





Computational imaging & learning: optimisation & data-driven learning II

Luca Calatroni

MaLGa center, DIBRIS, UniGe
MMS, Istituto Italiano di Tecnologia
Genoa, Italy

Mini-Corso Data Science @ UniPD
Università degli studi di Padova
May 11-14 2026

1. Deep unfolding/unrolling
 - From ISTA to LISTA
 - Unrolled primal-dual (no extrapolation)
 - To unroll or not to unroll?
2. Deep Equilibrium Models
 - DEQ-GD & DEQ-PGD
 - DEQ training
3. Outlook/conclusions

-  V. Monga, Y. Li, Y. Eldar, *Algorithm Unrolling: Interpretable, Efficient Deep Learning for Signal and Image Processing*, IEEE Signal Processing Magazine, 38 (2), 2021.
-  U. Kamilov, C. Bouman, G. Buzzard, B. Wohlberg, *Plug-and-Play Methods for Integrating Physical and Learned Models in Computational Imaging: Theory, algorithms, and applications*, IEEE Signal Processing Magazine 40 (1), 2023.
-  C. Schönlieb, Z. Shumaylov, *Data-driven approaches to inverse problems*, arXiv preprint: <https://arxiv.org/abs/2506.11732>, 2025.
-  D. Gilton, G. Ongie, R. Willett, *Deep equilibrium architectures for inverse problems in imaging*, IEEE Transactions on Computational Imaging, 7, 2021.

Deep unfolding/unrolling

Deep unfolding/unrolling

From ISTA to LISTA

$$\min_{\mathbf{u} \in \mathbb{R}^N} \frac{1}{2} \|\mathbf{A}\mathbf{u} - \mathbf{y}\|^2 + \lambda \|\mathbf{u}\|_1$$

¹Daubechies, Defrise, De Mol, 2004

$$\min_{\mathbf{u} \in \mathbb{R}^N} \frac{1}{2} \|\mathbf{A}\mathbf{u} - \mathbf{y}\|^2 + \lambda \|\mathbf{u}\|_1$$

Iterative Soft Thresholding Algorithm (ISTA)¹

PGD iterations: for $\mathbf{u}^0 \in \mathbb{R}^N$, $\tau \in (0, 2/L)$ and $k \geq 0$

$$\mathbf{u}^{k+1} = \text{prox}_{\tau\lambda\|\cdot\|_1}(\mathbf{u}^k - \tau\mathbf{A}^\top(\mathbf{A}\mathbf{u}^k - \mathbf{y})) = \mathcal{T}_{\tau\lambda}(\mathbf{u}^k - \tau\mathbf{A}^\top(\mathbf{A}\mathbf{u}^k - \mathbf{y}))$$

where $\mathcal{T}_{\tau\lambda}(z)$ is the n -dimensional *soft-thresholding* operator.

¹Daubechies, Defrise, De Mol, 2004

$$\min_{\mathbf{u} \in \mathbb{R}^N} \frac{1}{2} \|\mathbf{A}\mathbf{u} - \mathbf{y}\|^2 + \lambda \|\mathbf{u}\|_1$$

Iterative Soft Thresholding Algorithm (ISTA)¹

PGD iterations: for $\mathbf{u}^0 \in \mathbb{R}^N$, $\tau \in (0, 2/L)$ and $k = 0, \dots, K - 1$

$$\mathbf{u}^{k+1} = \text{prox}_{\tau\lambda\|\cdot\|_1}(\mathbf{u}^k - \tau\mathbf{A}^\top(\mathbf{A}\mathbf{u}^k - \mathbf{y})) = \mathcal{T}_{\tau\lambda}(\mathbf{u}^k - \tau\mathbf{A}^\top(\mathbf{A}\mathbf{u}^k - \mathbf{y}))$$

where $\mathcal{T}_{\tau\lambda}(z)$ is the n -dimensional *soft-thresholding* operator.

- In practice: iterate **till convergence** ($k \rightarrow +\infty$)
- What if we stop **at a certain K** ? (\approx max. computational budget)

¹Daubechies, Defrise, De Mol, 2004

A neural network interpretation of ISTA

Look at one iteration step:

$$\mathbf{u}^{k+1} = \mathcal{T}_{\tau\lambda}(\mathbf{u}^k - \tau\mathbf{A}^\top(\mathbf{A}\mathbf{u}^k - \mathbf{y})) = \underbrace{\mathcal{T}_{\tau\lambda}}_{\approx \text{ReLU}} \left(\underbrace{(\mathbf{I} - \tau\mathbf{A}^\top\mathbf{A})\mathbf{u}^k + \tau\mathbf{A}^\top\mathbf{y}}_{\text{linear}} \right)$$

... analogy with k -th layer of a NN, $k = 0, \dots, K - 1$.

A neural network interpretation of ISTA

Look at one iteration step:

$$\mathbf{u}^{k+1} = \mathcal{T}_{\tau\lambda}(\mathbf{u}^k - \tau\mathbf{A}^\top(\mathbf{A}\mathbf{u}^k - \mathbf{y})) = \underbrace{\mathcal{T}_{\tau\lambda}}_{\approx \text{ReLU}} \left(\underbrace{(\mathbf{I} - \tau\mathbf{A}^\top\mathbf{A})\mathbf{u}^k + \tau\mathbf{A}^\top\mathbf{y}}_{\text{linear}} \right)$$

... analogy with k -th layer of a NN, $k = 0, \dots, K - 1$.

