeNet

Pros

1. Easy to implement
2. Fast (useful in adv. training)
3. Bounded in uniform norm

Cons

1. Requires ε large (sometime visible)
2. Not always provide an Adv

Adversarial Training

Train the models by leveraging the following loss function

$$\tilde{\mathcal{L}}(f(x), l) = \alpha \cdot \mathcal{L}(f(x), l) + (1 - \alpha) \cdot \mathcal{L}(f(x + \delta^*), l)$$

Adversarial Loss

Minimal Adversarial Perturbation with Penalty

Minimum Problem with Constraints

$$\begin{aligned} \min_{\delta \in \mathbb{R}^n} \quad & \|\delta\|_p \quad * \\ \text{s.t.} \quad & \mathcal{K}(x + \delta) = l_{target} \\ & 0 \leq x + \delta \leq 1 \end{aligned}$$

Minimum Problem with Penalty

$$\begin{aligned} \min_{\delta \in \mathbb{R}^n} \quad & \boxed{c} \|\delta\|_p + \mathcal{L}(f(x + \delta), l_{target}) \\ \text{s.t.} \quad & 0 \leq x + \delta \leq 1 \end{aligned}$$

Penalty Constant.

Unbounded solution for $c=0$,
 $\delta(c) \sim 0$ for large values of c

Theorem.

Under convex assumption on $\mathcal{L}(f(x + \delta), l_{target})$ in the variable x we have that

1. For each c , the Box constrained L-BFGS has a solution $\delta(c)$
2. The penalty c can be optimized by line-search such that $\mathcal{K}(x + \delta(c)) = l_{target}$

Minimal Adversarial Perturbation with Penalty

BL-BFGS

Box constrained Limited memory Broyden-Fletcher-Goldfarb-Shanno
is a quasi-Newton method for minimum problem with box constraints. ^a

$$\begin{aligned} \min_{\delta \in \mathbb{R}^n} \quad & c \|\delta\|_p + \mathcal{L}(f(x + \delta), l_{target}) \\ \text{s.t.} \quad & 0 \leq x + \delta \leq 1 \end{aligned}$$

Pros.

1. Provide particularly low perturbations (approximate the Minimal Adversarial Perturbation)
2. Reliable. Almost always converge.

Cons.

1. Slow. Requires ~20k forwards of the model for each line search of c



First Adversarial Examples presented to the community. Attack of AlexNet ^b

^a Dimitri P. Bertsekas. "Constrained Optimization and Lagrange Multiplier Methods" (1996)
^b Christian Szegedy et al. "Intriguing properties of neural networks" (2013)

Including the Box Constraint in the Objective (CW)

Minimum Problem with Constraints

$$\begin{aligned} \min_{\delta \in \mathbb{R}^n} \quad & \|\delta\|_p && * && * \\ \text{s.t.} \quad & \mathcal{K}(x + \delta) = l_{target} \\ & 0 \leq x + \delta \leq 1 \end{aligned}$$

Change of Variable (Simplification 1)

Optimizing w instead of directly optimize over the adversarial perturbation \square

$$\delta = \frac{1}{2} (\tanh(w) + 1) - x$$

Remark!

The perturbation \square satisfies the Box.Con.

Improve Loss Function (Simplification 2)

$$\mathcal{L}(y, l) = \left[y_l - \max_{j \neq l} y_j \right]^+$$

Positive part of the “margin”
computed on the logits of f

Including the Box Constraint in the Objective (CW)

Minimum Problem with Constraints

$$\begin{aligned} \min_{\delta \in \mathbb{R}^n} \quad & \|\delta\|_p && * && * \\ \text{s.t.} \quad & \cancel{\mathcal{K}(x + \delta) = l_{target}} && && \\ & \cancel{0 \leq x + \delta \leq 1} && && \end{aligned}$$



Minimum Problem with Penalty

$$\begin{aligned} \min_{\delta \in \mathbb{R}^n} \quad & c \|\delta\|_p + \mathcal{L}(f(x + \delta), l_{target}) \\ \text{s.t.} \quad & 0 \leq x + \delta \leq 1 \end{aligned}$$

CW Minimum Problem (Unconstrained)

$$\min_{w \in \mathbb{R}^n} \underbrace{\left\| \frac{1}{2} (\tanh(w) + 1) - x \right\|_p}_{\delta} + c \cdot \underbrace{\mathcal{L}\left(f\left(\frac{1}{2}(\tanh(w) + 1)\right), l_{target}\right)}_{x+\delta}$$

Remark

The penalty c can be optimized by line-search such that $\mathcal{K}(x + \delta(c)) = l_{target}$

Including the Box Constraint in the Objective (CW)

CW Minimum Problem (Unconstrained)

$$\min_{w \in \mathbb{R}^n} \underbrace{\left\| \frac{1}{2} (\tanh(w) + 1) - x \right\|_p}_{\delta} + c \cdot \underbrace{\mathcal{L}\left(f\left(\frac{1}{2}(\tanh(w) + 1)\right), l_{target}\right)}_{x+\delta}$$

Pros.

