

A path-norm toolkit for modern networks: consequences, promises and challenges

Antoine Gonon, Nicolas Brisebarre, Elisa Riccietti, Rémi Gribonval

Univ Lyon, ENS de Lyon, UCBL, CNRS, Inria, LIP, F-69342 Lyon



SUMMARY

MODERN CHALLENGES

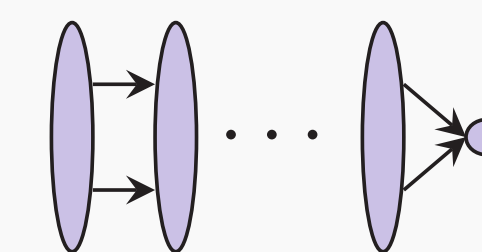
- 1 *Generalization* [1, 2]
- 2 *Robustness* [1]
- 3 Implicit bias [7]
- 4 Identifiability [5, 6]

Existing: "PATH-NORM", A PROMISING TOOL?

Theory: 😊 partial answers to 1-4

Practice: 😊 correlates with *generalization* [3, 4]
😊 easy to compute [4]

... only developed for toy networks



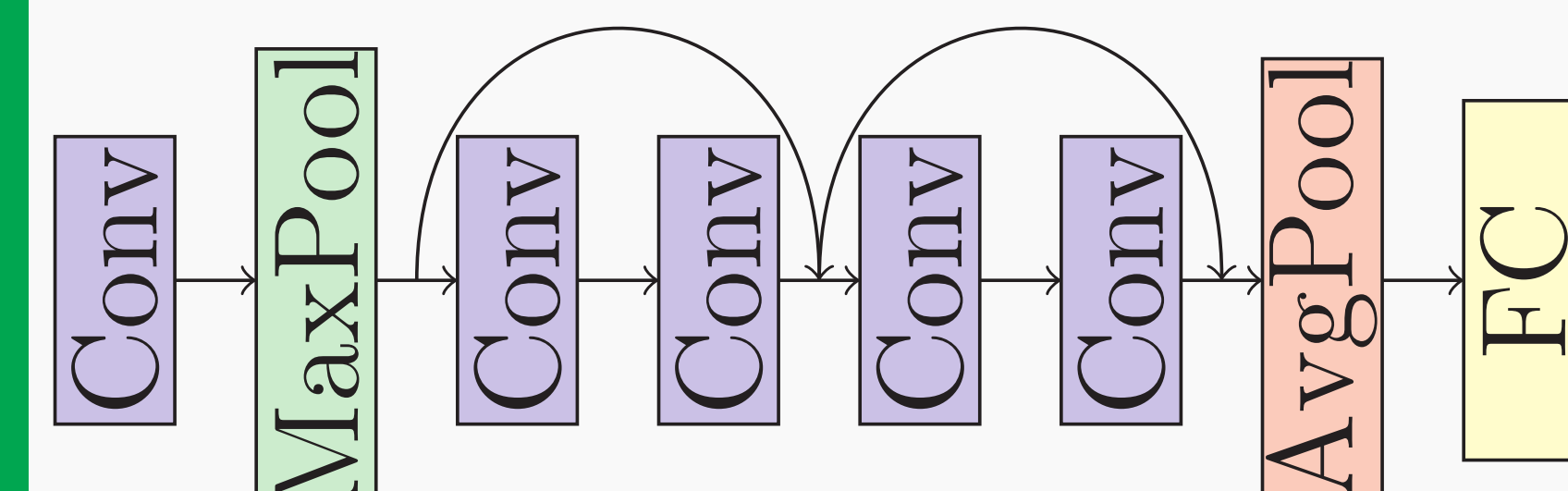
This work: TOOLS FOR MODERN RELU NETWORKS

RESNETS, VGGs, U-NETS, ReLU MOBILENETS, INCEPTION NETS, ALEXNET...

Theory: 😊😊 sharper *generalization* bound

Practice: 😊😊 PyTorch implementation
😊😊 first numerical assessment on ResNets/ImageNet

😊 now for modern networks



"PATH-NORM" PHILOSOPHY

$\theta = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$ network parameters (weight, biases)

Measure $\begin{cases} Generalization \\ Robustness \end{cases} \leq f(\|\theta\|)$

→ replace with

$\Phi(\theta) = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$ path-lifting

