



# Topics in Numerical Analysis II

## Computational Inverse Problems

Bangti Jin (b.jin@cuhk.edu.hk)

Chinese University of Hong Kong

November 20, 2023



# Outline

- 1 Deep learning for inverse problems: supervised approach
  - Supervised approaches to image reconstruction
  - Model adaptation
  - Deep equilibrium model



# Review

linear (nonlinear) inverse problems

$$Ax = y$$

- Hilbert space: truncated SVD, Tikhonov regularization, iterative methods (Landweber, Kaczmarz method, stochastic gradient descent)
- Banach space: sparsity regularization (theory + algorithm)
- other possibilities: other penalties, iterative methods in Banach space

What is beyond ?



# What is deep learning

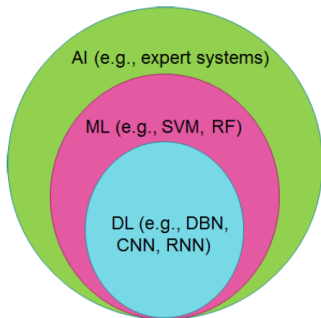
- gained a lot of attention recently and impressive results, e.g., image classification, speech recognition, segmentation, natural language processing, alphaGo, chatGPT, alphaFold, etc.
- computationally highly efficient
- limitation of current approaches to inverse problems: expressibility of **hand-crafted** prior (regularizer)
- instead: learn prior from data itself to improve the image reconstruction





## what is deep learning ?

- artificial intelligence: any techniques that enables computers to mimic human intelligence
- machine learning: a subset of AI using statistical techniques to enable machine to learn from data
- deep learning: a subset of ML where complicated tasks are performed by breaking them down into several layers

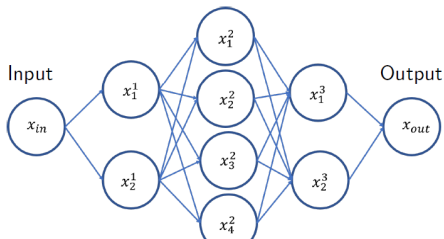




# convolutional neural networks (CNN)

- aim: to establish a mapping  $G_\theta(x_i) = x_o$  with input / output  $x_i, x_o \in X = \mathbb{R}^{n \times n}$  both being images
- construction: A fully convolutional neural network (CNN) is the composition of an affine linear operation, based on **convolution**, and a **nonlinear** function. The essential building block is

$$x_{i+1} = \sigma(b + \omega * x_i)$$





continuous convolution

$$(f * g)(t) = \int_{-\infty}^{\infty} f(t-s)g(s)ds$$

discrete convolution

$$(f * g)_j = \sum_{i=-\infty}^{\infty} f_{j-i}g_i$$

circular discrete convolution (nonzero for  $[0, N-1]$ )

$$(f * g)_j = \sum_{i=0}^{N-1} f_i g_{j-i \bmod N}$$

fast convolution using fast Fourier transform (FFT), of complexity  $\mathcal{O}(N \log N)$

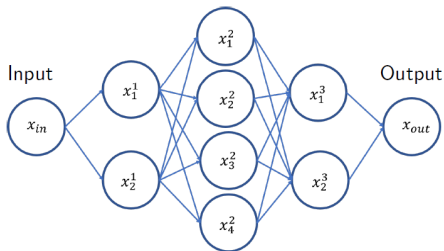


## CNN (1st layer)

Each layer consists of a convolutional layer (scalar + convolution) followed by a nonlinear function

$$x_i^1 = \sigma(b_i^1 + \omega_{i,1}^1 * x_{in}), \quad i = 1, 2$$

- $b$ : biases (scalar);  $\omega$ : a convolutional filter
- $\sigma$ : pointwise nonlinear function, typically  $\text{ReLU}(x) = \max(x, 0)$



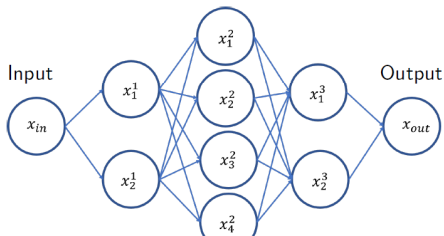


## CNN (2nd layer)

Each layer consists of a convolutional layer (scalar + convolution) followed by a nonlinear function

$$x_i^2 = \sigma \left( b_i^2 + \sum_{j=1}^2 \omega_{i,j}^2 * x_j^1 \right), \quad i = 1, \dots, 4$$

