

# Recent advances in adversarial machine learning: defense, transferable and camouflaged attacks

Xingjun Ma

School of Computing and Information Systems The University of Melbourne

**April 2020** 



#### Deep learning models are used everywhere



1



#### **Deep neural networks are vulnerable**



# Small perturbation can fool state-of-the-art ML models.

Szegedy et al. 2013, Goodfellow et al. 2014



# Security risks in medical diagnosis



Understanding Adversarial Attacks on Deep Learning Based Medical Image Analysis Systems **Ma** et al., Pattern Recognition, 2020.



#### Security threats to autonomous driving



#### Adversarial traffic signs all recognized as: 45km speed limit.



4

Evtimov et al. 2017





Original: What is the oncorhynchus also called? A: chum salmon Changed: What's the oncorhynchus also called? A: keta

Original: How long is the Rhine? A: 1,230 km Changed: How long is the Rhine??

A: more than 1,050,000



## Security risks in face or object recognition











#### How adversarial examples are crafted



8

Train a model:



Adversarial Attack:



### How adversarial examples are crafted

Model training:

$$\min_{\theta} \sum_{(x_i, y_i) \in D_{train}} L(f_{\theta}(x_i), y_i)$$

 $D_{train}$ : training data  $x_i$ : training sample  $y_i$ : class label L: loss function  $f_{\theta}$ : model

Adversarial attack:



• Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2014):

$$x' = x + \varepsilon \cdot \operatorname{sign} \nabla_x L(f_{\theta}(x), y)$$
 x': advs example



• Viewing DNN as a sequence of transformed spaces:



#### Non-linear explanation:

- Non-linear transformations leads to the existence of small "pockets" in the deep space:
- Regions of **low probability** (not naturally occurring).
- Densely scattered regions.
- Continuous regions.
- Close to normal data subspace.

Characterizing Adversarial Subspace Using Local Intrinsic Dimensionality. *Ma*, et al. ICLR 2018





- An illustrative example
  - $x \in [-1, 1), y \in [-1, 1), z \in [-1, 2)$
  - Binary classification
    - Class 1:  $z < x^2 + y^3$
    - Class 2:  $z \ge x^2 + y^3$
  - -x, y and z are increased by 0.01
    - $\rightarrow$  a total of 200 × 200 × 300
      - =  $1.2 \times 10^7$  points



- How many points are needed to reconstruct the decision boundary?
  - Training dataset: choose 80, 800, 8000, 80000 points randomly
  - Test dataset: choose 40, 400, 4000, 40000 points randomly
  - Boundary dataset (adversarial samples are likely to locate here):  $x^2 + y^3 - 0.1 < z < x^2 + y^3 + 0.1$



# Insufficient training data?

#### Test result

#### RBF SVMs

Size of the training dataset	Accuracy on its own test dataset	Accuracy on the test dataset with $4 \times 10^4$ points		Accuracy on the boundary dataset		
80	100	92.7		60.8		
800	99.0	97.4		74.9		
8000	99.5	99.6		94.1		
80000	99.9	99.9		99.9		98.9

Linear SVMs

Size of the training dataset	Accuracy on its own test dataset	Accuracy on the test dataset with $4 \times 10^4$ points	Accuracy on the boundary dataset
80	100	96.3	70 1
800	99.8	99.0	85.7
8000	99.9	99.8	97.3
80000	99.98	99.98	99.5

- 8000: 0.067% of 1.2 × 10<sup>7</sup>
- MNIST: 28 × 28 8-bit greyscale images, (2<sup>8</sup>)<sup>28×28</sup> ≈ 1.1 × 10<sup>1888</sup>
- $1.1 \times 10^{1888} \times 0.067\% \gg 6 \times 10^5$







• Viewing DNN as a stack of linear operations:

 $w^T x + b$ 

#### Linear explanation:

- Adversarial subspaces span a contiguous multidimensional space:
- Small changes at individual dimensions can sum up to significant change in final output:  $\sum_{i=0}^{n} x_i + \epsilon$ .
- Adversarial examples can always be found if  $\epsilon$  is large enough.



Goodfellow et al. 2014, 2016

-400

-20

-10

10

20

30



#### Training models on adversarial examples.





Adversarial training is a **min-max optimization** process:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \max_{\|x_i' - x_i\|_p \le \epsilon} L(f_{\theta}(x_i'), y_i)$$

attacking

*L*: loss,  $f_{\theta}$ : model,  $x_i$ : clean example,  $y_i$ : class,  $x'_i$ : adversarial example.

#### 1. Inner Maximization:

- This is to generate adversarial examples, by maximizing the loss *L*.
- It is a constrained optimization problem:  $||x_i' x_i||_p \leq \epsilon$ .

#### 2. Outer Minimization:

- A typical process to train a model, but on adversarial examples  $x'_i$  generated by the inner maximization.

On the Convergence and Robustness of Adversarial Training. Wang\*, Ma\*, et al., ICML 2019.

Mary et al. ICLR 2018.



#### Improving Adversarial Robustness Requires Revisiting Misclassified Examples

Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma and Quanquan Gu

ICLR 2020.



