

Efficient local linearity regularization to overcome catastrophic overfitting

Elias Abad Rocamora¹, Fanghui Liu², Grigorios G. Chrysos³,
Pablo M. Olmos⁴, Volkan Cevher¹



Vienna 05/2024



Background: Adversarial Training (AT)

Considering a classifier f_{θ} parametrized by θ and a distribution \mathcal{D} on inputs \mathbf{x} and labels y .

- We are interested in **fast Adversarial Training (AT)** (Madry et al., 2018) :

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \mathcal{L}(f_{\theta}(\mathbf{x} + \delta), y) \right]$$



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



Background: Adversarial Training (AT)

Considering a classifier f_{θ} parametrized by θ and a distribution \mathcal{D} on inputs \mathbf{x} and labels y .

- We are interested in **fast** Adversarial Training (AT) (Madry et al., 2018) :

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \mathcal{L}(f_{\theta}(\mathbf{x} + \delta), y) \right]$$

- Solving the inner max slows down training:

$$\delta_{t+1} = \text{proj}_{\|\delta\|_{\infty} \leq \epsilon} [\delta_{t+1} + \gamma \cdot \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(f_{\theta}(\mathbf{x} + \delta_t), y))]$$

10 PGD steps are typically employed $\Rightarrow \times 10$ overhead! 🤔



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

Background: Catastrophic Overfitting (CO)

- As a cheap alternative, 1 PGD step (FGSM) can be used:

$$\delta_{\text{FGSM}} = \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{f}_{\theta}(\mathbf{x}), y))$$



Wong et al., [Fast is better than free: Revisiting adversarial training](#) *ICLR*, 2020.

Andriuschenko and Flammarion, [Understanding and Improving Fast Adversarial Training](#) *NeurIPS*, 2020.



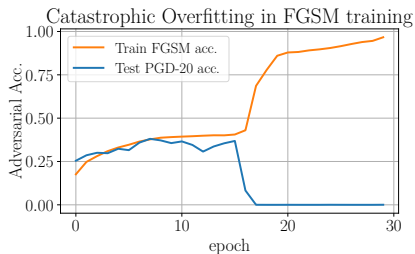
Background: Catastrophic Overfitting (CO)

- As a cheap alternative, 1 PGD step (FGSM) can be used:

$$\delta_{\text{FGSM}} = \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{f}_{\theta}(\mathbf{x}), y))$$

- Sadly, Catastrophic Overfitting (CO) appears:

- Example: PreActResNet18 on CIFAR10 at $\epsilon = 8/255$.
- Training results 100% robust to FGSM attacks and 0% robust to PGD-20 attacks 🤖



Wong et al., *Fast is better than free: Revisiting adversarial training* ICLR, 2020.

Andriuschenko and Flammarion, *Understanding and Improving Fast Adversarial Training* NeurIPS, 2020.



Background: Local linearity regularization

- Local linearity regularization can overcome CO. Let $\eta \sim \text{Unif}(-\epsilon, \epsilon)$:

Background: Local linearity regularization

- Local linearity regularization can overcome CO. Let $\boldsymbol{\eta} \sim \text{Unif}(-\epsilon, \epsilon)$:
- GradAlign:

$$1 - \text{cos-sim} [\nabla_{\mathbf{x}}\mathcal{L}(\mathbf{f}_{\theta}(\mathbf{x}), y), \nabla_{\mathbf{x}}\mathcal{L}(\mathbf{f}_{\theta}(\mathbf{x} + \boldsymbol{\eta}), y)]$$



Andriuschenko and Flammarion, Understanding and Improving Fast Adversarial Training *NeurIPS*, 2020.

- CURE:

$$\|\nabla_{\mathbf{x}}\mathcal{L}(\mathbf{f}_{\theta}(\mathbf{x}), y) - \nabla_{\mathbf{x}}\mathcal{L}(\mathbf{f}_{\theta}(\mathbf{x} + \boldsymbol{\eta}), y)\|_2$$



Moosavi-Dezfooli et al., Robustness via Curvature Regularization, and Vice Versa *CVPR*, 2019.

