Theoretically Principled Trade-off between Robustness and Accuracy

Hongyang Zhang, TTIC



Jiantao Jiao







Laurent Ghaoui



Michael I. Jordan



Simons Institute for the Theory of Computing Dec. 17th, 2019

Deep networks are unsafe



"Simons Institute" 87.7% confidence



=



"Simons Institute" 99.3% confidence human thinks

2





machine thinks



"TTIC" 99.3% confidence

Deep networks are unsafe



"panda"

Adversarial Noise



Adversarial Rotation

Adversarial Photographer

+

+



"gibbon"



"vulture"



"not hotdog"



=



"orangutan"



"hotdog"

[BCZOCG'18] Unrestricted Adversarial Example, 2018

Why are there adversarial examples?

• We use a wrong loss function







Non-Linear Case

Trade-off between Robustness and Accuracy

$$R_{rob}(f) := \mathbb{E}_{(X,Y)\sim D} \mathbb{1}\{\exists X' \in \mathbb{B}(X,\varepsilon) \ s. t. \ f(X')Y \le 0\}$$
$$R_{nat}(f) := \mathbb{E}_{(X,Y)\sim D} \mathbb{1}\{f(X)Y \le 0\}$$

• An example of trade-off:



Surrogate Loss

• Classification-calibrated loss ϕ :

$$H(\eta) := \min_{\alpha \in \mathbb{R}} (\eta \phi(\alpha) + (1 - \eta) \phi(-\alpha))$$
$$H^{-}(\eta) := \min_{\alpha : \alpha(2\eta - 1) \le 0} (\eta \phi(\alpha) + (1 - \eta) \phi(-\alpha))$$

Definition (classification-calibrated loss):

 ϕ is classification-calibrated loss, if for any $\eta \neq 1/2$, $H^{-}(\eta) > H(\eta)$.

Intuitive explanation:

- Think about η as $\eta(x) = \Pr[Y = +1 | X = x]$, and α as score of positive class by f
- Then $H(\eta) = \min_{f} R_{\phi}(f)$ $H^{-}(\eta) = \min_{f} R_{\phi}(f)$ s.t. *f* is inconsistent with Bayes optimal classifier
- Classification-calibrated loss: wrong classifier leads to larger loss for all $\eta(x)$

[BJM'06] Convexity, Classification, and Risk Bounds, 2006

Surrogate Loss



[BJM'06] Convexity, Classification, and Risk Bounds, 2006

7

Theorem 1 (Informal, upper bound, ZYJXGJ'19):

We have $R_{rob}(f) - R_{nat}^* \leq R_{\phi}(f) - R_{\phi}^* + \mathbb{E} \max_{X' \in \mathbb{B}(X,\varepsilon)} \phi(f(X')f(X)/\lambda).$

Proof Sketch:

• An important decomposition: $R_{rob}(f) = R_{nat}(f) + R_{bdy}(f)$ where $R_{bdy}(f) = \mathbb{E}_{(X,Y)\sim D} 1\{\exists X \in \varepsilon \text{ neighbour of } f \text{ s.t. } f(X)Y > 0\}$ $\int_{f^*}^{f^*} \int_{f^*}^{f^*} \int_{f^*}^{f^$

[ZYJXGJ'19] Theoretically Principled Trade-off between Robustness and Accuracy, ICML 2019

Theorem 1 (Informal, upper bound, ZYJXGJ'19):

We have $R_{rob}(f) - R_{nat}^* \leq R_{\phi}(f) - R_{\phi}^* + \mathbb{E} \max_{X' \in \mathbb{B}(X,\varepsilon)} \phi(f(X')f(X)/\lambda).$

Proof Sketch:

