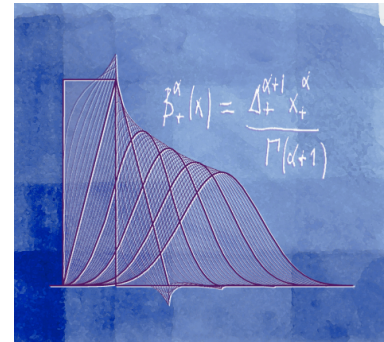


Splines and imaging: From compressed sensing to deep neural nets

Prof. Michael Unser, LIB

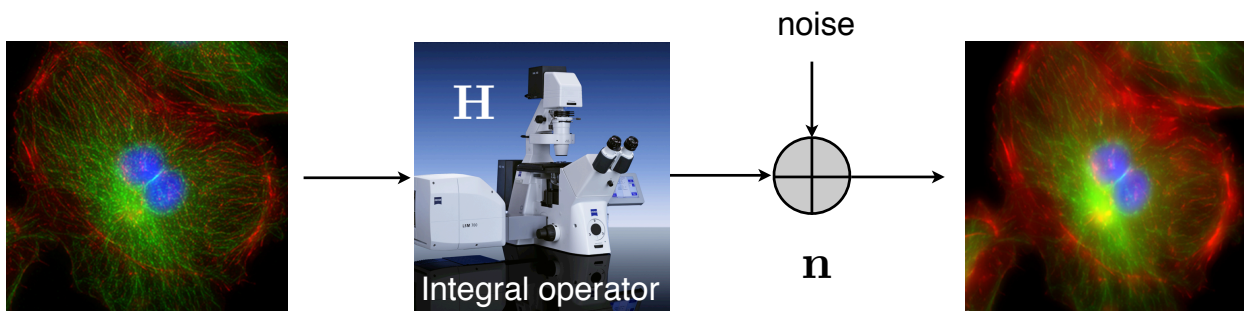


Deep Learning and Medical Imaging, IPAM, UCLA, January 27-31, 2020

Variational formulation of inverse problem

- Linear forward model

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



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

- Reconstruction as an optimization problem

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

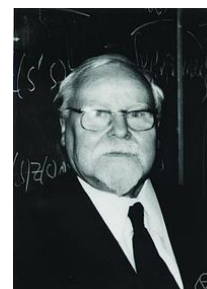
Linear inverse problems (20th century theory)

- Dealing with **ill-posed problems**: Tikhonov **regularization**

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

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

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



Andrey N. Tikhonov (1906-1993)

- Equivalent variational problem

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

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

Interpretation: “**filtered**” backprojection

3

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)

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

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

(Wahba 1990; Schölkopf 2001)

⇒ kernel estimator

4

OUTLINE

■ Introduction ✓

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

■ Continuous-domain theory of sparsity

- Splines and operators
- gTV regularization: representer theorem for CS

■ From compressed sensing to deep neural networks

- Unrolling forward/backward iterations: FBPCConv

■ Deep neural networks vs. deep splines

- Continuous piecewise linear (CPWL) functions / splines
- New representer theorem for deep neural networks



SWISS NATIONAL SCIENCE FOUNDATION



5

Part I: Continuous-domain theory of sparsity



L_1 splines

(Fisher-Jerome 1975)

gTV optimality of splines for inverse problems

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

6

Splines are analog, but intrinsically sparse

$L\{\cdot\}$: differential operator (translation-invariant)

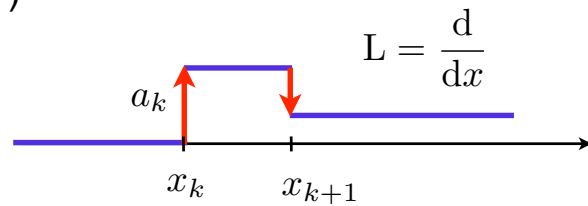
δ : Dirac distribution

Definition

The function $s(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$ (possibly of slow growth) is a **nonuniform L-spline** with knots $\{\mathbf{x}_k\}_{k \in S}$

$$\Leftrightarrow Ls = \sum_{k \in S} a_k \delta(\cdot - \mathbf{x}_k) = w : \text{ spline's innovation}$$

Spline theory: (Schultz-Varga, 1967)

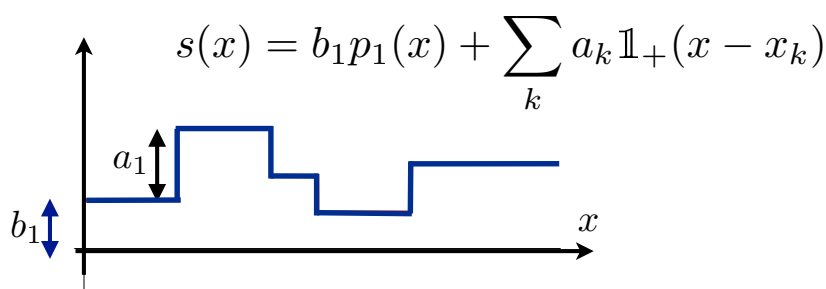
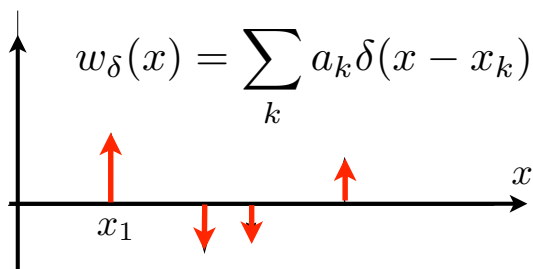


