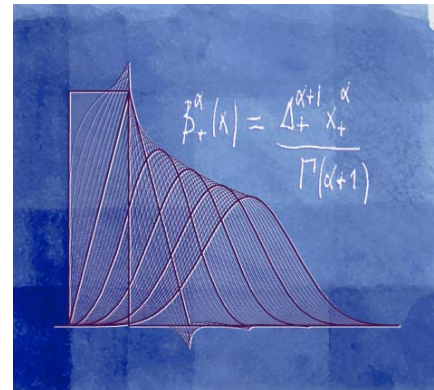




New representer theorems: From compressed sensing to deep learning

Michael Unser
Biomedical Imaging Group
EPFL, Lausanne, Switzerland

Joint work with
Julien Fageot, John-Paul Ward
and Kyong Jin



Mathematisches Kolloquium, Universität Wien, October 24, 2018, Wien, Austria

OUTLINE

■ Introduction

- Image reconstruction as an inverse problem
- Learning as an inverse problem

■ Prologue: discrete-domain regularization

■ Continuous-domain theory

- Splines and operators
- L_2 regularization (theory of RKHS) : classical representer theorem
- gTV regularization: representer theorem for CS

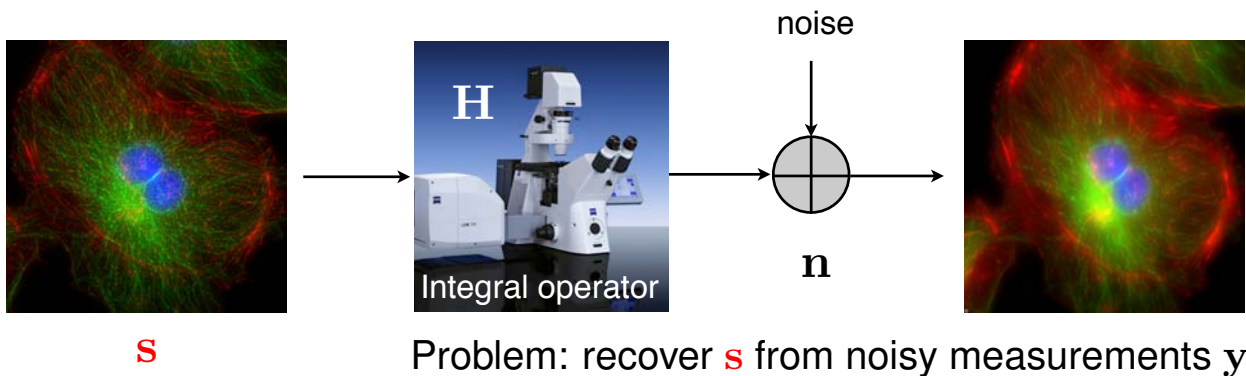
■ From compressed sensing to deep networks

- Unrolling forward/backward iterations: FBPCnv
- New representer theorem for deep neural networks

Variational formulation of inverse problem

Linear forward model

$$\mathbf{y} = \mathbf{H}\mathbf{s} + \mathbf{n}$$



Reconstruction as an optimization problem

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

3

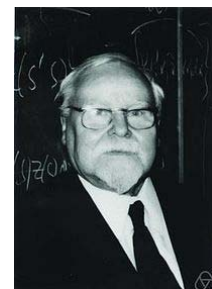
Linear inverse problems (20th century theory)

Dealing with ill-posed problems: Tikhonov regularization

$\mathcal{R}(\mathbf{s}) = \|\mathbf{L}\mathbf{s}\|_2^2$: regularization (or smoothness) functional

\mathbf{L} : regularization operator (i.e., Gradient)

$$\min_{\mathbf{s}} \mathcal{R}(\mathbf{s}) \quad \text{subject to} \quad \|\mathbf{y} - \mathbf{H}\mathbf{s}\|_2^2 \leq \sigma^2$$



Andrey N. Tikhonov (1906-1993)

Equivalent variational problem

$$\mathbf{s}^* = \arg \min \underbrace{\|\mathbf{y} - \mathbf{H}\mathbf{s}\|_2^2}_{\text{data consistency}} + \underbrace{\lambda \|\mathbf{L}\mathbf{s}\|_2^2}_{\text{regularization}}$$

Formal linear solution: $\mathbf{s} = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{L}^T \mathbf{L})^{-1} \mathbf{H}^T \mathbf{y} = \mathbf{R}_\lambda \cdot \mathbf{y}$

Interpretation: “filtered” backprojection

4

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 $f : \mathbb{R}^N \rightarrow \mathbb{R}$ such that $f(\mathbf{x}_m) \approx y_m$ for $m = 1, \dots, M$

- Introduce smoothness or **regularization** constraint (Poggio-Girosi 1990)

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

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

- Regularized least-squares fit

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

5

Unifying continuous-domain formulation

Unknown is a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$

- Regularization functional: $R(f) : \mathcal{B}(\mathbb{R}^d) \rightarrow \mathbb{R}^+$

Promotes smoothness (Sobolev norm) or sparsity (gTV)

- Native space $\mathcal{B}(\mathbb{R}^d)$

Banach vs. Hilbert space (RKHS)

$$\mathcal{B}(\mathbb{R}^d) = \{f : \mathbb{R}^d \rightarrow \mathbb{R} : R(f) < \infty\}$$

- Linear measurement operator $\mathbf{H} : \mathcal{B} \rightarrow \mathbb{R}^M$ **Linear functionals vs. point values**

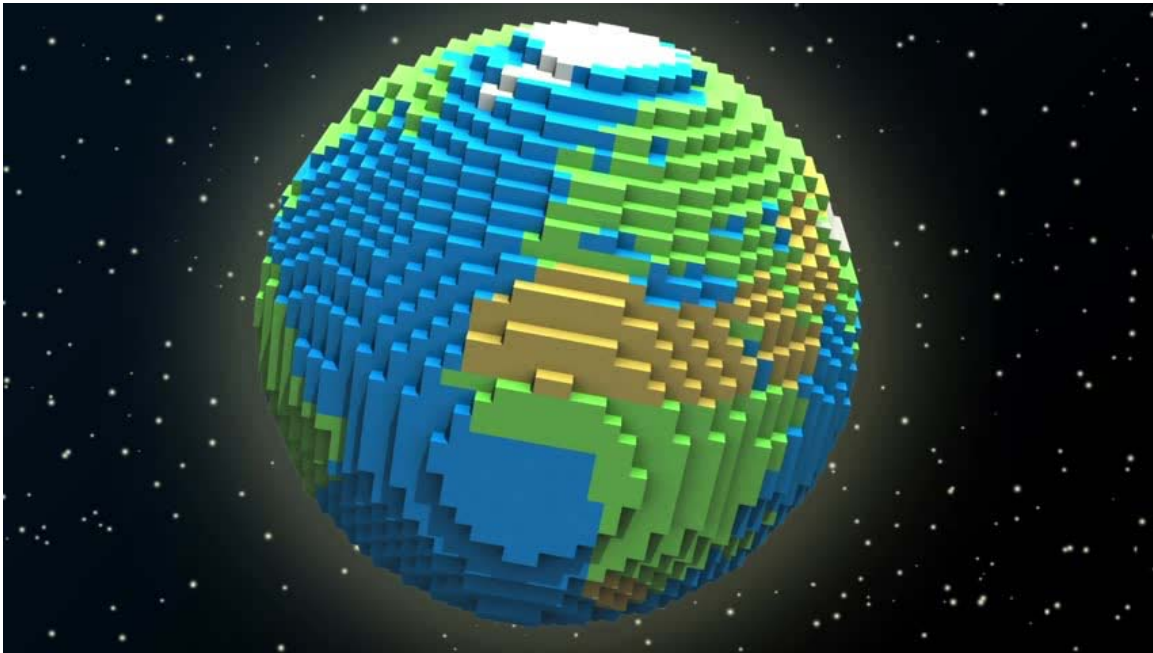
$$\mathbf{H} = (h_1, \dots, h_M) : f \mapsto (\langle h_1, f \rangle, \dots, \langle h_M, f \rangle)$$

