

# Regularization by Denoising: Clarifications and New Interpretations

Phil Schniter and Ted Reehorst



THE OHIO STATE UNIVERSITY

With support from: NSF CCF-1527162 and NSF CCF-1716388

SIAM Conference on Imaging Science (Bologna, Italy) — June 7, 2018

# Outline

- Introduction to RED
- Clarifications on RED
- New Interpretation of RED
- Fast RED Algorithms

# Inverse Problems in Imaging

Inverse problems in imaging:

Recover  $x^0$  from measurements  $y = \text{corrupted}(Ax^0)$

where  $A$  is a known linear operator.

In this talk, we'll focus on additive white Gaussian noise (AWGN):

Recover  $x^0$  from measurements  $y = Ax^0 + e$  with  $e \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ .

Other corruptions include loss of phase, quantization, Poisson arrivals...

# The Variational Approach and MAP Estimation

The **variational approach** to recovering  $\mathbf{x}$  solves an optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \{ \ell(\mathbf{x}; \mathbf{y}) + \lambda \rho(\mathbf{x}) \} \text{ with } \begin{cases} \ell(\mathbf{x}; \mathbf{y}) : \text{ loss function} \\ \rho(\mathbf{x}) : \text{ regularization} \\ \lambda > 0 : \text{ tuning parameter} \end{cases}$$

Can be interpreted as **Bayesian MAP** estimation:

$$\hat{\mathbf{x}}_{\text{map}} = \arg \min_{\mathbf{x}} \{ -\ln p(\mathbf{y}|\mathbf{x}) - \ln p(\mathbf{x}) \} \text{ with } \begin{cases} p(\mathbf{y}|\mathbf{x}) : \text{ likelihood} \\ p(\mathbf{x}) : \text{ prior} \end{cases}$$

AWGN likelihood implies **quadratic loss**  $\ell(\mathbf{x}; \mathbf{y}) = \frac{1}{2\sigma^2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2$ .

But how should we choose the regularization  $\rho(\cdot)$ ?

# Regularization by Denoising (RED)

Recently, Romano, Elad and Milanfar<sup>1</sup> proposed the RED regularization

$$\rho_{\text{red}}(\mathbf{x}) \triangleq \frac{1}{2} \mathbf{x}^\top (\mathbf{x} - \mathbf{f}(\mathbf{x})),$$

where  $\mathbf{f} : \mathbb{R}^N \rightarrow \mathbb{R}^N$  is an image denoising function (e.g., BM3D).

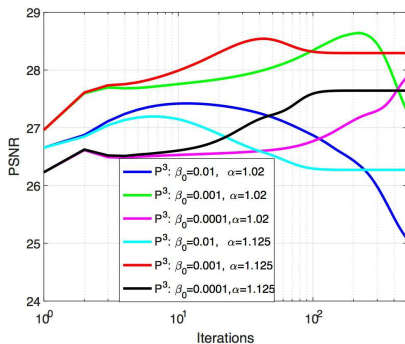
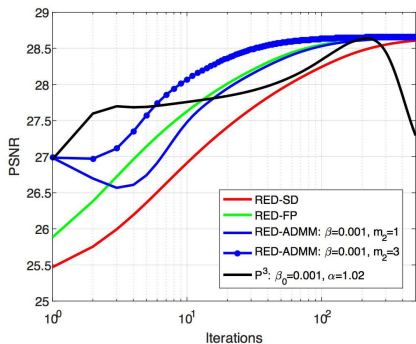
RED leads to a family of “plug-and-play” (PnP) algorithms, similar to those proposed by Bouman et al.<sup>2</sup> and Metzler et al.<sup>3</sup>, but with some advantages.

---

<sup>1</sup>Romano, Elad, Milanfar'17, <sup>2</sup>Venkatkrishnan, Bouman, Wolhberg'13, <sup>3</sup>Metzler, Maleki, Baraniuk'15

## RED versus PnP

Experiments in the RED paper<sup>1</sup> show advantages of RED algs over PnP:



Above represents super-resolution recovery averaged over 10 test images.

# Claims about RED

The RED paper<sup>1</sup> claims ...