# Adversarial risk:

$$\mathcal{R}(h_{\boldsymbol{\theta}}) = \frac{1}{n} \sum_{i=1}^{n} \max_{\mathbf{x}'_{i} \in \mathcal{B}_{\epsilon}(\mathbf{x}_{i})} \mathbb{1}(h_{\boldsymbol{\theta}}(\mathbf{x}'_{i}) \neq y_{i}),$$

Revisited adversarial risk (correctly- vs mis-classified):

$$\min_{\boldsymbol{\theta}} \mathcal{R}_{\text{misc}}(h_{\boldsymbol{\theta}}) := \frac{1}{n} \Big( \sum_{i \in \mathcal{S}_{h_{\boldsymbol{\theta}}}^+} \mathcal{R}^+(h_{\boldsymbol{\theta}}, \mathbf{x}_i) + \sum_{i \in \mathcal{S}_{h_{\boldsymbol{\theta}}}^-} \mathcal{R}^-(h_{\boldsymbol{\theta}}, \mathbf{x}_i) \Big)$$
$$= \frac{1}{n} \sum_{i=1}^n \Big\{ \mathbb{1}(h_{\boldsymbol{\theta}}(\hat{\mathbf{x}}_i') \neq y_i) + \mathbb{1}(h_{\boldsymbol{\theta}}(\mathbf{x}_i) \neq h_{\boldsymbol{\theta}}(\hat{\mathbf{x}}_i')) \cdot \mathbb{1}(h_{\boldsymbol{\theta}}(\mathbf{x}_i) \neq y_i) \Big\}$$



• Surrogate loss functions (existing methods and MART)

Defense Method	Loss Function
Standard	$ ext{CE}(\mathbf{p}(\hat{\mathbf{x}}',oldsymbol{ heta}),y)$
ALP	$ ext{CE}(\mathbf{p}(\hat{\mathbf{x}}',oldsymbol{ heta}),y) + \lambda \cdot \ \mathbf{p}(\hat{\mathbf{x}}',oldsymbol{ heta}) - \mathbf{p}(\mathbf{x},oldsymbol{ heta})\ _2^2$
CLP	$ ext{CE}(\mathbf{p}(\mathbf{x},oldsymbol{ heta}),y) + \lambda \cdot \ \mathbf{p}(\hat{\mathbf{x}}',oldsymbol{ heta}) - \mathbf{p}(\mathbf{x},oldsymbol{ heta})\ _2^2$
TRADES	$ ext{CE}( extbf{p}( extbf{x},oldsymbol{ heta}),y) + \lambda \cdot  ext{KL}ig( extbf{p}( extbf{x},oldsymbol{ heta})     extbf{p}(\hat{ extbf{x}}',oldsymbol{ heta})ig)$
MMA	$\operatorname{CE}(\mathbf{p}(\hat{\mathbf{x}}',\boldsymbol{\theta}),y) \cdot \mathbb{1}(h_{\boldsymbol{\theta}}(\mathbf{x})=y) + \operatorname{CE}(\mathbf{p}(\mathbf{x},\boldsymbol{\theta}),y) \cdot \mathbb{1}(h_{\boldsymbol{\theta}}(\mathbf{x})\neq y)$
MART	$\text{BCE}(\mathbf{p}(\hat{\mathbf{x}}', \boldsymbol{\theta}), y) + \lambda \cdot \text{KL}(\mathbf{p}(\mathbf{x}, \boldsymbol{\theta})    \mathbf{p}(\hat{\mathbf{x}}', \boldsymbol{\theta})) \cdot (1 - \mathbf{p}_y(\mathbf{x}, \boldsymbol{\theta}))$

• Semi-supervised extension of MART:

$$\mathcal{L}_{\text{semi}}^{\text{MART}}(\boldsymbol{\theta}) = \sum_{i \in \mathcal{S}_{\text{sup}}} \ell_{\text{sup}}^{\text{MART}}(\mathbf{x}_i, y_i; \boldsymbol{\theta}) + \gamma \cdot \sum_{i \in \mathcal{S}_{\text{unsup}}} \ell_{\text{unsup}}^{\text{MART}}(\mathbf{x}_i, y_i; \boldsymbol{\theta})$$



• White-box robustness: ResNet-18, CIFAR-10,  $\epsilon = 8/255$ 

	MNIST				CIFAR-10			
Defense	Natural	FGSM	$PGD^{20}$	$CW_{\infty}$	Natural	FGSM	$PGD^{20}$	$CW_{\infty}$
Standard	99.11	97.17	94.62	94.25	84.44	61.89	47.55	45.98
MMA	98.92	97.25	95.25	94.77	84.76	62.08	48.33	45.77
Dynamic	98.96	97.34	95.27	94.85	83.33	62.47	49.40	46.94
TRADES	<b>99.25</b>	96.67	94.58	94.03	82.90	62.82	50.25	48.29
MART	98.74	<b>97.8</b> 7	96.48	96.10	83.07	65.65	55.57	54.87

