

On BV Spaces, Splines, and Neural Networks

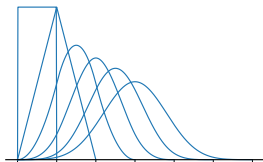
Rahul Parhi

Department of Electrical and Computer Engineering
University of Wisconsin–Madison

(joint work with Robert Nowak)

Analysis Seminar

November 16th, 2021



Outline

- 1 Introduction
 - Reconstructing functions from measurements
 - Variational methods for function reconstruction
 - Splines and variational methods in BV spaces
- 2 Neural networks and variational methods
 - New BV spaces related to neural networks
- 3 What are these new spaces?

UW–Madison is (was?) the mecca of splines

- **Isaac Schoenberg** (1903–1990) invented the spline in the 1940s.
⇒ Was at UW–Madison from 1966–1990.
- **Carl de Boor** wrote many influential papers and books about approximation theory and numerical algorithms with splines.
⇒ Was at UW–Madison from 1972–2003.
- **Grace Wahba** wrote many influential papers about smoothing noisy data with splines, making splines popular in statistics.
⇒ Was at UW–Madison from 1967–2018.
- **Amos Ron** has written many influential papers about approximation theory with splines.
⇒ Has been at UW–Madison since 1988.

Remark

People seem to have forgotten about splines...

Reconstructing functions from measurements

- A fundamental problem in science and engineering is to **reconstruct** a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ from **measurements**.

$$y_n = \langle h_n, f \rangle + \varepsilon_n, \quad n = 1, \dots, N$$

- $H\{f\} = (\langle h_1, f \rangle, \dots, \langle h_N, f \rangle) \in \mathbb{R}^N$ symbolizes the **linear** measurement process.
- $\varepsilon = (\varepsilon_1, \dots, \varepsilon_N) \in \mathbb{R}^N$ are perturbation or noise terms, typically zero-mean random variables.
- $\mathbf{y} = (y_1, \dots, y_N) \in \mathbb{R}^N$ denotes the (possibly noisy) data.

Remark

This is an **ill-posed inverse problem**.

How do we solve this problem?

- Solving this problem requires choosing/designing a model for your functions.
 - ⇒ Assume that $f \in \mathcal{X}'$, where \mathcal{X}' is some Banach space.
 - ⇒ It will be useful to think of \mathcal{X}' as a dual space.
- To remedy the ill-posed nature of the reconstruction problem, a **minimum-energy** requirement is often imposed to **regularize** the solution

$$\min_{f \in \mathcal{X}'} \|f\|_{\mathcal{X}'} \quad \text{s.t.} \quad \langle h_n, f \rangle = y_n, \quad n = 1, \dots, N.$$

or, if the data is noisy,

$$\min_{f \in \mathcal{X}'} \sum_{n=1}^N |y_n - \langle h_n, f \rangle|^2 + \lambda \|f\|_{\mathcal{X}'}^p,$$

where $\lambda > 0$ and $1 \leq p < \infty$.

⇒ $\|\cdot\|_{\mathcal{X}'}$ is a norm or seminorm that defines \mathcal{X}' and is often referred to as a **regularizer**.

Representer theorems

$$\min_{f \in \mathcal{X}'} \|f\|_{\mathcal{X}'} \quad \text{s.t.} \quad \langle h_n, f \rangle = y_n, \quad n = 1, \dots, N.$$

$$\min_{f \in \mathcal{X}'} \sum_{n=1}^N |y_n - \langle h_n, f \rangle|^2 + \lambda \|f\|_{\mathcal{X}'}^p,$$

- Representer theorems provide a **parametric representation**—ideally, a finite linear expansion in terms of some “basis” functions or **atoms**—that span the solution set.
⇒ Representer theorems provide a way to recast **infinite-dimensional** optimization problems as **finite-dimensional** optimization problems, providing a first step to designing a **numerical method** to solve the infinite-dimensional problem.
- This variational formulation for function reconstruction captures many classical methods such as **smoothing splines**, **locally adaptive splines**, and **wavelet methods**.

(Cubic) smoothing splines

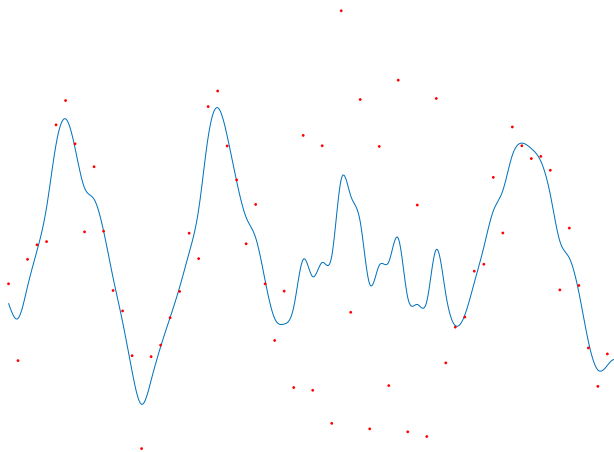
$$\min_{f \in \dot{H}^2(\mathbb{R})} \sum_{n=1}^N |y_n - f(x_n)|^2 + \lambda \|D^2 f\|_{L^2}^2$$

- Measurements are point evaluations $h_n = \delta(\cdot - x_n)$ for some $\{x_n\}_{n=1}^N \subset \mathbb{R}$.
- The solution to this problem is **unique** and a **cubic spline**

$$x \mapsto \sum_{n=1}^N v_k (x - x_n)_+^3 + c_1 x + c_0$$

- The number of **knots** is equal to the number of data.
- The **atoms** of the solution are $(\cdot - x_n)_+^3$.

