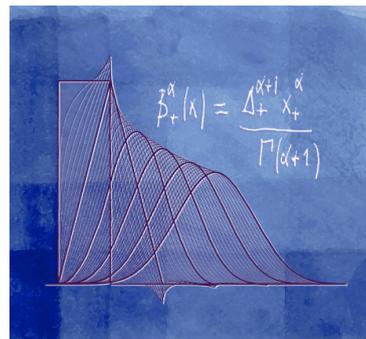


New representer theorems for inverse problems and machine learning

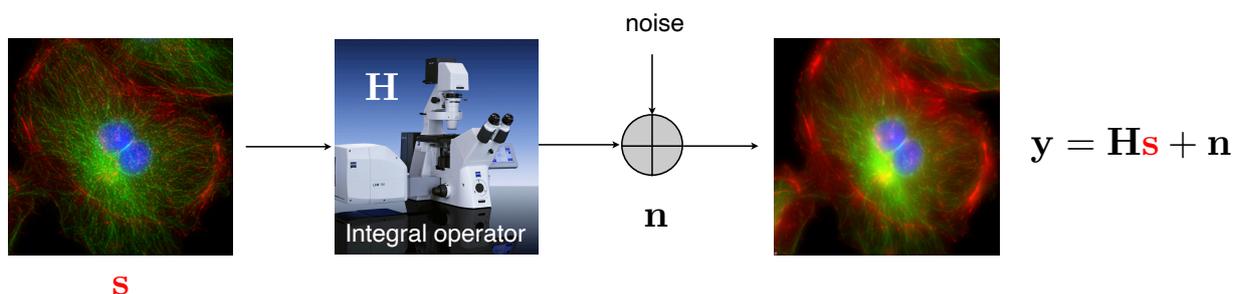
Michael Unser
 Biomedical Imaging Group
 EPFL, Lausanne, Switzerland



EPFL CIS – RIKEN AIP Seminar Series (virtual), March 2, 2022

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

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

Supervised learning as a (linear) inverse problem

but an infinite-dimensional one ...

Given the data points $(\mathbf{x}_m, y_m) \in \mathbb{R}^{N+1}$, 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

(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)

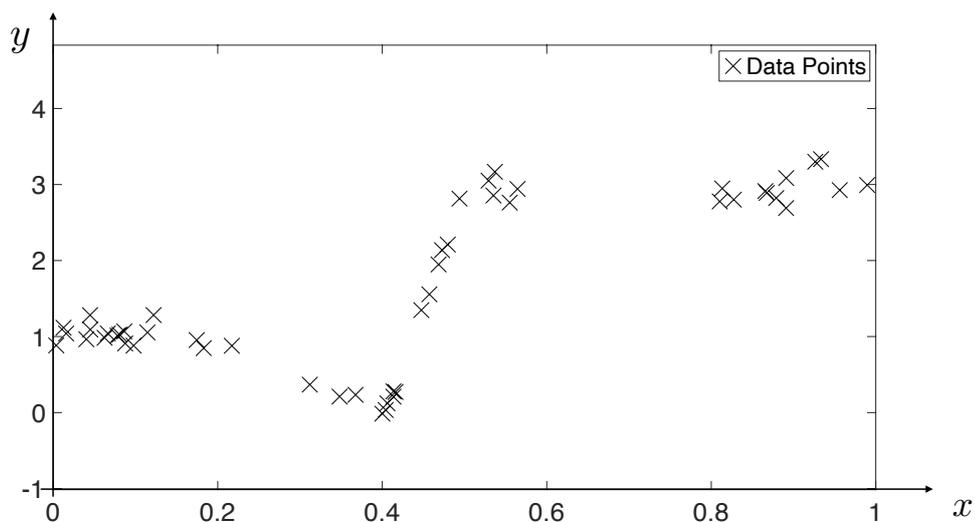
$$f_{\text{RKHS}} = \arg \min_{f \in \mathcal{H}} \left(\sum_{m=1}^M |y_m - f(\mathbf{x}_m)|^2 + \lambda R(f) \right) \quad \text{with} \quad R(f) = \|f\|_{\mathcal{H}}^2$$

\Rightarrow kernel estimator

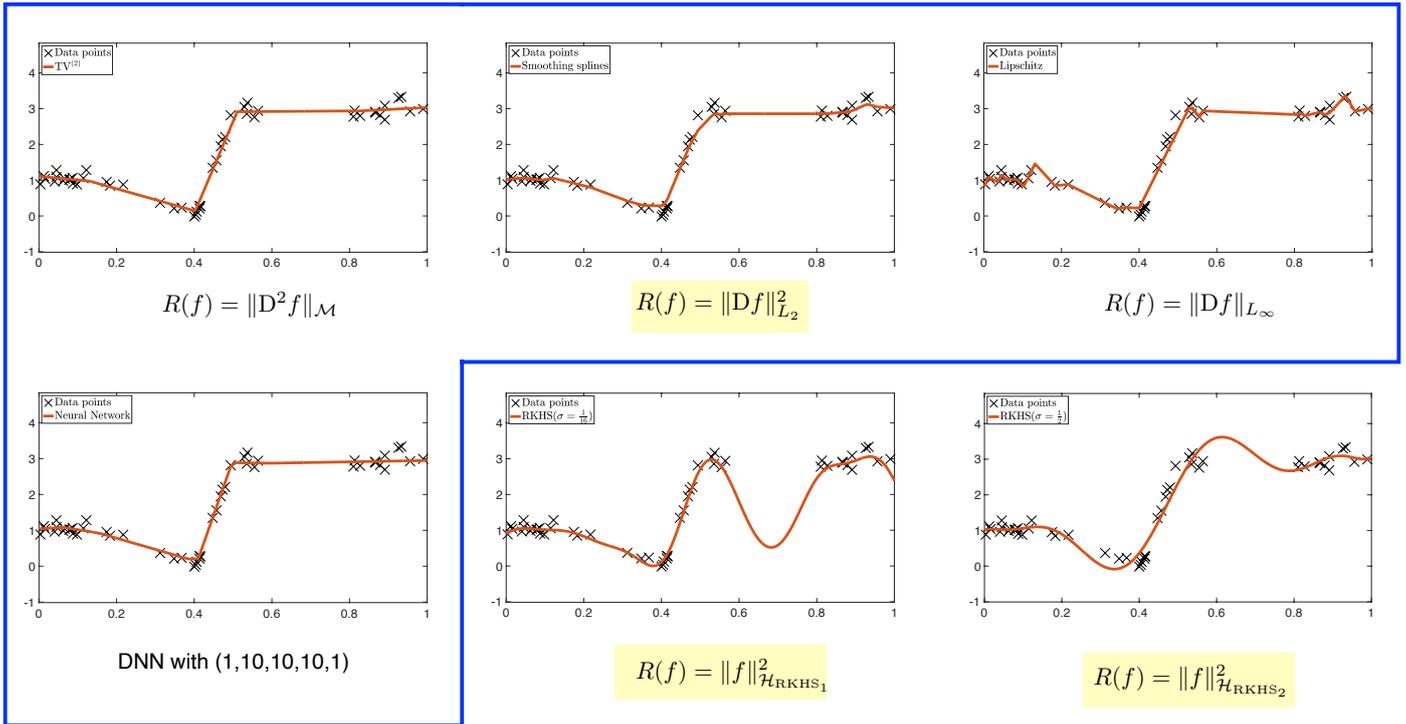
(Wahba 1990; Schölkopf 2001)

3

Can you learn the map $y = f(x)$?



