



מכון ויצמן למדע
WEIZMANN INSTITUTE OF SCIENCE

25/11/2021

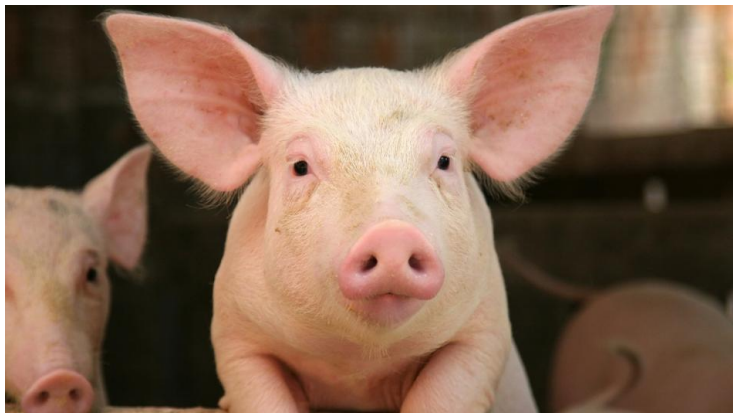


Introduction To Adversarial Examples

Niv Haim

Weizmann Institute

DL4CV Course Winter 2022 (20224182)



+



x 0.02

98.6% pig

=



99.0% airliner

(0.000000000000000000
000000000000005% pig...)

Biggio et al. 2013, "Evasion attacks against machine learning at test time"
Szegedy et al. 2014, "Intriguing properties of neural networks"
Goodfellow et al. 2015, "Explaining and Harnessing Adversarial Examples"

What is an Adversarial Example?

Perturbation Attack

- Originally coined by Szegedy et al., 2013:

“we find that applying an imperceptible non-random perturbation to a test image, it is possible to arbitrarily change the network’s prediction.

*... we term the so perturbed examples ‘**adversarial examples**’*”



98.6% pig

+



x 0.02

=



99.0% airliner

Outline

Today we will:

airliner

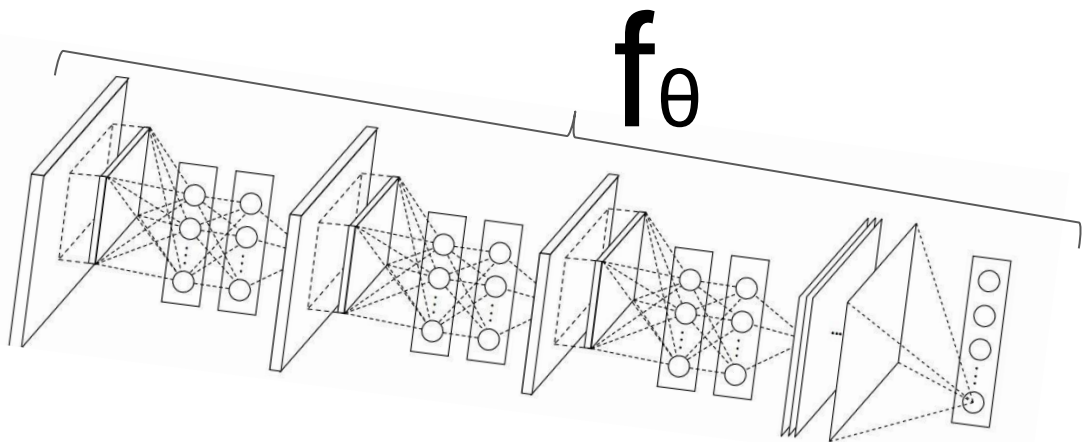


- See Adversarial Example
- Discuss what they are
- Learn how to generate them
- Learn how to (maybe) defend against them
- Learn about properties and advantages

Brief recap on training neural networks



Image by [Simon](#) from [Pixabay](#)



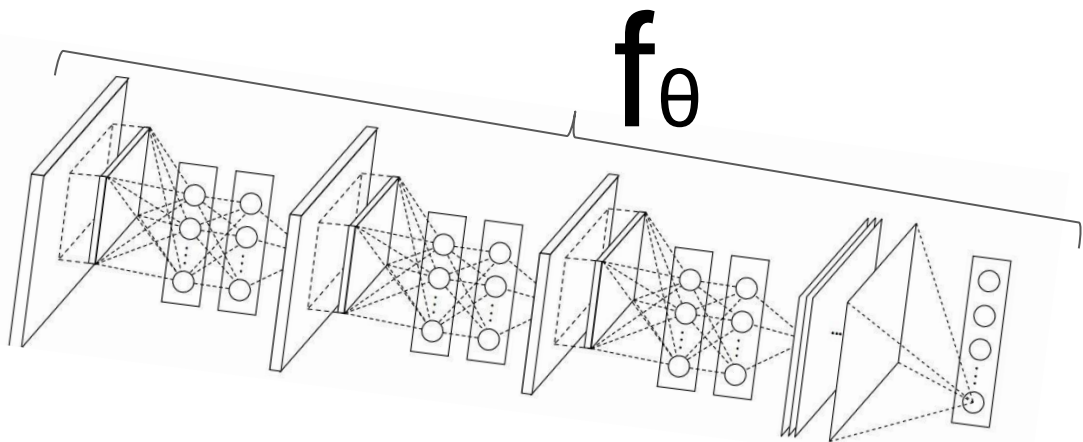
$$L(f_{\theta}(x), y)$$

purpose of loss:
How “well” we classify

Brief recap on training neural networks



Image by [Simon](#) from [Pixabay](#)



most common loss – CrossEntropy:

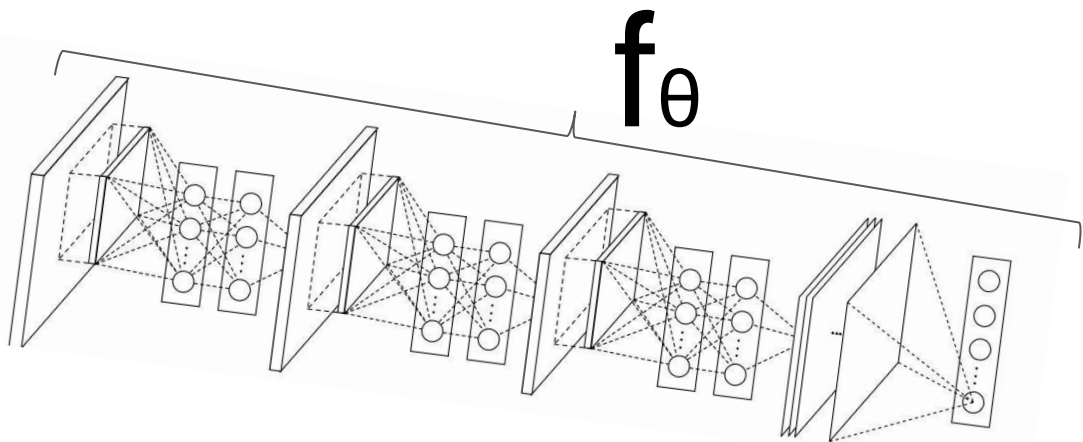
$$L(f_{\theta}(x), y)$$

$$L(f_{\theta}(x), y) = -\log \left(\frac{e^{f_{\theta}(x)_y}}{\sum_j e^{f_{\theta}(x)_j}} \right)$$

Brief recap on training neural networks



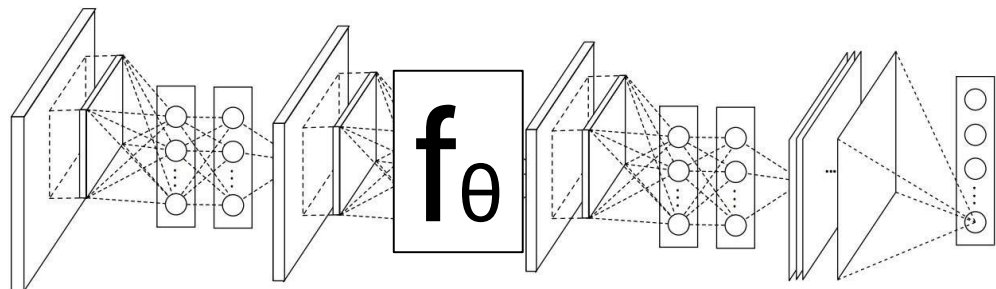
Image by [Simon](#) from [Pixabay](#)



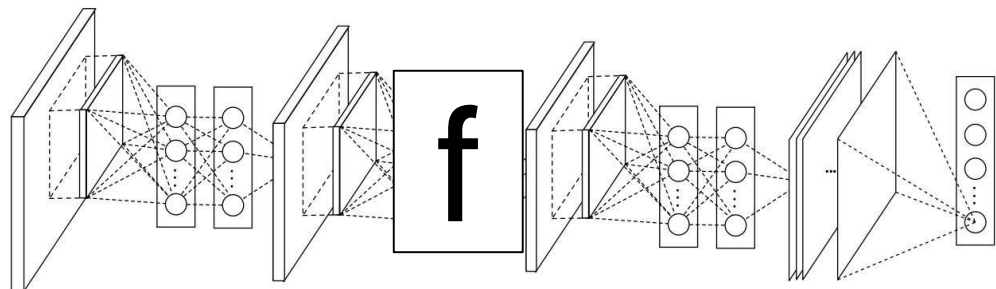
minimize loss:

$$L(f_{\theta}(x), y) \longrightarrow - \nabla_{\theta} L$$

Generating an Adversarial Example



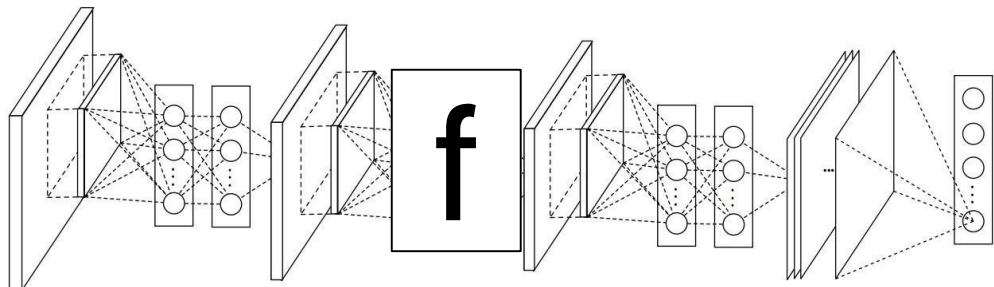
Generating an Adversarial Example



Generating an Adversarial Example



89.7% pig



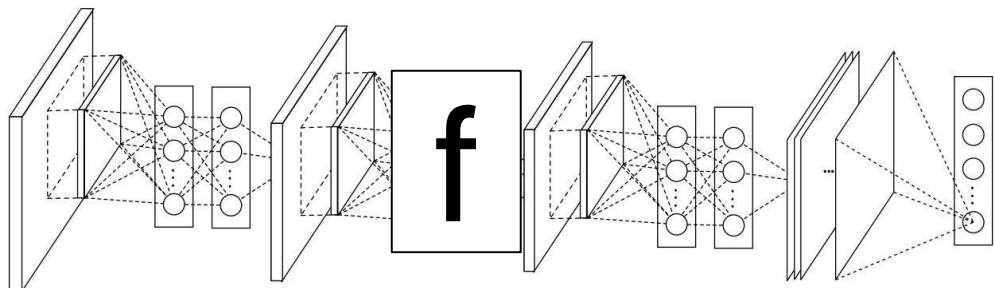
want to fool classifier
by changing δ

$$f(x + \delta) \neq y$$

Generating an Adversarial Example



89.7% pig



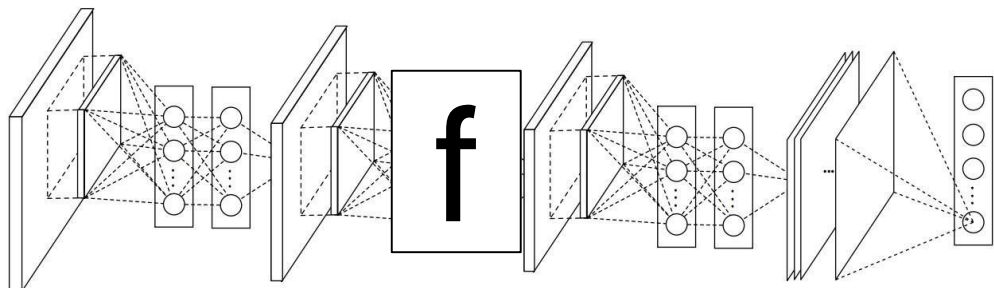
want to fool classifier \rightarrow d measures “badness”
by changing δ

$$d(f(x+\delta), y)$$

Generating an Adversarial Example



89.7% pig



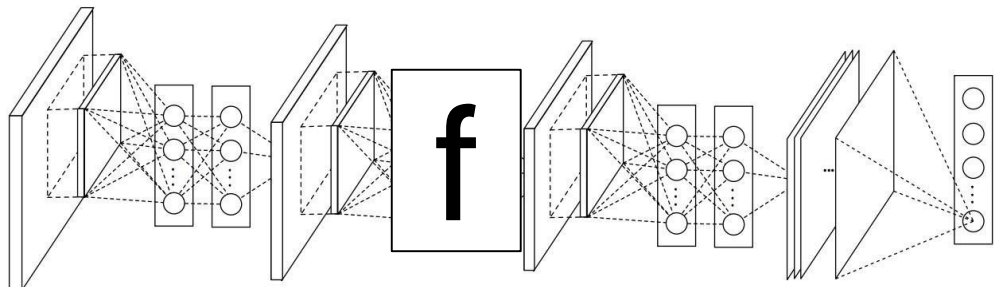
want to fool classifier → used L to ~~maximize “wellness”~~
maximize “badness”?

