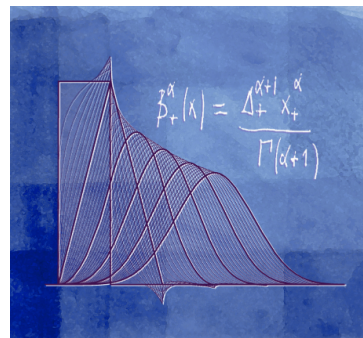


Part II: Variational Optimality of Neural Networks

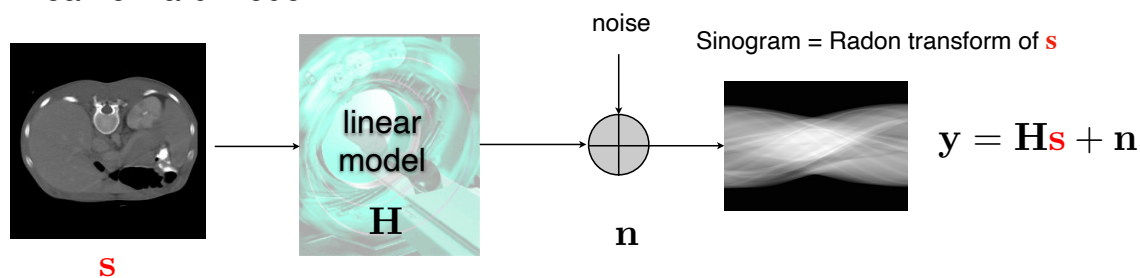
Michael Unser
 Biomedical Imaging Group
 EPFL, Lausanne, Switzerland



Summer School, Mathematics and Machine Learning for Image Analysis, Bologna, June 4-12, 2024

Variational formulation of inverse problems in imaging

Linear forward model



Problem: recover \mathbf{s} from noisy measurements \mathbf{y}

Regularization of ill-posed inverse problem

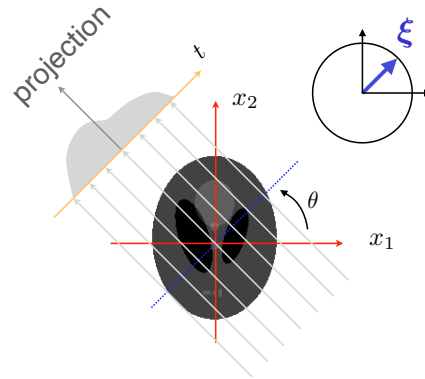
$$\mathbf{s}_{\text{rec}} = \arg \min_{\mathbf{s} \in \mathbb{R}^N} \underbrace{\|\mathbf{y} - \mathbf{H}\mathbf{s}\|_2^2}_{\text{data consistency}} + \underbrace{\lambda \|\mathbf{L}\mathbf{s}\|_p^p}_{\text{regularization}}, \quad p = 1, 2$$

The Radon transform and the FBP algorithm

Unit circle: $\mathbb{S}^1 = \{\xi \in \mathbb{R}^2 : \|\xi\| = 1\} = \{\xi = (\cos \theta, \sin \theta), \theta \in [0, 2\pi)\}$

- Radon transform of $s \in L_1(\mathbb{R}^2)$

$$\begin{aligned} R\{s\}(t, \xi) &= \int_{\mathbf{x} \in \mathbb{R}^2: \xi^T \mathbf{x} = t} s(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathbb{R}^2} \delta(t - \xi^T \mathbf{x}) s(\mathbf{x}) d\mathbf{x}, \quad (t, \xi) \in \mathbb{R} \times \mathbb{S}^1 \end{aligned}$$



- Reconstruction from $y(t, \xi) = R\{s\}(t, \xi)$: the **Filtered BackProjection** algorithm

$$s = R^* K_{\text{rad}} \{y\}$$

- K_{rad} : “**radial**” filtering in Radon space along the variable $t \in \mathbb{R}$.
Fourier symbol $\hat{K}_{\text{rad}}(\omega) \propto |\omega|$
- R^* : **backprojection** operator (the adjoint of R)

Supervised learning as a (linear) inverse problem

but an infinite-dimensional one ...

Given the data points $(\mathbf{x}_m, y_m) \in \mathbb{R}^N \times \mathbb{R}$, find $f : \mathbb{R}^N \rightarrow \mathbb{R}$ s.t. $f(\mathbf{x}_m) \approx y_m$ for $m = 1, \dots, M$

- Introduce smoothness or **regularization** constraint ($p = 2$) (Poggio-Girosi 1990)

$$R(f) = \|f\|_{\mathcal{H}}^2 = \|Lf\|_{L_2}^2 = \int_{\mathbb{R}^N} |Lf(\mathbf{x})|^2 d\mathbf{x}: \text{regularization functional}$$

$$\min_{f \in \mathcal{H}} R(f) \quad \text{subject to} \quad \sum_{m=1}^M |y_m - f(\mathbf{x}_m)|^2 \leq \sigma^2$$



- Regularized least-squares fit (theory of RKHS)

$$f_{\text{RKHS}} = \arg \min_{f \in \mathcal{H}} \left(\sum_{m=1}^M |y_m - f(\mathbf{x}_m)|^2 + \lambda R(f) \right) \quad \text{with} \quad R(f) = \|f\|_{\mathcal{H}}^2$$

⇒ kernel estimator

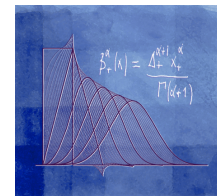
(Wahba 1990; Schölkopf 2001)

OUTLINE

- Connection with computational imaging ✓
- Variational formulation of learning: State-of-the art
 - Classical RKHS and kernel methods
 - Optimality results for shallow ReLU neural networks
- Radon-domain regularization for neural nets
 - Admissible regularization operator
 - Unifying representer theorem NEW
 - Laplacian revisited
 - Examples of admissible (operator, activation) pairs



(surprize?) ⇒ connection with *fractional splines*



5

Functions vs. distributions

■ Mathematical context

- $\mathcal{S}(\mathbb{R}^d)$: Schwartz's space of smooth and rapidly-decaying functions on \mathbb{R}^d

$$\varphi : \mathbb{R}^d \rightarrow \mathbb{R} \quad (\text{or } \mathbb{C})$$

$$\mathbf{x} \mapsto \varphi(\mathbf{x})$$



Laurent Schwartz (1915-2002)

- $\mathcal{S}'(\mathbb{R}^d)$: the space of **continuous linear functionals** on $\mathcal{S}(\mathbb{R}^d)$ = tempered distributions

$$f : \mathcal{S}(\mathbb{R}^d) \rightarrow \mathbb{R} \quad (\text{or } \mathbb{C})$$

$$\varphi \mapsto \langle f, \varphi \rangle = \int_{\mathbb{R}^d} f(\mathbf{x})\varphi(\mathbf{x})d\mathbf{x} \quad [\text{Formal or explicit (for locally-integrable functions)}]$$

- $L_2(\mathbb{R}^d)$: space of square-integrable functions on \mathbb{R}^d

- L_2 -norm: $\|f\|_{L_2} = \left(\int_{\mathbb{R}^d} |f(\mathbf{x})|^2 d\mathbf{x} \right)^{\frac{1}{2}}$

- $L_2(\mathbb{R}^d) = \{f : \mathbb{R}^d \rightarrow \mathbb{R} \text{ with } f \text{ measurable and } \|f\|_{L_2} < \infty\} = \overline{(\mathcal{S}(\mathbb{R}^d), \|\cdot\|_{L_2})}$