Interpretation: K -layer NN obtained by **unrolling/unfolding** the iterations:

$$\begin{aligned}\mathbf{u}^{k+1} &= \mathcal{T}_{\tau\lambda}(\mathbf{u}^k - \tau\mathbf{A}^\top(\mathbf{A}\mathbf{u}^k - \mathbf{y})) \\ &= \mathcal{T}_{\tau\lambda} \left(\mathcal{T}_{\tau\lambda} \left(\mathbf{u}^{k-1} - \tau\mathbf{A}^\top(\mathbf{A}\mathbf{u}^{k-1} - \mathbf{y}) \right) - \tau\mathbf{A}^\top \left(\mathbf{A} \left(\mathcal{T}_{\tau\lambda} \left(\mathbf{u}^{k-1} - \tau\mathbf{A}^\top(\mathbf{A}\mathbf{u}^{k-1} - \mathbf{y}) \right) - \mathbf{y} \right) \right) \right) \\ &= \dots \\ &= \mathcal{T}_{\tau\lambda} \left(\mathcal{T}_{\tau\lambda} \left(\dots \left(\mathbf{u}^0 - \tau\mathbf{A}^\top(\mathbf{A}\mathbf{u}^0 - \mathbf{y}) \right) \dots \right) \right)\end{aligned}$$

Same parameters $\{\mathbf{A}, \tau, \lambda\}$ shared across the network (recurrent NN).

The ISTA “network”

ISTA

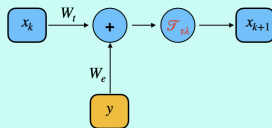
For $k = 0, \dots, K-1$

$$x_{k+1} = \mathcal{F}_{\tau\lambda}((\text{Id} - \tau A^T A)x_k + \tau A^T y)$$

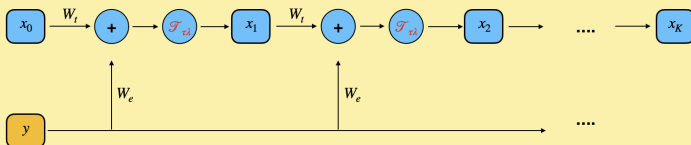
end

$$W_t := \text{Id} - \tau A^T A$$

$$W_e := \tau A^T$$



Stacking



Parameters of the network are: $\theta = \{W_t, W_e, \tau\}$.

Learned ISTA (LISTA) & generalisations

- Given a dataset $\{(\mathbf{u}_j, \mathbf{y}_j)\}$, $j = 1, \dots, J$, learn $\theta = \{\mathbf{W}_t, \mathbf{W}_e, \tau\}$ by minimising²:

$$\mathcal{L}(\mathbf{W}_t, \mathbf{W}_e, \tau) := \frac{1}{2J} \sum_{j=1}^J \|\mathbf{u}_j^K(\mathbf{W}_t, \mathbf{W}_e, \tau; \mathbf{y}_j) - \mathbf{u}_j\|^2$$

- $\mathbf{u}_j^K(\mathbf{W}_t, \mathbf{W}_e, \tau; \mathbf{y}_j)$ is computed via K iterations of ISTA.
- Backpropagation + SGD /ADAM for training

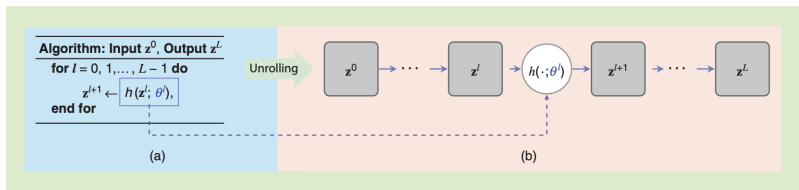
²Gregor, LeCun, '10

Learned ISTA (LISTA) & generalisations

- Given a dataset $\{(\mathbf{u}_j, \mathbf{y}_j)\}$, $j = 1, \dots, J$, learn $\theta = \{\mathbf{W}_t, \mathbf{W}_e, \tau\}$ by minimising²:

$$\mathcal{L}(\mathbf{W}_t, \mathbf{W}_e, \tau) := \frac{1}{2J} \sum_{j=1}^J \|\mathbf{u}_j^K(\mathbf{W}_t, \mathbf{W}_e, \tau; \mathbf{y}_j) - \mathbf{u}_j\|^2$$

- $\mathbf{u}_j^K(\mathbf{W}_t, \mathbf{W}_e, \tau; \mathbf{y}_j)$ is computed via K iterations of ISTA.
- Backpropagation + SGD /ADAM for training



General structure of algorithmic unrolling.

²Gregor, LeCun, '10

Let now the parameters change at each layer, so that $\theta = \{\mathbf{W}_t^k, \mathbf{W}_e^k, \tau_k, \lambda_k\}_{k=0}^{K-1}$.

Learn weight matrices + non-linearities end-to-end ³.

$$\mathcal{L}(\{\mathbf{W}_t^k, \mathbf{W}_e^k, \tau_k, \lambda_k\}_{k=0}^{K-1}) := \frac{1}{2J} \sum_{j=1}^J \|\mathbf{u}_j^K \left(\{\mathbf{W}_t^k, \mathbf{W}_e^k, \tau_k, \lambda_k\}_{k=0}^{K-1}; \mathbf{y}_j \right) - \mathbf{u}_j\|^2$$

- **Pro:** flexible and expressive. Same idea can be adapted to any iterative algorithm.
- **Con:** no clear interpretation from an optimisation viewpoint!

$$\mathbf{u}_j^K \left(\{\mathbf{W}_t^k, \mathbf{W}_e^k, \tau_k, \lambda_k\}_{k=0}^{K-1}; \mathbf{y}_j \right) \in \arg \min_{\mathbf{u}} ???$$