$$L(f(x+\delta), y)$$

Generating an Adversarial Example



89.7% pig



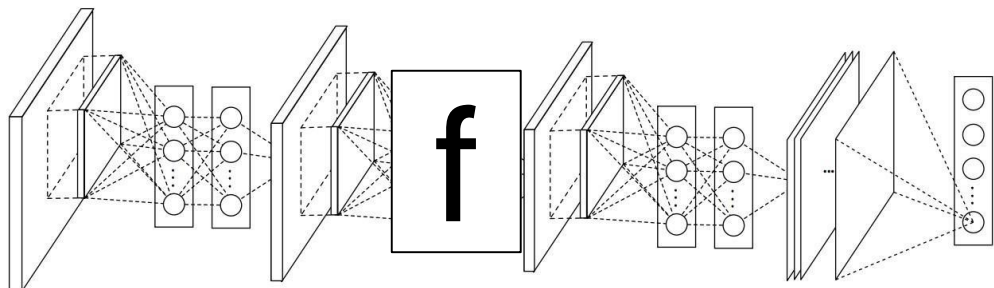
want to fool classifier \rightarrow maximize L w.r.t δ

$$L(f(x+\delta), y) \rightarrow \nabla_{\delta} L$$

Generating an Adversarial Example

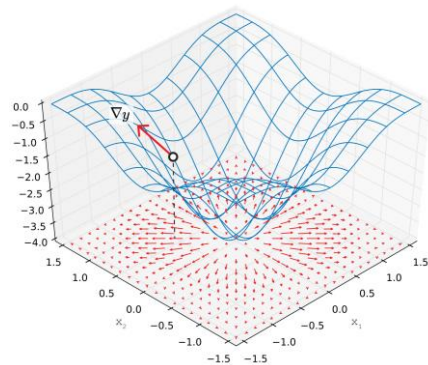


89.7% pig



want to fool classifier \rightarrow maximize L w.r.t δ

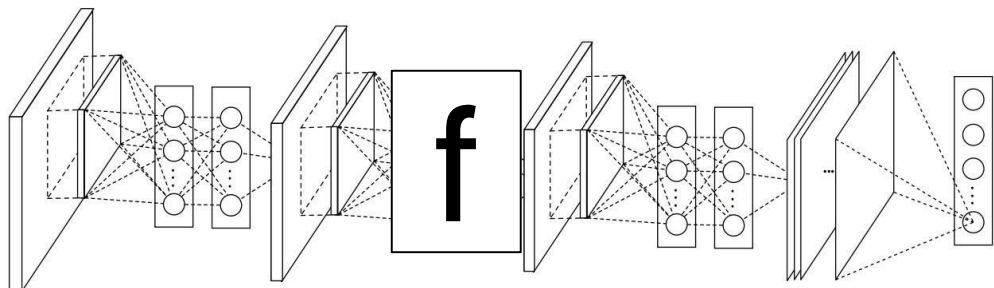
$$L(f(x+\delta), y) \rightarrow +\nabla_{\delta} L$$



Generating an Adversarial Example

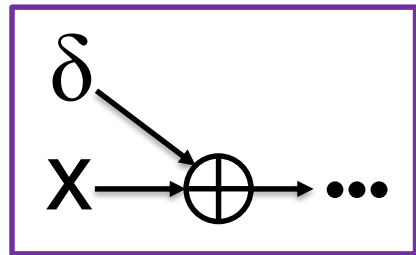


89.7% pig



want to fool classifier \rightarrow maximize L w.r.t x

$$L(f(\underbrace{x+\delta}_{\text{input}}), y) \longrightarrow + \nabla_{\underbrace{x}_{\text{input}}} L$$

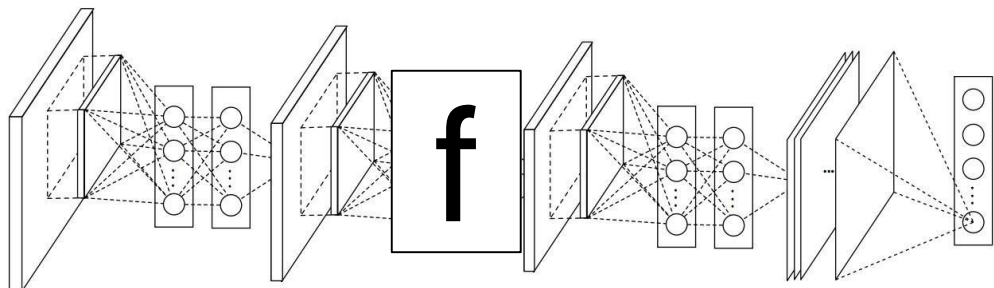


(just a technicality..)

Generating an Adversarial Example



89.7% pig



want to fool classifier \rightarrow maximize L w.r.t x

$$L(f(x+\delta), y) \rightarrow \delta = +\nabla_x L$$

Follow the gradient w.r.t x (the input image)



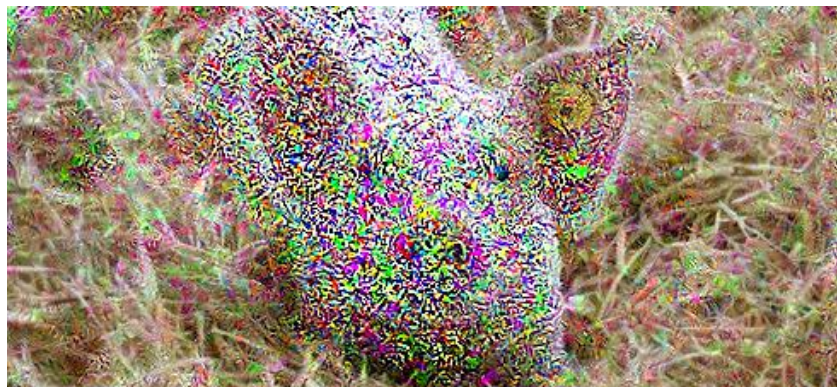
X (original): 89.7% pig



$X + \nabla_x L$: 68.6% hay

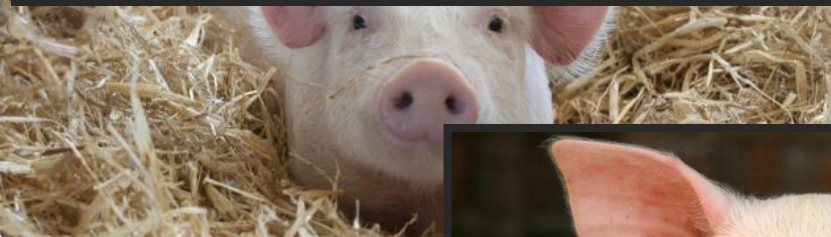


$X + 10 \times \nabla_x L$: 44.7% pig



$X + 100 \times \nabla_x L$: 44.8% fireguard

Did we generate an adversarial example? Need small δ ...



X (original):



: 68.6% hay



99.0% airliner



X + 10 \times $\nabla_x L$: 44.7% pig



X + 100 \times $\nabla_x L$: 44.8% fireguard

We want *small noise*



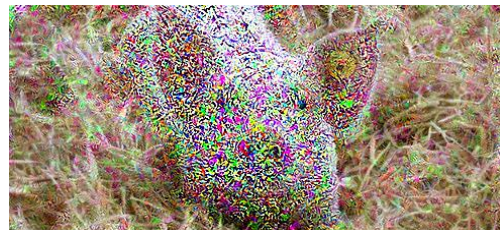
X

+



δ

=



$$\|\delta\|_{\infty} \leq 0.1$$

0.1	-0.1		
	0.1	0.05	-0.02
		-0.09	
			10^{-5}

What is small δ ?

$$\|\delta\|_{\infty} < \epsilon$$

small δ & $\delta = f(\nabla_x L)$?

“Enforcing $\|\nabla_x L\|_\infty < \epsilon$ ” :


ϵ	$-\epsilon$	ϵ	...
...	ϵ
...	...	$-\epsilon$...
...	...	0	...



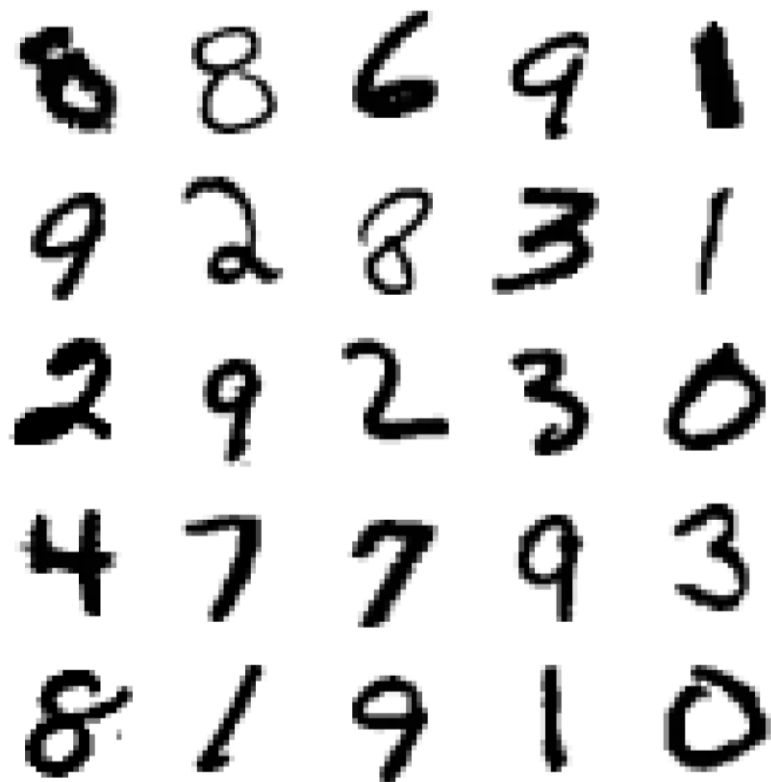
$$\delta = \epsilon \cdot \text{sgn}(\nabla_x L)$$

Fast Gradient Sign Method

a.k.a FGSM (Goodfellow et al. 2015)


$$\delta = \max_{\|\delta\|_\infty \leq \epsilon} L(f(x + \delta), y) \approx \max_{\|\delta\|_\infty \leq \epsilon} L(f(x), y) + \nabla_x L \delta$$

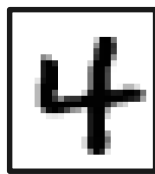
FGSM – example on MNIST



Classifier

```
model = nn.Sequential(  
    nn.Conv2d(1, 16, 4, stride=2, padding=1),  
    nn.ReLU(),  
    nn.Conv2d(16, 32, 4, stride=2, padding=1),  
    nn.ReLU(),  
    Flatten(),  
    nn.Linear(32 * 7 * 7, 100),  
    nn.ReLU(),  
    nn.Linear(100, 10)  
)
```

FGSM - simple but vicious



$$\mathbf{X}^{\text{adv}} = \mathbf{X} + \epsilon \text{sgn}(\nabla_{\mathbf{X}} L(\mathbf{X}, y_{\text{true}}))$$

Simple, Fast and Vicious

Test Error: 98.7%

FGSM ($\epsilon=0.1$) Error: 40.0%

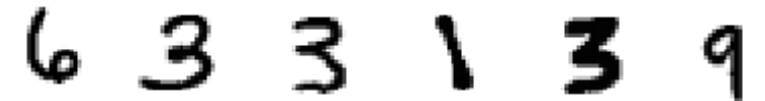
Pred: 4 Pred: 9 Pred: 9 Pred: 6 Pred: 4 Pred: 6



Pred: 2 Pred: 7 Pred: 7 Pred: 7 Pred: 2 Pred: 0



Pred: 6 Pred: 3 Pred: 3 Pred: 1 Pred: 3 Pred: 9



Pred: 4 Pred: 4 Pred: 7 Pred: 6 Pred: 7 Pred: 6



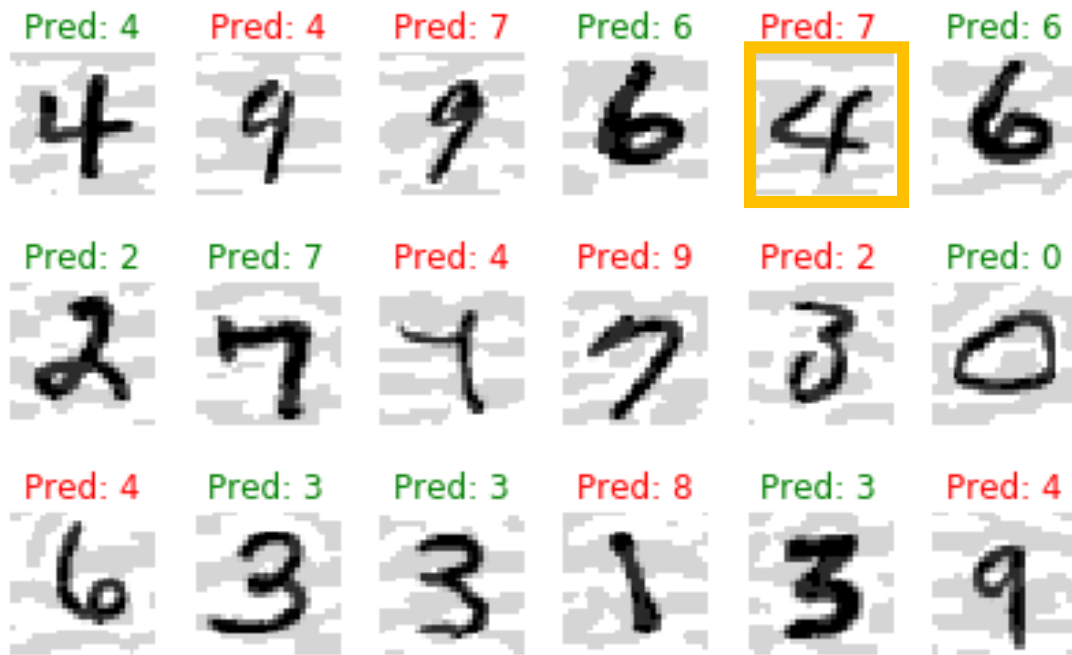
Pred: 2 Pred: 7 Pred: 4 Pred: 9 Pred: 2 Pred: 0



Pred: 4 Pred: 3 Pred: 3 Pred: 8 Pred: 3 Pred: 4



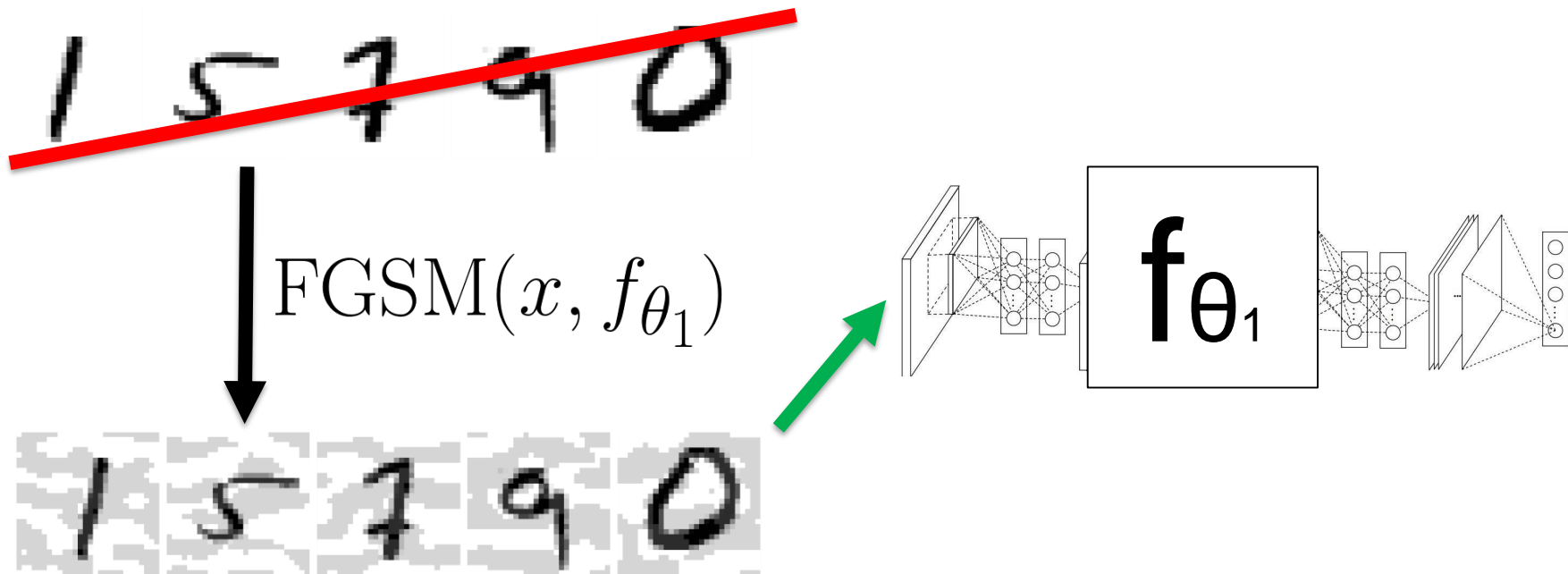
Adversarial Training



I want you to be 4!