(Cubic) smoothing splines



(Linear) locally adaptive splines

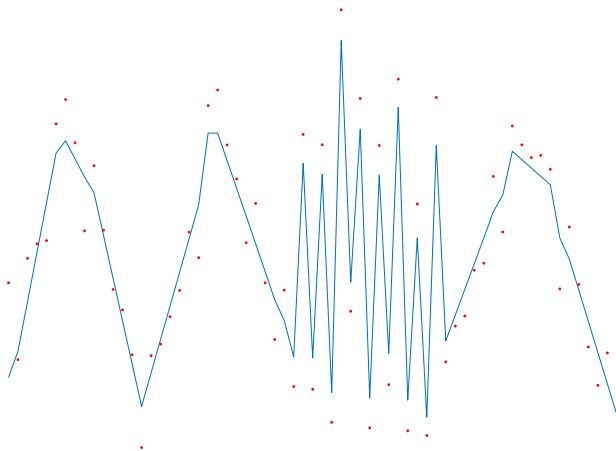
$$\min_{f: D} \sum_{n=1}^N |y_n - f(x_n)|^2 + \lambda \text{TV}(Df)$$

- Measurements are point evaluations $h_n = \delta(\cdot - x_n)$ for some $\{x_n\}_{n=1}^N \subset \mathbb{R}$.
- There **exists** a solution to this problem that is a **linear spline**

$$x \mapsto \sum_{k=1}^K v_k (x - t_k)_+ + c_1 x + c_0$$

- The number of **knots** is $K \leq N$.
- The **atoms** of the solution are $(\cdot - t_k)_+$.

(Linear) locally adaptive splines



Wavelet methods

- Wavelet atoms are translates and dilates of a **mother wavelet** function ψ

$$\left\{ \psi_{j,k}(x) = 2^{j/2} \psi(2^j x - k) \right\}_{j,k \in \mathbb{Z}}$$



- The Daubechies wavelets provide orthobases for $L^2(\mathbb{R})$.

- $\|f\|_{B_{p,q}^s} \asymp \left(\sum_j 2^{jsq} \|\Delta_j f\|_{L^p}^q \right)^{1/q}$

- Wavelet characterizations of $B_{p,q}^s(\mathbb{R})$ essentially replace Δ_j with

$$P_j f = \sum_{k \in \mathbb{Z}} \langle f, \psi_{j,k} \rangle \psi_{j,k}$$

(Db3) wavelet methods

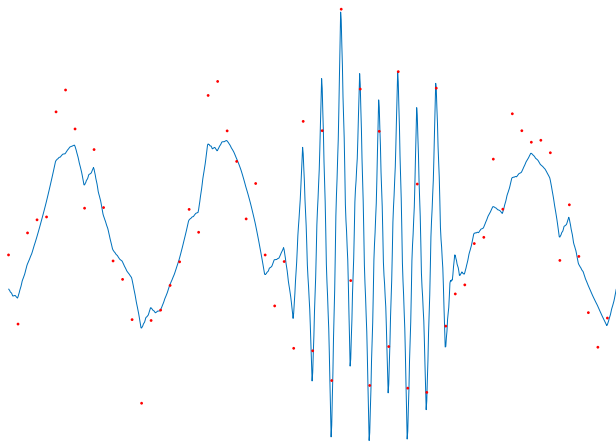
$$\min_{f \in B_{1,1}^2(\mathbb{R})} \sum_{n=1}^N |y_n - f(x_n)|^2 + \lambda \|f\|_{B_{1,1}^2}$$

- $\|\cdot\|_{B_{1,1}^2}$ is the db3 wavelet dependent norm.
- Measurements are point evaluations $h_n = \delta(\cdot - x_n)$ for some $\{x_n\}_{n=1}^N \subset \mathbb{R}$.
- There **exists** a solution to this problem that can be written in terms of a finite number of wavelet functions

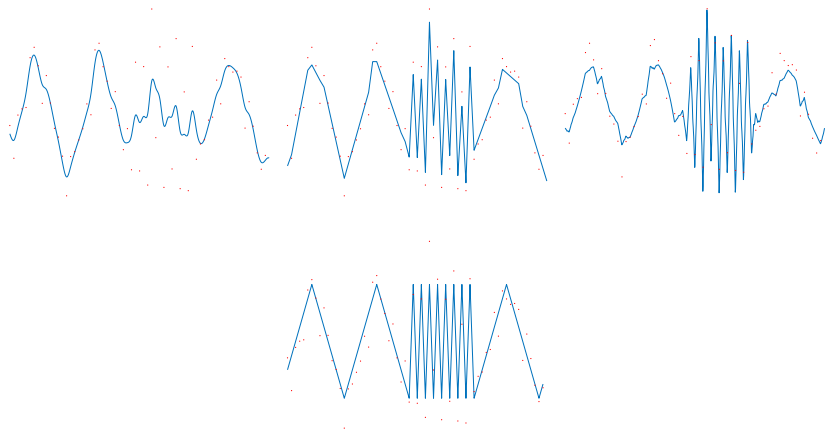
$$x \mapsto \sum_{(j,k) \in \mathcal{I}} c_{j,k} \psi_{j,k}(x)$$

- $|\mathcal{I}|$ depends on λ .
- The **atoms** of the solution are $\psi_{j,k}$.

(Db3) wavelet methods



Comparing these methods



- The data-generating function is not in $\dot{H}^2(\mathbb{R})$!
⇒ Choosing the right function space is **crucial** in getting accurate reconstructions.

Variational formulation for function reconstruction

- **Neural networks** are outperforming and replacing classical methods in many reconstruction tasks.
- Unlike classical methods, neural networks are not well-understood mathematically.
 - ⇒ Classical methods assume regularity of the underlying function and design a procedure to optimally reconstruct functions with that regularity.
 - ⇒ Neural network methods do not aim assume any kind of regularity a priori, but outperform classical methods in practice.
- It turns out that this variational formulation also captures neural networks.