³Monga, Li, Eldar, 2021

LISTA for image super-resolution

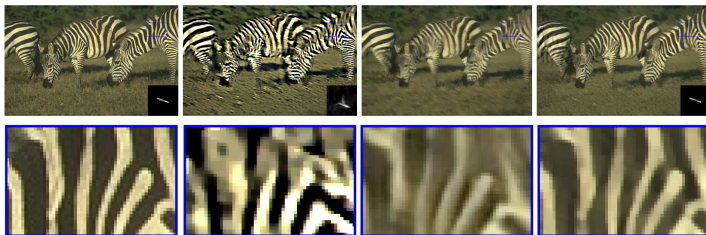


Figure 1: Left: GT, centre: standard approaches, Right: Unrolled version

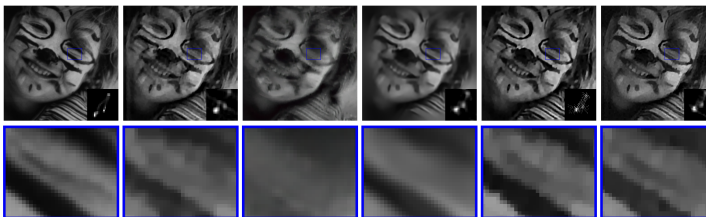


Figure 2: Left: GT, centre: standard approaches, Right: Unrolled version

Example: learn a customised TV prior via PGD

$$\mathbf{y} = \mathbf{A}\mathbf{u} + \mathbf{n}, \quad \min_{\mathbf{u}} \frac{1}{2} \|\mathbf{A}\mathbf{u} - \mathbf{y}\|^2 + \lambda \text{TV}_{\xi}(\mathbf{u})$$

where $\text{TV}_{\xi} = \|\mathbf{D}\mathbf{u}\|_{2,1,\xi}$. Consider K iterations of:

$$\mathbf{u}^{k+1} = \text{prox}_{\tau\lambda\text{TV}_{\epsilon}}(\mathbf{u}^k - \tau\mathbf{A}^{\top}(\mathbf{A}\mathbf{u}^k - \mathbf{y})) \quad \leftrightarrow \quad \mathbf{u}^{k+1} = \text{prox}_{\tau_k\lambda_k\text{TV}_{\epsilon}}(\mathbf{u}^k - \tau_k\mathbf{A}^{\top}(\mathbf{A}\mathbf{u}^k - \mathbf{y}))$$

Learn optimal $\boldsymbol{\theta} = (\tau_0, \dots, \tau_{K-1}, \lambda_0, \dots, \lambda_{K-1})$

Example: learn a customised TV prior via PGD

$$\mathbf{y} = \mathbf{A}\mathbf{u} + \mathbf{n}, \quad \min_{\mathbf{u}} \frac{1}{2} \|\mathbf{A}\mathbf{u} - \mathbf{y}\|^2 + \lambda \text{TV}_{\xi}(\mathbf{u})$$

where $\text{TV}_{\xi} = \|\mathbf{D}\mathbf{u}\|_{2,1,\xi}$. Consider K iterations of:

$$\mathbf{u}^{k+1} = \text{prox}_{\tau\lambda\text{TV}_{\epsilon}}(\mathbf{u}^k - \tau\mathbf{A}^{\top}(\mathbf{A}\mathbf{u}^k - \mathbf{y})) \quad \leftrightarrow \quad \mathbf{u}^{k+1} = \text{prox}_{\tau_k\lambda_k\text{TV}_{\epsilon}}(\mathbf{u}^k - \tau_k\mathbf{A}^{\top}(\mathbf{A}\mathbf{u}^k - \mathbf{y}))$$

Learn optimal $\boldsymbol{\theta} = (\tau_0, \dots, \tau_{K-1}, \lambda_0, \dots, \lambda_{K-1})$

Train using dataset $\{(\mathbf{u}_j, \mathbf{y}_j)\}$ via:

$$\boldsymbol{\theta} \in \arg \min_{\boldsymbol{\theta}} \frac{1}{2J} \sum_j \|\mathbf{u}_j - \mathbf{u}_j^K(\boldsymbol{\theta}; \mathbf{y}_j)\|^2$$

where $\mathbf{u}_j^K(\boldsymbol{\theta}; \mathbf{y}_j) = \text{PGD}_{\boldsymbol{\theta}}^K(\mathbf{y}_j)$.

DeepInv example https://deepinv.github.io/deepinv/auto_examples/unfolded/demo_custom_prior_unfolded.html.

Deep unfolding/unrolling

Unrolled primal-dual (no extrapolation)

Unrolling primal–dual algorithms

$$\mathbf{v}^{k+1} = \text{prox}_{\eta g^*}(\mathbf{v}^k + \eta \mathbf{K} \mathbf{u}^k)$$

$$\mathbf{u}^{k+1} = \text{prox}_{\tau f_{\mathbf{y}}}(\mathbf{u}^k - \tau \mathbf{K}^{\top} \mathbf{v}^{k+1})$$

Unrolling primal–dual algorithms

$$\mathbf{v}^{k+1} = \text{prox}_{\eta g^*}(\mathbf{v}^k + \eta \mathbf{K} \mathbf{u}^k)$$

$$\mathbf{u}^{k+1} = \text{prox}_{\tau f_{\mathbf{y}}}(\mathbf{u}^k - \tau \mathbf{K}^\top \mathbf{v}^{k+1})$$

Unrolled version: general form (simplified version: $\theta_{1,k} \equiv \theta_1$, $\theta_{2,k} \equiv \theta_2$)

$$\begin{aligned} \mathbf{v}^{k+1} &= \Psi_{\theta_{1,k}}(\mathbf{v}^k + \eta_k \mathbf{K} \mathbf{u}^k) \\ \mathbf{u}^{k+1} &= \Phi_{\theta_{2,k}}(\mathbf{u}^k - \tau_k \mathbf{K}^\top \mathbf{v}^{k+1}) \end{aligned}$$

- Step-sizes $\{\tau_k, \eta_k\}$ /prox operators $\Psi_{\theta_{1,k}}$ and $\Phi_{\theta_{2,k}}$ are learned end-to-end for all k .
- Output: $\mathbf{u}^K(\theta; \mathbf{y})$, with $\theta = \{\theta_{1,0}, \dots, \theta_{1,K-1}, \tau_0, \dots, \tau_{K-1}, \theta_{2,0}, \dots, \theta_{2,K-1}, \eta_0, \dots, \eta_{K-1}\}$.

→ K -layer NN inspired by a primal–dual solver, trained using samples $\{(\mathbf{u}_j, \mathbf{y}_j)\}$ by:

$$\min_{\theta} \frac{1}{2J} \sum_{j=1}^J \|\mathbf{u}^K(\theta; \mathbf{y}_j) - \mathbf{u}_j\|^2.$$

(Adler & Öktem, 2017) → [Tatiana's course for applications to CT](#).