4



5

OUTLINE

- **Introduction** ✓
 - Learning as an inverse problem
 - Teaser: Search of the best learner
- **Foundations of functional learning**
 - Banach spaces and duality mappings
 - **Unifying representer theorem** *NEW*
- **From classical to modern regularization-based techniques**
 - Kernel methods
 - Smoothing splines
 - **Sparse adaptive splines**
- **Deep neural networks vs. deep splines**
 - Continuous piecewise linear (CPWL) functions / splines
 - **Representer theorem for deep neural networks**



6

General notion of Banach space



Stefan Banach (1892-1945)

Normed space: vector space \mathcal{X} equipped with a norm $\|\cdot\|_{\mathcal{X}}$

Convergent sequence of functions (φ_i) in \mathcal{X} :

$$\lim_i \varphi_i = \varphi; \quad \text{i.e., } \lim_i \|\varphi - \varphi_i\|_{\mathcal{X}} = 0$$

Definition

A Banach space is a **complete normed** space \mathcal{X} ; that is, such that $\lim_i \varphi_i = \varphi \in \mathcal{X}$ for any convergent sequence (φ_i) in \mathcal{X} .

■ Generality of the concept

- Linear space of vectors $\mathbf{u} = (u_1, \dots, u_N) \in \mathbb{R}^N$
- Linear space of functions $u : \mathbb{R}^d \rightarrow \mathbb{R}$
- Space of linear functional $u : \mathcal{X} \rightarrow \mathbb{R}$
- Linear space of vector-valued functions $\mathbf{u} = (u_1, \dots, u_N) : \mathbb{R}^d \rightarrow \mathbb{R}^N$
- Linear space $\mathcal{L}(\mathcal{X}, \mathcal{Y})$ of bounded operators $U : \mathcal{X} \rightarrow \mathcal{Y}$

7

Dual of a Banach space

Dual of the Banach space $(\mathcal{X}, \|\cdot\|_{\mathcal{X}})$:

$\mathcal{X}' =$ space of linear functionals $g : f \mapsto \langle g, f \rangle \triangleq g(f) \in \mathbb{R}$ that are continuous on \mathcal{X}

\mathcal{X}' is a Banach space equipped with the **dual norm**:

$$\|g\|_{\mathcal{X}'} = \sup_{f \in \mathcal{X} \setminus \{0\}} \left(\frac{\langle g, f \rangle}{\|f\|_{\mathcal{X}}} \right)$$

■ Generic duality bound

$$\Rightarrow \|g\|_{\mathcal{X}'} \geq \frac{|\langle g, f \rangle|}{\|f\|_{\mathcal{X}}}, \quad f \neq 0$$

For any $f \in \mathcal{X}, g \in \mathcal{X}'$: $|\langle g, f \rangle| \leq \|g\|_{\mathcal{X}'} \|f\|_{\mathcal{X}}$

- Duals of L_p spaces: $(L_p(\mathbb{R}^d))' = L_{p'}(\mathbb{R}^d)$ with $\frac{1}{p} + \frac{1}{p'} = 1$ for $p \in (1, \infty)$

Hölder inequality: $|\langle f, \varphi \rangle| \leq \int_{\mathbb{R}^d} |f(\mathbf{r})\varphi(\mathbf{r})| \, d\mathbf{r} \leq \|f\|_{L_p} \|\varphi\|_{L_{p'}}$

8

Riesz conjugate for Hilbert spaces

- Duality bound for Hilbert spaces (equivalent to Cauchy-Schwarz inequality)

$$\text{For all } (u, v) \in \mathcal{H} \times \mathcal{H}': \quad |\langle u, v \rangle| \leq \|u\|_{\mathcal{H}} \|v\|_{\mathcal{H}'}$$



Frigyes Riesz (1880-1956)

- Definition

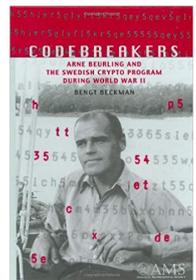
The **Riesz conjugate** of $u \in \mathcal{H}$ is the unique element $u^* \in \mathcal{H}'$ such that

$$\langle u, u^* \rangle = \langle u, u \rangle_{\mathcal{H}} = \|u\|_{\mathcal{H}}^2 = \|u\|_{\mathcal{H}} \|u^*\|_{\mathcal{H}'} \quad (\text{sharp duality bound})$$

- Properties

- Norm preservation: $\|u\|_{\mathcal{H}} = \|u^*\|_{\mathcal{H}'}$ +
- $u^* = \mathcal{R}^{-1}\{u\}$ (inverse Riesz map) (isometry)
- Invertibility: $u = (u^*)^* = \mathcal{R}\{u^*\}$ (\mathcal{H}')' = \mathcal{H} (reflexivity)
- Linearity: $(u_1 + u_2)^* = u_1^* + u_2^*$

Generalization: Duality mapping



Arne Beurling (1905-1986)

Definition

Let $(\mathcal{X}, \mathcal{X}')$ be a dual pair of Banach spaces. Then, the elements $f^* \in \mathcal{X}'$ and $f \in \mathcal{X}$ form a **conjugate pair** if

- $\|f^*\|_{\mathcal{X}'} = \|f\|_{\mathcal{X}}$ (norm preservation), and
- $\langle f^*, f \rangle_{\mathcal{X}' \times \mathcal{X}} = \|f^*\|_{\mathcal{X}'} \|f\|_{\mathcal{X}}$ (sharp duality bound).

For any given $f \in \mathcal{X}$, the set of admissible conjugates defines the **duality mapping**

$$J(f) = \{f^* \in \mathcal{X}' : \|f^*\|_{\mathcal{X}'} = \|f\|_{\mathcal{X}} \text{ and } \langle f^*, f \rangle_{\mathcal{X}' \times \mathcal{X}} = \|f^*\|_{\mathcal{X}'} \|f\|_{\mathcal{X}}\},$$

which is a non-empty subset of \mathcal{X}' . Whenever the duality mapping is single-valued (for instance, when \mathcal{X}' is strictly convex), one also defines the duality operator $J_{\mathcal{X}} : \mathcal{X} \rightarrow \mathcal{X}'$, which is such that $f^* = J_{\mathcal{X}}(f)$.

(Beurling-Livingston, 1962)

Properties of duality mapping

Theorem

Let $(\mathcal{X}, \mathcal{X}')$ be a dual pair of Banach spaces. Then, the following holds:

1. Every $f \in \mathcal{X}$ admits at least one conjugate $f^* \in \mathcal{X}'$.
2. For every $f \in \mathcal{X}$, the set $J(f)$ is convex and weak-* closed in \mathcal{X}' .
3. The duality mapping is **single-valued** if \mathcal{X}' is **strictly convex**; the latter condition is also necessary if \mathcal{X} is reflexive.

\mathcal{X} is **strictly convex** if, for all $f_1, f_2 \in \mathcal{X}$ such that $\|f_1\|_{\mathcal{X}} = \|f_2\|_{\mathcal{X}} = 1$ and $f_1 \neq f_2$, one has $\|\lambda f_1 + (1 - \lambda)f_2\|_{\mathcal{X}} < 1$ for any $\lambda \in (0, 1)$.

\mathcal{X} is **reflexive** if $\mathcal{X}'' = \mathcal{X}$.

11

Mother of all representer theorems

$$\arg \min_{f \in \mathcal{X}'} E(\mathbf{y}, \boldsymbol{\nu}(f)) + \psi(\|f\|_{\mathcal{X}'})$$



Lausanne, Christmas 2018

Mathematical assumptions:

- $(\mathcal{X}, \mathcal{X}')$ is a dual pair of Banach spaces.
- $\mathcal{N}_{\boldsymbol{\nu}} = \text{span}\{\nu_m\}_{m=1}^M \subset \mathcal{X}$ with the ν_m being linearly independent.
- $\boldsymbol{\nu} : \mathcal{X}' \rightarrow \mathbb{R}^M : f \mapsto (\langle \nu_1, f \rangle, \dots, \langle \nu_M, f \rangle)$ is the **linear measurement operator** (it is weak* **continuous** on \mathcal{X}' because $\nu_1, \dots, \nu_M \in \mathcal{X}$).
- $E : \mathbb{R}^M \times \mathbb{R}^M \rightarrow \mathbb{R}^+$ is a strictly-convex loss functional.
- $\psi : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is some arbitrary strictly-increasing convex function.