7

Spline synthesis: example

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

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



8

Spline synthesis: generalization

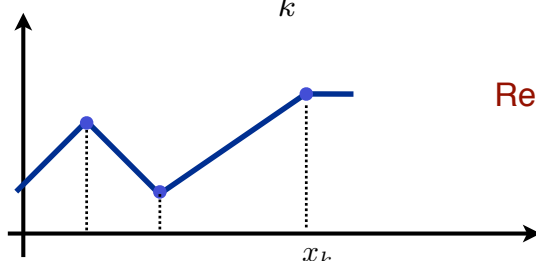
L: spline admissible operator (LSI)

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

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

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

9

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 real-valued **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 $f \in L_1(\mathbb{R}^d) \Rightarrow L_1(\mathbb{R}^d) \subseteq \mathcal{M}(\mathbb{R}^d)$

10

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

$\nu : \mathcal{M}_L \rightarrow \mathbb{R}^M$

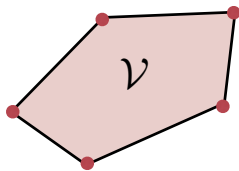
$$(P1') \quad \arg \min_{f \in \mathcal{M}_L(\mathbb{R}^d)} (F(\mathbf{y}, \nu(f)) + \lambda \|Lf\|_{\mathcal{M}}) \quad \text{with } \nu(f) = (\langle h_1, f \rangle, \dots, \langle h_M, f \rangle)$$

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)

11

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

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

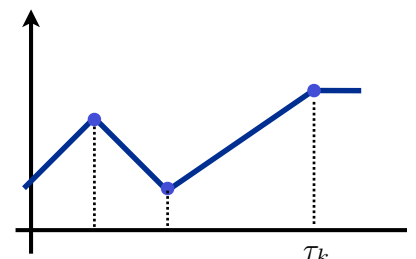
- Total 2nd-variation: $\text{TV}^{(2)}(s) = \sup_{\|\varphi\|_{\infty} \leq 1} \langle D^2 s, \varphi \rangle = \|D^2 s\|_{\mathcal{M}}$

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

- Generic form of the solution

$$s_{\text{spline}}(x) = b_1 + b_2 x + \sum_{k=1}^K a_k (x - \tau_k)_+$$

↗
no penalty



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

12

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: $(-\Delta)^{\frac{\gamma}{2}} \xleftrightarrow{\mathcal{F}} \|\boldsymbol{\omega}\|^\gamma$
 - ⇒ polyharmonic splines (Duchon 1977)
- Elliptical differential operators; e.g., $L = (-\Delta + \alpha I)^\gamma$
 - ⇒ Sobolev splines (Ward-U. 2014)

13

Recovery with sparsity constraints: discretization

- Constrained optimization formulation

Auxiliary **innovation** variable: $\mathbf{u} = \mathbf{L}\mathbf{s}$

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

- Augmented Lagrangian method

Quadratic penalty term: $\frac{\mu}{2} \|\mathbf{L}\mathbf{s} - \mathbf{u}\|_2^2$

Lagrange multiplier vector: $\boldsymbol{\alpha}$

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

(Ramani-Fessler, *IEEE TMI* 2011)

14

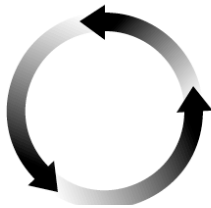
Discretization: compatible with CS paradigm

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

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 $k = 0, \dots, K$

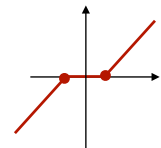


Linear step

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

Proximal step = pointwise non-linearity

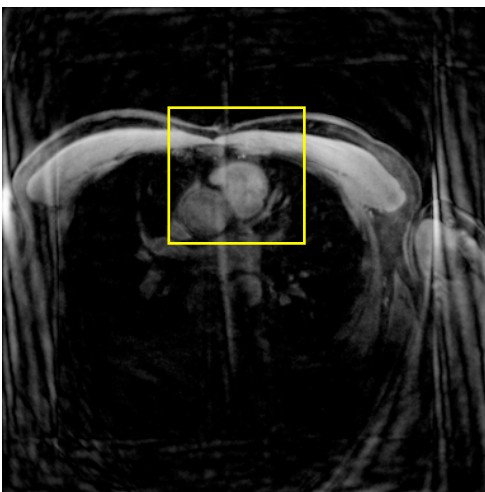
$$\mathbf{u}^{k+1} = \text{prox}_{|\cdot|} \left(\mathbf{L}\mathbf{s}^{k+1} + \frac{1}{\mu} \boldsymbol{\alpha}^{k+1}; \frac{\lambda}{\mu} \right)$$



