

Representer theorems for machine learning and inverse problems

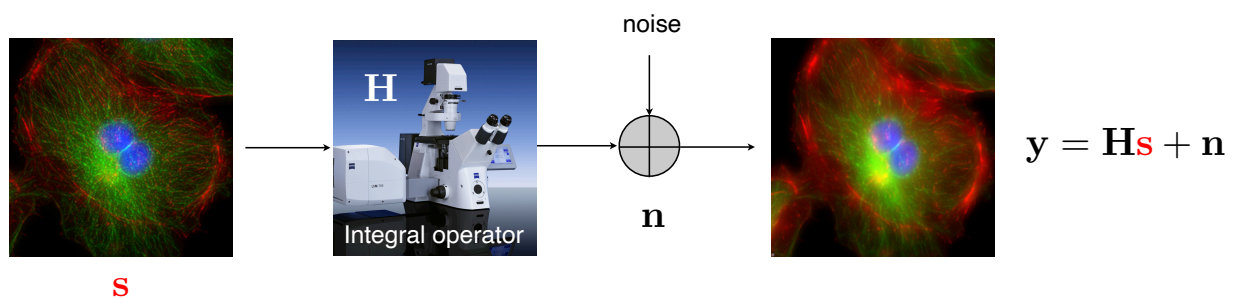
Michael Unser
Biomedical Imaging Group
EPFL, Lausanne, Switzerland



One World Seminar “Mathematical Methods for Arbitrary Data Sources”, June 8, 2020

Variational formulation of inverse problems

- Linear forward model



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

- Reconstruction as an optimization problem

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

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)

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

\Rightarrow kernel estimator
(Wahba 1990; Schölkopf 2001)

3

RKHS representer theorem for machine learning

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

(Poggio-Girosi 1990)

$r_{\mathcal{H}} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is the (unique) **reproducing kernel** for the RKHS \mathcal{H} 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 $\mathbf{x}_0 \in \mathbb{R}^d$

(Aronszajn, 1950)

Formal characterization: $r_{\mathcal{H}}(\cdot, \mathbf{x}_0) = R\{\delta(\cdot - \mathbf{x}_0)\} = (\delta(\cdot - \mathbf{x}_0))^*$ (Riesz conjugate)

Representer theorem for L_2 -regularization

The solution of (P2) has the generic parametric form: $f(\mathbf{x}) = \sum_{m=1}^M a_m r_{\mathcal{H}}(\mathbf{x}, \mathbf{x}_m)$

(de Boor 1966; Kimeldorf-Wahba 1971; Poggio-Girosi 1990)

4

RKHS representer theorem for machine learning

$$(P2') \quad \arg \min_{f \in \mathcal{H}} \left(\sum_{m=1}^M E(y_m, f(\mathbf{x}_m)) + \lambda \|f\|_{\mathcal{H}}^2 \right) \quad \text{with} \quad E : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R} \text{ convex}$$

$r_{\mathcal{H}} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is the (unique) **reproducing kernel** for the RKHS \mathcal{H} 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 $\mathbf{x}_0 \in \mathbb{R}^d$

(Aronszajn, 1950)

Formal characterization: $r_{\mathcal{H}}(\cdot, \mathbf{x}_0) = \mathbb{R}\{\delta(\cdot - \mathbf{x}_0)\} = (\delta(\cdot - \mathbf{x}_0))^*$ (Riesz conjugate)

Representer theorem for L_2 -regularization

The solution of (P2') has the generic parametric form: $f(\mathbf{x}) = \sum_{m=1}^M a_m r_{\mathcal{H}}(\mathbf{x}, \mathbf{x}_m)$

(de Boor 1966; Kimeldorf-Wahba 1971; Poggio-Girosi 1990; Schölkopf 2001)

Supports the theory of SVM, **kernel methods**, etc.

5

Is there a mother of all representer theorems ?

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

Classical representer theorem in machine learning:

- $\mathcal{X}' = \mathcal{H}$ is a reproducing kernel Hilbert space.
- $\nu : \mathcal{H}' \rightarrow \mathbb{R}^M : f \mapsto (f(\mathbf{x}_1), \dots, f(\mathbf{x}_M))$ is the sampling operator.

(de Boor 1966; Kimeldorf-Wahba 1971; Poggio-Girosi 1990; Schölkopf 2001)

Most general set-up:

- \mathcal{X} is a Banach space.
- $\nu : \mathcal{X} \rightarrow \mathbb{R}^M : f \mapsto (\langle \nu_1, f \rangle, \dots, \langle \nu_M, f \rangle)$ is a general linear measurement operator.
- $E : \mathbb{R}^M \times \mathbb{R}^M \rightarrow \mathbb{R}^+ \cup \{+\infty\}$ is a proper l.s.c. convex loss functional.
- $\psi : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is some arbitrary strictly-increasing convex function.