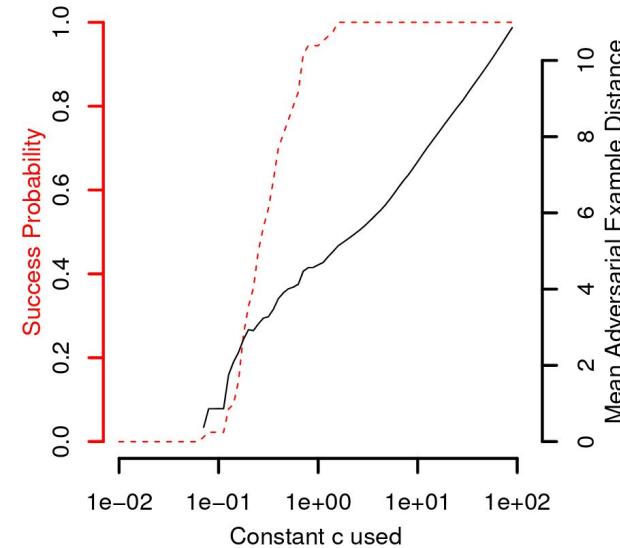
The first work that **directly** aims at solving MAP

MAP provides **Certifiable Guarantees** (next lecture)

Cons.

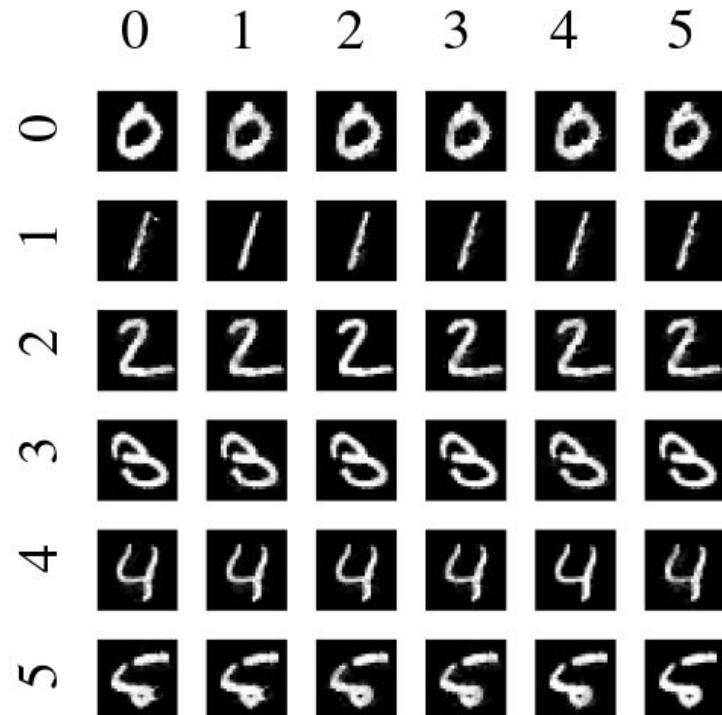
Forward required for each c : ~20k (Hence slow)

Binary steps required : ~ 12

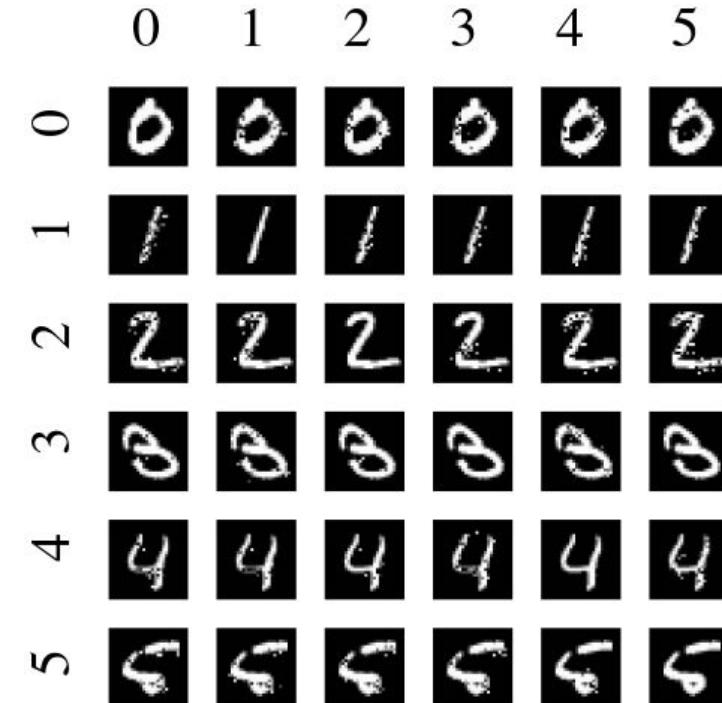


- A large penalty implies high success rate but farther adversarial examples
- A small penalty implies accurate solutions but small convergence rate.

