对抗机器学习: 攻击、防御与模型鲁棒性评估



京东智联云机器学习部负责人 京东 (上海) 人工智能研究院负责人



ImageNet Challenge





Accuracy is NOT the Sole Metric

Suppose you have a well-trained ImageNet classifier that achieves >97% accuracy.
 Should you always trust its prediction?

Accuracy is NOT the Sole Metric

Suppose you have a well-trained ImageNet classifier that achieves >97% accuracy.
 Should you always trust its prediction?









Any differences?



Accuracy is NOT the Sole Metric

Suppose you have a well-trained ImageNet classifier that achieves >97% accuracy.
 Should you always trust its prediction?



Adversarial Examples

A small perturbation could dramatically change the prediction





Trade-off between Accuracy & Robustness

Solely pursuing for high-accuracy rate may be misleading



Su et al., Is Robustness the Cost of Accuracy? A Comprehensive Study on the Robustness of 18 Deep Image Classification Models, ECCV 2018



Dilemma in Security-critical Tasks

Security-critical tasks require both accuracy and robustness



Autonomous Driving



Healthcare



Finance



Surveillance



Military



Law



Adversarial Machine Learning: a Fast Growing Area



https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html



Key Problems in Adversarial Machine Learning



Attack Settings: White and Black-Box Attacks



White-box Attack by Optimization-based Approaches



N. Carlini, D. Wagner. Towards Evaluating the Robustness of Neural Networks. IEEE Symposium on Security and Privacy, 2017 Chen et al. EAD: Elastic-Net Attacks to Deep Neural Networks via Adversarial Examples, AAAI 2018

Black-box Attacks with Soft Labels

■ The adversary has no access to the structure and parameters of deep neural networks

He can only query the model and get the probability outputs



Black-box Attacks by Zeroth Order Optimization (ZOO)

$$\mathbf{x}^* = \arg\min_{\mathbf{x}} \operatorname{Dist}(\mathbf{x}, \mathbf{x}_0) + c \cdot \mathcal{L}(\mathbf{x})$$

Cannot compute the gradient in the black-box setting

Chen et al., ZOO: Zeroth Order Optimization Based Black-box Attacks to Deep Neural Networks without Training Substitute Models. AISec@CCS 2017



Black-box Attacks by Zeroth Order Optimization (ZOO)

$$\mathbf{x}^* = \arg\min_{\mathbf{x}} \operatorname{Dist}(\mathbf{x}, \mathbf{x}_0) + c \cdot \mathcal{L}(\mathbf{x})$$

- Cannot compute the gradient in the black-box setting
- Symmetric difference quotient to estimate the gradient:

$$\hat{g}_i pprox rac{\partial \mathcal{L}(oldsymbol{x})}{\partial x_i} pprox rac{\mathcal{L}(oldsymbol{x} + \epsilon oldsymbol{e}_i) - \mathcal{L}(oldsymbol{x} - \epsilon oldsymbol{e}_i)}{2\epsilon}$$

- However, need O(d) queries to estimate a gradient
 - ImageNet: d = 299*299*3 > 268K
 - 100 iterations => 26.8 million queries

Chen et al., ZOO: Zeroth Order Optimization Based Black-box Attacks to Deep Neural Networks without Training Substitute Models. AlSec@CCS 2017



 \hat{g}_1

 $\mathbf{x} \leftarrow \mathbf{x} - \eta$

Black-box Attacks by Zeroth Order Optimization (ZOO)

■ Reduce number of queries:

- Stochastic Coordinate descent: update a small set of coordinates at each time
- Greedy approach: select important coordinates first
- Attack-space dimension reduction + hierarchical attack

■ The first black-box attack algorithm that achieves close to 100% attack success rate.

		MNIS	5T			
	Success Rate	Rate Avg. L ₂ Avg. Time (per atta				
White-box	100 %	2.00661	0.53 min			
Black-box (Substitute Model)	(Substitute Model) 26.74 % 5.272 0.80 min (+ 6.16					
Proposed black-box (ZOO)	98.9 %	1.987068	1.62 min			
		CIFAR	210			
	Success Rate	ess Rate Avg. L_2 Avg. Time (per attack)				
White-box	100 %	0.37974	0.16 min			
Black-box (Substitute Model)	5.3 %	5.7439	0.49 min (+ 7.81 min)			
Proposed Black-box (ZOO)	96.8 %	0.39879	3.95 min			

Chen et al., ZOO: Zeroth Order Optimization Based Black-box Attacks to Deep Neural Networks without Training Substitute Models. Alsec CL

Further Improves Query Efficiency

AutoZOOM:

1. Dispense with coordinate-wise estimation and instead propose a scaled random full gradient estimator.

2. An autoencoder trained offline with unlabeled data or a bilinear resizing operation for attack acceleration.

Reduced at least 93.2% query count 99.4% for ImageNet





Input-free Attack:

 Start with a gray color image;
 Shrink the dimension, then perturb a small region and tile it to cover the input image.

