

# Deep Regularization via Bi-Level Optimization

Srijith Nair and Philip Schniter

Supported by NSF grant CCF-1955587



THE OHIO STATE UNIVERSITY

2022 Asilomar Conf. Signals, Systems, and Computers

Paper ID: 1259

## Linear inverse problems

**Goal:** Recover signal  $\mathbf{x}$  from noisy measurements  $\mathbf{y}$ :

$$\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{w}.$$

- Applications: deblurring, superresolution, inpainting, computed tomography, magnetic resonance imaging (MRI), etc.

**Challenge:** Often,  $\mathbf{A}$  is not full-column rank

- Components of  $\mathbf{x}$  in  $\text{null}(\mathbf{A})$  are not measured!
- Accurate recovery requires leveraging prior information about  $\mathbf{x}$

## The variational method

A common approach is to recover the signal by solving an optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \{ \ell(\mathbf{x}; \mathbf{y}) + \lambda r(\mathbf{x}) \}$$

- The **data-fidelity term**  $\ell(\cdot; \mathbf{y})$  is typically chosen as the negative log likelihood, e.g.,  $\ell(\mathbf{x}; \mathbf{y}) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$  for white Gaussian  $\mathbf{w}$
- The **regularization**  $r(\cdot)$  is difficult to choose. We want  $r(\cdot)$  such that
  - the cost has no local minima (e.g.,  $r(\cdot)$  is convex)
  - the optimization is not expensive (e.g.,  $r(\cdot)$  is differentiable or proximable)
  - the solution is accurate (i.e.,  $\hat{\mathbf{x}} \approx \mathbf{x}_0$ )
- The regularization **weight**  $\lambda > 0$  should be tuned for best performance
- An example optimization **algorithm** is proximal gradient descent (PGD):

$$\hat{\mathbf{x}}_{\text{new}} = \text{prox}_{\ell\tau/\lambda}(\hat{\mathbf{x}} - \tau \nabla r(\hat{\mathbf{x}})), \quad \tau \in (0, \frac{1}{\text{Lip}(\nabla r)})$$

## Deep regularization

We propose to implement  $r_{\theta}(\cdot)$  using a **deep neural network** and compute the gradient  $\nabla r_{\theta}(\hat{\mathbf{x}})$  in PGD using **automatic differentiation**

Key questions:

- How should  $r_{\theta}(\cdot)$  be constructed?
- How should  $\theta$  be trained?

## Prior work

- RED algorithm:** [1]
  - Train a deep denoiser  $\mathbf{d}_{\theta}(\mathbf{x})$  via  $\hat{\theta} = \arg \min_{\theta} \mathbb{E} \{ \|\mathbf{x}_0 - \mathbf{d}_{\theta}(\mathbf{x}_0 + \mathbf{n})\|^2 \}$
  - Set  $r_{\theta}(\mathbf{x}) \triangleq \frac{1}{2} \mathbf{x}^{\top} (\mathbf{x} - \mathbf{d}_{\theta}(\mathbf{x}))$
  - Use the approximation  $\nabla r_{\theta}(\mathbf{x}) \approx \mathbf{x} - \mathbf{d}_{\theta}(\mathbf{x})$  in PGD
    - For the above to be exact, we need that  $\mathbf{d}_{\theta}(\cdot)$  is Jacobian-symmetric and locally homogeneous, but these properties are not satisfied by most practical denoisers [2]
- Input-convex regularization with adversarial training:** [3]
  - Construct  $r_{\theta}(\cdot)$  using an input-convex deep net [4]
  - Train  $r_{\theta}(\cdot)$  as a Wasserstein-based discriminator of real  $\mathbf{x}$  versus fake  $\mathbf{x}$
  - Use  $\nabla r_{\theta}(\cdot)$  in PGD
- Gradient-step approach:** [5, 6]
  - Construct a deep regularizer  $r_{\theta}(\cdot)$ 
    - e.g., for some deep net  $\mathbf{h}_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , construct  $r_{\theta}(\mathbf{x}) = \|\mathbf{h}_{\theta}(\mathbf{x})\|^2$  or  $r_{\theta}(\mathbf{x}) = \|\mathbf{x} - \mathbf{h}_{\theta}(\mathbf{x})\|^2$  or  $r_{\theta}(\mathbf{x}) = \mathbf{x}^{\top} (\mathbf{x} - \mathbf{h}_{\theta}(\mathbf{x}))$  or etc.
  - Train  $\mathbf{d}_{\theta}(\mathbf{x}) \triangleq \mathbf{x} - \nabla r_{\theta}(\mathbf{x})$  as a denoiser
  - Use  $\nabla r_{\theta}(\cdot)$  in PGD