Examples of attack with CW



ℓ^2 Minimal Adversarial Perturbation



ℓ^0 Minimal Adversarial Perturbation

DeepFool: A Fast and Accurate Method

The Binary Linear Case

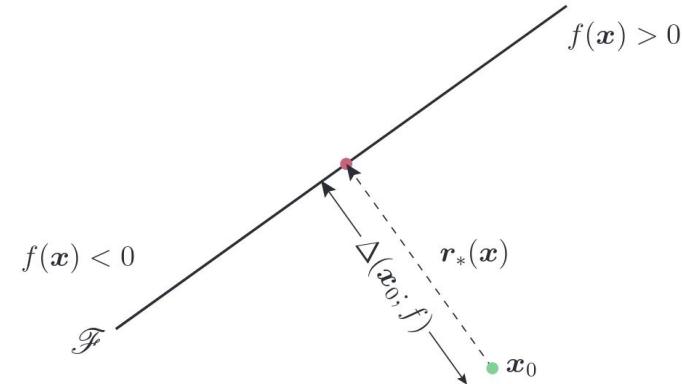
The following is an explicit solution.

If $f(x) = w^T x + b$

then

$$\delta = -\frac{f(x)}{\|w\|^2} w$$

Minimal Adversarial Perturbation
that fools the binary linear classifier



DeepFool: A Fast and Accurate Method

Extension to the Multiclass Case

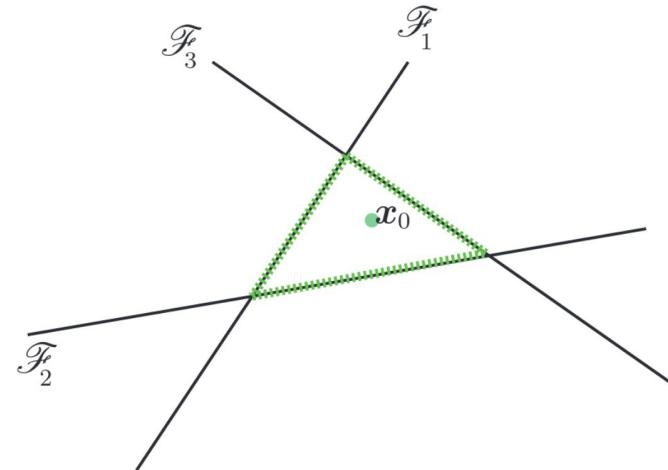
$$f(x) = Wx + b, \quad W = \begin{bmatrix} w^{(1)\top} \\ \vdots \\ w^{(C)\top} \end{bmatrix} \in \mathbb{R}^{C \times n}$$

Remark.

Minimal Adversarial Perturbations targeted for each class i

$$\delta^{(i)} = -\frac{f_i(x)}{\|w^{(i)}\|^2} w^{(i)}$$

The smallest corresponds to the **untargeted** attack.



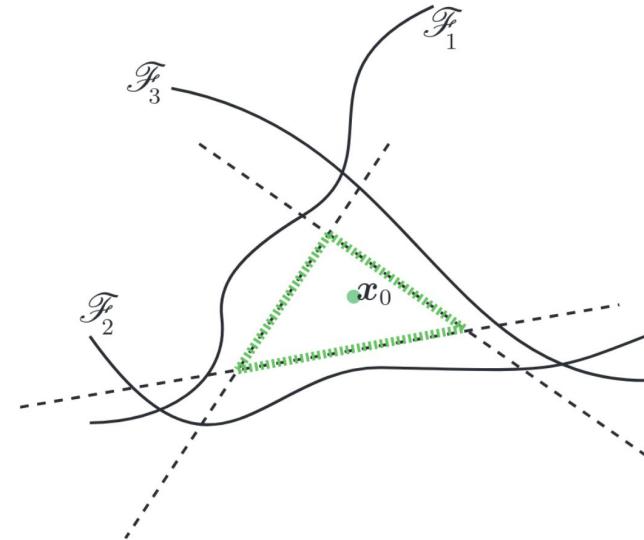
DeepFool: A Fast and Accurate Method

Not-Linear Multiclass Case

Algorithm 2 DeepFool: multi-class case

```
1: input: Image  $\mathbf{x}$ , classifier  $f$ .  
2: output: Perturbation  $\hat{\mathbf{r}}$ .  
3:  
4: Initialize  $\mathbf{x}_0 \leftarrow \mathbf{x}$ ,  $i \leftarrow 0$ .  
5: while  $\hat{k}(\mathbf{x}_i) = \hat{k}(\mathbf{x}_0)$  do  
6:   for  $k \neq \hat{k}(\mathbf{x}_0)$  do  
7:      $\mathbf{w}'_k \leftarrow \nabla f_k(\mathbf{x}_i) - \nabla f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_i)$   
8:      $f'_k \leftarrow f_k(\mathbf{x}_i) - f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_i)$   
9:   end for  
10:   $\hat{l} \leftarrow \arg \min_{k \neq \hat{k}(\mathbf{x}_0)} \frac{|f'_k|}{\|\mathbf{w}'_k\|_2}$   
11:   $\mathbf{r}_i \leftarrow \frac{|f'_l|}{\|\mathbf{w}'_l\|_2} \mathbf{w}'_l$   
12:   $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \mathbf{r}_i$   
13:   $i \leftarrow i + 1$   
14: end while  
15: return  $\hat{\mathbf{r}} = \sum_i \mathbf{r}_i$ 
```