- Regularized functional fit to the data **Arbitrary convex loss vs. least squares**

$$f_{\text{opt}} = \arg \min_{f \in \mathcal{B}} \left(\underbrace{\sum_{m=1}^M |y_m - \langle h_m, f \rangle|^2}_{E(\mathbf{y}, \mathbf{H}\{f\})} + \lambda R(f) \right) \quad (\text{Schölkopf 2001; Rosasco 2004})$$

6

Prologue: Discrete-domain regularization



7

Classical least-squares fit with l_2 regularization

- Linear measurement model:

$$y_m = \langle \mathbf{h}_m, \mathbf{x} \rangle + n[m], \quad m = 1, \dots, M$$

- System matrix of size $M \times N$: $\mathbf{H} = [\mathbf{h}_1 \cdots \mathbf{h}_M]^T$

$$\mathbf{x}_{\text{LS}} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_2^2$$

$$\Rightarrow \mathbf{x}_{\text{LS}} = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I}_N)^{-1} \mathbf{H}^T \mathbf{y}$$

$$= \mathbf{H}^T \mathbf{a} = \sum_{m=1}^M a_m \mathbf{h}_m \quad \text{where} \quad \mathbf{a} = (\mathbf{H}\mathbf{H}^T + \lambda \mathbf{I}_M)^{-1} \mathbf{y}$$

Interpretation: $\mathbf{x}_{\text{LS}} \in \text{span}\{\mathbf{h}_m\}_{m=1}^M$

Lemma

$$(\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I}_N)^{-1} \mathbf{H}^T = \mathbf{H}^T (\mathbf{H}\mathbf{H}^T + \lambda \mathbf{I}_M)^{-1}$$

8

Switch to l_1 regularization \Rightarrow sparsifying effect

- Linear measurement model:
 $y_m = \langle \mathbf{h}_m, \mathbf{x} \rangle + n[m], \quad m = 1, \dots, M$

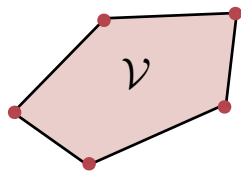
- System matrix of size $M \times N$: $\mathbf{H} = [\mathbf{h}_1 \cdots \mathbf{h}_M]^T$

$$(P1): \quad \mathcal{V} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_{\ell_1}$$

Representer theorem for unconstrained l_1 minimization
 The solution set \mathcal{V} of (P1) is convex, compact with extreme points of the form

$$\mathbf{x}_{\text{sparse}} = \sum_{k=1}^K a_k \mathbf{e}_{n_k} \quad \text{with} \quad K = \|\mathbf{x}_{\text{sparse}}\|_0 \leq M.$$

element of canonical basis with $[\mathbf{e}_n]_m = \delta_{m-n}$



If CS condition on \mathbf{H} is satisfied, then solution is unique

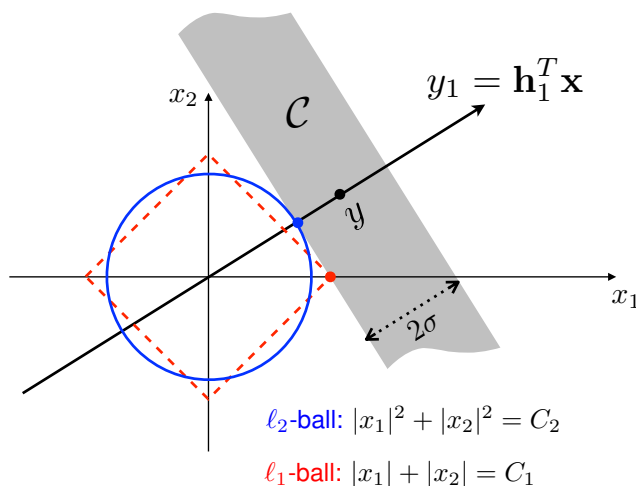
(U.-Fageot-Gupta *IEEE Trans. Info. Theory*, Sept. 2016)

Geometry of l_2 vs. l_1 minimization

- Prototypical inverse problem

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

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

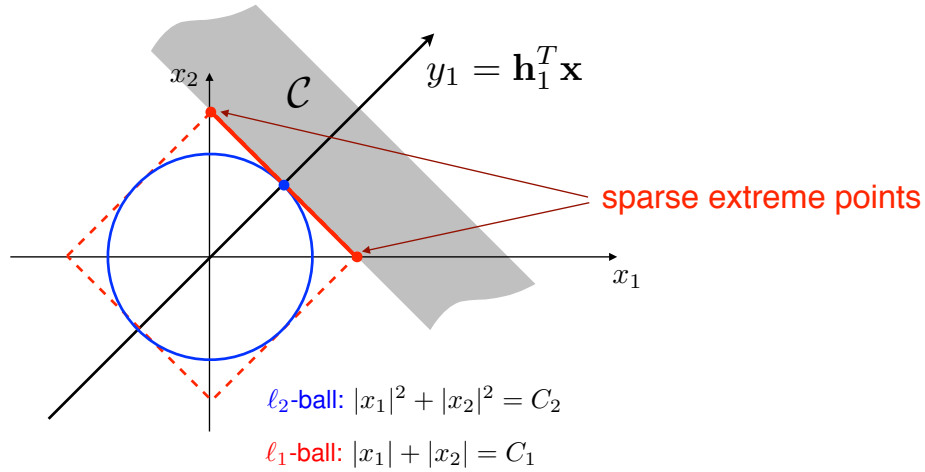


Geometry of l_2 vs. l_1 minimization

■ 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

11

Part II: Continuous-domain theory



12

Continuous-domain regularization (L_2 scenario)

Regularization functional: $\|Lf\|_{L_2}^2 = \int_{\mathbb{R}^d} |Lf(\mathbf{x})|^2 d\mathbf{x}$

L : suitable differential operator

■ Theory of reproducing kernel Hilbert spaces (Aronszajn 1950)

$$\langle f, g \rangle_{\mathcal{H}} = \langle Lf, Lg \rangle$$

■ Interpolation and approximation theory

■ Smoothing splines (Schoenberg 1964, Kimeldorf-Wahba 1971)

