Deeper Understanding of Deep Learning: Functional Analysis for Neural Networks



Robert Nowak University of Wisconsin-Madison https://nowak.ece.wisc.edu

Joint work with



Rahul Parhi



Liu Yang

Joe Shenouda

TAMIDS Seminar & Tutorial Series Fall 2023



TEXAS A&M Institute of Data Science

November 6, 2023

Here we are today: AI and Society



Al is one of the most powerful technologies of our time. President Biden has been clear that we must take bold action to harness the benefits and mitigate the risks of AI. The Biden-Harris Administration has acted decisively to protect safety and rights in the age of AI, so that everyone can benefit from the promise of AI.

> Develop standards, tools, and tests to help ensure that AI systems are safe, secure, and trustworthy

Modern Neural Network Architectures



Transformer

The **evolved transformer** <u>D So, Q Le, C Liang</u> - International Conference on Machine ..., 2019



Training Neural Networks

Neural Networks with Rectified Linear Units (ReLUs)





$$f(\boldsymbol{x}) = \sum_{j=1}^{K} \boldsymbol{v}_j (\boldsymbol{w}_j^T \boldsymbol{x})_+$$

T 7

"Training" Neural Networks

Consider the optimization

$$\min_{\boldsymbol{v},\boldsymbol{w}} \frac{1}{2} \sum_{i=1}^{n} \left\| \boldsymbol{y}_{i} - \sum_{j} \boldsymbol{v}_{j} \left(\boldsymbol{w}_{j}^{T} \boldsymbol{x}_{i} \right)_{+} \right\|_{2}^{2} + \underbrace{\operatorname{data fitting term} L(\boldsymbol{v},\boldsymbol{w})}_{\text{data fitting term}} \right\|_{2}^{2}$$

$$\frac{\lambda}{2} \sum_{j} \| \boldsymbol{v}_{j} \|_{2}^{2} + \| \boldsymbol{w}_{j} \|_{2}^{2}$$

regularize sum of squared weights

Gradient descent update for $oldsymbol{v}_j$ (same for $oldsymbol{w}_j)$

$$oldsymbol{v}_j \leftarrow oldsymbol{v}_j - \mu \Big(
abla_{oldsymbol{v}_j} L + \lambda oldsymbol{v}_j \Big) \ \leftarrow oldsymbol{v}_j - \mu
abla_{oldsymbol{v}_j} L - \mu \lambda oldsymbol{v}_j$$

Standard approach when training neural networks (possibly with a different loss)

Neural Balance

because ReLU is piecewise linear

$$v_j(w_j^T x)_+ = \alpha^{-1} v_j (\alpha w_j^T x)_+, \forall \alpha > 0$$
 "Neural Balance"
at the global minimum of weight decay objective $||v_j||_2 = ||w_j||_2$
Proof: The solution to the optimization

$$\min_{\alpha>0} \|\alpha \, \boldsymbol{w}\|_2^2 + \|\alpha^{-1} \boldsymbol{v}\|_2^2$$

is $\alpha = \sqrt{\|\boldsymbol{v}\|_2 / \|\boldsymbol{w}\|_2}$

at the global minimum $\|m{v}_j\|_2 = \|m{w}_j\|_2$ and therefore $rac{1}{2}ig(\|m{v}_j\|_2^2 + \|m{w}_j\|_2^2ig) \ = \ \|m{v}_j\|_2 \|m{w}_j\|_2$

Grandvalet, Y. (1998). Least absolute shrinkage is equivalent to quadratic penalization. In *International Conference on Artificial Neural Networks* (pp. 201-206).

Neyshabur, B., Tomioka, R., & Srebro, N. (2015). In Search of the Real Inductive Bias: On the Role of Implicit Regularization in Deep Learning. In *ICLR (Workshop)*.