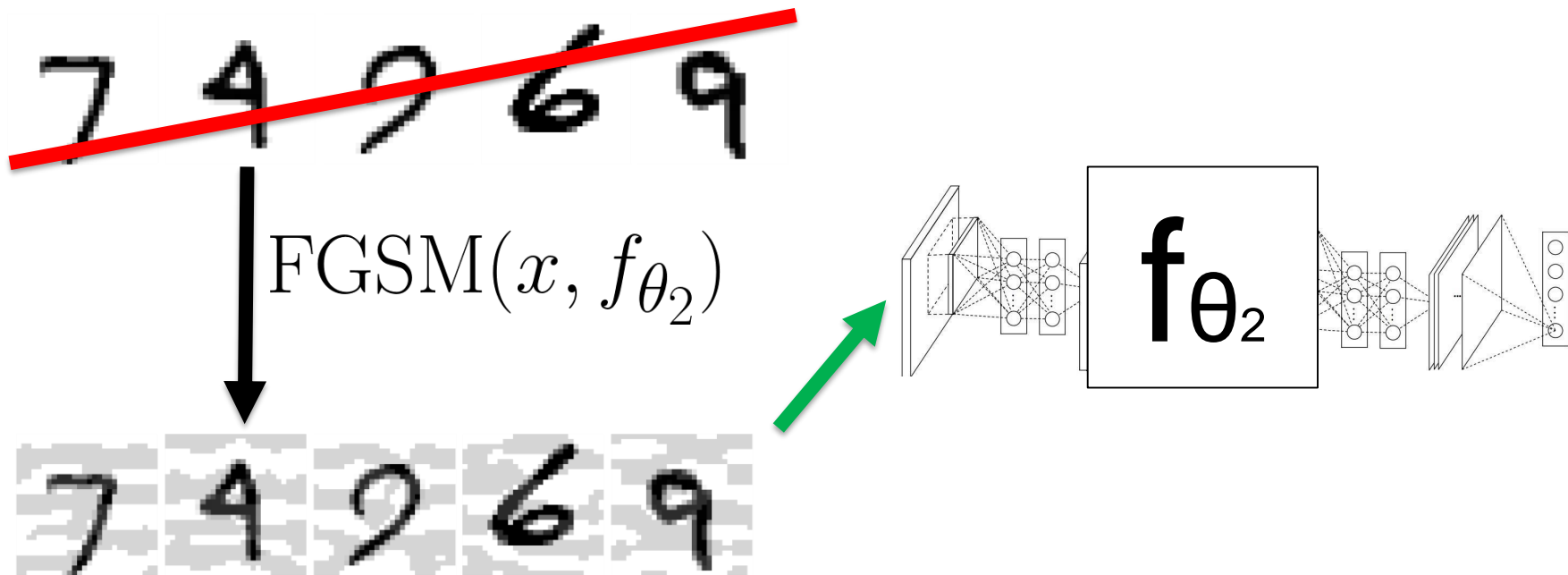
Adversarial Training

Train on adversarial examples (kind of augmentation)



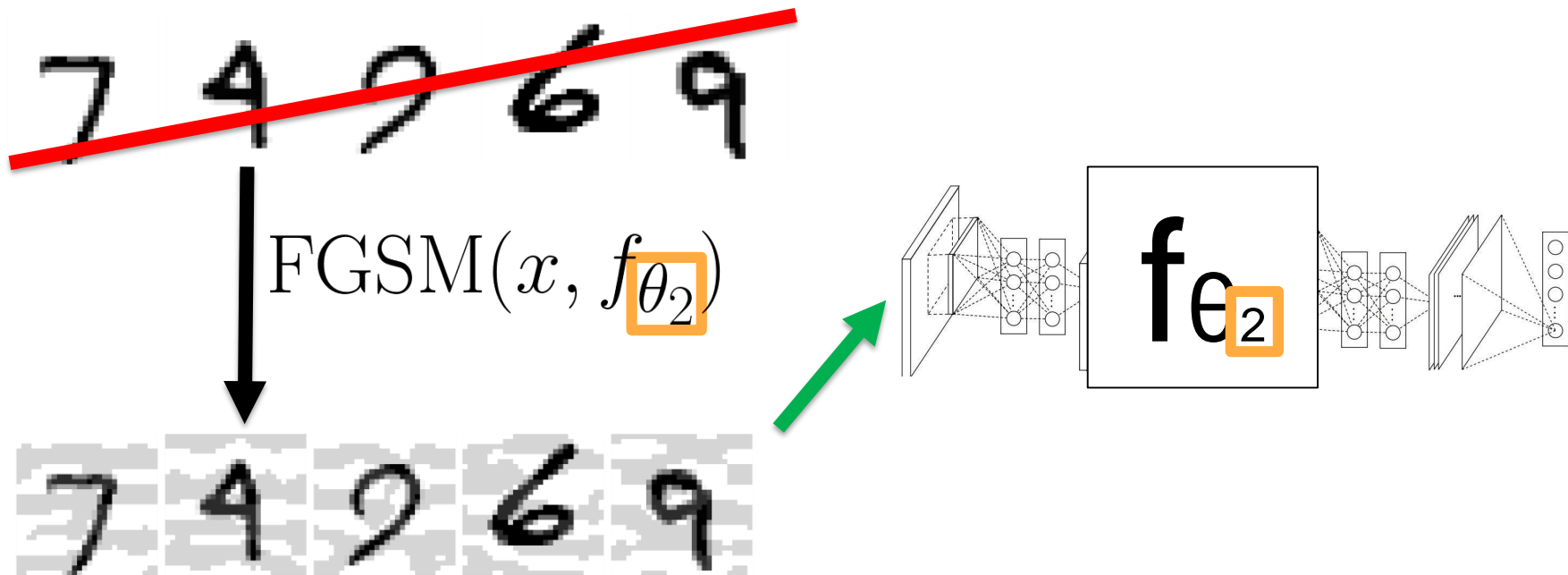
Adversarial Training

Train on adversarial examples (kind of augmentation)



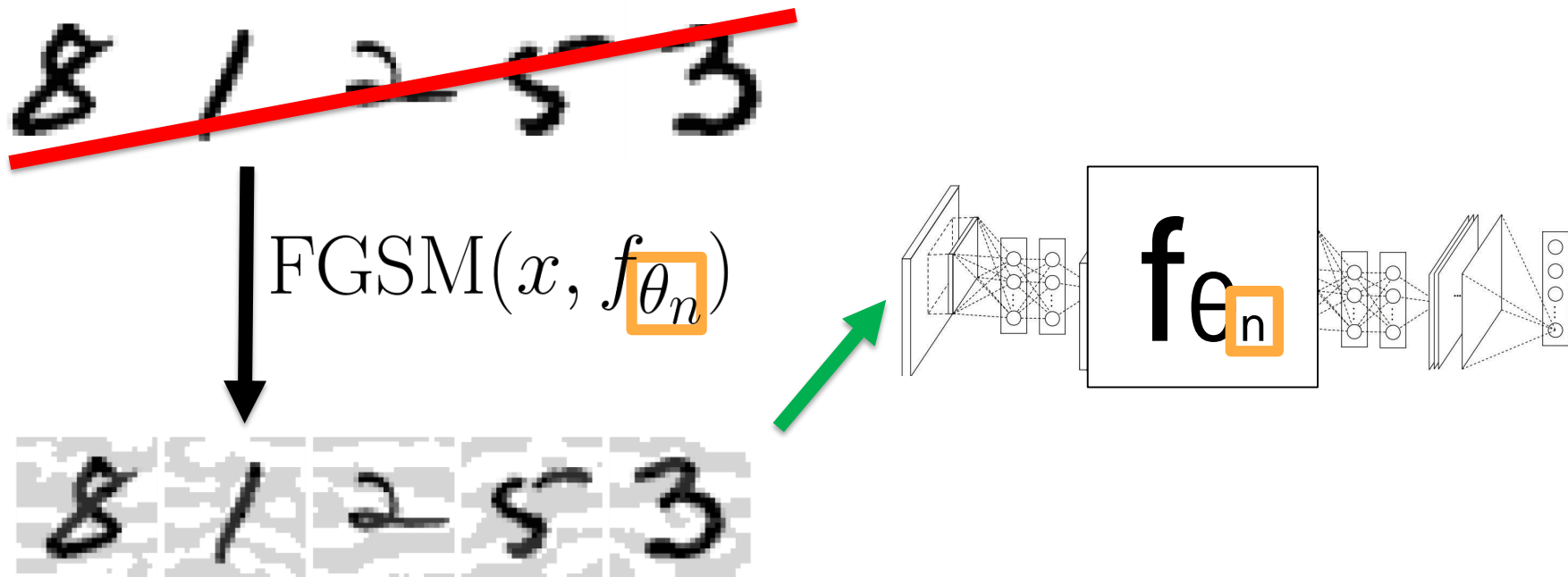
Adversarial Training

Train on adversarial examples (kind of augmentation)



Adversarial Training

Train on adversarial examples (kind of augmentation)



Adversarial Training - MNIST

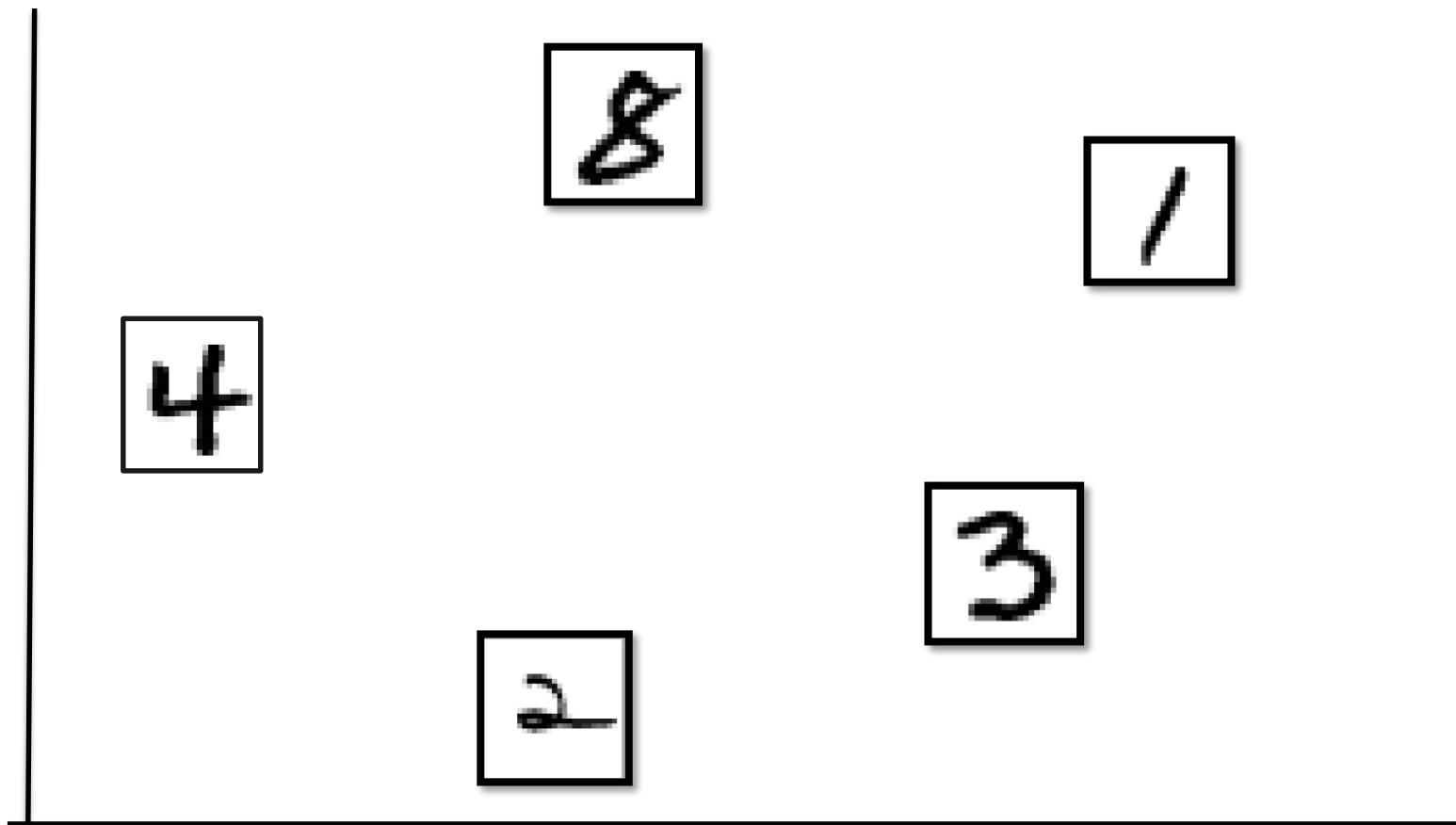
	Test Accuracy	FGSM Accuracy
Standard Training	98.7%	40.7%
Adv. Training (FGSM)	97.2%	94.0%

Did we solve the problem?

