

AI4Science: Striking the Best Data-Knowledge Tradeoff

Ju Sun (Computer Sci. & Eng., UMN)

May 23, 2025

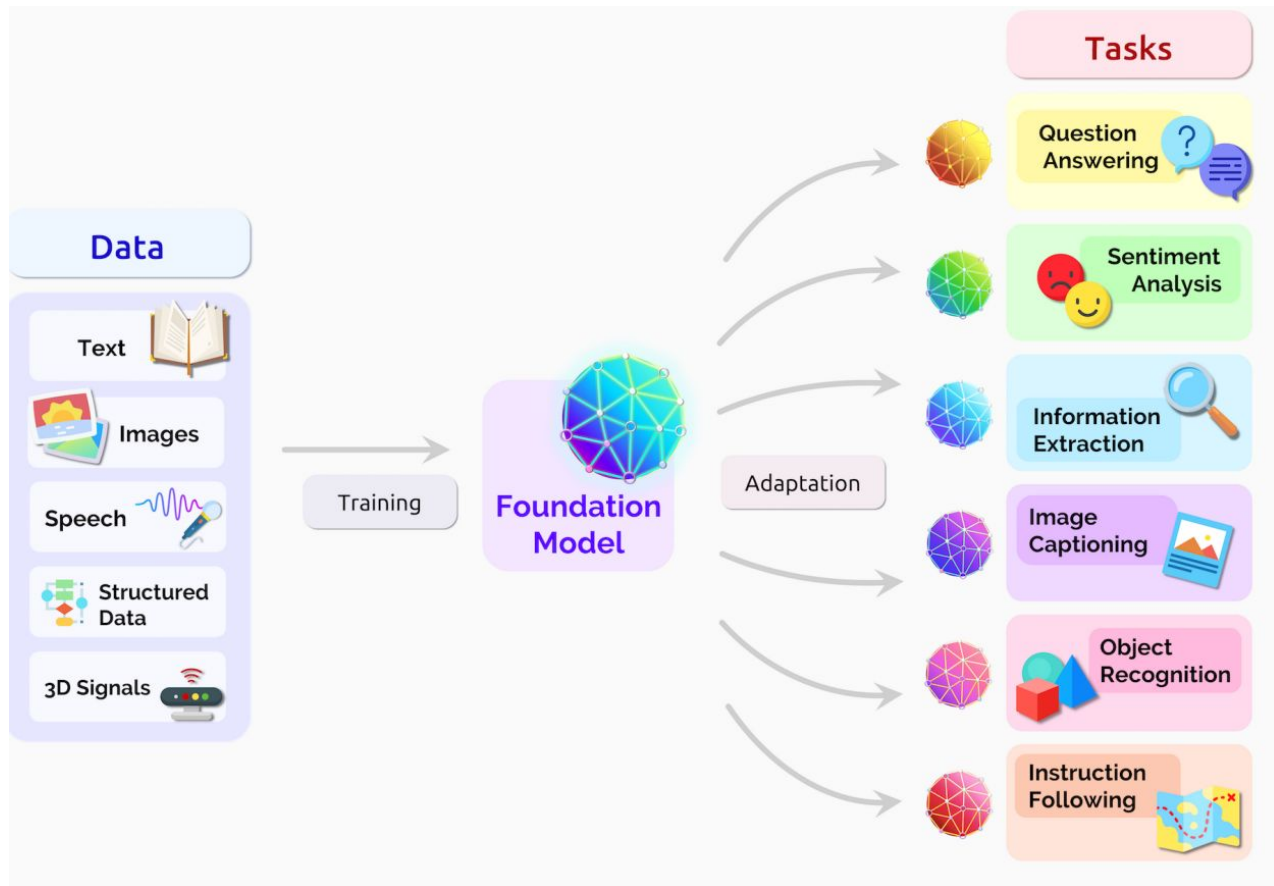
STROBE Seminar (@Physics & Astronomy, UCLA)



UNIVERSITY OF MINNESOTA

Driven to DiscoverSM

The “foundation model” movement



Credit: **On the Opportunities and Risks of Foundation Models**

<https://arxiv.org/abs/2108.07258>

CV/NLP domains are lucky

LAION 

Large-scale Artificial Intelligence Open Network

TRULY OPEN AI. 100% NON-PROFIT. 100%

LAION, as a non-profit organization, provides datasets and models to liberate machine learning research. By doing so, we encourage open public education and a more environmentally friendly use of resources by reusing existing datasets and models.

Re-LAION 5B release (30.08.2024)

TABLE 2: Statistics of commonly-used data sources.

| Corpora | Size | Source | Latest Update Time |
|------------------------------------|-------|--------------|--------------------|
| BookCorpus [158] | 5GB | Books | Dec-2015 |
| Gutenberg [159] | - | Books | Dec-2021 |
| C4 [82] | 800GB | CommonCrawl | Apr-2019 |
| CC-Stories-R [160] | 31GB | CommonCrawl | Sep-2019 |
| CC-NEWS [27] | 78GB | CommonCrawl | Feb-2019 |
| REALNEWS [161] | 120GB | CommonCrawl | Apr-2019 |
| OpenWebText [162] | 38GB | Reddit links | Mar-2023 |
| Pushift.io [163] | 2TB | Reddit links | Mar-2023 |
| Wikipedia [164] | 21GB | Wikipedia | Mar-2023 |
| BigQuery [165] | - | Codes | Mar-2023 |
| the Pile [166] | 800GB | Other | Dec-2020 |
| ROOTS [167] | 1.6TB | Other | Jun-2022 |

source: <https://arxiv.org/abs/2303.18223>

Not all fields are as lucky

Thrust B: How Should Domain Knowledge Be Incorporated into Supervised Machine Learning?

The central question for this thrust is “which knowledge should be leveraged in SciML, and how should this knowledge be included?” Any answers will naturally depend on the SciML task and computational budgets, thus mirroring standard considerations in traditional scientific computing.

Hard Constraints. One research avenue involves incorporation of domain knowledge through imposition of constraints that cannot be violated. These hard constraints could be enforced during training, replacing what typically is an unconstrained optimization problem with a constrained one. In general, such constraints could involve simulations or highly nonlinear functions of the training parameters. Therefore, there is a need to identify particular cases when constraint qualification conditions can be ensured as these conditions are necessary regularity conditions for constrained optimization [57–59]. Although incorporating constraints during training generally makes maximal use of training data, there may be additional opportunities to employ constraints at the time of prediction (e.g., by projecting predictions onto the region induced by the constraints).

Soft Constraints. A similar avenue for incorporating domain knowledge involves modifying the objective function (soft constraints) used in training. It is understood that ML loss function selection should be guided by the task and data. Therefore, opportunities exist for developing loss functions that incorporate domain knowledge and analyzing the resulting impact on solvability

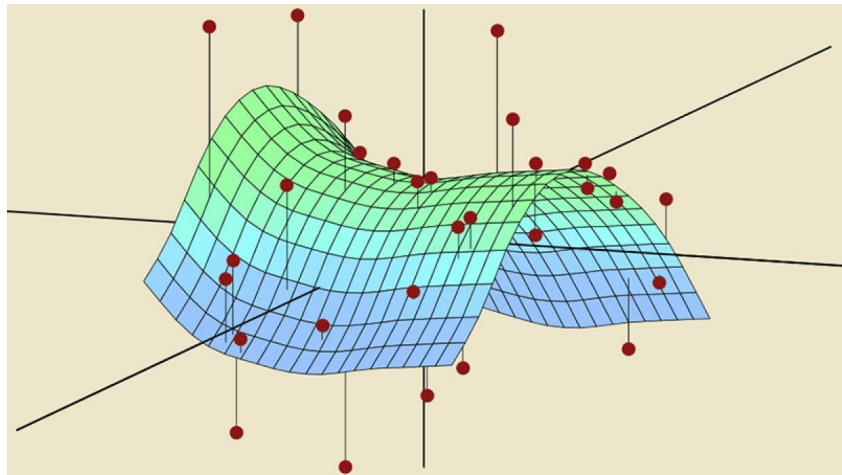


Ref <https://www.osti.gov/servlets/purl/1478744>

Domain-Aware Scientific Machine Learning

There's no free lunch!

(Self)-Supervised learning as data fitting



Typically, #data points we need grow **exponentially** with respect to dimension (i.e., **curse of dimensionality**)

Knowledge

Small-data AI

Building in prior knowledge is **crucial** for reducing the data complexity
e.g., “convolutional” layers



Large-data AI

Data

Today's talk:

several stories about data-knowledge tradeoffs

- Scientific inverse problems (SIPs)
 - Data-driven (data-rich) methods for SIPs
 - Single-instance (data-poor) methods for SIPs
- Principled computational tool for data-knowledge tradeoffs

Scientific Inverse Problems

Inverse problems

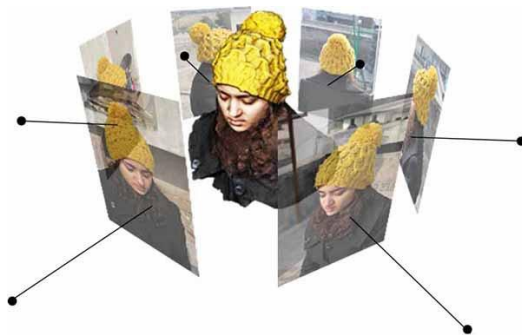
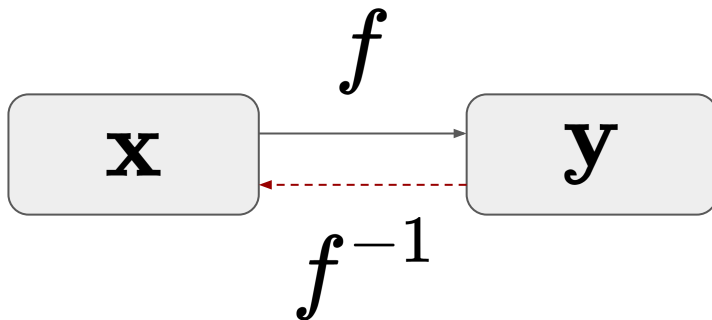
Inverse problem: given $\mathbf{y} = f(\mathbf{x})$, recover \mathbf{x}



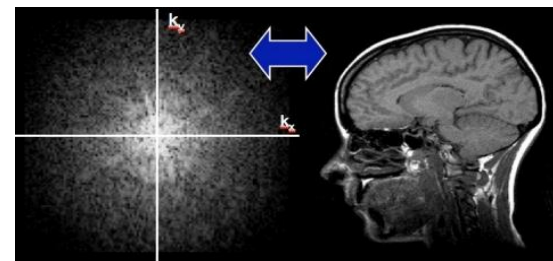
Image denoising



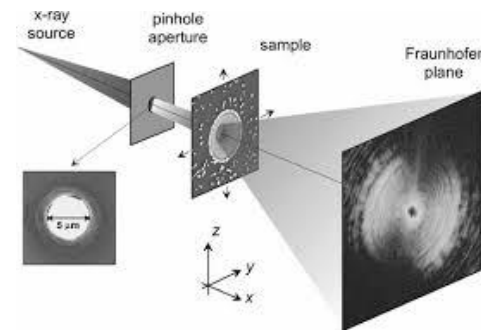
Image super-resolution



3D reconstruction



MRI reconstruction



Coherent diffraction imaging (CDI)

Traditional methods

Inverse problem: given $\mathbf{y} = f(\mathbf{x})$, recover \mathbf{x}

$$\min_{\mathbf{x}} \underbrace{\ell(\mathbf{y}, f(\mathbf{x}))}_{\text{data fitting}} + \lambda \underbrace{R(\mathbf{x})}_{\text{regularizer}} \quad \text{RegFit}$$

Limitations:

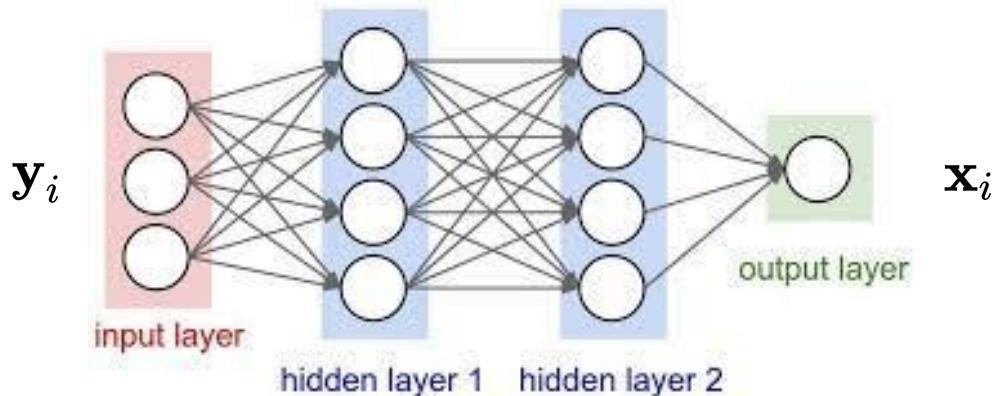
- Which ℓ ? (e.g., unknown/compound noise)
- Which R ? (e.g., structures not amenable to math description)
- Speed