- $b$ : biases (scalar);  $\omega$ : a convolutional filter
- $\sigma$ : pointwise nonlinear function, typically  $\text{ReLU}(x) = \max(x, 0)$



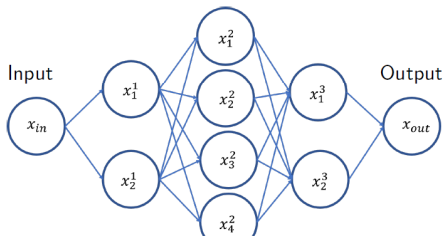


## CNN (3rd layer)

Each layer consists of a convolutional layer (scalar + convolution) followed by a nonlinear function

$$x_i^3 = \sigma \left( b_i^3 + \sum_{j=1}^4 \omega_{i,j}^3 * x_j^2 \right), \quad i = 1, 2$$

- $b$ : biases (scalar);  $\omega$ : a convolutional filter
- $\sigma$ : pointwise nonlinear function, typically  $\text{ReLU}(x) = \max(x, 0)$



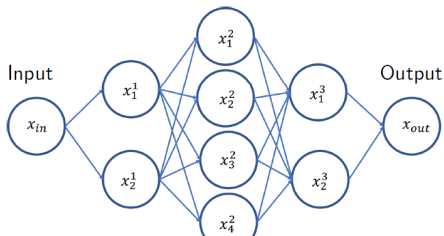


## CNN (last layer)

Each layer consists of a convolutional layer (scalar + convolution) followed by a nonlinear function

$$x_o = \sigma \left( b_1^4 + \sum_{j=1}^2 \omega_{1,j}^4 * x_j^3 \right)$$

- $b$ : biases (scalar);  $\omega$ : a convolutional filter
- $\sigma$ : pointwise nonlinear function, typically  $\text{ReLU}(x) = \max(x, 0)$





# training of CNN

- denote the CNN with given parameters  $\theta$  by  $G_\theta$
- With an input set  $\{x^i\}$ , and the desired outputs  $\{x_\dagger^i\}$ , seek to minimize a loss function, i.e.,

$$\text{loss}(\theta; x^i) = \|G_\theta(x^i) - x_\dagger^i\|_2^2$$

The set of (input, output) pairs  $\{(x^i, x_\dagger^i)\}_{i=1}^N$  is the training data

- training = optimization of  $\theta$ , by minimizing the loss  $\text{loss}(\theta; \mathbf{x})$  with respect to  $\theta$  over the training set

caveat:

- require a lot of training data
- training can be very expensive





the setting of general **discrete** inverse problems

$$y = A(x^\dagger) + \delta y$$

- $x^\dagger \in X$ : true unknown
- $y \in Y$ : measured data
- $A : X \rightarrow Y$ : forward operator (maybe imperfectly known only)
- $\delta y \in Y$ : data noise



## parameterised reconstruction of the unknown

- given the measurement  $y$  and operator  $A$ , aim to find a mapping  $A^\dagger : Y \rightarrow X$  s.t.

$$A^\dagger(y) \approx x^\dagger$$

- parameterise the mapping  $A^\dagger : Y \rightarrow X$  with some parameter  $\theta$  and find an optimal set of parameters  $\theta^*$  s.t.

$$A_{\theta^*}^\dagger(y) \approx x^\dagger$$

- Issue: What is the input ?
- Issue: How to construct the mapping ?



# gradient unrolling

given the classical minimization problem

$$x^* = \arg \min_{x \geq 0} \frac{1}{2} \|Ax - y\|_2^2 + \frac{\alpha}{2} \|x\|_2^2$$

This can be iteratively solved by a projected gradient descent

$$x_{k+1} = P_{x \geq 0}(x_k - \lambda((A^T A + \alpha I)x_k - A^T y))$$

which can be rewritten as

$$x_{k+1} = P_{x \geq 0}(((1 - \lambda\alpha) - \lambda A^T A)x_k + \lambda A^T y)$$



given the update equations

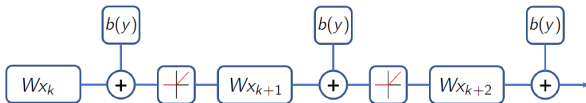
$$x_{k+1} = P_{x \geq 0} \left( \underbrace{((1 - \lambda\alpha) - \lambda A^\top A)}_W x_k + \underbrace{\lambda A^\top y}_{b(y)} \right)$$

Interpretation:

$Wx_k = ((1 - \lambda\alpha) - \lambda A^\top A)x_k$  : performs a filtering

and

$b(y) := \lambda A^\top y$  : data-dependent constant





## motivation: ISTA

given the minimization problem with a sparsity transform:

$$x^* = \arg \min_{x \geq 0} \frac{1}{2} \|Ax - y\|_2^2 + \alpha \|Wx\|_1$$

$W$ : orthonormal frame, we can solve this with the iterative thresholding algorithm (ISTA), step size  $\lambda$

$$x_{k+1} = S_{\alpha, W}(x_k - \lambda A^\top (Ax_k - y))$$

with the soft-thresholding operator  $S_\mu$  with  $\mu = \lambda\alpha$ ,

$$S_{\alpha, W}(x) = W^{-1} S_\mu(Wx)$$

This can be written in one step as

$$x_{k+1} = W^{-1}(S_\mu W(x_k - \lambda A^\top (Ax_k - y)))$$



Consider the one-step

$$x_{k+1} = W^{-1}(S_{\mu} W(x_k - \lambda A^{\top}(Ax_k - y)))$$

we now write  $a_k = Wx_k$  and apply  $W$  on both sides, we obtain

$$\begin{aligned} Wx_{k+1} &= WW^{-1}(S_{\mu} W(x_k - \lambda A^{\top}(Ax_k - y))) \\ \Rightarrow a_{k+1} &= S_{\mu} W(x_k - \lambda A^{\top}(Ax_k - y)) \\ \Rightarrow a_{k+1} &= S_{\mu}(Wx_k - \lambda WA^{\top}(Ax_k - y)) \\ \Rightarrow a_{k+1} &= S_{\mu}(a_k - \lambda WA^{\top}(AW^{-1}a_k - y)) \\ \Rightarrow a_{k+1} &= S_{\mu}((a_k - \lambda WA^{\top}AW^{-1})a_k - \lambda WA^{\top}y) \end{aligned}$$



given the update equation

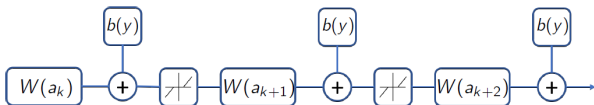
$$a_{k+1} = S_{\mu}((a_k - \lambda WA^T AW^{-1})a_k - \lambda WA^T y)$$

note that similarly

$$W(a_k) = (I - \lambda W^T A^T AW^{-1})a_k \text{ perform filtering}$$

and

$$b(y) = \lambda WA^T y \text{ data dependent constant}$$





## Interpretation: pure CNN

The general operation can be summarized as:

- iteratively filtering by  $W$
- adding a bias  $b(y)$
- applying a pointwise nonlinearity  $S$  (ReLU or soft-thresholding)

the procedure is identical with a layer in a CNN

$$x_i^{k+1} = \sigma(b_i^k + \sum_{j=1}^{m_k} \omega_{i,j}^k * x_j^k), \quad i = 1, \dots, n$$

letting  $Wx^k = \sum_{j=1}^{m_k} \omega_{i,j}^k * x_j^k$ , then we have simplified





# How to use CNN in inverse problems

approach (i): post-processing K. H. Jin, M. T. McCann, E. Froustey IEEE Trans. Imag. Proc.

2017

- Motivated by the unrolling, use a CNN as sort of post-processing of an initial reconstruction
- given a reconstruction operator  $A^\dagger : Y \rightarrow X$ , such as filtered backprojection, we compute

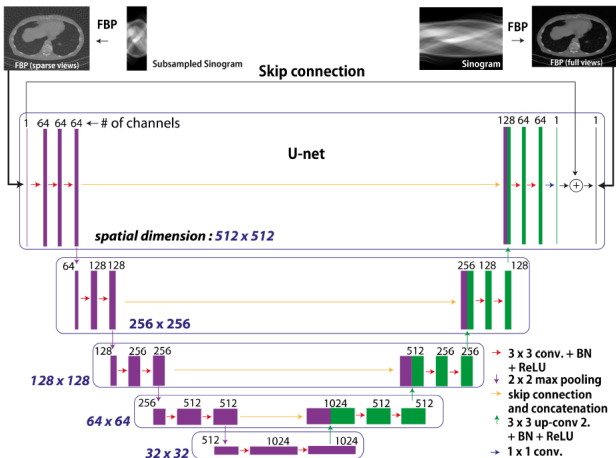
$$x^0 = A^\dagger y$$

and train a CNN  $f_\theta$  to clean the input images (i.e., remove noise and undersampling artefacts)