Learning TV-parameter maps via primal-dual unrolling⁴

$$\arg \min_{\mathbf{u} \in \mathbb{R}^N} \frac{1}{2} \|\mathbf{A}\mathbf{u} - \mathbf{y}\|^2 + \|\mathbf{NN}_{\theta}(\mathbf{z})\mathbf{D}\mathbf{u}\|_{2,1}$$

where $\mathbf{NN}_{\theta}(\mathbf{z})$ parametrises an adaptive spatio-temporal regularisation map Λ and $\mathbf{A} = \mathbf{S}\mathcal{F}$ is a subsampled Fourier transform (MRI).

⁴Kofler, Papafitsoros et al. '23

Deep unfolding/unrolling

To unroll or not to unroll?

$$\mathbf{u}^{k+1} = \text{Algo}_{\theta_k}(\mathbf{u}^k)$$

$$\mathbf{u}^{k+1} = \text{Algo}_{\theta_k}(\mathbf{u}^k)$$

Option 1: shared parameters

$$\theta_k \equiv \theta \quad \forall \text{ iteration/layer } k = 0, \dots, K - 1$$

- Recurrent NN structure
- Closer to classical optimisation
- Fewer parameters to estimate (step-sizes, but also inertia, pre-conditioners...)
- Better interpretability

Parameter sharing vs layer-specific parameters

$$\mathbf{u}^{k+1} = \text{Algo}_{\theta_k}(\mathbf{u}^k)$$

Option 1: shared parameters

$$\theta_k \equiv \theta \quad \forall \text{ iteration/layer } k = 0, \dots, K - 1$$

- Recurrent NN structure
- Closer to classical optimisation
- Fewer parameters to estimate (step-sizes, but also inertia, pre-conditioners...)
- Better interpretability

Option 2: layer-specific parameters

$$\theta_k \neq \theta_j, \quad \text{for iterations/layers } k \neq j$$

- (Significantly) more expressive
- Underlying objective changes at each layer
- Harder theoretical interpretation

$$\mathbf{u}^{k+1} = \text{Algo}_{\theta_k}(\mathbf{u}^k)$$

Option 1: shared parameters

$$\theta_k \equiv \theta \quad \forall \text{ iteration/layer } k = 0, \dots, K - 1$$

- Recurrent NN structure
- Closer to classical optimisation
- Fewer parameters to estimate (step-sizes, but also inertia, pre-conditioners...)
- Better interpretability

Option 2: layer-specific parameters

$$\theta_k \neq \theta_j, \quad \text{for iterations/layers } k \neq j$$

- (Significantly) more expressive
- Underlying objective changes at each layer
- Harder theoretical interpretation

Trade-off: structure vs expressivity

Computational cost and memory trade-offs

	Classical solver	Unrolled network
Iterations	Until convergence	Fixed number K
Training cost	None	Potentially large
Inference cost	Variable	Fixed
Memory usage	Low	High
Interpretability	High	Low/Moderate

- Reduced iterative runtime VS. offline training.
- Inference becomes predictable and fast.
- Memory footprint increases due to backpropagation.

Implicit bias: learning the algorithm

Classical optimisation:

$$\mathbf{u}^* \in \arg \min_{\mathbf{u}} F(\mathbf{u}) \quad \rightarrow \quad \mathbf{u}^* \approx \mathbf{u}^K = \text{Algo}^K(\text{params}, \mathbf{y}, \mathbf{u}^0)$$

Unrolling:

$$\mathbf{u}^K = \text{Algo}_{\theta}^K(\mathbf{y}, \mathbf{u}^0) \quad \rightarrow \quad \mathbf{u}^K \approx \arg \min \quad ??$$

Reasons:

- We are not learning a mapping $\mathbf{y} \mapsto \mathbf{u}$.
- We are learning:
 - Step-sizes
 - Metrics/preconditioners
 - Thresholding rules
 - Update directions

Interpretation:

- Learning modifies the *geometry* of optimisation.
- Learned weights can act as adaptive preconditioners.

Unrolling learns **how** to optimise, not **what** to optimise.

What happens as $K \rightarrow +\infty$?

Do not truncate at fixed depth, consider fixed-point eq. ($\mathbf{u}^{k+1} \approx \mathbf{u}^k$ at convergence)

$$\mathbf{u}^\infty = \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})$$

for **shared** parameters θ .

Deep Equilibrium Models: infinitely deep NNs defined by the structure of one layer.

- Do not store all layers
- Directly look at fixed-point equations
- Backpropagate via implicit differentiation

Deep Equilibrium Models

Deep Equilibrium Models (DEQs)

Main idea: unrolling with $K \rightarrow \infty$ + weight tying ($\theta_k \equiv \theta$ for all $k \geq 0$)⁵.

$$\mathbf{u}^{k+1} = \text{Algo}_{\theta}(\mathbf{u}^k; \mathbf{y}), \quad k = 0, \dots, +\infty$$

Fixed point condition:

$$\mathbf{u}^{\infty} = \mathbf{u}^{\infty}(\theta; \mathbf{y}) = \text{Algo}_{\theta}(\mathbf{u}^{\infty}; \mathbf{y})$$

⁵Bai, Kolter, Koltun, '19, Gilton, Ongie, Willett, '21

⁶Crockett, Fessler, '18, Franceschi et al, '18

Deep Equilibrium Models (DEQs)

Main idea: unrolling with $K \rightarrow \infty$ + weight tying ($\theta_k \equiv \theta$ for all $k \geq 0$)⁵.

$$\mathbf{u}^{k+1} = \text{Algo}_{\theta}(\mathbf{u}^k; \mathbf{y}), \quad k = 0, \dots, +\infty$$

Fixed point condition:

$$\mathbf{u}^{\infty} = \mathbf{u}^{\infty}(\theta; \mathbf{y}) = \text{Algo}_{\theta}(\mathbf{u}^{\infty}; \mathbf{y})$$