■ Thin-plate splines, radial basis functions (Duchon 1977)

■ Machine learning

■ Radial basis functions, kernel methods (Poggio-Girosi 1990)

■ Representer theorem(s) (Schölkopf-Smola 2001)

13

Splines are analog, but intrinsically sparse

$L\{\cdot\}$: admissible differential operator

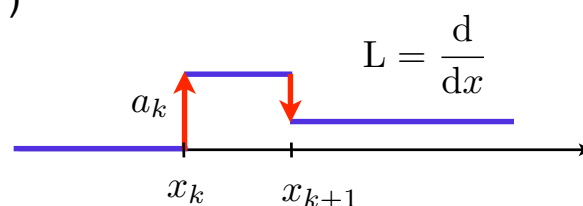
$\delta(\cdot - \mathbf{x}_0)$: Dirac impulse shifted by $\mathbf{x}_0 \in \mathbb{R}^d$

Definition

The function $s : \mathbb{R}^d \rightarrow \mathbb{R}$ is a (non-uniform) L -spline with knots $(\mathbf{x}_k)_{k=1}^K$ if

$$L\{s\} = \sum_{k=1}^K a_k \delta(\cdot - \mathbf{x}_k) = \mathbf{w}_{\delta} \quad : \quad \text{spline's innovation}$$

Spline theory: (Schultz-Varga, 1967)

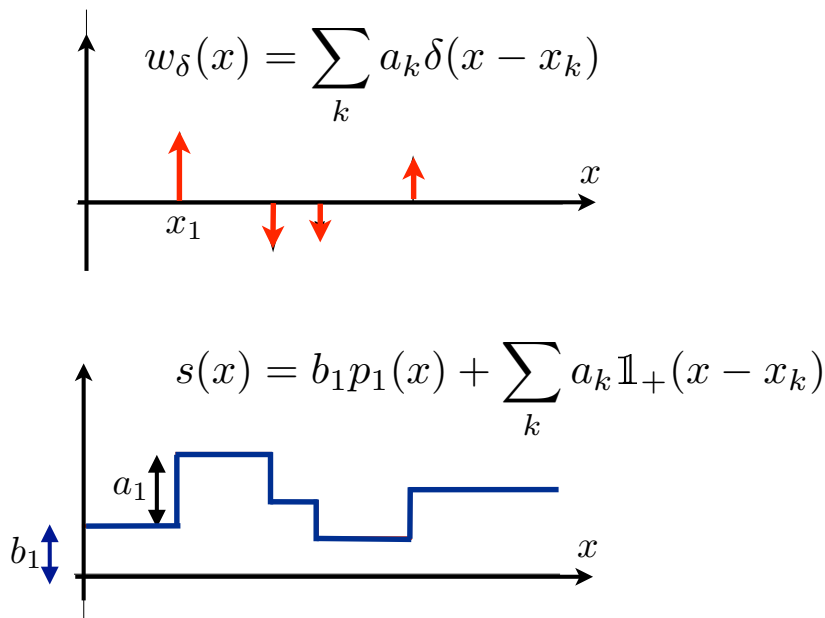


14

Spline synthesis: example

$$L = D = \frac{d}{dx} \quad \text{Null space: } \mathcal{N}_D = \text{span}\{p_1\}, \quad p_1(x) = 1$$

$$\rho_D(x) = D^{-1}\{\delta\}(x) = \mathbb{1}_+(x): \text{Heaviside function}$$



15

Spline synthesis: generalization

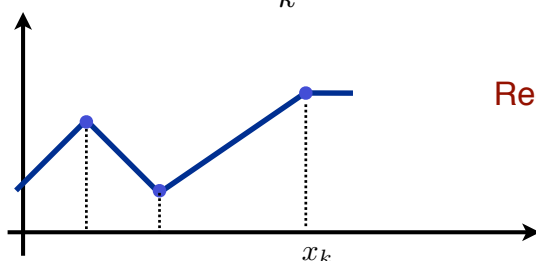
L: spline admissible operator (LSI)

$$\rho_L(\mathbf{x}) = L^{-1}\{\delta\}(\mathbf{x}): \text{Green's function of L}$$

$$\text{Finite-dimensional null space: } \mathcal{N}_L = \text{span}\{p_n\}_{n=1}^{N_0}$$

$$\text{Spline's innovation: } w_\delta(\mathbf{x}) = \sum_k a_k \delta(\mathbf{x} - \mathbf{x}_k)$$

$$\Rightarrow s(\mathbf{x}) = \sum_k a_k \rho_L(\mathbf{x} - \mathbf{x}_k) + \sum_{n=1}^{N_0} b_n p_n(\mathbf{x})$$



Requires specification of boundary conditions

16

RKHS representer theorem for L_2 regularization

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

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

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

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

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

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

(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.

17

L_2 representer theorem for variational splines

Theoretical difficulty: $\|f\|_{\mathcal{H}}^2 \rightarrow \|L f\|_{L_2}^2$ (only a semi-norm !)

$$(P2) \quad \arg \min_{f \in \mathcal{H}_L} \left(\sum_{m=1}^M |y_m - f(\mathbf{x}_m)|^2 + \lambda \|L f\|_{L_2(\mathbb{R}^d)}^2 \right)$$

$\rho_{L^*L}(\mathbf{x}) = (L^*L)^{-1}\{\delta\}(\mathbf{x})$: **Green's function** of (L^*L)

(Schoenberg 1964, Kimeldorf-Wahba 1971)

L_2 representer theorem for variational splines

The solution of (P2) is unique and of the form

$$f(\mathbf{x}) = \sum_{m=1}^M a_m \rho_{L^*L}(\mathbf{x} - \mathbf{x}_m) + \sum_{n=1}^{N_0} b_n p_n(\mathbf{x});$$

i.e., it is a (L^*L) -spline with knots at the $\{\mathbf{x}_m\}$.

