

Revisiting and Advancing Fast Adversarial Training through the Lens of Bi-level Optimization

Yihua Zhang^{1,*}, Guanhua Zhang^{2,*}, Prashant Khanduri³,

Mingyi Hong³, Shiyu Chang², Sijia Liu^{1,4}

¹Michigan State University, ²UC Santa Barbara, ³University of Minnesota, ⁴MIT-IBM Watson Lab

Poster Session: Hall E #528, Wednesday (Tonight)



PAPER



CODE

(Min-Max) Adversarial Training: Existing Principled Solution

- Nearly all existing work adopted the **Adversarial Training (AT)** framework [Madry et al. 2017], formulated as **min-max optimization**

Training over adversarially perturbed dataset

$$\text{minimize}_{\theta} \mathbb{E}_{(x,t) \sim D} \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \ell_{\text{tr}}(\theta; x + \delta, t) \right]$$

Sample-wise 'adversarial attack' generation

Assumption 1 (formulation level):

Attack type restriction: Must be the **opposite** of training objective.

Assumption 2 (computation level):

Each training step needs **multiple** gradient back-propagations to generate attacks.

(Min-Max) Adversarial Training: Existing Principled Solution

- In our paper, we focus on the following question:

How to advance the **algorithmic foundation** to advance Adversarial Training?

- Answer: Bi-level Optimization
- A properly designed formulation and solver will help scale up AT!

Bi-Level Optimization (BLO) Enables General AT Formulation

- Standard min-max formulation for adversarial training:

$$\min_{\theta} \mathbb{E}_{(x,t) \sim D} \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \ell_{\text{tr}}(\theta; x + \delta, t) \right]$$

- BLO-oriented adversarial training (AT)

Upper-level optimization $\min_{\theta} \ell_{\text{tr}}(\theta, \delta^*(\theta))$

Lower-level optimization s. t. $\delta^*(\theta) = \operatorname{argmin}_{\|\delta\|_{\infty} \leq \epsilon} \ell_{\text{atk}}(\theta, \delta)$

- Decouple Attack objective ℓ_{atk} from training objective ℓ_{tr}
- Why BLO? A possible framework of attack-agnostic robust training
A careful design of ℓ_{atk} can scale up adversarial training

Implicit Gradient --- The Tricky Part of BLO

- **Presence of implicit gradient (IG): the ‘fingerprint’ of BLO**

The upper-level gradient calculation:

$$\frac{d\ell_{\text{tr}}(\boldsymbol{\theta}, \boldsymbol{\delta}^*(\boldsymbol{\theta}))}{d\boldsymbol{\theta}} = \frac{\partial \ell_{\text{tr}}(\boldsymbol{\theta}, \boldsymbol{\delta}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} + \underbrace{\frac{d\boldsymbol{\delta}^*(\boldsymbol{\theta})^T}{d\boldsymbol{\theta}}}_{\text{IG}} \frac{\partial \ell_{\text{tr}}(\boldsymbol{\theta}, \boldsymbol{\delta}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\delta}}$$

IG: Gradient flow from lower-level to upper-level.

$\boldsymbol{\delta}^*(\boldsymbol{\theta})$ is an implicit function of $\boldsymbol{\theta}$

- The lower-level constraint makes BLO even harder!
- Properly designed lower-level objective makes it tractable!