6

CONTENT

- **Regularization in science and engineering** ✓
 - Inverse problems and learning
 - Representer theorem for RKHS
- **Unifying representer theorem**
 - Banach spaces and their duals
 - Duality mapping
 - Mother of all representer theorems
- **Applications:** Optimization in specific Banach spaces
 - Learning in RKHS
 - Tikhonov regularization
 - l_p -norm regularization
 - Sparse kernel expansions
 - Deep neural networks

7

General notion of Banach space

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

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

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



Stefan Banach (1892-1945)

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

8

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**:

$$\|g\|_{\mathcal{X}'} = \sup_{f \in \mathcal{X} \setminus \{0\}} \left(\frac{\langle g, f \rangle}{\|f\|_{\mathcal{X}}} \right) = \sup_{f \in \mathcal{X}: \|f\|_{\mathcal{X}} \leq 1} |\langle g, f \rangle|$$

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

9

Riesz conjugate for Hilbert spaces

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

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



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

■ $u^* = \mathbf{R}^{-1}\{u\}$ (inverse Riesz map)

■ Norm preservation: $\|u\|_{\mathcal{H}} = \|u^*\|_{\mathcal{H}'}$

■ Invertibility: $u = (u^*)^* = \mathbf{R}\{u^*\}$

■ Linearity: $(u_1 + u_2)^* = u_1^* + u_2^*$

+

(isometry)

$(\mathcal{H}')' = \mathcal{H}$ (reflexivity)

10

Generalization: Duality mapping

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

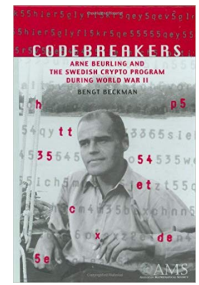
- $\|f^*\|_{\mathcal{X}'} = \|f\|_{\mathcal{X}}$ (**norm preservation**), and
- $\langle f^*, f \rangle_{\mathcal{X}' \times \mathcal{X}} = \|f^*\|_{\mathcal{X}'} \|f\|_{\mathcal{X}}$ (**sharp duality bound**).

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)



Arne Beurling (1905-1986)

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. $J(\lambda f) = \lambda J(f)$ for any $\lambda \in \mathbb{R}^+$ (homogeneity).
3. For every $f \in \mathcal{X}$, the set $J(f)$ is convex and weak-* closed in \mathcal{X}' .
4. The duality mapping is **single-valued** if \mathcal{X}' is **strictly convex**; the latter condition is also necessary if \mathcal{X} is reflexive.
5. When \mathcal{X} is **reflexive**, then the duality map is **bijective** if and only if both \mathcal{X} and \mathcal{X}' are **strictly convex**.

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

\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 $\lambda \in (0, 1)$.

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.

13

Mother of all representer theorems (Cont'd)

Theorem

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

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

is **non-empty, convex** and weak*-compact.

Any solution $f_0 \in S \subset \mathcal{X}'$ is a $(\mathcal{X}', \mathcal{X})$ -**conjugate of a common**

$$\nu_0 = \sum_{m=1}^M a_m \nu_m \in \mathcal{N}_{\boldsymbol{\nu}} \subset \mathcal{X}$$

with suitable weights $\mathbf{a} \in \mathbb{R}^M$; i.e., $S \subseteq J(\nu_0)$.

If the Banach space \mathcal{X} is **reflexive and strictly convex**, then the solution is **unique** with $f_0 = J_{\mathcal{X}}\{\nu_0\} \in \mathcal{X}'$ (Banach conjugate of ν_0). If \mathcal{X} is a Hilbert space, then $f_0 = \sum_{m=1}^M a_m \nu_m^*$ where ν_m^* is the Riesz conjugate of ν_m .

(Unser, ArXiv 2019)

14

CONTENT

- Regularization in science and engineering ✓
- Unifying representer theorem ✓
- **Applications:** Optimization in specific Banach spaces
 - Learning in RKHS
 - Tikhonov regularization
 - l_p -norm regularization
 - Sparsity promoting regularization
 - Sparse kernel expansions
 - Deep neural networks

15

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 $\delta(\cdot - \mathbf{x}) \in \mathcal{H}$ for any $\mathbf{x} \in \mathbb{R}^d$. The corresponding unique **Hilbert conjugate** $h(\cdot, \mathbf{x}) = (\delta(\cdot - \mathbf{x}))^* \in \mathcal{H}$ when indexed by \mathbf{x} 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)

16

Learning in RKHS (Cont'd)