- The parameters  $\theta$  are given by convolution filters  $W$  and biases  $b$



# standard architecture – U-net





## popularity of U-net

### **U-net: Convolutional networks for biomedical image segmentation**

[O Ronneberger](#), [P Fischer](#), [T Brox](#) - International Conference on Medical ..., 2015 - Sprin

... We demonstrate the application of the **u-net** to three different segmentation tasks. Th

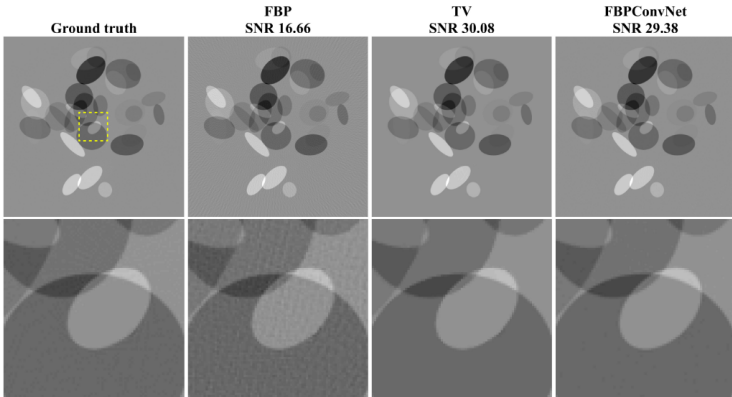
... The **u-net** (averaged over 7 rotated versions of the input data) achieves without any fi

☆ Save  Cite Cited by 52086 Related articles All 30 versions

(as of Nov. 2022)

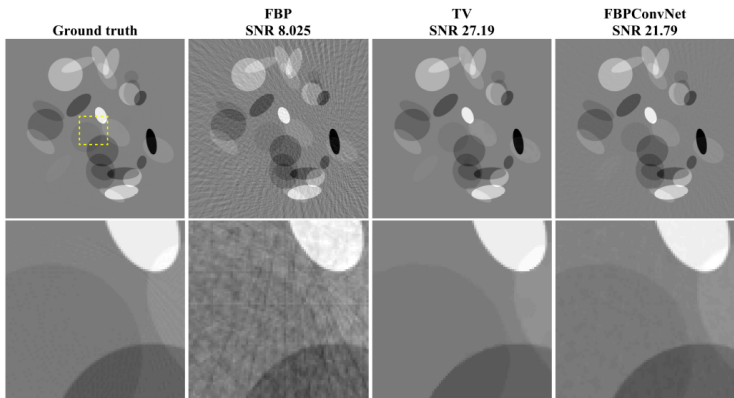


## ellipse data with 143 views



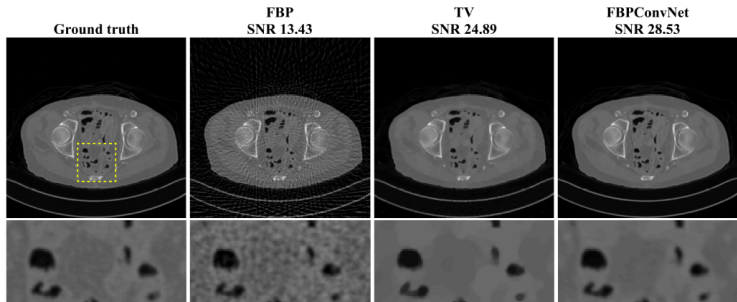


## ellipse data with 50 views





## reconstruction on biomedical data set with 50 views





# How to use CNN in inverse problems

approach (ii): learned model based reconstruction

- aim: model-based learning / reconstruction: in this approach, the forward and adjoint operators of the imaging problem are used directly in the recovery procedure.
- the general problem

$$x^* = \arg \min_{x \geq 0} \|Ax - y\|_2^2 + \alpha \mathcal{R}(x)$$

where  $\mathcal{R}$  denotes some regularisation term.