- Continuous and dense embeddings: $\mathcal{S}(\mathbb{R}^d) \xhookrightarrow{d} L_2(\mathbb{R}^d) \xhookrightarrow{d} \mathcal{S}'(\mathbb{R}^d)$

6

RKHS representer theorem for L_2 regularization ($p = 2$)

$$(P2) \quad \arg \min_{f \in \mathcal{H}} \left(\sum_{m=1}^M |y_m - f(\mathbf{x}_m)|^2 + \lambda \|f\|_{\mathcal{H}}^2 \right)$$

(deBoor 1966; Poggio-Girosi 1991)

$r_{\mathcal{H}} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is the (unique) **reproducing kernel** for the Hilbert $\mathcal{H} \subset \mathcal{S}'(\mathbb{R}^d)$ if

- $r_{\mathcal{H}}(\cdot, \mathbf{x}_0) \in \mathcal{H}$ for all $\mathbf{x}_0 \in \mathbb{R}^d$
- $f(\mathbf{x}_0) = \langle r_{\mathcal{H}}(\cdot, \mathbf{x}_0), f \rangle_{\mathcal{H}}$ for all $f \in \mathcal{H}$ and $\mathbf{x}_0 \in \mathbb{R}^d$

$$\Leftrightarrow \delta(\cdot - \mathbf{x}_0) \in \mathcal{H}'$$

Convex loss function: $E : \mathbb{R}^M \times \mathbb{R}^M \rightarrow \mathbb{R}$

Sample values: $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_M))$

$$(P2') \quad \arg \min_{f \in \mathcal{H}} (E(\mathbf{y}, \mathbf{f}) + \lambda \|f\|_{\mathcal{H}}^2)$$

(Schölkopf-Smola 2001)

Representer theorem for L_2 -regularization

The generic parametric form of the solution of (P2') is

$$f(\mathbf{x}) = \sum_{m=1}^M a_m r_{\mathcal{H}}(\mathbf{x}, \mathbf{x}_m)$$

Supports the theory of SVM, kernel methods, variational splines, etc.

7

And what about neural networks ?

Link with splines (gTV)

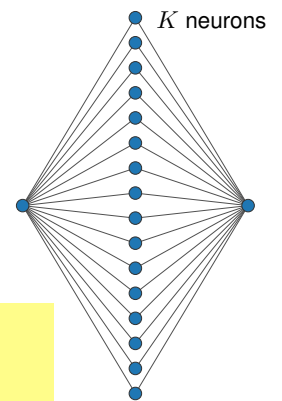
- **Shallow univariate** ReLU neural network with skip connection

$$f_{\theta}(x) = c_0 + c_1 x + \sum_{k=1}^K v_k (w_k x - b_k)_+$$

$$= c_0 + c_1 x + \sum_{k=1}^{K_0} a_k (x - \tau_k)_+$$

- Standard training with weight decay

$$(NN-1) : \arg \min_{\theta=(\mathbf{v}, \mathbf{w}, \mathbf{b}, \mathbf{c})} \sum_{m=1}^M |y_m - f_{\theta}(x_m)|^2 + \frac{\lambda}{2} \sum_{k=1}^K |v_k|^2 + |w_k|^2$$



Theorem

For any $K \geq K_0$ (with $K_0 < M$), the solution of (NN-1) is achieved by the **sparse adaptive spline**:

$$f_{\text{spline}} = \arg \min_{f \in \text{BV}^{(2)}(\mathbb{R})} \left(\sum_{m=1}^M |y_m - f(x_m)|^2 + \lambda \|D^2 f\|_{\mathcal{M}} \right).$$

(U.-Fageot-Ward, 2017; Savarese 2019; Parhi-Nowak 2020)

8

Proper continuous counterpart of ℓ_1 -norm



2



Johann Radon (1887-1956)

- Dual definition of ℓ_1 -norm (in finite dimensions only)

$$\|\mathbf{f}\|_{\ell_1} = \sum_{n=1}^N |f_n| = \sup_{\mathbf{u} \in \mathbb{R}^N: \|\mathbf{u}\|_{\infty} \leq 1} \langle \mathbf{f}, \mathbf{u} \rangle$$

- Space $C_0(\mathbb{R}^d)$ of functions on \mathbb{R}^d that are continuous, bounded, and decaying at infinity

$$C_0(\mathbb{R}^d) = \overline{(\mathcal{S}(\mathbb{R}^d), \|\cdot\|_{L_{\infty}})} \subset L_{\infty}(\mathbb{R}^d)$$

- Space of **bounded Radon measures** on \mathbb{R}^d

$$\mathcal{M}(\mathbb{R}^d) = (C_0(\mathbb{R}^d))' = \{f \in \mathcal{S}'(\mathbb{R}^d) : \|f\|_{\mathcal{M}} \triangleq \sup_{\varphi \in \mathcal{S}(\mathbb{R}^d): \|\varphi\|_{\infty} \leq 1} \langle f, \varphi \rangle < +\infty\}$$

- **Superset** of $L_1(\mathbb{R}^d)$

$$\forall f \in L_1(\mathbb{R}^d) : \|f\|_{\mathcal{M}} = \|f\|_{L_1} \Rightarrow L_1(\mathbb{R}^d) \subset \mathcal{M}(\mathbb{R}^d)$$

- **Extreme points** of unit ball in $\mathcal{M}(\mathbb{R}^d)$: $e_k = \pm \delta(\cdot - \boldsymbol{\tau}_k)$ with $\boldsymbol{\tau}_k \in \mathbb{R}^d$

9

Multi-dimensional extension via hyper-spherical measures

- Integral representation of infinite-width shallow neural network

$$f(\mathbf{x}) = \int_{\mathbb{R} \times \mathbb{S}^{d-1}} \sigma(\boldsymbol{\xi}^T \mathbf{x} - t) d\mu(t, \boldsymbol{\xi}) \mu(\cdot, \boldsymbol{\xi}) (\mathbf{x}) = \mathbf{R}^* \{ \sigma \otimes \mu(\cdot, \boldsymbol{\xi}) \} (\mathbf{x}) \quad \mathbf{R}^*: \text{Radon's backprojection operator}$$

- Hyper-spherical counterpart of spike deconvolution problem (Duval-Peyré 2014; Bach 2017)

$$\arg \min_{\mu \in \mathcal{M}(\mathbb{R} \times \mathbb{S}^{d-1})} \left(\sum_{m=1}^M |y_m - \mathbf{R}^* \{ \sigma \otimes \mu \} (\mathbf{x}_m)|^2 + \lambda \|\mu\|_{\mathcal{M}} \right)$$

$$\text{Existence of solutions of the form: } f(\mathbf{x}) = \sum_{k=1}^{K_0} a_k \sigma(\boldsymbol{\xi}_k^T \mathbf{x} - t_k)$$

- Reproducing kernel Banach space (RKBS) (Bartolucci-DeVito-Rosasco-Vigogna; ACHA 2023)

$$\begin{aligned} \mathcal{B} &= \{f_{\mu} : \mu \in \mathcal{M}(\Theta)\}, \\ f_{\mu}(x) &= \int_{\Theta} \rho(x, \theta) \beta(\theta) d\mu(\theta) \\ \|f\|_{\mathcal{B}} &= \inf \{ \|\mu\|_{\mathcal{M}} : f_{\mu} = f \} \end{aligned}$$

$$\begin{aligned} \arg \min_{f \in \mathcal{B}} \left(\sum_{m=1}^M |y_m - f(x_m)|^2 + \lambda \|f\|_{\mathcal{B}} \right) \\ \Rightarrow f(x) = \sum_{k=1}^{K_0} a_k \rho(x, \theta_k) \end{aligned}$$

10

OUTLINE

- Connection with computational imaging ✓
- Variational formulation of learning: State-of-the-art ✓
 - Classical RKHS and kernel methods
 - Optimality results for shallow ReLU neural networks
- **Radon-domain regularization** yields neural nets
 - Admissible regularization operators
 - Null space of polynomials
 - **Unifying representer theorem**
 - **Native** spaces
 - Example of admissible (operator, activation) pairs

11

Admissible regularization operator

■ Isotropic convolution operator

A linear operator $L : \mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}'(\mathbb{R}^d)$ that is **shift-invariant** and **isotropic** is uniquely characterized by its **radial frequency profile** $\widehat{L}_{\text{rad}} : \mathbb{R} \rightarrow \mathbb{R}$.