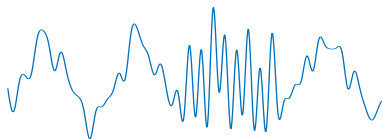
(P. and Nowak 2020, 2021a,b,c)

$$\min_{f \in \mathcal{X}'} \|f\|_{\mathcal{X}'} \quad \text{s.t.} \quad \langle h_n, f \rangle = y_n, \quad n = 1, \dots, N.$$

- ⇒ \mathcal{X}' is a **new**, not a classically studied function space.
- ⇒ \mathcal{X}' is a new kind of BV space that captures the regularity properties that neural networks are intrinsically optimal for.

Spatial inhomogeneity and sparsity

- Many real-world objects (signals, images, functions, etc.) are **spatially inhomogeneous** or exhibit **sparsity**.
 - ⇒ Spatial inhomogeneity is modeled via Besov spaces with $p < 2$.
 - ⇒ Sparsity is modeled via Besov spaces with $p = 1$ or with BV-type spaces.



- e.g., BV is a common model for natural images.

$$B_{1,1}^1(\Omega) \subset BV(\Omega) \subset B_{1,\infty}^1(\Omega)$$

- These kinds of spaces are interesting from an analysis perspective since they are **non-reflexive**.
- The spaces related to neural networks are non-reflexive.
 - ⇒ Before understanding neural networks, we first need to understand (locally adaptive) splines.

What are splines?

- A **polynomial spline** of degree $m - 1$ is a piecewise polynomial function that is C^{m-2} .
- The locations of the discontinuities of the $(m - 1)$ th derivative are referred to as the **knots** of the spline.
- A function f is a polynomial spline of degree $m - 1$ if

$$D^m f = \sum_{k=1}^K v_k \delta(\cdot - t_k),$$



where m is the **order** of the spline and $\{t_k\}_{k=1}^K$ are the **knots**.

- Therefore, every spline can be written as

$$x \mapsto \sum_{k=1}^K v_k \rho_m(x - t_k) + \sum_{n=0}^{m-1} c_n x^n,$$

where $\rho_m(x) = x_+^{m-1}/(m-1)!$ is a **Green's function** of D^m .
 $\implies \rho_m$ is the m th-order **truncated power function**.

Notation

- For $f : \mathbb{R} \rightarrow \mathbb{R}$, let $\text{TV}(f) = \sup_{x_1 < \dots < x_N} \sum_{n=0}^N |f(x_n) - f(x_{n-1})|$.
- $\mathcal{M}(\mathbb{R}) \subset \mathcal{S}'(\mathbb{R})$ be the space of finite Radon measures on \mathbb{R} .

$$\|u\|_{\mathcal{M}} = \sup_{\substack{\varphi \in C_0(\mathbb{R}) \\ \|\varphi\|_{\infty} = 1}} \langle u, \varphi \rangle$$

- $\|\cdot\|_{\mathcal{M}}$ is the total variation norm in the sense of measures.
- $\text{TV}(f) = \|Df\|_{\mathcal{M}}$
- Let $\text{TV}^m(f) = \text{TV}(D^{m-1}f) = \|D^m f\|_{\mathcal{M}}$ denote the m th-order total variation of f .
- Let

$$\text{BV}^m(\mathbb{R}) = \{f : \mathbb{R} \rightarrow \mathbb{R} : \text{TV}^m(f) < \infty\}$$

denote the space of functions of m th-order BV on \mathbb{R} .

\implies This is a **non-reflexive Banach space**.

Representer theorem for locally adaptive splines

Theorem (Fisher and Jerome 1975)

Consider the variational problem

$$\mathcal{V} := \arg \min_{f \in \text{BV}^m(\mathbb{R})} \|D^m f\|_{\mathcal{M}} \quad \text{s.t.} \quad \langle h_n, f \rangle = y_n, \quad n = 1, \dots, N,$$

where $h_n : \text{BV}^m(\mathbb{R}) \rightarrow \mathbb{R}$ is weak* continuous, i.e., h_n lies in the predual of $\text{BV}^m(\mathbb{R})$. Then, the solution set \mathcal{V} is nonempty, convex, and weak* compact. Moreover, the extreme points of \mathcal{V} are m th-order splines of the form

$$s(x) = \sum_{k=1}^K v_k \rho_m(x - t_k) + \sum_{n=0}^{m-1} c_n x^n, \quad K \leq N.$$

- These solutions are **sparse**: $\|D^m s\|_{\mathcal{M}} = \|\mathbf{v}\|_1$ and $K \leq N$.
- The knot locations $\{t_k\}_{k=1}^K$ are **adaptive**.

Proof sketch

Step 1: Understanding the Banach structure of $BV^m(\mathbb{R})$.

- $BV^m(\mathbb{R})$ is defined by the seminorm $f \mapsto \|D^m f\|_{\mathcal{M}}$
- The null space of this seminorm is the space of polynomials of degree $< m$, denoted $\mathcal{P}_{m-1}(\mathbb{R})$.
- We can equip $\mathcal{P}_{m-1}(\mathbb{R})$ with a **biorthogonal system** so that every $q \in \mathcal{P}_{m-1}(\mathbb{R})$ admits a unique representation

$$q = \sum_{k=1}^m \langle \phi_k, q \rangle p_k,$$

where $\{p_k\}_{k=1}^m$ is a basis for $\mathcal{P}_{m-1}(\mathbb{R})$ and $\phi_k \in \mathcal{P}_{m-1}(\mathbb{R})'$ such that $\langle \phi_k, p_n \rangle = \delta_{kn}$, e.g., choose $p_k(x) = \frac{x^{k-1}}{(k-1)!}$ and $\phi_k(x) = (-1)^{k-1} \delta^{(k-1)}(x)$.

Proof sketch (cont.)