- This leads to the gradient descent scheme

$$x^{k+1} = P_{x \geq 0}(x^k - \lambda(A^\top(Ax_k - y) + \alpha \nabla \mathcal{R}(x^k)))$$



including the model in the learned algorithm: ResNet

- without the data fit, we only have

$$x^{k+1} = P_{x \geq 0}(x^k - \lambda \alpha \nabla \mathcal{R}(x^k))$$

- This is similar to **ResNet** structure, with the convolutional part learning the action of  $\lambda \alpha \nabla \mathcal{R}(x^k)$







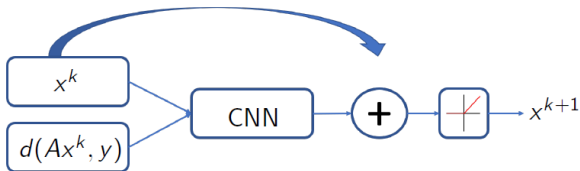
- add the model information (i.e., the data fit)

$$d(y, Ax) = \|Ax - y\|_2^2$$

- the update equation reads

$$x^{k+1} = P_{x \geq 0}(x^k - \lambda(\nabla d(y, Ax^k) + \alpha \nabla \mathcal{R}(x^k)))$$

- including the gradient of the data fit modifies the ResNet to a simple model based learned gradient descent network (each block  $G_{\theta_k}^k$ )





The algorithm can be summarized as J. Adler, O. Öktem, Inverse Problems 2017; IEEE

Trans. Med. Imag. 2018

---

### Algorithm 1 Learned gradient descent

---

```
1: function  $G_\theta(y)$ 
2:    $x^0 \leftarrow A^T y$ 
3:    $k \leftarrow 0$ 
4:   for  $k < k_{max}$  do
5:     Compute  $\nabla d(y, Ax^k) = A^*(Ax^k - y)$ 
6:      $x^{k+1} \leftarrow G_{\theta_k}^k(\nabla d(y, Ax^k), x^k)$ 
7:      $k \leftarrow k + 1$ 
8:   end for
9: end function(return  $x^{k_{max}}$ )
```

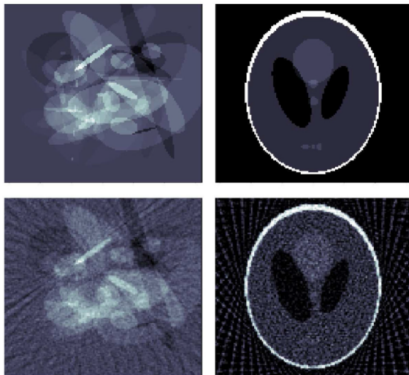
---



create a training set: train the network on phantoms of random ellipses by minimizing

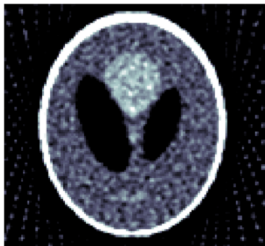
$$\text{loss}(\theta; x^K) = \|G_\theta(y) - x^\dagger\|_2^2$$

and test the performance on the Shepp Logan phantom





results by filtered back projection (FBP)





## results by total variation





## results by post-processing





results by learned gradient descent





results by learned primal dual







## quantitative comparative results

Method	PSNR (dB)	SSIM	Runtime (ms)	Parameters
FBP	19.75	0.597	4	1
TV	28.06	0.928	5 166	1
Learned U-Net	29.20	0.943	9	$10^7$
Learned Gradient Descent	32.02	-	-	-
Learned Primal-Dual	38.28	0.988	49	$2.4 \cdot 10^5$



# The peril of training data

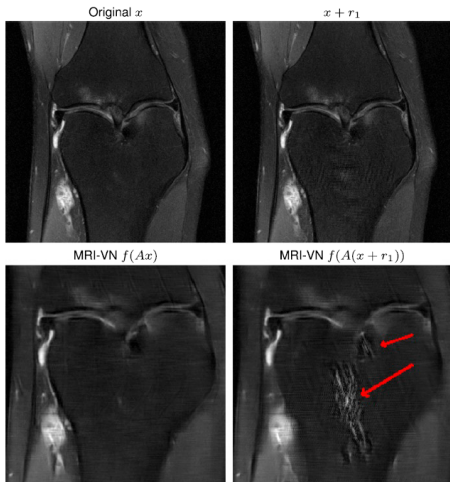
sensitivity of classification neural networks



K. Eykholt et al CVPR 2018



# The peril of training data





# The need for model adaptation

setting of the problem D. Gilton, G. Ongie, R. Willett. IEEE Trans. Comput. Imag. 2021

- Suppose we have neural network estimator  $\hat{x} = f_0(y)$  trained to solve

$$y = A_0x + \epsilon, \quad x \sim P_X, \epsilon \sim P_{N_0}$$

- concerned task:

$$y = A_1x + \epsilon', \quad x \sim P_X, \epsilon' \sim P_{N_1}$$

with  $A_1$  known, or partially unknown (e.g., additional param. in  $A_1$ )

- goal: can we adapt / modify the estimator  $\hat{x} = f_0(y)$  to the new task ?