- LLR:

$$\left| \mathcal{L}(\mathbf{f}_{\theta}(\mathbf{x} + \boldsymbol{\eta}), y) - \mathcal{L}(\mathbf{f}_{\theta}(\mathbf{x}), y) + \boldsymbol{\eta}^{\top} \nabla_{\mathbf{x}}\mathcal{L}(\mathbf{f}_{\theta}(\mathbf{x}), y) \right|$$



Qin et al., Adversarial Robustness through Local Linearization *NeurIPS*, 2019.

Background: Local linearity regularization

- Local linearity regularization can overcome CO. Let $\boldsymbol{\eta} \sim \text{Unif}(-\epsilon, \epsilon)$:
- GradAlign:

$$1 - \text{cos-sim} [\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), y), \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\eta}), y)]$$



Andriuschenko and Flammarion, Understanding and Improving Fast Adversarial Training *NeurIPS*, 2020.

- CURE:

$$\|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), y) - \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\eta}), y)\|_2$$



Moosavi-Dezfooli et al., Robustness via Curvature Regularization, and Vice Versa *CVPR*, 2019.

- LLR:

$$\left| \mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\eta}), y) - \mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), y) + \boldsymbol{\eta}^{\top} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), y) \right|$$



Qin et al., Adversarial Robustness through Local Linearization *NeurIPS*, 2019.

- The problem:** Differentiating gradients ($\nabla_{\boldsymbol{\theta}} \nabla_{\mathbf{x}} \cdot$) is costly.

Background: Local linearity regularization

- Local linearity regularization can overcome CO. Let $\boldsymbol{\eta} \sim \text{Unif}(-\epsilon, \epsilon)$:
- GradAlign:

$$1 - \text{cos-sim} [\nabla_{\mathbf{x}}\mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), y), \nabla_{\mathbf{x}}\mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\eta}), y)]$$



Andriuschenko and Flammarion, Understanding and Improving Fast Adversarial Training *NeurIPS*, 2020.

- CURE:

$$\|\nabla_{\mathbf{x}}\mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), y) - \nabla_{\mathbf{x}}\mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\eta}), y)\|_2$$



Moosavi-Dezfooli et al., Robustness via Curvature Regularization, and Vice Versa *CVPR*, 2019.

- LLR:

$$\left| \mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\eta}), y) - \mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), y) + \boldsymbol{\eta}^{\top} \nabla_{\mathbf{x}}\mathcal{L}(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), y) \right|$$



Qin et al., Adversarial Robustness through Local Linearization *NeurIPS*, 2019.

- The problem:** Differentiating gradients ($\nabla_{\boldsymbol{\theta}}\nabla_{\mathbf{x}}\cdot$) is costly.
- The challenge:** *Can we efficiently regularize local linearity?*

Our Method: ELLE

- If a function h is locally linear in a convex set \mathcal{X} , it must satisfy:

$$h((1-\alpha)\mathbf{x}_a + \alpha\mathbf{x}_b) = (1-\alpha)h(\mathbf{x}_a) + \alpha h(\mathbf{x}_b), \quad \forall \alpha \in [0, 1], \forall \mathbf{x}_a, \mathbf{x}_b \in \mathcal{X}$$

Our Method: ELLE

- If a function h is locally linear in a convex set \mathcal{X} , it must satisfy:

$$h((1-\alpha)\cdot\mathbf{x}_a + \alpha\cdot\mathbf{x}_b) = (1-\alpha)\cdot h(\mathbf{x}_a) + \alpha\cdot h(\mathbf{x}_b), \quad \forall \alpha \in [0, 1], \forall \mathbf{x}_a, \mathbf{x}_b \in \mathcal{X}$$