## Proposed regularizer construction

- We propose to construct the deep regularizer  $r_{\theta}(\cdot)$  as

$$r_{\theta}(\mathbf{x}) = \frac{1+\alpha}{2} \|\mathbf{x}\|^2 - s_{\theta}(\mathbf{x})$$

with  $\alpha \geq 0$  and  **$\beta$ -smooth** neural network  $s_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}$ , i.e.,

$$\|\nabla s_{\theta}(\mathbf{x}_1) - \nabla s_{\theta}(\mathbf{x}_2)\| \leq \beta \|\mathbf{x}_1 - \mathbf{x}_2\| \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$$

- If  $\beta = 1$ , the regularizer  $r_{\theta}(\cdot)$  is **convex** (or **strongly convex** when  $\alpha > 0$ )
- The structure of  $s_{\theta}(\cdot)$  is arbitrary
- We propose to enforce  $\beta$ -smoothness **adversarially**, by adding the following Lipschitz penalty while training the regularizer  $r_{\theta}(\cdot)$  [7]:

$$\max_i \{0, \max_{\delta} \hat{\beta}_{\theta}(\mathbf{x}_i, \delta) - \beta\}$$

where  $\{\mathbf{x}_i\}$  are the training samples and

$$\hat{\beta}_{\theta}(\mathbf{x}, \delta) \triangleq \frac{\|\nabla s_{\theta}(\mathbf{x}) - \nabla s_{\theta}(\mathbf{x} + \delta)\|}{\|\delta\|}$$

## Proposed training scheme

We propose a **bi-level training scheme** [8]:

- $\theta$  and  $\lambda$  are chosen to minimize the loss

$$\mathcal{L}(\theta, \lambda) \triangleq \sum_{i=1}^n \left[ \|\hat{\mathbf{x}}_i(\theta, \lambda) - \mathbf{x}_i\|^2 + \xi \max \{0, \max_{\delta_i} \hat{\beta}_{\theta}(\mathbf{x}_i, \delta_i) - \beta\} \right]$$

where

$$\hat{\mathbf{x}}_i(\theta, \lambda) = \arg \min_{\mathbf{x}} \{ \underbrace{\ell(\mathbf{x}; \mathbf{y}_i) + \lambda r_{\theta}(\mathbf{x})}_{\triangleq J_{\theta}(\mathbf{x}, \mathbf{y}_i)} \}$$

Note:

- The supervised L2 term  $\|\hat{\mathbf{x}}_i(\theta, \lambda) - \mathbf{x}_i\|^2$  is one of several options
- In practice, we alternate between updating  $\theta$ ,  $\lambda$ , and  $\{\delta_i\}$

## Bi-level optimization details

Details of the  $\theta$  optimization:

- For fixed  $\{\delta_i\}$ , define  $\rho(\theta) \triangleq \sum_{i=1}^n \max \{0, \max_{\delta_i} \hat{\beta}_{\theta}(\mathbf{x}_i, \delta_i) - \beta\}$
- Can show that

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = \sum_{i=1}^n \frac{\partial^2 J_{\theta}(\hat{\mathbf{x}}_i; \mathbf{y}_i)}{\partial \theta \partial \mathbf{x}^{\top}} \underbrace{\left[ \frac{\partial^2 J_{\theta}(\hat{\mathbf{x}}_i; \mathbf{y}_i)}{\partial \mathbf{x} \partial \mathbf{x}^{\top}} \right]^{-1}}_{\triangleq \gamma_i} (\mathbf{x}_i - \hat{\mathbf{x}}_i) + \xi \frac{\partial \rho(\theta)}{\partial \theta}$$

