





# Efficient Training of Low-Curvature Neural Networks

Suraj Srinivas<sup>\*</sup>, Kyle Matoba<sup>\*</sup>, Hima Lakkaraju, François Fleuret Harvard University, Idiap Research Institute & EPFL, Harvard University, University of Geneva

\*Equal contribution

## How to Penalize Curvature

Given function  $f : \mathbb{R}^n \to \mathbb{R}$ , define its curvature at **x** by

 $\mathcal{C}_f(x) = \frac{\|\nabla_x^2 f(\mathbf{x})\|_2}{\|\nabla_x f(\mathbf{x})\|_2 + \epsilon} \xrightarrow{\to} \text{spectral norm of Hessian} \xrightarrow{\to} \text{gradient norm}$ 

There are two broad strategies for curvature penalization

- Local Penalization: minimizes point-wise curvature, e.g.: [Dezfooli et al., 2019]
- 2. Global Penalization: minimizes upper-bound on curvature across data points, e.g.: [Dombrowski et al., 2022]

| Local Penalization          | Global Penalization                                |
|-----------------------------|--|
| + (Nearly) exact penalties  | - Loose bounds can lead to<br>inexact penalization |
| - Expensive                 | + Efficient!                                       |
| - Penalization only at data | + Penalization everywhere                          |

In this work, we choose to do global penalization, but find that using approximate local penalization also helps!

# Upper Bounding Normalized Curvature

Given a function  $f = f_L \circ f_{L-1} \circ \ldots \circ f_1$  with  $f_i : \mathbb{R}^{n_{i-1}} \to \mathbb{R}^{n_i}$ , the curvature  $C_f$  can be bounded by the sum of curvatures of individual layers  $\mathcal{C}_{f_i}(\mathbf{x})$ , i.e.,

$$\max_{\mathbf{x}} \mathcal{C}_f(\mathbf{x}) \leq \sum_{i=1}^{L} \max_{\mathbf{x}} \mathcal{C}_{f_i}(\mathbf{x}) \prod_{j=1}^{i} ||W_j||_2.$$

To minimize the bound, we

- Penalize curvature of activation functions,  $\max_{\mathbf{X}} \mathcal{C}_{f_i}(\mathbf{X})$
- 2. Penalize Lipschitz constants of linear layers,  $||W_i||_2$

Ingredients of LCNNs

(1) Centered  $\beta$ -softplus is defined as

$$s_{c}(\mathbf{x};\beta) = \frac{1}{\beta} \log \left( \frac{1 + \exp(\beta \mathbf{x})}{2} \right)$$

This has several nice properties:

- . Approximates ReLU:  $s_c(x; \beta \to \infty) = \text{ReLU}(x)$
- 2. Bounded curvature:  $\mathcal{C}_{S_c}(x) \leq \beta$
- 3. Stable at low curvature:  $s_c(x; \beta \to 0) = x/2$
- . Stable upon composition:  $s_c^n(x=0;\beta)=0$

Standard softplus diverges for both #3 and #4!

(2)  $\gamma$ -Batchnorm is defined as

$$\gamma \text{-BN}(x) \leftarrow \underbrace{\min(\gamma, ||\text{BN}(x)||_2)}_{\text{scaling factor} \leq \gamma} \times \frac{\text{BN}(x)}{||\text{BN}(x)||_2}$$

(3) "Real" spectral normalization of linear layers, which is defined in previous work [Ryu et al., 2019].

(4) <u>Penalization</u> of curvature by penalization of  $\beta$  and  $\gamma$ parameters, avoids the need to compute Hessian norms!

#### Advantages of LCNNs

Low curvature models have stable gradients,

$$\frac{\|\nabla_{\mathsf{x}} f(\mathbf{x} + \epsilon) - \nabla_{\mathsf{x}} f(\mathbf{x})\|_{2}}{\|\nabla_{\mathsf{x}} f(\mathbf{x})\|_{2}} \sim \|\epsilon\|_{2} \mathcal{C}_{f}(\mathbf{x})$$

...and low curvature is necessary for robustness!

$$\|f(\mathbf{x}+\epsilon)-f(\mathbf{x})\|_{2}\sim\|\epsilon\|_{2}\|\nabla_{\mathsf{x}}f(\mathbf{x})\|_{2}(1+\|\epsilon\|_{2}\mathcal{C}_{f}(\mathbf{x}))$$





### Evaluating Point-wise Curvature

We present results on Resnet-18 models trained on CIFAR100 dataset.

We find that LCNN-based models reduce point-wise curvature without loss in accuracy, even though they are designed to penalize global curvature.

| Model           | $\mathbb{E}_X \mathcal{C}_f(X)$ | Accuracy (%)                  |
|-----------------|---------------------------------|-------------------------------|
| Standard        | $270.89 \pm _{75.04}$           | $77.42 \pm 0.11$              |
| LCNNs           | $69.50{\scriptstyle~\pm~2.41}$  | $77.30 \pm 0.11$              |
| GradReg         | $89.47 \pm 5.86$                | $77.20 \pm 0.26$              |
| LCNNs + GradReg | $25.30 \pm 0.09$                | $77.29 \scriptstyle~\pm 0.07$ |
| HessianReg [1]  | $116.31 \pm \textbf{4.58}$      | $76.48 \pm 0.07$              |

### Evaluating Adversarial Robustness

We find that LCNN models provide adversarial robustness (to PGD advex and Autoattack) competitive with adversarial training without sacrificing accuracy.

| Model           | Clean acc. $(\%)$ | Adv. acc. (%)   |
|-----------------|-------------------|-----------------|
| Standard        | $77.42$ $\pm$ .10 | $16.11 \pm .21$ |
| LCNN            | $77.16 \pm .07$   | $17.66 \pm .18$ |
| GradReg         | 77.20 ± .26       | $38.09 \pm .47$ |
| LCNNs + GradReg | $77.29 \pm .26$   | $42.70 \pm .77$ |
| Adv. Training   | $76.96 \pm .26$   | $44.98 \pm .57$ |

#### References

[1] Dezfooli, Fawzi, Juesato, and Frossard. "Robustness via curvature regularization, and vice versa", CVPR 2019 [2] Ryu, Liu, Wang, Chen, Wang, and Yin. "Plug-and-play methods provably converge with properly trained denoisers", ICML 2019

[3] Dombrowski, Anders, Muller, and Kessel. "Towards robust explanations for deep neural networks". Pattern Recognition 2022