- 1 If  $\mathbf{f}(\cdot)$  is **locally homogeneous** (LH), i.e.,

$$\mathbf{f}((1 + \epsilon)\mathbf{x}) = (1 + \epsilon)\mathbf{f}(\mathbf{x}) \text{ for small } \epsilon,$$

and **differentiable**, then gradient of  $\rho_{\text{red}}(\mathbf{x}) \triangleq \frac{1}{2}\mathbf{x}^\top(\mathbf{x} - \mathbf{f}(\mathbf{x}))$  obeys

$$\nabla \rho_{\text{red}}(\mathbf{x}) = \mathbf{x} - \mathbf{f}(\mathbf{x}).$$

- 2 If the Jacobian  $J\mathbf{f}(\mathbf{x})$  is **strongly passive**, i.e.,

$$\|J\mathbf{f}(\mathbf{x})\|_2 \leq 1,$$

then the RED regularization  $\rho_{\text{red}}(\mathbf{x})$  is **convex**.

# Implications of RED Claims

- The convexity claim on  $\rho_{\text{red}}(\cdot)$  implies that minimization of

$$C_{\text{red}}(\mathbf{x}) \triangleq \frac{1}{2\sigma^2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 + \lambda\rho_{\text{red}}(\mathbf{x})$$

can be easily tackled by many algs (e.g., SD, ADMM, etc.).

- The gradient claim  $\nabla\rho_{\text{red}}(\mathbf{x}) = \mathbf{x} - \mathbf{f}(\mathbf{x})$  implies the minimizers obey

RED fixed-point condition:  $\frac{1}{\sigma^2} \mathbf{A}^\top (\mathbf{A}\hat{\mathbf{x}} - \mathbf{y}) + \lambda(\hat{\mathbf{x}} - \mathbf{f}(\hat{\mathbf{x}})) = \mathbf{0}$

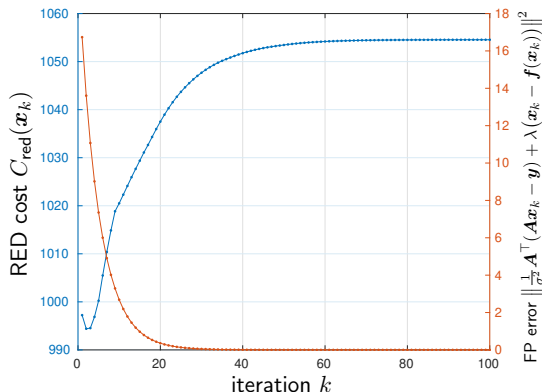
The RED algorithms find exactly these  $\hat{\mathbf{x}}$ .



# Mysterious Behavior

Surprisingly, the RED algorithms do not always behave as expected!

We expect SD to decrease the (convex) RED cost, but it is **increasing** it!



$$\text{RED-SD: } \mathbf{x}_{k+1} = \mathbf{x}_k - \mu \nabla C_{\text{red}}(\mathbf{x}_k)$$

# Clarifications on RED Gradient

It can be shown that...

- **differentiability** in  $\mathbf{f}(\cdot)$  implies

$$\nabla \rho_{\text{red}}(\mathbf{x}) \stackrel{\text{D}}{=} \mathbf{x} - \frac{1}{2} \mathbf{f}(\mathbf{x}) - \frac{1}{2} [J \mathbf{f}(\mathbf{x})]^\top \mathbf{x}.$$

- adding **local-homogeneity** (LH) gives

$$\nabla \rho_{\text{red}}(\mathbf{x}) \stackrel{\text{D,LH}}{=} \mathbf{x} - \frac{1}{2} [J \mathbf{f}(\mathbf{x})] \mathbf{x} - \frac{1}{2} [J \mathbf{f}(\mathbf{x})]^\top \mathbf{x}.$$

- adding **Jacobian symmetry** (JS) finally leads to

$$\nabla \rho_{\text{red}}(\mathbf{x}) \stackrel{\text{D,LH,JS}}{=} \mathbf{x} - \mathbf{f}(\mathbf{x}) \quad \dots \text{which yields the RED algorithms.}$$

So both LH and JS are needed to link RED cost to RED algs.

## Which Denoisers Yield Jacobian Symmetry?

Clear that these yield JS:

- Linear denoisers  $\mathbf{f}(\mathbf{x}) = \mathbf{W}\mathbf{x}$  with  $\mathbf{W} = \mathbf{W}^\top$ .
- Transform-domain-thresholding (TDT) denoisers  $\mathbf{f}(\mathbf{x}) = \mathbf{W}^\top \mathbf{g}(\mathbf{W}\mathbf{x})$ .
- MAP or MMSE denoisers under any assumed prior  $\mathbf{x} \sim \hat{p}_{\mathbf{x}}$ .