## transfer learning approach

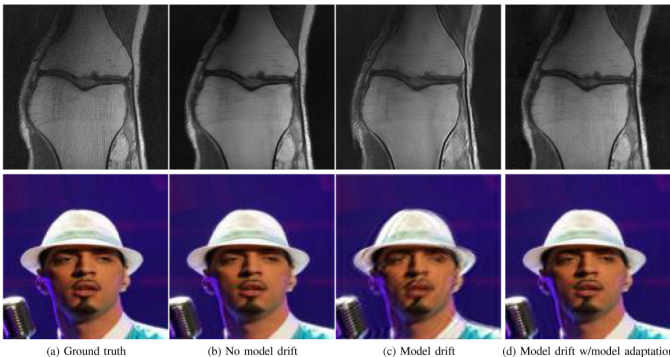
- naive approach: apply the neural network directly, with  $f(\cdot; \theta_0^*, A_0)$  by  $f(\cdot; \theta_0^*, A_1)$
- adaptation: estimate the image by  $f(y; \theta_1^*, A_1)$ , with  $\theta_1^*$  solving

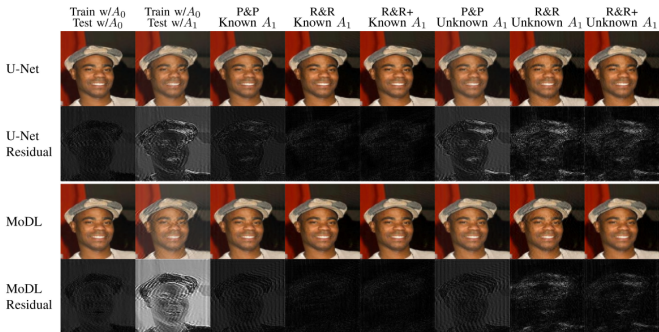
$$\min_{\theta} \underbrace{\|y - A_1 f(y; \theta, A_1)\|^2}_{\text{data consistency}} + \mu \underbrace{\|\theta - \theta_0^*\|^2}_{\text{constraint param.}}$$

with  $\mu > 0$  is a tunable parameter

- with training dataset, learning the model  $A$  as well ...

$$(\theta_1^*, A_1^*) = \arg \min_{\theta, A \in \mathcal{A}} \frac{1}{N} \sum_{i=1}^N \|y_i - A f(y_i; \theta, A)\|^2 + \mu \|\theta - \theta_0^*\|^2$$







# Deep equilibrium model

state of the art

- modern feedforward deep networks build on the concept of layers: each network is a stack of  $L$  transformations,  $L$  being the depth (determined by model designer)
- update: backpropagation through  $L$  layers via chain rule
- requires storing all intermed. values (activations) of these layers  
⇒ **memory requirement is proportional to  $L$**





What happens if you iterate more ?



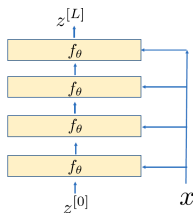
at iter 1, 10, 20, 30, 40

D. Gilton, G. Ongie, R. Willett. IEEE Trans. Comput. Imag. 2021



weight-tied deep sequence models:

$$z_{1:T}^{[i+1]} = f_{\theta}(z_{1:T}^{[i]}; x_{1:T}), \quad i = 0, \dots, L-1,$$
$$z_{1:T}^{[0]} = 0, G(x_{1:T}) = z_{1:T}^{[L]}$$



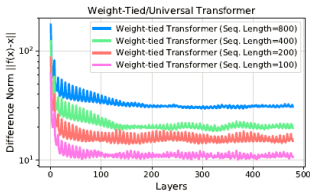
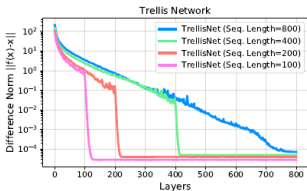
features

- reduce the model size, a form of regularization
- any deep network can be represented by a weight-tied deep network of equal depth and linear increase in width
- the network can be unrolled to any depth (improved feature extraction as depth increases)
- in practice only fixed layers, due to memory on training hardware



hypothesis: the weight-tied models converge to a fixed point, as the depth increases to infinity

$$z_{1:T}^* = \lim_{i \rightarrow \infty} z_{1:T}^{[i]} = \lim_{i \rightarrow \infty} f_{\theta}(z_{1:T}^{[i]}; x_{1:T}) = f_{\theta}(z_{1:T}^*; x_{1:T})$$