DL has changed
everything

DL methods for SIPs: the radical/simplistic way

Inverse problem: given $\mathbf{y} = f(\mathbf{x})$, recover \mathbf{x}

Learn the f^{-1} with a training set $\{(\mathbf{y}_i, \mathbf{x}_i)\}$



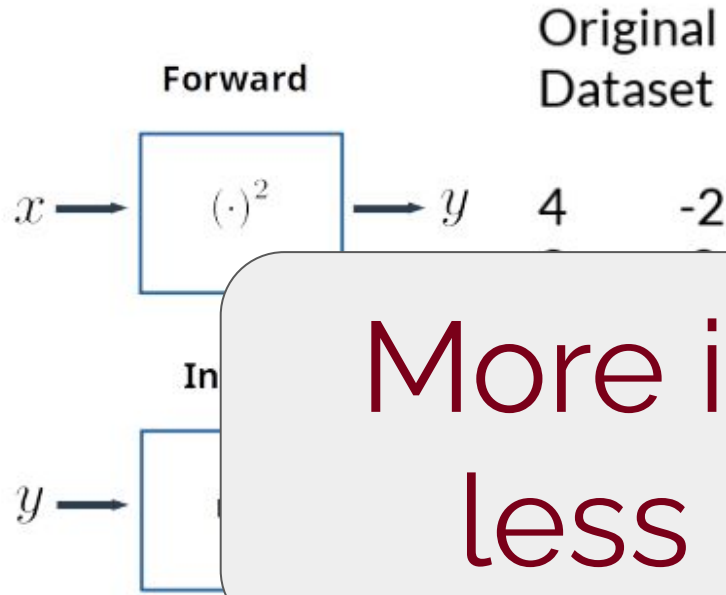
Limitations:

- Wasteful: not using f
- Representative data?
- Not always straightforward

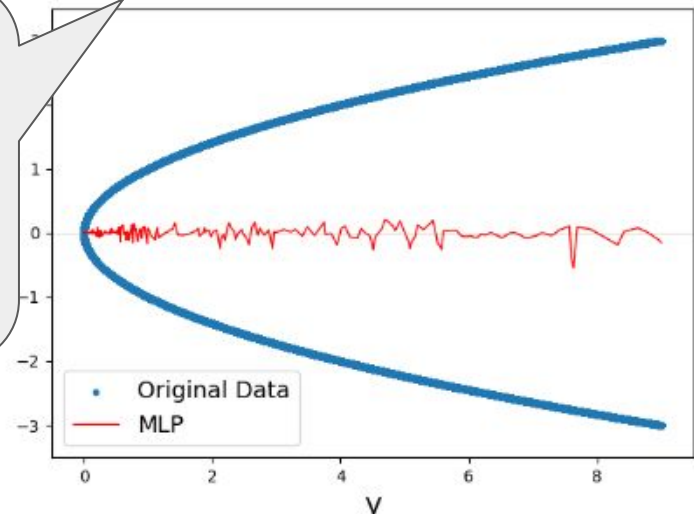
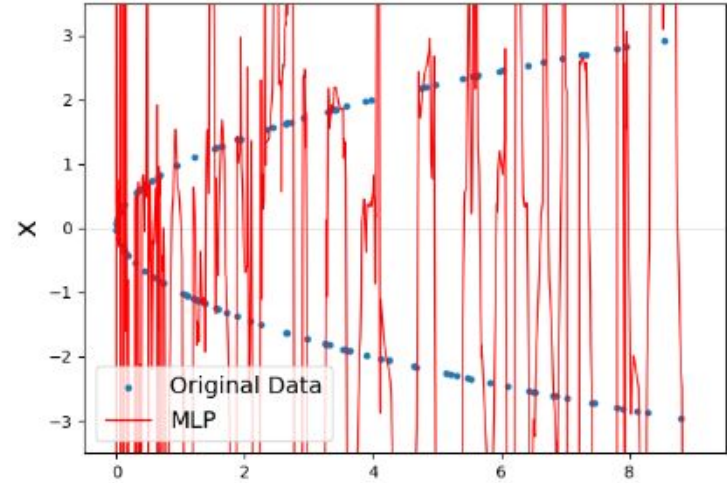
(see, e.g., Tayal et al. **Inverse Problems, Deep Learning, and Symmetry Breaking**.

<https://arxiv.org/abs/2003.09077>)

Story I: More could be less

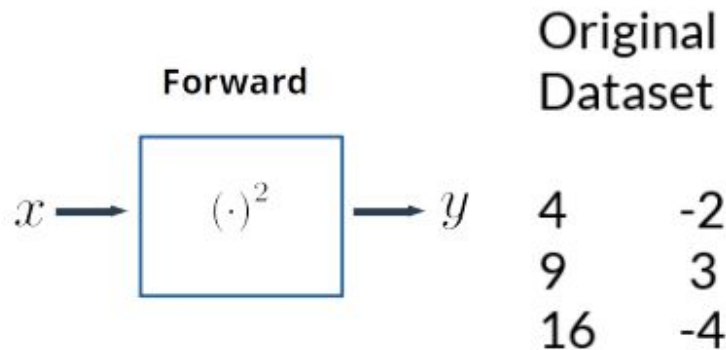


More is
less

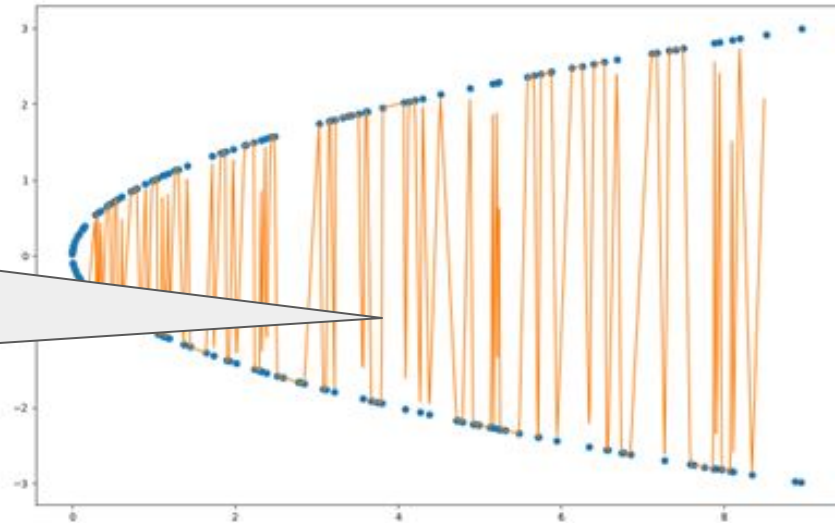


Why “more is less” here? **Forward symmetry:** $\{+\sqrt{y}, -\sqrt{y}\} \leftrightarrow y$

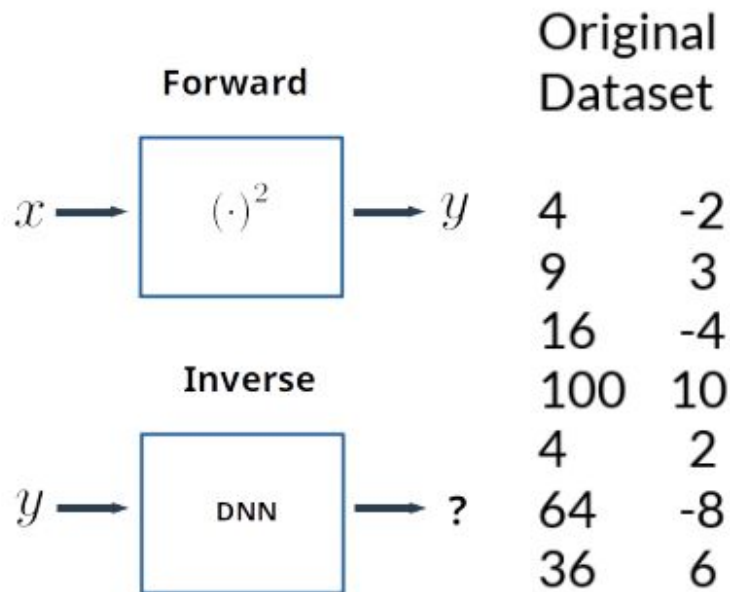
Implies: on dense training set, very close y 's can be mapped to very far away x 's different by signs



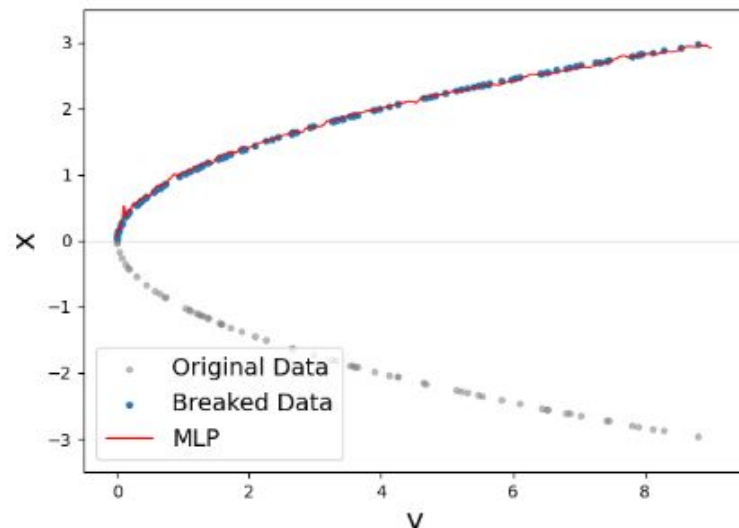
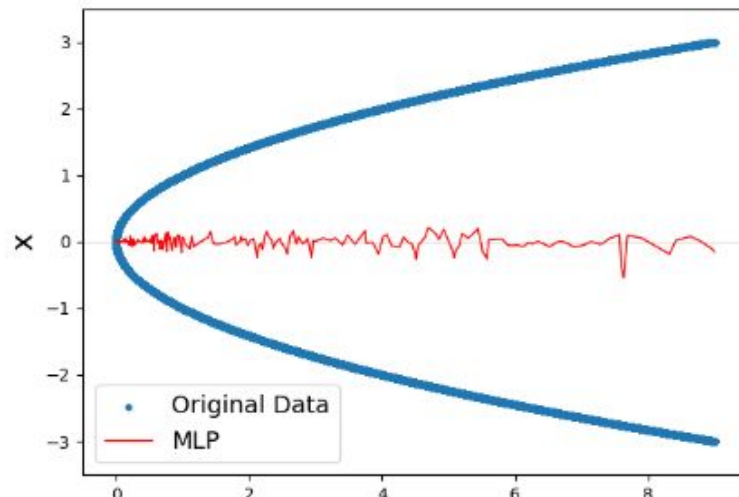
Highly oscillatory target function to learn by DNNs—difficult



Remedy: symmetry breaking



Fix all signs to be positive

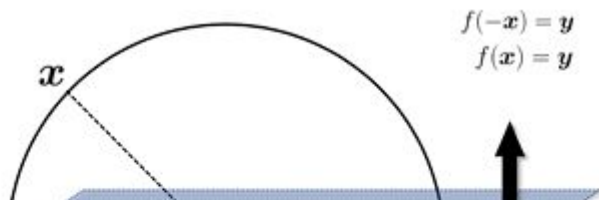


A slightly more complicated example

$$y = |\mathbf{A}\mathbf{x}|^2 \quad \mathbf{A} : \text{iid Gaussian} \quad (\text{Gaussian phase retrieval})$$

Forward symmetry: global sign

$$y = |\mathbf{A}\mathbf{x}|^2 = |\mathbf{A}(-\mathbf{x})|^2$$



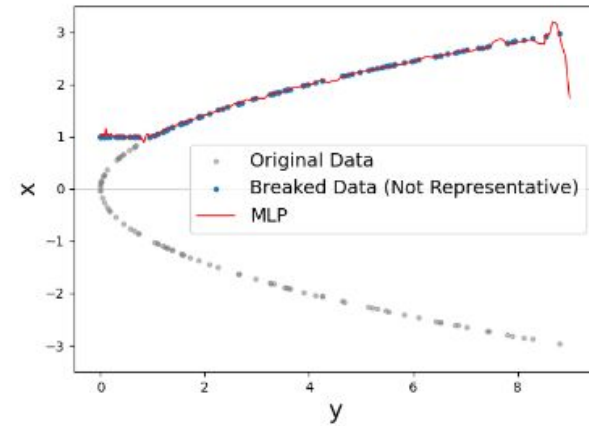
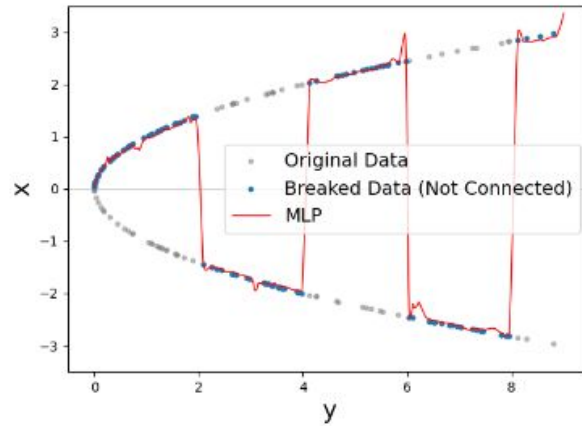
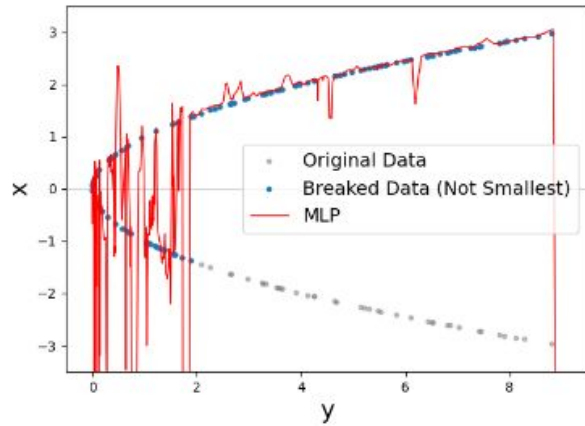
| Dim | Sample | After Symmetry Breaking | | Before Symmetry Breaking | |
|-----|--------|-------------------------|-------|--------------------------|-------|
| | | DNN | K-NN | DNN | K-NN |
| 5 | 2e4 | 4.08 | 11.82 | 85.37 | 68.26 |
| | 5e4 | 2.20 | 9.41 | 90.51 | 66.58 |
| | 1e5 | 1.30 | 7.98 | 96.66 | 66.18 |
| | 1e6 | 0.37 | 4.71 | 122.71 | 65.08 |

More is more

More is less

Symmetry-breaking principle

Symmetry breaking: a preprocessing step on the training set



Finding the smallest, connected, representative set

[Submitted on 18 Mar 2024]

What is Wrong with End-to-End Learning for Phase Retrieval?