Idea: Iterate the procedure



DeepFool: A Fast and Accurate Method

Empirical Robustness Estimation

$$\rho_{adv}(\mathcal{X}) = \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \frac{\|\delta(x)\|_2}{\|x\|_2}$$

Minimal Adversarial
Perturbation on x

Testing Dataset

Classifier	Test error	$\hat{\rho}_{adv}$ [DeepFool]	time
LeNet (MNIST)	1%	2.0×10^{-1}	110 ms
FC500-150-10 (MNIST)	1.7%	1.1×10^{-1}	50 ms
NIN (CIFAR-10)	11.5%	2.3×10^{-2}	1100 ms
LeNet (CIFAR-10)	22.6%	3.0×10^{-2}	220 ms
CaffeNet (ILSVRC2012)	42.6%	2.7×10^{-3}	510 ms*
GoogLeNet (ILSVRC2012)	31.3%	1.9×10^{-3}	800 ms*

Pros.

1. Fast, requires ~3C forwards
2. Provide Empirical MAP for the whole dataset (CW can not)

Cons.

1. Still no theoretical guarantees
2. Suboptimal solutions are found

Projected Gradient Descent (PGD)

Maximum Problem (Un-Targeted)

$$\max_{\delta \in \mathbb{R}^n} \mathcal{L}(f(x + \delta), l_{true})$$

$$\text{s.t. } \|\delta\|_p \leq \varepsilon$$

$$0 \leq x + \delta \leq 1$$

$$Q * (x + \delta) \in \{0, \dots, Q\}$$

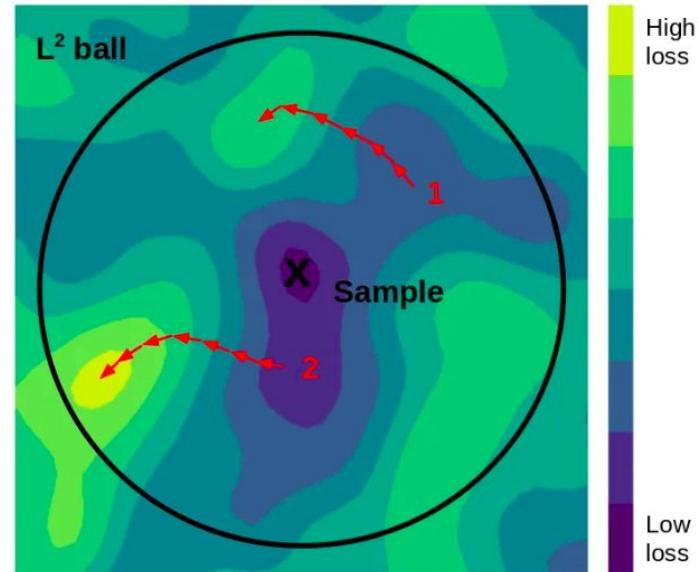
Iterating FGSM

Learning rate of the
Gradient Ascent

$$\tilde{x}^{(t)} = x^{(t-1)} + \alpha \operatorname{sgn} \left(\nabla_x \mathcal{L}(f(x^{(t-1)}), l_{true}) \right)$$

$$x^{(t)} = \Pi_{B_p(x, \varepsilon)}(\tilde{x}^{(t)})$$

Projection into the
 p -Ball of radius ε

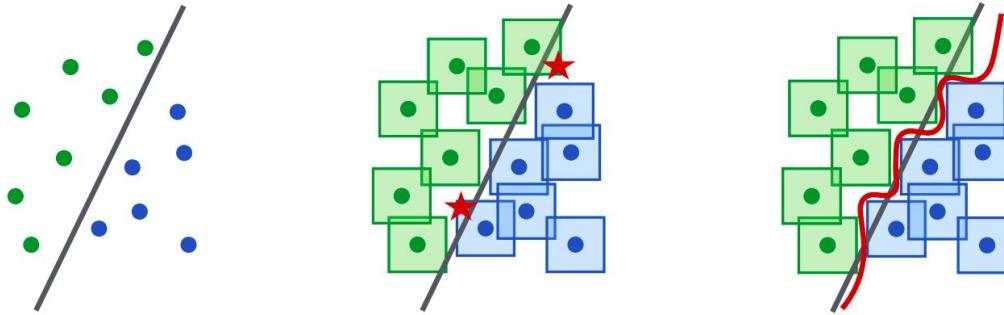


1. Take a random starting point
2. Iterate the SG(A) for few steps
3. Project into the \mathbb{I}^p Ball
4. Repeat 1,2,3

Adversarial Training by PGD

Minimization of the Empirical Worst Case

$$\min_{\theta} \rho(\theta), \quad \rho(\theta) = \mathbb{E}_{(x,l) \in \mathcal{X}} \left[\max_{\|\delta\|_p < \varepsilon} \mathcal{L}(f_\theta(x + \delta), l; \theta) \right]$$