Secret Sparsity of Weight Decay

gradient descent with weight decay aims to minimize sum of losses + sum of squared weights

$$\min_{\boldsymbol{v},\boldsymbol{w}} \frac{1}{2} \sum_{i=1}^{n} \left\| \boldsymbol{y}_{i} - \sum_{j} \boldsymbol{v}_{j} \left(\boldsymbol{w}_{j}^{T} \boldsymbol{x}_{i} \right)_{+} \right\|_{2}^{2} + \frac{\lambda}{2} \sum_{j} \left\| \boldsymbol{v}_{j} \right\|_{2}^{2} + \left\| \boldsymbol{w}_{j} \right\|_{2}^{2}$$
squared

same global minima as "path-norm" regularization

$$\min_{\boldsymbol{v}, \boldsymbol{w}} \frac{1}{2} \sum_{i=1}^{n} \| \boldsymbol{y}_{i} - \sum_{j} \boldsymbol{v}_{j} (\boldsymbol{w}_{j}^{T} \boldsymbol{x}_{i})_{+} \|_{2}^{2} + \lambda \sum_{j} \| \boldsymbol{v}_{j} \|_{2} \| \boldsymbol{w}_{j} \|_{2}$$
"path-norm"

same global minima as *constrained* ℓ^1 regularization

$$\min_{\boldsymbol{v},\boldsymbol{w}:\|\boldsymbol{w}_{j}\|_{2}=1} \frac{1}{2} \sum_{i=1}^{n} \|\boldsymbol{y}_{i} - \sum_{j} \boldsymbol{v}_{j} (\boldsymbol{w}_{j}^{T} \boldsymbol{x}_{i})_{+}\|_{2}^{2} + \lambda \sum_{j} \|\boldsymbol{v}_{j}\|_{2}$$

"multi-task" lasso

 ℓ^2 norm

Weight Decay Produces Sparse Networks



weight decay (eventually) produces a very sparse neural network ... after millions of GD steps!

equivalent ℓ^1 path-norm regularization (gradient descent + soft-thresholding) converges 15X faster training a 320 neuron ReLU network to fit 32 data points





Accurate, Robust, and Sparser Networks



architecture: three-layer, fully connected, dataset: MNIST

What Kinds of Functions Do Neural Networks Learn?

Shallow Networks and Neural Balance



Equivalent optimizations:

$$\begin{split} \min_{\boldsymbol{v},\boldsymbol{w}} \frac{1}{2} \sum_{i=1}^{n} \left\| y_{i} - \sum_{j} v_{j} \left(\boldsymbol{w}_{j}^{T} \boldsymbol{x}_{i} \right)_{+} \right\|_{2}^{2} &+ \frac{\lambda}{2} \sum_{j} |v_{j}|_{2}^{2} + \|\boldsymbol{w}_{j}\|_{2}^{2} \\ \min_{\boldsymbol{v},\boldsymbol{w}} \frac{1}{2} \sum_{i=1}^{n} \left\| y_{i} - \sum_{j} v_{j} \left(\boldsymbol{w}_{j}^{T} \boldsymbol{x}_{i} \right)_{+} \right\|_{2}^{2} &+ \lambda \sum_{j} |v_{j}| \|\boldsymbol{w}_{j}\|_{2} \\ \end{split}$$
 "path-norm"

Norm on the Space of Neural Network Functions

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \to \mathbb{R} : f(\boldsymbol{x}) = \sum_{j=1}^K v_j (\boldsymbol{w}_j^T \boldsymbol{x})_+, \ \boldsymbol{w}_j \in \mathbb{R}^d, \ v_j \in \mathbb{R} \right\}$$

The **path-norm** is a norm on \mathcal{F}

$$\|f\| := \sum_{j} |v_j| \| m{w}_j \|_2$$

Completion of \mathcal{F} with respect to $\|\cdot\|$ is the Banach space of all functions generated by "infinite-width" neural networks. Each such function f is expressed as an integral over neurons with respect to a finite measure ν

$$f(\boldsymbol{x}) = \int_{\boldsymbol{w} \in \mathbb{S}^{d-1}} (\boldsymbol{w}^T \boldsymbol{x})_+ d\nu(\boldsymbol{w})$$

Universal approximation bounds for superpositions of a sigmoidal function <u>AR Barron</u> - IEEE Transactions on Information theory, **1993** - ieeexplore.ieee.org

Breaking the curse of dimensionality with convex neural networks F Bach - The Journal of Machine Learning Research, 2017 - dl.acm.org

Path-Norm Related to Derivatives of f



Savarese, P., Evron, I., Soudry, D., & Srebro, N. (2019, June). How do infinite width bounded norm networks look in function space?. In *Conference on Learning Theory* (pp. 2667-2690). PMLR.