- An important decomposition: $R_{rob}(f) = R_{nat}(f) + R_{bdy}(f)$ where $R_{bdy}(f) = \mathbb{E}_{(X,Y)\sim D} \mathbb{1}\{\exists X \in \varepsilon \text{ neighbour of } f \text{ s. t. } f(X)Y > 0\}$
- $R_{rob}(f) R_{nat}^* = R_{nat}(f) R_{nat}^* + R_{bdy}(f)$
- $R_{nat}(f) R_{nat}^* \le R_{\phi}(f) R_{\phi}^*$ by [BJM'06]
- $R_{bdy}(f) \leq \mathbb{E} \max_{X' \in \mathbb{B}(X,\varepsilon)} \mathbb{1}(f(X')f(X) < 0) \leq \mathbb{E} \max_{X' \in \mathbb{B}(X,\varepsilon)} \phi(f(X')f(X)/\lambda)$

[BJM'06] Convexity, Classification, and Risk Bounds, 2006

[ZYJXGJ'19] Theoretically Principled Trade-off between Robustness and Accuracy, ICML 2019

Theorem 1 (Informal, upper bound, ZYJXGJ'19):

We have $R_{rob}(f) - R_{nat}^* \leq R_{\phi}(f) - R_{\phi}^* + \mathbb{E} \max_{X' \in \mathbb{B}(X,\varepsilon)} \phi(f(X')f(X)/\lambda).$

Theorem 2 (Informal, lower bound, ZYJXGJ'19):

There exist a data distribution, a classifier f, and an $\lambda > 0$ such that $R_{rob}(f) - R_{nat}^* \ge R_{\phi}(f) - R_{\phi}^* + \mathbb{E} \max_{X' \in \mathbb{B}(X,\varepsilon)} \phi(f(X')f(X)/\lambda).$

Theorem 1 (Informal, upper bound, ZYJXGJ'19):

We have $R_{rob}(f) - R_{nat}^* \leq R_{\phi}(f) - R_{\phi}^* + \mathbb{E} \max_{X' \in \mathbb{B}(X,\varepsilon)} \phi(f(X')f(X)/\lambda).$

• New Surrogate Loss (TRADES):



[ZYJXGJ'19] Theoretically Principled Trade-off between Robustness and Accuracy, ICML 2019