Fourier symbol of L : $\widehat{L}(\omega) = \widehat{L}_{\text{rad}}(\|\omega\|)$

Definition

An isotropic regularization operator with frequency profile $\widehat{L}_{\text{rad}}(\omega)$ is **spline-admissible** with a **polynomial null space** of degree n_0 (possibly trivial) if

1. $\widehat{L}_{\text{rad}}(\omega)$ **does not vanish** over \mathbb{R} , except for a **zero** of order $\gamma_0 \in (n_0, n_0 + 1]$ **at the origin**; that is, $|\widehat{L}_{\text{rad}}(\omega)|/|\omega|^{\gamma_0} = C_0$ as $\omega \rightarrow 0$.
2. **Ellipticity**: There exists an order $\gamma_1 > 1$, a constant $C_1 > 0$, and a radius $R_1 > 0$ such that $|\widehat{L}_{\text{rad}}(\omega)| \geq C_1|\omega|^{\gamma_1}$ for all $|\omega| > R_1$.

■ Example

$L = \Delta$ (Laplacian) with $\widehat{L}_{\text{rad}}(\omega) = -\omega^2$, $\gamma_0 = \gamma_1 = 2$, and $n_0 = 1$.

12

Null space of polynomials

Isotropic regularization operator L with frequency profile $\widehat{L}_{\text{rad}} : \mathbb{R} \rightarrow \mathbb{R}$

- Effect of a γ_0 th-order zero

$$\lim_{\omega \rightarrow 0} \frac{|\widehat{L}_{\text{rad}}(\omega)|}{|\omega|^{\gamma_0}} = C_0 \quad \Rightarrow \quad \text{annihilates polynomials of degree } n_0 = \lceil \gamma_0 - 1 \rceil$$

- The space \mathcal{P}_{n_0} of polynomials of degree n_0

- Taylor (or monomial) basis: $m_{\mathbf{k}}(\mathbf{x}) \triangleq \frac{\mathbf{x}^{\mathbf{k}}}{\mathbf{k}!}$
- $\mathcal{P}_{n_0} = \{p_0 = \sum_{|\mathbf{k}| \leq n_0} b_{\mathbf{k}} m_{\mathbf{k}} : \|p_0\|_{\mathcal{P}} \triangleq \| (b_{\mathbf{k}})_{|\mathbf{k}| \leq n_0} \|_2 < \infty\}$.

Proposition (Construction of biorthogonal basis)

There exists an isotropic window $\kappa_{\text{iso}} \in \mathcal{S}(\mathbb{R}^d)$ with $0 \leq \widehat{\kappa}_{\text{iso}}(\boldsymbol{\omega}) \leq 1$ and $\widehat{\kappa}_{\text{iso}}(\boldsymbol{\omega}) = 0$ for $\|\boldsymbol{\omega}\| \geq 1$ such that, for all $\mathbf{k}, \mathbf{n} \in \mathbb{N}^d$,

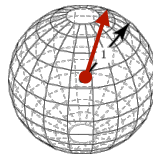
$$m_{\mathbf{n}}^* \triangleq (-1)^{|\mathbf{n}|} \partial^{\mathbf{n}} \kappa_{\text{iso}} \quad \text{and} \quad \langle m_{\mathbf{k}}, m_{\mathbf{n}}^* \rangle = \delta_{\mathbf{k}-\mathbf{n}} \quad (\text{biorthogonality})$$

- Dual space $\mathcal{P}'_{n_0} = \{p_0^* = \sum_{|\mathbf{k}| \leq n_0} b_{\mathbf{k}}^* m_{\mathbf{k}}^* : \|p_0^*\|_{\mathcal{P}'} \triangleq \| (b_{\mathbf{k}}^*) \|_2 < \infty\} \subset \mathcal{S}(\mathbb{R}^d)$

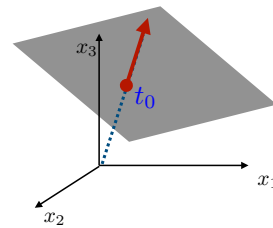
13

The Radon transform: Classical integral formulation

Unit sphere: $\mathbb{S}^{d-1} = \{\boldsymbol{\xi} \in \mathbb{R}^d : \|\boldsymbol{\xi}\| = 1\}$



Hyperplane $P_{\boldsymbol{\xi}_0, t_0} = \{\mathbf{x} \in \mathbb{R}^d : \boldsymbol{\xi}_0^T \mathbf{x} = t_0\}$



- Radon transform of $f \in L_1(\mathbb{R}^d)$

$$\begin{aligned} \mathbb{R}\{f\}(t, \boldsymbol{\xi}) &= \int_{\boldsymbol{\xi}^T \mathbf{x} = t} f(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathbb{R}^d} \delta(t - \boldsymbol{\xi}^T \mathbf{x}) f(\mathbf{x}) d\mathbf{x}, \quad (t, \boldsymbol{\xi}) \in \mathbb{R} \times \mathbb{S}^{d-1} \end{aligned}$$

- Backprojection operator: From Radon domain to Euclidean space

$$\mathbb{R}^*\{g\}(\mathbf{x}) = \int_{\mathbb{S}^{d-1}} \underbrace{g(\boldsymbol{\xi}^T \mathbf{x}, \boldsymbol{\xi})}_t d\boldsymbol{\xi}, \quad \mathbf{x} \in \mathbb{R}^d$$

- Fourier slice theorem

$$\mathbb{R}\{f\}(t, \boldsymbol{\xi}_0) = \frac{1}{2\pi} \int_{\mathbb{R}} \widehat{f}(\boldsymbol{\omega} \boldsymbol{\xi}_0) e^{i\boldsymbol{\omega} t} d\boldsymbol{\omega} = \mathcal{F}_{\boldsymbol{\omega} \rightarrow t}^{-1} \{\widehat{f}(\boldsymbol{\omega} \boldsymbol{\xi}_0)\} \{t\}$$

Fourier transform:

$$\widehat{f}(\boldsymbol{\omega}) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} f(\mathbf{x}) e^{-i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} d\mathbf{x}$$

14

New twist: Generic, Radon-domain regularization

Isotropic regularization operator L with frequency profile $\widehat{L}_{\text{rad}} : \mathbb{R} \rightarrow \mathbb{R}$

“Radonized” version of the operator:

$$L_R = K_{\text{rad}}RL$$

■ Role of each operator

- L : differential operator such as Laplacian, to penalize high frequency components
- R : Radon transform, to project in hyperspherical domain with $(t, \xi) \in \mathbb{R} \times \mathbb{S}^{d-1}$
- K_{rad} : isotropic Radon-domain filtering with $\widehat{K}_{\text{rad}}(\omega) \propto |\omega|^{d-1}$, to facilitate inversion