- Consider the operator D_{ϕ}^{-m} whose Schwartz kernel is given by

$$\begin{aligned}g_{\phi}(x, y) &= \rho_m(x - y) - \sum_{k=1}^m \langle \phi_k, \rho_m(\cdot - y) \rangle p_k(x) \\ &= \rho_m(x - y) - \text{proj}_{\mathcal{P}_{m-1}(\mathbb{R})} \{ \rho_m(\cdot - y) \}(x)\end{aligned}$$

- For all $\mu \in \mathcal{M}(\mathbb{R})$, this operator satisfies

$$\begin{aligned}D^m D_{\phi}^{-m} \mu &= \mu \\ \phi(D_{\phi}^{-m} \mu) &= \mathbf{0},\end{aligned}$$

where $\phi(f) = (\langle \phi_1, f \rangle, \dots, \langle \phi_m, f \rangle)$.

Proof sketch (cont.)

- Define

$$\text{BV}_\phi^m(\mathbb{R}) = \{f \in \text{BV}^m(\mathbb{R}) : \phi(f) = \mathbf{0}\} \cong \text{BV}^m(\mathbb{R}) / \mathcal{P}_{m-1}(\mathbb{R}).$$

Therefore, $\text{BV}^m(\mathbb{R}) = \text{BV}_\phi^m(\mathbb{R}) \oplus \mathcal{P}_{m-1}(\mathbb{R})$.

- Since we factored out the null space, $\text{BV}_\phi^m(\mathbb{R})$ is a Banach space when equipped with the norm $f \mapsto \|D^m f\|_{\mathcal{M}}$.
- Consider the mapping $D^m : \text{BV}_\phi^m(\mathbb{R}) \rightarrow \mathcal{M}(\mathbb{R})$ that maps $f \mapsto \mu = D^m f$. Since $\|D^m f\|_{\mathcal{M}} = \|\mu\|_{\mathcal{M}}$, this map is norm preserving and it is invertible with inverse given by $D_\phi^{-m} : \mathcal{M}(\mathbb{R}) \rightarrow \text{BV}_\phi^m(\mathbb{R})$.
 \implies i.e., $\text{BV}_\phi^m(\mathbb{R}) \cong \mathcal{M}(\mathbb{R})$.

Proof sketch (cont.)

- Every $f \in BV^m(\mathbb{R})$ admits the direct-sum decomposition

$$f = D_{\phi}^{-m} \mu + q,$$

where $\mu = D^m f \in \mathcal{M}(\mathbb{R})$ and $q = \text{proj}_{\mathcal{P}_{m-1}(\mathbb{R})} f \in \mathcal{P}_{m-1}(\mathbb{R})$.

- When equipped with the norm

$$\|f\|_{BV^m(\mathbb{R})} = \|D^m f\|_{\mathcal{M}} + \|\phi(f)\|_2,$$

$BV^m(\mathbb{R})$ is a Banach space, in particular, a non-reflexive Banach space (since $\mathcal{M}(\mathbb{R})$ is a non-reflexive Banach space).

\implies e.g., $\|f\|_{BV(\mathbb{R})} = \text{TV}(f) + |f(0)|$.

Proof sketch (cont.)

Step 2: Transform the variational problem over $BV^m(\mathbb{R})$ into a variational problem over $\mathcal{M}(\mathbb{R})$.

- Recall the variational problem

$$\min_{f \in BV^m(\mathbb{R})} \|D^m f\|_{\mathcal{M}} \quad \text{s.t.} \quad H\{f\} = \mathbf{y} \in \mathbb{R}^N.$$

- From the direct-sum decomposition $f = D_{\phi}^{-m} \mu + q$, this problem is equivalent to

$$\min_{\mu \in \mathcal{M}(\mathbb{R})} \|\mu\|_{\mathcal{M}} \quad \text{s.t.} \quad H D_{\phi}^{-m} \mu = \mathbf{y} - Hq \in \mathbb{R}^N.$$

Proof sketch (cont.)

- Classical problem known as **Radon measure recovery**:

$$\mathcal{V} := \arg \min_{\mu \in \mathcal{M}(\mathbb{R})} \|\mu\|_{\mathcal{M}} \quad \text{s.t.} \quad \mathcal{A}\mu = z,$$

where $\mathcal{A} : \mathcal{M}(\mathbb{R}) \rightarrow \mathbb{R}^N$ is weak* continuous and linear.

- \implies It is well-known (Fisher and Jerome 1975; Zuhovickiř 1948) that the solution set \mathcal{V} is nonempty, convex, and weak* compact, and that the extreme points take the form

$$\sum_{k=1}^K v_k \delta(\cdot - t_k),$$

where $t_k \in \mathbb{R}$, $k = 1, \dots, K$, and $K \leq N$.

- \implies weak* continuity of the measurement operator plays a crucial role in proving that solutions exist.

Proof sketch (cont.)

- Therefore, the solution set to

$$\min_{f \in \text{BV}^m(\mathbb{R})} \|D^m f\|_{\mathcal{M}} \quad \text{s.t.} \quad \mathbb{H}\{f\} = \mathbf{y} \in \mathbb{R}^N.$$

nonempty, convex, and weak* compact. Moreover, the extreme points are m th-order splines of the form

$$s(x) = \sum_{k=1}^K v_k \rho_m(x - t_k) + \sum_{n=0}^{m-1} c_n x^n,$$

where $K \leq N$.

- ⇒ Plug in the solution to the Radon measure recovery problem into the direct-sum decomposition. □

What's going on?

$$\min_{f \in \text{BV}^m(\mathbb{R})} \|D^m f\|_{\mathcal{M}} \quad \text{s.t.} \quad \mathbf{H}\{f\} = \mathbf{y} \in \mathbb{R}^N.$$

- The extreme points of the solution set are written as a superposition of extreme points of the the unit ball associated to $f \mapsto \|D^m f\|_{\mathcal{M}}$ (in the quotient space $\text{BV}^m(\mathbb{R})/\mathcal{P}_{m-1}(\mathbb{R})$).