Significant Experimental Results

Experiments ---- CIFAR10

Defense	Defense type	Under which attack	Dataset	Distance	$\mathcal{A}_{\mathrm{nat}}(f)$	$\mathcal{A}_{ m rob}(f)$		
[BRRG18]	gradient mask	[ACW18]	CIFAR10	$0.031 \ (\ell_{\infty})$	-	0%		
[MLW ⁺ 18]	gradient mask	[ACW18]	CIFAR10	$0.031~(\ell_{\infty})$	-	5%		
[DAL+18]	gradient mask	[ACW18]	CIFAR10	$0.031~(\ell_{\infty})$	-	0%		
[SKN +18]	gradient mask	[ACW18]	CIFAR10	$0.031~(\ell_{\infty})$	-	9%		
[NKM17]	gradient mask	[ACW18]	CIFAR10	$0.015~(\ell_{\infty})$	-	15%		
[WSMK18]	robust opt.	FGSM ²⁰ (PGD)	CIFAR10	$0.031 \ (\ell_{\infty})$	27.07%	23.54%		
[MMS ⁺ 18]	robust opt.	FGSM ²⁰ (PGD)	CIFAR10	$0.031~(\ell_{\infty})$	87.30%	47.04%		
$\min_{f} \max_{X' \in B_{\varepsilon}(X)} \phi(Yf(X')) \text{(by Madry et al.)}$								
TRADES $(1/\lambda = 1)$	regularization	FGSM ²⁰ (PGD)	CIFAR10	$0.031 \ (\ell_{\infty})$	88.64%	49.14%		
TRADES $(1/\lambda = 6)$	regularization	FGSM ²⁰ (PGD)	CIFAR10	$0.031~(\ell_{\infty})$	84.92%	56.61%		
TRADES $(1/\lambda = 6)$ $\min_{f} [\mathbb{I}]$	$E \phi(Yf(X))$	$FGSM^{20} (PGD) + \mathbb{E} \max_{X' \in B_{\varepsilon}(X)} q$	b(f(X)f)	$\left[\begin{array}{c} 0.031 \ (\ell_{\infty}) \end{array} \right]$	84.92% (OUR	56.61% S)		
TRADES $(1/\lambda = 6)$ min_f TRADES $(1/\lambda = 6)$	regularization $E \phi(Yf(X))$ regularization	$FGSM^{20} (PGD) + \mathbb{E} \max_{X' \in B_{\mathcal{E}}(X)} q$ $LBFGSAttack$	$\begin{array}{ } \text{CIFAR10} \\ b(f(X)f \\ \text{CIFAR10} \\ \end{array}$	$\frac{0.031 (\ell_{\infty})}{(X'))/\lambda}$ $0.031 (\ell_{\infty})$	84.92% (OUR 84.92%	S)		
TRADES $(1/\lambda = 6)$ f TRADES $(1/\lambda = 6)$ TRADES $(1/\lambda = 1)$	regularization $E \phi(Yf(X))$ regularization regularization	$FGSM^{20} (PGD) + \mathbb{E} \max_{X' \in B_{\mathcal{E}}(X)} q$ $LBFGSAttack$ $MI-FGSM$	CIFAR10 $b(f(X)f)$ $CIFAR10$ $CIFAR10$ $CIFAR10$	$ \begin{array}{c} 0.031 \left(\ell_{\infty}\right) \\ \hline \left(X'\right) \right) / \lambda \\ 0.031 \left(\ell_{\infty}\right) \\ 0.031 \left(\ell_{\infty}\right) \end{array} $	84.92% (OUR 84.92% 88.64%	56.61% S) 81.58% 51.26%		
TRADES $(1/\lambda = 6)$ min_f TRADES $(1/\lambda = 6)$ TRADES $(1/\lambda = 1)$ TRADES $(1/\lambda = 6)$	regularization $E \phi(Yf(X))$ regularization regularization regularization	FGSM ²⁰ (PGD) + $\mathbb{E} \max_{X' \in B_{\mathcal{E}}(X)} q$ LBFGSAttack MI-FGSM MI-FGSM	CIFAR10 b(f(X)f) CIFAR10 CIFAR10 CIFAR10	$ \begin{array}{c} 0.031 \ (\ell_{\infty}) \\ \hline (X')) / \lambda \\ 0.031 \ (\ell_{\infty}) \\ 0.031 \ (\ell_{\infty}) \\ 0.031 \ (\ell_{\infty}) \end{array} $	84.92% (OUR 84.92% 88.64% 84.92%	56.61% S) 81.58% 51.26% 57.95%		
TRADES $(1/\lambda = 6)$ f TRADES $(1/\lambda = 6)$ TRADES $(1/\lambda = 1)$ TRADES $(1/\lambda = 6)$ TRADES $(1/\lambda = 1)$	regularization $E \phi(Yf(X))$ regularization regularization regularization regularization	$FGSM^{20} (PGD) + E \max_{X' \in B_{\mathcal{E}}(X)} q$ $LBFGSAttack$ $MI-FGSM$ $MI-FGSM$ $C&W$	CIFAR10 b(f(X)f) CIFAR10 CIFAR10 CIFAR10 CIFAR10	$\begin{array}{c} 0.031 \ (\ell_{\infty}) \\ \hline (X')) / \lambda \\ 0.031 \ (\ell_{\infty}) \\ 0.031 \ (\ell_{\infty}) \\ 0.031 \ (\ell_{\infty}) \\ 0.031 \ (\ell_{\infty}) \end{array}$	84.92% (OUR 84.92% 88.64% 84.92% 88.64%	56.61% S) 81.58% 51.26% 57.95% 84.03%		