With only 1.7K queries on average, can perturb a gray image to any target class of ImageNet with a 100% success rate on InceptionV3.

Tu et al., AutoZOOM: Autoencoder-Based Zeroth Order Optimization Method for Attacking Black-Box Neural Networks. AAAI 2019 Du et al., Towards Query Efficient Black-box Attacks: An Input-free Perspective. ASec@CCS 2018

Further Improves Query Efficiency

Using an **active learning strategy** to significantly reduce the number of queries

Algorithm 2 Substitute DNN training with active learning INPUT: target oracle \tilde{O} , a maximum number ρ_{max} of training epochs, and an initial training set S_0 .

OUTPUT: a trained substitute model F.

1: Define architecture F;

2: for
$$\rho = 0; \rho < \rho_{max}; \rho + + do$$

3: **if**
$$\rho = 0$$
 then

4:
$$D \leftarrow \{(\mathbf{x}, \tilde{O}(\mathbf{x})) | \mathbf{x} \in S_{\rho}\}$$

5: **else**

6:
$$D_{add} \leftarrow \{(\mathbf{x}, \tilde{O}(\mathbf{x})) | \mathbf{x} \in S_{add}\};$$

7:
$$D \leftarrow [D, D_{add}]$$

- B: end if
- 9: train F with D;
- 10: craft S_{add} ;
- 11: Use Active Learning strategy to generate a new S_{add} ;

12: $S_{\rho+1} \leftarrow S_{\rho} \cup S_{add};$

13: **end for**

Reduce more than 90% of queries Obtain an accurate substitute model 85% similar with the target oracle.

Spanning Attack

Constrain a subspace spanned by an auxiliary small unlabeled dataset

$$\boldsymbol{a} = \boldsymbol{b}^{\mathsf{T}} \boldsymbol{S} = (\boldsymbol{b}^{\mathsf{T}} \boldsymbol{U}_N \boldsymbol{\Sigma}_N) \boldsymbol{V}_N^{\mathsf{T}}.$$
$$\boldsymbol{\Sigma}_N^{-1} \boldsymbol{U}_N^{\mathsf{T}} \boldsymbol{S} = \boldsymbol{V}_N^{\mathsf{T}}.$$

$$\boldsymbol{a} = \boldsymbol{b}^{\mathsf{T}} \boldsymbol{V}_N^{\mathsf{T}} = (\boldsymbol{b}^{\mathsf{T}} \boldsymbol{\Sigma}_N^{-1} \boldsymbol{U}_N^{\mathsf{T}}) \boldsymbol{S}$$



The reinforced attack typically requires less than 50% queries while improves success rates in the meantime.

Propose a novel **Frank-Wolfe based projection-free attack framework** for both white-box and black-box settings

$$\min_{\mathbf{x}\in\mathcal{X}} f(\mathbf{x}_t) + \langle \mathbf{x} - \mathbf{x}_t, \nabla f(\mathbf{x}_t) \rangle$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma_t \epsilon \cdot \operatorname{sign}(\mathbf{m}_t) - \gamma_t (\mathbf{x}_t - \mathbf{x}_{ori})$$

```
Algorithm 1 Frank-Wolfe White-box Attack Algorithm1: input: number of iterations T, step sizes \{\gamma_t\};2: \mathbf{x}_0 = \mathbf{x}_{ori}, \mathbf{m}_{-1} = \nabla f(\mathbf{x}_0)3: for t = 0, \dots, T - 1 do4: \mathbf{m}_t = \beta \cdot \mathbf{m}_{t-1} + (1 - \beta) \cdot \nabla f(\mathbf{x}_t)5: \mathbf{v}_t = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, \mathbf{m}_t \rangle // \text{LMO}6: \mathbf{d}_t = \mathbf{v}_t - \mathbf{x}_t7: \mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{d}_t8: end for9: output: \mathbf{x}_T
```

The proposed attack algorithms with momentum mechanism enjoy an $O(1/\sqrt{T})$ convergence rate in the nonconvex setting.