- Special case of generalized representer theorem

- $\mathcal{X} = \mathcal{H}'$, $\mathcal{X}' = \mathcal{H}''$ (Reflexive Banach space)
- $\nu_m = \delta(\cdot - \mathbf{x}_m)$ (Dirac sampling functionals)
- Additive loss: $E(\mathbf{y}, \mathbf{z}) = \sum_{m=1}^M E_m(y_m, z_m)$

- Key observation

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

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

- Implied form of unique solution = linear kernel expansion

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

(Schölkopf representer theorem, 2001)

17

2. Tikhonov regularization

\mathcal{H} : **Hilbert space** on \mathbb{R}^d with **Riesz map** $J_{\mathcal{H}'} = \mathbb{R} : \mathcal{H}' \rightarrow \mathcal{H}$

- Ill-posed linear inverse problem

Measurement functionals: $\nu_1, \dots, \nu_M \in \mathcal{H}'$

Goal: recover a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ from noisy measurements $y_m = \langle \nu_m, f \rangle + \epsilon_m$

- Formulation of reconstruction problem (penalized least-squares)

Given the data $\mathbf{y} \in \mathbb{R}^M$, find the function $f_0 : \mathbb{R}^d \rightarrow \mathbb{R}$ s.t.

$$f_0 = \arg \min_{f \in \mathcal{H}} \left(\sum_{m=1}^M |y_m - \langle \nu_m, f \rangle|^2 + \lambda \|f\|_{\mathcal{H}}^2 \right)$$

$\lambda \in \mathbb{R}^+$: adjustable regularization parameter.

18

Tikhonov regularization: closed-form solution

■ Application of generalized representer theorem

- $\mathcal{X} = \mathcal{H}'$, $\mathcal{X}' = \mathcal{H}'' = \mathcal{H}$ (Hilbert space)
- Measurement functionals: $\nu_m \in \mathcal{H}'$, $m = 1, \dots, M$
- Conjugate functions: $\varphi_m = \nu_m^* = \mathbf{R}\{\nu_m\} \in \mathcal{H}$
- $\psi(t) = \lambda|t|^2$ (convex)

$$\Rightarrow f_0 = \sum_{m=1}^M a_m \varphi_m \in \text{span}\{\varphi_m\}$$

■ Optimal discretization: "the miraculous simplification"

- System matrix $\mathbf{H} \in \mathbb{R}^{M \times M} = \mathbf{Gram\ matrix}$ (symmetric, positive-definite)

$$[\mathbf{H}]_{m,n} = \langle \nu_m, \varphi_n \rangle = \langle \nu_m, \nu_n^* \rangle = \langle \nu_m^*, \nu_n^* \rangle_{\mathcal{H}} = \langle \varphi_m, \varphi_n \rangle_{\mathcal{H}}$$

- $f = \sum_{m=1}^M a_m \varphi_m \Rightarrow \nu(f) = \mathbf{H}\mathbf{a}$, $\|f\|_{\mathcal{H}}^2 = \mathbf{a}^T \mathbf{H}\mathbf{a}$

$$\Rightarrow \mathbf{a}_{\text{opt}} = \arg \min_{\mathbf{a} \in \mathbb{R}^M} (\|\mathbf{y} - \mathbf{H}\mathbf{a}\|_2^2 + \lambda \|\mathbf{H}\mathbf{a}\|_2^2) = (\mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

19

3. ℓ_p -norm regularization

■ Finite-dimensional setup (CS)

- Goal: Recover $\mathbf{s} = (s_n) \in \mathbb{R}^N$ from a set of corrupted linear measurements
 $y_m = \mathbf{h}_m^T \mathbf{s} + \epsilon_m$, $m = 1, \dots, M$
- **Compressed sensing** scenario: $M \ll N$
- Strategy: Try to favor **sparse** solutions

■ Formulation of reconstruction task

- Data $\mathbf{y} \in \mathbb{R}^M$
- System matrix $\mathbf{H} = [\mathbf{h}_1 \ \mathbf{h}_2 \ \dots \ \mathbf{h}_M]^T \in \mathbb{R}^{M \times N}$
- Minimization problem with $p > 0$ small

$$\mathbf{s} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(E(\mathbf{y}, \mathbf{H}\mathbf{x}) + \lambda \|\mathbf{x}\|_{\ell_p}^p \right)$$

$\lambda \in \mathbb{R}^+$: adjustable regularization parameter

20

ℓ_p -norm regularization (Cont'd)