■ “Easy” case where L is invertible = trivial null space

- $L^{-1}L = \text{Id}$ on $\mathcal{S}'(\mathbb{R}^d)$
- L_R has a trivial null space
- Inversion of Radon transform: $R^*K_{\text{rad}}R = \text{Id}$ on $\mathcal{S}'(\mathbb{R}^d)$
 $\Rightarrow L_R^{-1} = L^{-1}R^*$ (Ludwig 1966)

■ Non-trivial null space

- L and L_R share the same null space: \mathcal{P}_{n_0}
- Canonical scenario: $n_0 = 1$ (affine maps)
- Makes the inversion process more difficult

15

Representer theorem for neural nets: Context

Given the data points $(\mathbf{x}_m, y_m) \in \mathbb{R}^d \times \mathbb{R}$, find $f : \mathbb{R}^d \rightarrow \mathbb{R}$ s.t. $f(\mathbf{x}_m) \approx y_m$ for $m = 1, \dots, M$

■ Variational formulation with Radon-domain regularization

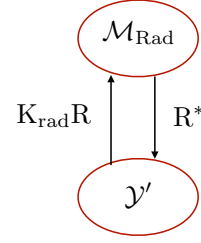
$$S = \arg \min_{f \in \mathcal{M}_{L_R}(\mathbb{R}^d)} \sum_{m=1}^M E(y_m, f(\mathbf{x}_m)) + \psi(\|L_R f\|_{\mathcal{M}(\mathbb{R} \times \mathbb{S}^{d-1})})$$

- **Regularization operator:** $L_R = K_{\text{rad}}RL : \mathcal{M}_{L_R}(\mathbb{R}^d) \rightarrow \mathcal{M}_{\text{Rad}}$ where L is admissible
- **Native space:** $\mathcal{M}_{L_R}(\mathbb{R}^d)$ = Banach space that is isometrically isomorphic to $\mathcal{M}_{\text{Rad}} \times \mathcal{P}_{n_0}$
- $\mathcal{M}_{\text{Rad}} = \mathcal{M}_{\text{even}}(\mathbb{R} \times \mathbb{S}^{d-1})$: Banach space of **Radon-compatible bounded measures**
- $E : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^+$ is a strictly-convex loss functional.
- $\psi : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is some arbitrary strictly-increasing convex function.

16

Theoretical ingredients

- Specification of Euclidean-to-Radon-domain isomorphisms
 - Problem: the Radon transform is not surjective on \mathcal{S} (resp., \mathcal{S}')
 - Identification of proper hyperspherical Banach subspace $\mathcal{M}_{\text{Rad}} \subset \mathcal{M}(\mathbb{R} \times \mathbb{S}^{d-1})$ over which \mathbb{R}^* is invertible



⇒ **Banach-space variant** of the distributional theory of the Radon transform

(Unser, JMLR 2023)

- Dealing with the non-trivial null space of $L_{\mathbb{R}}$

- Take inspiration from spline theory

(U.-Fageot-Ward, SIAM Rev 2017)

$$\text{Factoring out the null space: } (\text{Id} - \text{Proj}_{\mathcal{P}_{n_0}})\{f\} = f - \sum_{|\mathbf{k}| \leq n_0} \langle f, m_{\mathbf{k}}^* \rangle m_{\mathbf{k}}$$

- Abstract representer theorem for Banach space with semi-norm penalties (U.-Aziznejad, ACHA 2022)

- Prove that $\mathcal{M}_{L_{\mathbb{R}}}$ is a Banach space
- Establish weak* continuity of sampling functionals

(Unser FoCM in press)

17

Representer theorem for neural nets

$$S = \arg \min_{f \in \mathcal{M}_{L_{\mathbb{R}}}(\mathbb{R}^d)} \sum_{m=1}^M E(y_m, f(\mathbf{x}_m)) + \psi(\|L_{\mathbb{R}} f\|_{\mathcal{M}(\mathbb{R} \times \mathbb{S}^{d-1})}) \quad (1)$$

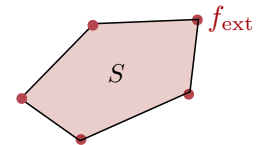
$$L_{\mathbb{R}} = K_{\text{rad}} \text{RL} : \mathcal{M}_{L_{\mathbb{R}}}(\mathbb{R}^d) \rightarrow \mathcal{M}(\mathbb{R} \times \mathbb{S}^{d-1})$$

Theorem

The solution set S of Problem (1) is non-empty and weak* compact. It is the weak* closure of the convex hull of its **extreme points**, which can all be written as

$$f_{\text{ext}}(\mathbf{x}) = p_0(\mathbf{x}) + \sum_{k=1}^{K_0} a_k \rho_{\text{rad}}(\boldsymbol{\xi}_k^{\top} \mathbf{x} - \tau_k)$$

with a fixed **activation function** $\rho_{\text{rad}} = \mathcal{F}^{-1}\{1/\widehat{L}_{\text{rad}}\}$, for some $K_0 \leq M - \dim \mathcal{P}_{n_0}$, $(a_k, \boldsymbol{\xi}_k, \tau_k) \in \mathbb{R} \times \mathbb{S}^{d-1} \times \mathbb{R}$ for $k = 1, \dots, K_0$, and a **null-space component** $p_0 \in \mathcal{P}_{n_0}$. The corresponding regularization cost (shared by all solutions) is $\|L_{\mathbb{R}} f_{\text{ext}}\|_{\mathcal{M}(\mathbb{R} \times \mathbb{S}^{d-1})} = \sum_{n=1}^{K_0} |a_n|$.



Special case of abstract rep theorem for direct sums (U.-Aziznejad, ACHA 2022)

18

Special case: Laplacian

- Properties of the (negative) Laplacian

- $(-\Delta)f(\mathbf{x}) = -\sum_{n=1}^d \frac{\partial^2}{\partial x_n^2} f(\mathbf{x})$
- Frequency symbol: $\|\boldsymbol{\omega}\|^2 \Rightarrow \widehat{\Delta}_{\text{rad}}(\omega) = \omega^2$ (radial profile)
- Annihilates all affine functions: $n_0 = 1$

- Outcome of representer theorem with $L_R = K_{\text{rad}}R(-\Delta)$

- Null space: $\mathcal{P}_1 = \{p_0(\mathbf{x}) = b_0 + \mathbf{b}^\top \mathbf{x} : (b_0, \mathbf{b}) \in \mathbb{R}^{d+1}\}$
- Activation function: $\rho_{\text{rad}}(t) = \mathcal{F}^{-1}\{\frac{1}{\omega^2}\}(t) = \frac{1}{2}|t| = t_+ - \frac{1}{2}t$

$$\Rightarrow \text{Shallow ReLU net: } f_{\text{ext}}(\mathbf{x}) = c_0 + \mathbf{c}^\top \mathbf{x} + \sum_{k=1}^{K_0} a_k (\boldsymbol{\xi}_k^\top \mathbf{x} - \tau_k)_+$$

sum of elementary (ReLU) ridges = **ridge spline** (Parhi-Nowak 2021)

19

Limit behaviour of multivariate 2-layer ReLU neural nets

- Shallow ReLU neural network $\mathbb{R}^d \rightarrow \mathbb{R}$ with skip connection

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = c_0 + \mathbf{c}_1^\top \mathbf{x} + \sum_{k=1}^K v_k (\mathbf{w}_k^\top \mathbf{x} - b_k)_+ = c_0 + \mathbf{c}_1^\top \mathbf{x} + \sum_{k=1}^{K_0} a_k (\boldsymbol{\xi}_k^\top \mathbf{x} - \tau_k)_+$$