Connection with **Bilevel Optimisation**⁶: DEQ as implicit solver of lower-level problem

$$\mathbf{u}^{\infty}(\theta; \mathbf{y}) \in \arg \min_{\mathbf{u} \in \mathbb{R}^N} F(\mathbf{u}; \mathbf{y}) \quad \text{or at least} \quad \mathbf{0} \in \partial F(\mathbf{u}^{\infty}(\theta; \mathbf{y}); \mathbf{y})$$

with training objective/upper level cost being, e.g.:

$$\min_{\theta} \frac{1}{2J} \sum_{j=1}^J \|\mathbf{u}^{\infty}(\theta; \mathbf{y}) - \mathbf{u}_j\|^2$$

- Lower level: solve fixed-point/minimisation problem (e.g., $F(\mathbf{u}) = f_{\mathbf{y}}(\mathbf{u}) + g_{\theta}(\mathbf{u})$).
- Upper level: optimise θ through that solution

⁵Bai, Kolter, Koltun, '19, Gilton, Ongie, Willett, '21

⁶Crockett, Fessler, '18, Franceschi et al, '18

Deep Equilibrium Models (DEQs)

Main idea: unrolling with $K \rightarrow \infty$ + weight tying ($\theta_k \equiv \theta$ for all $k \geq 0$)⁷.

$$\mathbf{u}^{k+1} = \text{Algo}_{\theta}(\mathbf{u}^k; \mathbf{y}), \quad k = 0, \dots, +\infty$$

Fixed point condition:

$$\mathbf{u}^{\infty} = \mathbf{u}^{\infty}(\theta; \mathbf{y}) = \text{Algo}_{\theta}(\mathbf{u}^{\infty}; \mathbf{y})$$

Sufficient condition for existence (and uniqueness) of fixed points? **Contractivity!**

$$\exists \rho < 1 \quad \text{s.t.} \quad \|\text{Algo}_{\theta}(\mathbf{u}_1; \mathbf{y}) - \text{Algo}_{\theta}(\mathbf{u}_2; \mathbf{y})\| \leq \rho \|\mathbf{u}_1 - \mathbf{u}_2\|, \quad \forall \mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^N.$$

If true: $\mathbf{u}^k \rightarrow \mathbf{u}^{\infty}$ (Banach's fixed point theorem).

⁷Bai, Kolter, Koltun, '19, Gilton, Ongie, Willett, '21

Deep Equilibrium Models (DEQs)

Main idea: unrolling with $K \rightarrow \infty$ + weight tying ($\theta_k \equiv \theta$ for all $k \geq 0$)⁷.

$$\mathbf{u}^{k+1} = \text{Algo}_{\theta}(\mathbf{u}^k; \mathbf{y}), \quad k = 0, \dots, +\infty$$

Fixed point condition:

$$\mathbf{u}^{\infty} = \mathbf{u}^{\infty}(\theta; \mathbf{y}) = \text{Algo}_{\theta}(\mathbf{u}^{\infty}; \mathbf{y})$$

Sufficient condition for existence (and uniqueness) of fixed points? **Contractivity!**

$$\exists \rho < 1 \quad \text{s.t.} \quad \|\text{Algo}_{\theta}(\mathbf{u}_1; \mathbf{y}) - \text{Algo}_{\theta}(\mathbf{u}_2; \mathbf{y})\| \leq \rho \|\mathbf{u}_1 - \mathbf{u}_2\|, \quad \forall \mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^N.$$

If true: $\mathbf{u}^k \rightarrow \mathbf{u}^{\infty}$ (Banach's fixed point theorem).

Gradient Descent (GD)

Let $F = f_{\mathbf{y}} + g$ be L -smooth,

$$\text{GD}(\mathbf{u}) = \mathbf{u} - \tau \nabla F(\mathbf{u}) = (\mathbf{I} - \tau \nabla f)(\mathbf{u})$$

If F is μ -strongly convex and $\tau \in (0, 2/L)$, then GD is a contraction with:

$$\rho = \max\{|1 - \tau\mu|, |1 - \tau L|\} < 1.$$

⁷Bai, Kolter, Koltun, '19, Gilton, Ongie, Willett, '21

Deep Equilibrium Models (DEQs)

Main idea: unrolling with $K \rightarrow \infty$ + weight tying ($\theta_k \equiv \theta$ for all $k \geq 0$)⁷.

$$\mathbf{u}^{k+1} = \text{Algo}_{\theta}(\mathbf{u}^k; \mathbf{y}), \quad k = 0, \dots, +\infty$$

Fixed point condition:

$$\mathbf{u}^{\infty} = \mathbf{u}^{\infty}(\theta; \mathbf{y}) = \text{Algo}_{\theta}(\mathbf{u}^{\infty}; \mathbf{y})$$

Sufficient condition for existence (and uniqueness) of fixed points? **Contractivity!**

$$\exists \rho < 1 \quad \text{s.t.} \quad \|\text{Algo}_{\theta}(\mathbf{u}_1; \mathbf{y}) - \text{Algo}_{\theta}(\mathbf{u}_2; \mathbf{y})\| \leq \rho \|\mathbf{u}_1 - \mathbf{u}_2\|, \quad \forall \mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^N.$$

If true: $\mathbf{u}^k \rightarrow \mathbf{u}^{\infty}$ (Banach's fixed point theorem).

Proximal Gradient Descent (PGD)

$F = f_{\mathbf{y}} + g$, with $f_{\mathbf{y}}$ L -smooth, g convex/non-smooth.

$$\text{PGD}(\mathbf{u}) = \text{prox}_{\tau g}(\mathbf{u} - \tau \nabla f_{\mathbf{y}}(\mathbf{u})).$$

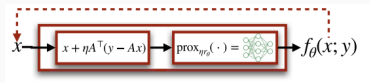
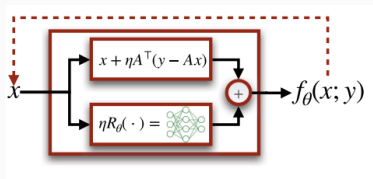
If $f_{\mathbf{y}}$ is strongly convex, then PGD is a contraction for suitably small choices of τ .