Weight Decay $\equiv TV^2(f)$ Regularization

weight decay is gradient of squared ℓ^2 regularization

$$\min_{\boldsymbol{v},\boldsymbol{w},\boldsymbol{b}} \frac{1}{2} \sum_{i=1}^{n} \left(y_i - \sum_j v_j \left(w_j x_i + b_j \right)_+ \right)^2 + \frac{\lambda}{2} \sum_j |v_j|^2 + |w_j|^2$$

equivalent to ℓ^1 regularization

$$\min_{\boldsymbol{v},\boldsymbol{w},\boldsymbol{b}} \frac{1}{2} \sum_{i=1}^{n} \left(y_i - \sum_j v_j \left(w_j x_i + b_j \right)_+ \right)^2 + \lambda \sum_j |v_j w_j|$$

equivalent to $TV^2(f)$ regularization

$$\min_{\boldsymbol{v},\boldsymbol{w},\boldsymbol{b}} \frac{1}{2} \sum_{i=1}^{n} \left(y_i - \sum_j v_j \left(w_j x_i + b_j \right)_+ \right)^2 + \lambda \mathsf{TV}^2(f)$$

remember: just think
$$f(x_i) \qquad \qquad \text{of this as } \int |f''(x)| \, dx$$

 BV^2 is the space of all functions with $\mathsf{TV}^2(f)\,<\,\infty$

Neural Network Representer Theorem (Parhi & N, 2020): For any dataset $\{x_i, y_i\}$ and any lower semicontinous loss function ℓ , there exists a solution to

$$\min_{f \in \mathsf{BV}^2} \sum_{i=1}^n \ell(y_i, f(x_i)) + \lambda \mathsf{TV}^2(f)$$

that is a sparse single hidden-layer ReLU neural network

$$f(x) = \sum_{k=1}^{K} v_k \underbrace{(w_k x - b_k)_+}_{\text{ReLU neurons}} + \underbrace{w_0 x + b_0}_{\text{skip connection}} \underbrace{K < n}_{\text{skip connection}}$$

- properties of differential operator $D^2 := \frac{d^2}{dx^2}$: Green's functions are ReLUs and null space = {linear functions}
- each $f \in BV^2$ can be expressed in terms of finite measure ν i.e., $D^2 f = \nu = \nu^+ \nu^-$
- minimizing total variation $\|\nu\|_{\mathcal{M}} = \int d\nu^+ + \int d\nu^-$ subject to linear constraints is the *measure recovery problem* \Rightarrow solution is sum of Dirac impulses Zuhovickii '48
- pseudoinverse (integrating twice) $\left(D^2\right)^+\delta$ = ReLU activation function

Non-Uniqueness of Solution



there are infinite number of neural network interpolators that minimize TV^2

Extension to Multivariate Case: Radon Transform



path-norm $(f) = \sum |v_k| \| \boldsymbol{w}_k \|_2 = \text{total variation of } \Lambda^{d-1} \mathscr{R} \Delta f$

Radon-domain second-order total variation

Ongie, G., Willett, R., Soudry, D., & Srebro, N. (2020). A Function Space View of Bounded Norm Infinite Width ReLU Nets: The Multivariate Case. In *International Conference on Learning Representations*.

Multidimensional ReLU Neurons

ReLU function: $\phi(w_1x_1 + w_2x_2 - b_0)$

Laplacian of ReLU: $\Delta \phi = \frac{\partial^2 \phi}{\partial x_1^2} + \frac{\partial^2 \phi}{\partial x_2^2}$

filtered Radon transform:

$$\Lambda \mathscr{R} \Delta \phi = \| \boldsymbol{w} \|_2 \, \delta(b - b_0) \, \delta(\theta - \tan^{-1} \frac{w_1}{w_2})$$



indicates slope, orientation and offset of ridge

The Banach Space $\mathscr{R}BV^2$

Radon domain TV²: \mathscr{R} TV² $(f) := {total variation of measure } \Lambda^{d-1} \mathscr{R} \Delta f$

$$\Delta^{d-1}\mathscr{R} = \text{filtered Radon transform}$$
$$\Delta = \sum_{k=1}^{d} \frac{\partial^2}{\partial x_k^2} = \text{Laplacian operator}$$

measures "sparsity" (L^1 -like norm) of second derivatives along each direction in \mathbb{R}^d

Second-Order Bounded Variation Space in Radon Domain:

$$\mathscr{R}\mathsf{BV}^2(\mathbb{R}^d) := \{f: \mathbb{R}^d \to \mathbb{R}, \mathscr{R}\mathsf{TV}^2(f) < \infty\}$$