$\leq g(\|\Phi(\theta)\|)$ "path-norm" based bound

😊 invariant under neuron-wise rescaling
😊 sharper than $\|\theta\|$

CONTRIBUTION 1: DEFINITION OF PATH-LIFTING AND PATH-ACTIVATIONS FOR MODERN NNS

$\begin{cases} \Phi(\theta) \text{ path-lifting (see paper)} \\ A(\theta, x) \text{ path-activations (binary matrix)} \end{cases}$ s.t. $R_\theta(x) := \text{output}(\theta, x) = \langle \Phi(\theta), A(\theta, x)x \rangle$

Corollary: Measure(*Robustness*)=Lipschitz constant= $\sup_{x \neq x'} \frac{\|R_\theta(x) - R_\theta(x')\|_1}{\|x - x'\|_\infty} \leq \|\Phi(\theta)\|_1 =:$ path-norm (easy to compute)

CONTRIBUTION 2: FIRST PATH-NORM BASED GENERALIZATION BOUND VALID FOR MODERN NETWORKS

Our bound: $C \times \|\Phi(\theta)\|_1$

$C = \frac{\sup_x \|x\|_\infty}{\sqrt{n}} L \sqrt{D \ln(K) + \ln(d_{\text{in}} d_{\text{out}})}$

- 😊 Valid for modern networks
- 😊 Sharper than $\prod_{k=1}^D \|M_k\|_{\infty \rightarrow \infty}$ for layered networks $x \mapsto M_D \text{ReLU}(\dots \text{ReLU}(M_1 x))$
- 😊 Improves on previous generalization bounds

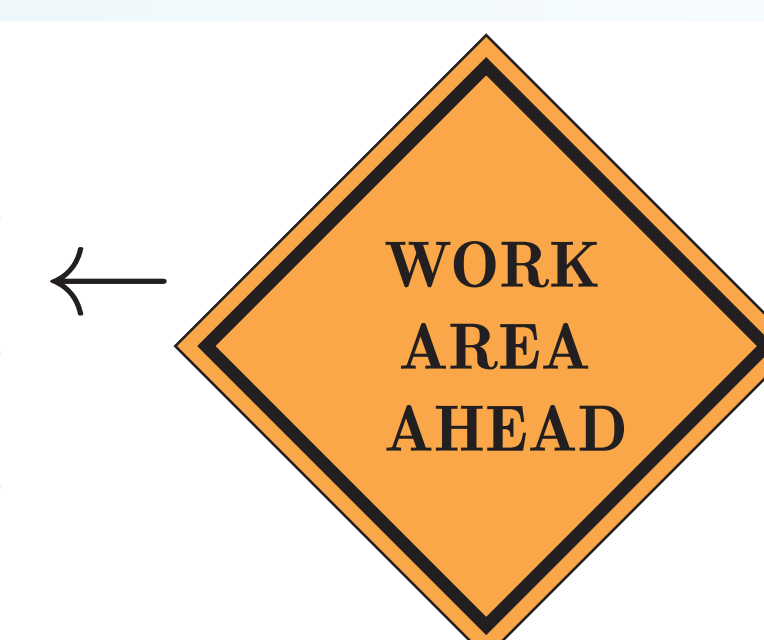
D = depth, K = k -max-pooling kernel size, L = loss Lipschitz constant, n = number of samples, $d_{\text{in}}/d_{\text{out}}$ = input/output dimension

CONTRIBUTION 3: FIRST ASSESSMENT OF THE PROMISES OF PATH-NORM ON MODERN NETWORKS

😊 Assessing the promises of path-norm for the first time on ResNets/ImageNet $C \simeq 0.1$

PyTorch pretrained ResNet18

$\ \Phi(\theta)\ _1$	1.3×10^{30}
$\ \Phi(\theta)\ _2$	2.5×10^2
$\ \Phi(\theta)\ _4$	7.2×10^{-6}



WHAT'S NEXT?

worst-case measure $\|\Phi(\theta)\|_1$



average-case complexity($\Phi(\theta)$)

REFERENCES

- [1] Neyshabur et al., Norm-Based Capacity Control in Neural Networks, COLT (2015)
- [2] Barron et al., Complexity, Statistical Risk, and Metric Entropy of Deep Nets Using Total Path Variation, preprint (2019)
- [3] Jiang et al., Fantastic Generalization Measures and Where to Find Them, ICLR (2020)
- [4] Dziugaite et al., In search of robust measures of generalization, NIPS (2020)
- [5] Bona-Pellissier et al., Local Identifiability of Deep ReLU Neural Networks: the theory, NIPS (2022)
- [6] Stock et al., An embedding of ReLU networks and an analysis of their identifiability, Constr. Approx. (2023)
- [7] Marcotte et al., Abide by the Law and Follow the Flow: Conservation Laws for Gradient Flows, NIPS (2023)