⁷Bai, Kolter, Koltun, '19, Gilton, Ongie, Willett, '21

Deep Equilibrium Models

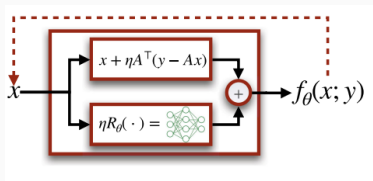
DEQ-GD & DEQ-PGD

DEQ-GD and DEQ-PGD



DEQ-GD ($\mathbf{NN}_\theta = \nabla g_\theta$) and DEQ-PGD ($\mathbf{NN}_\theta = \text{prox}_{g_\theta}$)

DEQ-GD and DEQ-PGD



DEQ-GD ($\mathbf{NN}_\theta = \nabla g_\theta$) and DEQ-PGD ($\mathbf{NN}_\theta = \text{prox}_{g_\theta}$)

DEQ-GD

Smooth $f_y + g \rightarrow \nabla f_y + \mathbf{NN}_\theta$. By triangle inequality:

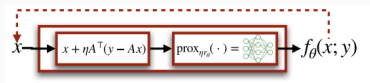
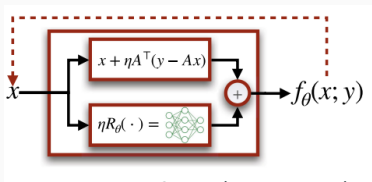
$$\|\text{GD}_\theta(\mathbf{u}_1) - \text{GD}_\theta(\mathbf{u}_2)\| \leq \|(\mathbf{I} - \tau \nabla f_y)(\mathbf{u}_1 - \mathbf{u}_2)\| + \tau \|\mathbf{NN}_\theta(\mathbf{u}_1) - \mathbf{NN}_\theta(\mathbf{u}_2)\|$$

If the following set of conditions apply:

- f_y is μ -strongly convex and L -smooth
- \mathbf{NN}_θ is $L_{\mathbf{NN}_\theta}$ -Lipschitz
- $\tau \in (0, 2/(L + L_{\mathbf{NN}_\theta}))$

then GD_θ is a contraction.

DEQ-GD and DEQ-PGD



DEQ-GD ($\mathbf{NN}_\theta = \nabla g_\theta$) and DEQ-PGD ($\mathbf{NN}_\theta = \text{prox}_{g_\theta}$)

DEQ-PGD

Non-smooth composite $f_y + g \rightarrow \mathbf{NN}_\theta \approx \text{prox}_g$, with \mathbf{NN}_θ $L_{\mathbf{NN}_\theta}$ -Lipschitz:

$$\begin{aligned} \|\text{PGD}_\theta(\mathbf{u}_1) - \text{PGD}_\theta(\mathbf{u}_2)\| &= \|\mathbf{NN}_\theta(I - \tau \nabla f_y)(\mathbf{u}_1) - \mathbf{NN}_\theta(I - \tau \nabla f_y)(\mathbf{u}_2)\| \\ &\leq L_{\mathbf{NN}_\theta} \|(I - \tau \nabla f_y)(\mathbf{u}_1 - \mathbf{u}_2)\|. \end{aligned}$$

If the following set of conditions apply:

- f_y is μ -strongly convex and L -smooth,
- $L_{\mathbf{NN}_\theta} \leq 1$ (e.g., non-expansive map),

then, for $\tau \in (0, 2/L)$, PGD_θ is a contraction.

Deep Equilibrium Models

DEQ training

Given samples $\{(\mathbf{u}_j, \mathbf{y}_j)\}$, train by minimising:

$$\min_{\theta} \frac{1}{2J} \sum_{j=1}^J \|\mathbf{u}_j^{\infty} - \mathbf{u}_j\|^2, \quad \mathbf{u}_j^{\infty} = \mathbf{u}_j^{\infty}(\theta; \mathbf{y}_j) = \text{Algo}_{\theta}(\mathbf{u}_j^{\infty}; \mathbf{y}_j)$$

1. **Forward pass:** compute \mathbf{u}_j^{∞} efficiently. (Fast) fixed-point iteration techniques (e.g., Anderson acceleration⁸)/iterative solvers with high-precision ($\epsilon \approx 10^{-9}$).
2. **Backward pass:** compute network updates. Supervised setting:

$$\nabla_{\theta} \left(\frac{1}{2J} \sum_j \|\mathbf{u}_j^{\infty} - \mathbf{u}_j\|^2 \right) = \frac{1}{2J} \sum_j \nabla_{\theta} \left(\|\mathbf{u}_j^{\infty} - \mathbf{u}_j\|^2 \right)$$

and for each term:

$$\nabla_{\theta} \left(\|\mathbf{u}_j^{\infty} - \mathbf{u}_j\|^2 \right) = \frac{d\mathbf{u}_j^{\infty}}{d\theta}{}^T \nabla_{\mathbf{u}} \|\mathbf{u}_j^{\infty} - \mathbf{u}_j\|^2 = 2 \frac{d\mathbf{u}_j^{\infty}}{d\theta}{}^T (\mathbf{u}_j^{\infty} - \mathbf{u}_j)$$

... where $d \cdot$ denotes the total derivative, as $\mathbf{u}^{\infty}(\theta) = \text{Algo}_{\theta}(\mathbf{u}^{\infty}(\theta); \mathbf{y})$.

⁸Walker, Ni, '11

Implicit differentiation (DEQ backward pass)

Fixed-point equation:

$$\mathbf{u}^\infty = \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})$$

Goal: compute $\frac{d\mathbf{u}^\infty}{d\theta}$ without backpropagating through all forward iterations.

