

The Spectral Bias of Polynomial Neural Networks

Moulik Choraria¹ Leello Daadi² Grigorios Chrysos²
Julien Mairal³ Volkan Cevher²

¹University of Illinois at Urbana-Champaign

²EPFL, Switzerland

³Univ. Grenoble-Alpes, Inria

Motivation

- Neural Networks (NNs) demonstrate a learning bias during training, wherein low-frequency functions are learned faster¹.
- Polynomial Neural Networks (PNNs) with multiplicative layers were recently shown to be effective at image generation and face recognition, where high-frequency information is critical².
- Inspired by these, we conduct a spectral analysis of PNNs, using the Π -Net parametrization.

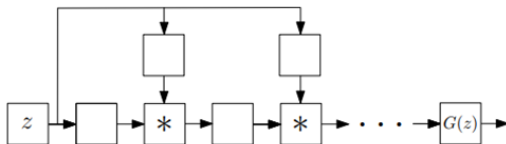


Figure: A Π -Net variant.

¹Nasim Rahaman et al. "On the Spectral Bias of Neural Networks". In: (2019).

²Grigorios G. Chrysos et al. " Π -nets: Deep Polynomial Neural Networks". In: 2020.

Setup (Neural Tangent Kernel Regime)

- The 2-layer feed-forward ReLU (σ) network in infinite width limit ($m \rightarrow \infty$): $f_{\mathbf{W}}(\mathbf{x}) = \sqrt{\frac{2}{m}} \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})$
- The corresponding Π -Net includes a multiplicative layer via the Hadamard product ($*$): $f_{\mathbf{W}}(\mathbf{x}) = \sqrt{\frac{2}{m}} \mathbf{W}_3 [\sigma(\mathbf{W}_1 \mathbf{x}) * \sigma(\mathbf{W}_2 \mathbf{x})]$
- The NTK for the 2-layer feed-forward network is composed of arc-cosine kernels: $\kappa(\mathbf{x}, \mathbf{x}') = 2\langle \mathbf{x}, \mathbf{x}' \rangle \kappa_1(\mathbf{x}, \mathbf{x}') + 2\kappa_2(\mathbf{x}, \mathbf{x}')$.
- The Π -Net NTK takes a product of arc-cosine kernels form: $\kappa_{\pi}(\mathbf{x}, \mathbf{x}') = 2(2\langle \mathbf{x}, \mathbf{x}' \rangle \kappa_1(\mathbf{x}, \mathbf{x}') + \kappa_2(\mathbf{x}, \mathbf{x}')) \kappa_2(\mathbf{x}, \mathbf{x}')$.

Result & Implications

- For the linear operator w.r.t. the 2-layer feed-forward NTK, eigenvalues $(\mu_k)_{k=1}^{\infty}$ for frequency k (even), decay as $\mu_k = \Omega(k^{-d-1})$ when $k \gg d$.³

Theorem (Ours)

Let $\{\mu_{\pi,1}, \mu_{\pi,2}, \dots\}$ denote the eigenvalues of the linear operator $L_{\kappa_{\pi}}$ w.r.t. the Π -kernel κ_{π} . For $k \gg d$ ($k \not\equiv 2 \pmod{4}$), $\mu_{\pi,k} = \Omega(k^{-d/2-2})$.

- Slower decay leads to larger RKHS (better approximation properties)
- Slower decay leads to faster learning for higher frequencies in kernel regression (translating to NNs in NTK regime).⁴

³Alberto Bietti and Julien Mairal. "On the Inductive Bias of Neural Tangent Kernels". In: (2019).

⁴Yuan Cao et al. "Towards Understanding the Spectral Bias of Deep Learning". In: (2020).