⇒ We showed that $\text{BV}_{\phi}^m(\mathbb{R}) \cong \text{BV}^m(\mathbb{R})/\mathcal{P}_{m-1}(\mathbb{R})$ is isometrically isomorphic to $\mathcal{M}(\mathbb{R})$.

⇒ It is well-known that the extreme points of the unit ball of $\mathcal{M}(\mathbb{R})$ take the form $\delta(\cdot - t_0)$, $t_0 \in \mathbb{R}$.

⇒ Therefore, the extreme points of $\text{BV}_{\phi}^m(\mathbb{R})$ take the form

$$\begin{aligned} D_{\phi}^{-m} \{\delta(\cdot - t_0)\} &= g_{\phi}(\cdot, t_0) \\ &= \rho_m(\cdot - t_0) - \text{proj}_{\mathcal{P}_{m-1}(\mathbb{R})} \{\rho_m(\cdot - t_0)\} \\ &\in [\rho_m(\cdot - t_0)] \end{aligned}$$

⇒ This unit ball has an interesting **geometry**.

A general representer theorem

Theorem

- Let $(\mathcal{X}, \mathcal{X}')$ be a dual pair of Banach spaces.
- Let $h_1, \dots, h_N \in \mathcal{X}$ be a set of measurement functionals.
- Let $H : \mathcal{X}' \rightarrow \mathbb{R}^N : f \mapsto (\langle h_1, f \rangle, \dots, \langle h_N, f \rangle)$ be a weak* continuous linear measurement operator.

Then, for any fixed $\mathbf{y} \in \mathbb{R}^N$, the solution set to

$$\arg \min_{f \in \mathcal{X}'} \|f\|_{\mathcal{X}'} \quad \text{s.t.} \quad H\{f\} = \mathbf{y}$$

is nonempty, convex, and weak* compact, and its extreme points take the form

$$\sum_{k=1}^K v_k u_k,$$

where $K \leq N$ and $u_k \in \text{Ext}\{f \in \mathcal{X}' : \|f\|_{\mathcal{X}'} \leq 1\}$.

A general representer theorem

Variants of this result have been proven by a number of authors:

- Claire Boyer et al. (2019). “On representer theorems and convex regularization”. In: **SIAM Journal on Optimization** 29.2, pp. 1260–1281
- Kristian Bredies and Marcello Carioni (2020). “Sparsity of solutions for variational inverse problems with finite-dimensional data”. In: **Calculus of Variations and Partial Differential Equations** 59.1, pp. 1–26
- Michael Unser (2021). “A unifying representer theorem for inverse problems and machine learning”. In: **Foundations of Computational Mathematics** 21.4, pp. 941–960

What is a neural network?

- Compositions of affine mappings and nonlinear mappings.
- A neural network $f : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_{L+1}}$ can be written as

$$\begin{cases} \mathbf{x}^{(0)} := \mathbf{x}, \\ \mathbf{x}^{(\ell)} := \boldsymbol{\rho}(\mathbf{A}^{(\ell)}\mathbf{x}^{(\ell-1)} - \mathbf{b}^{(\ell)}), \ell = 1, \dots, L, \\ f(\mathbf{x}) = \mathbf{A}^{(L+1)}\mathbf{x}^{(L)} - \mathbf{b}^{(L+1)}, \end{cases}$$

- $\implies L \in \mathbb{N}$ is the number of **hidden layers**.
- $\implies \mathbf{A}^{(\ell)} \in \mathbb{R}^{d_\ell \times d_{\ell-1}}$ are the **weights** of the neural network.
- $\implies \mathbf{b}^{(\ell)} \in \mathbb{R}^{d_\ell}$ are the **biases** of the neural network.
- $\implies \rho : \mathbb{R} \rightarrow \mathbb{R}$ is the **activation function**; ρ applies ρ entrywise.

Remark

When $L = 1$, the neural network is **shallow** and when $L > 1$, the neural network is **deep**.

What is a neural network?

- Today we will focus on **shallow** neural networks with **scalar** outputs.
- These are functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that can be written as

$$f(\mathbf{x}) = \mathbf{v}^\top \boldsymbol{\rho}(\mathbf{W}\mathbf{x} - \mathbf{b}) = \sum_{k=1}^K v_k \rho(\mathbf{w}_k^\top \mathbf{x} - b_k),$$

where $v_k \in \mathbb{R}$, $\mathbf{w}_k \in \mathbb{R}^d$, and $b_k \in \mathbb{R}$.

- A common choice for ρ is the truncated linear function $\rho(x) = x_+$, which we previously defined as ρ_2 .
- The atoms of a shallow neural network are functions of the form $\rho(\mathbf{w}_k^\top (\cdot) - b_k)$ and are often referred to as **neurons**.

Goal

Prove a representer theorem for neural networks.

Observations

- Recall the variational problem for locally adaptive splines

$$\min_{f \in \text{BV}^m(\mathbb{R})} \|D^m f\|_{\mathcal{M}} \quad \text{s.t.} \quad \langle h_n, f \rangle = y_n, \quad n = 1, \dots, N.$$

- The solution set is completely characterized by m th-order splines of the form

$$x \mapsto \sum_{k=1}^K v_k \rho_m(x - t_k) + \sum_{n=0}^{m-1} c_n x^n, \quad \rho_m(x) = \frac{x_+^{m-1}}{(m-1)!}.$$

- The operator D^m is a **sparsifying transform** to the atom $\rho_m(\cdot - t_k)$ since it **extracts the parameter** of the atom

$$D^m \rho_m(\cdot - t_k) = \delta(\cdot - t_k).$$

- In the proof of the representer theorem we were able to construct D_{ϕ}^{-m} , a **bounded, right inverse** of D^m .