- ELLE: Let $\mathbf{x}_a, \mathbf{x}_b \sim \text{Unif}(\mathbf{x} + \boldsymbol{\delta} : \|\boldsymbol{\delta}\|_\infty \leq \epsilon)$, $\alpha \sim \text{Unif}([0, 1])$ and $\mathbf{x}_c = (1 - \alpha) \cdot \mathbf{x}_a + \alpha \cdot \mathbf{x}_b$:

$$[(1 - \alpha) \cdot \mathcal{L}(\mathbf{f}_\theta(\mathbf{x}_a), y) + \alpha \cdot \mathcal{L}(\mathbf{f}_\theta(\mathbf{x}_b), y) - \mathcal{L}(\mathbf{f}_\theta(\mathbf{x}_c), y)]^2$$

Our Method: ELLE

- If a function h is locally linear in a convex set \mathcal{X} , it must satisfy:

$$h((1-\alpha)\cdot\mathbf{x}_a + \alpha\cdot\mathbf{x}_b) = (1-\alpha)\cdot h(\mathbf{x}_a) + \alpha\cdot h(\mathbf{x}_b), \quad \forall \alpha \in [0, 1], \forall \mathbf{x}_a, \mathbf{x}_b \in \mathcal{X}$$

- ELLE: Let $\mathbf{x}_a, \mathbf{x}_b \sim \text{Unif}(\mathbf{x} + \delta : \|\delta\|_\infty \leq \epsilon)$, $\alpha \sim \text{Unif}([0, 1])$ and $\mathbf{x}_c = (1 - \alpha) \cdot \mathbf{x}_a + \alpha \cdot \mathbf{x}_b$:

$$[(1 - \alpha) \cdot \mathcal{L}(\mathbf{f}_\theta(\mathbf{x}_a), y) + \alpha \cdot \mathcal{L}(\mathbf{f}_\theta(\mathbf{x}_b), y) - \mathcal{L}(\mathbf{f}_\theta(\mathbf{x}_c), y)]^2$$

- Our regularization term does not involve differentiating gradients 😊.

Our Method: ELLE

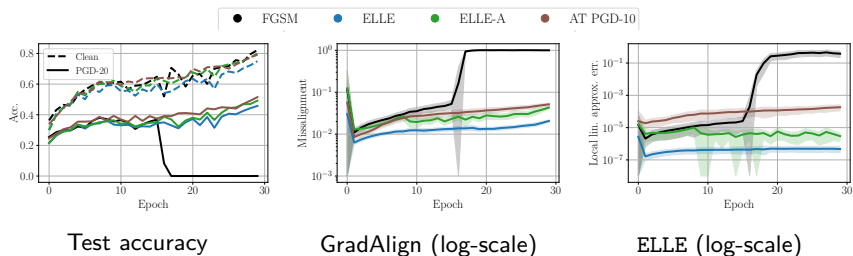


Figure: PreActResNet in CIFAR10 at $\epsilon = 8/255$

Our Method: ELLE

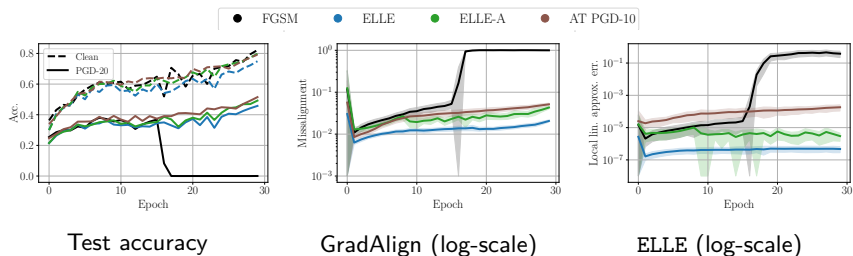


Figure: PreActResNet in CIFAR10 at $\epsilon = 8/255$

- Our local linearity metric follows the one of GradAlign at a reduced cost.

