# CS 499/579: TRUSTWORTHY ML ADVERSARIAL EXAMPLES: WHITE-BOX ATTACKS

Instructor: Sanghyun Hong

sanghyun.hong@oregonstate.edu





#### NOTES

- Call for actions
  - Homework 1 due
  - In-class presentation sign-ups
  - Term project team-up (by today)



- Research questions
  - What are the adversarial examples?
  - How can we find adversarial examples?
  - How can we exploit them in practice?
  - How can we defeat adversarial examples?



# How can we train models robust to adversarial inputs?

TOWARDS DEEP LEARNING MODELS RESISTANT TO ADVERSARIAL ATTACKS, MADRY ET AL., ICLR 2018

# HOW DID THE RESEARCH GO?

- Many attack proposals
  - FGSM
  - JSMA
  - DeepFool
  - DeepXplore<sup>1</sup>
  - C&W
  - ...
- Many defense proposals
  - Regularization ... broken
  - Defensive distillation ... broken
  - Adversarial training ... but with which attack?

- ...



Pei et al., DeepXplore: Automated Whitebox Testing of Deep Learning Systems, SOSP 2017

- Main research question
  - How can we train neural networks robust to adversarial examples?



## **R**EVISITING THE FORMULATION

- Test-time (evasion) attack
  - Suppose
    - A test-time input (*x*, *y*)
    - $(x, y) \sim D$ , D: data distribution;  $x \in \mathbb{R}^d$  and  $y \in [k]$ ;  $x \in [0, 1]$
    - A NN model f and its parameters heta
    - $L(\theta, x, y)$ : a loss function
  - Objective
    - Find an  $x^{adv} = x + \delta$  such that  $f(x^{adv}) \neq y$  while  $||\delta||_p \leq \varepsilon$



## **REVISITING THE FORMULATION**

- Test-time (evasion) attack
  - Suppose
    - A test-time input (*x*, *y*)
    - $(x, y) \sim D$ , D: data distribution;  $x \in \mathbb{R}^d$  and  $y \in [k]$ ;  $x \in [0, 1]$
    - A NN model f and its parameters  $\theta$
    - $L(\theta, x, y)$ : a loss function
  - Attacker's objective
    - Find an  $x^{adv} = x + \delta$  such that  $\max_{\delta \in S} L(\theta, x^{adv}, y)$  while  $||\delta||_p \le \varepsilon$



## **R**EVISITING THE FORMULATION

- Test-time (evasion) attack
  - Suppose
    - A test-time input (*x*, *y*)
    - $(x, y) \sim D$ , D: data distribution;  $x \in \mathbb{R}^d$  and  $y \in [k]$ ;  $x \in [0, 1]$
    - A NN model f and its parameters heta
    - $L(\theta, x, y)$ : a loss function
  - Attacker's objective
    - Find an  $x^{adv} = x + \delta$  such that  $\max_{\delta \in S} L(\theta, x^{adv}, y)$  while  $||\delta||_p \le \varepsilon$
  - Defender's objective
    - Train a neural network *f* robust to adversarial attacks
    - Find  $\theta$  such that  $\min_{\theta} \rho(\theta)$  where  $\rho(\theta) = \mathbb{E}_{(x,y)\sim D} [L(\theta, x^{adv}, y)]$



# **PUTTING ALL TOGETHER**

- (Models resilient to) test-time (evasion) attack
  - Suppose
    - A test-time input (*x*, *y*)
    - $(x, y) \sim D$ , D: data distribution;  $x \in \mathbb{R}^d$  and  $y \in [k]$ ;  $x \in [0, 1]$
    - A NN model f and its parameters heta
    - $L(\theta, x, y)$ : a loss function
  - Min-max optimization (between attacker's and defender's objectives)
    - Find  $\min_{\theta} \rho(\theta)$  where  $\rho(\theta) = \mathbb{E}_{(x,y)\sim D} \left[ \max_{\delta \in S} L(\theta, x + \delta, y) \right]$  while  $||\delta||_p \le \varepsilon$
    - s: a set of test-time samples