Example: $L = D^2$ with $\rho_{D^4}(x) \propto |x|^3 \Rightarrow f(x)$ is a cubic spline

18

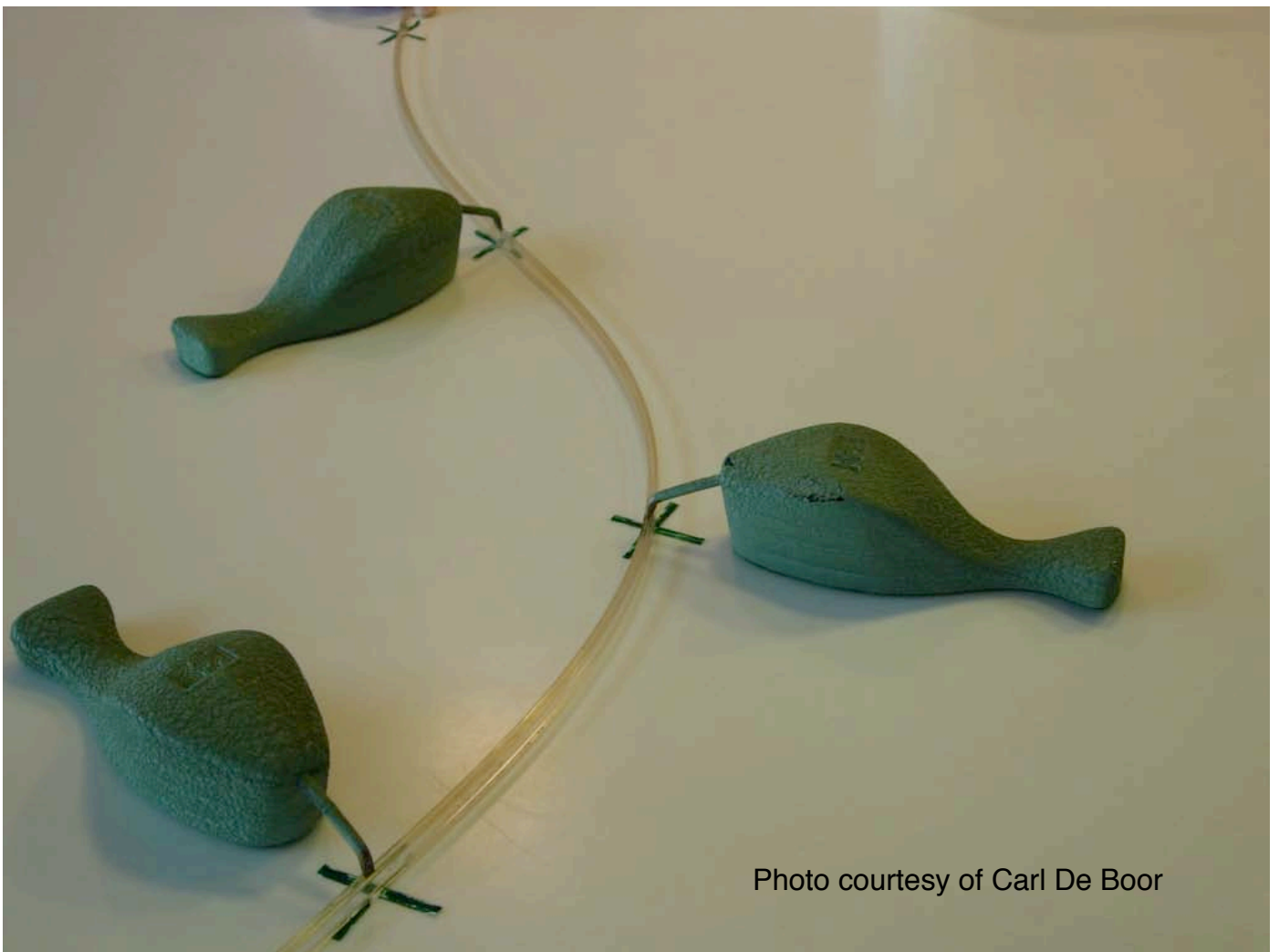


Photo courtesy of Carl De Boor

Quest for sparsity
in a continuous world

Sparsity and continuous-domain modeling

■ Compressed sensing (CS)

- Generalized sampling and infinite-dimensional CS (Adcock-Hansen 2011)
- Sampling: CS of analog signals (Eldar 2011)
- Recovery of Dirac impulses from Fourier measurements (Vetterli et al. 2002)
(Bredies 2013; Candes & Fernandez-Granda 2014; Duval-Peyré 2015)

■ Splines and approximation theory

- L_1 splines (Fisher-Jerome 1975)
- Locally-adaptive regression splines (Mammen-van de Geer 1997)
- Generalized TV (Steidl et al. 2005; Bredies et al. 2010)

■ Statistical modeling

- Sparse stochastic processes (Unser et al. 2011-2014)

21

Proper continuous counterpart of $\ell_1(\mathbb{Z}^d)$

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

$\mathcal{S}'(\mathbb{R}^d)$: Schwartz's space of tempered distributions

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

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

where $w : \varphi \mapsto \langle w, \varphi \rangle = \int_{\mathbb{R}^d} \varphi(\mathbf{r})w(\mathbf{r})d\mathbf{r}$

■ Equivalent definition of "total variation" norm

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

■ Basic inclusions

- $\delta(\cdot - \mathbf{x}_0) \in \mathcal{M}(\mathbb{R}^d)$ with $\|\delta(\cdot - \mathbf{x}_0)\|_{\mathcal{M}} = 1$ for any $\mathbf{x}_0 \in \mathbb{R}^d$
- $\|f\|_{\mathcal{M}} = \|f\|_{L_1(\mathbb{R}^d)}$ for all $f \in L_1(\mathbb{R}^d) \Rightarrow L_1(\mathbb{R}^d) \subseteq \mathcal{M}(\mathbb{R}^d)$

22

Representer theorem for gTV regularization

$$(P1) \quad \arg \min_{f \in \mathcal{M}_L(\mathbb{R}^d)} \left(\sum_{m=1}^M |y_m - \langle h_m, f \rangle|^2 + \lambda \|Lf\|_{\mathcal{M}} \right)$$

- L: spline-admissible operator with null space $\mathcal{N}_L = \text{span}\{p_n\}_{n=1}^{N_0}$
- gTV semi-norm: $\|L\{s\}\|_{\mathcal{M}} = \sup_{\|\varphi\|_{\infty} \leq 1} \langle L\{s\}, \varphi \rangle$
- Measurement functionals $h_m : \mathcal{M}_L(\mathbb{R}^d) \rightarrow \mathbb{R}$ (weak*-continuous)

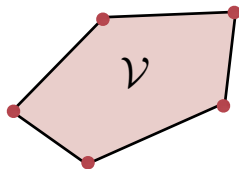
Convex loss function: $F : \mathbb{R}^M \times \mathbb{R}^M \rightarrow \mathbb{R}$ $\nu : \mathcal{M}_L \rightarrow \mathbb{R}^M$

$$(P1') \quad \arg \min_{f \in \mathcal{M}_L(\mathbb{R}^d)} (F(\mathbf{y}, \nu(f)) + \lambda \|Lf\|_{\mathcal{M}}) \quad \text{with } \nu(f) = (\langle h_1, f \rangle, \dots, \langle h_M, f \rangle)$$

Representer theorem for gTV-regularization
 The extreme points of (P1') are **non-uniform L-spline** of the form

$$f_{\text{spline}}(\mathbf{x}) = \sum_{k=1}^{K_{\text{knots}}} a_k \rho_L(\mathbf{x} - \mathbf{x}_k) + \sum_{n=1}^{N_0} b_n p_n(\mathbf{x})$$

with ρ_L such that $L\{\rho_L\} = \delta$, $K_{\text{knots}} \leq M - N_0$, and $\|Lf_{\text{spline}}\|_{\mathcal{M}} = \|\mathbf{a}\|_{\ell_1}$.