■ Application of general representer theorem

- $\mathcal{X} = (\mathbb{R}^N, \|\cdot\|_{\ell_q})$, $\mathcal{X}' = (\mathbb{R}^N, \|\cdot\|_{\ell_p})$ with $\frac{1}{p} + \frac{1}{q} = 1$
- Hölder inequality: $|\langle \mathbf{u}, \mathbf{v} \rangle| \leq \|\mathbf{u}\|_{\ell_p} \|\mathbf{v}\|_{\ell_q}$
- $\psi(x) = \lambda|x|^p$ is convex for $p \geq 1$
- $\|\cdot\|_{\ell_p}$ and $\|\cdot\|_{\ell_q}$ norms are strictly convex for $p \in (1, \infty) \Rightarrow$ **unique solution**
- Known q -to- p duality map: $[\mathbf{v}^*]_n = \frac{|v_n|^{q-1}}{\|\mathbf{v}\|_{\ell_q}^{q-2}} \text{sign}(v_n)$

■ Parametric form of the solution:

$$[\mathbf{s}]_n = \frac{[\mathbf{H}^T \mathbf{a}]_n^{q-1}}{\|\mathbf{H}^T \mathbf{a}\|_{\ell_q}^{q-2}} \text{sign}([\mathbf{H}^T \mathbf{a}]_n)$$

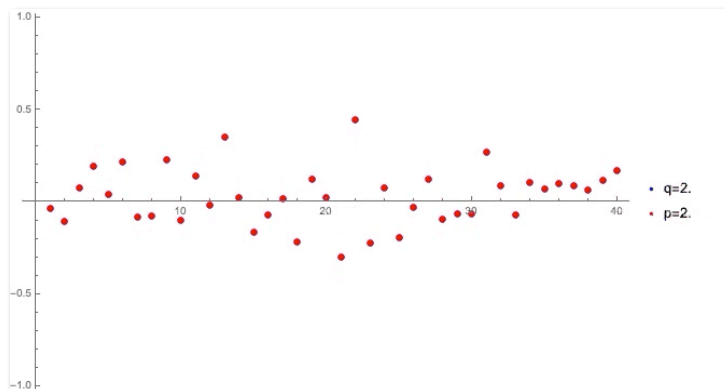
with parameter vector $\mathbf{a} \in \mathbb{R}^M$

21

Qualitative effect of Banach conjugation

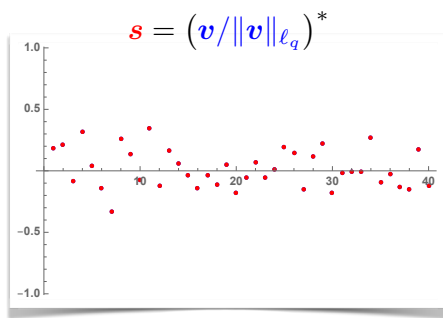
$$J_{\ell_q(\mathbb{Z})} : \ell_q(\mathbb{Z}) \rightarrow \ell_p(\mathbb{Z}) \quad v_n^* = \frac{|v_n|^{q-1}}{\|\mathbf{v}\|_{\ell_q}^{q-2}} \text{sign}(v_n)$$

$$\mathbf{s} = (\mathbf{v} / \|\mathbf{v}\|_{\ell_q})^*$$

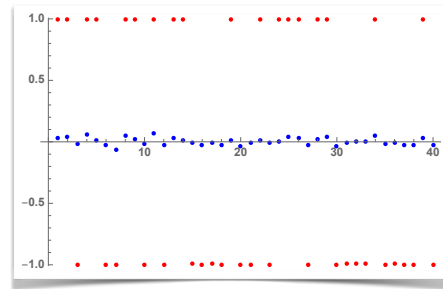


22

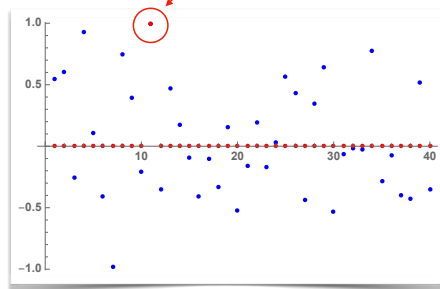
Qualitative effect of Banach conjugation



■ $(q, p) = (2, 2)$: identity



■ $(q, p) \rightarrow (1, \infty)$: saturation of v^*



■ $(q, p) \rightarrow (\infty, 1)$: sparsification of v^*

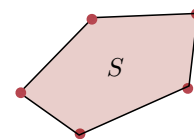
23

4. Sparsity promoting regularization

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

■ Cases where the solution set is not necessarily unique

- \mathcal{X}' is non-reflexive, non-strictly convex; e.g., $\mathcal{X}' = \ell_1(\mathbb{Z})$
- Representer theorem $\Rightarrow S$ is convex, weak* **compact**
- Krein-Milman theorem: S is the convex hull of its **extreme points**



Theorem

All extreme points f_0 of S can be expressed as

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

for some $1 \leq K_0 \leq M$ where the e_k are some **extreme points** of the unit “regularization” ball $B_{\mathcal{X}'} = \{f \in \mathcal{X}' : \|f\|_{\mathcal{X}'} \leq 1\}$ and $\mathbf{a} = (a_1, \dots, a_{K_0}) \in \mathbb{R}^{K_0}$.