Li, et al., Query-Efficient Black-Box Attack by Active Learning, ICDM 2018 Wang et al., Spanning Attack: Reinforce Black-box Attacks with Unlabeled Data, Machine Learning 2020 Chen et al., A Frank-Wolfe Framework for Efficient and Effective Adversarial Attacks, AAAI 2020



Black-box Attacks with Hard Labels

- The adversary has no access to the structure and parameters of deep neural networks
- He can only query the model and get the hard-label multi-class output



Optimization-based Hard-label Black-box Attack

Reformulate the attack optimization problem

 $\theta^* = \arg\min_{\boldsymbol{\theta}} g(\boldsymbol{\theta})$ Untargeted attack: $g(\boldsymbol{\theta}) = \arg\min_{\lambda>0} \left(f(\mathbf{x}_0 + \lambda \frac{\boldsymbol{\theta}}{\|\boldsymbol{\theta}\|}) \neq y_0 \right)$ Targeted attack: $g(\boldsymbol{\theta}) = \arg\min_{\lambda>0} \left(f(\mathbf{x}_0 + \lambda \frac{\boldsymbol{\theta}}{\|\boldsymbol{\theta}\|}) = t \right)$ g(

- Cannot compute the gradient of g
- However, can compute the function value of g via querying
- Binary search + fine-grained search



θ : the direction of adversarial example

Cheng et al., Query-Efficient Hard-label Black-box Attack: An Optimization-based Approach. ICLR 2019



Optimization-based Hard-label Black-box Attack

Algorithm for computing $g(\theta)$

Algorithm 1 Compute $g(\theta)$ locally

1: **Input:** Hard-label model f, original image x_0 , query direction θ , previous value v, increase/decrease ratio $\alpha = 0.01$, stopping tolerance ϵ (maximum tolerance of computed error) 2: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} / \|\boldsymbol{\theta}\|$ 3: if $f(\boldsymbol{x}_0 + v\boldsymbol{\theta}) = y_0$ then $v_{left} \leftarrow v, v_{right} \leftarrow (1 + \alpha)v$ 4: while $f(\boldsymbol{x}_0 + v_{right}\boldsymbol{\theta}) = y_0 \, \mathbf{do}$ 5: $v_{right} \leftarrow (1+\alpha)v_{right}$ 6: 7: else $v_{right} \leftarrow v, v_{left} \leftarrow (1 - \alpha)v$ 8: while $f(\boldsymbol{x}_0 + v_{left}\boldsymbol{\theta}) \neq y_0$ do 9: $v_{left} \leftarrow (1 - \alpha) v_{left}$ 10: 11: ## Binary Search within $[v_{left}, v_{riaht}]$ 12: while $v_{right} - v_{left} > \epsilon \, \mathbf{do}$ $v_{mid} \leftarrow (v_{right} + v_{left})/2$ 13: if $f(\boldsymbol{x}_0 + v_{mid}\boldsymbol{\theta}) = y_0$ then 14: 15: $v_{left} \leftarrow v_{mid}$ 16: else 17: $v_{right} \leftarrow v_{mid}$

18: return v_{right}

Zeroth-order optimization for minimizing $g(\theta)$

Algorithm 2 RGF for hard-label black-box attack
1: Input: Hard-label model f , original image x_0 , initial θ_0 .
2: for $t = 0, 1, 2, \dots, T$ do
3: Randomly choose u_t from a zero-mean Gaussian distribution
4: Evaluate $g(\boldsymbol{\theta}_t)$ and $g(\boldsymbol{\theta}_t + \beta \boldsymbol{u})$ using Algorithm 1
5: Compute $\hat{\boldsymbol{g}} = \frac{g(\boldsymbol{\theta}_t + \beta \boldsymbol{u}) - g(\boldsymbol{\theta}_t)}{\beta} \cdot \boldsymbol{u}$
6: Update $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \hat{\boldsymbol{g}}$
7: return $oldsymbol{x}_0 + g(oldsymbol{ heta}_T)oldsymbol{ heta}_T$

More than 4 times faster than Decision-attack

	Mì	NIST	CIFAR10				
	Avg L_2	# queries	queriesAvg L_2 # querie $30,103$ 0.2850 $55,552$ $58,508$ 0.2213 $140,572$ $92,018$ 0.2122 $316,792$				
	2.3158	30,103	0.2850	55,552			
Decision-attack (black-box)	2.0052	58,508	0.2213	140,572			
	1.8668	192,018	0.2122	316,791			
	1.8522	46,248	0.2758	61,869			
Opt-attack (black-box)	1.7744	57,741	0.2369	141,437			
	1.7114	73,293	0.2300	186,753			
C&W (white-box)	1.4178	-	0.1901	-			

Cheng etal., Query-Efficient Hard-label Black-box Attack: An Optimization-based Approach. ICLR 2019

Adversarial Examples in Image Captioning



Chen et al., Attacking Visual Language Grounding with Adversarial Examples: A Case Study on Neural Image Captioning. ACL 2018