Wenjie Zhang, Yuxiang Wan, Zhong Zhuang, Ju Sun

For nonlinear inverse problems that are prevalent in imaging science, symmetries in the forward model are common. When data-driven deep learning approaches are used to solve such problems, these intrinsic symmetries can cause substantial learning difficulties. In this paper, we explain how such difficulties arise and, more importantly, how to overcome them by preprocessing the training set before any learning, i.e., symmetry breaking. We take far-field phase retrieval (FFPR), which is central to many areas of scientific imaging, as an example and show that symmetric breaking can substantially improve data-driven learning. We also formulate the mathematical principle of symmetry breaking.

A version with careful mathematical analysis forthcoming ...

DL methods for SIPs: the middle way

Inverse problem: given $\mathbf{y} = f(\mathbf{x})$, recover \mathbf{x}

$$\min_{\mathbf{x}} \underbrace{\ell(\mathbf{y}, f(\mathbf{x}))}_{\text{data fitting}} + \lambda \underbrace{\mathbf{R}(\mathbf{x})}_{\text{regularizer}} \quad \text{RegFit}$$

Recipe: revamp numerical methods for RegFit with **pretrained/trainable DNNs**

DL methods for SIPs: the middle way

Algorithm unrolling

$$\min_{\mathbf{x}} \underbrace{\ell(\mathbf{y}, f(\mathbf{x}))}_{\text{data fitting}} + \lambda \underbrace{R(\mathbf{x})}_{\text{regularizer}}$$

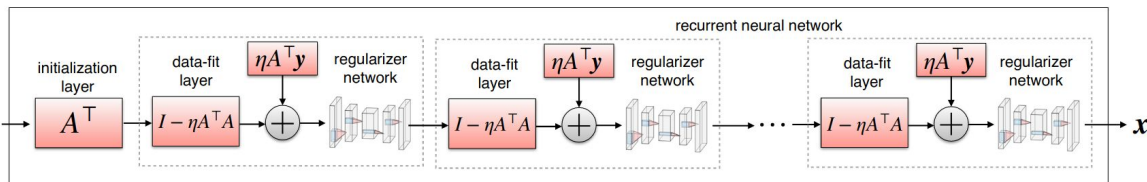
If R proximal friendly

$$\mathbf{x}^{k+1} = \mathcal{P}_R(\mathbf{x}^k - \eta \nabla^\top f(\mathbf{x}^k) \ell'(\mathbf{y}, f(\mathbf{x}^k)))$$

Idea: make \mathcal{P}_R trainable, using $\{(\mathbf{x}_i, \mathbf{y}_i)\}$

E.g.,

$$\ell(\mathbf{y}, f(\mathbf{x})) = \|\mathbf{y} - \mathbf{A} \mathbf{x}\|_2^2$$



DL methods for SIPs: the middle way

Using $\{\mathbf{x}_i\}$ only

$$\min_{\mathbf{x}} \underbrace{\ell(\mathbf{y}, f(\mathbf{x}))}_{\text{data fitting}} + \lambda \underbrace{R(\mathbf{x})}_{\text{regularizer}}$$

Plug-and-Play

$$\mathbf{x}^{k+1} = \mathcal{P}_R(\mathbf{x}^k - \eta \nabla^\top f(\mathbf{x}^k) \ell'(\mathbf{y}, f(\mathbf{x}^k)))$$

E.g. replace \mathcal{P}_R with pretrained denoiser

Deep generative models

Pretraining: $\mathbf{x}_i \approx G_\theta(\mathbf{z}_i) \quad \forall i$

Deployment: $\min_{\mathbf{z}} \ell(\mathbf{y}, f \circ G_\theta(\mathbf{z})) + \lambda R \circ G_\theta(\mathbf{z})$

DL methods for SIPs: a survey

Deep Learning Techniques for Inverse Problems in Imaging

Gregory Ongie^{*}, Ajil Jalal[†], Christopher A. Metzler[‡]
Richard G. Baraniuk[§], Alexandros G. Dimakis[¶], Rebecca Willett^{||}

April 2020

Abstract

Recent work in machine learning shows that deep neural networks can be used to solve a wide variety of inverse problems arising in computational imaging. We explore the central prevailing themes of this emerging area and present a taxonomy that can be used to categorize different problems and reconstruction methods. Our taxonomy is organized along two central axes: (1) whether or not a forward model is known and to what extent it is used in training and testing, and (2) whether or not the learning is supervised or unsupervised, i.e., whether or not the training relies on access to matched ground truth image and measurement pairs. We also discuss the tradeoffs associated with these different reconstruction approaches, caveats and common failure modes, plus open problems and avenues for future work.

Focuses on **linear**
inverse problems,
i.e., f linear

<https://arxiv.org/abs/2005.06001>

Limitations of middle ways:

- Representative data?
- Algorithm-sensitive
- Good initialization? (e.g.,
Manekar et al. **Deep Learning Initialized
Phase Retrieval.**
<https://sunju.org/pub/NIPS20-WS-DL4FPR.pdf>)

Other specialized surveys

Algorithm Unrolling: Interpretable, Efficient Deep Learning for Signal and Image Processing

Vishal Monga, *Senior Member, IEEE*, Yuelong Li, *Member, IEEE*, and Yonina C. Eldar, *Fellow, IEEE*

Focused on alg. unrolling

Untrained Neural Network Priors for Inverse Imaging Problems: A Survey

Deep Internal Learning:

Deep Learning from a Single Input

Understanding Untrained Deep Models for Inverse Problems: Algorithms and Theory

Tom Tirer *Member,*

**Focused on
single-instance methods**

Ismail Alkhouri, Evan Bell, Avrajit Ghosh, Shijun Liang, Rongrong Wang,

Theoretical Perspectives on Deep Learning Methods in Inverse Problems

Jonathan Scarlett, Reinhard Heckel, Miguel R. D. Rodrigues, Paul Hand, and Yonina C. Eldar

**Focused on theories for
linear IPs**

Story II: Don't be too Bayesian

$$d\mathbf{x} = -\beta_t/2 \cdot \mathbf{x}dt + \sqrt{\beta_t}d\mathbf{w},$$

Fixed forward diffusion process



Generative reverse denoising process

$$d\mathbf{x} = -\beta_t \left[\mathbf{x}/2 + \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt + \sqrt{\beta_t}d\bar{\mathbf{w}}.$$

$$\cong \boldsymbol{\varepsilon}_{\theta}^{(t)}(\mathbf{x})$$

Bayesian thinking

(Reverse SDE for DDPM) $d\mathbf{x} = -\beta_t [\mathbf{x}/2 + \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + \sqrt{\beta_t} d\bar{\mathbf{w}}$



Think of **conditional score function**

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}|\mathbf{y}) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) + \nabla_{\mathbf{x}} \log p_t(\mathbf{y}|\mathbf{x})$$



Conditional reverse SDE

$$d\mathbf{x} = [-\beta_t/2 \cdot \mathbf{x} - \beta_t(\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) + \nabla_{\mathbf{x}} \log p_t(\mathbf{y}|\mathbf{x}))] dt + \sqrt{\beta_t} d\bar{\mathbf{w}}$$

Interleaving methods

Algorithm 1 Template for interleaving methods

Input: # Diffusion steps T , measurement \mathbf{y}

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

2: **for** $i = T - 1$ to 0 **do**

3: $\hat{\mathbf{s}} \leftarrow \boldsymbol{\varepsilon}_{\theta}^{(i)}(\mathbf{x}_i)$

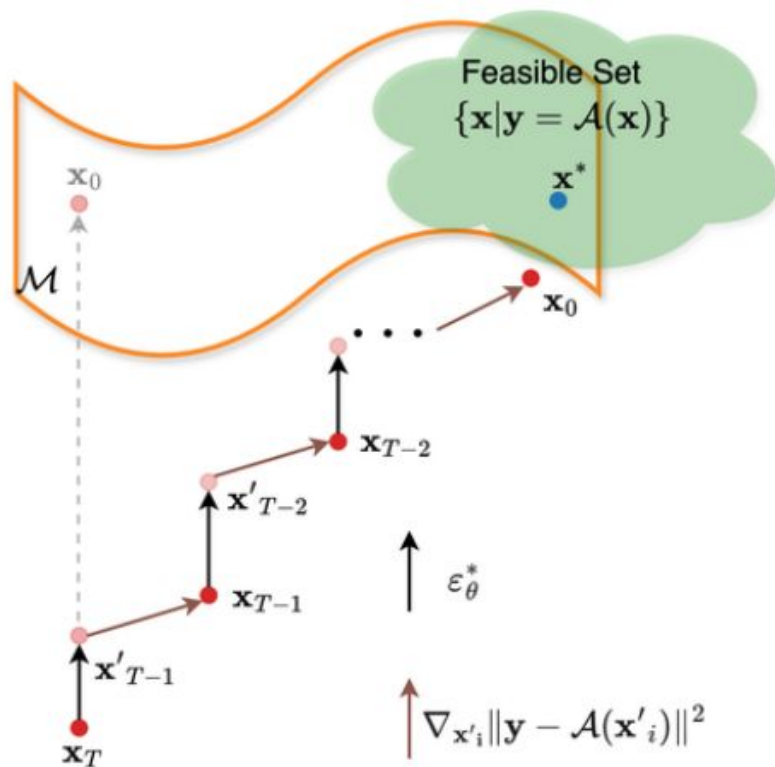
4: $\hat{\mathbf{x}}_0 \leftarrow \frac{1}{\sqrt{\bar{\alpha}_i}}(\mathbf{x}_i - \sqrt{1 - \bar{\alpha}_i}\hat{\mathbf{s}})$

5: $\mathbf{x}'_{i-1} \leftarrow$ DDIM reverse with $\hat{\mathbf{x}}_0$ and $\hat{\mathbf{s}}$