12

General representer theorem

Theorem

For any fixed $\mathbf{y} \in \mathbb{R}^M$, the solution set of the **generic** optimization problem

$$S = \arg \min_{f \in \mathcal{X}'} E(\mathbf{y}, \nu(f)) + \psi(\|f\|_{\mathcal{X}'})$$

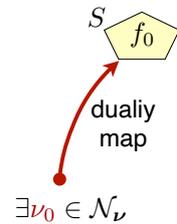
is **non-empty, convex** and weak*-compact, and all solutions $f_0 \in S \subset \mathcal{X}'$ are $(\mathcal{X}', \mathcal{X})$ -conjugate of a **common** $\nu_0 \in \mathcal{N}_\nu = \text{span}\{\nu_m\}_{m=1}^M \subset \mathcal{X}$.

The parametric form of the solution depends on the space type.

- 1) If \mathcal{X}' is a **Hilbert space** and ψ is strictly convex, then the solution is unique and it admits a **linear expansion** with coefficients $(a_m) \in \mathbb{R}^M$

$$f_0 = \sum_{m=1}^M a_m \varphi_m,$$

where $\varphi_m = J_{\mathcal{X}}\{\nu_m\} \in \mathcal{X}'$ with $J_{\mathcal{X}}$ the Riesz map $\mathcal{X} \rightarrow \mathcal{X}'$.



(Unser, FoCM 2021)

13

General representer theorem (Cont'd)

- 2) If \mathcal{X}' is a **strictly convex Banach space** and ψ is strictly convex, then the solution is unique and it admits the **representation** with $(a_m) \in \mathbb{R}^M$

$$f_0 = J_{\mathcal{X}} \left\{ \sum_{m=1}^M a_m \nu_m \right\},$$

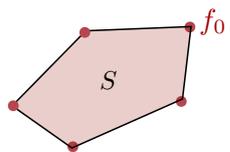
where $J_{\mathcal{X}}$ is the (nonlinear) duality operator $\mathcal{X} \rightarrow \mathcal{X}'$.

- 3) Otherwise, when \mathcal{X}' is **not strictly convex**, the solution set S is the convex hull of its **extreme points**, which can all be expressed as

$$f_0 = \sum_{k=1}^{K_0} c_k e_k,$$

for some $K_0 \leq M$, $c_1, \dots, c_{K_0} \in \mathbb{R}$, where $e_1, \dots, e_{K_0} \in \mathcal{X}'$ are some extreme points of the unit ball $B_{\mathcal{X}'} = \{x \in \mathcal{X}' : \|x\|_{\mathcal{X}'} \leq 1\}$.

(Unser, FoCM 2021)



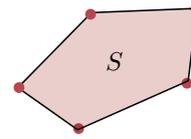
(Boyer-Chambolle-De Castro-Duval-De Gournay-Weiss, arXiv:1806.09810, 2019)

14

Extreme points

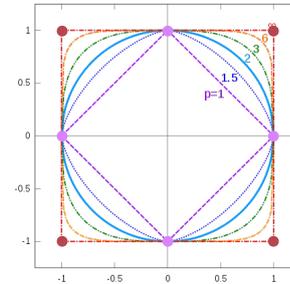
■ Definition

Let S be a convex set. Then, the point $x \in S$ is **extreme** if it cannot be expressed as a (non-trivial) convex combination of any other points in S .



■ Extreme points of unit ball in $\ell_p(\mathbb{Z})$

- $\ell_\infty(\mathbb{Z})$: $e_k[n] = \pm 1$
- $\ell_1(\mathbb{Z})$: $e_k = \pm \delta[\cdot - n_k]$ (Kronecker impulse)
- $\ell_p(\mathbb{Z})$ with $p \in (1, \infty)$: $e_k = u / \|u\|_{\ell_p}$ for any $u \in \ell_p(\mathbb{Z})$



⇒ sparse !!!

Definition of **strictly convexity** of a Banach space: all boundary points are extreme !!!

OUTLINE

- Introduction ✓
- Foundations of functional learning ✓
- From classical to **modern** regularization-based techniques
 - Learning in RKHS
 - Kernel methods of ML
 - Smoothing splines
 - **Sparse kernel methods**
 - **Sparse adaptive splines**
 - **Lipchitz splines**
- Deep neural networks vs. deep splines

NEW

1. Learning in reproducing kernel Hilbert space

Definition

A Hilbert space \mathcal{H} of functions on \mathbb{R}^d is called a **reproducing kernel Hilbert space** (RKHS) if $\delta(\cdot - \mathbf{x}_0) \in \mathcal{H}'$ for any $\mathbf{x}_0 \in \mathbb{R}^d$. The corresponding unique **Hilbert conjugate** $h(\cdot, \mathbf{x}_0) = (\delta(\cdot - \mathbf{x}_0))^* \in \mathcal{H}$ when indexed by \mathbf{x}_0 is called the **reproducing kernel** of \mathcal{H} .

Learning problem

Given the data $(\mathbf{x}_m, y_m)_{m=1}^M$ with $\mathbf{x}_m \in \mathbb{R}^d$, find the function $f_0 : \mathbb{R}^d \rightarrow \mathbb{R}$ s.t.

$$f_0 = \arg \min_{f \in \mathcal{H}} \left(\sum_{m=1}^M E_m(y_m, f(\mathbf{x}_m)) + \psi(\|f\|_{\mathcal{H}}) \right)$$

- $E_m : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ (strictly convex)
- $\psi : \mathbb{R} \rightarrow \mathbb{R}^+$ (strictly increasing and convex)

17

Learning in RKHS (Cont'd)

Special case of general representer theorem

- $\mathcal{X} = \mathcal{H}'$, $\mathcal{X}' = \mathcal{H}'' = \mathcal{H}$ (all Hilbert spaces are reflexive)
- $\nu_m = \delta(\cdot - \mathbf{x}_m)$ (Dirac sampling functionals)

- Additive loss: $E(\mathbf{y}, \mathbf{z}) = \sum_{m=1}^M E_m(y_m, z_m)$

specific of ML

Key observation

Reproducing kernel = Schwartz kernel of **Riesz map**

$$\mathbb{R} = J_{\mathcal{H}'} : \mathcal{H}' \rightarrow \mathcal{H} : \nu \mapsto \int_{\mathbb{R}^d} h(\cdot, \mathbf{y}) \nu(\mathbf{y}) d\mathbf{y} \quad \Rightarrow \quad \varphi_m = J_{\mathcal{H}'} \{ \delta(\cdot - \mathbf{x}_m) \} = h(\cdot, \mathbf{x}_m)$$

Implied form of unique solution = linear kernel expansion

$$f_0(\mathbf{x}) = \sum_{m=1}^M a_m \varphi_m(\mathbf{x}) = \sum_{m=1}^M a_m h(\mathbf{x}, \mathbf{x}_m)$$

(Schölkopf representer theorem, 2001)

18

2. Regularization with a LSI operator = kernel methods of ML

- Quadratic Tikhonov regularization functional

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

(Poggio-Girosi 1990)

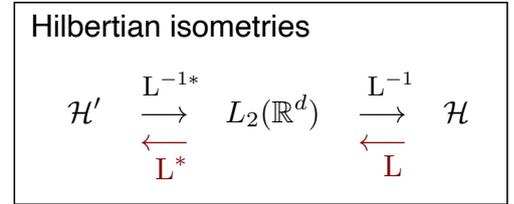
L: Linear shift-invariant (LSI), **invertible regularization operator**

