

On Fragile Features and Batch Normalization in Adversarial Training

Nils Philipp Walter David Stutz Bernt Schiele

Motivation

1. BN argued to **cause adversarial vulnerability**
2. **BN statistics** of clean adversarial examples **differ**
3. Non-robust Neural Networks contain robust sub-networks
4. Training BN and only BN (i.e. random filters) achieves decent accuracy

Research Question

Is it possible to **adversarially finetune** BN-layers to disable **fragile features** and boost **adversarial robustness**?

Experimental Setup

Compare adversarially fine-tuning BN (BN PARAMS) to standard baseline (NORMAL/ADV) and training BN and only BN (ADV BN ONLY):

Training	Model	All	BN Prms	BN stats	%prms
From scratch	NORMAL/ADV	✓	✓		100.00
	ADV BN ONLY	✗	✗		0.15
Fine-tune	BN STATS	✗	✗	✓	(0.15)
	BN PARAMS	✗	✓	✓	0.15

Implementation details:

- ResNet20 on CIFAR10
- Evaluation with AutoAttack, L_infty eps = 8/255
- AutoAugment Normal, crop+flips for Adv

Quantitive Results

Model	Test Error	Robust Test Error
NORMAL	4.11	99.8
ADV	17.67	58.4
+BN ONLY	43.33	98.0
+BN ONLY + log + conv1	35.22	98.1
BN STATS	6.33	99.6
BN PARAMS	29.67	70.5
BN PARAMS + log + conv1	26.89	68.2

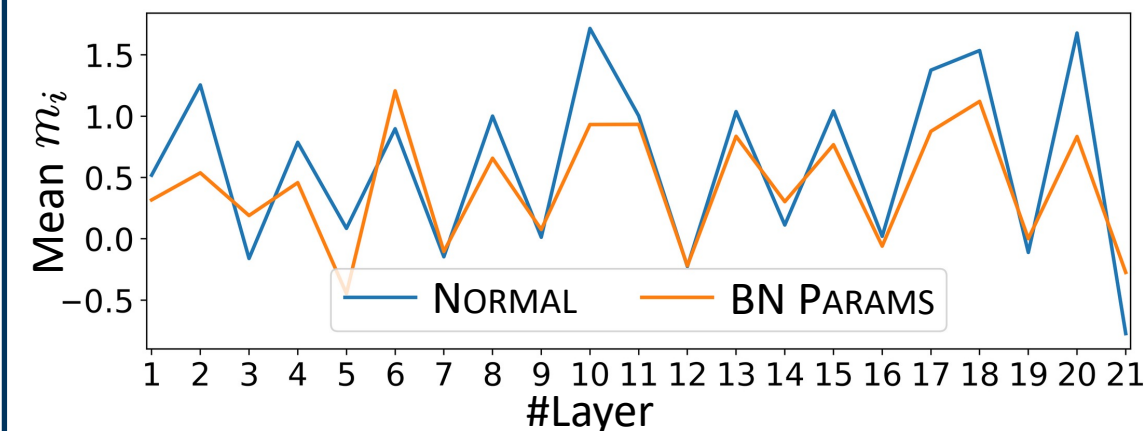
Hypothesis: BN PARAMS Turns Off (Fragile) Features

To compare BN parameters, rewrite BN as:

$$\text{BN}(z_i) = \left(\frac{\gamma_i}{\sqrt{\sigma_i^2 + \epsilon}} \right) z_i + \left(-\frac{\gamma_i \mu_i}{\sqrt{\sigma_i^2 + \epsilon}} + \beta_i \right)$$