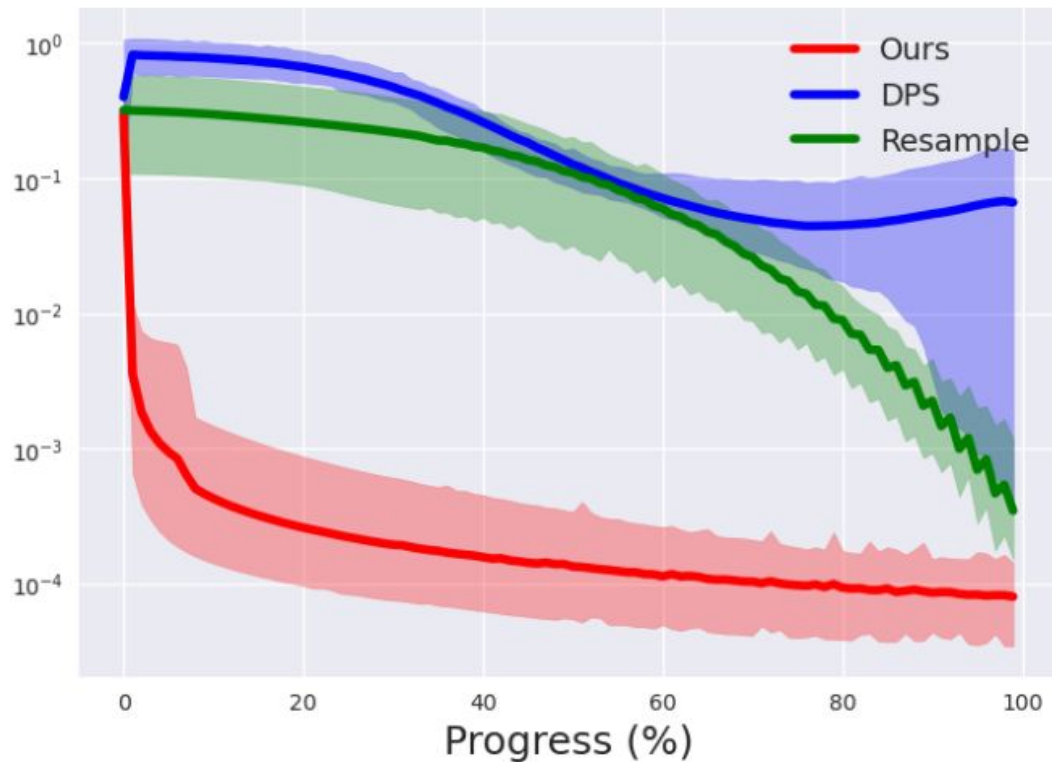
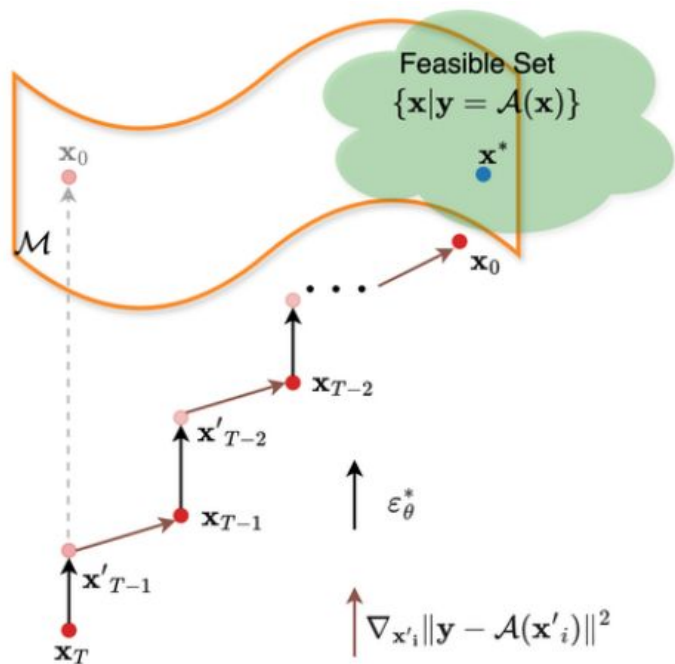
6: $\mathbf{x}_{i-1} \leftarrow$ (Approximately) Projection [39, 30, 33, 32, 40, 41, 34] or gradient update [20, 28, 19, 21, 29, 27, 26] with $\hat{\mathbf{x}}_0$ and \mathbf{x}'_{i-1} to get closer to $\{\mathbf{x} | \mathbf{y} = \mathcal{A}(\mathbf{x})\}$

7: **end for**

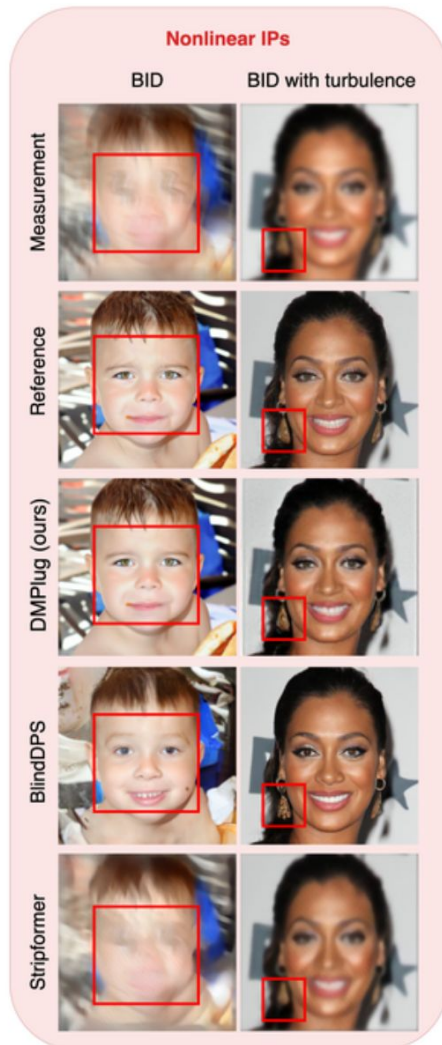
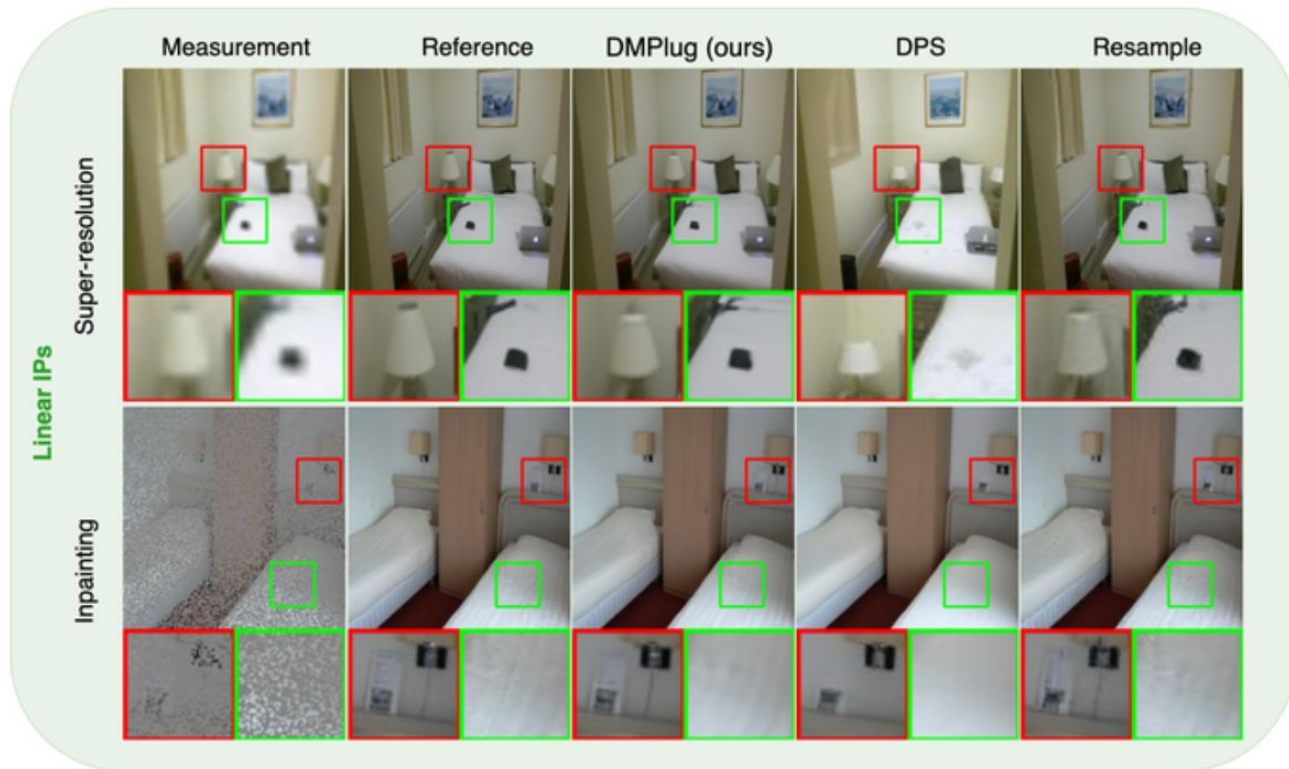
Output: Recovered object \mathbf{x}_0



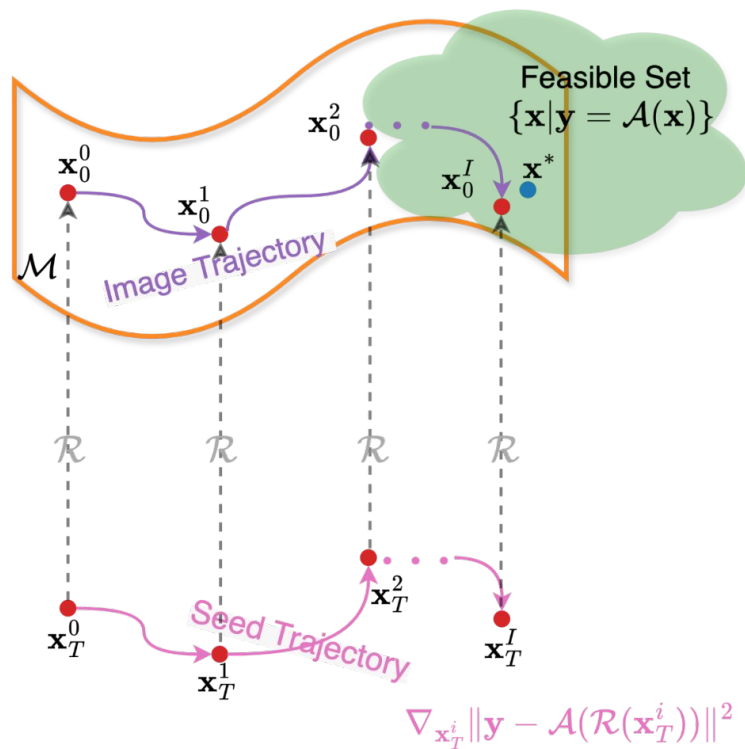
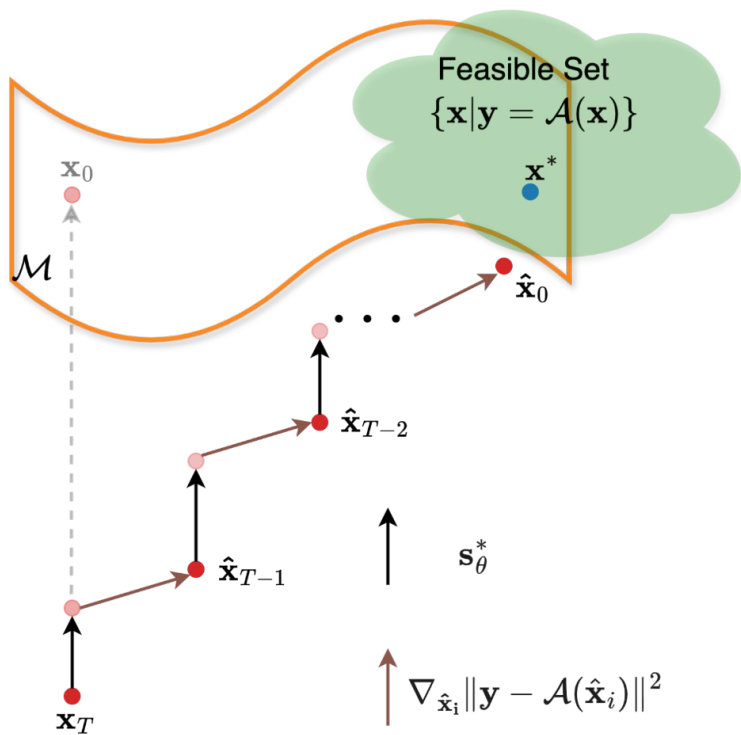
Feasibility crisis



Feasibility crisis



Explained in one picture (vs. our plugin idea)



On linear IPs

Table 1: (**Linear IPs**) **Super-resolution** and **inpainting** with additive Gaussian noise ($\sigma = 0.01$). (**Bold**: best, under: second best, **green**: performance increase, **red**: performance decrease)

| | Super-resolution ($4\times$) | | | | | | Inpainting (Random 70%) | | | | | |
|----------------------|---|--------------|--------------|---------------------------------------|--------------|--------------|---|--------------|--------------|---------------------------------------|--------------|--------------|
| | CelebA [65] (256×256) | | | FFHQ [66] (256×256) | | | CelebA [65] (256×256) | | | FFHQ [66] (256×256) | | |
| | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ |
| ADMM-PnP [68] | 0.217 | 26.99 | 0.808 | 0.229 | 26.25 | 0.794 | 0.091 | 31.94 | 0.923 | 0.104 | 30.64 | 0.901 |
| DMPS [29] | <u>0.070</u> | <u>28.89</u> | <u>0.848</u> | 0.076 | <u>28.03</u> | <u>0.843</u> | 0.297 | 24.52 | 0.693 | 0.326 | 23.31 | 0.664 |
| DDRM [32] | 0.226 | 26.34 | 0.754 | 0.282 | 25.11 | 0.731 | 0.185 | 26.10 | 0.712 | 0.201 | 25.44 | 0.722 |
| MCG [30] | 0.725 | 19.88 | 0.323 | 0.786 | 18.20 | 0.271 | 1.283 | 10.16 | 0.049 | 1.276 | 10.37 | 0.050 |
| ILVR [41] | 0.322 | 21.63 | 0.603 | 0.360 | 20.73 | 0.570 | 0.447 | 15.82 | 0.484 | 0.483 | 15.10 | 0.450 |
| DPS [19] | 0.087 | 28.32 | 0.823 | 0.098 | 27.44 | 0.814 | 0.043 | <u>32.24</u> | <u>0.924</u> | 0.046 | <u>30.95</u> | <u>0.913</u> |
| ReSample [20] | 0.080 | 28.29 | 0.819 | 0.108 | 25.22 | 0.773 | 0.039 | 30.12 | 0.904 | <u>0.044</u> | 27.91 | 0.884 |
| DMPlug (ours) | 0.067 | 31.25 | 0.878 | <u>0.079</u> | 30.25 | 0.871 | 0.039 | 34.03 | 0.936 | 0.038 | 33.01 | 0.931 |
| Ours vs. Best compe. | -0.003 | +2.36 | +0.030 | +0.003 | +2.22 | +0.028 | -0.000 | +1.79 | +0.012 | -0.006 | +2.06 | +0.018 |