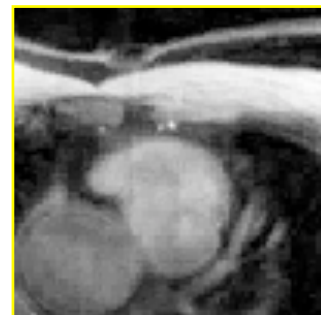
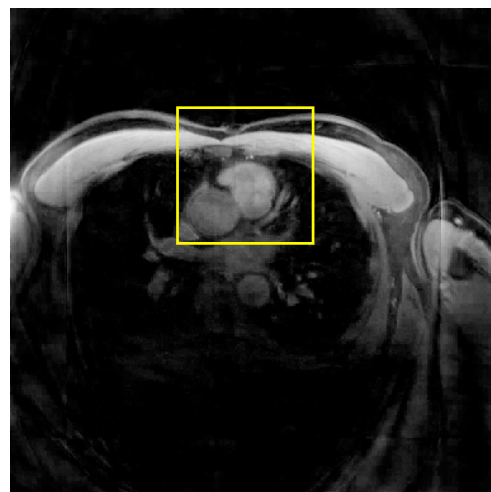
15

Example: ISMRM reconstruction challenge

L_2 regularization (Laplacian)



ℓ_1 / TV regularization



OUTLINE

- Introduction ✓
- Continuous-domain theory of sparsity ✓
- **From compressed sensing to deep neural networks**
 - Unrolling forward/backward iterations: FBPCConv
- **Deep neural networks vs. deep splines**
 - Continuous piecewise linear (CPWL) functions / splines
 - New representer theorem for deep neural networks

17

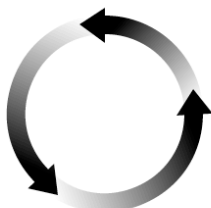
Structure of iterative reconstruction algorithm

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

ADMM

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

For $k = 0, \dots, K$

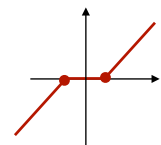


Linear step

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

Pointwise nonlinearity

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



18

Identification of convolution operators

Normal matrix: $\mathbf{A} = \mathbf{H}^T \mathbf{H}$ (symmetric)

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

■ Recognizing structured matrices

- \mathbf{L} : convolution matrix $\Rightarrow \mathbf{L}^T \mathbf{L}$: symmetric convolution matrix
- \mathbf{L}, \mathbf{A} : convolution matrices $\Rightarrow (\mathbf{A} + \lambda \mathbf{L}^T \mathbf{L})$: symmetric convolution matrix

- Applicable to

<ul style="list-style-type: none"> - deconvolution microscopy (Wiener filter) - parallel rays computer tomography (FBP) - MRI, including non-uniform sampling of k-space

■ Fast FFT-based implementation

■ Justification for use of convolution neural nets (CNN)

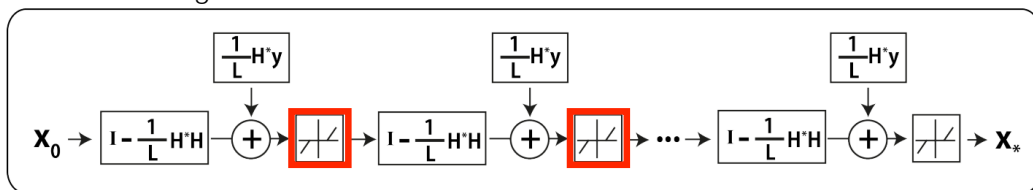
(see Theorem 1, Jin et al., *IEEE TIP* 2017)

Connection with deep neural networks

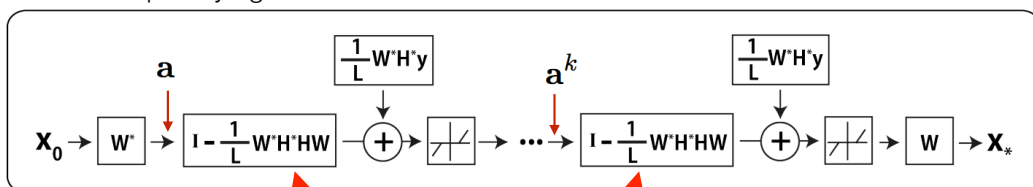
(Gregor-LeCun 2010)

Unrolled Iterative Shrinkage Thresholding Algorithm (ISTA)