S. Bai, J. Z. Kolter, and V. Koltun, NeurIPS 2019



eventual hidden layer values of an infinite depth network

$$z_{1:T}^* = f_{\theta}(z_{1:T}^*; x_{1:T})$$

**deep equilibrium model (DEM)**: compute the fixed point  $z_{1:T}^*$  directly using black-box root-finding, via implicit differentiation

- fixed point procedure

$$z_{1:T}^{[i+1]} = f_{\theta}(z_{1:T}^{[i]}; x_{1:T}), \quad i = 0, 1, \dots$$

- alternative method (quasi-Newton method)

S. Bai, J. Z. Kolter, and V. Koltun, NeurIPS 2019



equilibrium equation

$$f_{\theta}(z_{1:T}^*; x_{1:T}) = z_{1:T}^*$$

let  $(\cdot)$  be a parameter of  $f_{\theta}(z_{1:T}^*; x_{1:T})$  (e.g.,  $\theta$  or  $x_{1:T}$ ) implicit diff.  $\Rightarrow$

$$\frac{\partial z_{1:T}^*}{\partial(\cdot)} = \frac{\partial f_{\theta}(z_{1:T}^*; x_{1:T})}{\partial(\cdot)}$$

$$\frac{\partial z_{1:T}^*}{\partial(\cdot)} = \frac{df_{\theta}(z_{1:T}^*; x_{1:T})}{d(\cdot)} + \frac{\partial f_{\theta}(z_{1:T}^*; x_{1:T})}{\partial z_{1:T}^*} \frac{\partial z_{1:T}^*}{\partial(\cdot)}$$



$$\left( I - \frac{\partial f_{\theta}(z_{1:T}^*; x_{1:T})}{\partial z_{1:T}^*} \right) \frac{\partial z_{1:T}^*}{\partial (\cdot)} = \frac{df_{\theta}(z_{1:T}^*; x_{1:T})}{d(\cdot)}$$

fixed point equation:  $g_{\theta}(z_{1:T}^*) = f_{\theta}(z_{1:T}^*; x_{1:T}) - z_{1:T}^*$

$$J_{g_{\theta}}|_{z_{1:T}^*} = - \left( I - \frac{\partial f_{\theta}(z_{1:T}^*; x_{1:T})}{\partial z_{1:T}^*} \right)$$

$$\frac{\partial z_{1:T}^*}{\partial (\cdot)} = -(J_{g_{\theta}}^{-1}|_{z_{1:T}^*}) \frac{df_{\theta}(z_{1:T}^*; x_{1:T})}{d(\cdot)}$$

do not backprop through all layers



## backprop through equilibrium point

$$\frac{\partial \ell}{\partial (\cdot)} = \frac{\partial \ell}{\partial \mathbf{z}_{1:T}^*} \frac{\partial \mathbf{z}_{1:T}^*}{\partial (\cdot)} = -\frac{\partial \ell}{\partial \mathbf{z}_{1:T}^*} (\mathbf{J}_{g_\theta}^{-1} |_{\mathbf{z}_{1:T}^*}) \frac{d\mathbf{f}_\theta(\mathbf{z}_{1:T}^*; \mathbf{x}_{1:T})}{d(\cdot)}$$

- backward gradient through infinite stacking can be represented as one-step matrix multiplication (Jacobian)
- SGD

$$\theta^+ = \theta - \alpha \frac{\partial \ell}{\partial \theta} = \theta + \frac{\partial \ell}{\partial \mathbf{z}_{1:T}^*} (\mathbf{J}_{g_\theta}^{-1} |_{\mathbf{z}_{1:T}^*}) \frac{d\mathbf{f}_\theta(\mathbf{z}_{1:T}^*; \mathbf{x}_{1:T})}{d(\cdot)}$$

it does not depend on rootfinding procedure, structure of transformation, no intermediate value



## reliable estimation of equilibrium point

- use any black-box root-finding method, e.g., quasi-Newton (Broyden) method for

$$g_{\theta}(z_{1:T}^*; x_{1:T}) = f_{\theta}(z_{1:T}^*; x_{1:T}^*) - z_{1:T}^* \rightarrow 0.$$