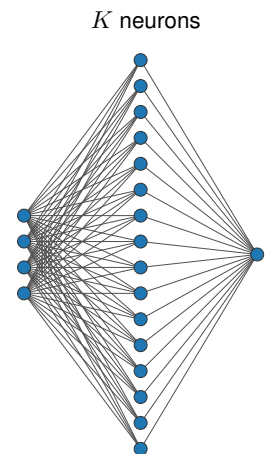
- Standard training with weight decay on $\mathbf{v} = (v_k)$ and $\mathbf{W} = [\mathbf{w}_1 \dots \mathbf{w}_K]$

$$(\text{NN-}d) : \arg \min_{\boldsymbol{\theta}=(\mathbf{v}, \mathbf{W}, \mathbf{b}, \mathbf{c})} \sum_{m=1}^M |y_m - f_{\boldsymbol{\theta}}(\mathbf{x}_m)|^2 + \frac{\lambda}{2} \sum_{k=1}^K |v_k|^2 + \|\mathbf{w}_k\|^2$$

Theorem

For any $K \geq K_0$ (with $K_0 < M$), the solution of (NN- d) is achieved by the **sparse ridge spline**:

$$f_{\text{ridge}} = \arg \min_{f \in \mathcal{M}_{\Delta_R}(\mathbb{R}^d)} \left(\sum_{m=1}^M |y_m - f(\mathbf{x}_m)|^2 + \lambda \|K_{\text{rad}}R\Delta f\|_{\mathcal{M}(\mathbb{R} \times \mathbb{S}^{d-1})} \right).$$



(Ongie et al. 2020; Parhi-Nowak 2021)

Delicate point: Proper delineation of the native space $\mathcal{M}_{\Delta_R}(\mathbb{R}^d)$

20

NATIVE SPACES

Hyper-spherical (test) functions and distributions

■ Test functions and tempered distributions

$$(\phi, g) \in \mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1}) \times \mathcal{S}'(\mathbb{R} \times \mathbb{S}^{d-1})$$

$$g : \phi \mapsto \langle g, \phi \rangle_{\text{Rad}} \in \mathbb{R}$$

For locally integrable functions $g : (t, \xi) \mapsto \mathbb{R}$:

$$\langle g, \phi \rangle_{\text{Rad}} = \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} g(t, \xi) \phi(t, \xi) dt d\xi$$

Special case: $d = 2$

$\xi = (\cos \theta, \sin \theta)$ with $d\xi = d\theta$ for $\theta \in [0, 2\pi]$

$$\langle g, \phi \rangle_{\text{Rad}} = \int_0^{2\pi} \int_{\mathbb{R}} g(t, \theta) \phi(t, \theta) dt d\theta$$

■ Radon transform and its adjoint

$$R : \mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1})$$

$$R^* : \mathcal{S}'(\mathbb{R} \times \mathbb{S}^{d-1}) \rightarrow \mathcal{S}'(\mathbb{R}^d)$$

R^* is the unique linear operator such that

$$\forall \varphi \in \mathcal{S}(\mathbb{R}^d) : \langle R^* \{g\}, \varphi \rangle = \langle g, R\{\varphi\} \rangle_{\text{Rad}}$$

Radon transform on $\mathcal{S}(\mathbb{R}^d)$

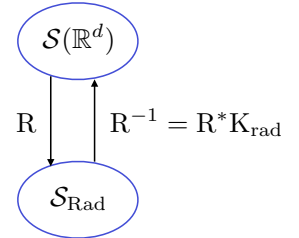
Theorem (Invertibility of Radon transform on $\mathcal{S}(\mathbb{R}^d)$)

1. R continuously maps $\mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}_{\text{Rad}} \subset \mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1})$
2. $R^*K_{\text{rad}}R = \text{Id}$ on $\mathcal{S}(\mathbb{R}^d) \Leftrightarrow R^{-1} = R^*K_{\text{rad}}$ on $\mathcal{S}_{\text{Rad}} = R(\mathcal{S}(\mathbb{R}^d))$

(Gelfand 1962; Helgason 1965; Ludwig 1966)

■ Radon-domain filtering operator

- K_{rad} : “radial” operator that acts along the Radon-domain variable t
- Radial frequency response: $\widehat{K}_{\text{rad}}(\omega) = c_d|\omega|^{d-1}$



23

Distributional theory of the (filtered) Radon transform

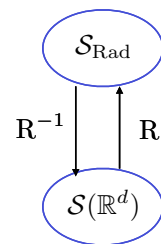
$$R : \mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1})$$

Difficulty: injective but not surjective !

Theorem (Variant on invertibility of Radon transform)

$\mathcal{S}_{\text{Rad}} = R(\mathcal{S}(\mathbb{R}^d))$ is a closed subspace of $\mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1})$. Moreover,

1. $R : \mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}_{\text{Rad}}$ is a **continuous bijection**, with $R^*K_{\text{rad}}R = \text{Id}$ on $\mathcal{S}(\mathbb{R}^d)$
2. $R^*K_{\text{rad}} : \mathcal{S}_{\text{Rad}} \rightarrow \mathcal{S}(\mathbb{R}^d)$ is a **continuous bijection**, with $RR^*K_{\text{rad}} = \text{Id}$ on \mathcal{S}_{Rad}
3. $R^* : \mathcal{S}'_{\text{Rad}} \rightarrow \mathcal{S}'(\mathbb{R}^d)$ is a **continuous bijection** with $K_{\text{rad}}RR^* = \text{Id}$ on $\mathcal{S}'_{\text{Rad}}$.

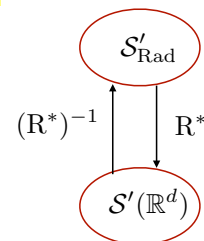


Theorem (Characterization of the range space)

Let $\phi \in \mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1})$. Then, $\phi \in \mathcal{S}_{\text{Rad}} \stackrel{\Delta}{=} \{\phi = R\{\varphi\} : \varphi \in \mathcal{S}(\mathbb{R}^d)\}$ iff.

1. Evenness: $\phi(t, \xi) = \phi(-t, -\xi)$.
2. $\Phi_k(\xi) = \int_{\mathbb{R}} \phi(t, \xi)t^k dt$ is a homogeneous polynomial in $\xi \in \mathbb{S}^{d-1}$ for any $k \in \mathbb{N}_0$

(Gelfand 1962; Helgason 1965; Ludwig 1966)



24

Banach space theory of the (filtered) Radon transform

$$\mathcal{X} = \overline{(\mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1}), \|\cdot\|_{\mathcal{X}})} \subset \mathcal{S}'(\mathbb{R} \times \mathbb{S}^{d-1})$$

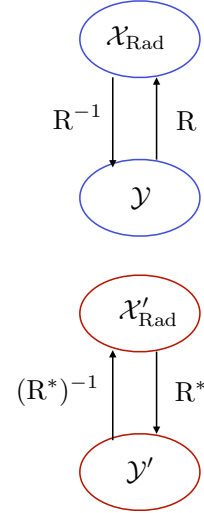
$$\mathcal{X}' = \{g \in \mathcal{S}'(\mathbb{R} \times \mathbb{S}^{d-1}) : \|g\|_{\mathcal{X}'} < \infty\} \quad \text{with} \quad \|g\|_{\mathcal{X}'} = \sup_{\phi \in \mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1}) : \|\phi\|_{\mathcal{X}} \leq 1} \langle g, \phi \rangle_{\text{Rad}}$$

$$\mathcal{X}_{\text{Rad}} = \overline{(\mathcal{S}_{\text{Rad}}, \|\cdot\|_{\mathcal{X}})} \subset \mathcal{X}$$

Theorem (Radon-compatible Banach isometries)