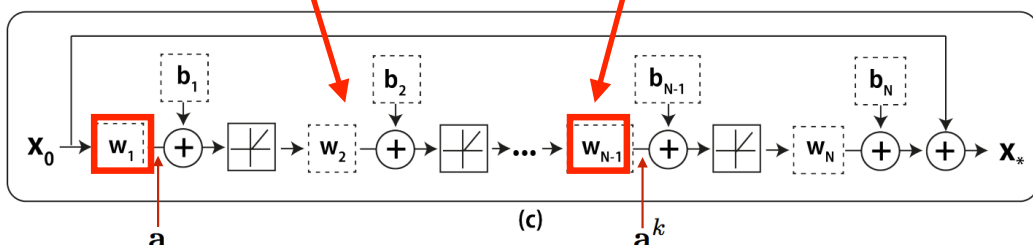
LISTA : learning-based ISTA



ISTA with sparsifying transformation (a)



FBPConvNet structures (b)



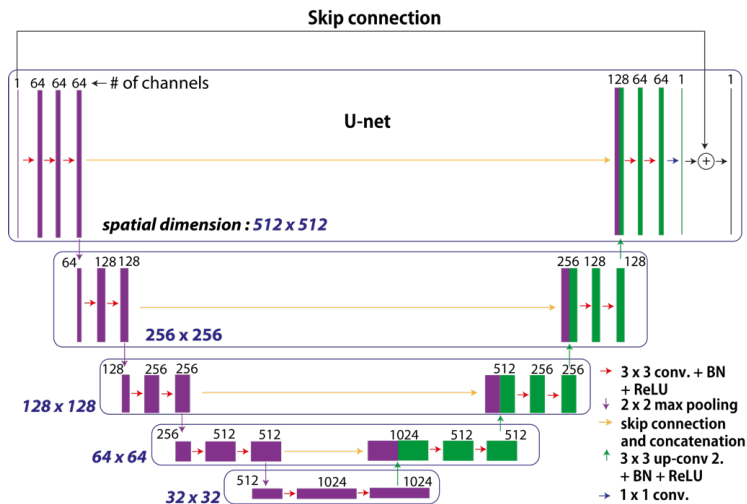
Recent appearance of Deep ConvNets

(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

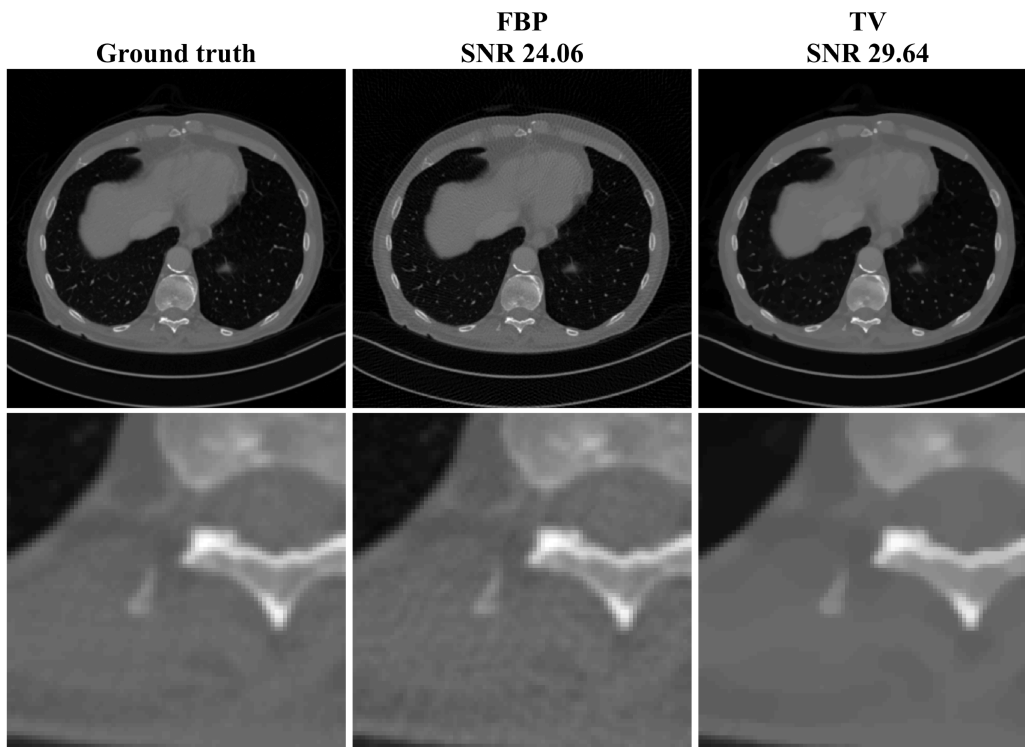
(Jin et al., IEEE TIP 2017)