- To compute  $\gamma_i$ , use CG (or similar) to solve  $\left[ \frac{\partial^2 J_{\theta}(\hat{\mathbf{x}}_i; \mathbf{y}_i)}{\partial \mathbf{x} \partial \mathbf{x}^{\top}} \right] \gamma_i = \mathbf{x}_i - \hat{\mathbf{x}}_i$
- Use auto-differentiation for gradients & Hessians (via JAX or FuncTorch)

Summary of each  $\theta$  update:

- For each  $i$ , use PGD to compute  $\hat{\mathbf{x}}_i = \arg \min_{\mathbf{x}} \{ \ell(\mathbf{x}; \mathbf{y}_i) + \lambda r_{\theta}(\mathbf{x}) \}$
- For each  $i$ , use CG to compute  $\gamma_i$  above
- Take gradient step  $\theta_{\text{new}} = \theta - \mu \frac{\partial \mathcal{L}(\theta)}{\partial \theta}$  using Adam

## Numerical experiments

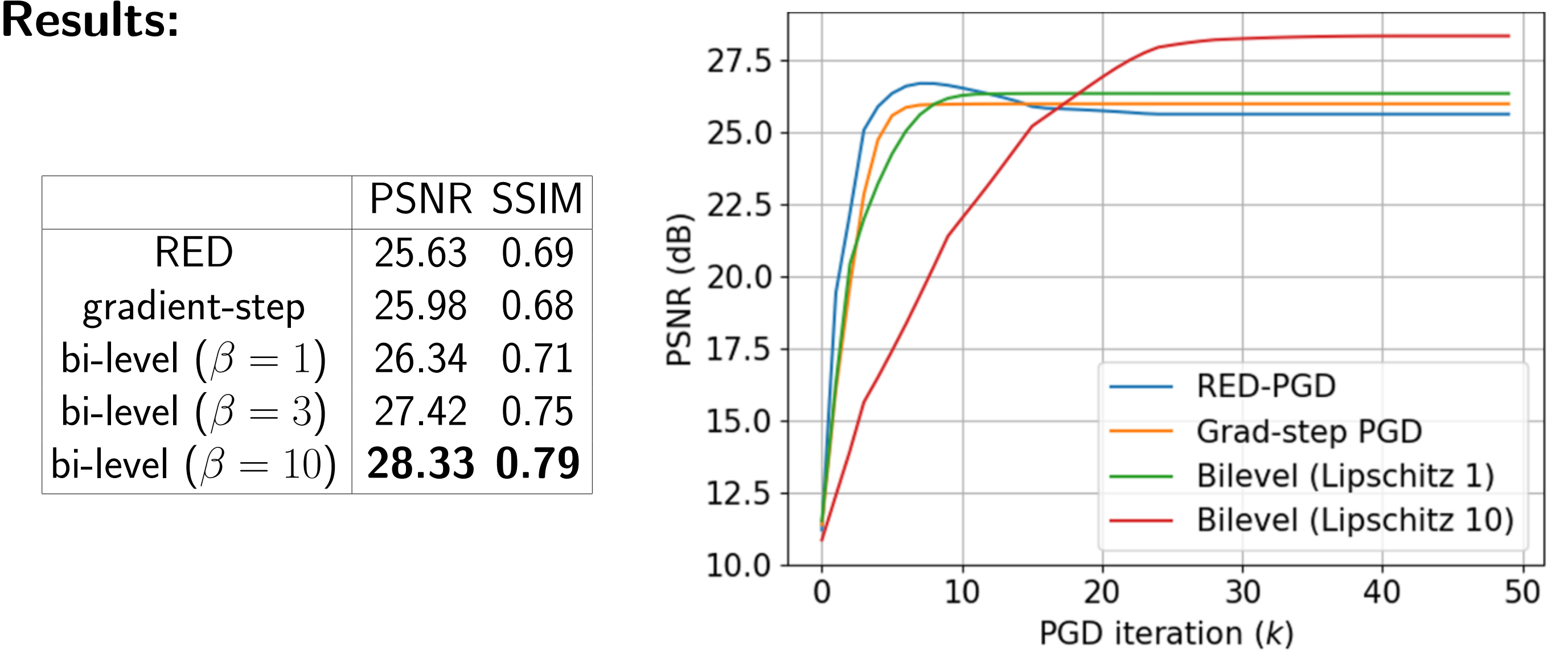
**Inverse problem:**

- Compressive sensing:** Recover  $\mathbf{x}$  from  $\mathbf{y} = \mathbf{A}\mathbf{x}$
- $\mathbf{x} \in \mathbb{R}^d$  is a  $d = 180^2$ -pixel grayscale image from BSD400
  - 370 training images (augmented to 1480 via rotation & flipping)
  - 30 testing images
- $\mathbf{A} \in \mathbb{R}^{\frac{d}{5} \times d}$  is a fast, structurally random matrix [9]

**Deep regularizer:**

- RED-style regularization:**  $r_{\theta}(\mathbf{x}) = \frac{1}{2} \mathbf{x}^{\top} (\mathbf{x} - \mathbf{h}_{\theta}(\mathbf{x}))$ 
  - Our scheme constrains  $\text{Lip}(\nabla s_{\theta})$  where  $s_{\theta}(\mathbf{x}) \triangleq \frac{1}{2} \|\mathbf{x}\|^2 - r_{\theta}(\mathbf{x})$
- $\mathbf{h}_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a **convolutional deep network**
  - 3 layers, 16 channels, SiLU activation, LayerNorm
  - This “simple” network acts as a proof-of-concept
- Recovery used **PGD with backtracking line-search** to automatically tune  $\tau$ 
  - Note: we increase the stepsize  $\tau$  on successful iterations