$\widehat{L}(\omega)$: **frequency response** of L

- Key observation

Reproducing kernel = Impulse response of $L^{-1}L^{-1*} = (L^*L)^{-1}$

$$\nu^* = J_{\mathcal{H}'}\{\nu\} = h * \nu \quad \text{where} \quad h = \mathcal{F}^{-1}\left\{\frac{1}{|\widehat{L}(\omega)|^2}\right\} \in L_1(\mathbb{R}^d)$$



- Parametric form of solution = expansion of **kernels centered on data points**

$$f_0(\mathbf{x}) = \sum_{m=1}^M a_m J_{\mathcal{H}'}\{\delta(\cdot - \mathbf{x}_m)\}(\mathbf{x}) = \sum_{m=1}^M a_m h(\mathbf{x} - \mathbf{x}_m)$$

19

3. Smoothing splines

$$f_0 = \arg \min_{f: \mathbb{R} \rightarrow \mathbb{R}} \left(\sum_{m=1}^M |f(x_m) - y_m|^2 + \lambda \int_{\mathbb{R}} \left| \frac{df(x)}{dx} \right|^2 dx \right)$$

(Schoenberg 1964; de Boor 1966)

- Smoothness regularization (spline semi-norm)

$$R(f) = \|Df\|_{L_2}^2 \quad \text{with} \quad D = \frac{d}{dx}; \quad \text{Null space} : \mathcal{N}_D = \{p(x) = a_0 : a_0 \in \mathbb{R}\}$$

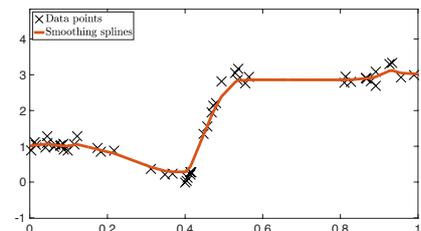
- Direct-sum RKHS topology: $L_{2,D}(\mathbb{R}) = \mathcal{H}_D \oplus \mathcal{N}_D$

D has a unique inverse only if one factors out the null space

$$\text{Impulse response of } (D^*D)^{-1}: \quad h(x) = \mathcal{F}^{-1}\left\{\frac{1}{|\omega|^2}\right\}(x) = \frac{1}{2}|x|$$

- Solution = linear spline with knots at x_1, \dots, x_M

$$f_0(x) = a_0 + \sum_{m=1}^M a_m |x - x_m|$$



20

4. Sparse kernel expansions

- Sparsity-promoting regularization functional

~~$$R(f) = \|Lf\|_{L_1} = \int_{\mathbb{R}^N} |Lf(x)| dx$$~~

L: Linear shift-invariant (LSI), **invertible regularization operator**

$\hat{L}(\omega)$: **frequency response** of L



Banach isometry

$$L_1(\mathbb{R}^d) \begin{array}{c} \xrightarrow{L^{-1}} \\ \xleftarrow{L} \end{array} L_{1,L}(\mathbb{R}^d)$$

Theoretical roadblock: The general representer theorem does not apply because **there exists no predual space** \mathcal{X} such that $L_1(\mathbb{R}^d) = \mathcal{X}'$.

The optimization problem is ill-defined and does not admit a solution !

21

Proper continuous counterpart of ℓ_1 -norm

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

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



Johann Radon (1887-1956)

- 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 $\mathcal{M}(\mathbb{R}^d)$: $e_k = \pm \delta(\cdot - \tau_k)$ with $\tau_k \in \mathbb{R}^d$

22

4. Sparse kernel expansions (2nd attempt)

■ Sparsity-promoting regularization functional

$$R(f) = \|Lf\|_{\mathcal{M}} = \sup_{\varphi \in C_0(\mathbb{R}^d): \|\varphi\|_{L^\infty} \leq 1} \langle Lf, \varphi \rangle$$

L: Linear shift-invariant (LSI), **invertible regularization operator**

$\hat{L}(\omega)$: **frequency response** of L

Impulse response of L^{-1} : $h = \mathcal{F}^{-1} \left\{ \frac{1}{\hat{L}(\omega)} \right\} \in L_1(\mathbb{R}^d)$

Banach isometry

$$\mathcal{M}(\mathbb{R}^d) \begin{array}{c} \xrightarrow{L^{-1}} \\ \xleftarrow{L} \end{array} \mathcal{M}_L(\mathbb{R}^d)$$

Extreme points: $e_k = L^{-1}\{\delta(\cdot - \tau_k)\}$

Corollary (3rd case of representer theorem)

(Aziznejad-U., SIAM 2021)

The extreme points f_0 of $S = \arg \min_{f \in \mathcal{M}_L(\mathbb{R}^d)} \left(\sum_{m=1}^M E_m(y_m, f(\mathbf{x}_m)) + \lambda \|Lf\|_{\mathcal{M}} \right)$ can all be expressed as

$$f_0(\mathbf{x}) = \sum_{k=1}^{K_0} a_k h(\mathbf{x} - \tau_k)$$

for some $K_0 \leq M$, $\tau_1, \dots, \tau_{K_0} \in \mathbb{R}^d$ and $\mathbf{a} = (a_k) \in \mathbb{R}^{K_0}$. Moreover, $\|Lf_0\|_{\mathcal{M}} = \sum_{k=1}^{K_0} |a_k| = \|\mathbf{a}\|_{\ell_1}$.

23

5. Sparse adaptive spline

$$f_0 = \arg \min_{f \in \mathcal{M}_{D^2}(\mathbb{R})} \left(\sum_{m=1}^M |f(x_m) - y_m|^2 + \lambda \|D^2 f\|_{\mathcal{M}} \right) \quad (\text{Mammen 1997; Unser 2017})$$

■ Sparsity-promoting regularization

$$R(f) = \|D^2 f\|_{\mathcal{M}} \quad \text{Null space: } \mathcal{N}_{D^2} = \{p(x) = b_0 + b_1 x : b_0, b_1 \in \mathbb{R}\}$$

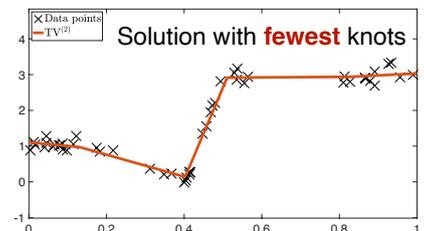
■ Direct-sum Banach topology: $\mathcal{M}_{D^2}(\mathbb{R}) = \mathcal{U}_{D^2} \oplus \mathcal{N}_{D^2}$

D^2 has a unique invertise only if one factors out the null space

Impulse response of D^{-2} (two-fold integrator): $h(x) = (x)_+ = \text{ReLU}(x)$

■ Solution = linear spline with (few) **adaptive** knots at $\tau_1, \dots, \tau_{K_0}$

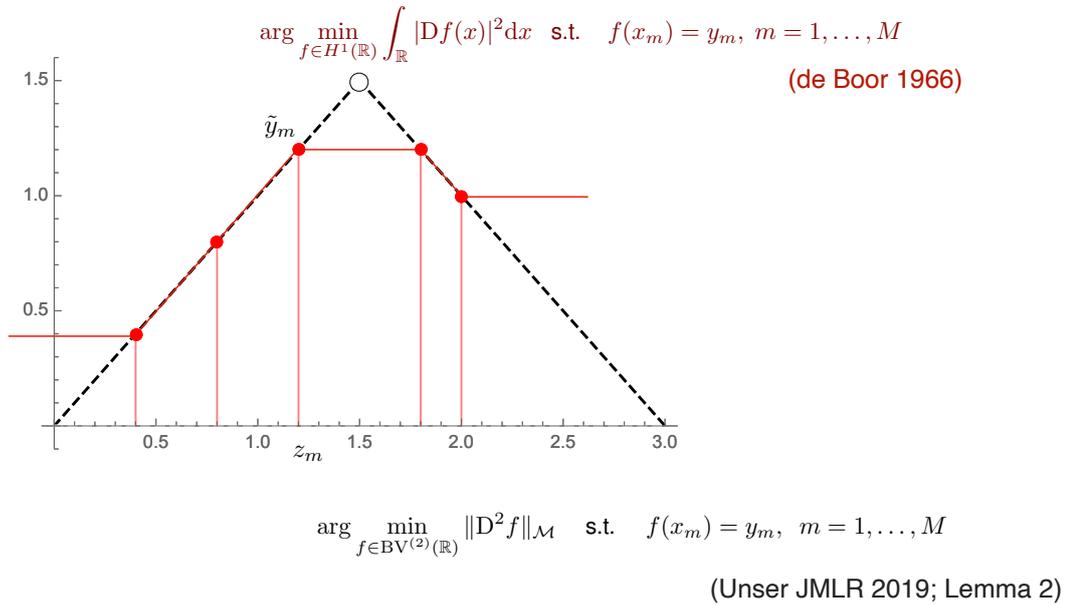
$$f_0(x) = b_0 + b_1 x + \sum_{k=1}^{K_0} a_k (x - \tau_k)_+$$