(U.-Fageot-Ward, *SIAM Review* 2017)

Example: 1D inverse problem with TV⁽²⁾ regularization

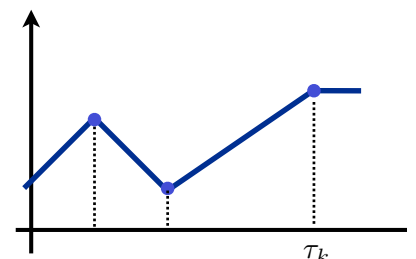
$$s_{\text{spline}} = \arg \min_{s \in \mathcal{M}_D^2(\mathbb{R})} \left(\sum_{m=1}^M |y_m - \langle h_m, s \rangle|^2 + \lambda \text{TV}^{(2)}(s) \right)$$

- Total 2nd-variation: $\text{TV}^{(2)}(s) = \sup_{\|\varphi\|_{\infty} \leq 1} \langle D^2 s, \varphi \rangle = \|D^2 s\|_{\mathcal{M}}$

$$L = D^2 = \frac{d^2}{dx^2} \quad \rho_{D^2}(x) = (x)_+ : \text{ReLU} \quad \mathcal{N}_{D^2} = \text{span}\{1, x\}$$

- Generic form of the solution

$$s_{\text{spline}}(x) = \underbrace{b_1 + b_2 x}_{\text{no penalty}} + \sum_{k=1}^K a_k (x - \tau_k)_+$$



with $K < M$ and free parameters b_1, b_2 and $(a_k, \tau_k)_{k=1}^K$

Other spline-admissible operators

- $L = D^n$ (pure derivatives)
 - ⇒ polynomial splines of degree $(n - 1)$ (Schoenberg 1946)
- $L = D^n + a_{n-1}D^{n-1} + \dots + a_0I$ (ordinary differential operator)
 - ⇒ exponential splines (Dahmen-Micchelli 1987)
- Fractional derivatives: $L = D^\gamma \xleftrightarrow{\mathcal{F}} (j\omega)^\gamma$
 - ⇒ fractional splines (U.-Blu 2000)
- Fractional Laplacian: $(-\Delta)^{\frac{\gamma}{2}} \xleftrightarrow{\mathcal{F}} \|\omega\|^\gamma$
 - ⇒ polyharmonic splines (Duchon 1977)
- Elliptical differential operators; e.g, $L = (-\Delta + \alpha I)^\gamma$
 - ⇒ Sobolev splines (Ward-U. 2014)

25

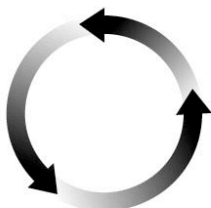
Discretization: compatible with CS paradigm

$$\mathbf{s}_{\text{sparse}} = \arg \min_{\mathbf{s} \in \mathbb{R}^K} \left(\frac{1}{2} \|\mathbf{y} - \mathbf{H}\mathbf{s}\|_2^2 + \lambda \|\mathbf{u}\|_1 \right) \text{ subject to } \mathbf{u} = \mathbf{L}\mathbf{s}$$

ADMM algorithm

$$\mathcal{L}_{\mathcal{A}}(\mathbf{s}, \mathbf{u}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{y} - \mathbf{H}\mathbf{s}\|_2^2 + \lambda \sum_n |[\mathbf{u}]_n| + \boldsymbol{\alpha}^T (\mathbf{L}\mathbf{s} - \mathbf{u}) + \frac{\mu}{2} \|\mathbf{L}\mathbf{s} - \mathbf{u}\|_2^2$$

For $k = 0, \dots, K$

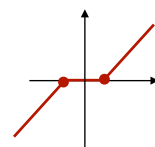


Linear step

$$\begin{aligned} \mathbf{s}^{k+1} &= (\mathbf{H}^T \mathbf{H} + \mu \mathbf{L}^T \mathbf{L})^{-1} (\mathbf{z}_0 + \mathbf{z}^{k+1}) \\ &\text{with } \mathbf{z}^{k+1} = \mathbf{L}^T (\mu \mathbf{u}^k - \boldsymbol{\alpha}^k) \\ \boldsymbol{\alpha}^{k+1} &= \boldsymbol{\alpha}^k + \mu (\mathbf{L}\mathbf{s}^{k+1} - \mathbf{u}^k) \end{aligned}$$

Proximal step = pointwise non-linearity

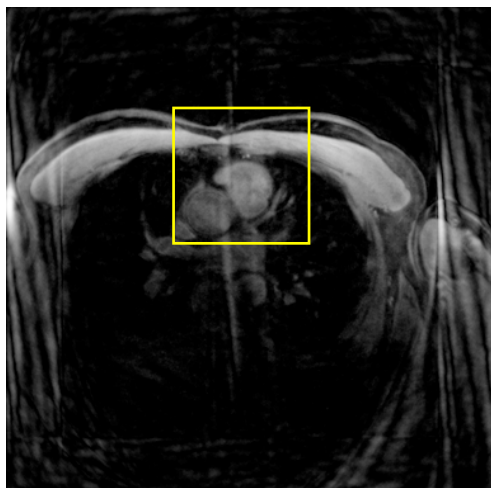
$$\mathbf{u}^{k+1} = \text{prox}_{|\cdot|} \left(\mathbf{L}\mathbf{s}^{k+1} + \frac{1}{\mu} \boldsymbol{\alpha}^{k+1}; \frac{\sigma^2}{\mu} \right)$$



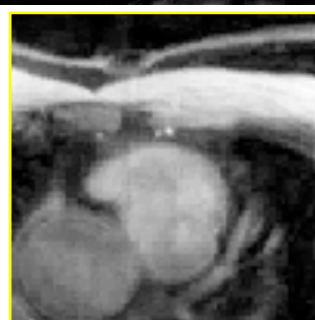
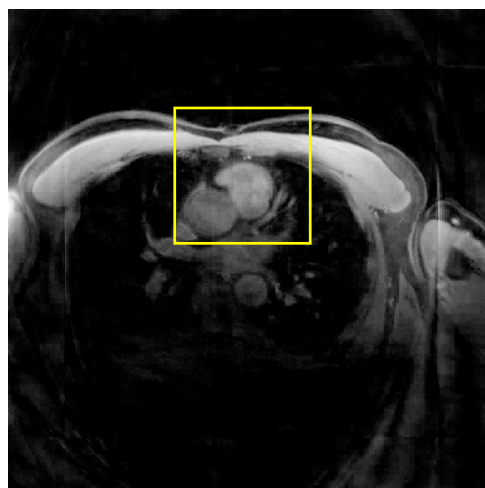
26

Example: ISMRM reconstruction challenge

L_2 regularization (Laplacian)



ℓ_1 / TV regularization



(Guerquin-Kern *IEEE TMI* 2011)

27

OUTLINE

- Linear inverse problems and regularization ✓
- Continuous-domain theory ✓
 - Splines and operators
 - Classical L_2 regularization: theory of RKHS
 - Minimization of gTV: the optimality of splines
- **From compressed sensing to deep networks**
 - Image recovery with sparsity constraints
 - FBPCnvNet
 - Representer theorem for deep neural networks

When is unrolled ADMM a deep ConvNet ?

