

Learning with Multiplicative Perturbations

Xiulong Yang & Shihao Ji Georgia State University



Background

• Adversarial Examples:



Perturbations are small, imperceptible to human.

Adversarial Training

• Improve robustness of DNNs against adversarial examples

$$\begin{aligned} \mathcal{L}_{\mathrm{AT}}(\boldsymbol{x}_{l}, y_{l}, \boldsymbol{r}_{\mathrm{adv}}, \boldsymbol{\theta}) &= D\left[h(y_{l}|\boldsymbol{x}_{l}), p(y|\boldsymbol{x}_{l} + \boldsymbol{r}_{\mathrm{adv}}, \boldsymbol{\theta})\right] \\ \text{with } \boldsymbol{r}_{\mathrm{adv}} &= \underset{\boldsymbol{r}; \|\boldsymbol{r}\| \leq \epsilon}{\operatorname{arg max}} D\left[h\left(y_{l}|\boldsymbol{x}_{l}\right), p\left(y|\boldsymbol{x}_{l} + \boldsymbol{r}, \boldsymbol{\theta}\right)\right], \end{aligned}$$

[Goodfellow et al. 2014]

 $\begin{aligned} \mathcal{L}_{\text{VAT}}(\boldsymbol{x}_*, \boldsymbol{r}_{\text{adv}}, \boldsymbol{\theta}) &= D\left[p(y | \boldsymbol{x}_*, \boldsymbol{\theta}), p(y | \boldsymbol{x}_* + \boldsymbol{r}_{\text{adv}}, \boldsymbol{\theta}) \right] \\ \text{with } \boldsymbol{r}_{\text{adv}} &= \underset{\boldsymbol{r}; \|\boldsymbol{r}\|_2 \leq \epsilon}{\operatorname{arg\,max}} D\left[p(y | \boldsymbol{x}_*, \boldsymbol{\theta}), p(y | \boldsymbol{x}_* + \boldsymbol{r}, \boldsymbol{\theta}) \right], \end{aligned}$

[Miyato et al. 2018]



• We propose a new type of adversarial perturbations:

Derive new loss functions:

xAT: xVAT:

$$\mathcal{L}_{xAT}(\boldsymbol{x}, \boldsymbol{z}_{xadv}, \boldsymbol{\theta}) = D\left[h(y|\boldsymbol{x}, \boldsymbol{\theta}), p(y|\boldsymbol{x} \odot \boldsymbol{z}_{xadv}, \boldsymbol{\theta})\right]$$

with $\boldsymbol{z}_{xadv} = \underset{\boldsymbol{z}}{\arg\max} D\left[h(y|\boldsymbol{x}, \boldsymbol{\theta}), p(y|\boldsymbol{x} \odot \boldsymbol{z}, \boldsymbol{\theta})\right]$

 $\mathcal{L}_{x\text{VAT}}(\boldsymbol{x}, \boldsymbol{z}_{x\text{adv}}, \boldsymbol{\theta}) = D\left[p(y|\boldsymbol{x}, \boldsymbol{\theta}), p(y|\boldsymbol{x} \odot \boldsymbol{z}_{x\text{adv}}, \boldsymbol{\theta})\right]$ with $\boldsymbol{z}_{x\text{adv}} = \underset{\boldsymbol{z}}{\arg\max} D\left[p(y|\boldsymbol{x}, \boldsymbol{\theta}), p(y|\boldsymbol{x} \odot \boldsymbol{z}, \boldsymbol{\theta})\right]$



• We use the L_0 -norm of z to regularize the learning:

$$egin{aligned} oldsymbol{z}_{ ext{xadv}} &= rg\max_{oldsymbol{z}} \Delta D(oldsymbol{z},oldsymbol{x},oldsymbol{ heta}) + \lambda \|oldsymbol{z}\|_{0} \ &= rg\max_{oldsymbol{z}} \Delta D(oldsymbol{z},oldsymbol{x},oldsymbol{ heta}) + \lambda \sum_{j=1}^P \mathbb{1}_{[z^j
eq 0]} \end{aligned}$$

• However, the discrete essence of z makes it undifferentiable.