Not clear that these yield JS:

- Pseudo-linear denoisers  $\mathbf{f}(\mathbf{x}) = \mathbf{W}(\mathbf{x})\mathbf{x}$  with non-linear  $\mathbf{W}(\cdot)$ .
- *Approximately* MAP or MMSE denoisers.

Most state-of-the-art denoisers fall into the 2nd category.

# Jacobian Symmetry Experiments

Avg JS error on suite of  $16 \times 16$  images:

	TDT	MF	NLM	BM3D	TNRD	DnCNN
$\frac{\ \widehat{Jf}(\mathbf{x}) - [\widehat{Jf}(\mathbf{x})]^\top\ _F^2}{\ \widehat{Jf}(\mathbf{x})\ _F^2}$	4.11e-21	1.35	0.118	0.186	0.0151	0.194

Avg gradient error on suite of  $16 \times 16$  images:

$\frac{\ \nabla \rho_{\text{red}}(\mathbf{x}) - \widehat{\nabla \rho_{\text{red}}}(\mathbf{x})\ ^2}{\ \widehat{\nabla \rho_{\text{red}}}(\mathbf{x})\ ^2}$	TDT	MF	NLM	BM3D	TNRD	DnCNN
$\nabla \rho_{\text{red}}(\mathbf{x})$ with D	3.39e-19	2.65e-15	6.17e-21	2.14e-13	5.42e-17	1.02e-12
$\nabla \rho_{\text{red}}(\mathbf{x})$ with D,LH,JS	0.565	0.966	0.913	1.00	0.957	0.852

Key points:

- 1 Large JS error for all but TDT.
- 2 Large gradient error under JS & LH assumptions for all denoisers!
- 3 Even TDT has large gradient error! Is LH the problem?

## Local Homogeneity Experiments

Avg LH error on suite of  $16 \times 16$  images:

	TDT	MF	NLM	BM3D	TNRD	DnCNN
$\frac{\ \mathbf{f}((1+\epsilon)\mathbf{x}) - (1+\epsilon)\mathbf{f}(\mathbf{x})\ ^2}{\ (1+\epsilon)\mathbf{f}(\mathbf{x})\ ^2}$	7.99e-10	0	5.60e-9	1.52e-13	5.09e-10	2.06e-9
$\frac{\ \widehat{J\mathbf{f}}(\mathbf{x})\mathbf{x} - \mathbf{f}(\mathbf{x})\ ^2}{\ \mathbf{f}(\mathbf{x})\ ^2}$	4.10e-4	2.14e-15	5.63e-3	0.214	2.60e-4	8.02e-3

Avg gradient error on suite of  $16 \times 16$  images:

$\frac{\ \nabla\rho_{\text{red}}(\mathbf{x}) - \widehat{\nabla\rho_{\text{red}}}(\mathbf{x})\ ^2}{\ \widehat{\nabla\rho_{\text{red}}}(\mathbf{x})\ ^2}$	TDT	MF	NLM	BM3D	TNRD	DnCNN
$\nabla\rho_{\text{red}}(\mathbf{x})$ with D	3.39e-19	2.65e-15	6.17e-21	2.14e-13	5.42e-17	1.02e-12
$\nabla\rho_{\text{red}}(\mathbf{x})$ with D,LH	0.565	6.09e-15	0.0699	0.344	0.139	1.20

Key points:

- It is important how LH is quantified.
- The RED gradient is **very sensitive to small imperfections in LH.**

# Implications of our Findings

We found:

- The RED algorithms solve a fixed-point equation corresponding to  $\nabla\rho(\mathbf{x}) = \mathbf{x} - \mathbf{f}(\mathbf{x})$ .
- $\mathbf{x} - \mathbf{f}(\mathbf{x})$  is very different from  $\nabla\rho_{\text{red}}(\mathbf{x})$  under practical  $\mathbf{f}(\cdot)$ , such as TDT, MF, NLM, BM3D, TNRD, and DnCNN.

Implication:

- $\rho_{\text{red}}(\cdot)$  does not explain the RED algorithms under practical  $\mathbf{f}(\cdot)$ .