Let $\|\varphi\|_{\mathcal{Y}} \triangleq \|\mathbb{R}\{\varphi\}\|_{\mathcal{X}}$. Then, $\mathbb{R} : (\mathcal{S}(\mathbb{R}^d), \|\cdot\|_{\mathcal{Y}}) \rightarrow (\mathcal{S}_{\text{Rad}}, \|\cdot\|_{\mathcal{X}})$ has a **unique isometric extension** $\mathbb{R} : \mathcal{Y} \rightarrow \mathcal{X}_{\text{Rad}}$ with $\mathcal{Y} = \overline{(\mathcal{S}(\mathbb{R}^d), \|\cdot\|_{\mathcal{Y}})}$. Moreover,

1. $\mathbb{R}^* K_{\text{rad}} : \mathcal{X}_{\text{Rad}} \rightarrow \mathcal{Y}$ is an **isometric bijection**, with $\mathbb{R}\mathbb{R}^* K_{\text{rad}} = \text{Id}$ on \mathcal{X}_{Rad}
2. $\mathbb{R}^* : \mathcal{X}'_{\text{Rad}} \rightarrow \mathcal{Y}'$ is an **isometric bijection** with $K_{\text{rad}}\mathbb{R}\mathbb{R}^* = \text{Id}$ on $\mathcal{X}'_{\text{Rad}}$.



(Unser, JMLR 2022)

25

Hyper-spherical functions and measures

■ Banach space of hyper-spherical bounded Radon measures

$$\mathcal{M}(\mathbb{R} \times \mathbb{S}^{d-1}) = (C_0(\mathbb{R} \times \mathbb{S}^{d-1}))' \quad \text{where} \quad C_0(\mathbb{R} \times \mathbb{S}^{d-1}) = \overline{(\mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1}), \|\cdot\|_{L_\infty})}$$

$$\text{Null space of } \mathbb{R}^*: \quad \ker(\mathbb{R}^*) = \{g \in \mathcal{M}(\mathbb{R} \times \mathbb{S}^{d-1}) : \langle \mathbb{R}\{\varphi\}, g \rangle_{\text{Rad}} = \langle \varphi, \mathbb{R}^*\{g\} \rangle = 0, \forall \varphi \in \mathcal{S}(\mathbb{R}^d)\}$$

Theorem (Inversion of backprojection operator) (Neumayer-U., Anal. and Appl. 2023)

The **quotient space** $\mathcal{M}_{\text{Rad}} = \mathcal{M}(\mathbb{R} \times \mathbb{S}^{d-1}) / \ker(\mathbb{R}^*)$ is a Banach space that is **isometrically isomorphic** to $\mathcal{M}_{\text{even}}(\mathbb{R} \times \mathbb{S}^{d-1})$. Consequently, $K_{\text{rad}}\mathbb{R}\mathbb{R}^* = \text{Id}$ on $\mathcal{M}_{\text{even}}$ (resp., \mathcal{M}_{Rad}).

$$C_{0,\text{even}}(\mathbb{R} \times \mathbb{S}^{d-1}) = \overline{(\mathcal{S}_{\text{Rad}}, \|\cdot\|_{L_\infty})} \subset C_0(\mathbb{R} \times \mathbb{S}^{d-1})$$

$$L_{p,\text{even}}(\mathbb{R} \times \mathbb{S}^{d-1}) = \overline{(\mathcal{S}_{\text{Rad}}, \|\cdot\|_{L_p})} \subset L_p(\mathbb{R} \times \mathbb{S}^{d-1}), \quad p \in (1, \infty)$$

26

Native space for Radon-domain regularization

■ Regularization functional: $\|L_{\mathbb{R}} f\|_{\mathcal{M}}$ with $L_{\mathbb{R}} = K_{\text{rad}} \mathbb{R} L : \mathcal{M}_{L_{\mathbb{R}}}(\mathbb{R}^d) \rightarrow \mathcal{M}_{\text{Rad}}$

■ Radon-domain \mathcal{M} -norm: $\|g\|_{\mathcal{M}} \triangleq \sup_{\phi \in \mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1}): \|\phi\|_{L_{\infty}} \leq 1} \langle g, \phi \rangle_{\text{Rad}}$,

■ Native space

$$\begin{aligned} \mathcal{M}_{L_{\mathbb{R}}}(\mathbb{R}^d) &= L_{\mathbb{R}}^{\dagger}(\mathcal{M}_{\text{Rad}}) \oplus \mathcal{P}_{n_0} \\ &= \{f = L_{\mathbb{R}}^{\dagger}\{w\} + p_0 : (w, p_0) \in \mathcal{M}_{\text{Rad}} \times \mathcal{P}_{n_0}\}, \end{aligned}$$

■ **Null space** of L and $L_{\mathbb{R}}$: $\mathcal{P} = \mathcal{P}_{n_0} = \text{span}\{m_{\mathbf{k}}\}_{|\mathbf{k}| \leq n_0}$ with $m_{\mathbf{k}}(\mathbf{x}) = \frac{\mathbf{x}^{\mathbf{k}}}{\mathbf{k}!}$

■ **Right-inverse** of $L_{\mathbb{R}}$ on \mathcal{M}_{Rad} : $L_{\mathbb{R}}^{\dagger} = (\text{Id} - \text{Proj}_{\mathcal{P}}) L^{-1} \mathbb{R}^*$

■ **Projector** $\text{Proj}_{\mathcal{P}} : \mathcal{S}'(\mathbb{R}^d) \rightarrow \mathcal{P} : f \mapsto \sum_{|\mathbf{k}| \leq n_0} \langle f, m_{\mathbf{k}}^* \rangle m_{\mathbf{k}}$ with $m_{\mathbf{k}}^* = (-1)^{|\mathbf{k}|} \partial^{\mathbf{k}} \kappa_{\text{iso}} \in \mathcal{S}(\mathbb{R}^d)$

27

Properties of native space

Right-inverse property: $\mathcal{U}' = L_{\mathbb{R}}^{\dagger}(\mathcal{M}_{\text{Rad}}) \Leftrightarrow L_{\mathbb{R}}(\mathcal{U}') = \mathcal{M}_{\text{Rad}}$

Theorem

(Unser FoCM in press)

Let L be an admissible operator with a polynomial null space $\mathcal{P} = \mathcal{P}_{n_0}$ (possibly trivial) of degree n_0 .

1. The space $\mathcal{M}_{L_{\mathbb{R}}}(\mathbb{R}^d) = \mathcal{U}' \oplus \mathcal{P}$ equipped with the composite norm $\|f\|_{\mathcal{M}_{L_{\mathbb{R}}}} = \|L_{\mathbb{R}}\{f\}\|_{\mathcal{M}} + \|\text{Proj}_{\mathcal{P}}\{f\}\|_{\mathcal{P}}$ is **complete** and **isomorphic** to $\mathcal{M}_{\text{Rad}} \times \mathcal{P}$
2. The operators $L_{\mathbb{R}} = K_{\text{rad}} \mathbb{R} L : \mathcal{M}_{L_{\mathbb{R}}} \rightarrow \mathcal{M}_{\text{Rad}}$ and $L_{\mathbb{R}}^{\dagger} = (\text{Id} - \text{Proj}_{\mathcal{P}}) L^{-1} \mathbb{R}^* : \mathcal{M}_{\text{Rad}} \rightarrow L_{\infty, -n_0}(\mathbb{R}^d)$ are continuous and have the following properties:

$$\begin{aligned} \forall w \in \mathcal{M}_{\text{Rad}} : L_{\mathbb{R}} L_{\mathbb{R}}^{\dagger}\{w\} &= w \\ \forall p_0 \in \mathcal{P} : L_{\mathbb{R}}\{p_0\} &= 0 \\ \forall f \in \mathcal{M}_{L_{\mathbb{R}}}(\mathbb{R}^d) : L_{\mathbb{R}}^{\dagger} L_{\mathbb{R}}\{f\} &= (\text{Id} - \text{Proj}_{\mathcal{P}})\{f\} = \text{Proj}_{\mathcal{U}'}\{f\}. \end{aligned}$$