- BLO-oriented adversarial training (AT)

$$\begin{aligned} & \text{minimize}_{\theta} \mathbb{E}_{x \sim D} [\ell_{\text{tr}}(\theta, \delta^*(\theta))] \\ & \text{subject to } \delta^*(\theta) = \operatorname{argmin}_{\delta \in C} \ell_{\text{atk}}(\theta, \delta) \end{aligned}$$

- BLO with customized lower-level attack objective

- **Linearization at \mathbf{z}** with quadratic regularization:

$$\ell_{\text{atk}}(\theta, \delta) = \langle \nabla_{\delta=\mathbf{z}} \ell_{\text{atk}}(\theta, \delta), \delta - \mathbf{z} \rangle + \left(\frac{\lambda}{2}\right) \|\delta - \mathbf{z}\|_2^2$$

- **Benefit:** Unique, computation-efficient, closed-form lower-level minimizer

$$\delta^*(\theta) = \operatorname{Proj}_C(\mathbf{z} - (1/\lambda) \nabla_{\delta} \ell_{\text{atk}}(\theta, \mathbf{z}))$$

Lower-level linearization leads to one-step PGD attack (**No SIGN**)!

Fast BAT

- **Derivation of implicit gradient (IG)** $\frac{d\delta^*(\theta)}{d\theta}$
 - **Key idea:** Extract implicit functions that involves IG from KKT conditions of lower-level problem
 - **Why is KKT tractable?** In robust training, the lower-level constraint $\delta \in \mathcal{C}$ is linear

Theorem 1 [Zhang et al., 2021]: With Hessian-free assumption, $\nabla_{\delta\delta} \ell_{\text{atk}}(\theta, \delta) = 0$

$$\frac{d\delta^*(\theta)^\top}{d\theta} = -(1/\lambda) \nabla_{\theta\delta} \ell_{\text{atk}}(\theta, \delta^*) \mathbf{H}_{\mathcal{C}}, \text{ with } \mathbf{H}_{\mathcal{C}} := \begin{bmatrix} 1_{p_1 < \delta_1^* < q_1} \mathbf{e}_1 & \cdots & 1_{p_d < \delta_d^* < q_d} \mathbf{e}_d \end{bmatrix}$$

$1_{p < \delta < q}$ is an indicator function, $p_i = \max\{-\epsilon, -x_i\}$, $q_i = \{\epsilon, 1 - x_i\}$

Fast BAT

- Fast Bi-level Adversarial Training (Fast BAT)

$$\text{minimize}_{\theta} E_{x \sim D}[\ell_{\text{tr}}(\theta, \delta^*(\theta))]$$

$$\text{subject to } \delta^*(\theta) = \operatorname{argmin}_{\delta \in \mathcal{C}} \langle \nabla_{\delta} \ell_{\text{atk}}(\theta, z), \delta - z \rangle + \left(\frac{\lambda}{2}\right) \|\delta - z\|_2^2$$

- Fast BAT algorithm:

- ❖ Fix θ , obtain lower-level solution $\delta^*(\theta)$

$$\delta^*(\theta) = \operatorname{Proj}_{\mathcal{C}}(z - (1/\lambda) \nabla_{\delta} \ell_{\text{atk}}(\theta, z))$$

(Single-step perturbation)

- ❖ Fix δ , obtain upper-level model update by SGD

$$\theta \leftarrow \theta - \alpha \frac{d\ell_{\text{tr}}(\theta, \delta^*(\theta))}{d\theta}$$

(IG-involved model updating)

$$\frac{d\ell_{\text{tr}}(\theta, \delta^*(\theta))}{d\theta} = \nabla_{\theta} \ell_{\text{tr}}(\theta, \delta^*(\theta)) + \underbrace{\frac{d\delta^*(\theta)^{\top}}{d\theta}}_{\text{IG}} \nabla_{\delta} \ell_{\text{tr}}(\theta, \delta^*(\theta))$$

$$\frac{d\delta^*(\theta)^{\top}}{d\theta} = -(1/\lambda) \nabla_{\theta \delta} \ell_{\text{atk}}(\theta, \delta^*) \mathbf{H}_{\mathcal{C}} \quad \text{(Theorem 1)}$$

Fast BAT vs. Linearization Type

- Fast BAT with gradient sign-based linearization

$$\text{minimize}_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim D} [\ell_{\text{tr}}(\boldsymbol{\theta}, \boldsymbol{\delta}^*(\boldsymbol{\theta}))]$$

$$\text{subject to } \boldsymbol{\delta}^*(\boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\delta} \in \mathcal{C}} \langle \text{sign}(\nabla_{\boldsymbol{\delta}} \ell_{\text{atk}}(\boldsymbol{\theta}, \mathbf{z})), \boldsymbol{\delta} - \mathbf{z} \rangle + \left(\frac{\lambda}{2}\right) \|\boldsymbol{\delta} - \mathbf{z}\|_2^2$$

- Why gradient sign?