Observations

- The atoms of a neural network take the form

$$\mathbf{x} \mapsto \rho(\mathbf{w}_0^\top \mathbf{x} - b_0), \quad \rho(x) = x_+$$

- Univariate neural networks with $\rho(x) = x_+$ are **linear splines**

$$\implies \mathbb{D}^2 \left\{ \sum_{k=1}^K v_k \rho(w_k(\cdot) - b_k) \right\} = \sum_{k=1}^K v_k |w_k| \delta(\cdot - b_k/w_k).$$

- The parameters of an atom are $\mathbf{w}_0 \in \mathbb{R}^d$ and $b_0 \in \mathbb{R}$.

$$\implies \text{WLOG, suppose } \mathbf{w}_0 \in \mathbb{S}^{d-1} = \{ \mathbf{v} \in \mathbb{R}^d : \|\mathbf{v}\|_2 = 1 \}.$$

$$\implies \rho(\mathbf{w}_0^\top \mathbf{x} - b_0) = \|\mathbf{w}_0\|_2 \rho\left(\left[\frac{\mathbf{w}_0}{\|\mathbf{w}_0\|_2} \right]^\top \mathbf{x} - \frac{b_0}{\|\mathbf{w}_0\|_2} \right)$$

Question

Question

Does there exist a sparsifying transform \mathbb{R} such that

$$\mathbb{R} \rho(\mathbf{w}_0^T(\cdot) - b_0) = \delta(\cdot - (\mathbf{w}_0, b_0)),$$

where $(\mathbf{w}_0, b_0) \in \mathbb{S}^{d-1} \times \mathbb{R}$?

Answer

Yes, and it involves **second-order differentiation** and the **Radon transform**.

The Radon transform

- For a sufficiently nice function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, the **Radon transform** is given by

$$\mathcal{R}\{f\}(\zeta) = \int_{\zeta} f,$$

$\implies \zeta \in \mathbb{P}^d$ is a $(d-1)$ -dimensional hyperplane in \mathbb{R}^d .

- Every hyperplane can be defined as the solution set to the equation $\gamma^T \mathbf{x} = t$, where $(\gamma, t) \in \mathbb{S}^{d-1} \times \mathbb{R}$.

$\implies (\gamma, t)$ and $(-\gamma, -t)$ are associated to the same hyperplane.

- We can identify $\mathcal{R}\{f\}$ with an even function on $\mathbb{S}^{d-1} \times \mathbb{R}$ given by

$$\mathcal{R}\{f\}(\gamma, t) = \int_{\gamma^T \mathbf{x} = t} f(\mathbf{x}) \, ds(\mathbf{x})$$

- The Radon transform as an integral operator whose Schwartz kernel is given by the distribution

$$k(\mathbf{x}, (\gamma, t)) = \delta(\gamma^T \mathbf{x} - t)$$

Inverting the dual Radon transform

- For a sufficiently nice function $\Phi : \mathbb{S}^{d-1} \times \mathbb{R} \rightarrow \mathbb{R}$, the **dual Radon transform** is given by

$$\mathcal{R}^*\{\Phi\}(\mathbf{x}) = \int_{\mathbb{S}^{d-1}} \Phi(\boldsymbol{\gamma}, \boldsymbol{\gamma}^T \mathbf{x}) \, d\sigma(\boldsymbol{\gamma})$$

- Given sufficiently nice $\Phi : \mathbb{S}^{d-1} \times \mathbb{R} \rightarrow \mathbb{R}$ such that $\Phi(\boldsymbol{\gamma}, t) = \Phi(-\boldsymbol{\gamma}, -t)$,

$$2(2\pi)^{d-1}\Phi = \Lambda^{d-1} \mathcal{R} \mathcal{R}^* \Phi = \mathcal{R} \mathcal{R}^* \Lambda^{d-1}\Phi,$$

where $\Lambda^{d-1} = (-\partial_t^2)^{\frac{d-1}{2}}$.

\implies In CT imaging parlance, Λ^{d-1} is known as the **backprojection filter** (or ramp filter).

The sparsifying transform

Lemma (Ongie et al., 2020, P. and Nowak, 2021)

Consider $r_{(\mathbf{w}_0, b_0)} = \rho(\mathbf{w}_0^\top(\cdot) - b_0)$, where $\rho(x) = x_+$. Then,

$$c_d \Lambda^{d-1} \mathcal{R} \Delta r_{(\mathbf{w}_0, b_0)} = \frac{\delta(\cdot - (\mathbf{w}_0, b_0)) + \delta(\cdot + (\mathbf{w}_0, b_0))}{2}$$

where $c_d := \frac{1}{2(2\pi)^{d-1}}$.

Proof

- Notice that

$$\Delta r_{(\mathbf{w}_0, b_0)} = \Delta \rho(\mathbf{w}_0^\top(\cdot) - b_0) = \delta(\mathbf{w}_0^\top(\cdot) - b_0),$$

which is the Schwartz kernel of the Radon transform.

- Next, consider the even test function ψ .
- Therefore,

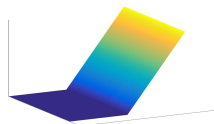
$$\begin{aligned} c_d \langle \Lambda^{d-1} \mathcal{R} \Delta r_{(\mathbf{w}_0, b_0)}, \psi \rangle &= c_d \langle \Delta r_{(\mathbf{w}_0, b_0)}, \mathcal{R}^* \Lambda^{d-1} \psi \rangle \\ &= (c_d \mathcal{R} \mathcal{R}^* \Lambda^{d-1} \psi)(\mathbf{w}_0, b_0) \\ &= \psi(\mathbf{w}_0, b_0) \end{aligned}$$



What's going on?