(Debarre arXiv 2020)

24

Comparison of linear interpolators



6. Lipschitz splines

$$f_0 = \arg \min_{f \in W_{\infty}^1(\mathbb{R})} \left(\sum_{m=1}^M |f(x_m) - y_m|^2 + \lambda \|Df\|_{L_{\infty}} \right)$$

- Lipschitz boundedness constraint

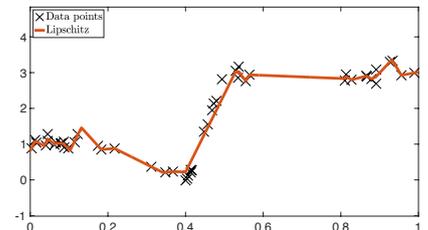
$$R(f) = \|Df\|_{L_{\infty}} \quad \text{Null space} : \mathcal{N}_D = \{p(x) = b_0 : b_0 \in \mathbb{R}\}$$

Extreme points of unit ball in $L_{\infty}(\mathbb{R})$: e_k such that $e_k(x) = \pm 1$

- Direct-sum Banach topology: $W_{\infty}^1(\mathbb{R}) = \mathcal{U}_D \oplus \mathcal{N}_D$

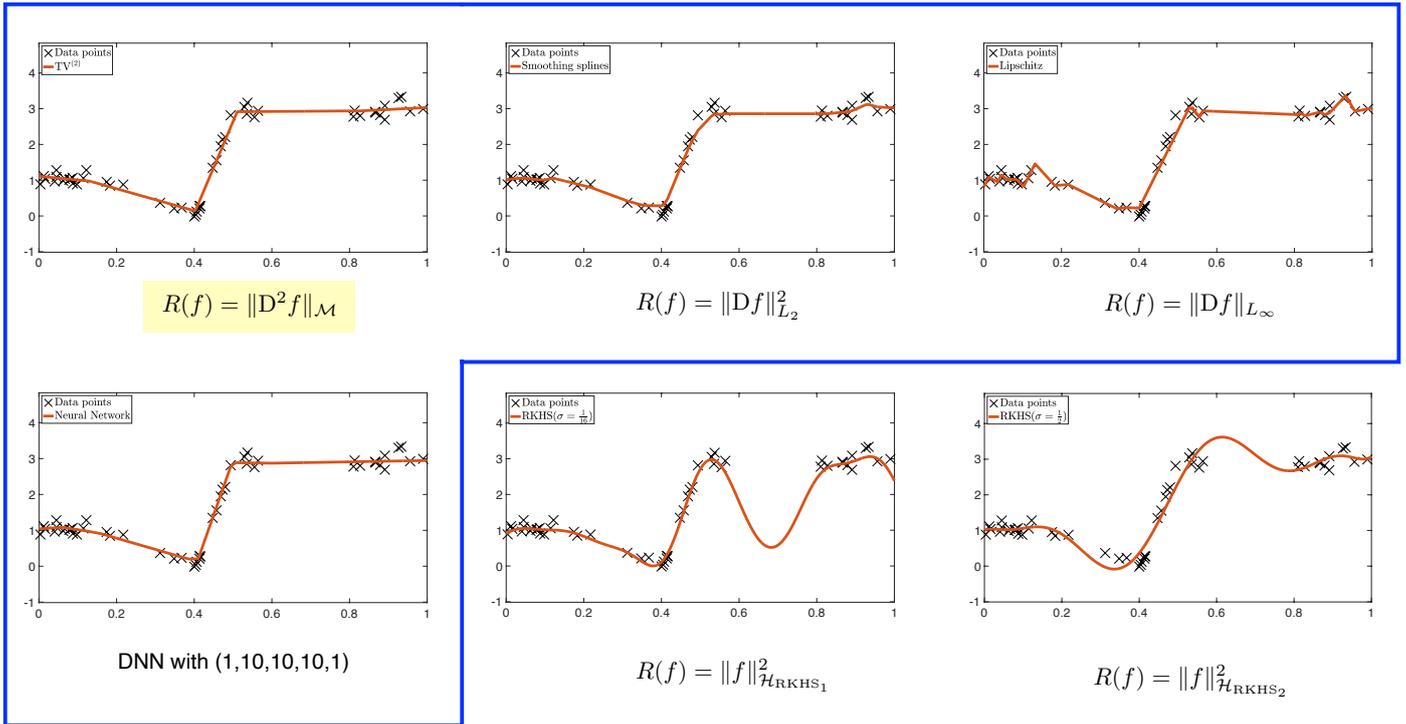
D has a unique inverse only if one factors out the null space

$$u_k = D^{-1}e_k(x) = \int_{-\infty}^x e_k(t) dt + C_k : \text{linear spline with binary slope } (\pm 1)$$



- Solution = **linear spline** with with many oscillations (non-unique)

$$f_0(x) = b_0 + \sum_{k=1}^{K_0} a_k u_k(x) \quad \text{(Aziznejad et al., ArXiv 2022)}$$



27

OUTLINE

- Introduction ✓
- Foundations of functional learning ✓
- From classical to modern regularization-based techniques ✓
- **Deep neural networks vs. deep splines**
 - Background
 - Continuous piecewise linear (CPWL) functions / splines
 - Variational formulation of shallow nets
 - Representer theorem for deep neural networks



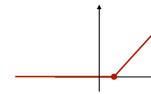
28

Deep neural networks and splines

■ Preferred choice of activation function: ReLU

- ReLU works nicely with dropout / ℓ_1 -regularization
- Networks with hidden ReLU are easier to train
- State-of-the-art performance

$$\text{ReLU}(x; b) = (x - b)_+$$



(Glorot *ICAI*S 2011)

(LeCun-Bengio-Hinton *Nature* 2015)

■ Deep nets as Continuous PieceWise-Linear maps

- ReLU \Rightarrow CPWL
- CPWL \Rightarrow Deep ReLU network

(Montufar *NIPS* 2014)

(Strang *SIAM News* 2018)

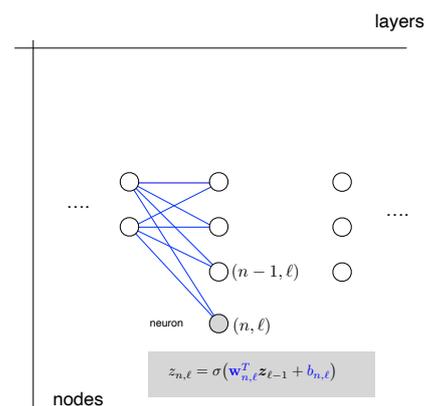
■ Deep ReLU nets = hierarchical splines

- ReLU is a piecewise-linear spline

(Poggio-Rosasco 2015)