TRADES $(1/\lambda = 6)$ min_f TRADES $(1/\lambda = 6)$ TRADES $(1/\lambda = 1)$ TRADES $(1/\lambda = 6)$ TRADES $(1/\lambda = 1)$ TRADES $(1/\lambda = 6)$	regularization $E \phi(Yf(X))$ regularization regularization regularization regularization regularization	$FGSM^{20} (PGD)$ $+ \mathbb{E} \max_{X' \in B_{\mathcal{E}}(X)} q$ $LBFGSAttack$ $MI-FGSM$ $MI-FGSM$ $C\&W$ $C\&W$	CIFAR10 b(f(X)f) CIFAR10 CIFAR10 CIFAR10 CIFAR10 CIFAR10	$\begin{array}{c} 0.031 \ (\ell_{\infty}) \\ \hline (X') \)/\lambda \\ 0.031 \ (\ell_{\infty}) \end{array}$	84.92% (OUR 84.92% 88.64% 84.92% 88.64% 84.92%	56.61% S) 81.58% 51.26% 57.95% 84.03% 81.24%		
TRADES $(1/\lambda = 6)$ f TRADES $(1/\lambda = 6)$ TRADES $(1/\lambda = 1)$ TRADES $(1/\lambda = 6)$ TRADES $(1/\lambda = 1)$ TRADES $(1/\lambda = 6)$ [SKC18]	regularization $E \phi(Yf(X))$ regularization regularization regularization regularization regularization gradient mask	$FGSM^{20} (PGD)$ $+ \mathbb{E} \max_{X' \in B_{\mathcal{E}}(X)} q$ $LBFGSAttack$ $MI-FGSM$ $MI-FGSM$ $C&W$ $C&W$ $C&W$ $[ACW18]$	CIFAR10 b(f(X)f) CIFAR10 CIFAR10 CIFAR10 CIFAR10 CIFAR10 MNIST	$\begin{array}{c} 0.031 \ (\ell_{\infty}) \\ \hline (X') \)/\lambda \\ 0.031 \ (\ell_{\infty}) \\ 0.005 \ (\ell_{2}) \end{array}$	84.92% (OUR 84.92% 88.64% 84.92% 88.64% 84.92% -	56.61% S) 81.58% 51.26% 57.95% 84.03% 81.24% 55%		
TRADES $(1/\lambda = 6)$	regularization $E \phi(Yf(X))$ regularization regularization regularization regularization regularization gradient mask robust opt.	$FGSM^{20} (PGD)$ $+ \mathbb{E} \max_{X' \in B_{\mathcal{E}}(X)} q$ $LBFGSAttack$ $MI-FGSM$ $MI-FGSM$ $C&W$ $C&W$ $C&W$ $[ACW18]$ $FGSM^{40} (PGD)$	CIFAR10 b(f(X)f) CIFAR10 CIFAR10 CIFAR10 CIFAR10 CIFAR10 MNIST MNIST	$\begin{array}{c} 0.031 \ (\ell_{\infty}) \\ \hline (X') \)/\lambda \\ 0.031 \ (\ell_{\infty}) \\ 0.005 \ (\ell_{2}) \\ 0.3 \ (\ell_{\infty}) \end{array}$	84.92% (OUR 84.92% 88.64% 84.92% 88.64% 84.92% - 99.36%	56.61% S) 81.58% 51.26% 57.95% 84.03% 81.24% 55% 96.01%		
TRADES $(1/\lambda = 6)$	regularization $E \phi(Yf(X))$ regularization regularization regularization regularization gradient mask robust opt. regularization	$FGSM^{20} (PGD)$ $+ \mathbb{E} \max_{X' \in B_{\mathcal{E}}(X)} q$ $LBFGSAttack$ $MI-FGSM$ $MI-FGSM$ $C&W$ $C&W$ $C&W$ $[ACW18]$ $FGSM^{40} (PGD)$ $FGSM^{40} (PGD)$	CIFAR10 b(f(X)f) CIFAR10 CIFAR10 CIFAR10 CIFAR10 CIFAR10 MNIST MNIST MNIST	$\begin{array}{c} 0.031 \ (\ell_{\infty}) \\ \hline (X') \)/\lambda \\ 0.031 \ (\ell_{\infty}) \\ 0.005 \ (\ell_{2}) \\ 0.3 \ (\ell_{\infty}) \\ 0.3 \ (\ell_{\infty}) \end{array}$	84.92% (OUR 84.92% 88.64% 84.92% 88.64% 84.92% - 99.36% 99.48%	56.61% 81.58% 51.26% 57.95% 84.03% 81.24% 55% 96.01% 96.07%		