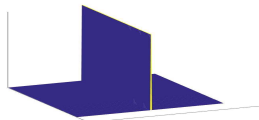
- Neural network atom

$$\Rightarrow \rho(\mathbf{w}_0^T(\cdot) - b_0), (\mathbf{w}_0, b_0) \in \mathbb{S}^{d-1} \times \mathbb{R}$$



- Laplacian of atom

$$\Rightarrow \Delta\{\rho(\mathbf{w}_0^T(\cdot) - b_0)\} = \delta(\mathbf{w}_0^T(\cdot) - b_0)$$



- Filtered Radon transform of Laplacian of atom

$$\Rightarrow (\Lambda^{d-1} \mathcal{R} \Delta)\{\rho(\mathbf{w}_0^T(\cdot) - b_0)\}(\boldsymbol{\gamma}, t) = \delta((\boldsymbol{\gamma}, t) - (\mathbf{w}_0, b_0)).$$

Other activation functions

- Due to the **intertwining relations** of the Laplacian and the Radon transform, we can write

$$c_d \Lambda^{d-1} \mathcal{R} \Delta = c_d \partial_t^2 \Lambda^{d-1} \mathcal{R}$$

- The operator $R_m := c_d \partial_t^m \Lambda^{d-1} \mathcal{R}$ sparsifies atoms of the form $\rho_m(\mathbf{w}_0^\top(\cdot) - b_0)$, where $\rho_m(x) = x_+^{m-1}/(m-1)!$.

\implies In particular,

$$R_m \rho_m(\mathbf{w}_0^\top(\cdot) - b_0) = \frac{\delta(\cdot - (\mathbf{w}_0, b_0)) + (-1)^m \delta(\cdot + (\mathbf{w}_0, b_0))}{2}$$

Higher-order BV spaces in the Radon domain

- Recall

$$\text{BV}^m(\mathbb{R}) = \{f : \mathbb{R} \rightarrow \mathbb{R} : D^m f \in \mathcal{M}(\mathbb{R})\}$$

Consider the space

$$\mathcal{R}\text{BV}^m(\mathbb{R}^d) := \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R} : \begin{array}{l} c_d \partial_t^m \Lambda^{d-1} \mathcal{R} f \in \mathcal{M}(\mathbb{S}^{d-1} \times \mathbb{R}), \\ \text{ess sup}_{\mathbf{x} \in \mathbb{R}^d} |f(\mathbf{x})| (1 + \|\mathbf{x}\|_2)^{-(m-1)} < \infty \end{array} \right\}$$

- When $d = 1$, one can show that $\mathcal{R}\text{BV}^m(\mathbb{R}) = \text{BV}^m(\mathbb{R})$.
 - \implies Define $\mathcal{R}\text{TV}^m(f) = c_d \|\partial_t^m \Lambda^{d-1} \mathcal{R} f\|_{\mathcal{M}}$.
 - \implies When $d = 1$, $\mathcal{R}\text{TV}^m(f) = \text{TV}^m(f)$.

Representer theorem for neural networks

Theorem (P. and Nowak, 2021)

Consider the variational problem

$$\arg \min_{f \in \mathcal{R}BV^m(\mathbb{R}^d)} c_d \|\partial_t^m \Lambda^{d-1} \mathcal{R} f\|_{\mathcal{M}} \quad \text{s.t.} \quad \langle h_n, f \rangle = y_n, \quad n = 1, \dots, N,$$

where $h_n : \mathcal{R}BV^m(\mathbb{R}^d) \rightarrow \mathbb{R}$ is weak* continuous, i.e., h_n lies in the predual of $\mathcal{R}BV^m(\mathbb{R}^d)$. Then, the solution set is nonempty, convex, and weak* compact. Moreover, the extreme points are shallow neural networks of the

$$s(\mathbf{x}) = \sum_{k=1}^K v_k \rho_m(\mathbf{w}_k^\top \mathbf{x} - b_k) + c(\mathbf{x}), \quad K \leq N,$$

where $c \in \mathcal{P}(\mathbb{R}^d)$ is a polynomial of degree $< m$.

Representer theorem for neural networks

- Proof is similar to the proof for locally adaptive splines.
- Similar results hold for **deep neural networks**. (P. and Nowak 2021c)
 - ⇒ Define $\mathcal{R}BV^m(\mathbb{R}^d; \mathbb{R}^D)$ and consider compositions of such functions.

What are $\mathcal{R}BV^m$ -spaces?

- Let $\Omega \subset \mathbb{R}^d$ be a bounded domain with nice boundary (e.g., Lipschitz). Then,

$$\mathcal{R}BV^m(\Omega) := \{f : \Omega \rightarrow \mathbb{R} : \exists g \in \mathcal{R}BV^m(\mathbb{R}^d) \text{ s.t. } g|_{\Omega} = f\}$$

- When $d = 1$, $\mathcal{R}BV^m(\Omega) = BV^m(\Omega)$.
 - $\implies B_{1,1}^m(\Omega) \subset BV^m(\Omega) \subset B_{1,\infty}^m(\Omega)$
 - \implies The **best** K -term approximation rate for $f \in BV^m(\Omega)$ is

$$\|f - f_K\|_{L^2(\Omega)} \lesssim K^{-m},$$

and is achieved by keeping the K largest Daubechies wavelet coefficients or with free knot spline approximation.

- When $d > 1$,
 - $\implies W^{d-1+m,1}(\Omega) \subset \mathcal{R}BV^m(\Omega)$. (P. and Nowak 2021b)
 - \implies We also know about the **best** approximation rates for $\mathcal{R}BV^m(\Omega)$.

Approximation properties of $\mathcal{R}BV^m(\Omega)$

Suppose $\Omega = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 \leq 1\}$.

- Every $f \in \mathcal{R}BV^m(\Omega)$ admits an integral representation

$$f(\mathbf{x}) = \int_{\mathbb{S}^{d-1} \times [-1,1]} \rho_m(\mathbf{w}^\top \mathbf{x} - b) \, d\mu(\mathbf{w}, b) + c(\mathbf{x}),$$

where $\rho_m(x) = x_+^{m-1}/(m-1)!$, $\mu \in \mathcal{M}(\mathbb{S}^{d-1} \times [-1,1])$, and $c(\cdot)$ is a polynomial of degree $< m$.