3. Embeddings: $\mathcal{S}(\mathbb{R}^d) \hookrightarrow \mathcal{M}_{L_{\mathbb{R}}}(\mathbb{R}^d) \hookrightarrow L_{\infty, -n_0}(\mathbb{R}^d) \xrightarrow{d} \mathcal{S}'(\mathbb{R}^d)$.

$$L_{\infty, -n_0}(\mathbb{R}^d) = \{f : \mathbb{R}^d \rightarrow \mathbb{R} \text{ s.t. } \sup_{\mathbf{x} \in \mathbb{R}^d} (1 + \|\mathbf{x}\|)^{-n_0} |f(\mathbf{x})| < \infty\}$$

28

Schwartz kernel of pseudo-inverse operator

Adjoint pair of pseudo-inverse operators

$$L_R^\dagger = (\text{Id} - \text{Proj}_{\mathcal{P}})L^{-1}R^* : \mathcal{S}(\mathbb{R} \times \mathbb{S}^{d-1}) \rightarrow \mathcal{S}'(\mathbb{R}^d) \quad (\text{Right-inverse of } L_R)$$

$$L_R^{\dagger*} = RL^{-1*}(\text{Id} - \text{Proj}_{\mathcal{P}'}) : \mathcal{S}'(\mathbb{R}^d) \rightarrow \mathcal{S}'(\mathbb{R} \times \mathbb{S}^{d-1}) \quad (\text{Left-inverse of } L_R^*)$$

Theorem (Generalized impulse response of $L_R^{\dagger*}$)

(Unser FoCM in press)

Let L be an admissible operator with a polynomial null space $\mathcal{P} = \mathcal{P}_{n_0}$ (possibly trivial) of degree n_0 , a frequency profile $\widehat{L}_{\text{rad}} : \mathbb{R} \rightarrow \mathbb{R}$, and $\rho_{\text{rad}}(t) = \mathcal{F}^{-1}\{1/\widehat{L}_{\text{rad}}\}(t)$. Then,

$$\begin{aligned} h(\mathbf{x}_0, (t, \boldsymbol{\xi})) &\triangleq L_R^{\dagger*}\{\delta(\cdot - \mathbf{x}_0)\}(t, \boldsymbol{\xi}) \\ &= \rho_{\text{rad}}(t - \boldsymbol{\xi}^\top \mathbf{x}_0) - \sum_{n=0}^{n_0} \frac{(-\boldsymbol{\xi}^\top \mathbf{x}_0)^n}{n!} (\kappa_{\text{rad}} * \partial^n \rho_{\text{rad}})(t) \end{aligned}$$

with $h(\mathbf{x}_0, \cdot) \in C_0(\mathbb{R} \times \mathbb{S}^{d-1})$ and

$$\sup_{(\mathbf{x}_0, \boldsymbol{\xi}) \in \mathbb{R}^d \times \mathbb{S}^{d-1}} (1 + |\boldsymbol{\xi}^\top \mathbf{x}_0|)^{-n_0} \|h(\mathbf{x}_0; (\cdot, \boldsymbol{\xi}))\|_{L_q(\mathbb{R})} < \infty.$$

for any $q \in [2, \infty]$.

29

Kernel of stable right-inverse of “radonized” Laplacian

Integral representation

$$\Delta_R^\dagger : w \mapsto f(\mathbf{x}) = \int_{\mathbb{R}} \int_{\mathbb{S}^{d-1}} h(\mathbf{x}; t, \boldsymbol{\xi}) w(t, \boldsymbol{\xi}) d\xi dt \quad \text{where}$$

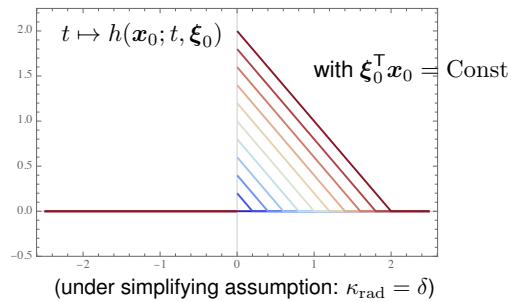
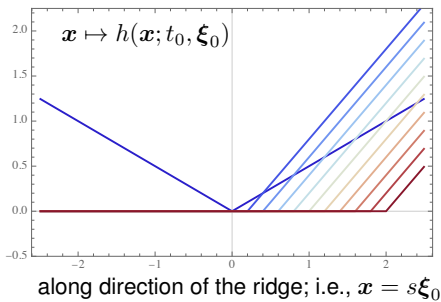
$$h(\mathbf{x}; t, \boldsymbol{\xi}) = \frac{1}{2} |\boldsymbol{\xi}^\top \mathbf{x} - t| - (\kappa_{\text{rad}} * \frac{1}{2} |\cdot|)(t) + (\boldsymbol{\xi}^\top \mathbf{x}) (\kappa_{\text{rad}} * \frac{1}{2} \text{sign})(t),$$

where $\kappa_{\text{rad}} \in \mathcal{S}_{\text{even}}(\mathbb{R})$ is a normalized radial smoothing kernel.

Continuity bound (non-trivial):

$$|\Delta_R^\dagger\{w\}(\mathbf{x})| \leq \|h(\mathbf{x}; \cdot)\|_{L_\infty} \|w\|_{\mathcal{M}} \leq (1 + \|\mathbf{x}\|) \|w\|_{\mathcal{M}}$$

← Dictionary elements
= extreme points of unit ball



30

Admissible (activation, operator) pairs

$\rho_{\text{rad}}(t)$	$\tilde{\rho}_{\text{rad}}(t)$	\hat{L}_{rad} or $\tilde{\hat{L}}_{\text{rad}}$	n_0 (null space)
Exponential			
$e^{- t }$	N/A	$1 + \omega ^2$	-1 (trivial)
Classical sigmoids (Barron 1993)			
	$\frac{\tanh(\frac{t}{2})}{2} = \frac{1}{2} + \frac{1}{1+e^{-t}}$	$\frac{\sinh(\pi\omega)}{\pi}$	0 (bias)
	$\frac{\arctan(t)}{\pi}$	$\omega e^{ \omega }$	0
Ridge splines (of degree $n \in \mathbb{N}$)			
$\frac{1}{2} t $ (or ReLU)	$t \log t $	$ \omega ^2$	1 (affine)
$\propto t^{2n} \log t $	$\frac{\text{sign}(t) t ^{2n}}{(2n)!}$	$ \omega ^{2n+1}$	$2n \geq 2$ (even)
$\frac{1}{2} \frac{ t ^{2n+1}}{(2n+1)!}$	$\propto t^{2n+1} \log t $	$ \omega ^{2n+2}$	$2n + 1 \geq 1$ (odd)
Fractional splines (degree $\alpha \in \mathbb{R}^+ \setminus \mathbb{N}$)			
$\frac{ t ^\alpha \sin(\frac{\alpha\pi}{2})}{\pi\Gamma(\alpha)}$	$\frac{\text{sign}(t) t ^\alpha \cos(\frac{\alpha\pi}{2})}{\pi\Gamma(\alpha)}$	$ \omega ^{\alpha+1}$	$\lceil \alpha \rceil$