**Results:**



**Key points:**

- As  $\beta = \text{Lip}(\nabla s)$  increases, the final-PSNR increases and the convergence speed decreases
- Setting  $\beta = 1$  yields a convex regularizer
- The proposed bi-level schemes attain the highest PSNRs

## References

- Y. Romano, M. Elad, and P. Milanfar, “The little engine that could: Regularization by denoising (RED),” *SIAM J. Imag. Sci.*, vol. 10, no. 4, pp. 1804–1844, 2017.
- E. T. Reehorst and P. Schniter, “Regularization by denoising: Clarifications and new interpretations,” *IEEE Trans. Comput. Imag.*, vol. 5, pp. 52–67, Mar. 2019.
- S. Lunz, O. Öktem, and C.-B. Schönlieb, “Adversarial regularizers in inverse problems,” *Proc. Neural Inf. Process. Syst. Conf.*, vol. 31, 2018.
- B. Amos, L. Xu, and J. Z. Kolter, “Input convex neural networks,” in *Proc. Int. Conf. Mach. Learn.*, pp. 146–155, 2017.
- R. Cohen, Y. Blau, D. Freedman, and E. Rivlin, “It has potential: Gradient-driven denoisers for convergent solutions to inverse problems,” in *Proc. Neural Inf. Process. Syst. Conf.*, vol. 34, pp. 18152–18164, 2021.
- S. Hurault, A. Leclaire, and N. Papadakis, “Gradient step denoiser for convergent plug-and-play,” in *Proc. Int. Conf. on Learn. Rep.*, 2022.
- L. Bungert, R. Raab, T. Roith, L. Schwin, and D. Tenbrinck, “CLIP: Cheap Lipschitz training of neural networks,” in *Intl. Conf. Scale Space & Variational Methods in Comp. Vis.*, pp. 307–319, 2021.
- K. G. Samuel and M. F. Tappen, “Learning optimized MAP estimates in continuously-valued MRF models,” in *Proc. IEEE Conf. Comp. Vision Pattern Recog.*, pp. 477–484, 2009.
- T. T. Do, L. Gan, N. H. Nguyen, and T. D. Tran, “Fast and efficient compressive sensing using structurally random matrices,” *IEEE Trans. Signal Process.*, vol. 60, pp. 139–154, Jan. 2012.