3. A large brown teddy bear laying on top of a bed.

Adversarial Examples in Sequence-to-Sequence Models



DATASET	SUCCESS RATE	BLEU	# CHANGED
GIGAWORD	86.0%	0.828	2.17
DUC2003	85.2%	0.774	2.90
DUC2004	84.2%	0.816	2.50

Attacking Text Summarization

Method	SUCCESS%	BLEU	# CHANGED
NON-OVERLAP	89.4%	0.349	3.5
1-keyword	100.0%	0.705	1.8
2-keyword	91.0 %	0.303	4.0
3-keyword	69.6%	0.205	5.3

SOURCE INPUT SEQ	UNDER NATO THREAT TO END HIS PUNISHING OFFENSIVE AGAINST ETHNIC ALBANIAN SEP- ARATISTS IN KOSOVO , PRESIDENT SLOBODAN MILOSEVIC OF YUGOSLAVIA HAS ORDERED
	MOST UNITS OF HIS ARMY BACK TO THEIR BARRACKS AND MAY WELL AVOID AN ATTACK BY
	THE ALLIANCE, MILITARY OBSERVERS AND DIPLOMATS SAY
ADV INPUT SEQ	UNDER NATO THREAT TO END HIS PUNISHING OFFENSIVE AGAINST ETHNIC ALBANIAN SEPA-
	RATISTS IN KOSOVO , PRESIDENT SLOBODAN MILOSEVIC OF YUGOSLAVIA HAS jean-sebastien
	MOST UNITS OF HIS ARMY BACK TO THEIR BARRACKS AND MAY WELL AVOID AN ATTACK BY
	THE ALLIANCE, MILITARY OBSERVERS AND DIPLOMATS SAY.
SOURCE OUTPUT SEQ	MILOSEVIC ORDERS ARMY BACK TO BARRACKS
ADV OUTPUT SEQ	nato may not attack kosovo

Cheng et al., Seq2Sick: Evaluating the Robustness of Sequence-to-Sequence Models with Adversarial Examples. AAAI 2020



Attacking Machine Translation

Adversarial Examples in Visual Question Answering



Fukui et al., Multimodal Compact Bilinear Pooling for Visual Question Answering and Visual Grounding. EMNLP 2016 Zeng et al., Adversarial Attacks Beyond the Image Space. CVPR 2019



Adversarial Examples in Reinforcement Learning





Adversarial Examples in Speech Recognition



Carlini et al., Audio Adversarial Examples: Targeted Attacks on Speech-to-Text. IEEE Symposium on Security and Privacy Workshops 2018 Yakura & Sakuma, Robust Audio Adversarial Example for a Physical Attack. IJCAI 2019 Taori et al., Targeted adversarial examples for black box audio systems. IEEE Security and Privacy Workshops 2019



Adversarial Examples in Shallow Learning







Figure 3: Model outputs for individual adversarial examples.

7 8

Research

Vidnerová & Nerud., Vulnerability of Machine Learning Models to Adversarial Examples. ITAT, 2016 Biggio et al., Security evaluation of support vector machines in adversarial environments. Support Vector Machines Applications 20 27

Adversarial Examples in Physical World



Brown et al., Adversarial Patch. 2017

Athalye et al., Synthesizing Robust Adversarial Examples. ICML 2018

Eykholt et al., Physical Adversarial Examples for Object Detectors WOOT 2018

Thys, Ranst, Goedemé., Fooling Automated Surveillance Cameras: Adversarial Patches to Attack Person Detection. 309 (1

Xũ et al., Evading Real-Time Person Detectors by Adversarial T-shirt. ECV Regence

Sharif, et al., A General Framework for Adversarial Examples with Objectives. ACM Transactions on Privacy and Security 2019

Why Do Adversarial Examples Matter?

Whenever there is an AI model, there is (almost) a way to generate adversarial examples



Key Problems in Adversarial Machine Learning



(Incomplete) List of Defense Strategies

Secondary Classification Defensive Distillation Normalization **Gradient Regularization Feature Squeezing Adversarial Training Activation Pruning Distributional Detection Model Compression PCA** Detection **Gradient Shattering Stochastic Gradients**



Adversarial Training



- Add adversarial examples to the training set, with their correct labels
- Robust optimization problem:

$$\theta^* = \arg\min_{\theta} \mathbb{E}_{x \in \mathcal{X}} \left[\max_{\delta \in [-\epsilon, \epsilon]^N} \ell(x + \delta; F_{\theta}) \right]$$

- (One of the) strongest defense mechanism so far
- Not scalable enough/ vulnerable to blind-spot attack/ high sample complexity . . .