Banach space representer theorems for neural networks and ridge splines <u>R Parhi</u>, <u>RD Nowak</u> - The Journal of Machine Learning Research, 2021 - dl.acm.org **Neural Network Representer Theorem (Parhi & N, 2020)**: For any dataset $\{x_i, y_i\}$ and any lower semicontinous loss function ℓ , there exists a solution to

$$\min_{f \in \mathscr{R}\mathsf{B}\mathsf{V}^2} \sum_{i=1}^n \ell(y_i, f(\boldsymbol{x}_i)) + \lambda \mathscr{R}\mathsf{T}\mathsf{V}^2(f)$$

with a representation in the form of a single hidden-layer neural network

$$f(\boldsymbol{x}) = \sum_{k=1}^{K} v_k (\boldsymbol{w}_k^T \boldsymbol{x} - b_k)_+ + (\boldsymbol{w}_0^T \boldsymbol{x} + b_0) (K < n)$$

ReLU neurons skip connection



implication: nonparametric problem reduces to training finite-dimensional neural network with weight decay

Neural Networks Adapt to Directional Smoothness



linear methods cannot adapt to directional smoothness in function/data

Minimax Optimality of ReLU Networks

Let $\Omega \subset \mathbb{R}^d$ be a bounded domain. A function $f \in \mathscr{R}BV^2(\Omega)$ can be approximated by a finite-width m ReLU neuron network f_m satisfying

$$||f - f_m||_{L^2}^2 = O(m^{-(1+3/d)}) = O(m^{-1})$$

Estimating f in this class from n nicely distributed samples on $\Omega \subset \mathbb{R}^d$

$$\mathbb{E}\|f - f_m\|_{L^2(\Omega)}^2 = \widetilde{O}\left(m^{-(1+3/d)} + \frac{m}{n}\right) = \widetilde{O}\left(n^{-\frac{d+3}{2d+3}}\right) = O\left(n^{-1/2}\right)$$

$$\overset{\text{approx.}}{\overset{\text{error}}{\overset{\text{or curse}}{\overset{\text{or curse}}{\overset{\overset{o}}{\overset{\overset{o}}{\overset{\overset{o}}{\overset{\overset{o}}{\overset{o}}{\overset{\overset{o}}}{\overset{\overset{o}}{\overset{\overset{o}}{\overset{\overset{o}}}{\overset{\overset{o}}{\overset{\overset{o}}{\overset{\overset{o}}}{\overset{\overset{o}}{\overset{\overset{o}}{\overset{\overset{o}}}{\overset{\overset{o}}{\overset{\overset{o}}{\overset{\overset{o}}}{\overset{\overset{o}}}{\overset{\overset{o}}{\overset{\overset{o}}}{\overset{\overset{o}}{\overset{\overset{o}}}{\overset{\overset{o}}{\overset{\overset{o}}}{\overset{\overset{o}}}{\overset{\overset{o}}{\overset{\overset{o}}{\overset{\overset{o}}}{\overset{\overset{o}}{\overset{\overset{o}}{\overset{\overset{o}}}{\overset{\overset{o}}}$$

this is the minimax rate for $\mathscr{R}\mathsf{BV}^2(\Omega)$

In contrast, for any linear method (thin-plate spline, RKHS, neural tangent kernel) the *linear* minimax rate is $n^{-\frac{3}{d+3}}$

Parhi & N, '21

Data Fitting and Extrapolation



neural networks learn and extrapolate very differently than classical multivariate estimation techniques and kernel methods in general

What Kinds of Functions Do Deep Neural Networks Learn?

Vector-Valued Networks



equivalent optimizations:

$$\min_{\boldsymbol{v},\boldsymbol{w}} \frac{1}{2} \sum_{i=1}^{n} \left\| \boldsymbol{y}_{i} - \sum_{j} \boldsymbol{v}_{j} \left(\boldsymbol{w}_{j}^{T} \boldsymbol{x}_{i} \right)_{+} \right\|_{2}^{2} + \frac{\lambda}{2} \sum_{j} \left\| \boldsymbol{v}_{j} \right\|_{2}^{2} + \left\| \boldsymbol{w}_{j} \right\|_{2}^{2}$$
$$\min_{\boldsymbol{v},\boldsymbol{w}: \| \boldsymbol{w}_{j} \|_{2} = 1} \frac{1}{2} \sum_{i=1}^{n} \left\| \boldsymbol{y}_{i} - \sum_{j} \boldsymbol{v}_{j} \left(\boldsymbol{w}_{j}^{T} \boldsymbol{x}_{i} \right)_{+} \right\|_{2}^{2} + \lambda \sum_{j} \| \boldsymbol{v}_{j} \|_{2}$$