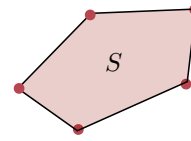
(Boyer-Chambolle-De Castro-Duval-De Gournay-Weiss, arXiv:1806.09810, 2019)

24

Extreme points

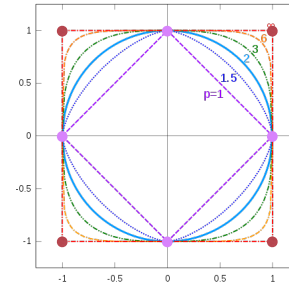
Definition

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 !!!

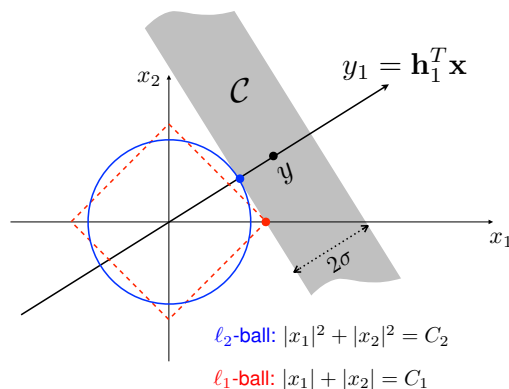
Definition of **strictly convexity**: all boundary points are extreme !!!

Geometry of ℓ_2 vs. ℓ_1 minimization

Prototypical inverse problem

$$\min_{\mathbf{x}} \{ \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\ell_2}^2 + \lambda \|\mathbf{x}\|_{\ell_2}^2 \} \Leftrightarrow \min_{\mathbf{x}} \|\mathbf{x}\|_{\ell_2} \text{ subject to } \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\ell_2}^2 \leq \sigma^2$$

$$\min_{\mathbf{x}} \{ \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\ell_2}^2 + \lambda \|\mathbf{x}\|_{\ell_1} \} \Leftrightarrow \min_{\mathbf{x}} \|\mathbf{x}\|_{\ell_1} \text{ subject to } \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\ell_2}^2 \leq \sigma^2$$

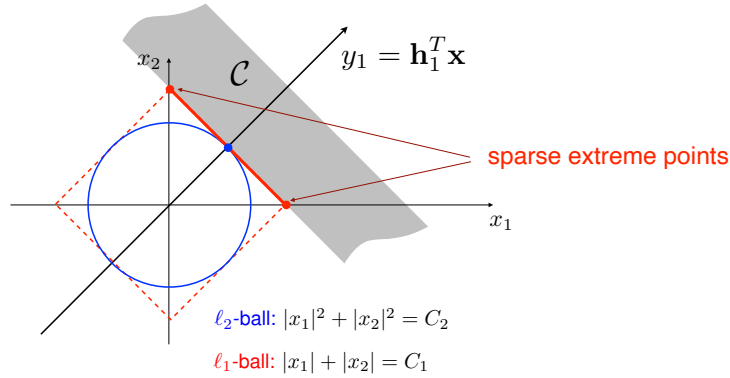


Geometry of l_2 vs. l_1 minimization

Prototypical inverse problem

$$\min_{\mathbf{x}} \{ \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\ell_2}^2 + \lambda \|\mathbf{x}\|_{\ell_2}^2 \} \Leftrightarrow \min_{\mathbf{x}} \|\mathbf{x}\|_{\ell_2} \text{ subject to } \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\ell_2}^2 \leq \sigma^2$$

$$\min_{\mathbf{x}} \{ \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\ell_2}^2 + \lambda \|\mathbf{x}\|_{\ell_1} \} \Leftrightarrow \min_{\mathbf{x}} \|\mathbf{x}\|_{\ell_1} \text{ subject to } \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\ell_2}^2 \leq \sigma^2$$



Configuration for **non-unique** ℓ_1 solution

5. Sparse kernel expansions

Context

- $\mathcal{S}(\mathbb{R}^d)$: Schwartz's space of smooth and rapidly decaying functions on \mathbb{R}^d
- $\mathcal{S}'(\mathbb{R}^d)$: the space of tempered distributions
- Regularization operator $L : \mathcal{S}'(\mathbb{R}^d) \xrightarrow{c} \mathcal{S}'(\mathbb{R}^d)$
- Inverse operator $L^{-1} : \mathcal{S}'(\mathbb{R}^d) \xrightarrow{c} \mathcal{S}'(\mathbb{R}^d)$
- Bivariate kernel: $h : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$