Interpretability



(a) clean example

(b) adversarial example by boundary attack with random spatial transformation





 (b) adversarial example by boundary attack with random spatial transformation



the class

of bicycle

(c) clean example



(d) adversarial example by boundary attack with random spatial transformation



(e) clean example



(f) adversarial example by boundary attack with random spatial transformation



(c) clean example



(e) clean example



(d) adversarial example by boundary attack with random spatial transformation



(f) adversarial example by boundary attack with random spatial transformation

the class of bird

Competition: NeurIPS 2018 Adversarial Vision Challenge



geted

Attack

- Evaluation criterion
 - 400+ teams, ~2,000 submissions
 - Tiny ImageNet dataset
 - Model Track and Attack Track
 - Participants in the two tracks play against each other

Competition: NeurIPS 2018 Adversarial Vision Challenge



📱 Final Result



Recent Developments of TRADES

- Acceleration [SNG+19,ZZL+19]:
 - Achieve 30x speed-up, almost as fast as natural training
- Semi-supervised learning/unlabel data [CRS+19,SFK+19]:
 - TRADES + self-training (500K) improves robustness by +5% on CIFAR10
- Applications [JHC+19]:
 - 1st place in Glue leaderboard (up until Dec. 9th) in NLP --- SMART
- Theoretical understanding (upcoming):
 - Benefits of local Lipschitzness
 - Provable certification of TRADES by random smoothing

[SNG+19] Adversarial training for free, 2019. [ZZL-[CRS+19] Unlabel data improves adversarial robustness, 2019. [SFK+19] Are labels requires for improving adversarial robustness?, 2019. [JHC+19] SMART, 2019.

[ZZL+19] You only propagate once, 2019.

Conclusions

- Adversarial Robustness
 - Trade-off matters in the adversarial defense
 - Matching upper and lower bounds on $R_{rob}(f) R_{nat}^*$
 - New surrogate loss for adversarial defense
 - Winners of NeurIPS 2018 Adversarial Vision Challenge
 - Some recent developments

Thank You

Trade-off between Robustness and Accuracy

• Our goal: Find a classifier \hat{f} such that $R_{rob}(\hat{f}) \leq OPT + \delta$

OPT: =
$$\min_{f} R_{rob}(f)$$
, s.t. $R_{nat}(f) \le R_{nat}^* + \delta$
suffice to show $R_{rob}(f) - R_{nat}^* \le \delta$

PyTorch Package

New Surrogate Loss:

$\min_{f} \left[\mathbb{E} \phi \left(Y f(X) \right) + \mathbb{E} \max_{X' \in B_{\varepsilon}(X)} \phi (f(X) f(X') / \lambda) \right]$

Natural training:



Adversarial training by TRADES:

To apply TRADES, cd into the directory, put 'trades.py' to the directory.

from trades import trades_loss

```
def train(args, model, device, train_loader, optimizer, epoch):
    model.train()
   for batch idx, (data, target) in enumerate(train loader):
        data, target = data.to(device), target.to(device)
        optimizer.zero grad()
        # calculate robust loss - TRADES loss
        loss = trades loss(model=model,
                           x_natural=data,
                           y=target,
                           optimizer=optimizer,
                           step size=args.step size,
                           epsilon=args.epsilon,
                           perturb steps=args.num steps,
                           batch size=args.batch size,
                           beta=args.beta,
                           distance='l inf')
        loss.backward()
        optimizer.step()
```

Link: https://github.com/yaodongyu/TRADES

Goodle

Unrestricted Adversarial Examples Challenge Duil Dassing

In the Unrestricted Adversarial Examples Challenge, attackers submit arbitrary adversarial inputs, and defenders are expected to assign low confidence to difficult inputs while retaining high confidence and accuracy on a clean, unambiguous test set. You can learn more about the motivation and structure of the contest in our recent paper

This repository contains code for the warm-up to the challenge, as well as the public proposal for the contest. We are currently accepting defenses for the warm-up.