A Black Box Attack

Model acts as an Oracle^a

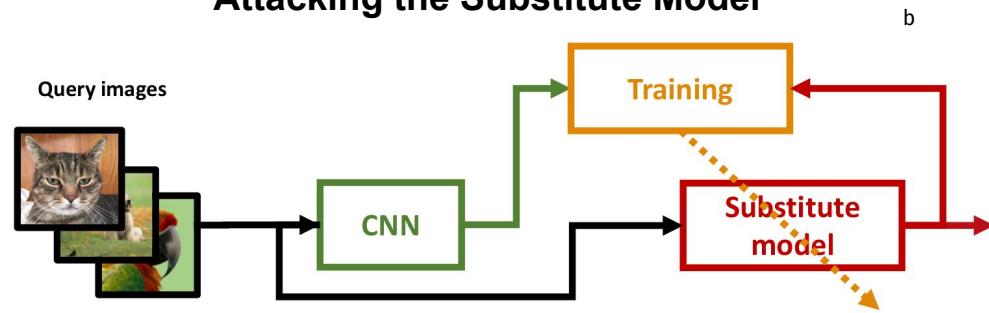
The model \mathcal{K} is given without:

1. Training dataset
2. Model Architecture and Weights

Procedure

1. Initial Collection
2. Architecture Selection
3. Substitute Training
4. Adversarial sample Crafting

Attacking the Substitute Model



A new model is trained, and a white-box attack is performed on the new one; the procedure leverages the **TRANSFERABILITY PROPERTY** of Adversarial Examples

^a Nicolas Papernot et al. "Practical Black-Box Attacks against Machine Learning" (2017)

^b Federico Nesti. "An introduction to Adversarial Examples", lecture of Deep Learning Course (2022)

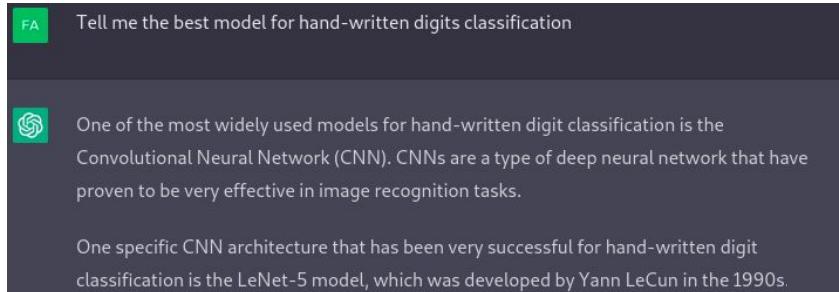
A Black Box Attack: Substitute DNN and Training

Initial Collection

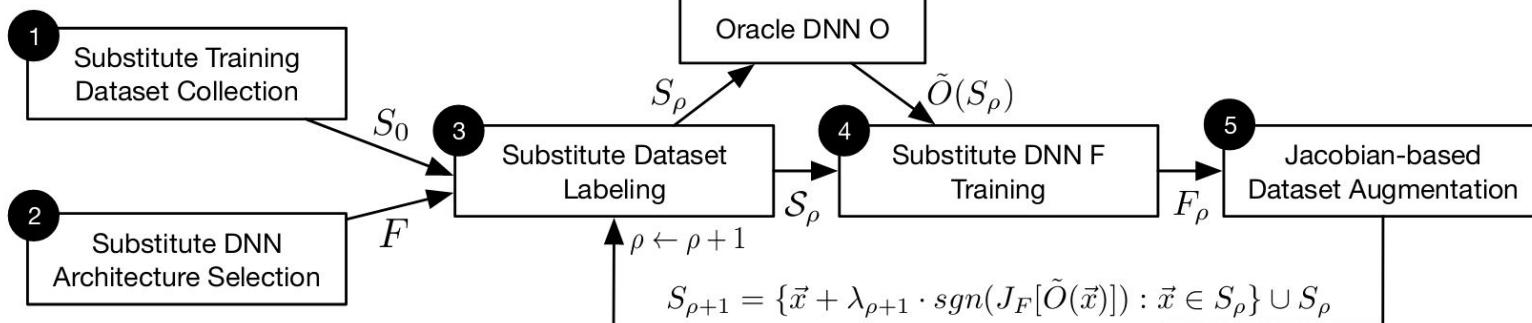


Hand-crafted images of 10 digits. At least one image per class has to be chosen in order to start the procedure.