Our Method: ELLE

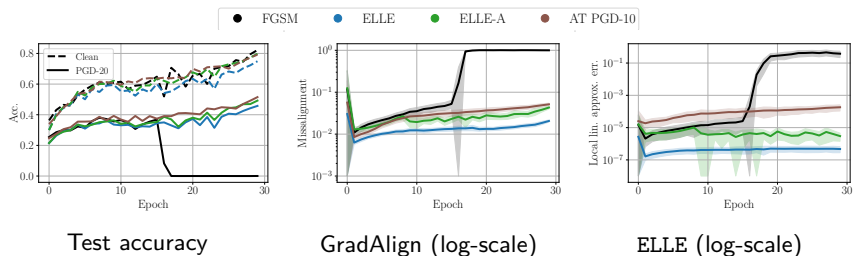
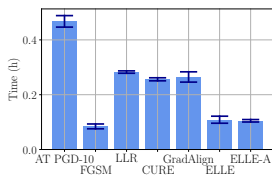


Figure: PreActResNet in CIFAR10 at $\epsilon = 8/255$

- Our local linearity metric follows the one of GradAlign at a reduced cost.
- When regularized (● and ●), CO is overcome.

Comparison with other regularization terms

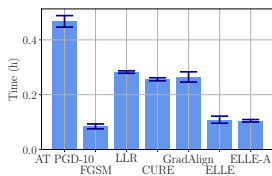


Runtime comparison

Method	8		16	
	AA	Clean	AA	Clean
LLR	42.18 ± (0.20)	75.02 ± (0.09)	16.92 ± (0.20)	42.81 ± (9.62)
CURE	43.60 ± (0.17)	77.74 ± (0.11)	<u>18.25</u> ± (0.45)	52.49 ± (0.04)
GradAlign	44.66 ± (0.21)	80.50 ± (0.07)	17.46 ± (1.71)	44.35 ± (15.32)
ELLE	42.78 ± (0.95)	<u>80.13</u> ± (0.32)	18.28 ± (0.17)	59.73 ± (0.16)
ELLE-A	<u>44.32</u> ± (0.04)	79.81 ± (0.10)	18.03 ± (0.15)	<u>59.21</u> ± (1.23)
AT PGD-10	46.95 ± (0.11)	79.11 ± (0.08)	24.77 ± (0.26)	59.64 ± (0.46)

PreActResNet18 in CIFAR10

Comparison with other regularization terms



Method	8		16	
	AA	Clean	AA	Clean
LLR	42.18 ± (0.20)	75.02 ± (0.09)	16.92 ± (0.20)	42.81 ± (9.62)
CURE	43.60 ± (0.17)	77.74 ± (0.11)	<u>18.25</u> ± (0.45)	52.49 ± (0.04)
GradAlign	44.66 ± (0.21)	80.50 ± (0.07)	17.46 ± (1.71)	44.35 ± (15.32)
ELLE	42.78 ± (0.95)	<u>80.13</u> ± (0.32)	18.28 ± (0.17)	59.73 ± (0.16)
ELLE-A	<u>44.32</u> ± (0.04)	79.81 ± (0.10)	18.03 ± (0.15)	<u>59.21</u> ± (1.23)
AT PGD-10	46.95 ± (0.11)	79.11 ± (0.08)	24.77 ± (0.26)	59.64 ± (0.46)

Runtime comparison

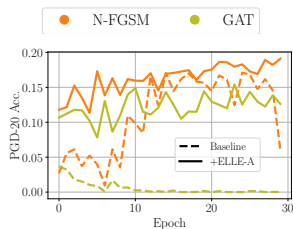
PreActResNet18 in CIFAR10

- ELLE attains comparable performance with negligible overhead.