Warm-up & Contest Timeline







Defense	Submitted by	Clean data	Common corruptions	Spatial grid attack	SPSA attack	Boundary attack	Submission Date
Pytorch ResNet50 (trained on bird-or- bicycle extras)	TRADESv2	100.0%	100.0%	99.5%	100.0%	95.0%	Jan 17th, 2019 (EST)
Keras ResNet (trained on ImageNet)	Google Brain	100.0%	99.2%	92.2%	1.6%	4.0%	Sept 29th, 2018
Pytorch ResNet (trained on bird-or- bicycle extras)	Google Brain	98.8%	74.6%	49.5%	2.5%	8.0%	Oct 1st, 2018







Defense	Submitted by	Clean data	Common corruptions	Spatial grid attack	SPSA attack	Boundary attack	Submission Date
Pytorch ResNet50 (trained on bird-or- bicycle extras)	TRADESv2	100.0%	100.0%	99.5%	100.0%	95.0%	Jan 17th, 2019 (EST)
Keras ResNet (trained on ImageNet)	Google Brain	100.0%	99.2%	92.2%	1.6%	4.0%	Sept 29th, 2018
Pytorch ResNet (trained on bird-or- bicycle extras)	Google Brain	98.8%	74.6%	49.5%	2.5%	8.0%	Oct 1st, 2018







Defense	Submitted by	Clean data	Common corruptions	Spatial grid attack	SPSA attack	Boundary attack	Submission Date
Pytorch ResNet50 (trained on bird-or- bicycle extras)	TRADESv2	100.0%	100.0%	99.5%	100.0%	95.0%	Jan 17th, 2019 (EST)
Keras ResNet (trained on ImageNet)	Google Brain	100.0%	99.2%	92.2%	1.6%	4.0%	Sept 29th, 2018
Pytorch ResNet (trained on bird-or- bicycle extras)	Google Brain	98.8%	74.6%	49.5%	2.5%	8.0%	Oct 1st, 2018







Defense	Submitted by	Clean data	Common corruptions	Spatial grid attack	SPSA attack	Boundary attack	Submission Date
Pytorch ResNet50 (trained on bird-or- bicycle extras)	TRADESv2	100.0%	100.0%	99.5%	100.0%	95.0%	Jan 17th, 2019 (EST)
Keras ResNet (trained on ImageNet)	Google Brain	100.0%	99.2%	92.2%	1.6%	4.0%	Sept 29th, 2018
Pytorch ResNet (trained on bird-or- bicycle extras)	Google Brain	98.8%	74.6%	49.5%	2.5%	8.0%	Oct 1st, 2018







Defense	Submitted by	Clean data	Common corruptions	Spatial grid attack	SPSA attack	Boundary attack	Submission Date
Pytorch ResNet50 (trained on bird-or- bicycle extras)	TRADESv2	100.0%	100.0%	99.5%	100.0%	95.0%	Jan 17th, 2019 (EST)
Keras ResNet (trained on ImageNet)	Google Brain	100.0%	99.2%	92.2%	1.6%	4.0%	Sept 29th, 2018
Pytorch ResNet (trained on bird-or- bicycle extras)	Google Brain	98.8%	74.6%	49.5%	2.5%	8.0%	Oct 1st, 2018

Future Directions about Robustness

- Computational and Statistical Theory
 - Understand the optimization principal of new surrogate loss
 - (Tight) sample complexity of adversarial learning
- Applications of AI Security
 - Robotics, autonomous cars
 - Medical diagnose
- Extensions with other frameworks
 - Self-supervised/semi-supervised learning
 - Neural ODE