#### SADDLE POINT PROBLEM: INNER MAXIMIZATION AND OUTER MINIMIZATION



#### **INNER MAXIMIZATION USING THE FIRST-ORDER ADVERSARY**

• Revisit FGSM (Fast Gradient Sign Method)

```
x + \varepsilon \operatorname{sgn}(\nabla_x L(\theta, x, y)).
```

- FGSM can be viewed as a simple one-step toward maximizing the loss (inner part)



#### **INNER MAXIMIZATION**

• Revisit FGSM (Fast Gradient Sign Method)

 $x + \varepsilon \operatorname{sgn}(\nabla_x L(\theta, x, y)).$ 

- FGSM can be viewed as a simple one-step toward maximizing the loss (inner part)
- PGD (Projected Gradient Descent)

$$x^{t+1} = \Pi_{x+S} \left( x^t + \alpha \operatorname{sgn}(\nabla_x L(\theta, x, y)) \right).$$
FGSM

- Multi-step adversary; much stronger than FGSM attack



• PGD (Projected Gradient Descent)

$$x^{t+1} = \Pi_{x+\mathcal{S}} \left( x^t + \alpha \operatorname{sgn}(\nabla_x L(\theta, x, y)) \right).$$

- Multi-step adversary; much stronger than FGSM attack
- Hyper-parameters
  - *t*: number of iterations
  - *α*: step-size
  - $\varepsilon$ : perturbation bound  $|x^* x|_p$
- Notation: PGD-t, bounded by  $\varepsilon$ , used the step-size of  $\alpha$



# **OUTER MINIMIZATION**

PGD (Projected Gradient Descent)

$$x^{t+1} = \Pi_{x+\mathcal{S}} \left( x^t + \alpha \operatorname{sgn}(\nabla_x L(\theta, x, y)) \right).$$

- Multi-step adversary; much stronger than FGSM attack
- Adversarial training
  - Make a model do correct prediction on adversarial examples
  - Training procedure
    - At each iteration of training
    - Craft PGD-t adversarial examples
    - Update the model towards making it correct on those adv examples



- Findings
  - (1, 3) PGD increases the loss values in a fairly consistent way
  - (2, 4) Models trained with PGD attacks are resilient to the same attacks





#### • Findings

University

- PGD increases the loss values in a fairly consistent way
- Models trained with PGD attacks are resilient to the same attacks
- Final loss of PGD attacks are concentrated (both for defended/undefended models)



- Why adversarial training (AT) works?
  - Capacity is crucial for the robustness: robust models need complex decision boundary
  - Capacity alone helps: high-capacity models show more robustness w/o AT





• ... Cont'd

Oregon State University

- Capacity is crucial for the robustness: robust models need complex decision boundary
- Capacity alone helps: high-capacity models show more robustness w/o AT
- AT with weak attacks (like FGSM) can't defeat a strong one like PGD
- (optional) Robustness may be at odds with accuracy



# SUMMARY

- Bottom-line
  - PGD is a strong attack we can use
  - Training a model with PGD can make it resilient to the first-order adversary
  - To achieve such robustness, we need sufficient model complexity



# CS 499/579: TRUSTWORTHY ML ADVERSARIAL ATTACKS: TRANSFERABILITY

Instructor: Sanghyun Hong

sanghyun.hong@oregonstate.edu





#### **ADVERSARIAL EXAMPLES ATTACKS**

- Test-time (evasion) attack
  - Given a test-time sample *x*
  - Craft an adversarial example  $x^*$  that fools the target neural network