Architecture Selection

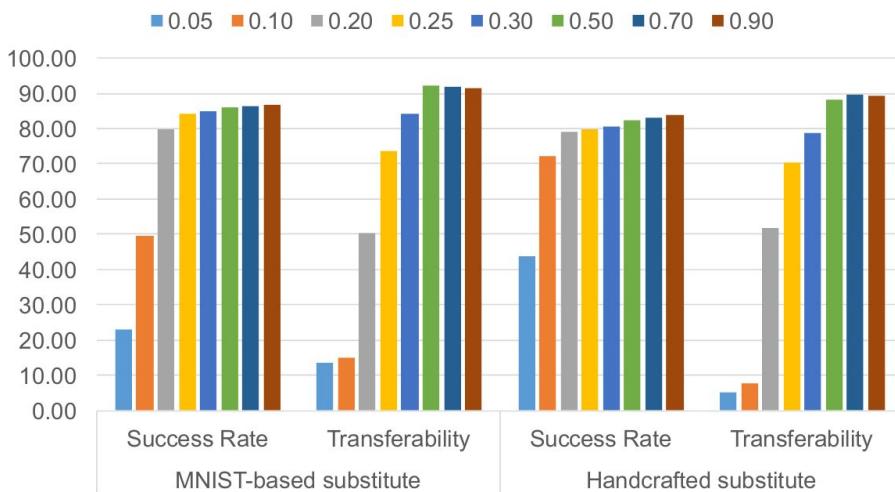


Training Schema

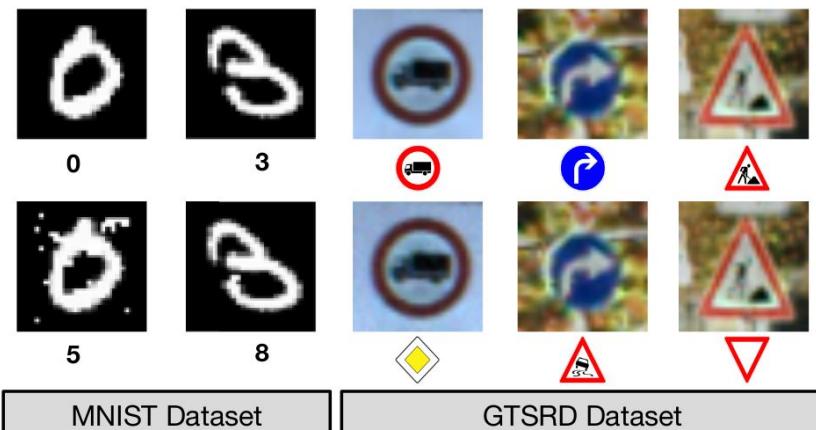


A Black Box Attack: Evaluation

Success rate and Transferability



Graphical Evaluation



For both models $\epsilon=0.2$

Defense Strategies

Attacker-Defensor Paradigm

1. Adversarial training: Augmenting the training data with adversarial examples.
2. Defensive distillation: Training a second DNN (**distilled model**) to mimic the output of the first DNN on the training set. The distilled model is, empirically, more robust.
3. Gradient masking: Modifying the DNN architecture to mask the gradients in certain layers (avoid white-box adversarial attacks).
4. Input Preprocessing: Modifying the input data. E.g., adding noise or blurring the image, to make it more difficult for an attacker to create effective adversarial examples.
5. Randomized defenses: Adding randomness to the DNN at different stages, such as the weights, the architecture, or the input. (**next lecture**)
6. Lipschitz Regularization: Crafting models with a bounded Lipschitz constant to certify the robustness against adversarial examples. (**next lecture**)

Adversarial Training

Adding Adversarial Examples at training stage to improve the robustness.

FGSM

Goodfellow et al. proposed a new loss function

$$\tilde{\mathcal{L}}(f(x), l) = \alpha \cdot \mathcal{L}(f(x), l) + (1 - \alpha) \cdot \mathcal{L}(f(x + \delta^*), l)$$

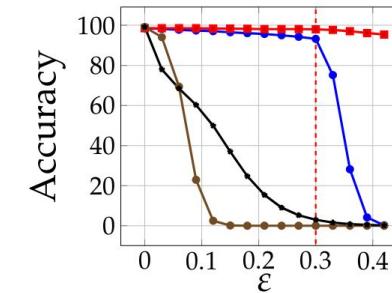
Adversarial Loss

Projected Gradient Descent

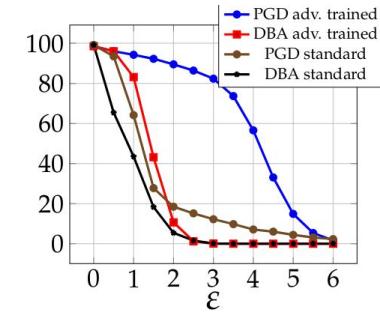
Madry et al. Improved the strategies with a finer estimation

$$\min_{\theta} \rho(\theta), \quad \rho(\theta) = \mathbb{E}_{(x,l) \in \mathcal{X}} \left[\max_{\|\delta\|_p < \varepsilon} \mathcal{L}(f_\theta(x + \delta), l; \theta) \right]$$

Remark! The Robustness evaluation is empirical (no guaranty)



(a) MNIST, ℓ_∞ -norm



(b) MNIST, ℓ_2 -norm

Goodfellow et al. "Explaining and Harnessing Adversarial Examples" (2014)