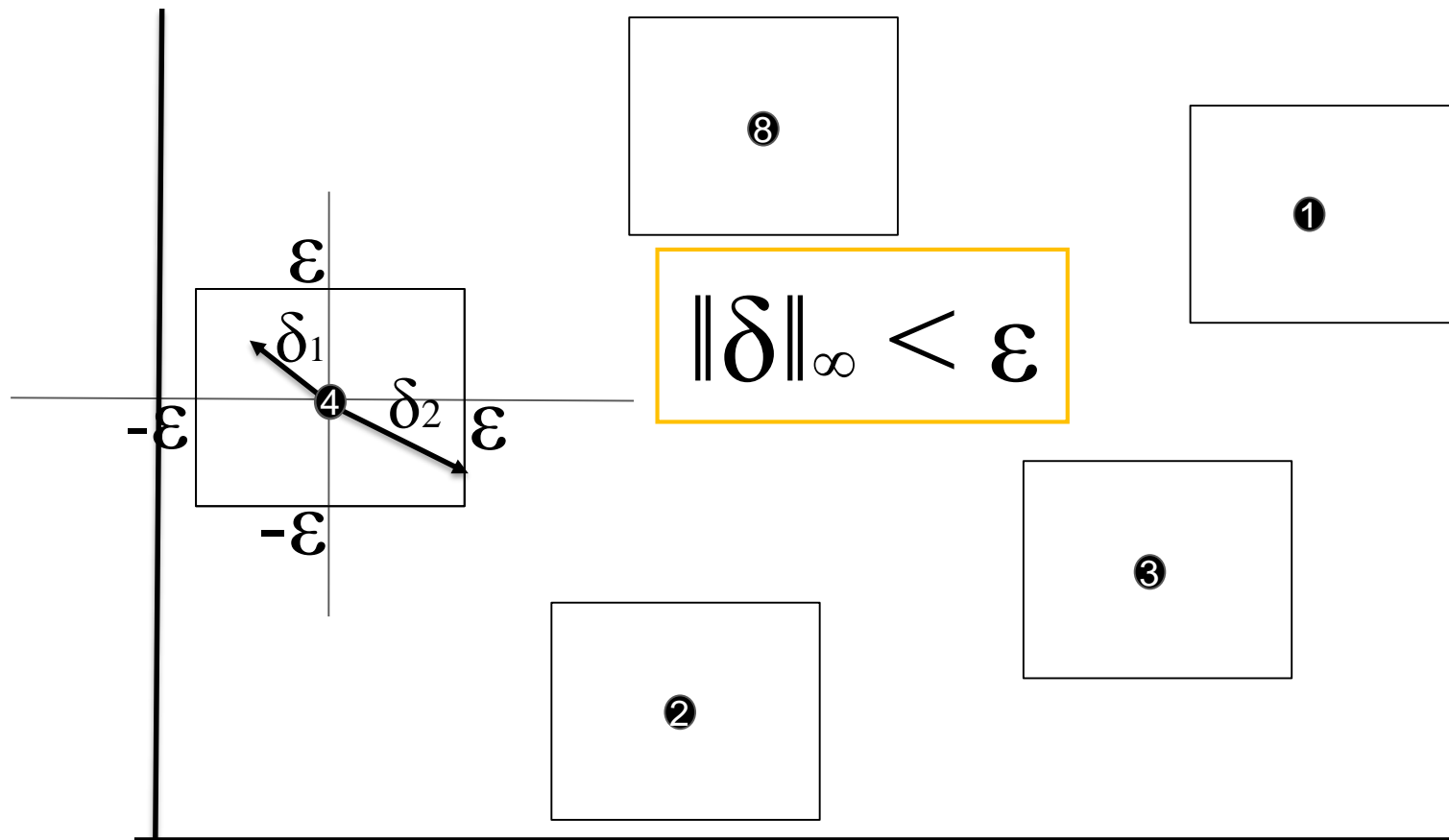
Outline

- See Adversarial Example
- Discuss what they are
- How to attack: FGSM
- How to defend: Adversarial training (AT)
 - Next: a better picture of AT (pictorially/optimization)
- Learn about properties and advantages

Perturbation Attack (pictorially)

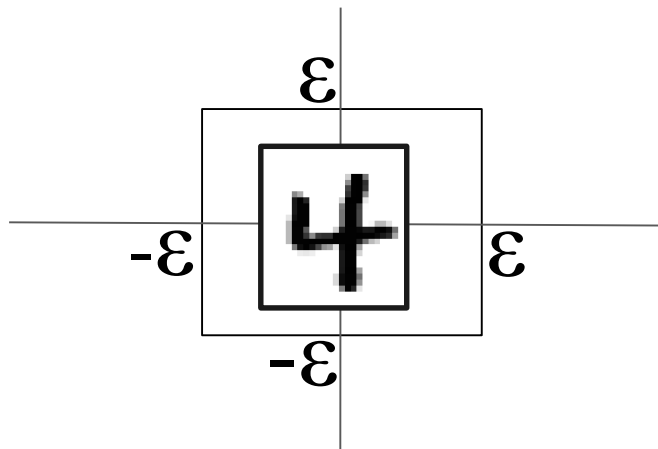


Perturbation Attack (pictorially)



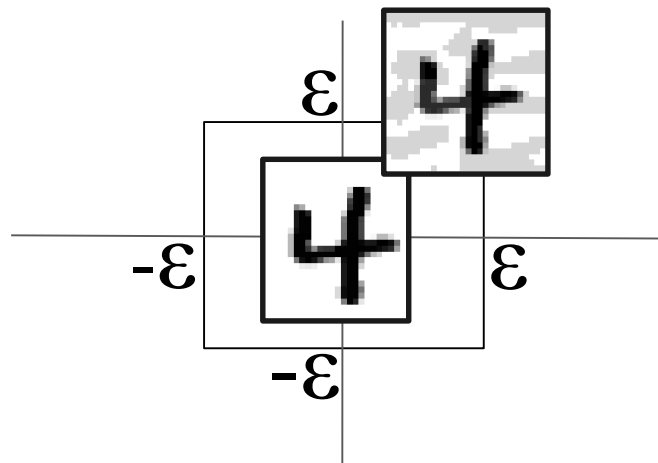
Perturbation Attack (pictorially)

FGSM



Perturbation Attack (pictorially)

FGSM

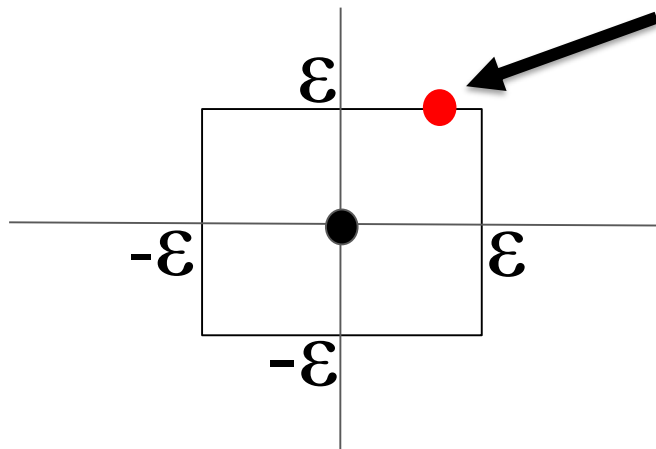


$$\mathbf{X}^{\text{adv}} = \mathbf{X} + \epsilon \text{sgn}(\nabla_{\mathbf{X}} L(\mathbf{X}, y_{\text{true}}))$$

Perturbation Attack (pictorially)

FGSM

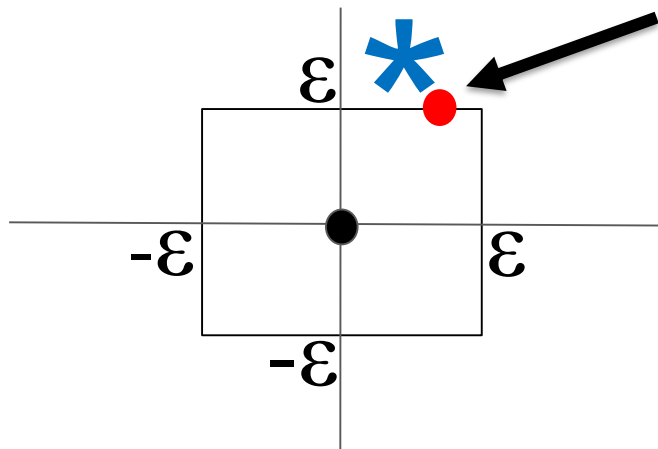
Possible AE (found by FGSM)




Perturbation Attack (pictorially)

FGSM

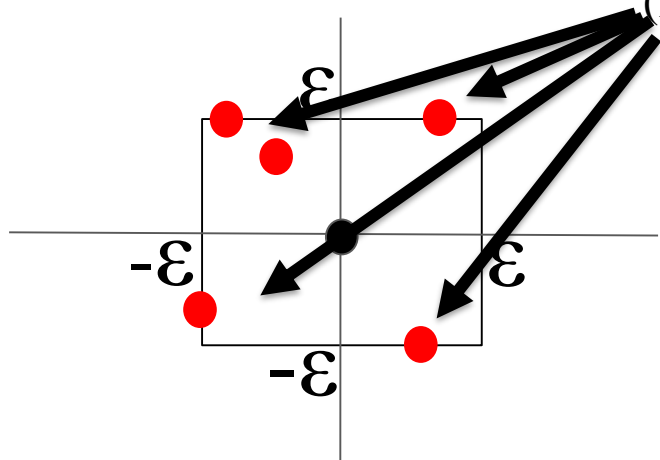
Possible AE (found by FGSM)



 dot should have been lying on one of the corners..

Perturbation Attack (pictorially)

Possible AEs
(need to be found)



“The Game” of AT:

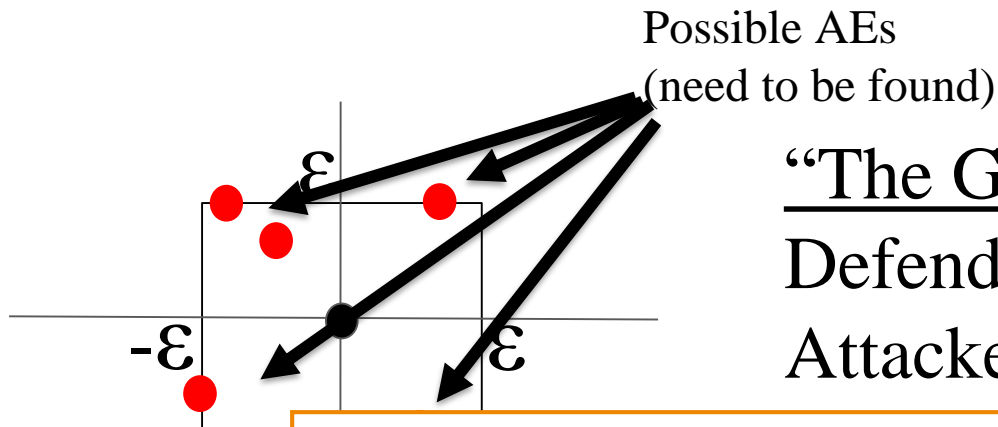
Defender: defend in box

Attacker: find AE in box

Coming
Up next:



Perturbation Attack (optimization)



“The Game” of AT:

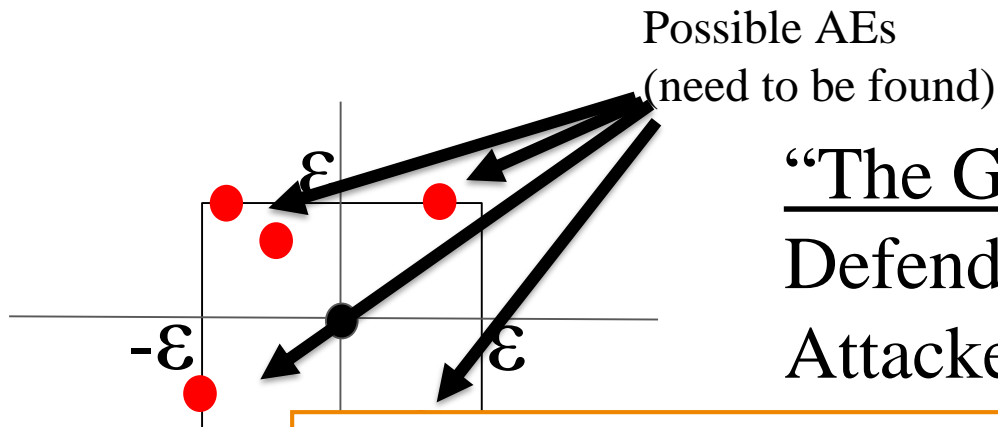
Defender: defend in box

Attacker: find AE in box

- Adversarial Training as a min-max optimization problem:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} [\max_{\delta \in S} L(f_{\theta}(x + \delta), y)]$$

Perturbation Attack (optimization)



“The Game” of AT:

Defender: defend in box

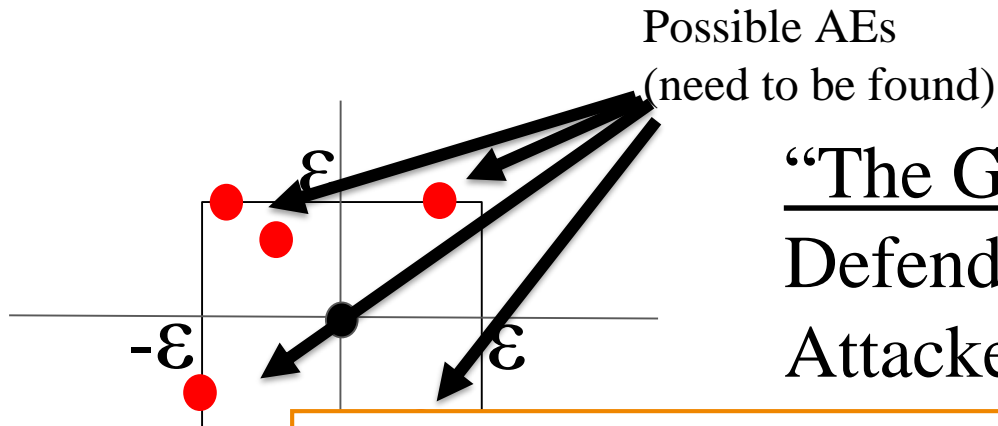
Attacker: find AE in box

- Adversarial Training as a min-max optimization problem:

Standard Loss

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} [L(f_{\theta}(x), y)]$$

Perturbation Attack (optimization)



“The Game” of AT:

Defender: defend in box

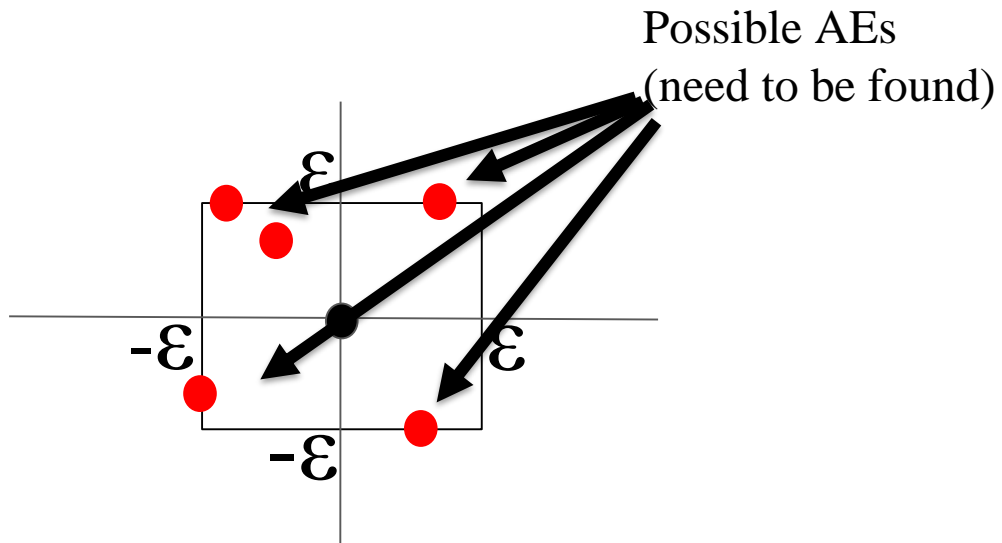
Attacker: find AE in box

- Adversarial Training as a min-max optimization problem:

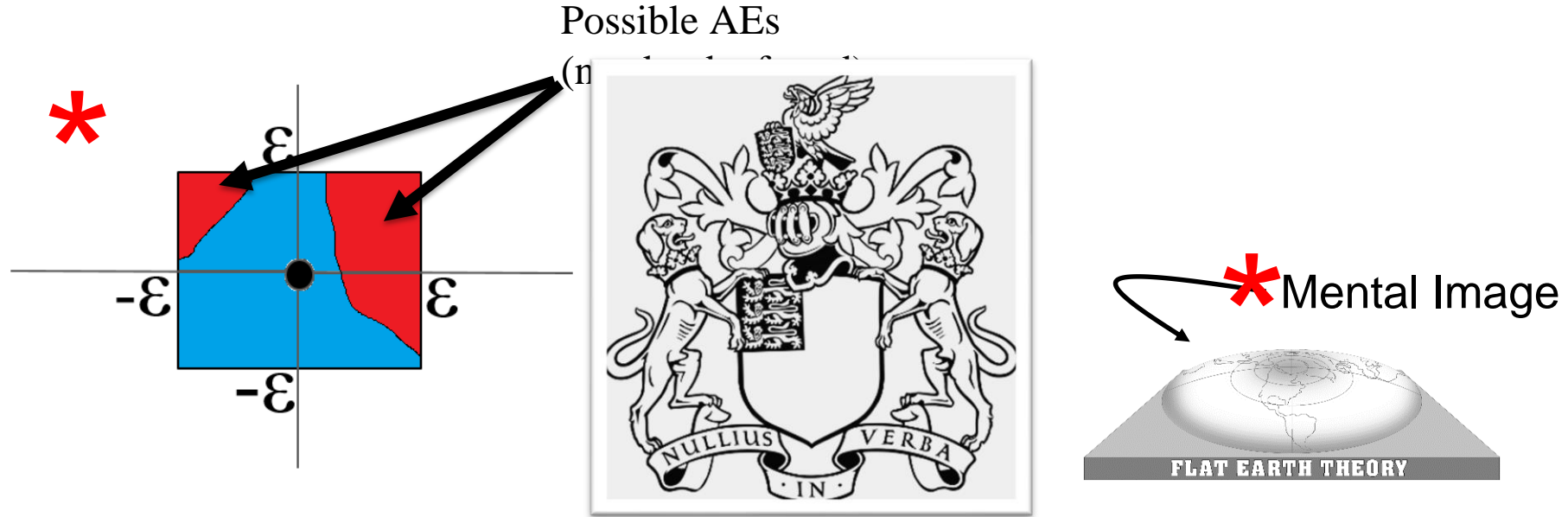
Adversarial Loss

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} \left[\max_{\delta \in S} L(f_{\theta}(x + \delta), y) \right]$$

Perturbation Attack (illustrations)



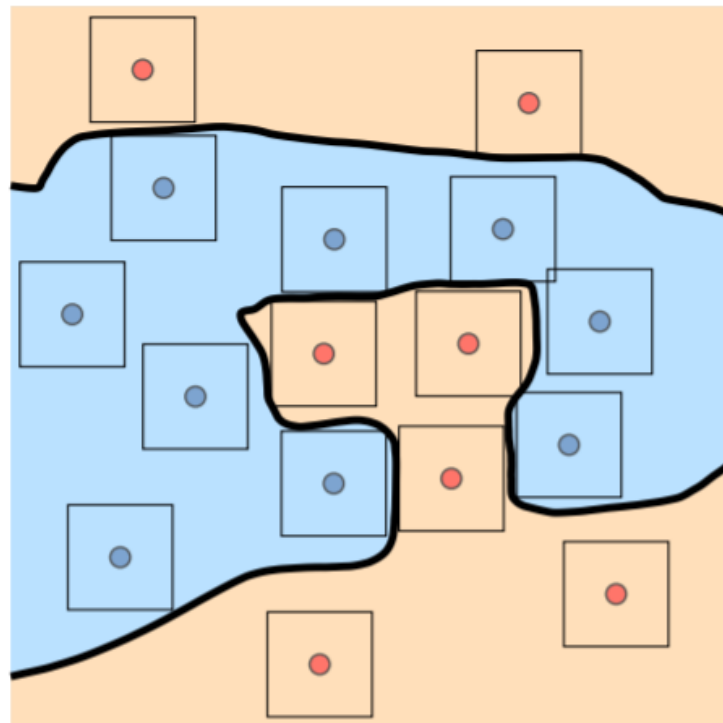
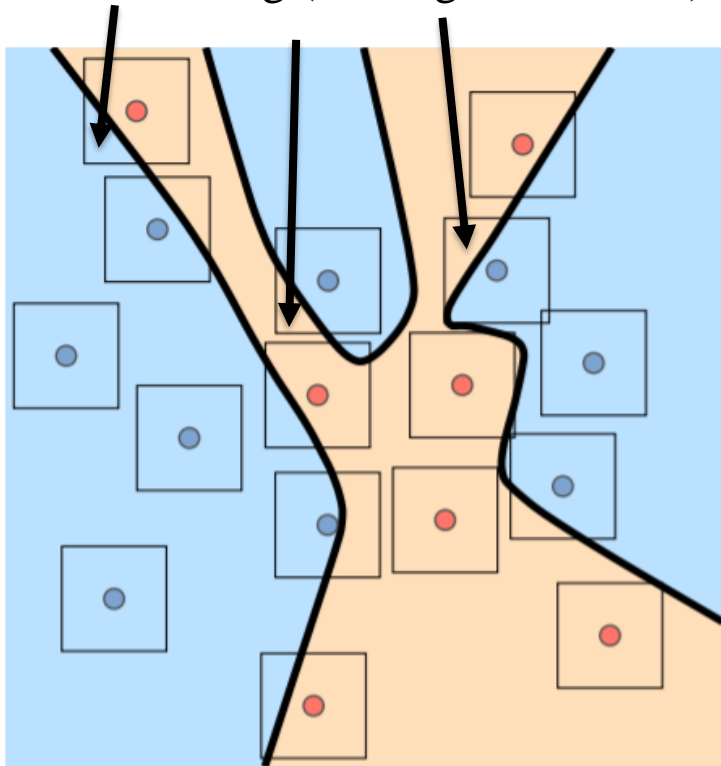
Perturbation Attack (illustrations)



Mental image alert! (“experimental” mental images could be horribly misleading)

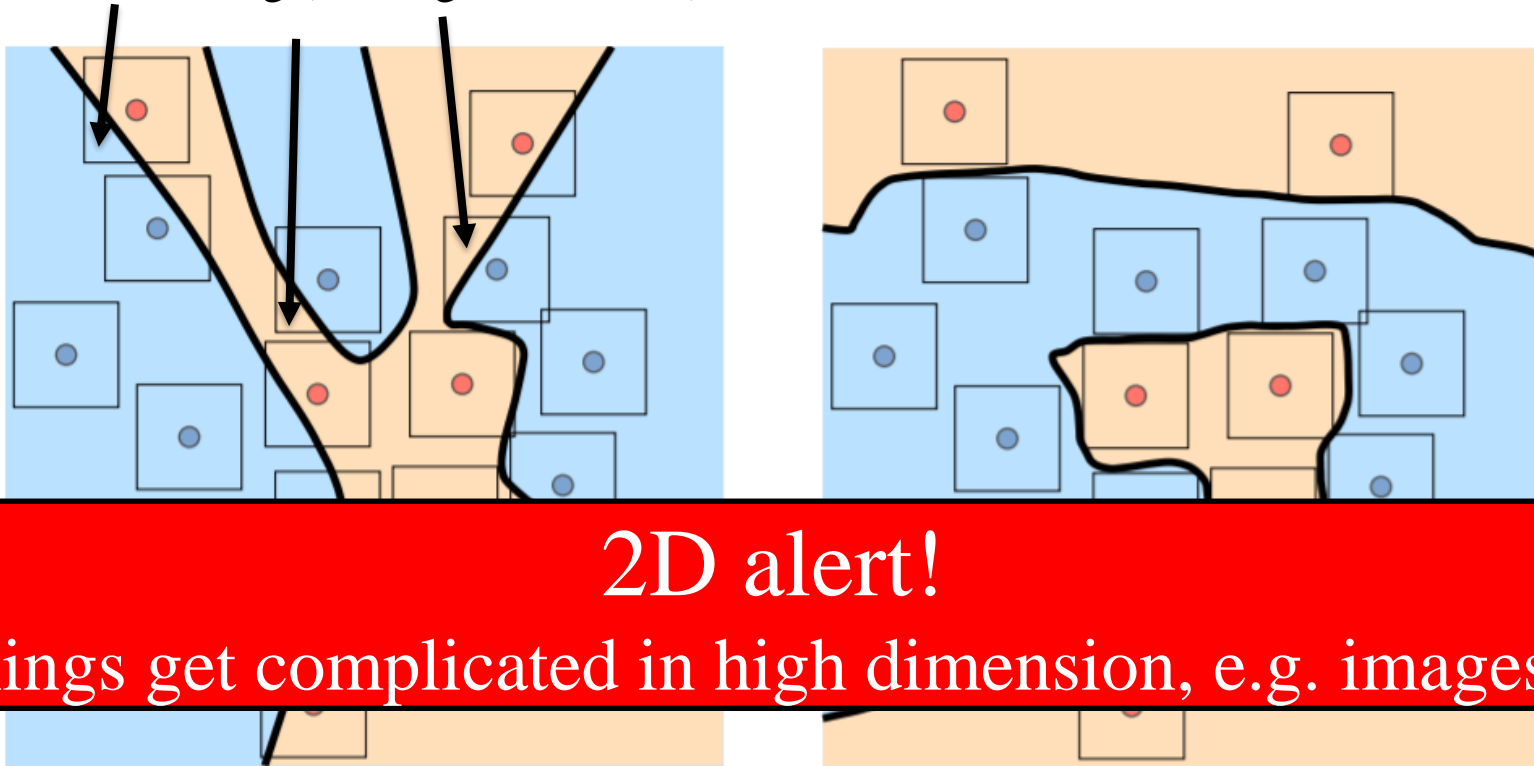
Perturbation Attack (better illustrations)

AEs lurking (waiting to be found)



Perturbation Attack (better illustrations)

AEs lurking (waiting to be found)

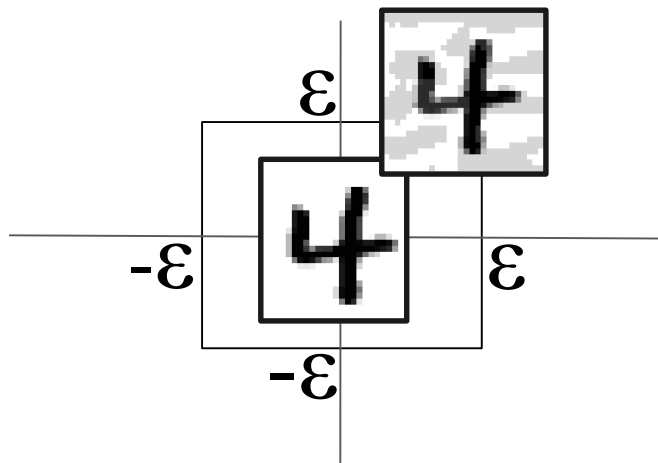


2D alert!

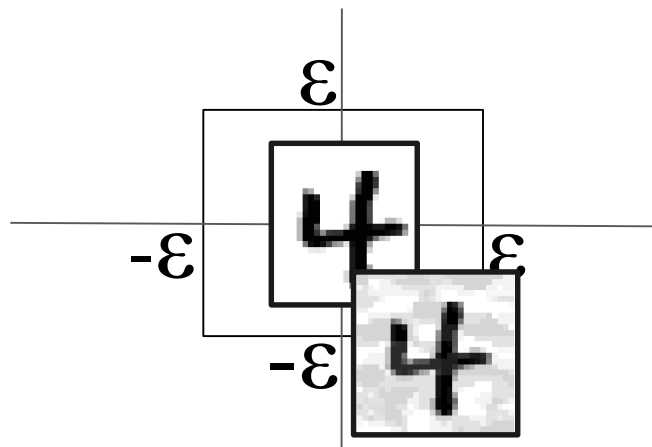
(Things get complicated in high dimension, e.g. images...)

PGD (a.k.a Iterated-GSM)

FGSM



PGD



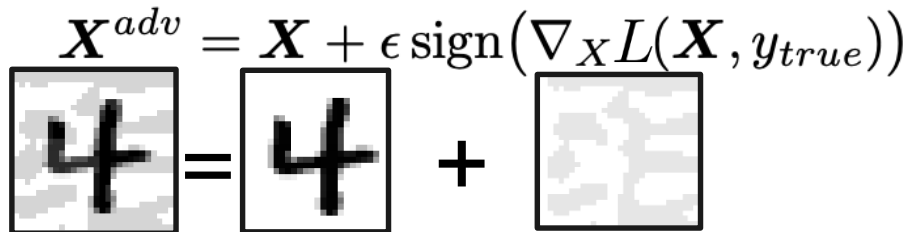
$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} [\max_{\delta \in S} L(f_{\theta}(x + \delta), y)]$$

PGD (a.k.a Iterated-GSM)

Attack Model:

$$\mathcal{S} = \{\delta \mid \|\delta\|_{\infty} < \epsilon\}$$

FGSM:

$$\mathbf{X}^{adv} = \mathbf{X} + \epsilon \text{sign}(\nabla_{\mathbf{X}} L(\mathbf{X}, y_{true}))$$


PGD (a.k.a Iterated-GSM)

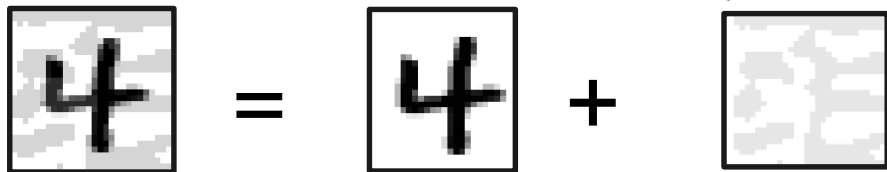
Attack Model:

$$\mathcal{S} = \{\delta \mid \|\delta\|_{\infty} < \varepsilon\}$$

PGD:

$$\mathbf{X}_0^{adv} = \mathbf{X},$$

$$\mathbf{X}_{N+1}^{adv} = \mathbf{X}_N^{adv} + \alpha \text{sign}(\nabla_{\mathbf{X}} L(\mathbf{X}_N^{adv}, y_{true}))$$



PGD (a.k.a Iterated-GSM)

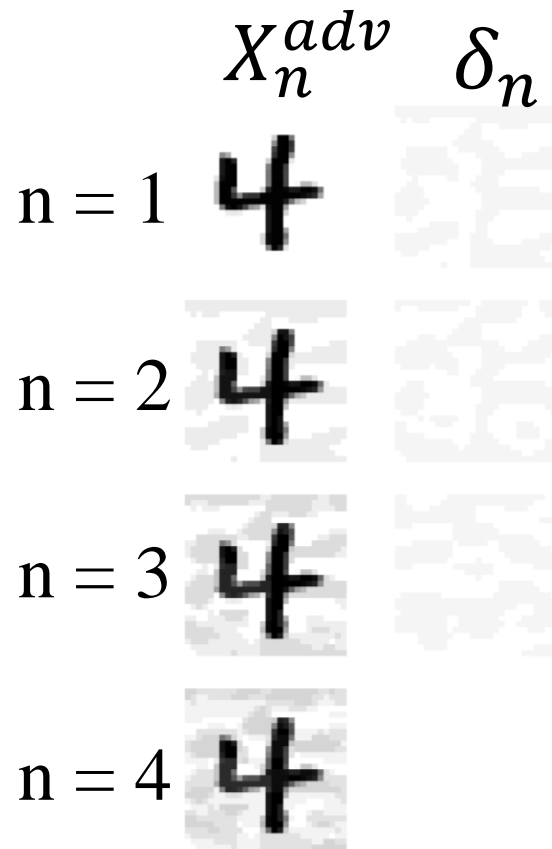
Attack Model:

$$\mathcal{S} = \{\delta \mid \|\delta\|_\infty < \epsilon\}$$

PGD:

$$\mathbf{X}_0^{adv} = \mathbf{X},$$

$$\mathbf{X}_{N+1}^{adv} = \text{Clip}_{\mathbf{X}, \epsilon} \left\{ \mathbf{X}_N^{adv} + \alpha \text{sign}(\nabla_{\mathbf{X}} L(\mathbf{X}_N^{adv}, y_{true})) \right\}$$



Adversarial Training

	Test Accuracy	FGSM Accuracy	PGD Accuracy
Standard Training	98.7%	40.7%	7.3%
Adv. Training (FGSM)	97.2%	94.0%	90.0%

What can we do to defend?