Dimension-free approximation rates

- Let $\mathcal{D} := \{g\}_{g \in \mathcal{D}}$ be a dictionary of atoms such that $g \in L^2(\Omega)$. If $\mathcal{D} \subset L^2(\Omega)$ is compact, then there exists

$$f_K = \sum_{k=1}^K \alpha_k g_k, \quad g_k \in \mathcal{D}$$

such that

$$\left\| \int_{\mathcal{D}} g(\cdot) \mu(g(\cdot)) - f_K \right\|_{L^2(\Omega)} \lesssim K^{-\frac{1}{2}}$$

- \implies This rate **does not grow** with the input dimension d .
- \implies Gilles Pisier (1981). “Remarques sur un résultat non publié de B. Maurey”. In: **Séminaire Analyse fonctionnelle (dit, pp. 1–12**
- In our problem, $\mathcal{D} = \{\rho_m(\mathbf{w}^\top(\cdot) - b)\}_{(\mathbf{w}, b) \in \mathbb{S}^{d-1} \times [-1, 1]}$.

Approximation properties of $\mathcal{R} \text{BV}^m(\Omega)$

- Given $f \in \mathcal{R} \text{BV}^m(\Omega)$, there exists

$$f_K(\mathbf{x}) = \sum_{k=1}^K v_k \rho_m(\mathbf{w}_k^\top \mathbf{x} - b_k) + c(\mathbf{x})$$

such that

$$\|f - f_K\|_{L^2(\Omega)} \lesssim K^{-\frac{1}{2} - \frac{2m-1}{2d}} \lesssim K^{-\frac{1}{2}}.$$

This is the **best** rate.

(Bach 2017; P. and Nowak 2021b; Siegel and Xu 2021)

- Compare this to the **best** K -term approximation rates in $H^m[0, 1]^d$, which scales as

$$\|f - f_K\|_{L^2[0,1]^d} \lesssim K^{-\frac{m}{d}}$$

and is achieved by truncated Fourier series approximation.

\Rightarrow This rate **grows exponentially** with the input dimension d .

Takeaway messages

- The intrinsic function spaces of neural networks are $\mathcal{R}BV^m$ -spaces.
 - ⇒ These are **not** classically studied function spaces. Perhaps explaining the lack of understanding of neural networks in practice.
- These spaces are **small** in the sense of their “dimension-free” approximation rates compared to classical function spaces, e.g., Sobolev spaces.
- When $m = 1$, $\mathcal{R}TV(\cdot) := \mathcal{R}TV^1(\cdot)$ is a **new** notion of multivariate total variation, different than $TV(f) = \|\nabla f\|_{\mathcal{M}}$.
 - ⇒ When $m > 1$, $\mathcal{R}TV^m(\cdot)$ provides a way to defined **higher-order** variants of multivariate total variation.

Open problems

- How do $\mathcal{R} BV^m$ -spaces relate to classical spaces, e.g., Besov, Triebel-Lizorkin, etc.?
 - ⇒ **Ridgelet analysis** is probably the right tool to answer this question.
- What other kinds of representer theorems can you get if you replace \mathcal{R} with a **generalized Radon transform** in $c_d \partial_t^m \Lambda^{d-1} \mathcal{R}$.
 - ⇒ Requires understanding the invertibility of the dual transform.

References I



Bach, Francis (2017). "Breaking the curse of dimensionality with convex neural networks". In: [Journal of Machine Learning Research](#).



Boyer, Claire, Antonin Chambolle, Yohann De Castro, Vincent Duval, Frédéric De Gournay, and Pierre Weiss (2019). "On representer theorems and convex regularization". In: [SIAM Journal on Optimization](#) 29.2, pp. 1260–1281.



Bredies, Kristian and Marcello Carioni (2020). "Sparsity of solutions for variational inverse problems with finite-dimensional data". In: [Calculus of Variations and Partial Differential Equations](#) 59.1, pp. 1–26.



Fisher, S. D. and Joseph W. Jerome (1975). "Spline solutions to L^1 extremal problems in one and several variables". In: [Journal of Approximation Theory](#).



Ongie, Greg, Rebecca Willett, Daniel Soudry, and Nathan Srebro (2020). "A Function Space View of Bounded Norm Infinite Width ReLU Nets: The Multivariate Case". In: [International Conference on Learning Representations](#).



P., Rahul and Robert D. Nowak (2020). "The role of neural network activation functions". In: [IEEE Signal Processing Letters](#).



— (2021a). "Banach Space Representer Theorems for Neural Networks and Ridge Splines". In: [Journal of Machine Learning Research](#).



— (2021b). "Near-Minimax Optimal Estimation With Shallow ReLU Neural Networks". In: [arXiv preprint arXiv:2109.08844](#).



— (2021c). "What Kinds of Functions do Deep Neural Networks Learn? Insights from Variational Spline Theory". In: [arXiv preprint arXiv:2105.03361](#).

References II



Pisier, Gilles (1981). "Remarques sur un résultat non publié de B. Maurey". In: [Séminaire Analyse fonctionnelle](#) (dit, pp. 1–12.



Siegel, Jonathan W. and Jinchao Xu (2021). "Characterization of the Variation Spaces Corresponding to Shallow Neural Networks". In: [arXiv preprint arXiv:2106.15002](#).



Unser, Michael (2021). "A unifying representer theorem for inverse problems and machine learning". In: [Foundations of Computational Mathematics](#) 21.4, pp. 941–960.



Zuhovickĩ, S. (1948). "Remarks on problems in approximation theory". In: [Mat. Zbirnik KDU](#), pp. 169–183.