21

CT data

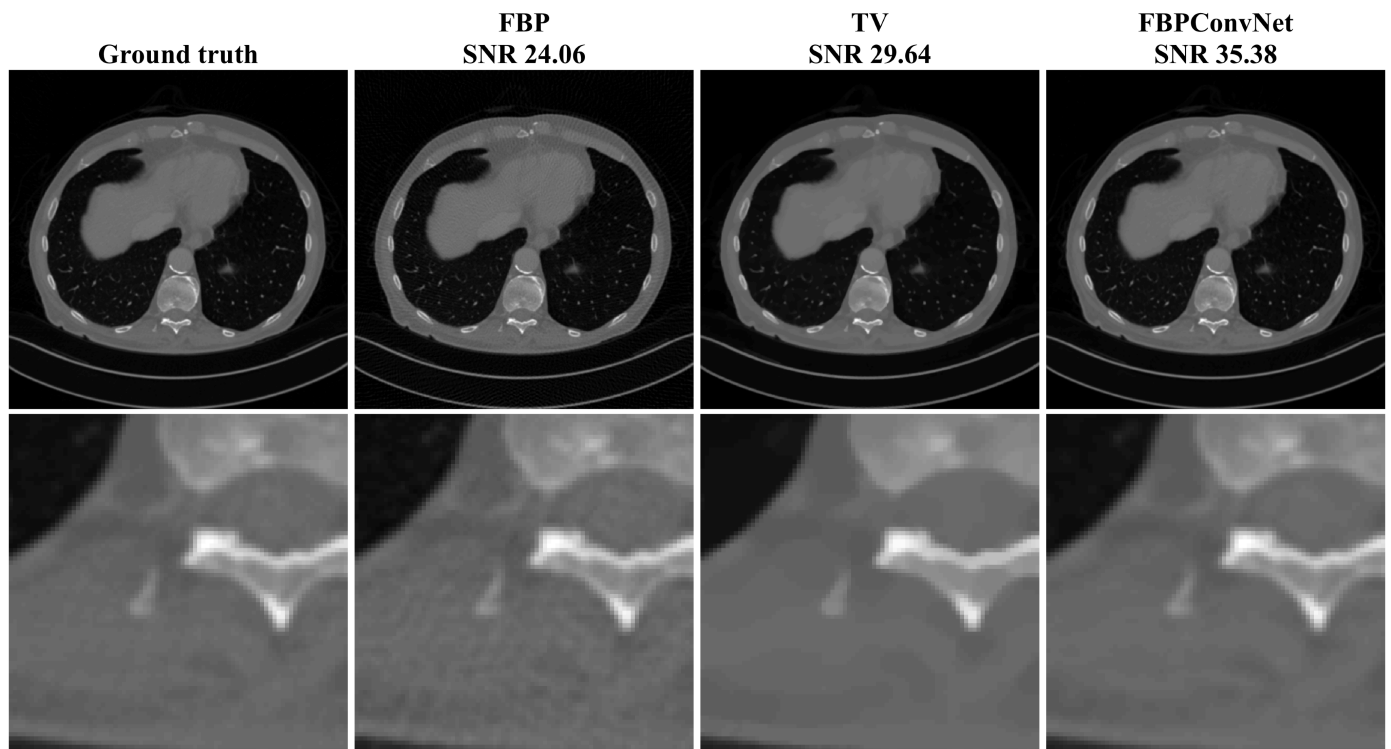
Dose reduction by 7: 143 views



Reconstructed from
from 1000 views

CT data

Dose reduction by 7: 143 views



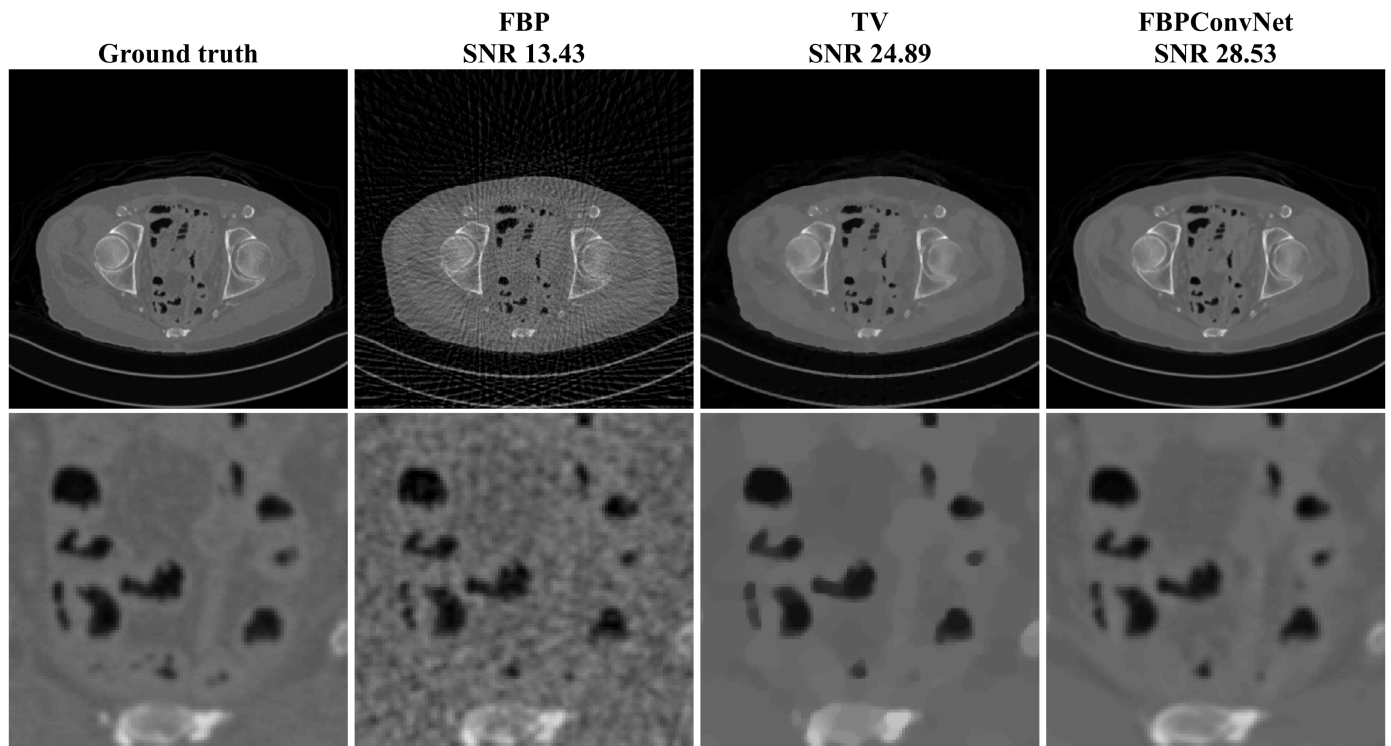
Reconstructed from
from 1000 views

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



CT data

Dose reduction by 20: 50 views



Reconstructed from
from 1000 views

(Jin-McCann-Froustey-Unser, *IEEE Trans. Im Proc.*, 2017)



OUTLINE

- Introduction ✓
- Continuous-domain theory of sparsity ✓
- From compressed sensing to deep neural networks ✓
- **Deep neural networks vs. deep splines**
 - Background
 - Continuous piecewise linear (CPWL) functions / splines
 - New representer theorem for deep neural networks

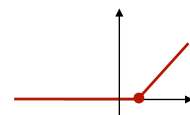


25

Deep neural networks and splines

- Preferred choice of activation function: ReLU