#### **A**DVERSARIAL ATTACKS

• Example: An adversary wants to upload NSFW image to the cloud





## WHITE-BOX ADVERSARIAL ATTACKS

• Example: An adversary wants to upload NSFW image to the cloud



- FGSM, C&W, PGD, ...: the attacker has complete access to the target model



#### **BLACK-BOX ADVERSARIAL ATTACKS**

• Example: An adversary wants to upload NSFW image to the cloud





# (TRANSFER-BASED) BLACK-BOX ADVERSARIAL ATTACK

• Example: An adversary wants to upload NSFW image to the cloud



- Transfer-based attacks<sup>12</sup> : craft adv. examples on a transfer prior





# (OPTIMIZATION-BASED) BLACK-BOX ADVERSARIAL ATTACK

• Example: An adversary wants to upload NSFW image to the cloud



- Transfer-based attacks<sup>12</sup> : craft adv. examples on a transfer prior
- Optimization-based attacks<sup>3</sup> : craft them iteratively with query outputs and a transfer prior

Goodfellow et al., Explaining and Harnessing Adversarial Examples, ICLR 2015
 Madry et al., Towards Deep Learning Models Resistant to Adversarial Attacks, ICLR 2018
 Cheng et al., Improving Black-box Adversarial Attacks with a Transfer-based Prior, NeurIPS 2019



#### **TODAY WE TALK ABOUT TRANSFER-BASED ATTACKS**

DELVING INTO TRANSFERABLE ADVERSARIAL EXAMPLES AND BLACK-BOX ATTACKS, LIU ET AL., ICLR 2017

- Research questions
  - How well do adversarial examples transfer?
  - How practical are the transfer-based attacks?
  - What factors influence the transferability?
  - How can we reduce the transferability?



- Empirical evaluation
  - Train two models on a dataset
  - Craft adversarial examples on a model A (targeted and non-targeted)
  - Measure the success of these examples on the other model B
- Setup
  - Choose 100 images randomly from the ImageNet test-set
  - Use ResNet-50/-101/-152, GoogleNet, and VGG-16 models
  - Matching rate and distortion ( $l_2$ -distance)
- Adversarial attacks
  - Optimization-based approach (similar to C&W)
  - Fast Gradient-based approach (similar to PGD)

•	Results	from	non-targeted	attacks	(Top-5 acc.)
---	---------	------	--------------	---------	--------------

RMSD	ResNet-152	ResNet-101	ResNet-50	VGG-16	GoogLeNet
22.83	0%	13%	18%	19%	11%
23.81	19%	0%	21%	21%	12%
22.86	23%	20%	0%	21%	18%
22.51	22%	17%	17%	0%	5%
22.58	39%	38%	34%	19%	0%
	RMSD22.8323.8122.8622.5122.58	RMSDResNet-15222.830%23.8119%22.8623%22.5122%22.5839%	RMSDResNet-152ResNet-10122.830%13%23.8119%0%22.8623%20%22.5122%17%22.5839%38%	RMSDResNet-152ResNet-101ResNet-5022.830%13%18%23.8119%0%21%22.8623%20%0%22.5122%17%17%22.5839%38%34%	RMSDResNet-152ResNet-101ResNet-50VGG-1622.830%13%18%19%23.8119%0%21%21%22.8623%20%0%21%22.5122%17%17%0%22.5839%38%34%19%

Panel A: Optimization-based approach

RMSD	ResNet-152	ResNet-101	ResNet-50	VGG-16	GoogLeNet
23.45	4%	13%	13%	20%	12%
23.49	19%	4%	11%	23%	13%
23.49	25%	19%	5%	25%	14%
23.73	20%	16%	15%	1%	7%
23.45	25%	25%	17%	19%	1%
	RMSD23.4523.4923.4923.7323.45	RMSDResNet-15223.454%23.4919%23.4925%23.7320%23.4525%	RMSDResNet-152ResNet-10123.454%13%23.4919%4%23.4925%19%23.7320%16%23.4525%25%	RMSDResNet-152ResNet-101ResNet-5023.454%13%13%23.4919%4%11%23.4925%19%5%23.7320%16%15%23.4525%25%17%	RMSDResNet-152ResNet-101ResNet-50VGG-1623.454%13%13%20%23.4919%4%11%23%23.4925%19%5%25%23.7320%16%15%1%23.4525%25%17%19%