Feedforward deep neural network

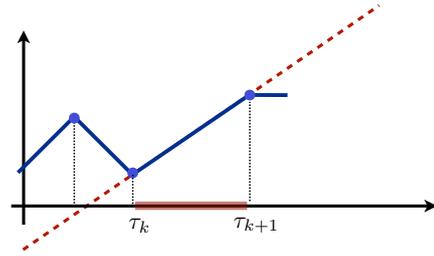
- Layers: $\ell = 1, \dots, L$
- Deep structure descriptor: (N_0, N_1, \dots, N_L)
- Neuron or node index: (n, ℓ) , $n = 1, \dots, N_\ell$
- Activation function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ (ReLU)
- Linear step: $\mathbb{R}^{N_{\ell-1}} \rightarrow \mathbb{R}^{N_\ell}$
 $f_\ell : \mathbf{x} \mapsto \mathbf{f}_\ell(\mathbf{x}) = \mathbf{W}_\ell \mathbf{x} + \mathbf{b}_\ell$
- Nonlinear step: $\mathbb{R}^{N_\ell} \rightarrow \mathbb{R}^{N_\ell}$
 $\sigma_\ell : \mathbf{x} \mapsto \sigma_\ell(\mathbf{x}) = (\sigma(x_1), \dots, \sigma(x_{N_\ell}))$



Learned

$$\mathbf{f}_{\text{deep}}(\mathbf{x}) = (\sigma_L \circ \mathbf{f}_L \circ \sigma_{L-1} \circ \dots \circ \sigma_2 \circ \mathbf{f}_2 \circ \sigma_1 \circ \mathbf{f}_1)(\mathbf{x})$$

Continuous-PieceWise Linear (CPWL) functions



1D: Non-uniform spline de degree 1

Partition: $\mathbb{R} = \bigcup_{k=0}^K P_k$ with $P_k = [\tau_k, \tau_{k+1})$, $\tau_0 = -\infty < \tau_1 < \dots < \tau_K < \tau_{K+1} = +\infty$.

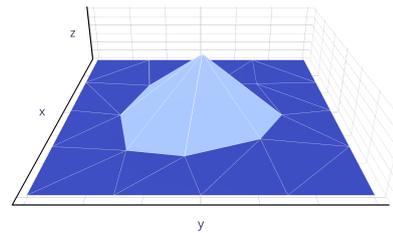
The function $f_{\text{spline}} : \mathbb{R} \rightarrow \mathbb{R}$ is a piecewise-linear spline with knots τ_1, \dots, τ_K if

- (i) for $x \in P_k$: $f_{\text{spline}}(x) = f_k(x) \triangleq a_k x + b_k$ with $(a_k, b_k) \in \mathbb{R}^2$, $k = 0, \dots, K$
- (ii) f_{spline} is continuous $\mathbb{R} \rightarrow \mathbb{R}$

$$\blacksquare f_{\text{spline}}(x) = \tilde{b}_0 + \tilde{b}_1 x + \sum_{k=1}^K \tilde{a}_k (x - \tau_k)_+ \quad \text{with } \tilde{b}_0, \tilde{b}_1 \in \mathbb{R}, (\tilde{a}_k) \in \mathbb{R}^K.$$

31

CPWL functions in high dimensions



Multidimensional generalization

Partition of domain into a finite number of non-overlapping **convex polytopes**; i.e.,

$$\mathbb{R}^N = \bigcup_{k=1}^K P_k \quad \text{with } \mu(P_{k_1} \cap P_{k_2}) = 0 \text{ for all } k_1 \neq k_2$$

The function $f_{\text{CPWL}} : \mathbb{R}^N \rightarrow \mathbb{R}$ is **continuous piecewise-linear** with partition P_1, \dots, P_K

- (i) for $\mathbf{x} \in P_k$: $f_{\text{CPWL}}(\mathbf{x}) = f_k(\mathbf{x}) \triangleq \mathbf{a}_k^T \mathbf{x} + b_k$ with $\mathbf{a}_k \in \mathbb{R}^N$, $b_k \in \mathbb{R}$, $k = 1, \dots, K$
- (ii) f_{CPWL} is continuous $\mathbb{R}^N \rightarrow \mathbb{R}$

The vector-valued function $\mathbf{f}_{\text{CPWL}} = (f_1, \dots, f_M) : \mathbb{R}^N \rightarrow \mathbb{R}^M$ is a CPWL if each component function $f_m : \mathbb{R}^N \rightarrow \mathbb{R}$ is CPWL.

32

Algebra of CPWL functions

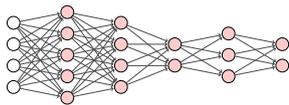
- any linear combination of (vector-valued) CPWL functions $\mathbb{R}^N \rightarrow \mathbb{R}^{N'}$ is CPWL, and,
- the composition $\mathbf{f}_2 \circ \mathbf{f}_1$ of any two CPWL functions with compatible domain and range—i.e., $\mathbf{f}_2 : \mathbb{R}^{N_1} \rightarrow \mathbb{R}^{N_2}$ and $\mathbf{f}_1 : \mathbb{R}^{N_0} \rightarrow \mathbb{R}^{N_1}$ —is CPWL $\mathbb{R}^{N_0} \rightarrow \mathbb{R}^{N_2}$.

Sketch of proof: The continuity property is preserved through composition. The composition of two affine transforms is an affine transform, including the scenario where the domain is partitioned.

- The max (resp. min) pooling of two (or more) CPWL functions is CPWL.

33

Implication for deep ReLU neural networks



 spline



$$\mathbf{f}_{\text{deep}}(\mathbf{x}) = (\sigma_L \circ \mathbf{f}_L \circ \sigma_{L-1} \circ \cdots \circ \sigma_2 \circ \mathbf{f}_2 \circ \sigma_1 \circ \mathbf{f}_1)(\mathbf{x})$$

- Each scalar neuron activation, $\sigma_{n,\ell}(x) = \text{ReLU}(x)$, is CPWL.
- Each layer function $\sigma_\ell \circ \mathbf{f}_\ell(\mathbf{x}) = (\mathbf{W}_\ell \mathbf{x} + \mathbf{b}_\ell)_+$ is CPWL
- The whole feedforward network $\mathbf{f}_{\text{deep}} : \mathbb{R}^{N_0} \rightarrow \mathbb{R}^{N_L}$ is CPWL
- This holds true as well for deep architectures that involve Max pooling for dimension reduction
- The CPWL also remains valid for more complicated neuronal responses as long as they are CPWL; that is, **linear splines**.

34

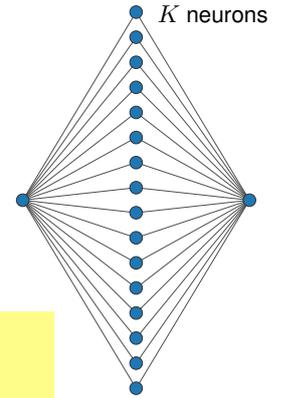
Limit behaviour of univariate shallow ReLU neural nets

- 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 (DNN-1) is achieved by the **sparse adaptive spline**:

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

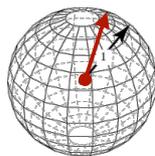
Arguments for the proof:

- Scale invariance of ReLU architecture: For any $\gamma > 0$, the map $(v_k, w_k) \mapsto (\gamma v_k, w_k/\gamma)$ does not affect f_{θ} .
- At the optimum of (NN-1), $|w_k| = |v_k|$, for $k = 1, \dots, K$ and $\text{TV}^{(2)}(f_{\theta}) = \sum_{k=1}^K |a_k|$ with $a_k = v_k |w_k|$.