Madry, et al., Towards Deep Learning Models Resistant to Adversarial Attacks. ICLR 2018



Further Improves Adversarial Training

Adversarial Distributional Training for Robust Deep Learning

Zhijie Deng^{*1} Yinpeng Dong^{*1} Tianyu Pang¹ Hang Su¹ Jun Zhu¹

Boosting Adversarial Training with Hypersphere Embedding

Tianyu Pang^{*}, Xiao Yang^{*}, Yinpeng Dong, Kun Xu, Hang Su, Jun Zhu Department of Computer Science and Technology Tsinghua University, Beijing, China {pty17, yangxiao19, dyp17}@mails.tsinghua.edu.cn kunxu.thu@gmail.com, {suhangss, dcszj}@mail.tsinghua.edu.cn

Convergence of Adversarial Training in Overparametrized Neural Networks

Ruiqi Gao^{1,*} Tianle Cai^{1,*} Haochuan Li² Liwei Wang³ Cho-Jui Hsieh⁴ Jason D. Lee⁵ ¹School of Mathematical Sciences, Peking University ²Department of EECS, Massachusetts Institute of Technology ³Key Laboratory of Machine Perception, MOE, School of EECS, Peking University ⁴Department of Computer Science, University of California, Los Angeles ⁵Department of Electrical Engineering, Princeton University **Theoretically Principled Trade-off between Robustness and Accuracy**

Hongyang Zhang¹² Yaodong Yu³ Jiantao Jiao⁴ Eric P. Xing¹⁵ Laurent El Ghaoui⁴ Michael I. Jordan⁴

Unlabeled Data Improves Adversarial Robustness

Yair Carmon* Stanford University yairc@stanford.edu Aditi Raghunathan* Stanford University aditir@stanford.edu

Ludwig Schmidt UC Berkeley ludwig@berkeley.edu

Percy Liang Stanford University pliang@cs.stanford.edu John C. Duchi Stanford University jduchi@stanford.edu

You Only Propagate Once: Accelerating Adversarial Training via Maximal Principle

Dinghuai Zhang*, Tianyuan Zhang* Peking University {zhangdinghuai, 1600012888}@pku.edu.cn Yiping Lu* Stanford University yplu@stanford.edu

Zhanxing Zhu[†]

School of Mathematical Sciences, Peking University Center for Data Science, Peking University Beijing Institute of Big Data Research zhanxing.zhu@pku.edu.cn

Bin Dong[†]

Dynamic Adversarial Training

First-Order Stationary Condition for constrained optimization (FOSC)

 $c(\mathbf{x}^k) = \max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x} - \mathbf{x}^k, \nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}^k) \rangle$

- The lower the score, the better convergence quality of the adversarial example \mathbf{x}^k
- Inner maximization: early (later) stages prefer low (high) convergence quality AEs

Algorithm 1 Dynamic Adversarial Training **Input:** Network h_{θ} , training data S, initial model parameters θ^0 , step size η_t , mini-batch \mathcal{B} , maximum FOSC value c_{max} , training epochs T, FOSC control epoch T', PGD step K, PGD step size α , maximum perturbation ϵ . for t = 0 to T - 1 do $c_t = \max(c_{\max} - t \cdot c_{\max}/T', 0)$ for each batch $\mathbf{x}_{\mathcal{B}}^{0}$ do $V = \mathbb{1}_{\mathcal{B}}$ # control vector of all elements is 1 while $\sum V > 0$ & k < K do $\mathbf{x}_{\mathcal{B}}^{k+\overline{1}} = \mathbf{x}_{\mathcal{B}}^{k} + V \cdot \alpha \cdot \operatorname{sign}(\nabla_{\mathbf{x}} \ell(h_{\theta}(\mathbf{x}_{\mathcal{B}}^{k}), y))$ $\mathbf{x}_{\boldsymbol{\beta}}^{k} = clip(\mathbf{x}_{\boldsymbol{\beta}}^{k}, \mathbf{x}_{\boldsymbol{\beta}}^{0} - \epsilon, \mathbf{x}_{\boldsymbol{\beta}}^{0} + \epsilon)$ $V = \mathbb{1}_{\mathcal{B}}(c(\mathbf{x}_{1\dots\mathcal{B}}^{k}) < c_{t})$ # The element of V becomes 0 at which FOSC is smaller than c_t end while $\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta_t \mathbf{g}(\boldsymbol{\theta}^t) \quad \# \mathbf{g}(\boldsymbol{\theta}^t) : stochastic gradient$ end for end for