Laurent Schwartz (1915-2002)

$$L^{-1}\{\varphi\} = \int_{\mathbb{R}^d} h(\cdot, \mathbf{y})\varphi(\mathbf{y})d\mathbf{y}$$

Schwartz kernel

Native Banach space for $(L, \mathcal{M}(\mathbb{R}^d))$

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

Isometry with space of Radon measures



Johann Radon (1887-1956)

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

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

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

- Basic isometries

$$L : \mathcal{M}_L(\mathbb{R}^d) \rightarrow \mathcal{M}(\mathbb{R}^d)$$

$$L^{-1} : \mathcal{M}(\mathbb{R}^d) \rightarrow \mathcal{M}_L(\mathbb{R}^d)$$

$$L^{-1} : \varphi \mapsto \int_{\mathbb{R}^d} h(\cdot, \mathbf{y}) \varphi(\mathbf{y}) d\mathbf{y}$$

- Extreme points** of unit ball in $\mathcal{M}_L(\mathbb{R}^d)$:

$$u_k = L^{-1}\{e_k\} = \pm L^{-1}\{\delta(\cdot - \tau_k)\} = \pm h(\cdot, \tau_k)$$

Sparse kernel expansions (Cont'd)

$$L^{-1} : \varphi \mapsto \int_{\mathbb{R}^d} h(\cdot, \mathbf{y}) \varphi(\mathbf{y}) d\mathbf{y}$$

$$S = \arg \min_{f \in \mathcal{M}_L(\mathbb{R}^d)} \left(\sum_{m=1}^M E_m(y_m, f(\mathbf{x}_m)) + \lambda \|Lf\|_{\mathcal{M}} \right)$$

Theorem

All extreme points f_0 of S can be expressed as

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

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

Special case: Translation-invariant kernels

Linear-shift invariant (LSI) setting

- LSI operator L with **frequency response** $\widehat{L}(\omega) = \mathcal{F}\{L\{\delta\}\}(\omega)$
- LSI inverse operator $L^{-1} : \varphi \mapsto h_{\text{LSI}} * \varphi$
- Translation-invariant kernel: $h(\mathbf{x}, \boldsymbol{\tau}) = h_{\text{LSI}}(\mathbf{x} - \boldsymbol{\tau})$

■ Determination of the kernel:
$$h_{\text{LSI}}(\mathbf{x}) = \mathcal{F}^{-1} \left\{ \frac{1}{\widehat{L}(\omega)} \right\}(\mathbf{x})$$

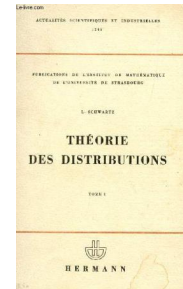
Determination of the regularization operator

$$\widehat{L}(\omega) = \frac{1}{\widehat{h}_{\text{LSI}}(\omega)}$$

$$L : \mathcal{S}'(\mathbb{R}^d) \xrightarrow{c.} \mathcal{S}'(\mathbb{R}^d) \Leftrightarrow \widehat{L}(\omega) \text{ smooth and slowly growing}$$

Example of admissible kernels:

$$h_{\text{LSI}}(\mathbf{x}) = \exp(-\|\mathbf{x}\|^\alpha) \quad \text{with } \alpha \in (0, 2)$$



31

RKHS vs. sparse kernel expansions (LSI)

$$\min_{f \in L_{2,L}(\mathbb{R}^d)} \left(\sum_{m=1}^M E_m(y_m, f(\mathbf{x}_m)) + \lambda \|Lf\|_{L_2}^2 \right)$$

$$\Rightarrow f_{\text{RKHS}}(\mathbf{x}) = \sum_{m=1}^M a_m h_{\text{PD}}(\mathbf{x} - \mathbf{x}_m)$$

$$\text{Quadratic energy: } \|Lf_{\text{RKHS}}\|_{L_2}^2 = \mathbf{a}^T \mathbf{G} \mathbf{a}$$

Positive-definite kernel:

$$h_{\text{PD}}(\mathbf{x}) = \mathcal{F}^{-1} \left\{ \frac{1}{|\widehat{L}(\omega)|^2} \right\}(\mathbf{x})$$

$$\min_{f \in \mathcal{M}_L(\mathbb{R}^d)} \left(\sum_{m=1}^M E_m(y_m, f(\mathbf{x}_m)) + \lambda \|Lf\|_{\mathcal{M}} \right)$$

$$\Rightarrow f_{\text{sparse}}(\mathbf{x}) = \sum_{k=1}^{K_0} a_k h_{\text{LSI}}(\mathbf{x} - \boldsymbol{\tau}_k)$$

$$\text{Sparsity-promoting energy: } \|Lf_{\text{sparse}}\|_{\mathcal{M}} = \|\mathbf{a}\|_{\ell_1}$$