On nonlinear IPs

Table 2: (**Nonlinear IP**) **Nonlinear deblurring** with additive Gaussian noise ($\sigma = 0.01$). (**Bold**: best, under: second best, **green**: performance increase, **red**: performance decrease)

| | CelebA 65 (256×256) | | | FFHQ 66 (256×256) | | | LSUN 67 (256×256) | | |
|------------------------|---------------------------------------|--------------|--------------|-------------------------------------|--------------|--------------|-------------------------------------|--------------|--------------|
| | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ |
| BKS-styleGAN 69 | 1.047 | 22.82 | 0.653 | 1.051 | 22.07 | 0.620 | 0.987 | 20.90 | 0.538 |
| BKS-generic 69 | 1.051 | 21.04 | 0.591 | 1.056 | 20.76 | 0.583 | 0.994 | 18.55 | 0.481 |
| MCG 30 | 0.705 | 13.18 | 0.135 | 0.675 | 13.71 | 0.167 | 0.698 | 14.28 | 0.188 |
| ILVR 41 | 0.335 | 21.08 | 0.586 | 0.374 | 20.40 | 0.556 | 0.482 | 18.76 | 0.444 |
| DPS 19 | 0.149 | 24.57 | 0.723 | 0.130 | 25.00 | 0.759 | 0.244 | 23.46 | 0.684 |
| ReSample 20 | <u>0.104</u> | <u>28.52</u> | <u>0.839</u> | <u>0.104</u> | <u>27.02</u> | <u>0.834</u> | <u>0.143</u> | <u>26.03</u> | <u>0.803</u> |
| DMPlug (ours) | 0.073 | 31.61 | 0.882 | 0.057 | 32.83 | 0.907 | 0.083 | 30.74 | 0.882 |
| Ours vs. Best compe. | -0.031 | +3.09 | +0.043 | -0.047 | +5.79 | +0.073 | -0.060 | +4.71 | +0.079 |

The paper (NeurIPS'24)

[Submitted on 27 May 2024]

DMPlug: A Plug-in Method for Solving Inverse Problems with Diffusion Models

Hengkang Wang, Xu Zhang, Taihui Li, Yuxiang Wan, Tiancong Chen, Ju Sun

Pretrained diffusion models (DMs) have recently been popularly used in solving inverse problems (IPs). The existing methods mostly interleave iterative steps in the reverse diffusion process and iterative steps to bring the iterates closer to satisfying the measurement constraint. However, such interleaving methods struggle to produce final results that look like natural objects of interest (i.e., manifold feasibility) and fit the measurement (i.e., measurement feasibility), especially for nonlinear IPs. Moreover, their capabilities to deal with noisy IPs with unknown types and levels of measurement noise are unknown. In this paper, we advocate viewing the reverse process in DMs as a function and propose a novel plug-in method for solving IPs using pretrained DMs, dubbed DMPlug. DMPlug addresses the issues of manifold feasibility and measurement feasibility in a principled manner, and also shows great potential for being robust to unknown types and levels of noise. Through extensive experiments across various IP tasks, including two linear and three nonlinear IPs, we demonstrate that DMPlug consistently outperforms state-of-the-art methods, often by large margins especially for nonlinear IPs. The code is available at [this https URL](https://github.com/HengkangWang/DMPlug).

<https://arxiv.org/abs/2405.16749>

DL methods for SIPS: the **economic/surprising** way

Deep image prior (DIP) $\mathbf{x} \approx G_{\theta}(\mathbf{z})$ G_{θ} (and \mathbf{z}) trainable

$$\min_{\mathbf{x}} \underbrace{\ell(\mathbf{y}, f(\mathbf{x}))}_{\text{data fitting}} + \lambda \underbrace{R(\mathbf{x})}_{\text{regularizer}}$$



$$\min_{\theta} \ell(\mathbf{y}, f \circ G_{\theta}(\mathbf{z})) + \lambda R \circ G_{\theta}(\mathbf{z})$$

No extra training data!

Ulyanov et al. **Deep image prior**. IJCV'20. <https://arxiv.org/abs/1711.10925>

Contrasting: Deep generative models

Pretraining: $\mathbf{x}_i \approx G_{\theta}(\mathbf{z}_i) \quad \forall i$

Deployment: $\min_{\mathbf{z}} \ell(\mathbf{y}, f \circ G_{\theta}(\mathbf{z})) + \lambda R \circ G_{\theta}(\mathbf{z})$

Deep image prior (DIP)

DIP's cousin(s)

$$\mathbf{x} \approx G_{\theta}(\mathbf{z}) \quad G_{\theta} \text{ (and } \mathbf{z}) \text{ trainable}$$

Idea: (visual) objects as continuous functions

Neural implicit representation (NIR)

$$\mathbf{x} \approx \mathcal{D} \circ \bar{\mathbf{x}} \quad \mathcal{D} : \text{discretization} \quad \bar{\mathbf{x}} : \text{continuous function}$$

Physics-informed neural networks (PINN)

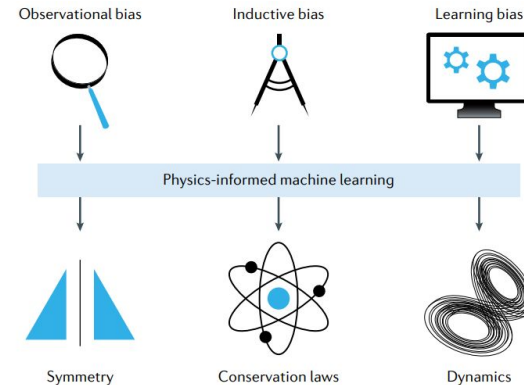
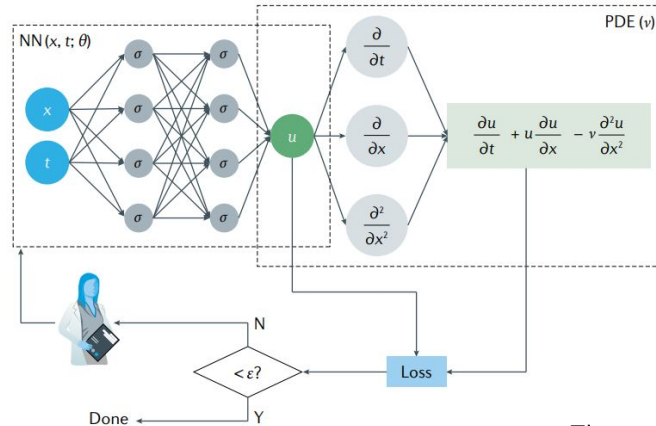


Figure credit: <https://www.nature.com/articles/s42254-021-00314-5>

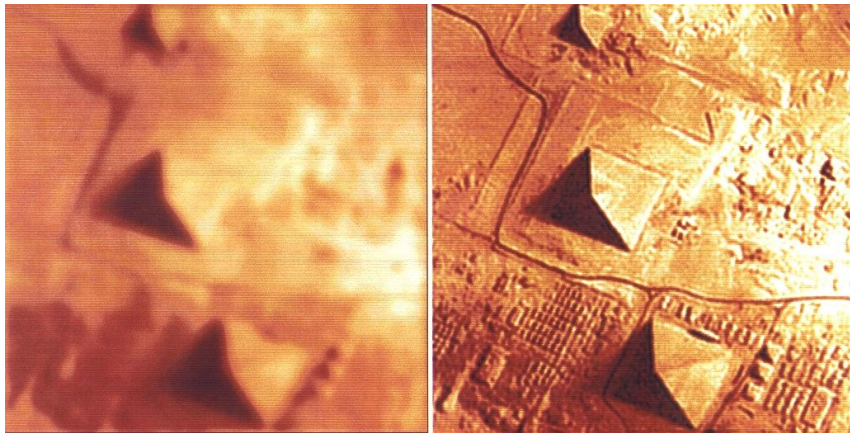
Story III: We benefit from DL even with a single data point

Blind image deblurring (BID)

$$\underbrace{\mathbf{y}}_{\text{blurry and noisy image}} = \underbrace{\mathbf{k}}_{\text{blur kernel}} * \underbrace{\mathbf{x}}_{\text{clean image}} + \underbrace{\mathbf{n}}_{\text{noise}}$$

Given \mathbf{y} ,
recover \mathbf{x} (and/or \mathbf{k})

Also **Blind Deconvolution**



Landmark surveys

- 1996: Kundur and Hatzinakos. **Blind image deconvolution**. <https://doi.org/10.1109/79.489268>
- 2011: Levin et al. **Understanding blind deconvolution algorithms**. <https://doi.org/10.1109/TPAMI.2011.148>
- 2012: Kohler et al. **Recording and playback of camera shake: Benchmarking blind deconvolution with a real-world database**. https://doi.org/10.1007/978-3-642-33786-4_3
- 2016: Lai et al. **A comparative study for single image blind deblurring**. <https://doi.org/10.1109/CVPR.2016.188>
- 2021: Koh et al. **Single image deblurring with neural networks: A comparative survey** <https://doi.org/10.1016/j.cviu.2020.103134>
- 2022: Zhang et al. **Deep image blurring: A survey** <https://doi.org/10.1007/s11263-022-01633-5>

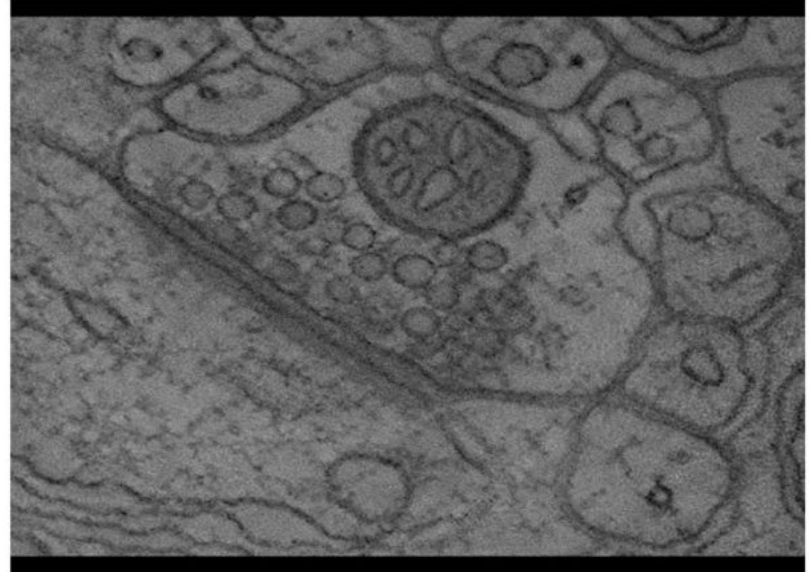
See also: **Awesome Deblurring**

<https://github.com/subeeshvasu/Awesome-Deblurring>

Key challenge of data-driven approach:

obtaining sufficiently expressive data (Koh et al'21. Zhang et al'22)

Untouched practical questions



Key question addressed in this paper: How do we solve blind image deblurring without knowing: (1) the size of the blur kernel, (2) the type and level of noise, and (3) whether it is blur / noise only or both ?