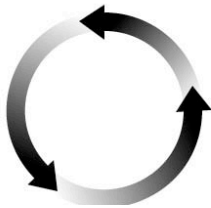
Answer: when $\mathbf{H}^T \mathbf{H}$ and \mathbf{L} are both convolutions

$$\mathcal{L}_{\mathcal{A}}(\mathbf{s}, \mathbf{u}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{y} - \mathbf{H}\mathbf{s}\|_2^2 + \sigma^2 \sum_n |[\mathbf{u}]_n| + \boldsymbol{\alpha}^T (\mathbf{L}\mathbf{s} - \mathbf{u}) + \frac{\mu}{2} \|\mathbf{L}\mathbf{s} - \mathbf{u}\|_2^2$$

ADMM algorithm

Initialization	
$\mathbf{z}_0 = \mathbf{H}^T \mathbf{y}$	$\mathbf{u}^0 = \mathbf{0}$
$\mathbf{s}^0 = \mathbf{0}$	$\boldsymbol{\alpha}^0 = \mathbf{0}$

For $k = 0, \dots, K$



Linear step = Convolutions

$$\mathbf{s}^{k+1} = (\mathbf{H}^T \mathbf{H} + \mu \mathbf{L}^T \mathbf{L})^{-1} (\mathbf{z}_0 + \mathbf{z}^{k+1})$$

with $\mathbf{z}^{k+1} = \mathbf{L}^T (\mu \mathbf{u}^k - \boldsymbol{\alpha}^k)$

$$\boldsymbol{\alpha}^{k+1} = \boldsymbol{\alpha}^k + \mu (\mathbf{L}\mathbf{s}^{k+1} - \mathbf{u}^k)$$

Proximal step = pointwise non-linearity

$$\mathbf{u}^{k+1} = \text{prox}_{|\cdot|} (\mathbf{L}\mathbf{s}^{k+1} + \frac{1}{\mu} \boldsymbol{\alpha}^{k+1}; \frac{\sigma^2}{\mu})$$

29

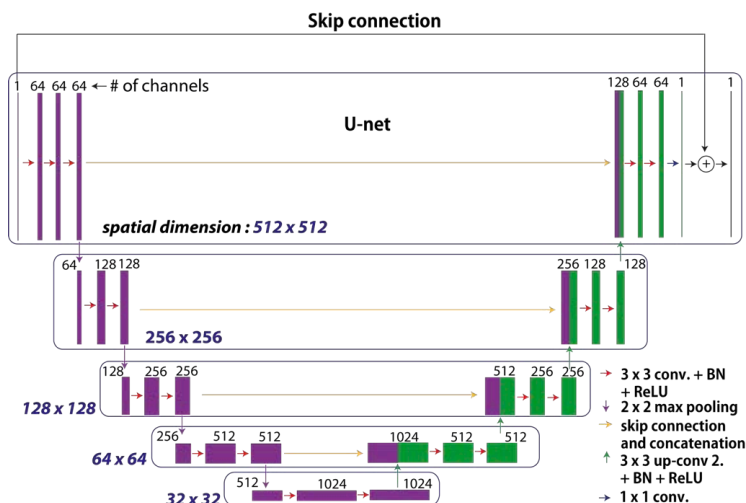
Recent appearance of Deep ConvNets

(Jin et al. 2016; Adler-Öktem 2017; Chen et al. 2017; ...)

■ CT reconstruction based on Deep ConvNets

- Input: Sparse view FBP reconstruction
- Training: Set of 500 high-quality full-view CT reconstructions
- Architecture: U-Net with skip connection

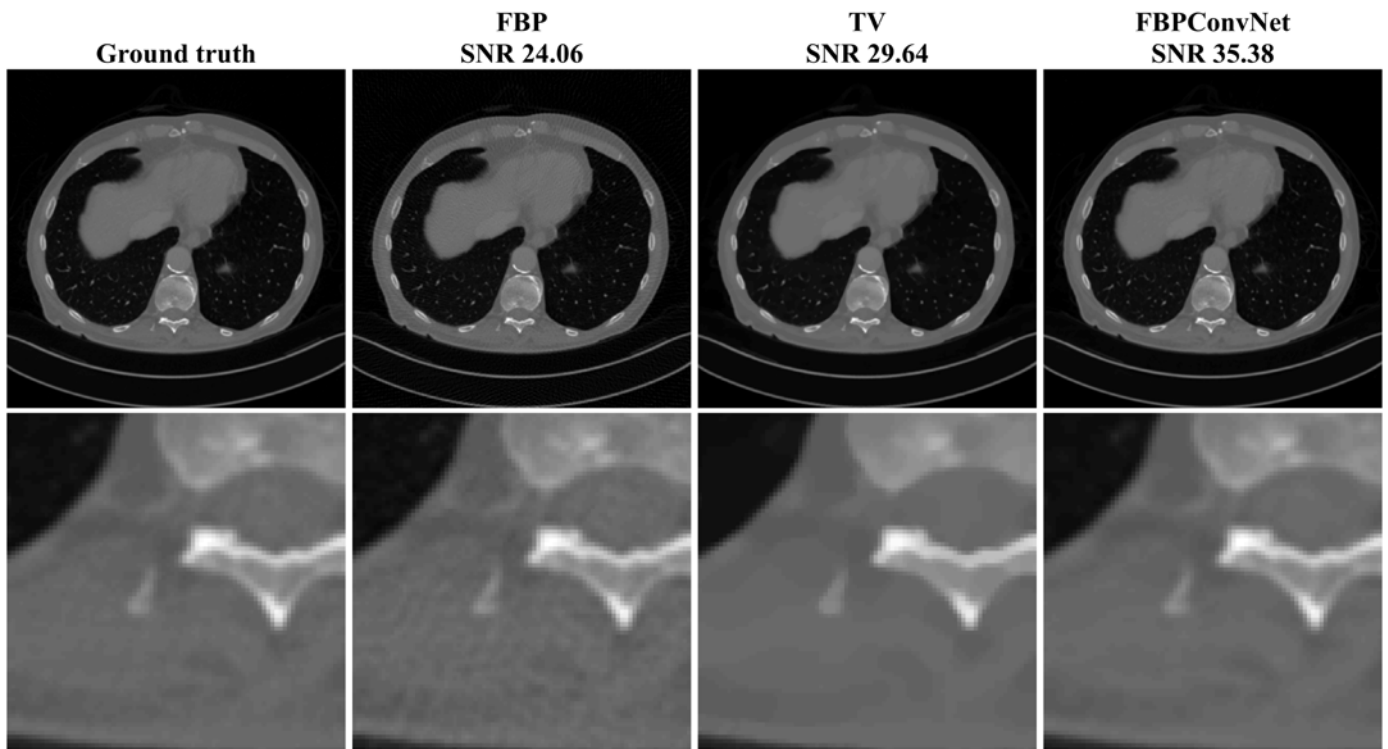
(Jin et al., IEEE TIP 2017)