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



- 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

- ReLU \Rightarrow CPWL

(Montufar *NIPS* 2014)

- CPWL \Rightarrow Deep ReLU network

(Strang *SIAM News* 2018)

- Deep ReLU nets = hierarchical splines

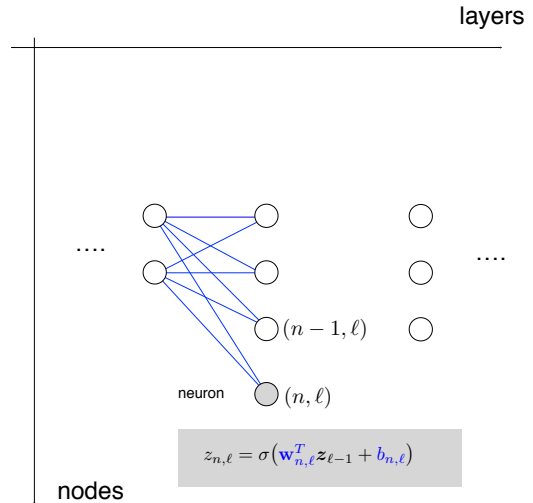
- ReLU is a piecewise-linear spline

(Poggio-Rosasco 2015)

26

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

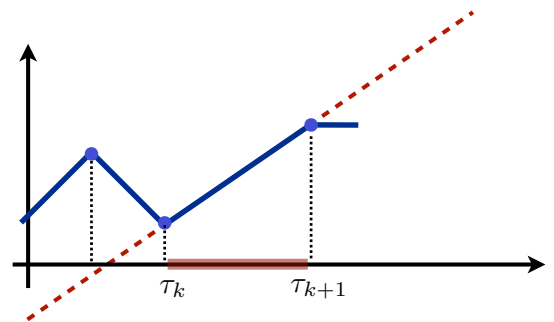


Learned

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

Continuous-PieceWise Linear (CPWL) functions

- 1D: Non-uniform spline de degree 1

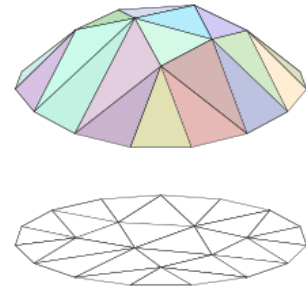


Partition: $\mathbb{R} = \bigcup_{k=0}^K P_k$ with $P_k = [\tau_k, \tau_{k+1})$, $\tau_0 = -\infty < \tau_1 < \dots < \tau_K < \tau_{K+1} = +\infty$.