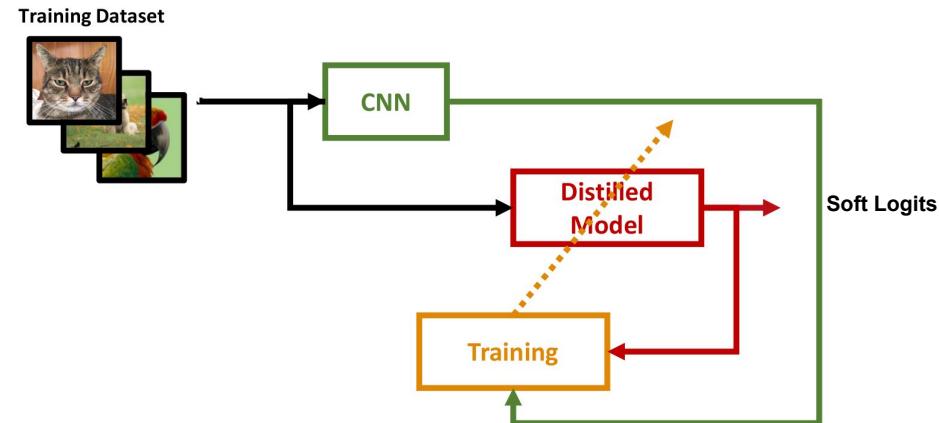
Aleksander Madry et al. "Towards Deep Learning Models Resistant to Adversarial"

Defensive Distillation

Training a Smaller Model by leveraging Cross-Entropy with high temperature

Soft Training

The **distilled model** is trained not with hard labels (one-hot encoded vectors), but with soft labels (probability values of the prediction of the original network)



Remark. The strategy is broken by Carlini-Wagner in 2017.

^a Nicolas Papernot et al. "Distillation as a defense to adversarial perturbations against deep neural networks" (2016)
^b Federico Nesti. "An introduction to Adversarial Examples", lecture of Deep Learning Course (2022)

Defensive Distillation

Cross-Entropy with Temperature

$$\mathcal{L}(y, l; T) = -\log \left(\frac{e^{\frac{y_l}{T}}}{\sum_{i=1}^C e^{\frac{y_i}{T}}} \right)$$

Temperature

High temperature increases the **uncertainty**,
and hence improves the (empirical) **robustness**
and the **generalization** of the distilled model

$$T = 1$$

(regular softmax)



$$T = 2$$



$$T \rightarrow \infty$$



$$y_{soft}(x) = \begin{bmatrix} 0,007 \\ 0,011 \\ 0,062 \\ 0,920 \end{bmatrix}$$

$$y_{soft}(x) = \begin{bmatrix} 0,059 \\ 0,076 \\ 0,178 \\ 0,687 \end{bmatrix}$$

$$y_{soft}(x) = \begin{bmatrix} 0,25 \\ 0,25 \\ 0,25 \\ 0,25 \end{bmatrix}$$

Gradient Masking

Avoid white-box attacks and evaluation based on the gradient by a.e flat activation function

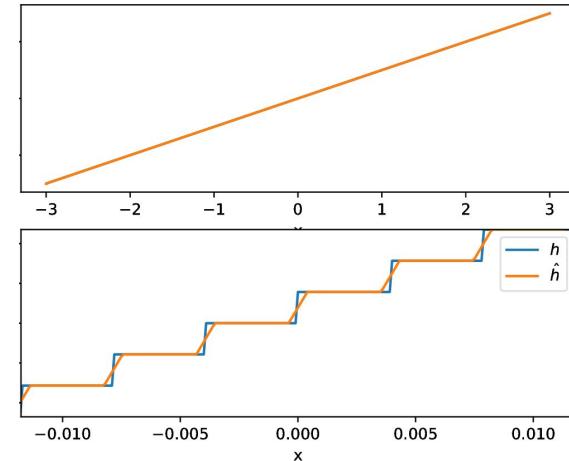
Remark. The strategy is broken by Papernot with black box attacks

Fun fact

Discussion with Ian Goodfellow and Our Clarifications

We received some inquiries on [Ian Goodfellow's comment "Gradient Masking Causes CLEVER to Overestimate Adversarial Perturbation Size"](#) on our paper. We thank Ian for the discussion but the comments are inappropriate and not applicable to our paper. CLEVER is intended to be a tool for network designer and to evaluate network robustness in the "white-box" setting. Especially, the argument that on digital computers all functions are not Lipschitz continuous and behave like a staircase function (where the gradient is zero almost everywhere) is incorrect. Under the white-box setting, gradients can be computed via automatic differentiation, which is well supported by mature packages like TensorFlow. See [our reply and discussions with Ian Goodfellow on gradient masking and implementation on digital computers](#).

CLEVER is a verification strategy that leverages the gradient to estimate the robustness of a sample

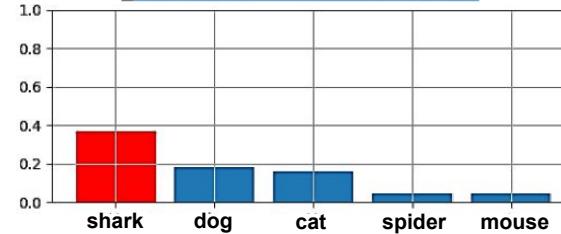
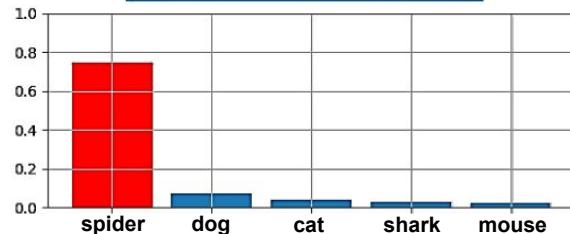


Gradient of a model with such a activation function is 0 almost-everywhere

Input Transformations



Translation



Other Transformations

Rotations, Mirror, Blur, Gaussian Noise

Conclusion

Adversarial Examples for Biological Vision Systems

Polyphemus Moth



Defense mechanisms that involves the **large eyespots** on its hindwings.