30

CT data

Dose reduction by 7: 143 views



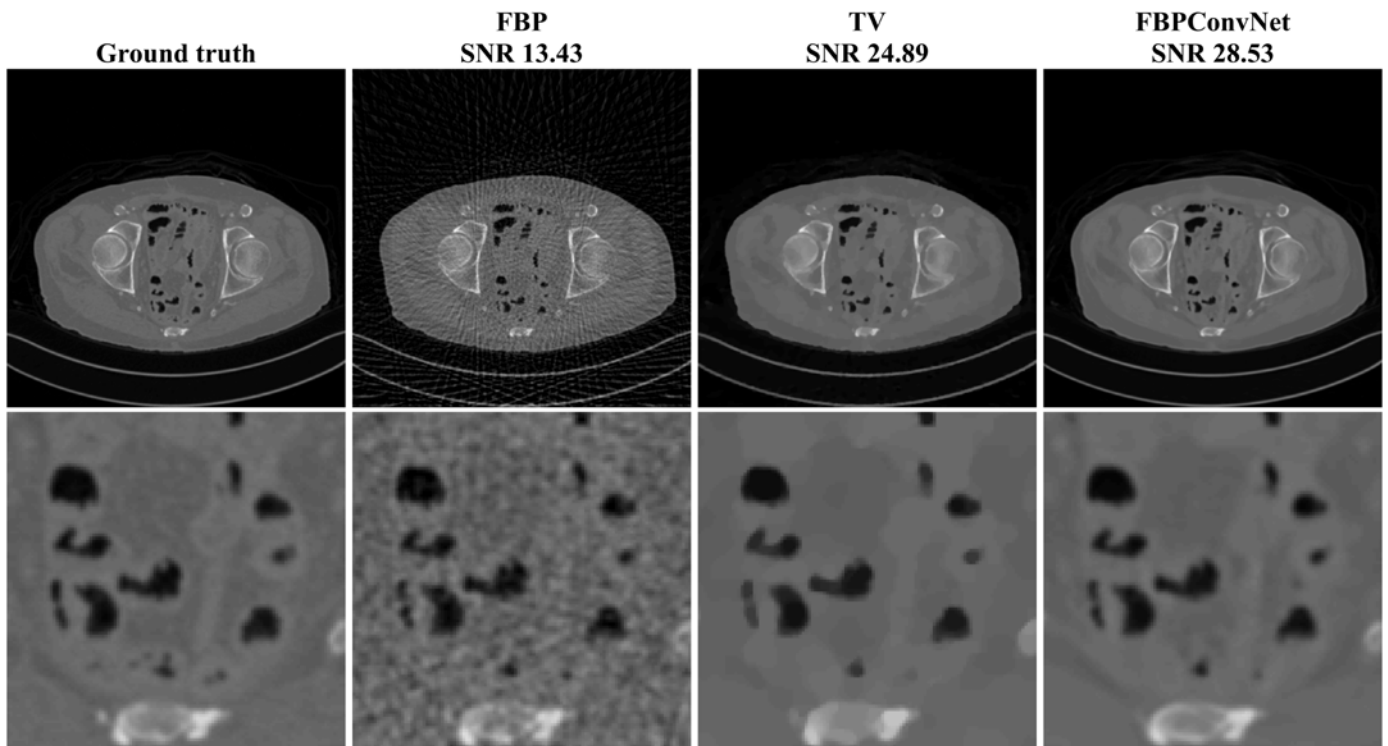
Reconstructed from
from 1000 views

(Jin et al., *IEEE Trans. Im Proc.*, 2017)



CT data

Dose reduction by 20: 50 views



Reconstructed from
from 1000 views

(Jin et al., *IEEE Trans. Im Proc.*, 2017)





Finale:

Representer theorem for deep learning

33

Deep neural networks and **splines**

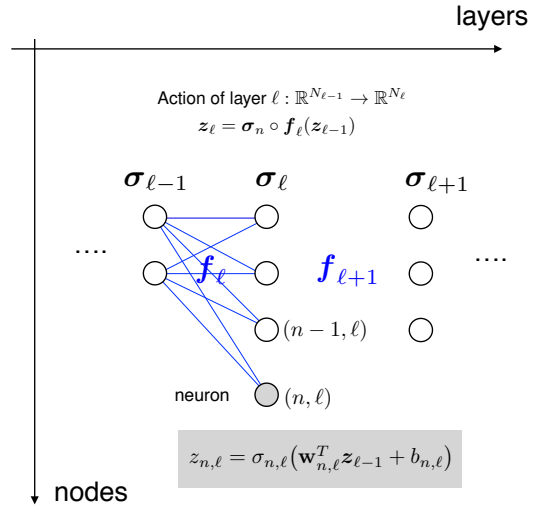
- Preferred choice of activation function: ReLU
 - ReLU works nicely with dropout / ℓ_1 -regularization (Glorot *ICAI*S 2011)
 - Networks with hidden ReLU are easier to train
 - State-of-the-art performance (LeCun-Bengio-Hinton *Nature* 2015)

- Deep nets as Continuous PieceWise-Linear maps
 - MaxOut \Rightarrow CPWL (Goodfellow *PMLR* 2013)
 - ReLU \Rightarrow CPWL (Montufar *NIPS* 2014)
 - CPWL \Rightarrow Deep ReLU network (Wang-Sun *IEEE-IT* 2005)

34

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) , $n = 1, \dots, N_\ell$
- Activations functions: $\sigma_{n,\ell} : \mathbb{R} \rightarrow \mathbb{R}$
- 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: $\mathbb{R}^{N_\ell} \rightarrow \mathbb{R}^{N_\ell}$
 $\sigma_\ell : \mathbf{x} \mapsto \sigma_\ell(\mathbf{x}) = (\sigma_{1,\ell}(x_1), \dots, \sigma_{N_\ell,\ell}(x_{N_\ell}))$



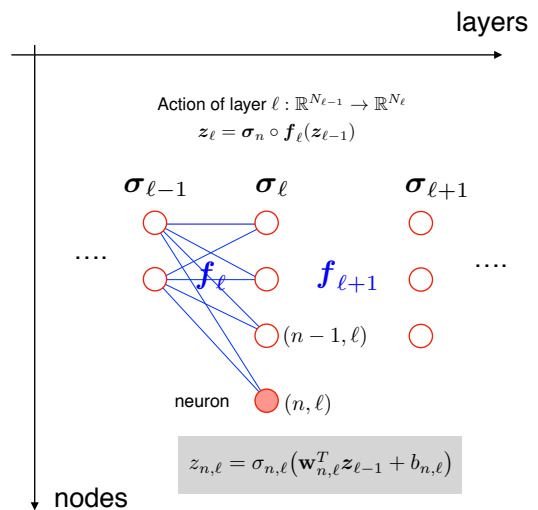
Conventional design: $\sigma_{n,\ell} = \sigma$

