
Adversarial Robustness of Convolutional Networks

Fergus O’Hanlon

Mathematical Institute

University of Oxford

Oxford, OX2 6GG

fergus.ohanlon@st-annes.ox.ac.uk

Abstract

Convolutional Neural Networks (CNNs) have been immensely successful in the past few decades. In particular, the introduction of Residual Networks (ResNets) was significantly influential. However, one of the flaws with CNNs is that they are susceptible to adversarial attacks. It is known that there is a relationship between natural accuracy and robust accuracy. Furthermore, increasing the width of ResNet is known to increase performance. After presenting "Convolutional Normalization", we explore the effect of normalization techniques on natural and robust accuracy and investigate their relationship with width for ResNet.

1 Introduction

Convolutional Neural Networks (CNNs) have been a major focus in the field of deep learning since the introduction of LeNet-5 [1]. Major breakthroughs in image, video, speech, and audio processing have been made in the past couple of decades with CNNs [2]. Since then, it has been found that certain small perturbations to a neural network’s input can cause it to misclassify the output [3]. Such perturbations are found by optimizing the input to maximise the prediction error and the perturbed examples are so called "adversarial examples". Normalization plays an important role in modern network architecture. For example, batch normalization (BN) [4] is an essential component for training deep CNNs. Convolutional Normalization (ConvNorm) is a new technique specific for CNNs, which this report aims to present and explore.

2 Framework

2.1 Deep Networks

A deep network $f(\cdot)$ of L layers is a composition of a series of mappings: $f(\mathbf{x}) = f^{L-1} \circ \dots \circ f^0(\mathbf{x})$. Each layer $f^\ell(\cdot)$ consists of a linear transform $\mathcal{A}^\ell(\cdot)$ followed by a nonlinear activation function $\varphi(\cdot)$. The output of each layer is given by the following recursion:

$$\mathbf{z}^{\ell+1} := f^\ell(\mathbf{z}^\ell) = \varphi \circ \mathcal{A}^\ell(\mathbf{z}^\ell)$$

for $\ell = 0, 1, \dots, L-1$, with $\mathbf{z}_0 = \mathbf{x}$. The goal of deep learning is to train a network which can be captured by the following minimisation problem:

$$\min_{\theta \in \Theta} \mathcal{L}(\theta; \{\mathbf{x}^i, \mathbf{y}^i\}_{i=1}^m) := \frac{1}{m} \sum_{i=1}^m \ell(f(\mathbf{x}^i, \theta), \mathbf{y}^i)$$

where θ denotes the networks parameters in $\{\mathcal{A}^\ell(\cdot)\}_{\ell=0}^{L-1}$ and ℓ is a chosen loss function.

2.2 Convolutional Neural Networks and Convolution Operators

CNNs consist of convolutional layers which implement the linear transform by convolving the input. Suppose that the input data \mathbf{x} has C output channels, $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_C)$. The linear operator for the ℓ 'th layer, $\mathcal{A}^\ell(\cdot) : \mathbb{R}^{C_\ell \times m} \mapsto \mathbb{R}^{C_{\ell+1} \times m}$ is a convolution operation with $C_{\ell+1}$ output channels,

$$\mathbf{z}^{\ell+1} = (\mathbf{z}_1^{\ell+1}, \dots, \mathbf{z}_{C_{\ell+1}}^{\ell+1}), \quad \mathbf{z}_k^{\ell+1} = \varphi \left(\sum_{j=1}^{C_\ell} \mathbf{a}_{kj}^\ell * \mathbf{z}_j^\ell \right) \quad (k = 1, \dots, C_{\ell+1})$$

Now, since any linear convolution can be reduced to a circular convolution via zero-padding, we consider $*$ as a circular convolution such that

$$\mathbf{y} = \mathbf{a} * \mathbf{x} = \mathbf{C}_\mathbf{a} \cdot \mathbf{x}$$

where $\mathbf{C}_\mathbf{a} \in \mathbb{R}^{C \times C}$ is a circulant matrix whose first column is \mathbf{a} . $\mathbf{C}_\mathbf{a}$ can be decomposed via the discrete Fourier Transform (DFT) matrix \mathbf{F} :

$$\mathbf{C}_\mathbf{a} = \mathbf{F}^* \text{diag}(\hat{\mathbf{a}}) \mathbf{F}, \quad \hat{\mathbf{a}} = \mathbf{F} \mathbf{a}$$

This means that the convolution $\mathbf{a} * \mathbf{x}$ can be computed efficiently in the frequency domain using the fast fourier transform (FFT).