Implicit differentiation (DEQ backward pass)

Fixed-point equation:

$$\mathbf{u}^\infty = \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})$$

Rewrite as a root-finding constraint:

$$\mathbf{F}(\mathbf{u}, \theta) := \mathbf{u} - \text{Algo}_\theta(\mathbf{u}; \mathbf{y}), \quad \mathbf{F}(\mathbf{u}^\infty, \theta) = \mathbf{0}.$$

Take total derivative w.r.t. θ :

$$\left. \frac{\partial \mathbf{F}}{\partial \mathbf{u}} \right|_{\mathbf{u}^\infty} \frac{d\mathbf{u}^\infty}{d\theta} + \left. \frac{\partial \mathbf{F}}{\partial \theta} \right|_{\mathbf{u}^\infty} = \mathbf{0}.$$

Implicit differentiation (DEQ backward pass)

Fixed-point equation:

$$\mathbf{u}^\infty = \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})$$

Rewrite as a root-finding constraint:

$$\mathbf{F}(\mathbf{u}, \theta) := \mathbf{u} - \text{Algo}_\theta(\mathbf{u}; \mathbf{y}), \quad \mathbf{F}(\mathbf{u}^\infty, \theta) = \mathbf{0}.$$

Take total derivative w.r.t. θ :

$$\left. \frac{\partial \mathbf{F}}{\partial \mathbf{u}} \right|_{\mathbf{u}^\infty} \frac{d\mathbf{u}^\infty}{d\theta} + \left. \frac{\partial \mathbf{F}}{\partial \theta} \right|_{\mathbf{u}^\infty} = \mathbf{0}.$$

Compute the two Jacobians:

$$\frac{\partial \mathbf{F}}{\partial \mathbf{u}} = \mathbf{I} - \frac{\partial \text{Algo}_\theta(\mathbf{u}; \mathbf{y})}{\partial \mathbf{u}}, \quad \frac{\partial \mathbf{F}}{\partial \theta} = -\frac{\partial \text{Algo}_\theta(\mathbf{u}; \mathbf{y})}{\partial \theta}.$$

Implicit differentiation (DEQ backward pass)

Fixed-point equation:

$$\mathbf{u}^\infty = \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})$$

Rewrite as a root-finding constraint:

$$\mathbf{F}(\mathbf{u}, \theta) := \mathbf{u} - \text{Algo}_\theta(\mathbf{u}; \mathbf{y}), \quad \mathbf{F}(\mathbf{u}^\infty, \theta) = \mathbf{0}.$$

Take total derivative w.r.t. θ :

$$\left. \frac{\partial \mathbf{F}}{\partial \mathbf{u}} \right|_{\mathbf{u}^\infty} \frac{d\mathbf{u}^\infty}{d\theta} + \left. \frac{\partial \mathbf{F}}{\partial \theta} \right|_{\mathbf{u}^\infty} = \mathbf{0}.$$

Evaluate at \mathbf{u}^∞ and substitute:

$$\left(\mathbf{I} - \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \mathbf{u}} \right) \frac{d\mathbf{u}^\infty}{d\theta} = \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \theta}.$$

Implicit differentiation (DEQ backward pass)

Fixed-point equation:

$$\mathbf{u}^\infty = \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})$$

Rewrite as a root-finding constraint:

$$\mathbf{F}(\mathbf{u}, \theta) := \mathbf{u} - \text{Algo}_\theta(\mathbf{u}; \mathbf{y}), \quad \mathbf{F}(\mathbf{u}^\infty, \theta) = \mathbf{0}.$$

Take total derivative w.r.t. θ :

$$\left. \frac{\partial \mathbf{F}}{\partial \mathbf{u}} \right|_{\mathbf{u}^\infty} \frac{d\mathbf{u}^\infty}{d\theta} + \left. \frac{\partial \mathbf{F}}{\partial \theta} \right|_{\mathbf{u}^\infty} = \mathbf{0}.$$

Evaluate at \mathbf{u}^∞ and substitute:

$$\left(\mathbf{I} - \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \mathbf{u}} \right) \frac{d\mathbf{u}^\infty}{d\theta} = \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \theta}.$$

Hence,

$$\frac{d\mathbf{u}^\infty}{d\theta} = \left(\mathbf{I} - \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \mathbf{u}} \right)^{-1} \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \theta}.$$

Key point: no unrolling, only Jacobians evaluated at the fixed point.

Implicit differentiation and Jacobian-Free Backpropagation

$$\frac{d\mathbf{u}^\infty}{d\boldsymbol{\theta}} = \left(\mathbf{I} - \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \mathbf{u}} \right)^{-1} \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \boldsymbol{\theta}}$$

If

$$\left\| \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \mathbf{u}} \right\|_2 < 1,$$

then the inverse admits a Neumann series expansion:

$$(\mathbf{I} - \mathbf{J}_\theta)^{-1} = \sum_{k=0}^{\infty} \mathbf{J}_\theta^k, \quad \mathbf{J}_\theta := \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \mathbf{u}}$$

Hence

$$\frac{d\mathbf{u}^\infty}{d\boldsymbol{\theta}} = \sum_{k=0}^{\infty} \mathbf{J}_\theta^k \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \boldsymbol{\theta}}$$

Implicit differentiation and Jacobian-Free Backpropagation

$$\frac{d\mathbf{u}^\infty}{d\theta} = \left(\mathbf{I} - \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \mathbf{u}} \right)^{-1} \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \theta}$$

If

$$\left\| \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \mathbf{u}} \right\|_2 < 1,$$

then the inverse admits a Neumann series expansion:

$$(\mathbf{I} - \mathbf{J}_\theta)^{-1} = \sum_{k=0}^{\infty} \mathbf{J}_\theta^k, \quad \mathbf{J}_\theta := \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \mathbf{u}}$$

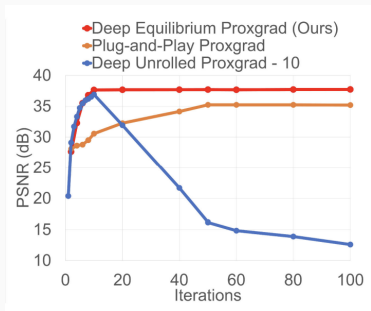
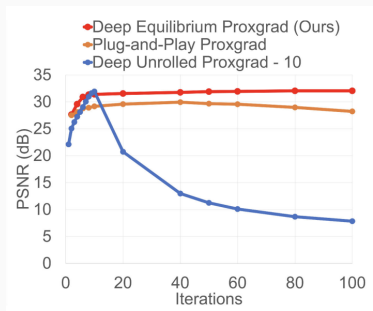
Hence

$$\frac{d\mathbf{u}^\infty}{d\theta} = \sum_{k=0}^{\infty} \mathbf{J}_\theta^k \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \theta} \underbrace{\approx}_{k=0} \mathbf{I} \frac{\partial \text{Algo}_\theta(\mathbf{u}^\infty; \mathbf{y})}{\partial \theta}$$