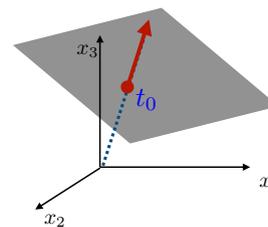
(Savarese 2019; Parhi-Nowak 2020)

The Radon transform and the FBP algorithm

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



Hyperplane $P_{\xi_0, t_0} = \{x \in \mathbb{R}^d : \xi_0^T x = t_0\}$



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

$$\mathbb{R}\{f\}(t, \xi) = \int_{\mathbb{R}^d} \delta(t - \xi^T x) f(x) dx, \quad (t, \xi) \in \mathbb{R} \times \mathbb{S}^{d-1}$$

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

$$f = \mathbb{R}^* \mathbb{K}_{\text{rad}} \{g\}$$

- \mathbb{K}_{rad} : "radial" filtering in Radon space along the variable t .
Frequency response: $\widehat{K}_{\text{rad}}(\omega) \propto |\omega|^{d-1}$
- \mathbb{R}^* : backprojection operator (the adjoint of \mathbb{R})

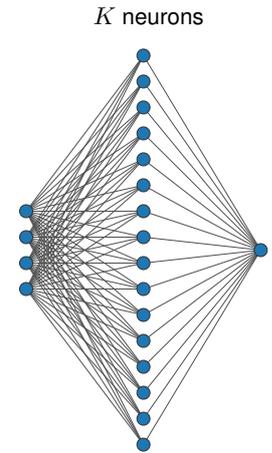
Limit behaviour of multivariate 2-layer ReLU neural nets

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

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

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

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



Theorem

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

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

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

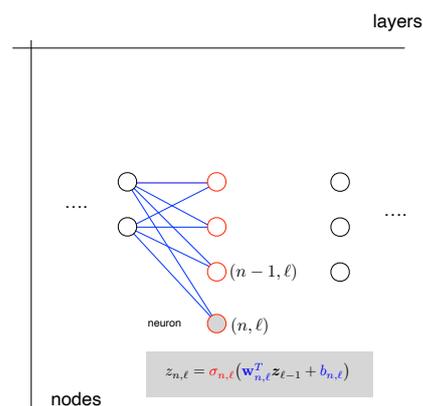
- $\mathcal{M}_{\text{Rad}}(\mathbb{R} \times \mathbb{S}^{d-1})$: space of bounded Radon-compatible measures
 $\mathcal{M}_{\text{Rad}} \subset \mathcal{S}'_{\text{Rad}} = \mathbf{K}_{\text{rad}} \mathbf{R}(\mathcal{S}'(\mathbb{R}^d))$
- $\mathcal{M}_{\Delta_R}(\mathbb{R}^d) =$ Banach space that is isometrically isomorphic to $\mathcal{M}_{\text{Rad}} \times \{c_0 + \mathbf{c}_1^T \mathbf{x}\}$
- Regularization operator $\Delta_R = \mathbf{K}_{\text{rad}} \mathbf{R} \Delta : \mathcal{M}_{\Delta_R}(\mathbb{R}^d) \rightarrow \mathcal{M}_{\text{Rad}}$

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

37

Refinement: free-form activation functions

- Layers: $\ell = 1, \dots, L$
- Deep structure descriptor: (N_0, N_1, \dots, N_L)
- Neuron or node index: (n, ℓ) , $n = 1, \dots, N_{\ell}$
- Activation function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ (ReLU)
- Linear step: $\mathbb{R}^{N_{\ell-1}} \rightarrow \mathbb{R}^{N_{\ell}}$
 $\mathbf{f}_{\ell} : \mathbf{x} \mapsto \mathbf{f}_{\ell}(\mathbf{x}) = \mathbf{W}_{\ell} \mathbf{x} + \mathbf{b}_{\ell}$
- Nonlinear step: $\mathbb{R}^{N_{\ell}} \rightarrow \mathbb{R}^{N_{\ell}}$
 $\boldsymbol{\sigma}_{\ell} : \mathbf{x} \mapsto \boldsymbol{\sigma}_{\ell}(\mathbf{x}) = (\sigma_{n,\ell}(x_1), \dots, \sigma_{N_{\ell},\ell}(x_{N_{\ell}}))$



$$\mathbf{f}_{\text{deep}}(\mathbf{x}) = (\boldsymbol{\sigma}_L \circ \mathbf{f}_L \circ \boldsymbol{\sigma}_{L-1} \circ \dots \circ \boldsymbol{\sigma}_2 \circ \mathbf{f}_2 \circ \boldsymbol{\sigma}_1 \circ \mathbf{f}_1)(\mathbf{x})$$

Joint learning / training ?

38

Constraining activation functions

■ Regularization functional

- Should not penalize simple solutions (e.g., identity or linear scaling)
- Should impose differentiability (for DNN to be trainable via backpropagation)
- Should favor simplest CPWL solutions; i.e., with “sparse 2nd derivatives”

■ Second total-variation of $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

$$\text{TV}^{(2)}(\sigma) \triangleq \|\text{D}^2\sigma\|_{\mathcal{M}} = \sup_{\varphi \in \mathcal{S}(\mathbb{R}) : \|\varphi\|_{\infty} \leq 1} \langle \text{D}^2\sigma, \varphi \rangle$$

■ Native space for $(\mathcal{M}(\mathbb{R}), \text{D}^2)$

$$\text{BV}^{(2)}(\mathbb{R}) = \{f : \mathbb{R} \rightarrow \mathbb{R} : \|\text{D}^2f\|_{\mathcal{M}} < \infty\}$$

39

Representer theorem for deep neural networks

Theorem (TV⁽²⁾-optimality of deep spline networks)

(Unser, JMLR 2019)

- neural network $\mathbf{f} : \mathbb{R}^{N_0} \rightarrow \mathbb{R}^{N_L}$ with **deep structure** (N_0, N_1, \dots, N_L)
 $\mathbf{x} \mapsto \mathbf{f}(\mathbf{x}) = (\sigma_L \circ \ell_L \circ \sigma_{L-1} \circ \dots \circ \ell_2 \circ \sigma_1 \circ \ell_1)(\mathbf{x})$
- **normalized** linear transformations $\ell_\ell : \mathbb{R}^{N_{\ell-1}} \rightarrow \mathbb{R}^{N_\ell}$, $\mathbf{x} \mapsto \mathbf{U}_\ell \mathbf{x}$ with weights
 $\mathbf{U}_\ell = [\mathbf{u}_{1,\ell} \dots \mathbf{u}_{N_\ell,\ell}]^T \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$ such that $\|\mathbf{u}_{n,\ell}\| = 1$
- **free-form** activations $\sigma_\ell = (\sigma_{1,\ell}, \dots, \sigma_{N_\ell,\ell}) : \mathbb{R}^{N_\ell} \rightarrow \mathbb{R}^{N_\ell}$ with $\sigma_{1,\ell}, \dots, \sigma_{N_\ell,\ell} \in \text{BV}^{(2)}(\mathbb{R})$

Given a series data points $(\mathbf{x}_m, \mathbf{y}_m)$ $m = 1, \dots, M$, we then define the training problem

$$\arg \min_{(\mathbf{U}_\ell), (\sigma_{n,\ell} \in \text{BV}^{(2)}(\mathbb{R}))} \left(\sum_{m=1}^M E(\mathbf{y}_m, \mathbf{f}(\mathbf{x}_m)) + \mu \sum_{\ell=1}^L R_\ell(\mathbf{U}_\ell) + \lambda \sum_{\ell=1}^L \sum_{n=1}^{N_\ell} \text{TV}^{(2)}(\sigma_{n,\ell}) \right) \quad (1)$$

- $E : \mathbb{R}^{N_L} \times \mathbb{R}^{N_L} \rightarrow \mathbb{R}^+$: arbitrary convex error function
- $R_\ell : \mathbb{R}^{N_\ell \times N_{\ell-1}} \rightarrow \mathbb{R}^+$: convex cost

If solution of (1) exists, then it is achieved by a **deep spline network** with activations of the form

$$\sigma_{n,\ell}(x) = b_{1,n,\ell} + b_{2,n,\ell}x + \sum_{k=1}^{K_{n,\ell}} a_{k,n,\ell}(x - \tau_{k,n,\ell})_+$$

with adaptive parameters $K_{n,\ell} \leq M - 2$, $\tau_{1,n,\ell}, \dots, \tau_{K_{n,\ell},n,\ell} \in \mathbb{R}$, and $b_{1,n,\ell}, b_{2,n,\ell}, a_{1,n,\ell}, \dots, a_{K_{n,\ell},n,\ell} \in \mathbb{R}$.