2.3 Adversarial Attacks

2.3.1 Fast Gradient Sign Method

The following method to generate adversarial examples is referred to as the "fast gradient sign method" (FGSM) [5]:

$$\tilde{\mathbf{x}} := \epsilon \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}, \mathbf{y}))$$

Where $\mathcal{L}(\boldsymbol{\theta}, \mathbf{x}, \mathbf{y})$ is the cost, $\boldsymbol{\theta}$ is the parameter vector, and \mathbf{y} is the target associated with \mathbf{x} .

2.3.2 Projected Gradient Descent

Where FGSM is a one-step scheme, [6] established a more powerful, multi-step scheme by considering projected gradient descent (PGD) on the negative loss function:

$$\mathbf{x}^{t+1} = \Pi_{\mathbf{x}+\mathcal{S}} (\mathbf{x}^t + \alpha \text{sign} \nabla_{\mathbf{x}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}, \mathbf{y}))$$

where $\Pi_{\mathbf{x}+\mathcal{S}}$ denotes projecting perturbations into the set \mathcal{S} and α is the step size.

3 Convolutional Normalization

ConvNorm was proposed by [7] which exploits the convolution structures of CNNs and exploits translation-invariance properties of convolutional operators. The technique improves training and robustness through reducing the network's Lipschitzness [8] by normalizing each output channel using preconditioning.

To introduce this method, consider a vanilla CNN using 1-stride and a layer ℓ that has C_I and C_O input and output channels, respectively. The output of the k -th channel before the nonlinear activation can be written as (omitting the layer number ℓ):

$$\mathbf{z}_{k,out} = \sum_{j=1}^{C_I} \mathbf{a}_{kj} * \mathbf{z}_{j,in} = [\mathbf{C}_{\mathbf{a}_{k1}} \quad \mathbf{C}_{\mathbf{a}_{k2}} \quad \dots \quad \mathbf{C}_{\mathbf{a}_{kC_I}}] \cdot \begin{bmatrix} \mathbf{z}_{1,in} \\ \mathbf{z}_{2,in} \\ \vdots \\ \mathbf{z}_{C_I,in} \end{bmatrix} = \mathbf{A}_k \cdot \mathbf{z}_{in}$$

where $\mathbf{A}_k := [\mathbf{C}_{\mathbf{a}_{k1}} \quad \mathbf{C}_{\mathbf{a}_{k2}} \quad \dots \quad \mathbf{C}_{\mathbf{a}_{kC_I}}]$ and $\mathbf{z}_{in} := [\mathbf{z}_{1,in} \quad \mathbf{z}_{2,in} \quad \dots \quad \mathbf{z}_{C_I,in}]^\top$

The idea of convnorm is to multiply each output channel $\mathbf{z}_{k,out}$ ($k = 1, \dots, C_O$) by a preconditioning matrix \mathbf{P}_k , as defined in equation (1) below. This is inspired by [9] which showed that normalizing

the output z eliminates bad local minimizers and dramatically improves the optimization landscape for learning the kernel \mathbf{a} .

$$\mathbf{P}_k = \left(\sum_{j=1}^{C_I} \mathbf{C}_{\mathbf{a}_{kj}} \mathbf{C}_{\mathbf{a}_{kj}}^\top \right)^{-1/2} = \left(\mathbf{A}_k \mathbf{A}_k^\top \right)^{-1/2} \quad (1)$$

The resulting normalized output is thus given by

$$\tilde{\mathbf{z}}_{k,out} := \mathbf{P}_k \mathbf{z}_{k,out} = (\mathbf{A}_k \mathbf{A}_k^\top)^{-1/2} \mathbf{A}_k \cdot \mathbf{z}_{in} = \mathbf{Q}_k(\mathbf{A}_k) \cdot \mathbf{z}_{k,in}$$

where $\mathbf{Q}_k(\mathbf{A}_k) := (\mathbf{A}_k \mathbf{A}_k^\top)^{-1/2} \mathbf{A}_k$. ConvNorm is essentially a reparameterization of the kernels $\{\mathbf{a}_{kj}\}_{j=1}^{C_I}$ for the k -th channel. Expanding the definition of $\mathbf{C}_{\mathbf{a}_{kj}}$ in equation (1) provides:

$$\begin{aligned} \mathbf{P}_k &= \left(\sum_{j=1}^{C_I} \mathbf{F}^* \text{diag}(\hat{\mathbf{a}}_{kj}) \mathbf{F} \mathbf{F}^* \text{diag}(\hat{\mathbf{a}}_{kj}) \mathbf{F} \right)^{-1/2} \\ &= \mathbf{F}^* \left(\text{diag} \left(\sum_{j=1}^{C_I} |\hat{\mathbf{a}}_{kj}|^2 \right) \right)^{\odot -1/2} \mathbf{F} = \mathbf{F}^* (\text{diag}(\mathbf{F} \mathbf{v}_k)) \mathbf{F} = \mathbf{C}_{\mathbf{v}_k} \end{aligned}$$

with $\mathbf{v}_k := \mathbf{F}^{-1} \left(\sum_{i=1}^{C_I} |\hat{\mathbf{a}}_{ki}|^2 \right)^{\odot -1/2}$ and \odot denoting entrywise operation. Rewriting $\mathbf{Q}_k(\mathbf{A}_k)$ as:

$$\mathbf{Q}_k(\mathbf{A}_k) = \mathbf{P}_k \mathbf{A}_k = [\mathbf{C}_{\mathbf{v}_k * \mathbf{a}_{k1}} \quad \dots \quad \mathbf{C}_{\mathbf{v}_k * \mathbf{a}_{kC_I}}]$$

means that the reparameterization can be written in terms of convolutions which can be efficiently computed using FFT. Specifically, for each $j = 1, \dots, C_I$:

$$\mathbf{C}_{\mathbf{v}_k * \mathbf{a}_{kj}} = \mathbf{F}^* \text{diag}(\hat{\mathbf{g}}(\mathbf{a}_{kj})) \mathbf{F}, \quad \hat{\mathbf{g}}(\mathbf{a}_{kj}) = \hat{\mathbf{a}}_{kj} \odot \left(\sum_{i=1}^{C_I} |\hat{\mathbf{a}}_{ki}|^2 \right)^{\odot -1/2}$$

The matrix \mathbf{A}_k can't be normalized to exact orthogonal since it is overcomplete. However it can be normalized to tight frame with $\mathbf{Q}_k \mathbf{Q}_k^\top = \mathbf{I}$ and so $\|\mathbf{Q}_k\|$ can be normalized in each channel to unity.

Applying the normalization over every channel provides the overall ConvNorm normalization as follows:

$$\tilde{\mathbf{z}}_{out} = \begin{bmatrix} \mathbf{P}_1 \mathbf{z}_{1,out} \\ \vdots \\ \mathbf{P}_{C_O} \mathbf{z}_{C_O,out} \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{Q}_1 \\ \vdots \\ \mathbf{Q}_{C_O} \end{bmatrix}}_{:= \mathbf{Q}} \mathbf{z}_{in} \quad (2)$$

Proposition 3.1. The spectral norm of \mathbf{Q} induced in (2) can be bounded by

$$\|\mathbf{Q}\| \leq \sqrt{\sum_{k=1}^{C_O} \|\mathbf{Q}_k\|^2}$$

and so the ConvNorm actually reduces $\|\mathbf{Q}\|$, improving the Lipschitzness in each layer.

Proof.

$$\sigma_1^2(\mathbf{Q}) = \lambda_1(\mathbf{Q}^\top \mathbf{Q}) = \lambda_1 \left(\sum_{i=1}^{C_O} \mathbf{Q}_i^\top \mathbf{Q}_i \right) \leq \sum_{i=1}^{C_O} \lambda_1(\mathbf{Q}_i^\top \mathbf{Q}_i) = \sum_{i=1}^{C_O} \sigma_1^2(\mathbf{Q}_i)$$

Taking the square root of both sides provides the desired result. \square

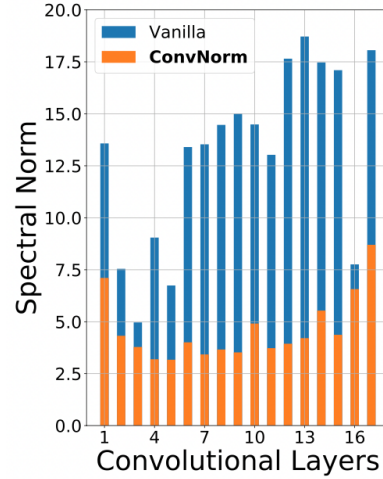


Figure 1: **Spectral norm for each layer on ResNet18 [7].** ConvNorm reduces the spectral norm

4 Experiments and Results

The ConvNorm repository [7] was imported and modified to implement the WideResNet architecture [10]. A new script was created based on the existing model for ResNet and adapted to implement the 'basic-wide' block [10]. Two normalization techniques were focused on, namely BN [4] and the previously discussed ConvNorm. Similar to [7], the training procedure described in [11] was used for robust training. Experiments were performed to investigate the effectiveness of the normalization techniques on WideResNet16 and explore how they perform against the width factor k .