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

35

Deep neural net with optimized activations

- Layers: $\ell = 1, \dots, L$
- Deep structure descriptor: (N_0, N_1, \dots, N_L)
- Neuron or node index: (n, ℓ) , $n = 1, \dots, N_\ell$
- Activations functions: $\sigma_{n,\ell} : \mathbb{R} \rightarrow \mathbb{R}$
- 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: $\mathbb{R}^{N_\ell} \rightarrow \mathbb{R}^{N_\ell}$
 $\sigma_\ell : \mathbf{x} \mapsto \sigma_\ell(\mathbf{x}) = (\sigma_{1,\ell}(x_1), \dots, \sigma_{N_\ell,\ell}(x_{N_\ell}))$



New **adaptive** design: $x \mapsto \sigma_{n,\ell}(x)$ s.t. $\text{TV}^{(2)}(\sigma_{n,\ell})$ minimum

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

36

New representer theorem for deep neural networks

(Unser, *arXiv:1802.09210*, Feb 2018)

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

- 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 of M data points $\mathbf{y}_m \approx \mathbf{f}(\mathbf{x}_m)$, 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 $b_{1,n,\ell}, b_{2,n,\ell}, a_{1,n,\ell}, \dots, a_{K_{n,\ell},n,\ell} \in \mathbb{R}$.

37

Outcome of representer theorem

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

- its number $0 \leq K = K_{n,\ell}$ 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 $\mathbf{b}_{n,\ell} = (b_{1,n,\ell}, b_{2,n,\ell}) \in \mathbb{R}^2$,
 $\mathbf{a}_{n,\ell} = (a_{1,n,\ell}, \dots, a_{K_{n,\ell},n,\ell}) \in \mathbb{R}^K$.

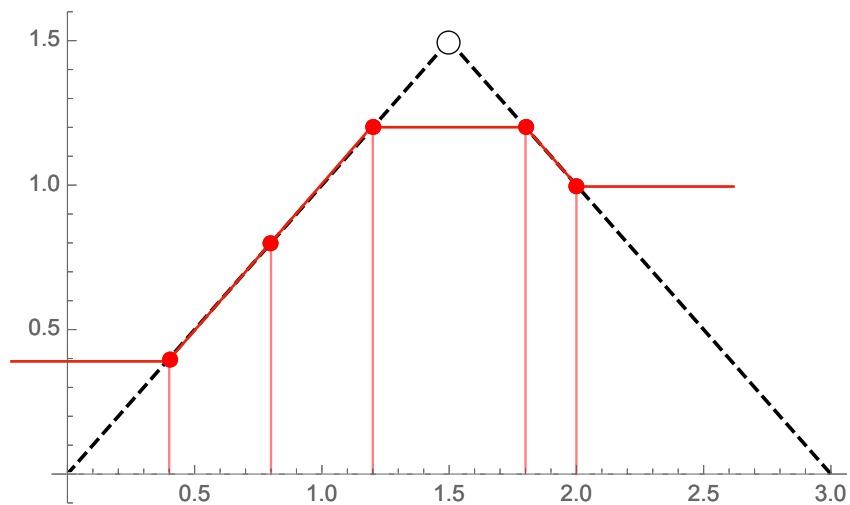
These parameters (including the number of knots) are **data-dependent** and need to be adjusted automatically during **training**.

- Link with ℓ_1 minimization techniques

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

38

Comparison of linear interpolators



39

Deep spline networks: Discussion

- Global optimality achieved with **spline activations**

- State-of-the-art ReLU networks $(K_{n,\ell} = 1, \mathbf{b}_{n,\ell} = \mathbf{0})$

- No need to normalize:

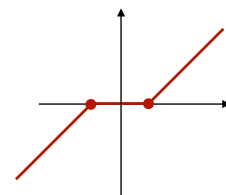
$$(\mathbf{w}_{n,\ell}^T \mathbf{x} - z_{n,\ell})_+ = (a_{n,\ell} \mathbf{u}_{n,\ell}^T \mathbf{x} - z_{n,\ell})_+ = a_{n,\ell} (\mathbf{u}_{n,\ell}^T \mathbf{x} - \tau_{n,\ell})_+$$

- Key features

- Produces a global mapping $x \mapsto \mathbf{f}(x)$ that is **continuous** and **piecewise-linear**
 - Direct control of complexity (number of knots): adjustment of λ
 - Ability to suppress unnecessary layers

- Backward compatibility

- Linear regression: $\lambda \rightarrow \infty \Rightarrow K_{n,\ell} = 0$
 - Compressed sensing / ℓ_1 minimization



40

SUMMARY: Controlling smoothness vs. sparsity

- New findings resonate with what is known in discrete setting
 - l_2 solution lives in a **fixed** subspace of dimension M
 - Tikhonov solution is intrinsically “**blurred**”
 - Minimization of l_1 favors sparse solutions (independently of sensing matrix)
- Specificities of continuous-domain formulation $s \in \mathcal{X}$
 - Functional model: class of signals + physics $s \mapsto \mathbf{z} = \mathbf{H}\{s\}$
 - Smoothing-splines: minimum “spline” energy $(L^*L)\{s_{\text{smooth}}\} = \sum_{m=1}^M a_m h_m$
 - L-splines = signals with “sparsest” innovation $L\{s_{\text{sparse}}\} = \sum_{k=1}^K a_k \delta(\cdot - \mathbf{x}_k)$
- Practical implications
 - Infinite-dimensional optimization is feasible (parametric form of solution)
 - gTV regularization favors **sparse** innovations with **adaptive** knots
 - Non-uniform L-splines: **universal** solutions of linear inverse problems

and deep neural networks ...

41

Acknowledgments

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

- Dr. Pouya Tafti
- Prof. Arash Amini
- Dr. Emrah Bostan
- Dr. Masih Nilchian
- Dr. Ulugbek Kamilov
- Dr. Cédric Vonesch
-



and collaborators ...

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



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

42

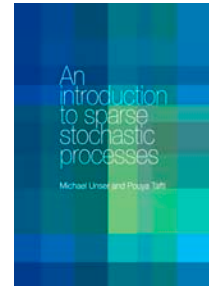
References

■ Theoretical results on sparsity-promoting regularization

- M. Unser, J. Fageot, H. Gupta, “Representer Theorems for Sparsity-Promoting ℓ_1 Regularization,” *IEEE Trans. Information Theory*, Vol. 62, No. 9, pp. 5167-5180.
- 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.

■ Theory of sparse stochastic processes

- M. Unser and P. Tafti, ***An Introduction to Sparse Stochastic Processes***, Cambridge University Press, 2014.
Preprint, available at <http://www.sparseprocesses.org>.
- **For splines:** see chapter 6
- E. Bostan, U.S. Kamilov, M. Nilchian, M. Unser, “Sparse Stochastic Processes and Discretization of Linear Inverse Problems,” *IEEE Trans. Image Processing*, vol. 22, no. 7, pp. 2699-2710, 2013.



■ Deep 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.
- M.T. McCann, K.H. Jin, M. Unser, “Convolutional Neural Networks for Inverse Problems in Imaging—A Review,” *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 85-95, Nov. 2017.
- M. Unser, “A representer theorem for deep neural networks,” preprint, Feb 2018, arXiv:1802.09210.