Combination with other methods

		ϵ		16		26	
Method	Model	AA	Clean	AA	Clean	AA	Clean
GAT	PRN	0.54 \pm (0.53)	78.52 \pm (0.25)	0.01 \pm (0.00)	84.35 \pm (0.34)		
GAT+ELLE-A		13.83 \pm (4.63)	65.71 \pm (2.48)	6.56 \pm (2.79)	58.37 \pm (2.39)		
N-FGSM		20.59 \pm (0.21)	61.24 \pm (0.26)	10.96 \pm (0.26)	37.73 \pm (0.32)		
N-FGSM+ELLE-A		20.48 \pm (0.57)	61.21 \pm (0.14)	12.03 \pm (1.02)	26.77 \pm (3.25)		
AT PGD-10		24.77 \pm (0.26)	59.64 \pm (0.46)	14.42 \pm (0.00)	34.90 \pm (0.61)		
GAT	WRN	0.95 \pm (0.09)	84.09 \pm (0.10)	0.00 \pm (0.00)	89.04 \pm (0.28)		
GAT+ELLE-A		17.30 \pm (1.20)	64.27 \pm (3.56)	6.74 \pm (2.08)	42.04 \pm (10.52)		
N-FGSM		20.54 \pm (0.18)	63.83 \pm (1.24)	3.31 \pm (2.58)	27.96 \pm (9.36)		
N-FGSM+ELLE-A		21.28 \pm (0.07)	64.25 \pm (0.20)	12.22 \pm (0.25)	33.50 \pm (0.14)		
AT PGD-10		26.77 \pm (0.28)	64.97 \pm (0.09)	14.61 \pm (0.10)	36.30 \pm (0.62)		

CIFAR10 *Short* schedule



WRN in CIFAR10 at $\epsilon = \frac{26}{255}$



Sriramanan et al., Guided adversarial attack for evaluating and enhancing adversarial defenses *NeurIPS*, 2020.



de Jorge et al., Make some noise: Reliable and efficient single-step adversarial training *NeurIPS*, 2022.

Combination with other methods

		ϵ		16		26	
Method	Model	AA	Clean	AA	Clean		
GAT	PRN	0.54 \pm (0.53)	78.52 \pm (0.25)	0.01 \pm (0.00)	84.35 \pm (0.34)		
GAT+ELLE-A		13.83 \pm (4.63)	65.71 \pm (2.48)	6.56 \pm (2.79)	58.37 \pm (2.39)		
N-FGSM		20.59 \pm (0.21)	61.24 \pm (0.26)	10.96 \pm (0.26)	37.73 \pm (0.32)		
N-FGSM+ELLE-A		20.48 \pm (0.57)	61.21 \pm (0.14)	12.03 \pm (1.02)	26.77 \pm (3.25)		
AT PGD-10		24.77 \pm (0.26)	59.64 \pm (0.46)	14.42 \pm (0.00)	34.90 \pm (0.61)		
GAT	WRN	0.95 \pm (0.09)	84.09 \pm (0.10)	0.00 \pm (0.00)	89.04 \pm (0.28)		
GAT+ELLE-A		17.30 \pm (1.20)	64.27 \pm (3.56)	6.74 \pm (2.08)	42.04 \pm (10.52)		
N-FGSM		20.54 \pm (0.18)	63.83 \pm (1.24)	3.31 \pm (2.58)	27.96 \pm (9.36)		
N-FGSM+ELLE-A		21.28 \pm (0.07)	64.25 \pm (0.20)	12.22 \pm (0.25)	33.50 \pm (0.14)		
AT PGD-10		26.77 \pm (0.28)	64.97 \pm (0.09)	14.61 \pm (0.10)	36.30 \pm (0.62)		

CIFAR10 *Short* schedule

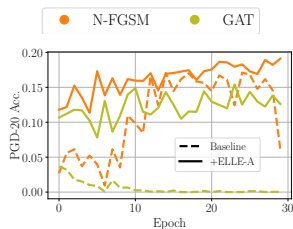
- Plugging our regularization term stabilizes and improves the performance of popular single-step AT methods.



Sriramanan et al., Guided adversarial attack for evaluating and enhancing adversarial defenses *NeurIPS*, 2020.



de Jorge et al., Make some noise: Reliable and efficient single-step adversarial training *NeurIPS*, 2022.



WRN in CIFAR10 at $\epsilon = \frac{26}{255}$

Thanks

Thanks for your attention!

contact: elias.abadrocamora@epfl.ch

supported by:



HASLERSTIFTUNG



The
Alan Turing
Institute