• We adopt the <u>Stochastic Variational Optimization</u> and the <u>Hard</u> <u>Concrete Gradient Estimator</u> techniques for optimization.

$$\log \alpha_{\text{xadv}} = \underset{\log \alpha}{\arg \max} \mathbb{E}_{\boldsymbol{u} \sim \mathcal{U}(0,1)} \left[\Delta D(g(f(\log \alpha, \boldsymbol{u})), \boldsymbol{x}, \boldsymbol{\theta}) \right]$$
$$+ \lambda \sum_{j=1}^{P} \sigma \left(\log \alpha^{j} - \beta \log \frac{-\gamma}{\zeta} \right)$$
(12)

with

- $$\begin{split} f(\log \boldsymbol{\alpha}, \boldsymbol{u}) = &\sigma\left((\log \boldsymbol{u} \log(1 \boldsymbol{u}) + \log \boldsymbol{\alpha})/\beta\right)(\zeta \gamma) + \gamma,\\ g(\cdot) &= \min(1, \max(0, \cdot)), \end{split}$$
- Generate the mask: $\hat{z}_{xadv} = g(f(\log \alpha_{xadv}, u)), \qquad u \sim \mathcal{U}(0, 1).$

Transductive vs. Inductive Training



Fig. 2. The pipeline of the transductive and inductive implementations of multiplicative adversarial training.

Shrink or Expand





Fig. 3. The effect of ϵ on different perturbations. (a) shows that the additive perturbations are on the surface of a ball with the radius ϵ . (b) demonstrates that our multiplicative perturbations are distributed within the rectangle.

Efficient Computing

• Both AT and VAT resort to optimizing additive perturbations and classifier parameter alternatively in two steps.



(a) Additive Perturbation Pipeline

• xAT/xVAT can update them simultaneously in one step.



Experiments: Semi-supervised learning

TEST ACCURACIES OF SEMI-SUPERVISED LEARNING ON MNIST, SVHN AND CIFAR-10. THE RESULTS ARE AVERAGED OVER 5 RUNS.

Method	Tes MNIST N_l =100	st Accurac SVHN N_l =1000	y (%) CIFAR-10 N _l =4000
GAN with feature match [22]	99.07	91.89	81.37
CatGAN [23]	98.09	-	80.42
Ladder Networks [24]	98.94	-	79.60
П-model [14]	-	94.57	83.45
Mean Teacher [16]	-	94.79	82.26
VAT [6]	98.64	94.23	85.18
xVAT (Transductive)	98.02	93.99	85.82
xVAT (Inductive)	97.82	94.22	86.59

Experiments: Supervised learning

TEST ACCURACIES OF SUPERVISED LEARNING ON CIFAR-10 AND CIFAR-100. THE RESULTS ARE AVERAGED OVER 5 RUNS.

Method	Test Accuracy (%) CIFAR-10 CIFAR-100		
Baseline (MLE) [14]	93.24	73.58	
П-model [14]	94.44	73.68	
Temporal ensembling [14]	94.40	73.70	
AT, $L_{\infty}(\text{ours})^*$	93.90	74.04	
VAT [6]	94.19	75.02	
xAT (Inductive) xVAT (Inductive)	93.70 93.88	74.62 75.30	

THE TRAINING SPEEDS OF VAT AND XVAT ON THE FOUR BENCHMARK DATASETS. THE RESULTS ARE AVERAGED OVER 5 RUNS.

Method	Seconds per epoch			
	MNIST	SVHN	CIFAR-10	CIFAR-100
VAT (ours)*	4.31	54.3	51.3	51.5
xVAT (Transductive)	4.54	36.6	34.1	39.3
xVAT (Inductive)	4.33	35.7	33.6	34.4

Multiplicative vs. Additive Perturbations

The multiplicative perturbations are (1) More perceptible
(2) More interpretable

	x	$\log \alpha$	Zxadv	$x_{\rm xadv}$	radv	x _{adv}	
MNIST	3	3			A.	3	-0.8
	1	1			#72	1	- 0.4
	3	2			a Est	2	0.2
51			5255 T-5255				LT 10
	51	-51			195	31	- 0.8
		》在 全部	N			- 0.6	
30110	100	1			1.5.5.2	The second	- 0.4
	1	13.		1-2	1000	1	- 0.2
ey	ep		通貨用	Can St.		0	0.0
CIFAR-10	E/SUDY	1000			100	123	□ ¹⁰
	2.3	32 (S)	2.0			- 0.8	
	ALC: NO.		S = S			- 0.6	
	100	STRATES		183		- 0.4	
	8					- 0.2	
		的深刻		14		0.0	



Histograms of weights shows that xVAT learns a denser classifier from multiplicative perturbations with more non-zero weights than VAT and MLE, which may indicate the adversarially trained DNNs need more capacities (active neurons) to against multiplicative perturbations.



Fig. 6. Histograms of the classifier weights learned by MLE, VAT and xVAT on CIFAR-100. The histograms are computed from different CNN layers.



Thank you!