Admissible kernel:

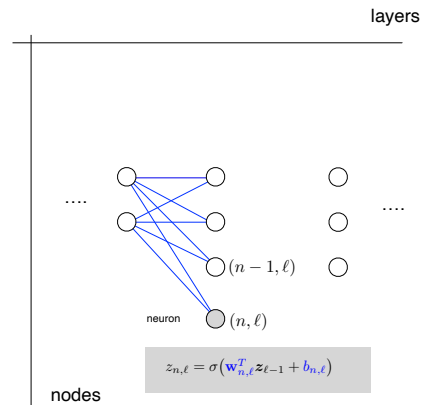
$$h_{\text{LSI}}(\mathbf{x}) = \mathcal{F}^{-1} \left\{ \frac{1}{\widehat{L}(\omega)} \right\}(\mathbf{x})$$

Adaptive parameters: $K_0 \leq M, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{K_0} \in \mathbb{R}^d$

32

6. 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: $\mathbb{R}^{N_\ell} \rightarrow \mathbb{R}^{N_\ell}$
 $\sigma_\ell : \mathbf{x} \mapsto \sigma_\ell(\mathbf{x}) = (\sigma(x_1), \dots, \sigma(x_{N_\ell}))$

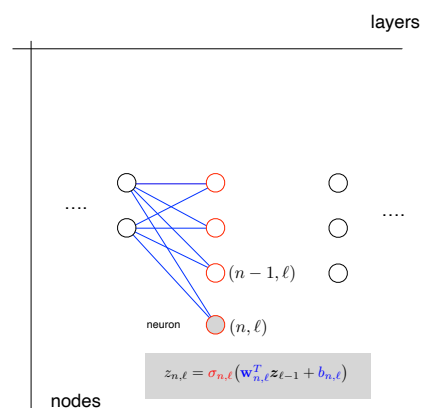


Learned

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

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: $\mathbb{R}^{N_{\ell-1}} \rightarrow \mathbb{R}^{N_\ell}$
 $f_\ell : \mathbf{x} \mapsto \mathbf{f}_\ell(\mathbf{x}) = \mathbf{W}_\ell \mathbf{x} + \mathbf{b}_\ell$
- Nonlinear step: $\mathbb{R}^{N_\ell} \rightarrow \mathbb{R}^{N_\ell}$
 $\sigma_\ell : \mathbf{x} \mapsto \sigma_\ell(\mathbf{x}) = (\sigma_{n,\ell}(x_1), \dots, \sigma_{n,\ell}(x_{N_\ell}))$



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

Joint learning / training ?

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 $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

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

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

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

35

Representer theorem for deep neural networks

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

(U. 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 $b_{1,n,\ell}, b_{2,n,\ell}, a_{1,n,\ell}, \dots, a_{K_{n,\ell},n,\ell} \in \mathbb{R}$.

36

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)

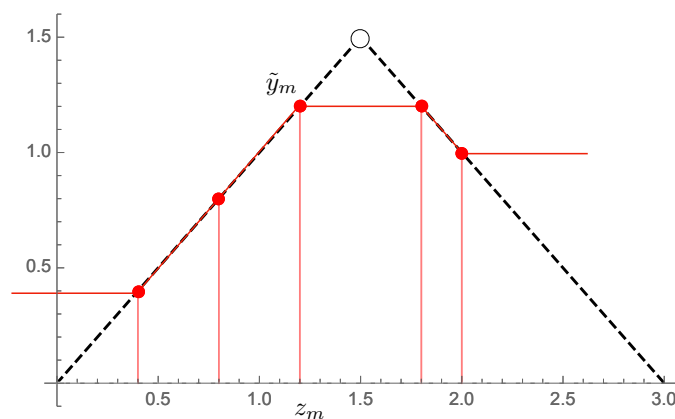
- Adaptive-piecewise linear (APL) networks ($K_{n,\ell} = 5$ or $7, \mathbf{b}_{n,\ell} = \mathbf{0}$)

(Agostinelli et al. 2015)

37

Comparison of linear interpolators

$$\arg \min_{f \in H^1(\mathbb{R})} \int_{\mathbb{R}} |Df(x)|^2 dx \quad \text{s.t.} \quad f(x_m) = y_m, \quad m = 1, \dots, M$$



$$\arg \min_{f \in BV^2(\mathbb{R})} \|D^2 f\|_{\mathcal{M}} \quad \text{s.t.} \quad f(x_m) = y_m, \quad m = 1, \dots, M$$