- Broyden iterations

$$z_{1:T}^{i+1} = z_{1:T}^i - \alpha B g_{\theta}(z_{1:T}^i; x_{1:T}), \quad i = 0, 1, \dots$$

with step size  $\alpha$ ,  $B \approx J_{g_{\theta}}^{-1}|_{z_{1:T}^i}$  (Jacobian approx.)





overall computational strategy:

- forward pass:  $z_{1:T}^* = \text{RootFind}(g_\theta; x_{1:T})$
- backward pass (implicit diff.):

$$\frac{\partial \ell}{\partial (\cdot)} = -\frac{\partial \ell}{\partial z_{1:T}^*} (J_{g_\theta}^{-1} |_{z_{1:T}^*}) \frac{df_\theta(z_{1:T}^*; x_{1:T})}{d(\cdot)}$$

- accelerating DEM: to compute  $-\frac{\partial \ell}{\partial z_{1:T}^*} (J_{g_\theta}^{-1} |_{z_{1:T}^*})$   
solve a linear system with quasi-Newton method

$$(J_{g_\theta}^\top |_{z_{1:T}^*}) y^\top + \left( \frac{\partial \ell}{\partial z_{1:T}^*} \right)^\top = 0.$$

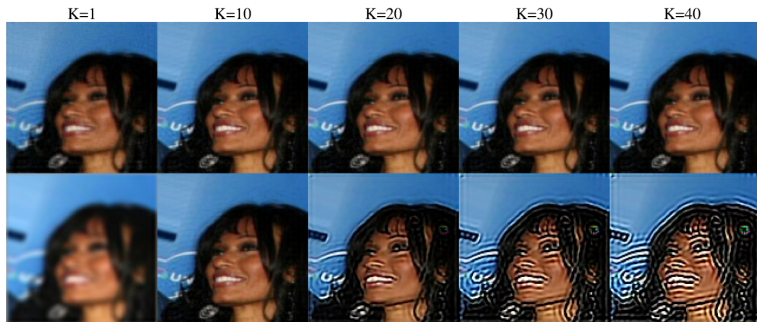


## properties

- memory cost:  $Z_{1:T}$ ,  $X_{1:T}$ ,  $f_\theta$  (and matrix-vector product)
- the analysis is independent of  $\theta$ , but should be stable and constrained
- universality of "single-layer" DEQs



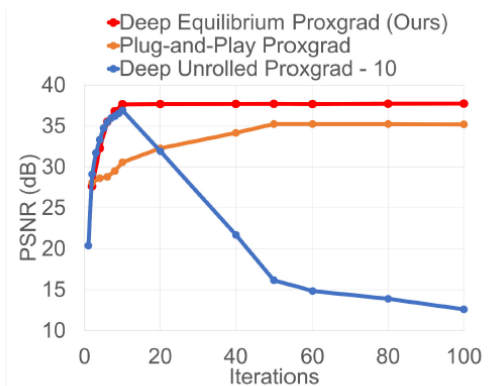
## performance on image recovery



D. Gilton, G. Ongie, R. Willett. IEEE Trans. Comput. Imag. 2021



## performance on image recovery



(b) Deblurring (2)

D. Gilton, G. Ongie, R. Willett. IEEE Trans. Comput. Imag. 2021



## preliminary theory for equilibrium methods Obmann and Haltmeier2023

$$A^*(Ax - y) + \alpha G(x) = 0$$

(i)  $x \mapsto \langle G(x), x - z \rangle$  is coercive (ii)  $G$  is weak to weak continuous

- (stability) Let  $\alpha > 0$  and  $(y_k) \subset Y$  with  $y_k \rightarrow y$ . Then, any sequence  $(x_k)$  satisfying  $T_\alpha(x_k, y_k) = 0$  has a weakly convergent subseq., and the limit is a solution of  $T_\alpha(x, y) = 0$ .
- (consistency) Let  $y \in \mathcal{R}(A)$ ,  $(y_k) \subset Y$  satisfy  $\|y_k - y\| = \delta_k \rightarrow 0$ , and let  $\alpha_k$  satisfy  $\alpha_k(\delta_k) \rightarrow 0$  and  $\delta_k^2/\alpha_k(\delta_k) \rightarrow 0$ . Any sequence  $(x_k)$  with  $T_{\alpha_k}(x_k, y_k) = 0$  has at least one weak cluster point, which is a solution to the problem with exact data.
- convergence rate ...



There are many exciting opportunities in developing algorithms and theories

- the choice of architecture is crucial
- the choice of training data is crucial
- the choice of training algorithm is crucial
- the specific physics / mathematical properties are important
- the theory is scarce
- ...