Learning Sinusoids

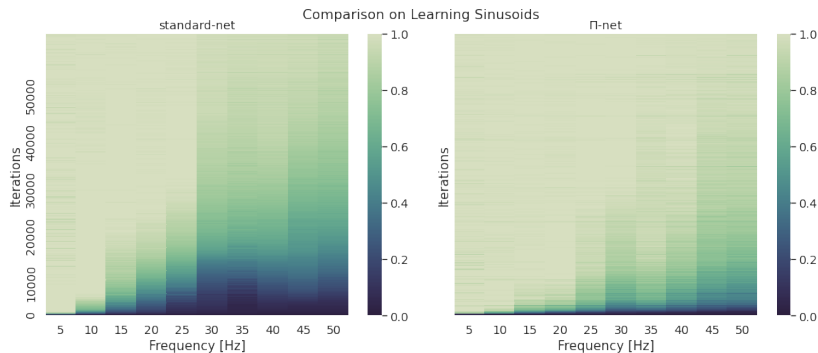


Figure: Comparison between rate of learning sinusoids with fully connected networks (heat map depicts fraction of learning for each frequency over iterations) yields that Π -Nets with multiplicative layers learn higher frequencies much faster than feed-forward networks.

Denoising Images

- We consider the task of denoising an image using a U-Net structure, in the Deep Image Prior⁵ setting to observe this spectral bias in presence of noise.

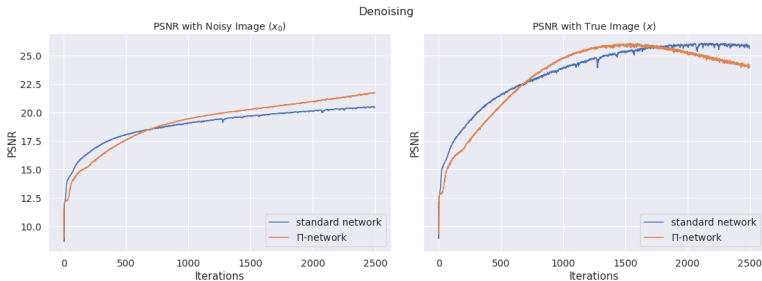


Figure: Denoising comparison shows that Π -Nets reach the peak PSNR w.r.t. the true image much faster and show a lesser impedance towards learning high frequency noise, compared to standard Deep networks.

⁵Victor Lempitsky, Andrea Vedaldi, and Dmitry Ulyanov. “Deep Image Prior”. In: 2018.

Label Noise in Classification

- We consider binary classification in MNIST (with two digits) with fully connected networks in the presence of label noise, as in Rahaman et al.

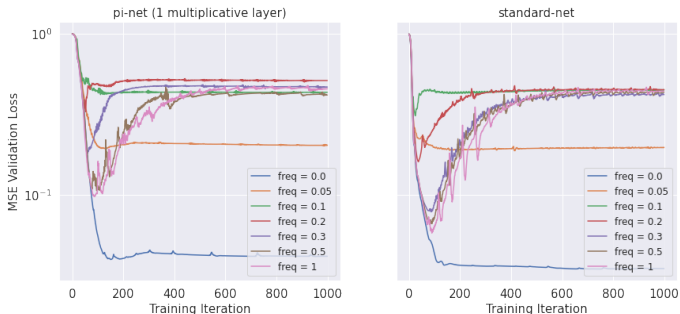


Figure: We observe that Π -Nets are much faster in picking up label noise of higher frequencies, leading to a much smaller “dip”, thus demonstrating a stronger tendency to learn more complex decision boundaries.

- **Extensions to general polynomials:** New proof techniques (beyond NTKs) are needed for formulating more general polynomials.
- **On Generalization in PNNs:** Deviation from the “overparametrized networks learn low-complexity functions” picture of generalization or that simple decision boundaries lead to adversarial robustness.
- **On the Empirical front:** Considering SGD (instead of GD in our experiments) can lead to different results. Additionally, in training of Wasserstein GANs, it would be interesting to relate the moment-matching perspective to spectral bias in polynomial networks.

The Spectral Bias of Polynomial Neural Networks

Contact author: moulikc2@illinois.edu