Adversarial Training

	Test Accuracy	FGSM Accuracy	PGD Accuracy
Standard Training	98.7%	40.7%	7.3%
Adv. Training (FGSM)	97.2%	94.0%	90.0%
Adv. Training (PGD)	98.0%	96.1%	95.9%

Adversarial Training – Other Datasets

CIFAR10 (ResNet50)	Test	PGD ($\epsilon = \frac{8}{255}$)
Standard Training	95.25%	0.00%
Adv. Training (PGD 8/255)	87.03%	53.29%

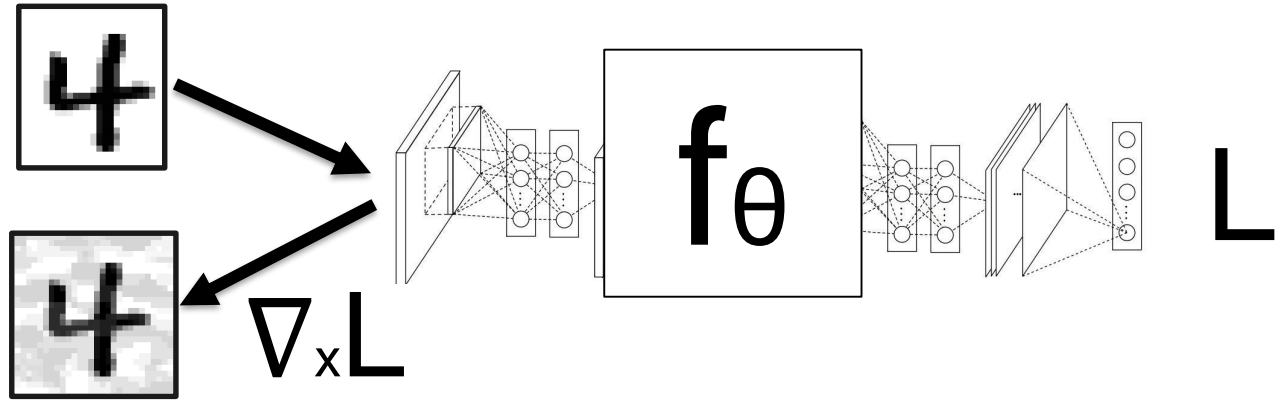
ImageNet (ResNet50)	Test	PGD ($\epsilon = \frac{8}{255}$)
Standard Training	76.13%	0.01%
Adv. Training (PGD 8/255)	47.91%	19.52%

Outline

- See Adversarial Example
- Discuss what they are
- How to attack: FGSM, PGD
- How to defend: Adversarial training (AT)
- Optimization view of AT
- Next: Black-Box attacks
- Learn about properties and advantages

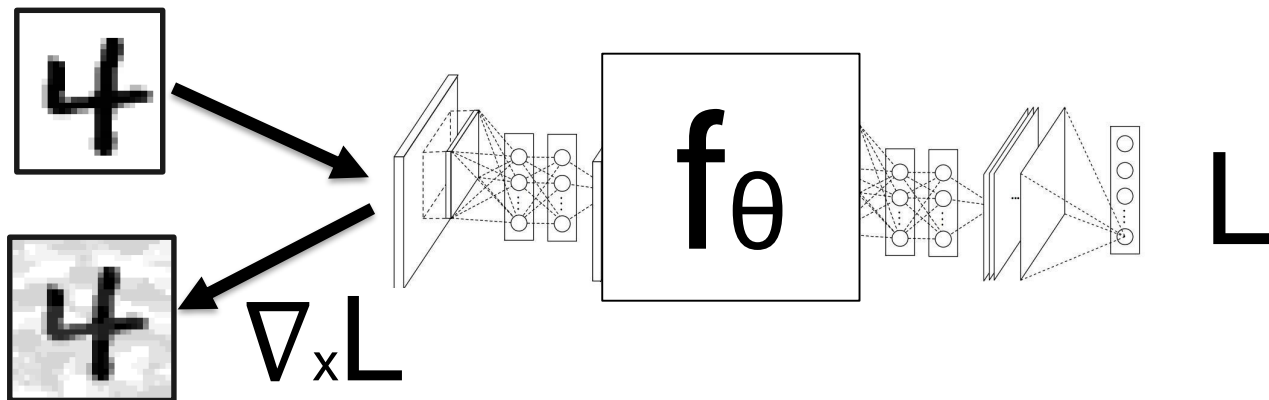
Black-Box Attacks

“White-Box”
(FGSM,
PGD, etc.)



Black-Box Attacks

“White-Box”

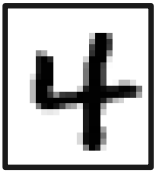


“Black-Box”



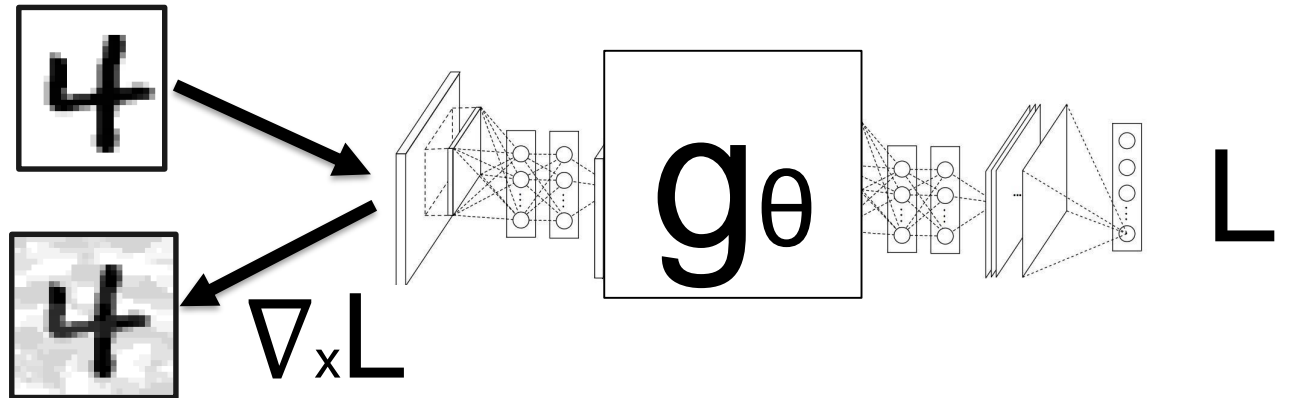
Black-Box Attacks

“Black-Box”



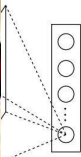
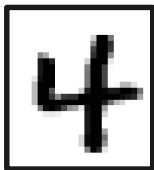
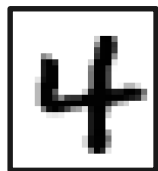
Black-Box Attacks

“Black-Box”



Black-Box Attacks

“Black-Box”



$\nabla_{\mathbf{x}} L$

Black-Box Attacks - Transferability

- Test set Accuracy

	ResNet-50	ResNet-101	ResNet-152	GoogLeNet	VGG-16
Top-5 accuracy	91.0%	91.7%	92.1%	89.0%	88.3%

- Accuracy under FGSM attack

	ResNet-152	ResNet-101	ResNet-50	VGG-16	GoogLeNet
ResNet-152	32%	55%	53%	47%	36%

Black-Box Attacks - Transferability

- Test set Accuracy

	ResNet-50	ResNet-101	ResNet-152	GoogLeNet	VGG-16
Top-5 accuracy	91.0%	91.7%	92.1%	89.0%	88.3%

- Accuracy under FGSM attack

	ResNet-152	ResNet-101	ResNet-50	VGG-16	GoogLeNet
ResNet-152	32%				
ResNet-101		33%			
ResNet-50			29%		
VGG-16				5%	
GoogLeNet					11%

White-Box
FGSM

Black-Box Attacks - Transferability

- Test set Accuracy

	ResNet-50	ResNet-101	ResNet-152	GoogLeNet	VGG-16
Top-5 accuracy	91.0%	91.7%	92.1%	89.0%	88.3%

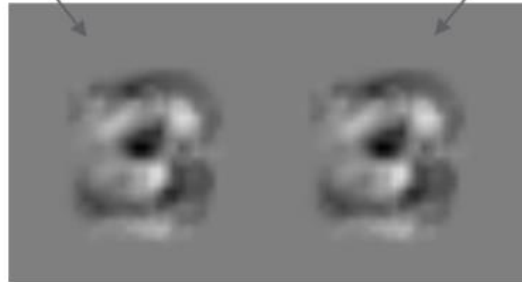
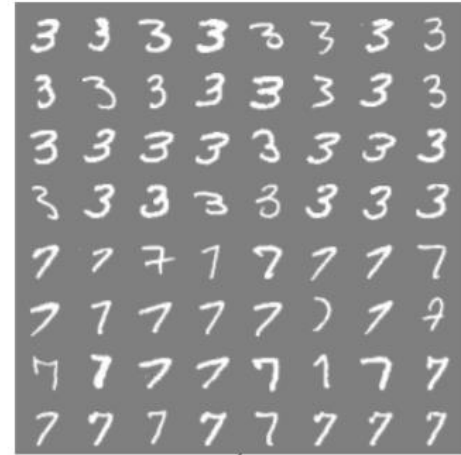
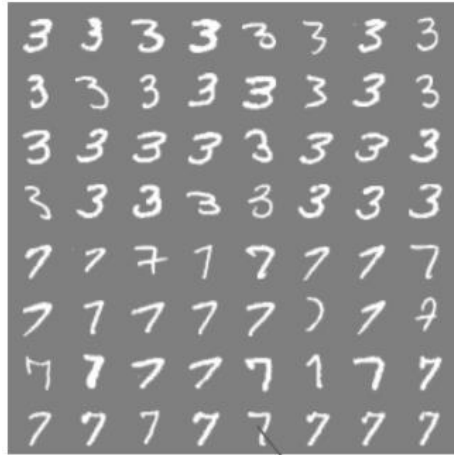
- Accuracy under FGSM attack

	ResNet-152	ResNet-101	ResNet-50	VGG-16	GoogLeNet
ResNet-152		55%	53%	47%	36%
ResNet-101	56%		50%	46%	40%
ResNet-50	59%	53%		47%	38%
VGG-16	42%	39%	41%		21%
GoogLeNet	71%	74%	62%	53%	

Black-Box

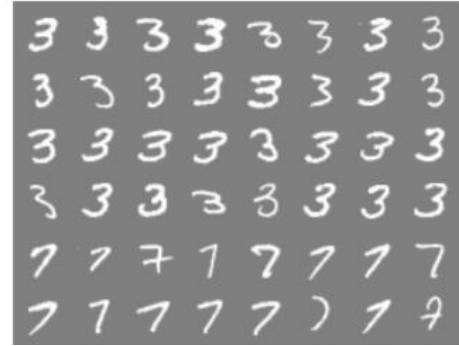
Black-Box Attacks - Transferability

- Possible reason:



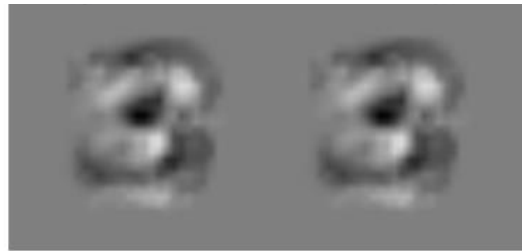
Black-Box Attacks - Transferability

- Possible reason:



Adversarial Examples comes from the data:

Ilyas et al. 2019, "Adversarial Examples Are Not Bugs, They Are Features"



Outline

- See Adversarial Example
- Discuss what they are
- How to attack: FGSM, PGD
- How to defend: Adversarial training (AT)
- Optimization view of AT
- Black-Box attacks (transferability)
- **Next: Summary**
- Surprising “advantages” of AE

Adversarial Examples – The Bigger Picture

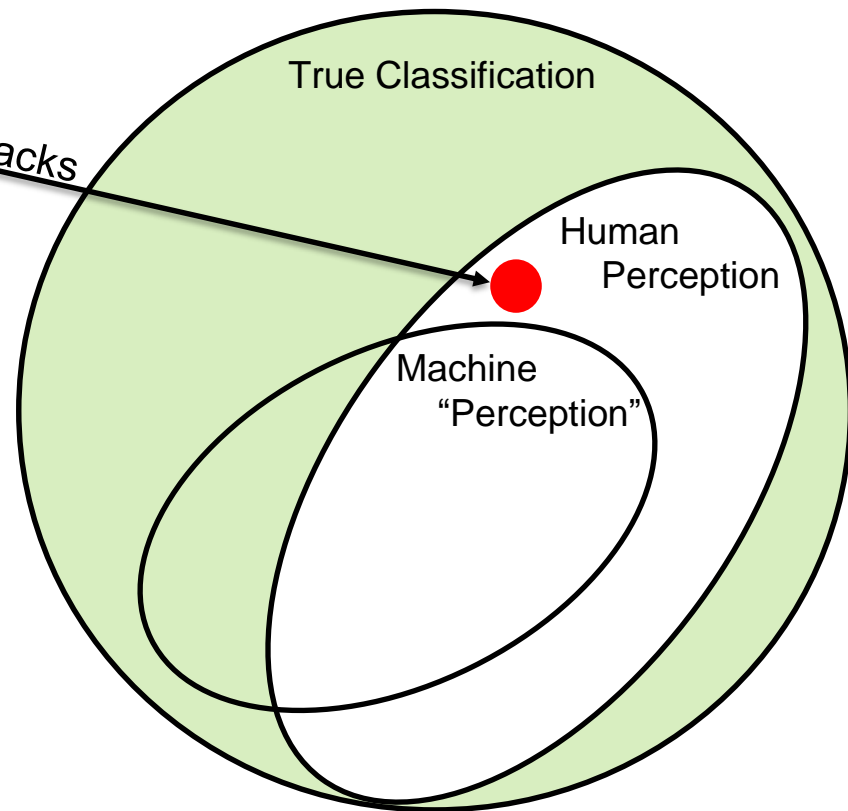
Is this surprising?