For every experiment, the models were trained on the CIFAR-10 dataset [12] for 40 epochs using cross-entropy loss and AdamW [12] for optimization with an initial learning rate of 0.001. 10% of the training set is used for validation and the validation set is treated as a held-out test set. PGD and FGSM attacks are performed with total $\epsilon = 8/255$ where PGD is implemented with 10 steps of size $\alpha = 2/255$. Computations were executed with a NVIDIA Tesla P100 PCIe 16 GB GPU.

4.1 Adjusting to Computational Limitations

WideResNet is much more efficient than ResNet due to the parallel nature of GPU computations [10]. We utilize this computational discount and chose this architecture to increase our capacity to explore more networks. With limited computational resources, AdamW was chosen as an optimizer with the aim of reducing training time [13].

4.2 Results

4.2.1 Natural Accuracy

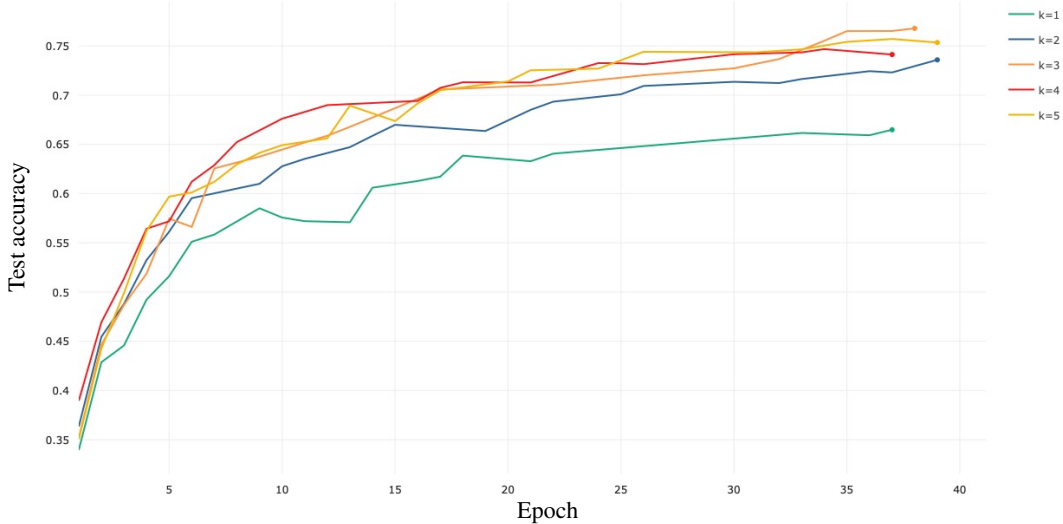


Figure 2: **Test accuracy of WideResNet-16 using various widening factors k .** Increasing the width of the model increases the natural accuracy. (Note: some of the data for the later epochs is missing due to upload rate limits of comet.ml)

4.2.2 Robust Accuracy

5 Discussion and Conclusion

Through the implementation of WideResNet, the relationship between ConvNorm and network width was explored. The empirical data suggests that ConvNorm improves natural accuracy but there is a significant decrease in robustness. While there may be a trade-off between natural accuracy and robustness [14] [15], BN also increases natural accuracy but doesn't exhibit the same proportional drop in robustness. While these initial results may indicate that ConvNorm does not improve