natural norm for vector-valued neural network functions

Vector-Valued Function Spaces

$$\min_{\boldsymbol{v}, \boldsymbol{w}: \|\boldsymbol{w}_j\|_2 = 1} \frac{1}{2} \sum_{i=1}^n \|\boldsymbol{y}_i - \sum_j \boldsymbol{v}_j (\boldsymbol{w}_j^T \boldsymbol{x}_i)_+\|_2^2 + \lambda \sum_j \|\boldsymbol{v}_j\|_2$$

continuum width networks:

$$f(\boldsymbol{x}) = \begin{bmatrix} f_1(\boldsymbol{x}) \\ f_2(\boldsymbol{x}) \\ \vdots \\ f_d(\boldsymbol{x}) \end{bmatrix} = \int_{\boldsymbol{w} \in \mathbb{S}^{d-1}} (\boldsymbol{w}^T \boldsymbol{x})_+ d\boldsymbol{\nu}(\boldsymbol{w})$$

finite vector-valued measure

$$\begin{split} \|f\| &:= \|\boldsymbol{\nu}\| = \text{total variation of } \ell^2 \text{-norm of } \boldsymbol{\nu} \\ &= \sum_j \|\boldsymbol{v}_j\|_2 \text{ for finite width networks} \\ & \text{weight decay encourages a small norm} \end{split}$$

Vector-Valued Variation Spaces and Bounds on Neural Network Widths J. Shenouda, R. Parhi, K. Lee, RN (2023)

Deep Representer Theorem

Let $\mathscr{R}BV_{deep}^2$ denote space of functions generated by taking compositions of functions from the vector-valued variation space

Multilayer Representer Theorem (Parhi & N 2021, Shenouda et al 2023): For any dataset $\{x_i, y_i\}$ and any lower semicontinous loss function ℓ , there exists a solution to

$$\min_{f \in \mathscr{R}\mathsf{B}\mathsf{V}^2_{\mathsf{deep}}} \sum_{i=1}^n \ell(y_i, f(\boldsymbol{x}_i)) + \lambda \mathscr{R}\mathsf{T}\mathsf{V}^2_{\mathsf{deep}}(f)$$

with a representation in the form of a sparse multilayer ReLU neural network with $< nd_l$ neurons in *l*-th layer

implication: nonparametric problem reduces to training multi-layer neural network with weight decay

The Effect of Weight Decay

 $\|f\| := \sum \|v_j\|_2$ for finite width networks



weight decay encourages functions with strong variations in only a few directions (sparse weights)

weight decay encourages output functions to "share" neurons

10 🕄

Another Consequence: Tight Bounds on Widths

consider one ReLU layer within a deep network



Layer Width Theorem : Let Φ and Ψ be matrices composed of the post-activation features and the outputs of a ReLU layer in a deep net, trained to minimize the weight decay objective. Then there exists a representation of this layer with K neurons where K satisfies

 $K \leq \mathsf{rank}(\mathbf{\Phi}) \times \mathsf{rank}(\mathbf{\Psi}) \leq n \, d$

can be a significant improvement on best prior bound of n^2 (Jacot et al '22) relevant because in practice the ranks often become small at deeper layers

Sparsifying Trained Networks

VGG-19 trained to minimize cross-entropy loss + weight decay on CIFAR10



Theory: $K \leq {\rm rank}({\bf \Phi}) \times {\rm rank}({\bf \Psi}) \approx 10 \times 10$ neurons are sufficient

	V	$\widehat{\mathbf{V}}$	
Active Neurons	512	47	
Test Accuracy	93.92%	93.88%	
Train Loss	0.0104	0.0112	

multi-task lasso finds an equivalent solution (wrt training data) with just 47 neurons!

Summary

ReLU neural networks are optimal solutions to data-fitting problems in new function spaces:

- Radon-domain bounded variation space
- Banach space, not Hilbert space
- immune to curse-of-dimensionality
- solutions are sparse / narrow
- adaptive to spatial and directional varying smoothness

Weight Decay is equivalent to "multi-task lasso" regularization and thus promotes solutions that have strong variation in limited directions and encourages outputs to share neurons

Thanks!

This work was supported by the following agencies