Theorem 2: With sign-based linearization, Fast BAT simplifies to alternating optimization (**without** involving computation of **implicit gradients**)

$$\text{Lower-level: } \boldsymbol{\delta}^*(\boldsymbol{\theta}) = \operatorname{Proj}_{\mathcal{C}}(\mathbf{z} - (1/\lambda) \operatorname{sign}(\nabla_{\boldsymbol{\delta}} \ell_{\text{atk}}(\boldsymbol{\theta}, \mathbf{z})))$$

$$\text{Upper-level: } \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \left(\frac{\partial \ell_{\text{tr}}(\boldsymbol{\theta}, \boldsymbol{\delta}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} + \mathbf{0} \right)$$

Fast AT
[Wong et al., 2020]

Fast BAT + gradient sign-based linearization => Fast AT

Numerical Experiments of Fast BAT on CIFAR-10

Metrics: Standard Accuracy, Robust Accuracy (PGD/AutoAttack)

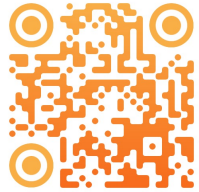
Baselines: Fast-AT, Fast-AT with Gradient Alignment, 2-step AT

CIFAR-10, PARN-18 trained with $\epsilon = 8/255$					
Method	SA (%)	RA-PGD (%)		RA-AA (%)	
		$\epsilon = 8$	$\epsilon = 16$	$\epsilon = 8$	$\epsilon = 16$
FAST-AT	82.39 ± 0.44	45.49 ± 0.41	9.56 ± 0.26	41.87 ± 0.15	7.91 ± 0.06
FAST-AT-GA	79.71 ± 0.44	47.27 ± 0.42	11.57 ± 0.32	43.24 ± 0.27	9.48 ± 0.15
PGD-2-AT	81.97 ± 0.41	44.62 ± 0.39	9.39 ± 0.32	41.73 ± 0.20	7.54 ± 0.25
FAST-BAT	79.97 ± 0.12	48.83 ± 0.17	14.00 ± 0.21	45.19 ± 0.12	11.51 ± 0.20
CIFAR-10, PARN-18 trained with $\epsilon = 16/255$					
FAST-AT	44.15 ± 7.27	37.17 ± 0.74	21.83 ± 1.32	31.66 ± 0.27	12.49 ± 0.33
FAST-AT-GA	58.29 ± 1.32	43.86 ± 0.67	26.01 ± 0.16	38.69 ± 0.56	17.97 ± 0.33
PGD-2-AT	68.04 ± 0.30	48.79 ± 0.31	24.30 ± 0.46	41.59 ± 0.22	15.40 ± 0.29
FAST-BAT	68.16 ± 0.25	49.05 ± 0.12	27.69 ± 0.16	43.64 ± 0.26	18.79 ± 0.24

- Fast BAT improves baselines in both SA and RA, and mitigates catastrophic overfitting!
- Improvement becomes more significant when facing stronger attack ($\epsilon = 16/255$)



PAPER



CODE

Poster Session: Hall E #528, Wednesday (Tonight)

thank you

tusind tak
謝謝 dakujem vám
ngiyabonga
merc
dziękuję
baie dankie
gracias
obrigada
obrigado
teşekkür ederim
شكرا
tack så mycket
suksema
danke
gràcies
tänan
dank u
mahalo
molte grazie
धन्यवाद