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

- (i) : f_{spline} is continuous $\mathbb{R} \rightarrow \mathbb{R}$
- (ii) : for $x \in P_k$: $f_{\text{spline}}(x) = f_k(x) \triangleq a_k x + b_k$ with $(a_k, b_k) \in \mathbb{R}^2$, $k = 0, \dots, K$
- $f_{\text{spline}}(x) = \tilde{b}_0 + \tilde{b}_1 x + \sum_{k=1}^K \tilde{a}_k (x - \tau_k)_+$ with $\tilde{b}_0, \tilde{b}_1 \in \mathbb{R}$, $(\tilde{a}_k) \in \mathbb{R}^K$.

CPWL functions in high dimensions



■ Multidimensional generalization

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

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

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

- (i) : f_{CPWL} is continuous $\mathbb{R}^N \rightarrow \mathbb{R}$
- (ii) : 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 $\mathbf{a}_k \in \mathbb{R}^N, b_k \in \mathbb{R}, k = 1, \dots, K$

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.

29

Algebra of CPWL functions

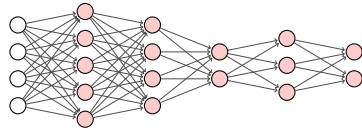
- any linear combination of (vector-valued) CPWL functions $\mathbb{R}^N \rightarrow \mathbb{R}^{N'}$ is CPWL, and,
- the composition $\mathbf{f}_2 \circ \mathbf{f}_1$ of any two CPWL functions with compatible domain and range—i.e., $\mathbf{f}_2 : \mathbb{R}^{N_1} \rightarrow \mathbb{R}^{N_2}$ and $\mathbf{f}_1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

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

30

Implication for deep ReLU neural networks

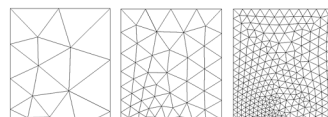


$$\mathbf{f}_{\text{deep}}(\mathbf{x}) = (\sigma_L \circ \mathbf{f}_L \circ \sigma_{L-1} \circ \dots \circ \sigma_2 \circ \mathbf{f}_2 \circ \sigma_1 \circ \mathbf{f}_1)(\mathbf{x})$$

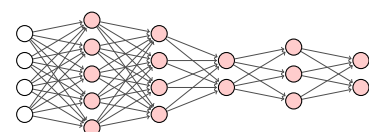
- Each scalar neuron activation, $\sigma_{n,\ell}(x) = \text{ReLU}(x)$, is CPWL.
- Each layer function $\sigma_\ell \circ \mathbf{f}_\ell(\mathbf{x}) = (\mathbf{W}_\ell \mathbf{x} + \mathbf{b}_\ell)_+$ is CPWL
- The whole feedforward network $\mathbf{f}_{\text{deep}} : \mathbb{R}^{N_0} \rightarrow \mathbb{R}^{N_L}$ 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**.

31

CPWL functions: further properties



- The CPWL model has universal approximation properties (as one increases the number of regions)
- Any CPWL function $\mathbb{R}^N \rightarrow \mathbb{R}$ can be implemented via a deep ReLU network with no more than $\log_2(N + 1) + 1$ layers

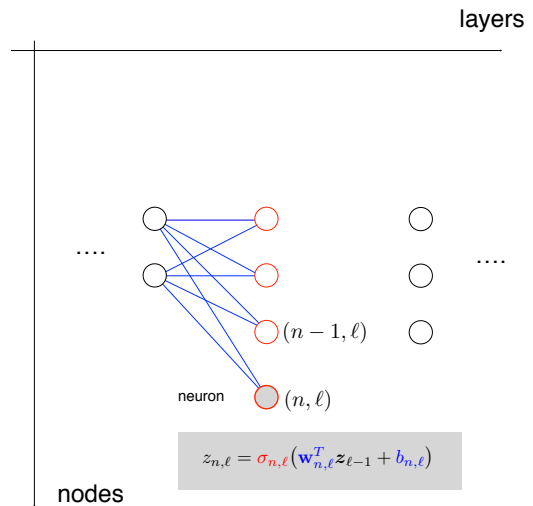


(Arora ICLR 2018)

32

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



$$\mathbf{f}_{\text{deep}}(\mathbf{x}) = (\sigma_L \circ f_L \circ \sigma_{L-1} \circ \dots \circ \sigma_2 \circ f_2 \circ \sigma_1 \circ f_1)(\mathbf{x})$$

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

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

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

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

equipped with the norm $\|f\|_{\text{BV}^{(2)}} \triangleq \|\mathbf{D}^2 f\|_{\mathcal{M}} + |f(0)| + |f(1) - f(0)|$

Representer theorem for deep neural networks

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

(U. arXiv:1802.09210, Feb 2018)

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

35

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

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

- Link with ℓ_1 minimization techniques

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

36

Optimality results

Lemma 1 (TV⁽²⁾-optimality of piecewise-linear interpolants)