Chihuahua or Muffin



Conclusion

Introduction, Definitions and Taxonomy

Adversarial Examples Generation

Defense Mechanisms

Thanks for the attention

Fabio Brau



Scuola Superiore Sant'Anna, Pisa



fabio.brau@santannapisa.it



<https://www.linkedin.com/in/fabio-brau/>



<http://retis.santannapisa.it/~f.brau>

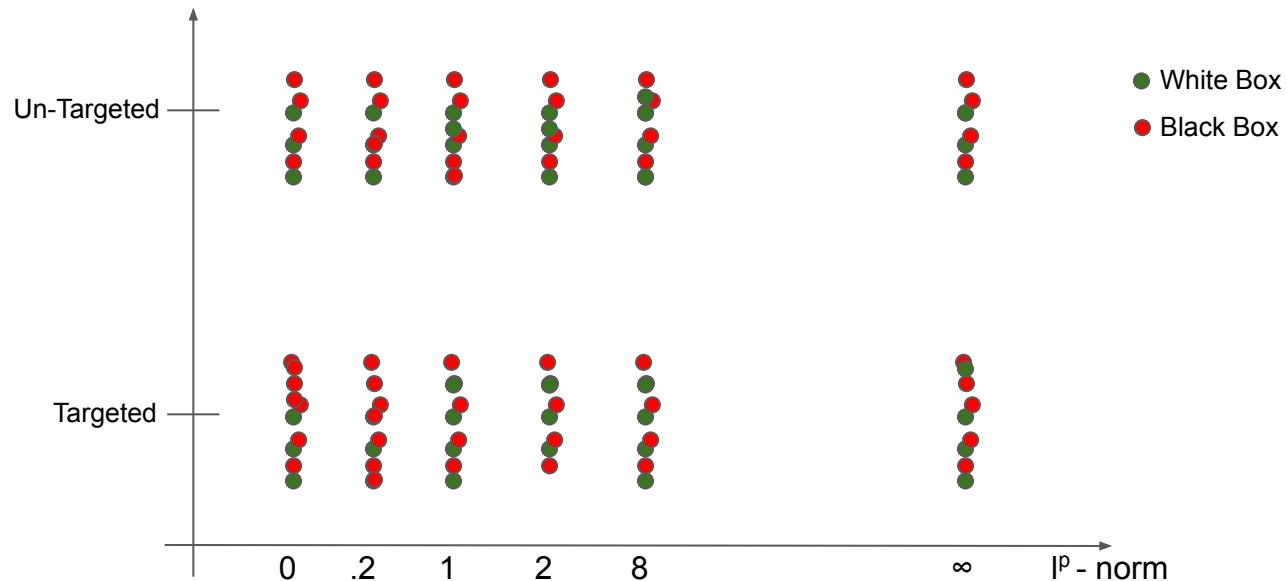


ISTITUTO
DI TELECOMUNICAZIONI,
INFORMATICA
E FOTONICA



Sant'Anna
Scuola Universitaria Superiore Pisa

Taxonomy Map



Generation of Adversarial Examples (Un-Targeted)

Adversarial Generation follows two main strategies

Loss Dependent

Increase the Loss Function on the correct class of the input x

$$\max_{\delta \in \mathbb{R}^n} \mathcal{L}(f(x + \delta), l_{true})$$

$$\text{s.t. } \|\delta\|_p \leq \varepsilon$$

$$0 \leq x + \delta \leq 1$$

Box Constraint

$$Q * (x + \delta) \in \{0, \dots, Q\}$$

Quantized
Constraint

Minimum Distance

Searching for the Closest Sample to x that is classified in a different class

$$\min_{\delta \in \mathbb{R}^n} \|\delta\|_p$$

$$\text{s.t. } \mathcal{K}(x + \delta) \neq l_{true}$$

$$0 \leq x + \delta \leq 1$$

$$Q * (x + \delta) \in \{0, \dots, Q\}$$

Generation of Adversarial Examples (Targeted)

Adversarial Generation follows two main strategies

Loss Dependent

Decrease the Loss Function on the target class

$$\max_{\delta \in \mathbb{R}^n} \mathcal{L}(f(x + \delta), l_{true})$$

$$\text{s.t. } \|\delta\|_p \leq \varepsilon$$

$$0 \leq x + \delta \leq 1$$

Box Constraint

$$Q * (x + \delta) \in \{0, \dots, Q\}$$

Minimum Distance

Searching for the Closest Sample to x that is classified in the target class

$$\min_{\delta \in \mathbb{R}^n} \|\delta\|_p$$

$$\text{s.t. } \mathcal{K}(x + \delta) \neq l_{true}$$

$$0 \leq x + \delta \leq 1$$

Quantized
Constraint

$$Q * (x + \delta) \in \{0, \dots, Q\}$$