Table 1: Examples of admissible symmetric and anti-symmetric activation functions with their corresponding regularization operator. The anti-

Anti-symmetric extension

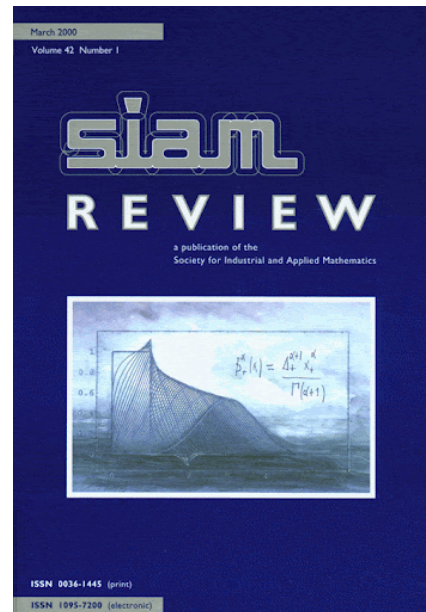
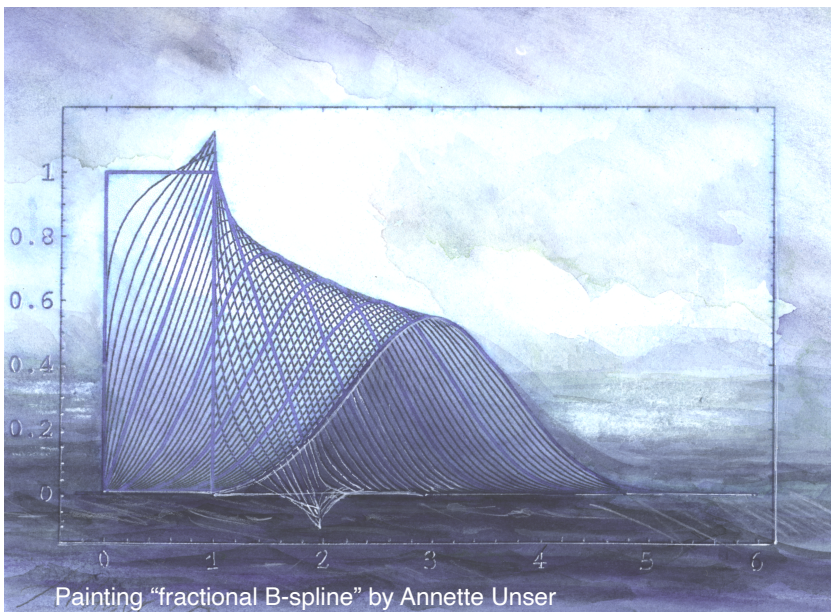
Replace symmetric filtering operator K_{rad} by anti-symmetric $\tilde{K}_{\text{rad}} = H_{\text{rad}}K_{\text{rad}}$.

H_{rad} : Hilbert transform

Bottom line

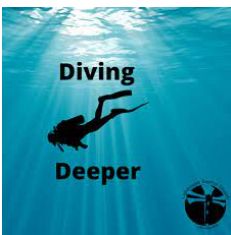
Choice of L (resp., L_R) and symmetry fixes activation $\sigma = \rho_{\text{rad}}$ and vice versa

M. Unser, T. Blu, "Fractional Splines and Wavelets," *SIAM Review*, March 2000.



CONCLUSION: Return of the spline

- Foundations of **functional learning**
 - Functional optimization in Banach spaces (enabled by representer theorems)
 - **Hilbert spaces**: the tools of classical ML
 - **Non-reflexive Banach spaces**: for sparsity-promoting regularization (e.g., CS)
 - **Isotropy + Radon transform**: The key for obtaining pointwise nonlinearities
- Splines and machine learning
 - Traditional kernel methods are closely related to splines ... and the same holds true for ReLU nets ...
 - Sparsity-promoting regularization offer promising perspectives
 - Radon-domain regularization \Rightarrow *Unifies* **Shallow neural nets** and **RBF methods**



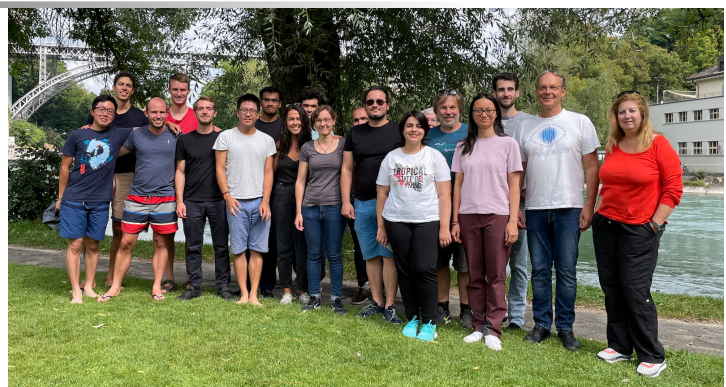
- Functional composition = **hierarchical splines**
 - Deep ReLU neural nets are **high-dimensional** piecewise-linear **splines**
 - **Free-form** activations with $TV^{(2)}$ -regularization \Rightarrow **Deep splines**

33

ACKNOWLEDGMENTS

Many thanks to (former) members of EPFL's Biomedical Imaging Group

- Prof. Rahul Parhi
- Dr. Shayan Aziznejad
- Dr. Julien Fageot
- Prof. John Paul Ward
- Dr. Thomas Debarre
- Dr. Mike McCann
- Dr. Harshit Gupta
- Prof. Kyong Jin
- Dr. Fangshu Yang
- Dr. Emrah Bostan
- Prof. Ulugbek Kamilov
-



and collaborators ...

- Prof. Demetri Psaltis
- Prof. Marco Stampanoni
- Prof. Carlos-Oscar Sorzano
- Prof. Luigi Ambrosio
-



34

References

<http://bigwww.epfl.ch/>

■ Sparse adaptive splines

- M. Unser, J. Fageot, J.P. Ward, "Splines Are Universal Solutions of Linear Inverse Problems with Generalized-TV Regularization," *SIAM Review*, vol. 59, No. 4, pp. 769-793, 2017.
- T. Debarre, Q. Denoyelle, M. Unser, J. Fageot, "Sparsest Continuous Piecewise-Linear Representation of Data," *Journal of Computational and Applied Mathematics*, vol. 406, paper no. 114044, pp. 1-30, 2022.

■ Representer theorems

- M. Unser, "A Representer Theorem for Deep Neural Networks," *Journal of Machine Learning Research*, vol. 20, no. 110, pp. 1-30, Jul. 2019.
- M. Unser, "A Unifying Representer Theorem for Inverse Problems and Machine Learning," *Foundations of Computational Mathematics*, vol. 21, pp. 941–960, 2021.
- M. Unser, S. Aziznejad, "Convex optimization in sums of Banach spaces," *Applied and Computational Harmonic Analysis*, vol. 56, no. 1, pp. 1-25, 2022.

■ Neural networks and the Radon transform

- M. Unser, "Ridges, Neural Networks, and the Radon Transform", *Journal of Machine Learning Research*, vol. 24, no. 37, pp. 1-33, 2023.
- S. Neumayer, M. Unser, "Explicit Representations for Banach Subspaces of Lizorkin Distributions," *Analysis and Applications* vol. 21, no. 5, pp. 1223–1250, September 2023. Preprint arXiv:2203.05312 [math.FA]
- M. Unser, "Unifying Variational Formulation of Supervised Learning: From Kernel Methods to Neural Networks," *Foundations of Computational Mathematics* (in press). Preprint arXiv:2206.14625 [cs.LG]