airliner



test+noise

Perturbation Attacks



Adversarial Examples – The Bigger Picture

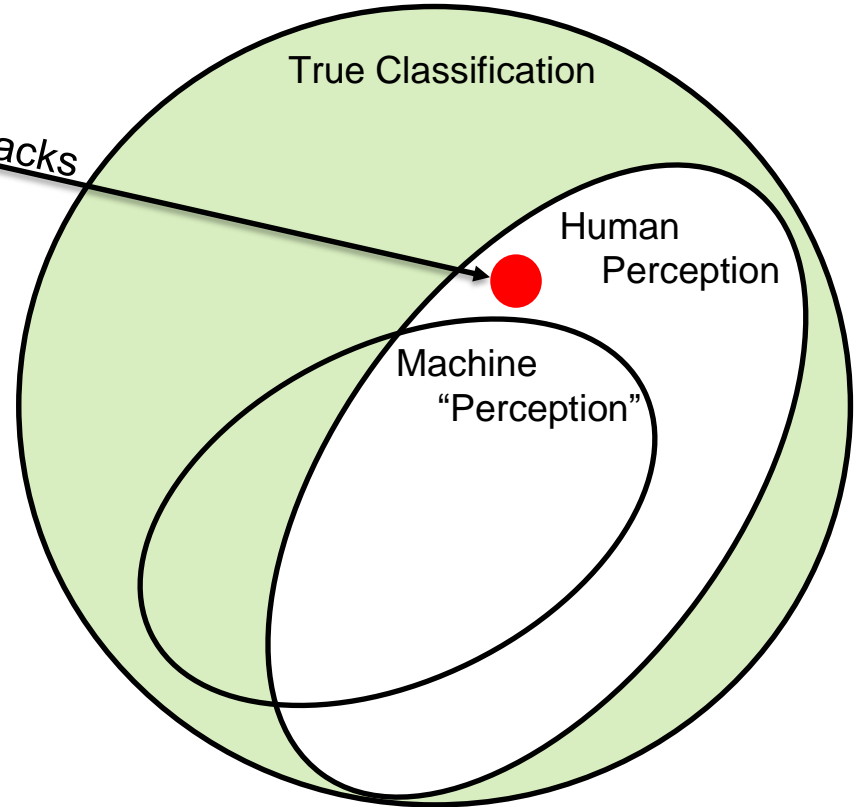
Inputs that fool a computer, but not a human

airliner



test+noise

Perturbation Attacks



Adversarial Examples – The Bigger Picture

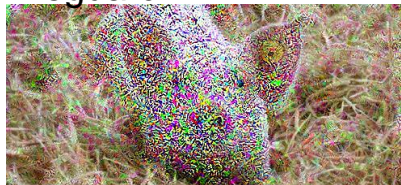
Inputs that fool a computer, but not a human

airliner

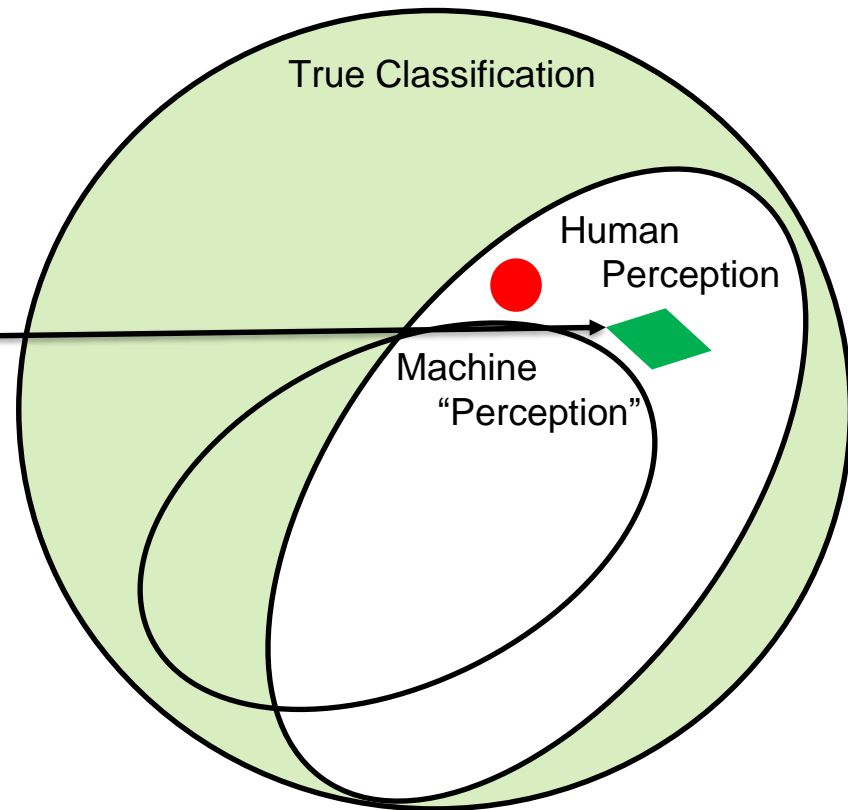


test+noise

fireguard



"noisy" image



Adversarial Examples – The Bigger Picture

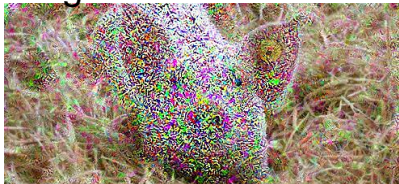
Inputs that fool a computer, but not a human

airliner



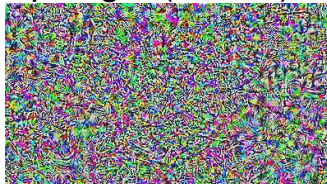
test+noise

fireguard

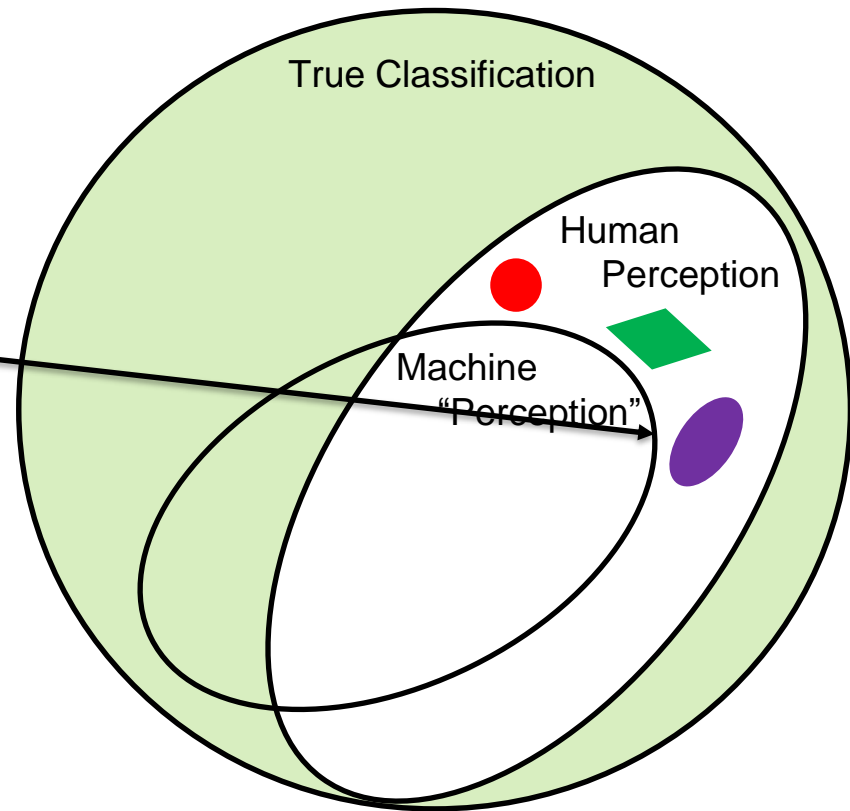


"noisy" image

spotlight (26.7%)



noise



Adversarial Examples – The Bigger Picture

Inputs that fool a computer, but not a human

airliner



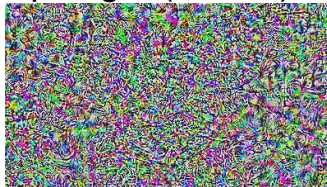
test+noise

fireguard



“noisy” image

spotlight (26.7%)

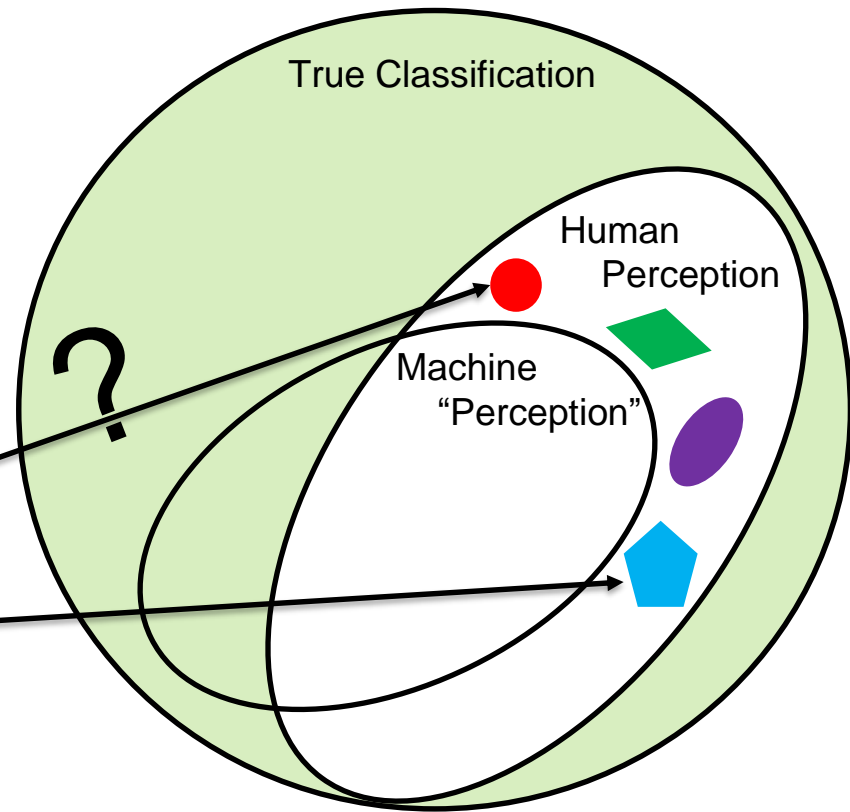


noise

cat ?



model failure



Adversarial Examples – The Bigger Picture

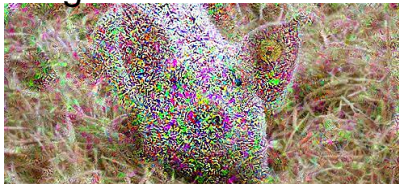
Inputs that fool a computer, but not a human

airliner



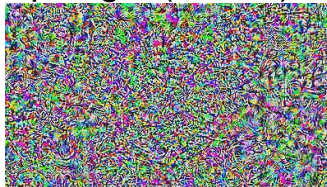
test+noise

fireguard



"noisy" image

spotlight (26.7%)



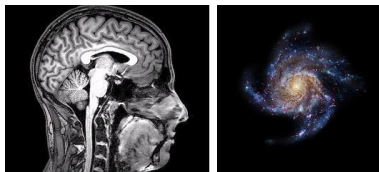
noise

cat ?

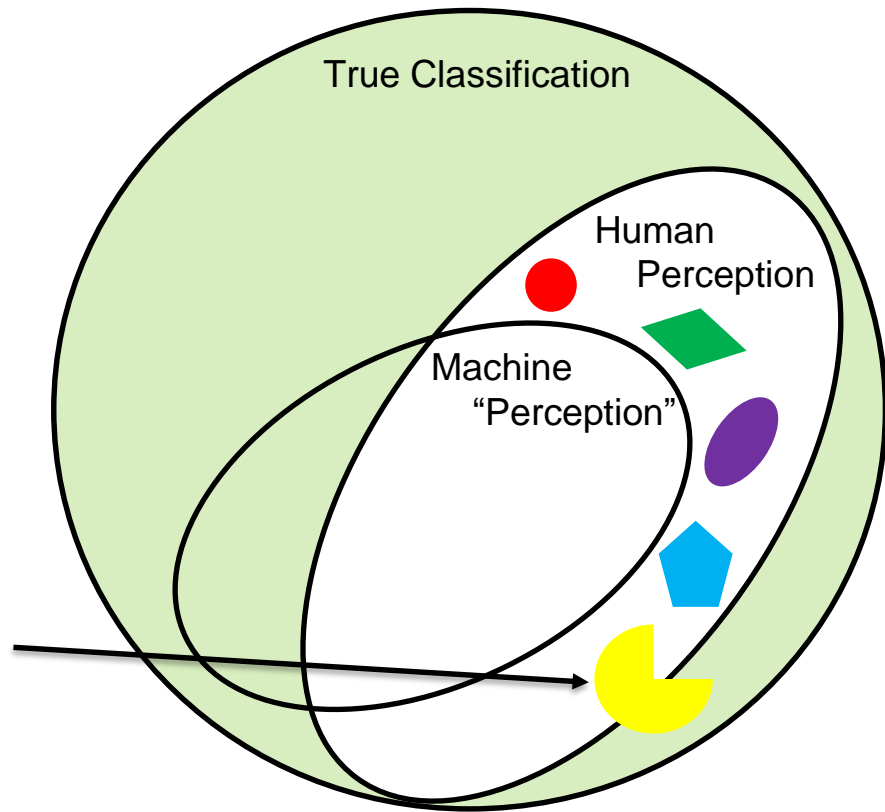


model failure

???



out-of-distribution



The Bigger Picture: Failure modes in machine learning

Intentionally-motivated failures

Unintended failures

Attack

Overview

Perturbation attack Attacker modifies the query to get appropriate response

Poisoning attack Attacker contaminates the training phase of ML systems to get intended result

Model Inversion Attacker recovers the secret features used in the model by through careful queries

Membership Inference Attacker can infer if a given data record was part of the model's training dataset or not

Model Stealing Attacker is able to recover the model through carefully-crafted queries

Reprogramming ML system Repurpose the ML system to perform an activity it was not programmed for

Adversarial Example in Physical Domain Attacker brings adversarial examples into physical domain to subvert ML system e.g: 3d printing special eyewear to fool facial

Failure

Overview

Reward Hacking Reinforcement Learning (RL) systems act in unintended ways because of mismatch between state reward and true reward

Side Effects RL system disrupts the environment as it tries to attain its goal

Distributional shifts The system is tested in one kind of environment, but is unable to adapt to changes in other kinds of environment

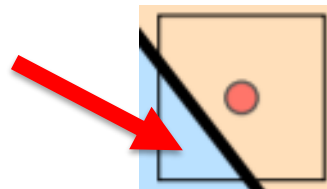
Natural Adversarial Examples Without attacker perturbations, the ML system fails owing to hard negative mining

Common Corruption The system is not able to handle common corruptions and perturbations such as tilting, zooming noisy images.