Double DIPs

$$\underbrace{\mathbf{y}}_{\text{blurry and noisy image}} = \underbrace{\mathbf{k}}_{\text{blur kernel}} * \underbrace{\mathbf{x}}_{\text{clean image}} + \underbrace{\mathbf{n}}_{\text{noise}}$$

$$\min_{\mathbf{k}, \mathbf{x}} \underbrace{\ell(\mathbf{y}, \mathbf{k} * \mathbf{x})}_{\text{data fitting}} + \lambda_k \underbrace{R_k(\mathbf{k})}_{\text{regularizing } \mathbf{k}} + \lambda_x \underbrace{R_x(\mathbf{x})}_{\text{regularizing } \mathbf{x}}$$

Idea: parameterize both \mathbf{k} and \mathbf{x} as DIPs

- CNN + CNN (Wang et al'19, <https://doi.ieeecomputersociety.org/10.1109/ICCVW.2019.00127>; Tran et al'21, <https://arxiv.org/abs/2104.00317>)
- MLP + CNN (SelfDeblur; Ren et al'20, <https://arxiv.org/abs/1908.02197>)

Still problematic with

- 1) kernel size over-specification
- 2) substantial noise

A glance of our modifications

Over-specify \mathbf{k}

Over-specify \mathbf{x}

\mathbf{k} ~half of the image sizes

Handle bounded shift



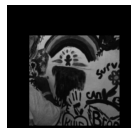
Ground Truth



Over-specified \mathbf{x}



Exact-specified \mathbf{x}



x



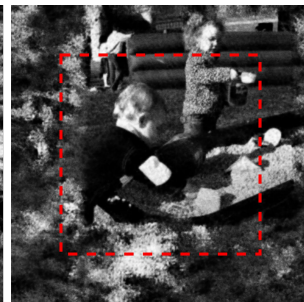
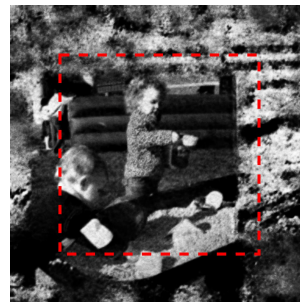
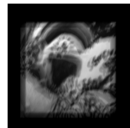
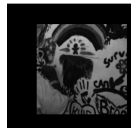
$*$

k

$=$



y



$$\min_{\theta_k, \theta_x} \|y - G_{\theta_k}(z_k) * G_{\theta_x}(z_x)\|_2^2 + \lambda \frac{\|\nabla G_{\theta_x}(z_x)\|_1}{\|\nabla G_{\theta_x}(z_x)\|_2}$$

ℓ_1/ℓ_2 vs ℓ_1

Table 1: ℓ_1/ℓ_2 vs TV for noise: mean and (std).

| | Low Level | | High Level | |
|-----------------|--------------|----------------|--------------|-----------------|
| | PSNR | λ | PSNR | λ |
| $\frac{L1}{L2}$ | 32.64 (0.69) | 0.0001 (0.018) | 27.74 (0.23) | 0.0002 (0.0019) |
| TV | 31.12 (0.52) | 0.002 (0.07) | 24.34 (0.78) | 0.02 (0.10) |

SelfDeblur vs our method



Clean



Blurry and noisy



SelfDeblur



Ours



Clean



Blurry and noisy

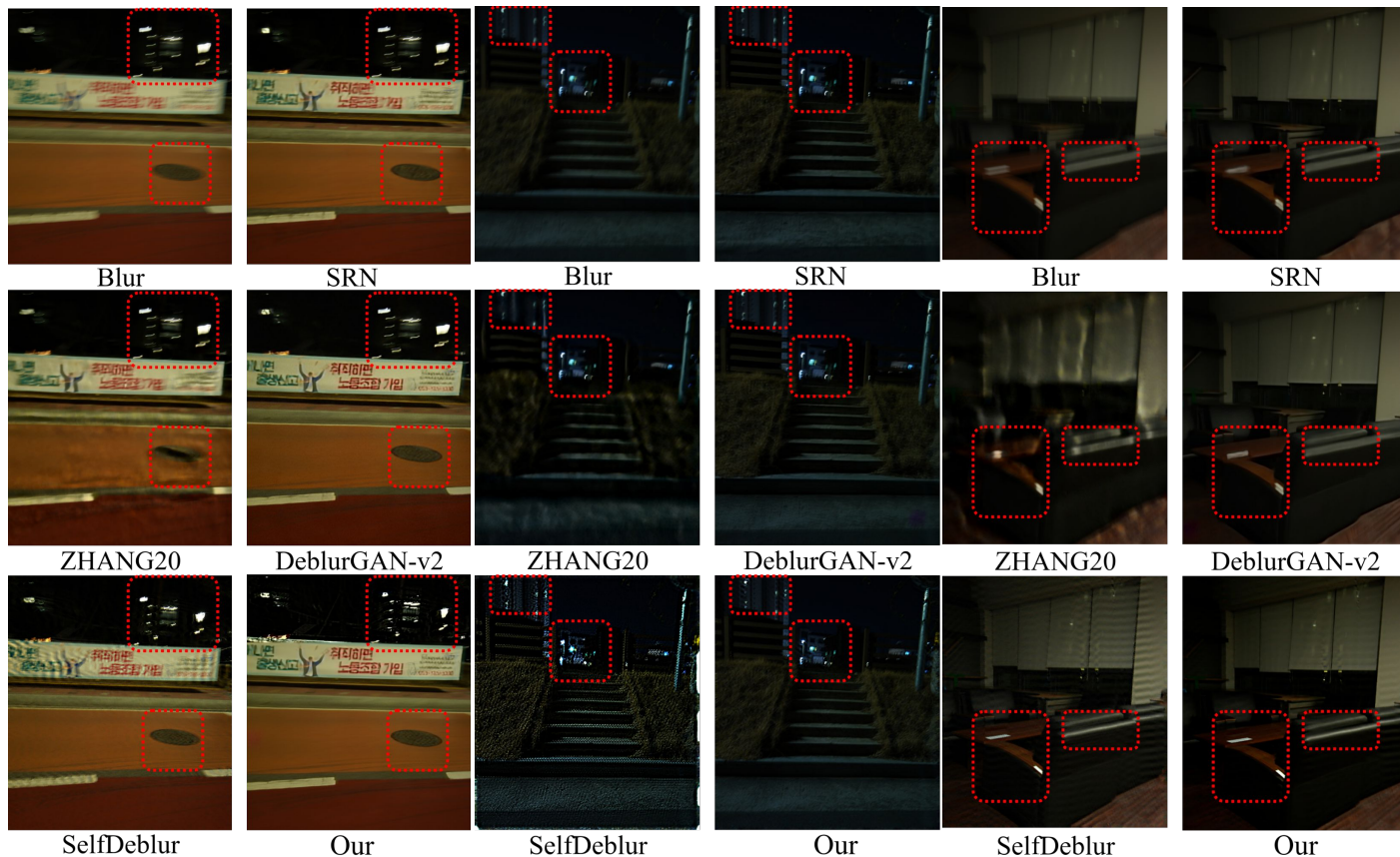


SelfDeblur



Ours

Real world results

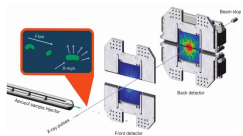


Difficult cases

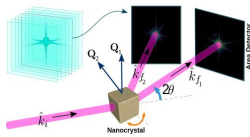
- 1) High depth contrast
- 2) High brightness contrast

**Outperform SOTA
data-driven methods!**

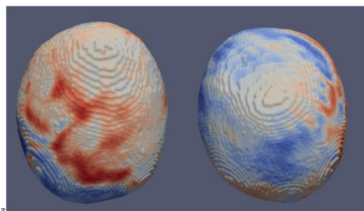
Breakthroughs in imaging



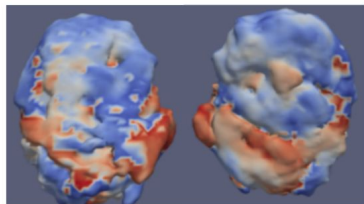
Coherent Diffraction Imaging



Bragg Coherent Diffraction Imaging



Our



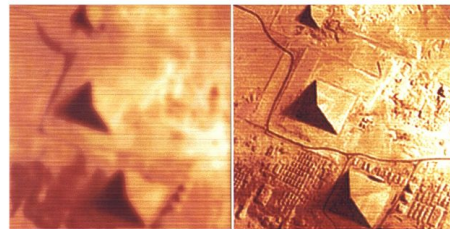
HIO+ER with Shrinkwrap

$$\underbrace{\mathbf{y}}_{\text{blurry and noisy image}} = \underbrace{\mathbf{k}}_{\text{blur kernel}} * \underbrace{\mathbf{x}}_{\text{clean image}} + \underbrace{\mathbf{n}}_{\text{noise}}$$

Mostly due to optical deficiencies (e.g., defocus) and motions

Given \mathbf{y} ,
recover \mathbf{x} (and/or \mathbf{k})

Also **Blind Deconvolution**



First PR method that won in a double-blind test, and systematic evaluation, beating a 40-years old legacy

Practical Phase Retrieval Using Double Deep Image Priors

Zhong Zhuang, David Yang, Felix Hofmann, David Barmherzig, Ju Sun

First BID method that works with unknown kernel size AND substantial noise

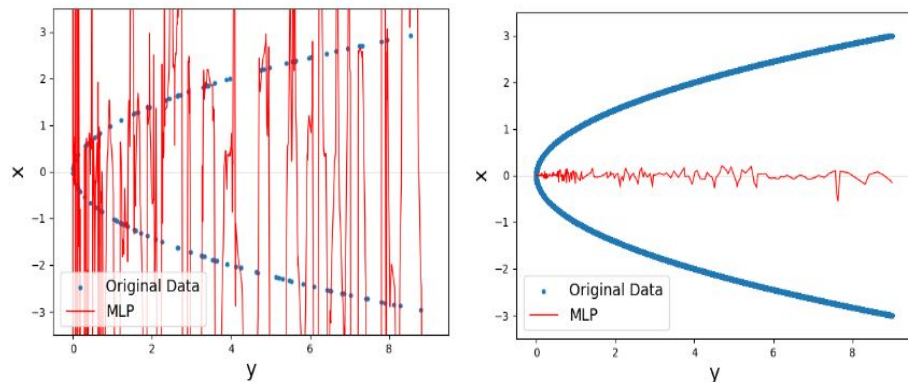
Blind Image Deblurring with Unknown Kernel Size and Substantial Noise

Zhong Zhuang, Taihui Li, Hengkang Wang, Ju Sun

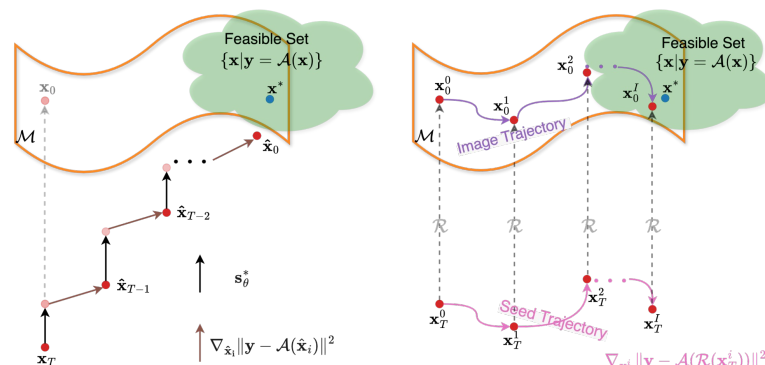
Related papers

- Li et al. **Self-Validation: Early Stopping for Single-Instance Deep Generative Priors** (BMVC'21) <https://arxiv.org/abs/2110.12271>
- Wang et al. **Early Stopping for Deep Image Prior** <https://arxiv.org/abs/2112.06074> (TMLR'23)
- Zhuang et al. **Blind Image Deblurring with Unknown Kernel Size and Substantial Noise.** <https://arxiv.org/abs/2208.09483> (IJCV'24)
- Zhuang et al. **Practical Phase Retrieval Using Double Deep Image Priors.** <https://arxiv.org/abs/2211.00799> (Electronic Imaging'24)
- Li et al. **Deep Random Projector: Toward Efficient Deep Image Prior.** (CVPR'23)

Data-driven methods for SIPs



Story I: More could be less



Story II: Don't be too Bayesian

Single-instance methods for SIPs

Story III: Benefit from DL with a single data point

Deep image prior (DIP) $\mathbf{x} \approx G_{\theta}(\mathbf{z})$ G_{θ} (and \mathbf{z}) trainable

$$\min_{\mathbf{x}} \underbrace{\ell(\mathbf{y}, f(\mathbf{x}))}_{\text{data fitting}} + \lambda \underbrace{R(\mathbf{x})}_{\text{regularizer}}$$

$$\min_{\theta} \ell(\mathbf{y}, f \circ G_{\theta}(\mathbf{z})) + \lambda R \circ G_{\theta}(\mathbf{z})$$

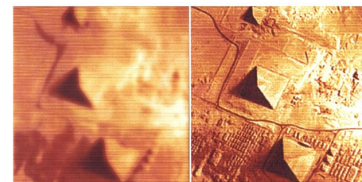
No extra training data!

$$\underbrace{\mathbf{y}}_{\text{blurry and noisy image}} = \underbrace{\mathbf{k}}_{\text{blur kernel}} * \underbrace{\mathbf{x}}_{\text{clean image}} + \underbrace{\mathbf{n}}_{\text{noise}}$$

Mostly due to optical deficiencies (e.g., defocus) and motions

Given \mathbf{y} ,
recover \mathbf{x} (and/or \mathbf{k})

Also **Blind Deconvolution**



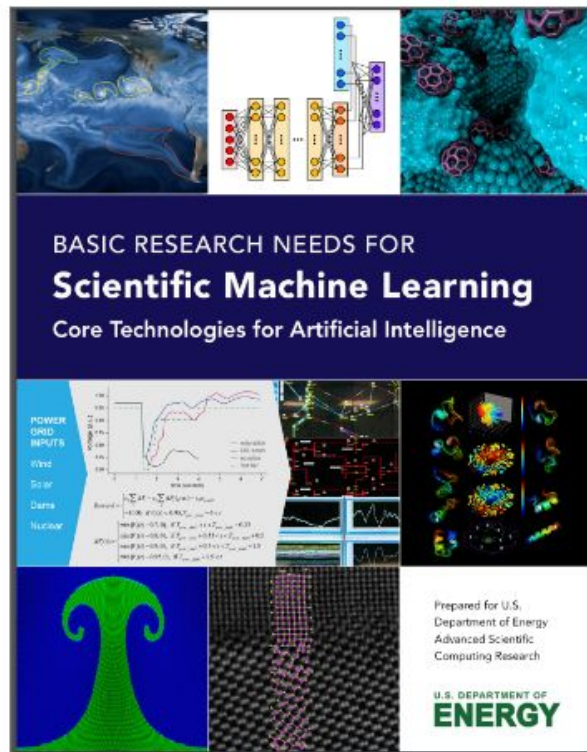
Principled data-knowledge tradeoff

Thrust B: How Should Domain Knowledge Be Incorporated into Supervised Machine Learning?