Panel B: Fast gradient approach



- More distortion leads to successful attacks?
  - Setup: VGG-16 to ResNet-152



**Oregon State** 

#### • Results from targeted attacks (Matching rate)

	RMSD	ResNet-152	ResNet-101	ResNet-50	VGG-16	GoogLeNet
ResNet-152	23.13	100%	2%	1%	1%	1%
ResNet-101	23.16	3%	100%	3%	2%	1%
ResNet-50	23.06	4%	2%	100%	1%	1%
VGG-16	23.59	2%	1%	2%	100%	1%
GoogLeNet	22.87	1%	1%	0%	1%	100%

- What if we use just random perturbations? Does *not* transfer



- Take aways
  - Non-targeted adversarial attacks transfer
  - Targeted adversarial attacks does not transfer well
  - Sub-research question: How we can make targeted attacks transferable?



#### **IMPROVING TRANSFERABILITY OF TARGETED ATTACKS**

#### • "Ensemble" (Used optimization-based attacks)

	RMSD	ResNet-152	ResNet-101	ResNet-50	VGG-16	GoogLeNet
-ResNet-152	30.68	38%	76%	70%	97%	76%
-ResNet-101	30.76	75%	43%	69%	98%	73%
-ResNet-50	30.26	84%	81%	46%	99%	77%
-VGG-16	31.13	74%	78%	68%	24%	63%
-GoogLeNet	29.70	90%	87%	83%	99%	11%

#### - What about non-targeted attacks?

	RMSD	ResNet-152	ResNet-101	ResNet-50	VGG-16	GoogLeNet
-ResNet-152	17.17	0%	0%	0%	0%	0%
-ResNet-101	17.25	0%	1%	0%	0%	0%
-ResNet-50	17.25	0%	0%	2%	0%	0%
-VGG-16	17.80	0%	0%	0%	6%	0%
-GoogLeNet	17.41	0%	0%	0%	0%	5%



#### **IMPROVING TRANSFERABILITY OF TARGETED ATTACKS**

- Why does ensemble work?
  - Hypothesis: it makes computed gradients are aligned to that of the target model
  - Evaluation approach
    - Compute the gradients of inputs from the models
    - Compute the cosine similarity between the gradients from two different models
  - Results

	ResNet-152	ResNet-101	ResNet-50	VGG-16	GoogLeNet
ResNet-152	1.00	_		—	_
ResNet-101	0.04	1.00	_	_	
ResNet-50	0.03	0.03	1.00	—	—
VGG-16	0.02	0.02	0.02	1.00	_
GoogLeNet	0.01	0.01	0.01	0.02	1.00



- Method
  - Craft adversarial examples on ImageNet models
  - Use them to fool the object recognition service in Clarifai.com (You can do as well)
- Setup
  - Choose 100 images randomly from the ImageNet test-set
  - Use models: ResNet-50/-101, GoogleNet and VGG-16
  - Matching rate
- Attacks
  - Optimization-based approach (similar to C&W)



- Transfer attack results
  - Non-targeted:
    - Most attacks transfer (= fooled Clarifai.com)
      - 57% AEs crafted on VGG-16 transfer
      - 76% AEs crafted on the ensemble transfer
  - Targeted:
    - Misclassification towards a target label
      - 2% AEs crafted on VGG-16 transfer
      - 18% AEs crafted on the ensemble transfer



#### • Transfer attack results

original image	true label	Clarifai.com results of original image	target label	targeted adversarial example	Clarifai.com results of targeted adversarial example
	viaduct	bridge, sight, arch, river, sky	window screen		window, wall, old, decoration, design
	hip, rose hip, rosehip	fruit, fall, food, little, wildlife	stupa, tope		Buddha, gold, temple, celebration, artistic



# **Thank You!**

Instructor: Sanghyun Hong

https://secure-ai.systems/courses/MLSec/F23