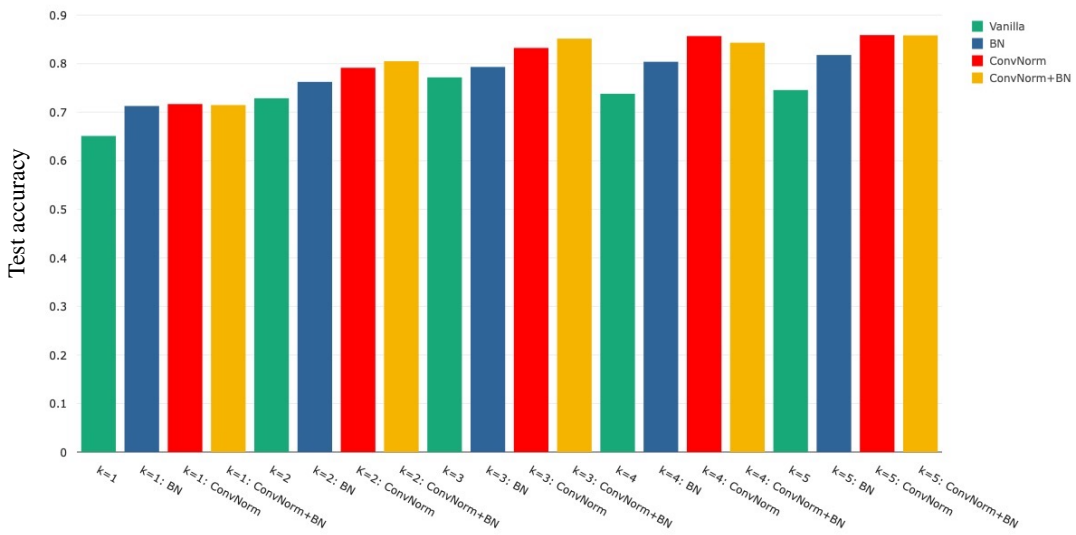


Figure 3: **Top-1 test accuracy of WideResNet-16- k with various normalization techniques.** ConvNorm and BN both improve the natural accuracy for increased widening factors.

Table 1: **Comparison of normalization methods under different adversarial attacks.**

| Top-N | Attack | k | Vanilla | BN | ConvNorm | BN+ConvNorm |
|-------------|--------|-----|---------|-------|----------|-------------|
| Precision@1 | PGD | 1 | 29.64 | 37.05 | 38.00 | 37.03 |
| | | 2 | 35.23 | 38.09 | 41.72 | 1.17 |
| | | 3 | 32.95 | 41.18 | 0.26 | 3.19 |
| | | 4 | 34.33 | 40.67 | 8.03 | 1.31 |
| | | 5 | 33.82 | 40.36 | 2.81 | 2.44 |
| | FGSM | 1 | 36.27 | 41.92 | 42.50 | 41.66 |
| | | 2 | 47.30 | 44.20 | 47.41 | 75.73 |
| | | 3 | 55.17 | 47.42 | 71.27 | 81.98 |
| | | 4 | 43.21 | 48.80 | 66.40 | 83.58 |
| | | 5 | 43.31 | 47.36 | 87.79 | 74.45 |
| Precision@5 | PGD | 1 | 85.49 | 90.12 | 90.65 | 90.80 |
| | | 2 | 88.93 | 91.47 | 91.93 | 63.53 |
| | | 3 | 89.43 | 92.19 | 46.56 | 68.35 |
| | | 4 | 87.55 | 92.29 | 78.07 | 56.49 |
| | | 5 | 88.78 | 92.80 | 68.16 | 62.81 |
| | FGSM | 1 | 88.32 | 91.34 | 91.75 | 91.79 |
| | | 2 | 92.57 | 92.45 | 93.28 | 98.57 |
| | | 3 | 94.98 | 93.42 | 98.05 | 99.04 |
| | | 4 | 90.81 | 93.78 | 97.80 | 99.29 |
| | | 5 | 91.75 | 93.89 | 99.44 | 98.35 |

The model trained with ConvNorm exhibits an increase in robustness against FGSM attacks as the widen factor k increases. However, it performs badly against the stronger PGD attacks for larger widths. The robustness of the vanilla model fluctuates slightly for both attacks as the width increases. BN almost always results in an increase of robustness in comparison to the vanilla model.

robustness of WideResNet, further experiments are worthwhile before conclusion. Possible avenues of further investigation include training using: CIFAR-100, more epochs, SGD with a learning rate schedule.

References

- [1] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [2] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- [3] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [4] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [5] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [6] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- [7] Sheng Liu, Xiao Li, Yuexiang Zhai, Chong You, Zhihui Zhu, Carlos Fernandez-Granda, and Qing Qu. Convolutional normalization: Improving deep convolutional network robustness and training. *arXiv preprint arXiv:2103.00673*, 2021.
- [8] Robert M Gray. Toeplitz and circulant matrices: A review. 2006.
- [9] Qing Qu, Xiao Li, and Zhihui Zhu. A nonconvex approach for exact and efficient multichannel sparse blind deconvolution. *arXiv preprint arXiv:1908.10776*, 2019.
- [10] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- [11] A Shafahi, M Najibi, MA Ghiasi, Z Xu, J Dickerson, C Studer, LS Davis, G Taylor, and T Goldstein. Adversarial training for free!” in advances in neural information processing systems. 2019.
- [12] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [13] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [14] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- [15] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning*, pages 7472–7482. PMLR, 2019.