• White-box robustness: WideResNet-34-10, CIFAR-10,  $\epsilon = 8/255$ 

		FGSM		$PGD^{20}$		$PGD^{100}$		$\mathrm{CW}_\infty$	
Defense	Natural	Best	Last	Best	Last	Best	Last	Best	Last
Standard	87.30	56.10	56.10	52.68	49.31	51.55	49.03	50.73	48.47
Dynamic	84.51	63.53	63.53	55.03	51.70	54.12	50.07	51.34	49.27
TRADES	84.22	64.70	64.70	56.40	53.16	55.68	51.27	51.98	51.12
MART	84.17	67.51	67.51	58.56	57.39	57.88	55.04	54.58	54.53



- White-box robustness: unlabled data, CIFAR-10,  $\epsilon = 8/255$ •
- a) WideResNet-34-8 with 100K unlabeled data b) WideResNet-28-10 with 500K unlabeled data

 $PGD^{20}$ 

62.76

63.10

65.04

Defense	Natural	$PGD^{20}$	 Defense	Natural
UAT++	86.04	59.41	 UAT++	86.21
RST	88.24	59.60	RST	89.70
MART	86.68	61.88	MART	86.30



#### Skip Connections Matter: on the Transferability of Adversarial Examples Generated with ResNets

Dongxian Wu, Yisen Wang, Shu-Tao Xia, James Bailey and Xingjun Ma. ICLR 2020.



• Gradient backpropagation with skip connections



Skip the gradients incrases transferability!



• New attack method: skip gradient method (SGM)

$$oldsymbol{x}_{adv}^{t+1} = \Pi_{\epsilon} \Big( oldsymbol{x}_{adv}^t + lpha \cdot ext{sign} ig( rac{\partial \ell}{\partial oldsymbol{z}_L} \prod_{i=0}^{L-1} [\gamma] rac{\partial f_{i+1}}{\partial oldsymbol{z}_i} + 1) rac{\partial oldsymbol{z}_0}{\partial oldsymbol{x}} ig) \Big)$$

Breaking down a network f according to its L residual blocks.

-----

	RN18	<b>RN34</b>	<b>RN50</b>	<b>RN101</b>	RN152	DN121	DN169	DN201
PGD	23.23±0.69	$24.38{\pm}0.41$	$22.80{\pm}0.55$	$22.98{\pm}0.83$	$26.56{\pm}0.75$	30.71±0.60	$30.90{\pm}0.31$	36.01±0.59
SGM	28.92±0.45	43.43±0.32	36.71±0.55	38.38±0.53	$\textbf{44.84}{\pm 0.14}$	57.38±0.14	60.45±0.42	65.48±0.23

ImageNet, target: Inception V3,  $\epsilon = 16/255$ 



Combined with existing methods: the success rates (%) of attacks crafted on source model DN201 against 7 target models.

Attack \Target	VGG19	RN152	DN201	SE154	IncV3	IncV4	IncRes
MI	75.09	76.39	<b>99.84</b>	64.38	59.62	54.85	50.05
MI+SGM	+12.01	+13.24	99.52	+17.16	+21.88	+15.57	+18.35
DI	78.11	78.18	99.81	61.75	60.04	56.15	49.00
DI+SGM	+12.28	+13.76	99.52	+20.92	+17.66	+15.78	+20.20
MI+DI	87.16	87.28	99.76	79.80	76.68	75.20	71.05
MI+DI+SGM	93.00	93.92	99.42	89.86	85.72	81.23	80.50



# Adversarial Camouflage: Hiding Adversarial Examples with Natural Styles Ranjie Duan, **Xingjun Ma**, Yisen Wang, James Bailey, Kai Qin, Yun Yang *CVPR 2020.*



#### Adversarial camouflage



Camouflage adversarial examples with customized styles.



# Adversarial camouflage



## Making large perturbations look natural: Adversarial attack + style transfer





(a)  $RP_2$  (b) AdvCam (c) AdvPatch (d) AdvCam

## A visually comparison to existing attacks



#### Adversarial camouflage



Revolver --> Toilet tissue

Minivan --> Traffic light

Scabbard --> Purse

Attacking the background is what makes the attack stealthy and ubiquitous.

Examples of camouflaged digital attacks



### Adversarial camouflage



Traffic sign -> Barbershop



Tree -> Street sign

Examples of camouflaged physical-world attacks



# Using adversarial camouflage to protect privacy



Here is an adversarial pikachu

This is a dog to Google Image Search.

pikachu.jpg × bull terrier	۵ 🎙 ۹		sign i
Q All Images 🤣 Shopping : More	Settings Tools		SafeSear
About 2 results (0.30 seconds)   Image size:   548 × 548   No other sizes of this image found.   Possible related search: bull terrier		Bull < Terrier Dog breed	
www.akc.org > Dog Breeds <b>Bull Terrier Dog Breed Information - American Kenne</b> Feb 12, 2020 - Right breed for you? <b>Bull Terrier</b> information including person grooming, pictures, videos, and the AKC breed standard.	I Club onality, history,	The Bull Terrier is a b terrier family. There is version of this breed known as the Miniatu Wikipedia	reed of dog in the also a miniature which is officially re Bull Terrier.



# Thank you!