Incomplete Testing The ML system is not tested in the realistic conditions that it is meant to operate in.

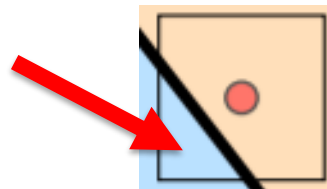
Adversarial Examples - Summary

- Remember the bigger picture (many failures)
- Hard to attack (need to find AE in box)
- Harder to defend (need to prove: no AEs in all box)



Adversarial Examples - Summary

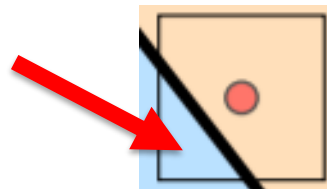
- Remember the bigger picture (many failures)
- Hard to attack (need to find AE in box)



- Harder to defend (need to **prove**: no AEs in all box)

Adversarial Examples - Summary

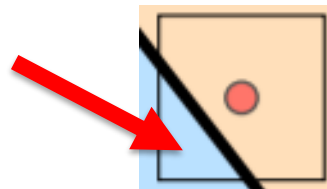
- Remember the bigger picture (many failures)
- Hard to attack (need to find AE in box)



- Harder to defend (need to **prove**: very hard to find AE in box)

Adversarial Examples - Summary

- Remember the bigger picture (many failures)
- Hard to attack (need to find AE in box)



- Harder to defend (need to **Evaluate**: very hard to find AE in box)
- Coming next: Robustness beyond security

Outline

- See Adversarial Example
- Discuss what they are
- How to attack: FGSM, PGD
- How to defend: Adversarial training (AT)
- Optimization view of AT
- Black-Box attacks (transferability)
- Summary (“security”)
- Surprising “advantages” of AE (beyond security)

Follow the gradient w.r.t x (the input image)



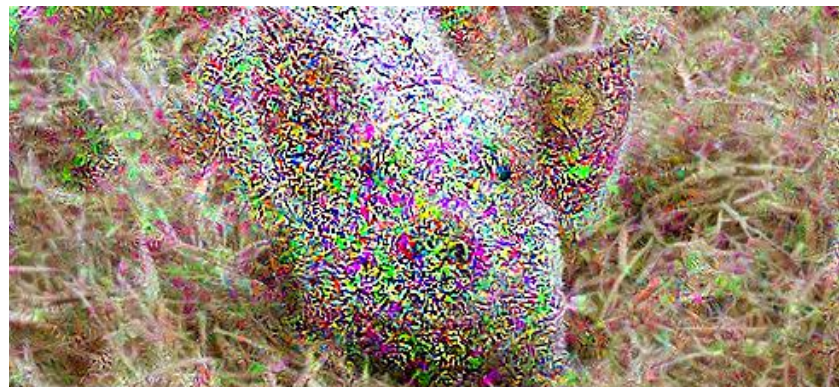
X (original): 89.7% pig



$X + \nabla_x L$: 68.6% hay



$X + 10 \times \nabla_x L$: 44.7% pig



$X + 100 \times \nabla_x L$: 44.8% fireguard

Follow $\nabla_x L(f(x), y)$ of Robust Model

Original



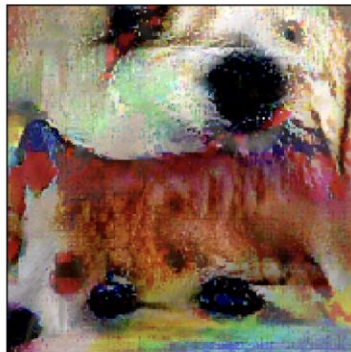
primate

Standard



dog

l_∞ -trained



dog

l_2 -trained



dog



bird



turtle



dog



cat

Image synthesis with Robust Classifier

cliff



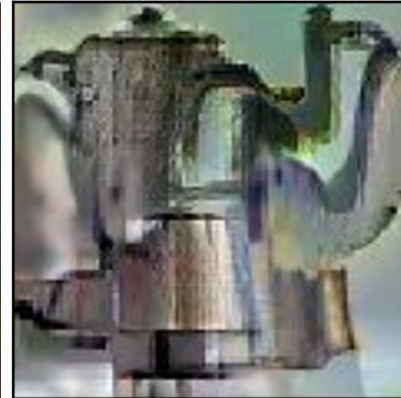
anemone fish



mashed potato



coffee pot



house finch



armadillo



chow



jigsaw



Norwich terrier

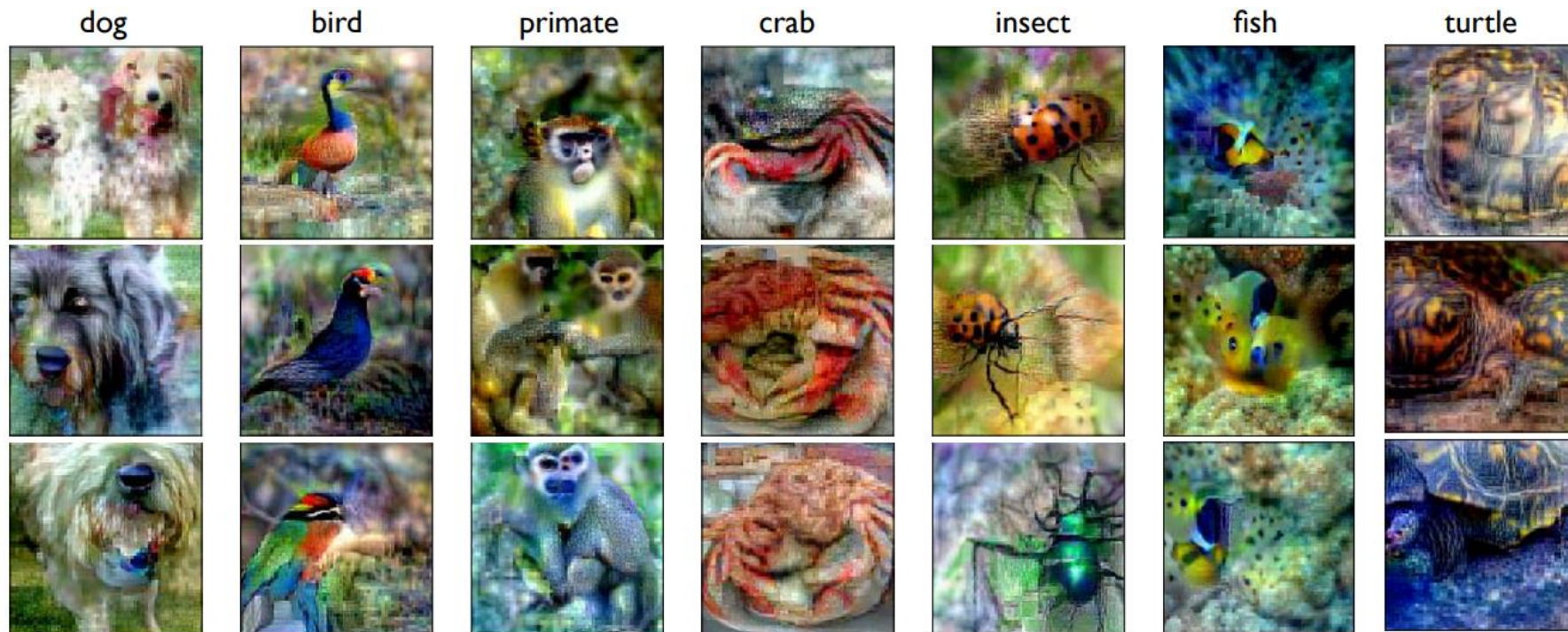


notebook



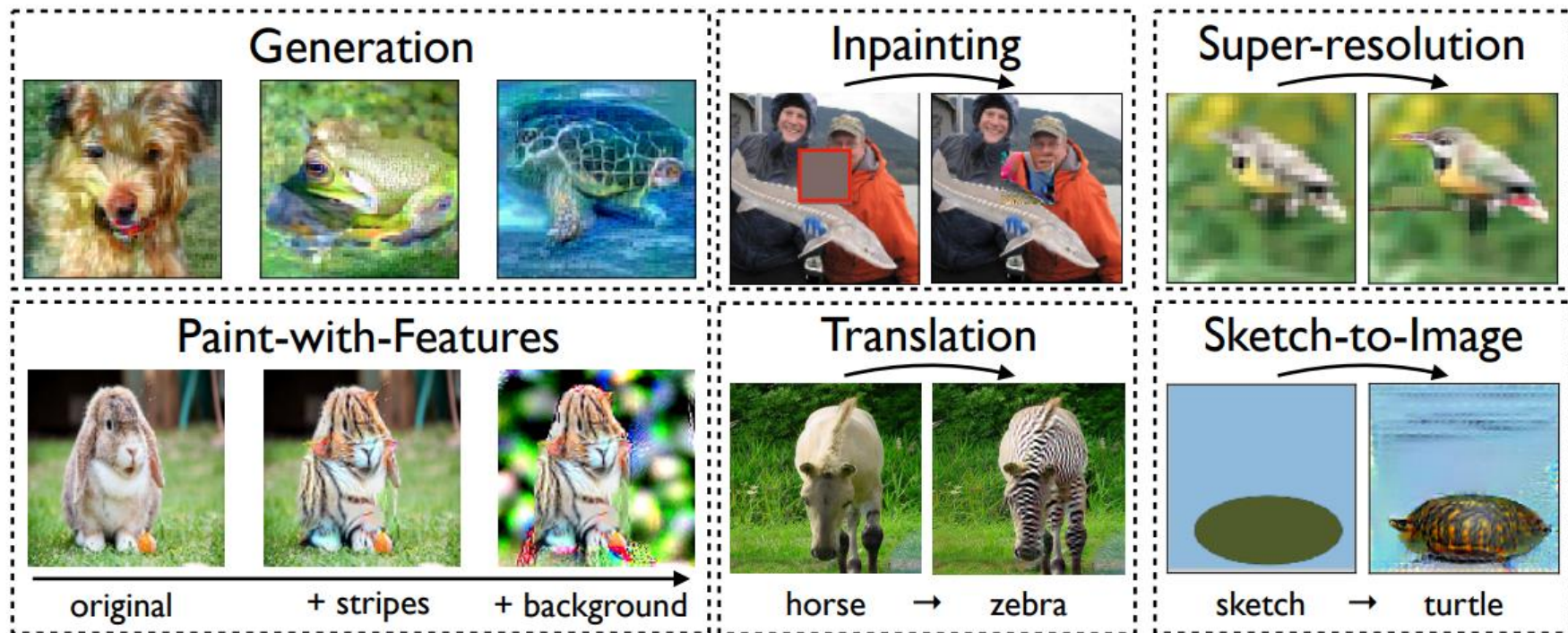
Santurkar et al. 2019, "Image Synthesis with a Single (Robust) Classifier"

Image synthesis with Robust Classifier



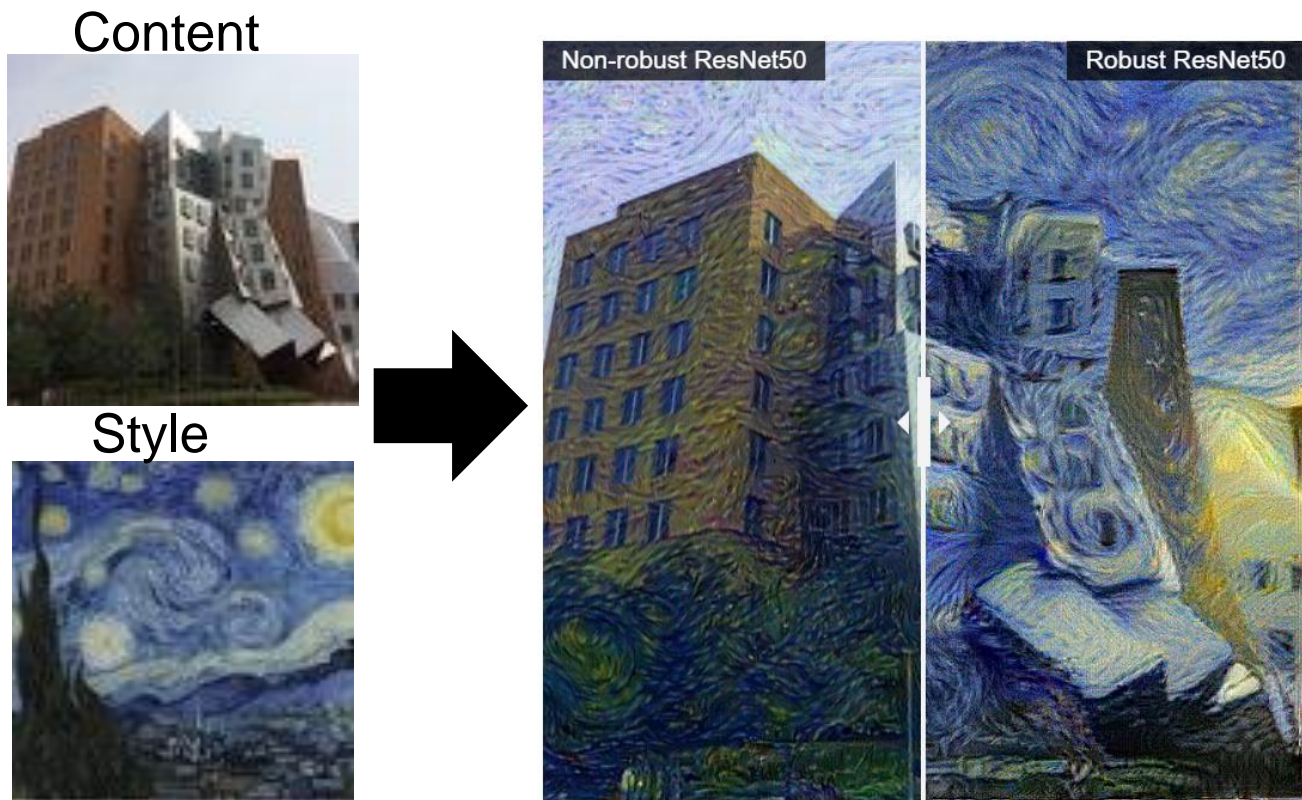
Santurkar et al. 2019, "Image Synthesis with a Single (Robust) Classifier"

Image synthesis with Robust Classifier



Santurkar et al. 2019, "Image Synthesis with a Single (Robust) Classifier"

Style Transfer with Robust Model





What have we learnt today?

- Saw a few Adversarial Examples
- Discussed what they are
- How to attack: FGSM, PGD
- How to “defend”: Adversarial training (AT)
- Optimization view of AT
- Black-Box attacks (transferability)
- Security-wise summary
- Surprising Visual properties of robust models (beyond security)

Monday:



Sequences
(RNN, Attention, ViT)

25/11/2021

Defence Against the Adversarial Examples



Niv Haim

Weizmann Institute

DL4CV Course Winter 2022 (20224182)