The central question for this thrust is “which knowledge should be leveraged in SciML, and how should this knowledge be included?” Any answers will naturally depend on the SciML task and computational budgets, thus mirroring standard considerations in traditional scientific computing.

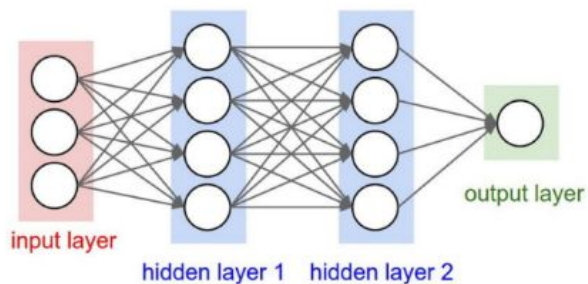
Hard Constraints. One research avenue involves incorporation of domain knowledge through imposition of constraints that cannot be violated. These hard constraints could be enforced during training, replacing what typically is an unconstrained optimization problem with a constrained one. In general, such constraints could involve simulations or highly nonlinear functions of the training parameters. Therefore, there is a need to identify particular cases when constraint qualification conditions can be ensured as these conditions are necessary regularity conditions for constrained optimization [57–59]. Although incorporating constraints during training generally makes maximal use of training data, there may be additional opportunities to employ constraints at the time of prediction (e.g., by projecting predictions onto the region induced by the constraints).

Soft Constraints. A similar avenue for incorporating domain knowledge involves modifying the objective function (soft constraints) used in training. It is understood that ML loss function selection should be guided by the task and data. Therefore, opportunities exist for developing loss functions that incorporate domain knowledge and analyzing the resulting impact on solvability



When DL meets constraints

Artificial neural networks



Unconstrained optimization

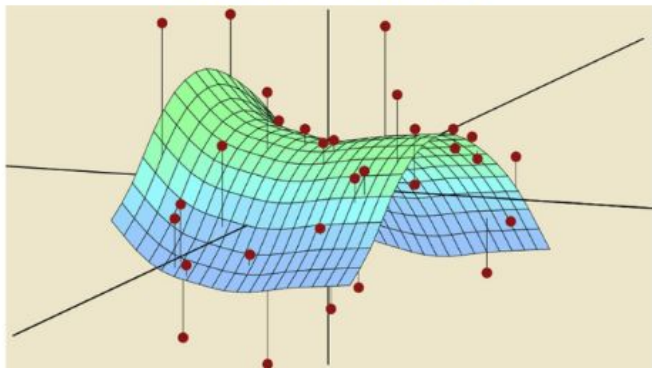
$$\min_{w'_i s, b'_i s} \frac{1}{n} \sum_{i=1}^n \ell[y_i, \{\text{NN}(w_1, \dots, w_k, b_1, \dots, b_k)\}(x_i)]$$
$$\min_x f(x)$$

“Solved”

Constrained optimization

$$\min_x f(x) \quad \text{s. t. } g(x) \leq 0$$

largely “unsolved”



used to approximate nonlinear functions

GAPS

$$\min_x f(x) \quad \text{s.t. } g(x) \leq 0$$

largely “unsolved”



An imaginary chat between a PhD student working in deep learning (**DLP**) and a PhD student working in optimization (**OP**)

DLP: Man, I've solved a constrained DL problem recently

OP: Oh, that's a hard problem

DLP: Really? I actually did it

OP: How?

DLP: My problem is $\min_x f(x), \text{s.t. } g(x) \leq 0$. I put $g(x)$ as a penalty and then call ADAM

OP: Are you sure it works?

DLP: Yes, the performance is improved and my paper is published at ICML

OP: Why don't you try augmented Lagrangian methods?

DLP: No implementation in Pytorch. Is it possible we work out some theory about my method?

OP: I think it's hard. It's not convex

DL with nontrivial constraints: many pitfalls

- **Robustness evaluation**
- Imbalanced learning
- Topology optimization

Deep Learning with Nontrivial Constraints: Methods and Applications

Chuan He¹, Ryan Devera¹, Wenjie Zhang¹, Ying Cui², Zhaosong Lu³ and Ju Sun¹

¹Computer Science and Engineering, University of Minnesota

²Industrial Engineering and Operations Research, University of California, Berkeley

³Industrial and Systems Engineering, University of Minnesota

{he000233, dever120, zhan7867}@umn.edu, yingcui@berkeley.edu, {zhaosong, jusun}@umn.edu

Robustness evaluation: penalty methods for complicated d (perceptual attack)

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ \text{s. t. } & d(\mathbf{x}, \mathbf{x}') \leq \epsilon, \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$

$$d(\mathbf{x}, \mathbf{x}') \doteq \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2$$

where $\phi(\mathbf{x}) \doteq [\hat{g}_1(\mathbf{x}), \dots, \hat{g}_L(\mathbf{x})]$

perceptual distance

Projection onto the constraint is complicated

Penalty methods

$$\max_{\tilde{\mathbf{x}}} \mathcal{L}(f(\tilde{\mathbf{x}}), y) - \lambda \max(0, \|\phi(\tilde{\mathbf{x}}) - \phi(\mathbf{x})\|_2 - \epsilon)$$

Solve it for each fixed λ and then increase λ

Algorithm 2 Lagrangian Perceptual Attack (LPA)

```
1: procedure LPA(classifier network  $f(\cdot)$ , LPIPS distance  $d(\cdot, \cdot)$ , input  $\mathbf{x}$ , label  $y$ , bound  $\epsilon$ )
2:    $\lambda \leftarrow 0.01$ 
3:    $\tilde{\mathbf{x}} \leftarrow \mathbf{x} + 0.01 * \mathcal{N}(0, 1)$   $\triangleright$  initialize perturbations with random Gaussian noise
4:   for  $i$  in  $1, \dots, S$  do  $\triangleright$  we use  $S = 5$  iterations to search for the best value of  $\lambda$ 
5:     for  $t$  in  $1, \dots, T$  do  $\triangleright T$  is the number of steps
6:        $\Delta \leftarrow \nabla_{\tilde{\mathbf{x}}} [\mathcal{L}(f(\tilde{\mathbf{x}}), y) - \lambda \max(0, d(\tilde{\mathbf{x}}, \mathbf{x}) - \epsilon)]$   $\triangleright$  take the gradient of (5)
7:        $\hat{\Delta} = \Delta / \|\Delta\|_2$   $\triangleright$  normalize the gradient
8:        $\eta = \epsilon * (0.1)^{t/T}$   $\triangleright$  the step size  $\eta$  decays exponentially
9:        $m \leftarrow d(\tilde{\mathbf{x}}, \tilde{\mathbf{x}} + h\hat{\Delta})/h$   $\triangleright m \approx$  derivative of  $d(\tilde{\mathbf{x}}, \cdot)$  in the direction of  $\hat{\Delta}$ ;  $h = 0.1$ 
10:       $\tilde{\mathbf{x}} \leftarrow \tilde{\mathbf{x}} + (\eta/m)\hat{\Delta}$   $\triangleright$  take a step of size  $\eta$  in LPIPS distance
11:    end for
12:    if  $d(\tilde{\mathbf{x}}, \mathbf{x}) > \epsilon$  then
13:       $\lambda \leftarrow 10\lambda$   $\triangleright$  increase  $\lambda$  if the attack goes outside the bound
14:    end if
15:  end for
16:   $\tilde{\mathbf{x}} \leftarrow \text{PROJECT}(d, \tilde{\mathbf{x}}, \mathbf{x}, \epsilon)$ 
17:  return  $\tilde{\mathbf{x}}$ 
18: end procedure
```

Problem with penalty methods

| Method | cross-entropy loss | | margin loss | |
|-------------|--------------------|------------------|-------------|------------------|
| | Viol. (%) ↓ | Att. Succ. (%) ↑ | Viol. (%) ↓ | Att. Succ. (%) ↑ |
| Fast-LPA | 73.8 | 3.54 | 41.6 | 56.8 |
| LPA | 0.00 | 80.5 | 0.00 | 97.0 |
| PPGD | 5.44 | 25.5 | 0.00 | 38.5 |
| PWCF (ours) | 0.62 | 93.6 | 0.00 | 100 |

LPA, Fast-LPA: penalty methods

PPGD: Projected gradient descent

Penalty methods tend to encounter

large constraint violation (i.e., infeasible solution, known in optimization theory) or **suboptimal solution**

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ \text{s. t. } & d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n \end{aligned}$$

$$\begin{aligned} & d(\mathbf{x}, \mathbf{x}') \doteq \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2 \\ \text{where } & \phi(\mathbf{x}) \doteq [\hat{g}_1(\mathbf{x}), \dots, \hat{g}_L(\mathbf{x})] \end{aligned}$$

PWCF, an optimizer with a principled stopping criterion on **stationarity** & **feasibility**

Key algorithm



<http://www.timmitchell.com/software/GRANSO/>

Nonconvex, nonsmooth, constrained

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad \text{s.t.} \quad c_i(\mathbf{x}) \leq 0, \quad \forall i \in \mathcal{I}; \quad c_i(\mathbf{x}) = 0, \quad \forall i \in \mathcal{E}.$$

Penalty sequential quadratic programming (P-SQP)

$$\begin{aligned} \min_{d \in \mathbb{R}^n, s \in \mathbb{R}^p} \quad & \mu(f(x_k) + \nabla f(x_k)^\top d) + e^\top s + \frac{1}{2} d^\top H_k d \\ \text{s.t.} \quad & c(x_k) + \nabla c(x_k)^\top d \leq s, \quad s \geq 0, \end{aligned}$$

Ref: **Curtis, Frank E., Tim Mitchell, and Michael L. Overton.** "A BFGS-SQP method for nonsmooth, nonconvex, constrained optimization and its evaluation using relative minimization profiles." *Optimization Methods and Software* 32.1 (2017): 148-181.

Algorithm highlights

Steering strategy for the penalty parameter

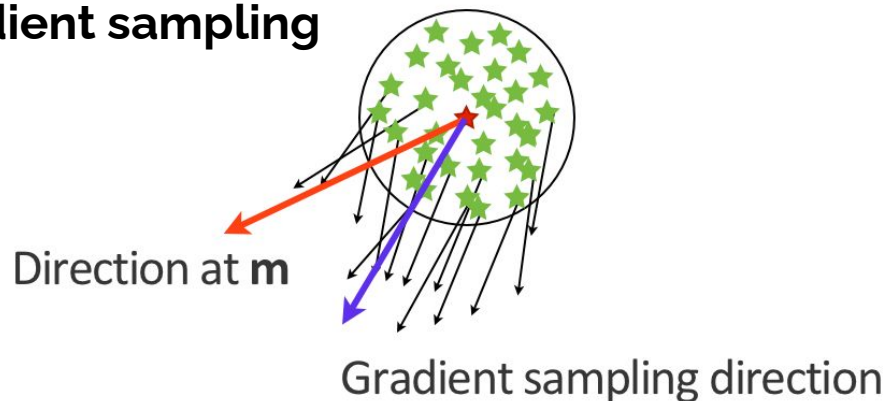
If feasibility improvement is insufficient : $l_\delta(d_k; x_k) < c_v v(x_k)$

Stationarity based on (approximate) gradient sampling

$$G_k := [\nabla f(x^k) \quad \nabla f(x^{k,1}) \quad \dots \quad \nabla f(x^{k,m})]$$

$$\min_{\lambda \in \mathbb{R}^{m+1}} \quad \frac{1}{2} \|G_k \lambda\|_2^2$$

$$\text{s.t. } \mathbf{1}^T \lambda = 1, \quad \lambda \geq 0$$



Key take-away



- Principled stopping criterion and line search, to obtain a **solution with certificate** (stationarity & feasibility check)
- Quasi-newton style method for fast convergence, i.e., **reasonable speed and high-precision solution**

Limitations of GRANSO



```
% Gradient of inner product with respect to A
f_grad      = imag((conj(Bty)*Cx.)/(y'*x));
f_grad      = f_grad(:);
```

```
% Gradient of inner product with respect to A
ci_grad     = real((conj(Bty)*Cx.)/(y'*x));
ci_grad     = ci_grad(:);
```

analytical gradients required

```
p          = size(B,2);
m          = size(C,1);
X          = reshape(x,p,m);
```

vector variables only

Lack of Auto-Differentiation

Lack of GPU Support

No native support of tensor variables

⇒ impossible to do deep learning with GRANSO

GRANSO meets PyTorch



NCVX PyGRANSO
Documentation

Search the docs ...