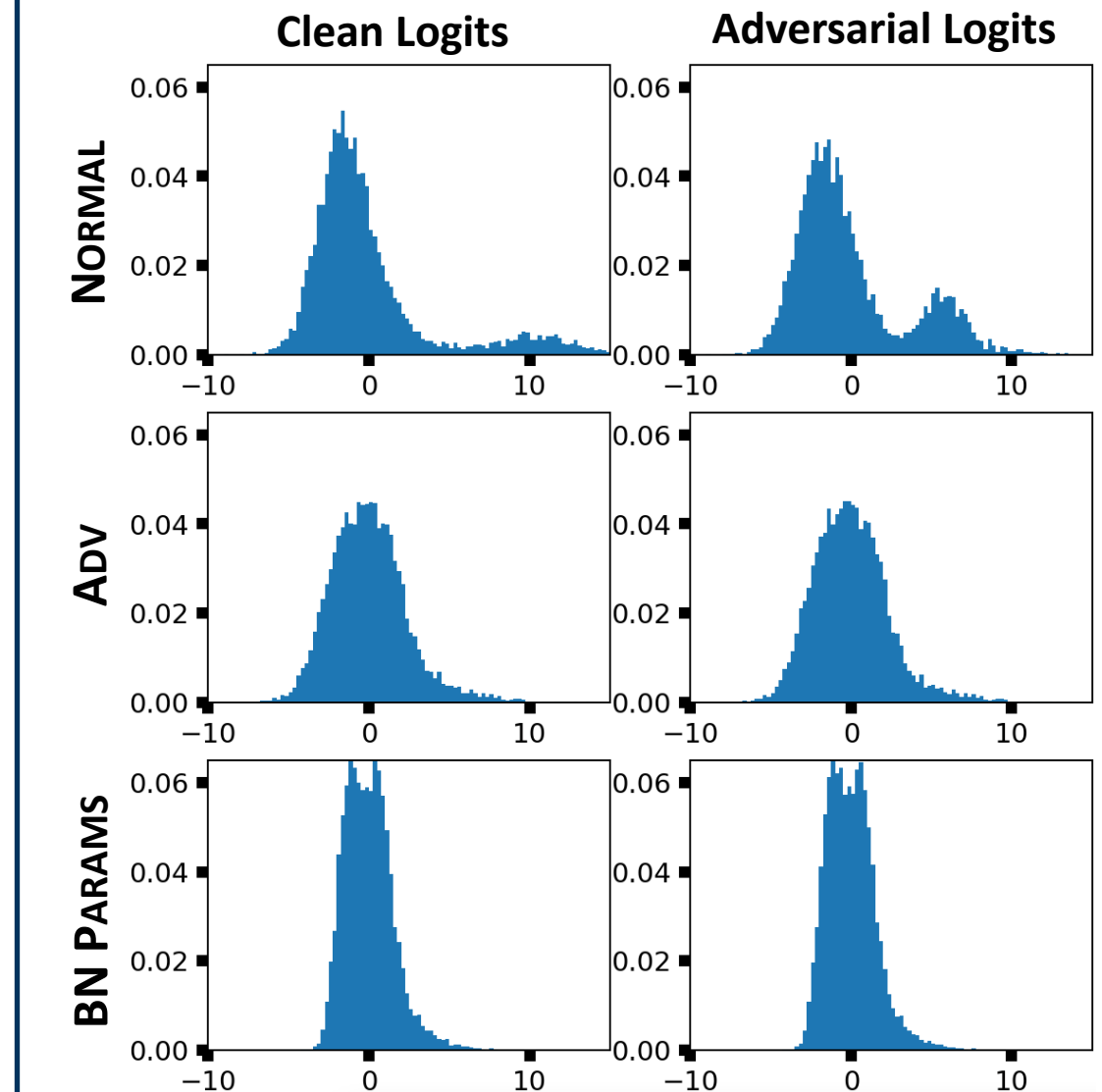
$$= m_i \cdot z_i + b_i$$

→ Analysis can be reduced to m_i and b_i



- BN PARAMS seems to disable part of the fragile features
- Average difference in mean weight is -0.177

BN PARAMS Mimics Logits of Adv



Conclusion

- Only adjusting statistics has no effect
- BN PARAMS obtains non-trivial adversarial robustness
- Indicates that NORMAL learns some robust features and fragile ones can be ignored
- Training BN on random features (ADV BN ONLY) does not provide adversarial robustness