Theorem 1. Suppose Assumptions 1, 2 and 3 hold. Let $\Delta = L_S(\theta^0) - \min_{\theta} L_S(\theta)$. If the step size of the outer minimization is set to $\eta_t = \eta = \min(1/L, \sqrt{\Delta/L\sigma^2 T})$. Then the output of Algorithm 1 satisfies

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\|\nabla L_S(\boldsymbol{\theta}^t)\|_2^2 \right] \le 4\sigma \sqrt{\frac{L\Delta}{T}} + \frac{5L_{\theta x}^2 \delta}{\mu}$$

where $L = (L_{\theta x} L_{x\theta} / \mu + L_{\theta \theta}).$

		M	NIST		CIFAR-10					
Defense	FGSM	PGD-10	PGD-20	$C\&W_{\infty}$	FGSM	PGD-10	PGD-20	$C\&W_{\infty}$		
Standard	96.12	95.73	95.73	97.20	65.65	65.80	65.60	66.12		
Curriculum	96.59	95.87	96.09	97.52	71.25	71.44	71.13	71.94		
Dynamic	97.60	97.01	96.97	98.36	71.95	72.15	72.02	72.85		

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

Adversarial Training with Misclassified Examples

Answer the following question: are the adversarial examples generated from misclassified and correctly classified examples, equally important for adversarial robustness?



Treat misclassified examples separately

Table 1: Loss function comparison with existing work. The adversarial example $\hat{\mathbf{x}}'$ is generated by (10) for all defense methods except TRADES and MMA. The adversarial example in TRADES is generated by maximizing its regularization term (KL-divergence), and the adversarial example in MMA is generated by solving (10) with different perturbation limit (*i.e.*, ϵ).

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

Algorithm 1 Misclassification Aware adveRsarial Training (MART) 1: Input: Training data $\{\mathbf{x}_i, y_i\}_{i=1,...,n}$, outer iteration number T_O , inner iteration number T_I , maximum perturbation ϵ , step size for inner optimization η_I , step size for outer optimization η_O 2: Initialization: Standard random initialization of h_{θ} 3: for $t = 1, ..., T_O$ do Uniformly sample a minibatch of training data $B^{(t)}$ for $\mathbf{x}_i \in B^{(t)}$ do 5: $\mathbf{x}'_i = \mathbf{x}_i + \epsilon \cdot \xi$, with $\xi \sim \mathcal{U}(-1, 1)$ # \mathcal{U} is a uniform distribution 6: for $s = 1, \ldots, T_I$ do 7: $\mathbf{x}'_i \leftarrow \Pi_{\mathcal{B}_{\epsilon}(\mathbf{x}_i)} \big(\mathbf{x}'_i + \eta_I \cdot \operatorname{sign}(\nabla_{\mathbf{x}'_i} \operatorname{CE}(\mathbf{p}(\mathbf{x}'_i, \boldsymbol{\theta}), y_i)) \big)$ $\# \Pi(\cdot)$ is the projection operator 8: end for 9: 10: $\hat{\mathbf{x}}'_i \leftarrow \mathbf{x}'_i$ 11: end for $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta_O \sum_{\mathbf{x}_i \in B^{(t)}} \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{x}_i, y_i, \hat{\mathbf{x}}'_i; \boldsymbol{\theta})$ 12: 13: end for 14: **Output:** Robust classifier h_{θ}

			•			· ·				
		MN	IST		CIFAR-10					
Defense	FGSM	PGD^{10}	PGD^{20}	CW_{∞}	FGSM	PGD^{10}	PGD^{20}	CW_{∞}		
Standard	96.12	95.73	95.47	96.34	79.98	80.27	80.01	80.85		
MMA	96.11	95.94	95.81	96.87	80.28	80.52	80.48	81.32		
Dynamic	97.60	96.25	95.82	97.03	81.37	81.71	81.38	82.05		
TRADES	97.49	96.03	95.73	97.20	81.52	81.73	81.53	82.11		
MART	97.77	96.96	96.97	98.36	82.75	82.93	82.70	82.95		

Table 3: Black-box robustness (accuracy (%) on black-box test attacks) on MNIST and CIFAR-10.