38

Conclusion

- Unifying result that supports all known “representer” theorems
 - Classical methods based on **quadratic minimization**
 - Kernel-based methods for RKHS (Poggio-Girosi 1990; Schölkopf 2001)
 - Tikhonov regularization (Tikhonov 1977; Gupta 2018)
 - Optimization in **reflexive and strictly-convex** Banach spaces
 - L_p splines (de Boor 1976; ...)
 - Reproducing kernel Banach spaces (Zhang-Xu-Zhang 2009; Zhang-Zhang 2012)
 - Modern **sparsity-based** optimization
 - ℓ_1 -minimization for compressed sensing (Donoho 2006; Candes 2006; Baraniuk 2007)
 - Total variation minimization for the recovery of spikes (Candes Fernandez-Grada 2013; Duval-Peyré 2015)
 - L-splines are optimum solutions for inverse problems with generalized total-variation regularization (Unser-Fageot-Ward 2017; Flinth-Weiss 2018; Bredies-Carioni 2020)
 - Optimality of deep ReLU networks (Unser 2019)

39

Conclusion (Cont'd)

- Remarkable level of generality \Rightarrow opens up new perspectives
 - Fundamental ingredients for applicability
 - **Banach space** that is matched to the problem at hand
 - Knowledge of **dual mapping** vs. **extreme points**

No need for Fréchet derivatives or sub-gradients !!!

- Sparse kernel expansions: Open computational challenge
 - Efficient algorithm for **displacing/removing** kernels

40

Acknowledgments

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

- Dr. Julien Fageot
- Prof. John Paul Ward
- Harshit Gupta
- Shayan Aziznejad
- Thomas Debarre
- Dr. Emrah Bostan
- Prof. Ulugbek Kamilov
- Prof. Matthieu Guerquin-Kern

....



and collaborators ...

- Prof. Demetri Psaltis
- Prof. Marco Stampanoni
- Prof. Carlos-Oscar Sorzano
- Dr. Arne Seitz
-



Selected references

■ Reference books

- A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*, Kluwer, 2004.
- W. Rudin, *Functional Analysis*, McGraw-Hill Book Co., 1991.
- L. Schwartz. *Théorie des Distributions*. Hermann, Paris, 1966.

■ Selected papers

- M. Unser, "A unifying representer theorem for inverse problems and machine learning," arXiv:1903.00687 (2019).
- M. Unser, J. Fageot, and H. Gupta, Representer theorems for sparsity-promoting ℓ_1 regularization, *IEEE Trans. Information Theory*, 62 (2016), pp. 5167–5180.
- M. Unser, J. Fageot, and J. P. Ward, "Splines are universal solutions of linear inverse problems with generalized-TV regularization," *SIAM Review*, 59 (2017), pp. 769–793.
- H. Gupta, J. Fageot, and M. Unser, "Continuous-domain solutions of linear inverse problems with Tikhonov versus generalized TV regularization," *IEEE Trans. Signal Processing*, 66 (2018), pp. 4670–4684.
- C. Boyer, A. Chambolle, Y. De Castro, V. Duval, F. De Gournay, and P. Weiss, "On representer theorems and convex regularization", *SIAM Journal of Optimization*, 29 (2019) pp. 1260-1281.
- S. Aziznejad and M. Unser, "Multiple-kernel regression with sparsity constraint", arXiv:1811.00836, 2018.
- M. Unser, "A representer theorem for deep neural networks," *J. Machine Learning Research*, 20 (2019), no. 110, pp. 1-30.

- 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 :

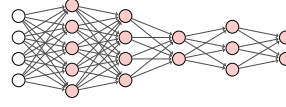
- Initialization (input): $\tilde{\mathbf{y}}_{m,0} = \mathbf{x}_m$.

- 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

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

43

Tikhonov regularization (Exact solution)

$$f_0 = \arg \min_{f \in \mathcal{H}} \left(\sum_{m=1}^M |y_m - \langle \nu_m, f \rangle|^2 + \lambda \|f\|_{\mathcal{H}}^2 \right)$$

and

$$f_0 = \text{span}\{\varphi_m\}_{m=1}^M$$

- Equivalent finite-dimensional problem

$$\mathbf{a}_0 = \arg \min_{\mathbf{a} \in \mathbb{R}^M} (\|\mathbf{y} - \mathbf{H}\mathbf{a}\|^2 + \lambda \mathbf{a}^T \mathbf{H}\mathbf{a})$$

- Closed-form solution

$$\mathbf{a}_0 = (\mathbf{H}\mathbf{H} + \lambda \mathbf{H})^{-1} \mathbf{H}\mathbf{y} = (\mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

\mathbf{H} invertible $\Leftrightarrow \nu_m$ are linearly independent

44