Introduction

Installation

Settings

Examples



Home



NCVX Package

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \text{ s.t. } c_i(\mathbf{x}) \leq 0, \forall i \in \mathcal{I}; c_i(\mathbf{x}) = 0, \forall i \in \mathcal{E}$$

First general-purpose solver for constrained DL problems

NCVX: A General-Purpose Optimization Solver for Constrained Machine and Deep Learning

Buyun Liang, Tim Mitchell, Ju Sun

Example 1: Support Vector Machine (SVM)

Soft-margin SVM

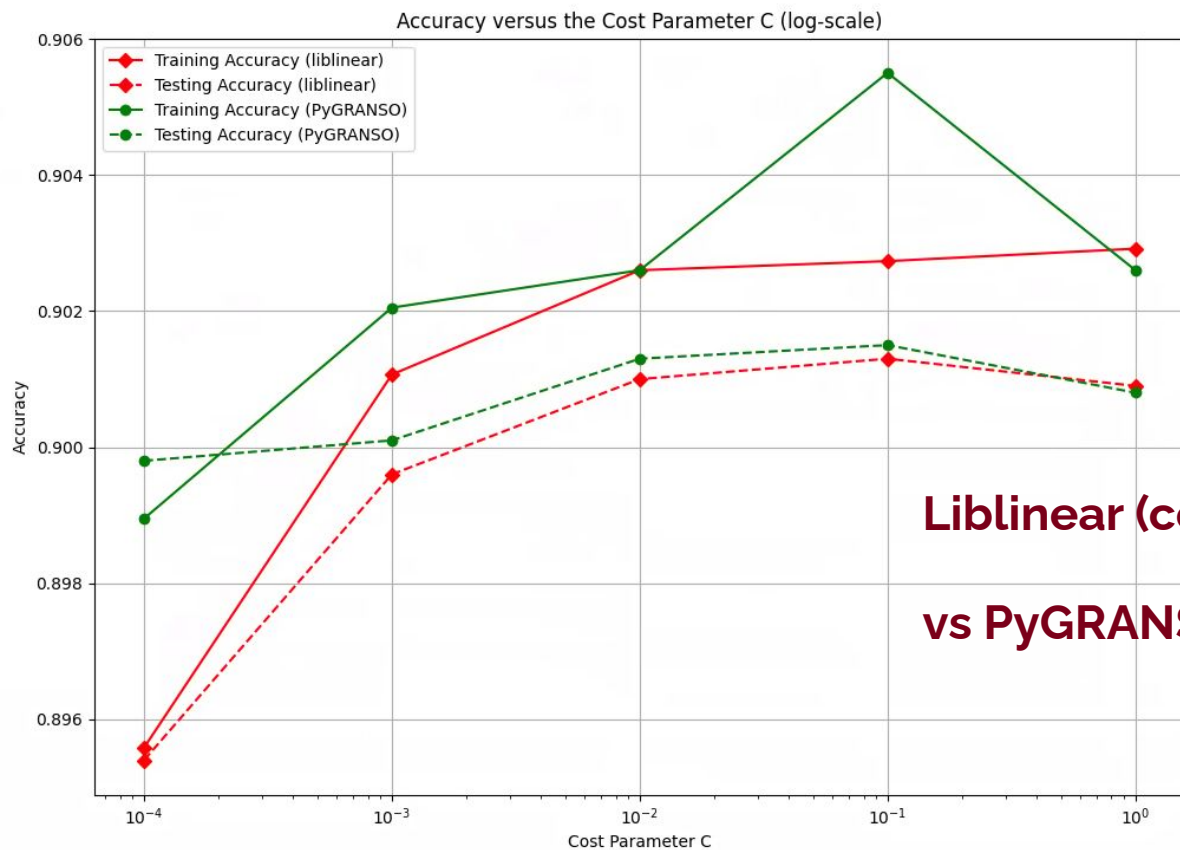
$$\min_{\mathbf{w}, b, \zeta} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \zeta_i$$

s.t. $y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \zeta_i, \zeta_i \geq 0 \quad \forall i = 1, \dots, n$

```
def comb_fn(X_struct):
    # obtain optimization variables
    w = X_struct.w
    b = X_struct.b
    zeta = X_struct.zeta
    # objective function
    f = 0.5*w.T@w + C*torch.sum(zeta)
    # inequality constraints
    ci = pygransoStruct()
    ci.c1 = 1 - zeta - y*(x@w+b)
    ci.c2 = -zeta
    # equality constraint
    ce = None
    return [f,ci,ce]

# specify optimization variables
var_in = {"w": [d,1], "b": [1,1], "zeta": [n,1]}
# pygranso main algorithm
soln = pygranso(var_in,comb_fn)
```

Binary classification (odd vs even digits) on MNIST dataset



**Liblinear (coordinate descent)
vs PyGRANSO**

Example 2: Robustness—min formulation

$$\begin{aligned} \min_{\mathbf{x}'} \quad & d(\mathbf{x}, \mathbf{x}') \\ \text{s. t.} \quad & \max_{\ell \neq c} f_{\boldsymbol{\theta}}^{\ell}(\mathbf{x}') \geq f_{\boldsymbol{\theta}}^c(\mathbf{x}') \\ & \mathbf{x}' \in [0, 1]^n \end{aligned}$$

```
def comb_fn(X_struct):
    # obtain optimization variables
    x_prime = X_struct.x_prime
    # objective function
    f = d(x, x_prime)
    # inequality constraints
    ci = pygransoStruct()
    f_theta_all = f_theta(x_prime)
    fy = f_theta_all[:, y] # true class output
    # output except true class
    fi = torch.hstack((f_theta_all[:, :y], f_theta_all[:, y+1:]))
    ci.c1 = fy - torch.max(fi)
    ci.c2 = -x_prime
    ci.c3 = x_prime-1
    # equality constraint
    ce = None
    return [f, ci, ce]
# specify optimization variable (tensor)
var_in = {"x_prime": list(x.shape)}
# pygranso main algorithm
soln = pygranso(var_in, comb_fn)
```

CIFAR10 dataset

Compared with FAB [iterative constraint
linearization + projected gradient]

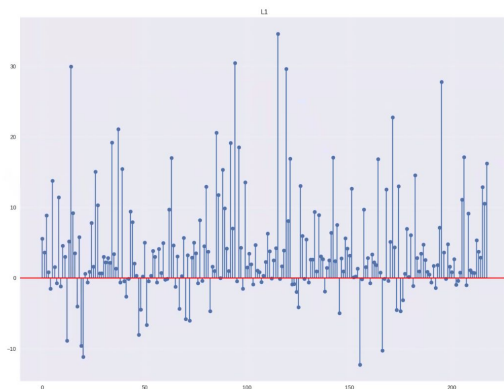
<https://github.com/fra31/auto-attack>

$$\min_{x'} d(x, x')$$

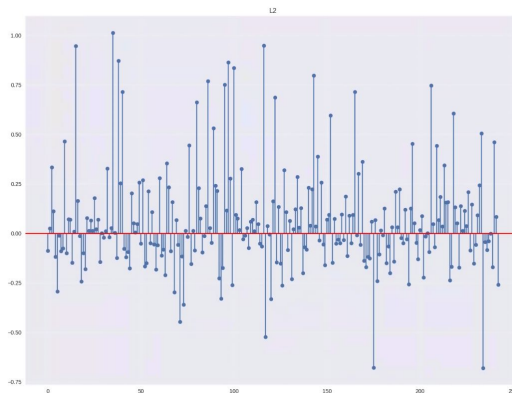
$$\text{s. t. } \max_{\ell \neq c} f_{\theta}^{\ell}(x') \geq f_{\theta}^c(x')$$

$$x' \in [0, 1]^n$$

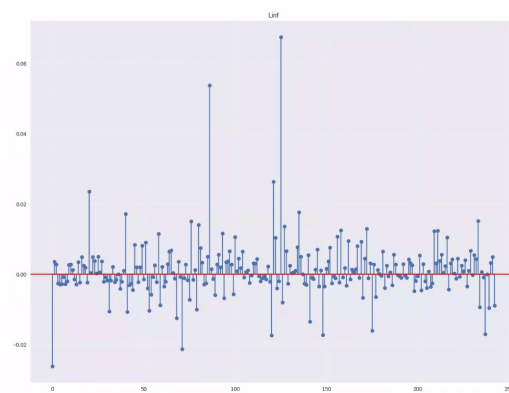
X-axis: image index; Y-axis: PyGRANSO radius - FAB radius



L1 attack



L2 attack



Linf attack

<https://ncvx.org/>

Many
others

NCVX PyGRANSO
Documentation

Search the docs ...

Introduction

Installation

Settings

Examples

- Rosenbrock
- Eigenvalue Optimization
- Dictionary Learning
- Nonlinear Feasibility Problem
- Sphere Manifold
- Trace Optimization
- Robust PCA
- Generalized LASSO
- Logistic Regression
- LeNet5
- Perceptual Attack
- Orthogonal RNN

Highlights



Home



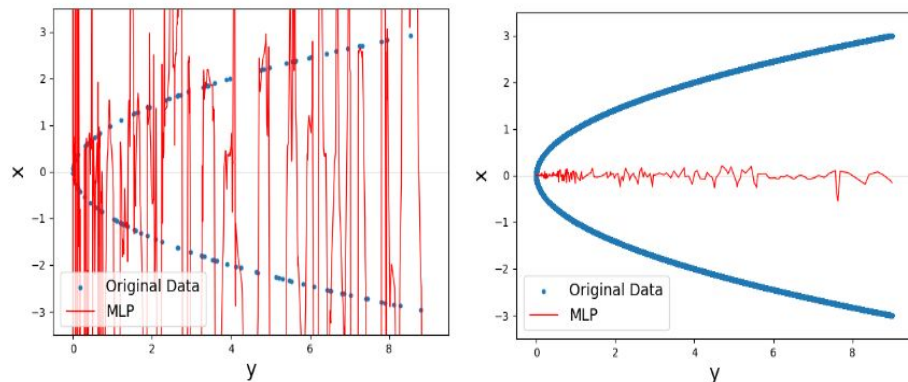
NCVX Package

NCVX (NonConVeX) is a user-friendly and scalable python software package targeting general nonsmooth NCVX problems with nonsmooth constraints. **NCVX** is being developed by **GLOVEX** at the Department of Computer Science & Engineering, University of Minnesota, Twin Cities.

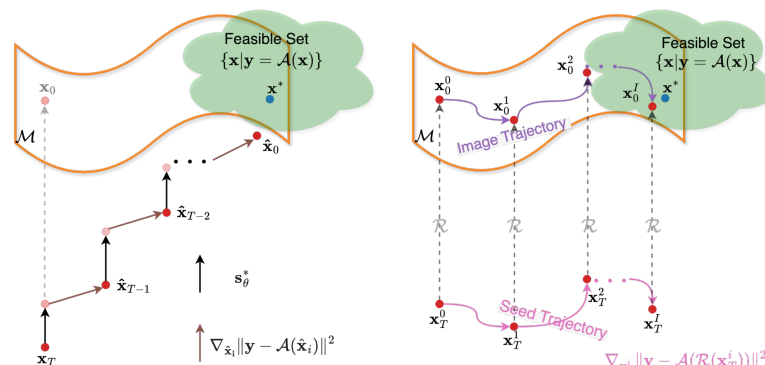
The initial release of **NCVX** contains the solver **PyGRANSO**, a PyTorch-enabled port of **GRANSO** incorporating auto-differentiation, GPU acceleration, tensor input, and support for new QP solvers. As a highlight, **PyGRANSO** can solve general constrained deep learning problems, the first of its kind.



Data-driven methods for SIPs



Story I: More could be less



Story II: Don't be too Bayesian

Single-instance methods for SIPs

Story III: Benefit from DL with a single data point

Deep image prior (DIP) $\mathbf{x} \approx G_{\theta}(\mathbf{z})$ G_{θ} (and \mathbf{z}) trainable

$$\min_{\mathbf{x}} \underbrace{\ell(\mathbf{y}, f(\mathbf{x}))}_{\text{data fitting}} + \lambda \underbrace{R(\mathbf{x})}_{\text{regularizer}}$$

$$\min_{\theta} \ell(\mathbf{y}, f \circ G_{\theta}(\mathbf{z})) + \lambda R \circ G_{\theta}(\mathbf{z})$$

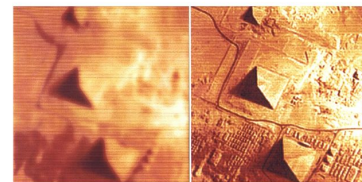
No extra training data!

$$\underbrace{\mathbf{y}}_{\text{blurry and noisy image}} = \underbrace{\mathbf{k}}_{\text{blur kernel}} * \underbrace{\mathbf{x}}_{\text{clean image}} + \underbrace{\mathbf{n}}_{\text{noise}}$$

Mostly due to optical deficiencies (e.g., defocus) and motions

Given \mathbf{y} ,
recover \mathbf{x} (and/or \mathbf{k})

Also **Blind Deconvolution**



A (the?) tool for DL with nontrivial constraints

GRANVISO + PyTorch

 PyGRANSO

NCVX PyGRANSO
Documentation

🔍 Search the docs ...

Introduction

Installation

Settings

Examples



Home



NCVX Package

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \text{ s.t. } c_i(\mathbf{x}) \leq 0, \forall i \in \mathcal{I}; c_i(\mathbf{x}) = 0, \forall i \in \mathcal{E}$$

First general-purpose solver for constrained DL problems

**NCVX: A General-Purpose Optimization Solver for
Constrained Machine and Deep Learning**

Buyun Liang, Tim Mitchell, Ju Sun