40

Outcome of representer theorem

Each neuron (fixed index (n, ℓ)) is characterized by

- its number $0 \leq K_{n,\ell}$ of knots (ideally, much smaller than M);
- the location $\{\tau_k = \tau_{k,n,\ell}\}_{k=1}^{K_{n,\ell}}$ of these knots (ReLU biases);
- the expansion coefficients $\mathbf{b}_{n,\ell} = (b_{1,n,\ell}, b_{2,n,\ell}) \in \mathbb{R}^2$,
 $\mathbf{a}_{n,\ell} = (a_{1,n,\ell}, \dots, a_{K,n,\ell}) \in \mathbb{R}^K$.

These parameters (including the number of knots) are **data-dependent** and adjusted automatically during training.

■ Link with ℓ_1 minimization techniques

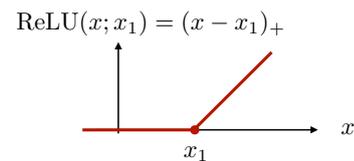
$$\text{TV}^{(2)}\{\sigma_{n,\ell}\} = \sum_{k=1}^{K_{n,\ell}} |a_{k,n,\ell}| = \|\mathbf{a}_{n,\ell}\|_1$$

41

Deep spline networks: Discussion

- Global optimality achieved with **spline activations**
- Justification of popular schemes / Backward compatibility

- Standard ReLU networks ($K_{n,\ell} = 1$, $\mathbf{b}_{n,\ell} = \mathbf{0}$)



(Glorot *ICAI*S 2011)

(LeCun-Bengio-Hinton *Nature* 2015)

- Linear regression: $\lambda \rightarrow \infty \Rightarrow K_{n,\ell} = 0$

- State-of-the-art Parametric ReLU networks ($K_{n,\ell} = 1$)
1 ReLU + linear term (per neuron)

(He et al. *CVPR* 2015)

- Adaptive-piecewise linear (APL) networks ($K_{n,\ell} = 5$ or 7 , $\mathbf{b}_{n,\ell} = \mathbf{0}$)

(Agostinelli et al. 2015)

42

Deep spline networks (Cont'd)

■ Key features

- Direct control of complexity (number of knots): adjustment of λ
- Ability to suppress unnecessary layers

■ Challenges

- Adaptive knots: more difficult optimization problem \Rightarrow In need of novel training algorithms
- Optimal allocation of knots
 ℓ_1 -minimization with knot deletion mechanism (even for single layer)
- Finding the tradeoff: more complex activations vs. deeper architectures

43

CONCLUSION: Return of the spline

■ Foundations of functional learning

- Functional optimization in Banach spaces (enabled by representer theorem)
- **Hilbert spaces**: the tools of classical ML
- **Non-convex Banach spaces**: for sparsity-promoting regularization (e.g., CS)

■ Splines and machine learning

- Traditional kernel methods are closely related to splines (with one knot/kernel per data point)
- Sparse variants offer promising perspectives
- Deep ReLU neural nets are **high-dimensional** piecewise-linear **splines**
- Functional optimization for the streamlining of neuronal architectures
- **Free-form** activations with TV-regularization \Rightarrow **Deep splines**

44

ACKNOWLEDGMENTS

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

- Dr. Julien Fageot
- Shayan Aziznejad
- 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 Jianwei Ma
-



References

■ 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," arXiv:2003.10112, 2020.

■ Representer theorems

- M. Unser, "A Unifying Representer Theorem for Inverse Problems and Machine Learning," *Foundations of Computational Mathematics*, vol. 21, pp. 941-960, 2021.
- S. Aziznejad, M. Unser, "Multi-Kernel Regression with Sparsity Constraints," *SIAM Journal on Mathematics of Data Science*, vol. 3, no. 1, pp. 201-224, 2021.

■ Neural networks

- K.H. Jin, M.T. McCann, E. Froustey, M. Unser, "Deep Convolutional Neural Network for Inverse Problems in Imaging," *IEEE Trans. Image Processing*, vol. 26, no. 9, pp. 4509-4522, Sep. 2017.
- H. Gupta, K.H. Jin, H.Q. Nguyen, M.T. McCann, M. Unser, "CNN-Based Projected Gradient Descent for Consistent CT Image Reconstruction," *IEEE Trans. Medical Imaging*, vol. 37, no. 6, pp. 1440-1453, 2018.
- M. Unser, "A Representer Theorem for Deep Neural Networks," *J. Machine Learning Research*, vol. 20, no. 110, pp. 1-30, Jul. 2019.

- Preprints and demos: <http://bigwww.epfl.ch/>

Sketch of proof

$$\min_{(\mathbf{U}_\ell), (\tilde{\sigma}_{n,\ell} \in \text{BV}^{(2)}(\mathbb{R}))} \left(\sum_{m=1}^M E(\mathbf{y}_m, \mathbf{f}(\mathbf{x}_m)) + \mu \sum_{\ell=1}^L R_\ell(\mathbf{U}_\ell) + \lambda \sum_{\ell=1}^L \sum_{n=1}^{N_\ell} \text{TV}^{(2)}(\tilde{\sigma}_{n,\ell}) \right)$$

Optimal solution $\tilde{\mathbf{f}} = \tilde{\sigma}_L \circ \tilde{\ell}_L \circ \tilde{\sigma}_{L-1} \circ \dots \circ \tilde{\ell}_2 \circ \tilde{\sigma}_1 \circ \tilde{\ell}_1$ with optimized weights $\tilde{\mathbf{U}}_\ell$ and neuronal activations $\tilde{\sigma}_{n,\ell}$.

Apply “optimal” network $\tilde{\mathbf{f}}$ to each data point \mathbf{x}_m :

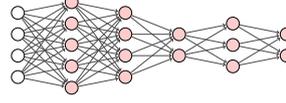
- Initialization (input): $\tilde{\mathbf{y}}_{m,0} = \mathbf{x}_m$.

- For $\ell = 1, \dots, L$

$$\mathbf{z}_{m,\ell} = (z_{1,m,\ell}, \dots, z_{N_\ell,m,\ell}) = \tilde{\mathbf{U}}_\ell \tilde{\mathbf{y}}_{m,\ell-1}$$

$$\tilde{\mathbf{y}}_{m,\ell} = (\tilde{y}_{1,m,\ell}, \dots, \tilde{y}_{N_\ell,m,\ell}) \in \mathbb{R}^{N_\ell}$$

$$\text{with } \tilde{y}_{n,m,\ell} = \tilde{\sigma}_{n,\ell}(z_{n,m,\ell}) \quad n = 1, \dots, N_\ell.$$



$$\Rightarrow \tilde{\mathbf{f}}(\mathbf{x}_m) = \tilde{\mathbf{y}}_{m,L}$$

This fixes two terms of minimal criterion: $\sum_{m=1}^M E(\mathbf{y}_m, \tilde{\mathbf{y}}_{m,L})$ and $\sum_{\ell=1}^L R_\ell(\tilde{\mathbf{U}}_\ell)$.

$\tilde{\mathbf{f}}$ achieves global optimum

$$\Leftrightarrow \tilde{\sigma}_{n,\ell} = \arg \min_{f \in \text{BV}^{(2)}(\mathbb{R})} \|D^2 f\|_{\mathcal{M}} \quad \text{s.t.} \quad f(z_{n,m,\ell}) = \tilde{y}_{n,m,\ell}, \quad m = 1, \dots, M$$