Wang et al., Improving Adversarial Robustness Requires Revisiting Misclassified Examples, ICLR 202

Adversarial Robustness Leaderboard



https://robustbench.github.io/

Key Problems in Adversarial Machine Learning

h









Robustness Estimation: CLEVER Score

CLEVER: Cross Lipschitz Extreme Value for nEtwork Robustness

First attack-independent robustness metric that can be applied to any neural network classifier

Theorem 3.1 (Formal guarantee on lower bound β_L). Let $x_0 \in \mathbb{R}^d$ and $f : \mathbb{R}^d \to \mathbb{R}^K$ be a multiclass classifier with continuously differentiable components f_i and let $c = \operatorname{argmax}_{1 \le i \le K} f_i(x_0)$ be the class which f predicts for x_0 . For all $\delta \in \mathbb{R}^d$ with

$$\|\boldsymbol{\delta}\|_p \le \min_{j \ne c} \frac{f_c(\boldsymbol{x_0}) - f_j(\boldsymbol{x_0})}{L_q^j},$$

argmax_{1 \le i \le K} $f_i(x_0 + \delta) = c$ holds with $\frac{1}{p} + \frac{1}{q} = 1, 1 \le p, q \le \infty$ and L_q^j is the Lipschitz constant for the function $f_c(x) - f_j(x)$ in ℓ_p norm.

Remark: $\beta_L = \min_{j \neq c} \frac{f_c(\boldsymbol{x}_0) - f_j(\boldsymbol{x}_0)}{L_q^j}$ is a lower bound of minimum distortion. $L_q^j = \max \|\nabla g(\boldsymbol{x})\|_q$, where $g(\boldsymbol{x}_0) = f_c(\boldsymbol{x}_0) - f_j(\boldsymbol{x}_0)$

Weng, et al. Evaluating the Robustness of Neural Networks: An Extreme Value Theory Approach. ICLR 2018

Δ

(1)

CLEVER score ≈

Minimal distortion Δ

Robustness Estimation: CLEVER Score

Efficiently estimate the Lipschitz constant by sampling around input + extreme value theory

```
Algorithm 1: CLEVER-t, compute CLEVER score for targeted attack
   Input: a K-class classifier f(x), data example x_0 with predicted class c, target class j, batch size
             N_b, number of samples per batch N_s, perturbation norm p, maximum perturbation R
   Result: CLEVER Score \mu \in \mathbb{R}_+ for target class j
1 S \leftarrow \{\emptyset\}, g(\boldsymbol{x}) \leftarrow f_c(\boldsymbol{x}) - f_j(\boldsymbol{x}), q \leftarrow \frac{p}{n-1}.
<sup>2</sup> for i \leftarrow 1 to N_b do
        for k \leftarrow 1 to N_{\circ} do
             randomly select a point \boldsymbol{x}^{(i,k)} \in B_n(\boldsymbol{x_0}, R)
4
             compute b_{ik} \leftarrow \|\nabla q(\boldsymbol{x}^{(i,k)})\|_q via back propagation
5
        end
        S \leftarrow S \cup \{\max_k \{b_{ik}\}\}
7
s end
         \leftarrow MLE of location parameter of reverse Weibull distribution on S
9 \hat{a}_W
10 \mu \leftarrow \min(\frac{g(\boldsymbol{x}_0)}{\hat{a}}, R)
   Algorithm 2: CLEVER-u, compute CLEVER score for un-targeted attack
   Input: Same as Algorithm 1, but without a target class j
   Result: CLEVER score \nu \in \mathbb{R}_+ for un-targeted attack
1 for j \leftarrow 1 to K, j \neq c do
        \mu_i \leftarrow \text{CLEVER-t}(f, \boldsymbol{x_0}, c, j, N_b, N_s, p, R)
2
3 end
```

- 4 $\nu \leftarrow \min_{i} \{\mu_i\}$
- CLEVER score enables robustness comparison between
 - different models
 - different datasets
 - different neural network architectures
 - different defense mechanisms

Demo: http://bigcheck.mybluemix.net/



Weng, et al. Evaluating the Robustness of Neural Networks: An Extreme Value Theory Approach. ICLR 2018

Robustness Evaluation Metrics



Robustness Verification: Nearest Neighbor Classifiers

Robustness verification for **ReLU network** (Katz et al., 2017) and tree ensemble (Kantchelian et al., 2016) are **NP-complete**. How about nearest neighbor classifiers?



Finding the minimum perturbation to make it closest to X_j $\epsilon^{(j)} = \min_{\delta} \frac{1}{2} \|\delta\|^2$ s.t. $\|\boldsymbol{z} + \boldsymbol{\delta} - \boldsymbol{x}_j\|^2 \le \|\boldsymbol{z} + \boldsymbol{\delta} - \boldsymbol{x}_i\|^2, \forall i \neq j$

A quadratic programming problem with linear constraint, which is polynomial time solvable.

$$oldsymbol{\epsilon}^{(j)} = \min_{oldsymbol{\delta}} rac{1}{2} \|oldsymbol{\delta}\|^2 \quad ext{s.t.} \quad \mathcal{A}^{(j)}oldsymbol{\delta} + oldsymbol{b}^{(j)} \geq 0$$

where
$$A_{i,:}^{(j)} = \mathbf{x}_j - \mathbf{x}_i$$
, $b_i^{(j)} = \frac{\|\mathbf{z} - \mathbf{x}_i\|^2 - \|\mathbf{z} - \mathbf{x}_j\|^2}{2}$.