Jacobian-Free Backpropagation (Fung, Heaton et al. '22)

- Backprop through infinite layers without explicit Jacobian storage.
- Approximate hypergradient direction (not the steepest, but still a descent one)

Performance of DEQs - test time



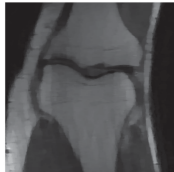
Performance of DEQs - test time



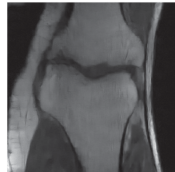
(a) Ground truth



(b) IFFT ($A^T y$), PSNR = 24.53 dB dB



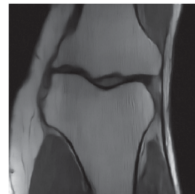
(c) TV-Regularized, PSNR =



(d) PnP-Prox, PSNR = 29.24 dB dB



(e) DU-PROX, PSNR = 31.02 dB dB



(f) DE-PROX, PSNR = 32.09 dB dB

Outlook: DEQs beyond convexity/contractivity

$$\mathbf{u}^* \arg \min_{\mathbf{u} \in \mathbb{R}^N} f_{\mathbf{y}}(\mathbf{u}) + g_{\theta}(\mathbf{u}), \quad g_{\theta}(\mathbf{u}) \text{ parameterised by } \mathbf{NN}_{\theta}$$

Example (RED): $g_{\theta}(\mathbf{u}) = \frac{\lambda}{2} \mathbf{u}^{\top} (\mathbf{u} - \mathbf{NN}_{\theta}(\mathbf{u}))$.

Outlook: DEQs beyond convexity/contractivity

$$\mathbf{u}^* \arg \min_{\mathbf{u} \in \mathbb{R}^N} f_{\mathbf{y}}(\mathbf{u}) + g_{\theta}(\mathbf{u}), \quad g_{\theta}(\mathbf{u}) \text{ parameterised by } \text{NN}_{\theta}$$

Example (RED): $g_{\theta}(\mathbf{u}) = \frac{\lambda}{2} \mathbf{u}^{\top} (\mathbf{u} - \text{NN}_{\theta}(\mathbf{u}))$.

More expressive networks and general data terms:

- No strongly convex/ L -smooth $f_{\mathbf{y}}$ (e.g., $\text{KL}(\mathbf{y}; \mathbf{A}\mathbf{u})$)
- Non-convex/non-smooth $g_{\theta}(\mathbf{u}) \rightarrow \text{prox}_{g_{\theta}}$? Non-expansiveness?

Generalisations are possible (see, e.g., Daniele, Villa, Vaiter, Calatroni, '26 for KL) beyond L -smoothness, convexity⁹.

Well-posed DEQ layer

Let $f_{\theta} : \mathbb{R}^N \rightarrow \mathbb{R}^N$. We say f_{θ} is well-posed if:

- $\text{Fix}(f_{\theta}) = \{\mathbf{u}^{\infty} : f_{\theta}(\mathbf{u}^{\infty}) = \mathbf{u}^{\infty}\} \neq \emptyset$.
- for a given $\mathbf{u}^0 \in \mathbb{R}^N$ and having defined:

$$\mathbf{u}^{k+1} = f_{\theta}(\mathbf{u}^k; \mathbf{u}^0), \quad \forall k \geq 0,$$

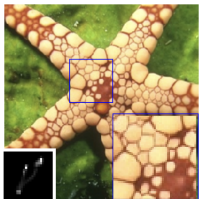
then $\mathbf{u}^k \rightarrow \mathbf{u}^{\infty} \in \text{Fix}(f_{\theta})$.

Convergence is possible through conditions implying only existence, not uniqueness.

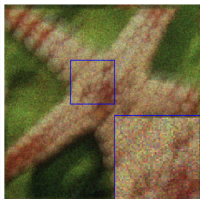
⁹Bolte et al., '18

Is it worth it?

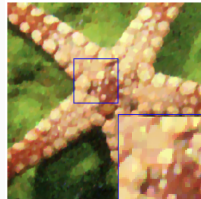
Yes! More expressive network + interpretability/convergence.



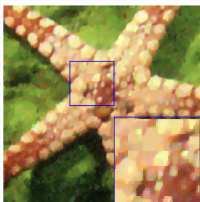
u



y , PSNR=13.36



u_{TV}^* , PSNR=22.39



u_{PnP}^* , PSNR=24.15



u_{DEQ}^* , PSNR=24.47

Outlook/conclusions

- DL can be integrated within optimisation algorithms improving performance, and preserving (sometimes) theoretical guarantees
- The understanding of these strategies passes through convex/non-convex calculus/operator theory
- **Main use:** estimating priors π_U (hard!)

Conclusions and outlook

- DL can be integrated within optimisation algorithms improving performance, and preserving (sometimes) theoretical guarantees
- The understanding of these strategies passes through convex/non-convex calculus/operator theory
- **Main use:** estimating priors π_U (hard!)

Some **challenging** open questions:

- **Plug & play:** theory is pretty developed, but adaptation to non-quadratic data models $D(\mathbf{y}; \mathbf{A}\mathbf{u})$ is scarce, Beyond denoisers?
- **Unrolling:** convergence guarantees are often missing, not clear variational interpretation, how to assess quality of the minimiser?
- **DEQs:** conditions for well-posedness in non-convex/non-contractive regimes + guarantees for training.

Thank you!
Questions?

luca.calatroni@unige.it