Consider a series of scalar data points $(x_m, y_m), m = 1, \dots, M$ with $M > 2$ and $x_1 \neq x_2$. Then, the extremal points of the interpolation problem

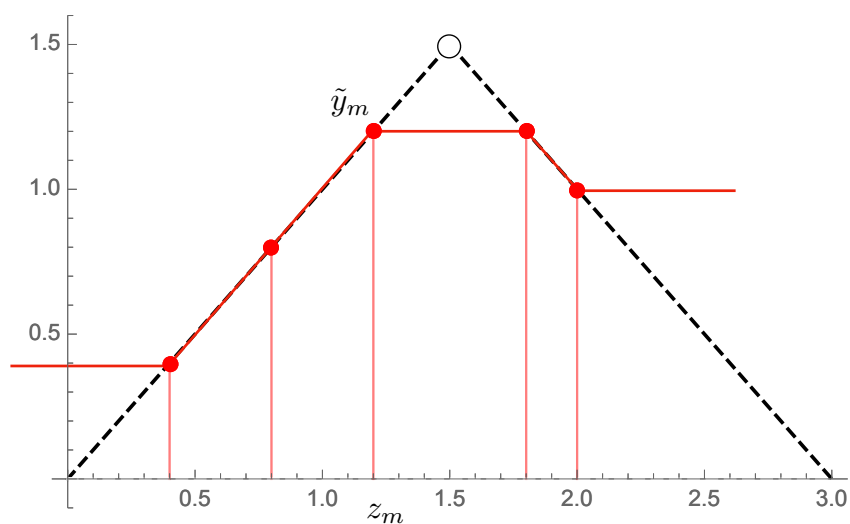
$$\arg \min_{f \in \text{BV}^{(2)}(\mathbb{R})} \|D^2 f\|_{\mathcal{M}} \quad \text{s.t.} \quad f(x_m) = y_m, \quad m = 1, \dots, M$$

are nonuniform splines of degree 1 with no more than $(M - 2)$ adaptive knots.

(U., JMLR 2019; Appendix C)

37

Comparison of linear interpolators



38

Spline interpolants: RKHS vs sparse

Lemma 1 (TV⁽²⁾-optimality of piecewise-linear interpolants)

Consider a series of scalar data points $(x_m, y_m), m = 1, \dots, M$ with $M > 2$ and $x_1 \neq x_2$. Then, the extremal points of the interpolation problem

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

are nonuniform splines of degree 1 with no more than $(M - 2)$ adaptive knots.

Proposition 2 (Sobolev optimality of piecewise-linear interpolation)

Let $H^1(\mathbb{R}) = \{f : \mathbb{R} \rightarrow \mathbb{R} : \|Df\|_{L_2}^2 + |f(0)|^2 < \infty\}$. Given a series of distinct data points $(x_m, y_m), m = 1, \dots, M$, the interpolation problem

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

has a unique piecewise-linear solution that can be written as

$$s_2(x) = b_1 + \sum_{m=1}^M a_m (x - x_m)_+.$$

39

Deep spline networks: Discussion

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

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

No need to normalize:

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

- 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 $(K_{n,\ell} = 5 \text{ or } 7, \mathbf{b}_{n,\ell} = \mathbf{0})$
(Agostinelli et al. 2015)

40

Deep spline networks (Cont'd)



■ Key features

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

■ Generalizations

- Broad family of cost functionals
- Cases where a subset of network components is fixed
- Generalized forms of regularization: $\psi(\text{TV}^{(2)}(\sigma_{n,\ell}))$

■ Challenges

⇒ In need for more powerful training algorithms

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

41

Acknowledgments

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

- Dr. Julien Fageot
- Prof. John Paul Ward
- Dr. Mike McCann
- Dr. Kyong Jin
- Harshit Gupta
- Dr. Ha Nguyen
- 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
-



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

42

References

■ Optimality of 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.

■ Image reconstruction with sparsity constraints (CS)

- M. Guerquin-Kern, M. Häberlin, K.P. Pruessmann, M. Unser, "A Fast Wavelet-Based Reconstruction Method for Magnetic Resonance Imaging," *IEEE Transactions on Medical Imaging*, vol. 30, no. 9, pp. 1649-1660, 2011.
- 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.
- 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, pp. 1-30, Jul. 2019.

43

Sketch of proof

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

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 :

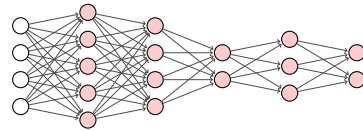
- Initialization (input): $\tilde{\mathbf{y}}_{m,0} = \mathbf{x}_m$.

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

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

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

$$\text{with } \tilde{y}_{n,m,\ell} = \tilde{\sigma}_{n,\ell}(z_{n,m,\ell}) \quad n = 1, \dots, N_\ell. \quad \Rightarrow \quad \tilde{\mathbf{f}}(\mathbf{x}_m) = \tilde{\mathbf{y}}_{m,L}$$



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

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

44