O(n) quadratic problems, each QP has O(n) variables to solve, total time complexity >O(n^3) Speed up the algorithm by

- greedy coordinate ascent algorithm
- a screening rule to remove variables in each dual QP problem

 removing unimportant subproblems without solving them

ConvNet



The first robustness verification for nearest neighbor models

Wang et al., Evaluating the Robustness of Nearest Neighbor Classifiers: A Primal-Dual Perspective 2019

1-NN

Robustness Verification: Distance Metric Learning

The first adversarial verification method and the first certified defense for distance metric learning.

Compute a lower bound of the minimal adversarial perturbation of Mahalanobis K-NN

Theorem 1 (Robustness verification for Mahalanobis K-NN). Given a Mahalanobis K-NN classifier parameterized by a neighbor parameter K, a training dataset S and a positive semi-definite matrix M, for any instance $(\mathbf{x}_{test}, y_{test})$ we have

$$\epsilon^*(\boldsymbol{x}_{test}, y_{test}; \boldsymbol{M}) \ge \underset{j:y_j \neq y_{test}}{kth \min} \underset{i:y_i = y_{test}}{kth \max} \tilde{\epsilon}(\boldsymbol{x}_i, \boldsymbol{x}_j, \boldsymbol{x}_{test}; \boldsymbol{M}),$$
(10)

where kth max and kth min select the k-th maximum and k-th minimum respectively with k = (K+1)/2.

where
$$\tilde{\epsilon}(x^+, x^-, x; M) = \frac{d_M(x, x^-) - d_M(x, x^+)}{2\sqrt{(x^+ - x^-)^\top M^\top M(x^+ - x^-)}}.$$





Train a robust Mahalanobis distance with small certified and empirical robust errors

Algorithm 1: Adversarially robust metric learning (ARML)							
Input: Training data S , number of epochs T .							
Output: Positive semi-definite matrix M .							
1 Initialize G and M as identity matrices ;							
2 for $t = 0 \dots T - 1$ do							
3 Update G with the gradient							
$\mathbb{E}_{(\boldsymbol{x},\boldsymbol{y})\in\mathbb{S}}\nabla_{\boldsymbol{G}}\ell\left(\tilde{\epsilon}\left(\mathrm{randnear}_{\boldsymbol{M}}^{+}(\boldsymbol{x}),\mathrm{randnear}_{\boldsymbol{M}}^{-}(\boldsymbol{x}),\boldsymbol{x};\boldsymbol{G}^{\top}\boldsymbol{G}\right)\right);$							
4 Update M with the constraint $M = G^{\top}G$;							
5 end							

Table 2: Certified robust errors (left) and empirical robust errors (right) of Mahalanobis *K*-NN. The best (minimum) robust errors among all methods are in bold.

			Certified robust errors					Empirical robust errors					
	ℓ_2 -radius	0.000	0.500	1.000	1.500	2.000	2.500	0.000	0.500	1.000	1.500	2.000	2.500
	Euclidean	0.038	0.134	0.360	0.618	0.814	0.975	0.031	0.063	0.104	0.155	0.204	0.262
NO HOT	NCA	0.030	0.175	0.528	0.870	0.986	1.000	0.027	0.063	0.120	0.216	0.330	0.535
MNIST	LMNN	0.040	0.669	1.000	1.000	1.000	1.000	0.036	0.121	0.336	0.775	0.972	1.000
	ITML	0.106	0.731	0.943	1.000	1.000	1.000	0.084	0.218	0.355	0.510	0.669	0.844
	LFDA	0.237	1.000	1.000	1.000	1.000	1.000	0.215	1.000	1.000	1.000	1.000	1.000
	ARML (Ours)	0.034	0.101	0.276	0.537	0.760	0.951	0.032	0.055	0.077	0.109	0.160	0.213

ARML is more robust both provably (in terms of the certified robust error) **and empirically** (in terms of the empirical robust error).

45Wang et al., Provably Robust Metric Learning. NeurIPS 2020

Adversarial ML Problems from Industry's Point of View

Accuracy is not the sole metric to grade an AI model, neither is robustness

- The industry needs accurate, effective, robust, and sometimes fair and interpretable AI models.
- Scalability is one of the key problems

- Defense in the real-world
 - Robustness of robust models

Plug and Play robust module is necessary for protecting AI models that have been deployed



Thanks!

We Strive to Build Trustworthy Al Welcome to join us!

yijinfeng@jd.com