A bigger problem:

- For non-JS  $\mathbf{f}(\cdot)$ , can show that **there exists no explicit regularizer  $\rho(\cdot)$  for which  $\nabla\rho(\mathbf{x}) = \mathbf{x} - \mathbf{f}(\mathbf{x})$ , i.e., explaining the RED algorithms!**

# How to Explain the RED Algorithms?

The RED algorithms assume  $\nabla\rho(\mathbf{x}) = \mathbf{x} - \mathbf{f}(\mathbf{x})$  and work very well.

Can we justify this  $\nabla\rho(\mathbf{x})$ ?

Even when  $\mathbf{f}(\cdot)$  is not LH and/or JS?

Yes! Using [score matching](#). We explain this in 3 steps:

- 1 regularization by log-likelihood (RLL),
- 2 RLL as kernel density estimation (KDE),
- 3 score matching.

## Regularization by Log-Likelihood (RLL)

- Consider noisy pseudo-measurements

$$\mathbf{r} = \mathbf{x}^0 + \mathcal{N}(0, \nu \mathbf{I}).$$

- Suppose we adopt the prior pdf  $\hat{p}_{\mathbf{x}}$ . Then the likelihood of  $\mathbf{r}$  is

$$\hat{p}_{\mathbf{r}}(\mathbf{r}; \nu) = \int_{\mathbb{R}^N} \mathcal{N}(\mathbf{r}; \mathbf{x}, \nu \mathbf{I}) \hat{p}_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}. \quad \text{“Gaussian blurred prior”}$$

- Define the **RLL regularization** as

$$\rho_{\text{LL}}(\mathbf{r}; \nu) \triangleq -\nu \ln \hat{p}_{\mathbf{r}}(\mathbf{r}; \nu)$$

- Then it can be shown using Tweedie's formula<sup>4</sup> that

$$\nabla \rho_{\text{LL}}(\mathbf{r}; \nu) = \mathbf{r} - \hat{\mathbf{f}}_{\text{mmse}, \nu}(\mathbf{r}),$$

which is **consistent with the RED algorithms!**

---

<sup>4</sup>Robbins'56



# RLL as Kernel Density Estimation

- Given training data  $\{\mathbf{x}_t\}_{t=1}^T$ , consider the **empirical prior**

$$\hat{p}_{\mathbf{x}}(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \delta(\mathbf{x} - \mathbf{x}_t).$$

- A better match to the true  $p_{\mathbf{x}}$  is obtained via **KDE** or **Parzen windowing**:

$$\tilde{p}_{\mathbf{x}}(\mathbf{x}; \nu) = \frac{1}{T} \sum_{t=1}^T \mathcal{N}(\mathbf{x}; \mathbf{x}_t, \nu \mathbf{I}). \quad \text{“blurred empirical prior”}$$

- Using this  $\tilde{p}_{\mathbf{x}}$  for MAP/variational optimization yields

$$\begin{aligned} \hat{\mathbf{x}} &= \arg \min_{\mathbf{x}} \frac{1}{2\sigma^2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 - \ln \tilde{p}_{\mathbf{x}}(\mathbf{x}; \nu) \\ &= \arg \min_{\mathbf{x}} \frac{1}{2\sigma^2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 + \lambda \rho_{\text{LL}}(\mathbf{x}; \nu) \text{ for } \lambda = \frac{1}{\nu}. \end{aligned}$$

So **RLL arises naturally in non-parametric estimation via KDE!**

# Score-Matching by Denoising

- The above RLL/KDE framework encompasses **only JS** denoisers  $\mathbf{f}(\cdot)$ . We now generalize.

- First note that, for large # of examples  $T$ , gradient is very expensive:

$$\nabla \ln \tilde{p}_{\mathbf{x}}(\mathbf{x}; \nu) = \frac{\hat{\mathbf{f}}_{\text{mmse}, \nu}(\mathbf{x}) - \mathbf{x}}{\nu} \quad \text{with} \quad \hat{\mathbf{f}}_{\text{mmse}, \nu}(\mathbf{x}) = \frac{\sum_{t=1}^T (\mathbf{x}_t - \mathbf{x}) \mathcal{N}(\mathbf{x}; \mathbf{x}_t, \nu \mathbf{I})}{\sum_{t=1}^T \mathcal{N}(\mathbf{x}; \mathbf{x}_t, \nu \mathbf{I})}.$$

- Practical idea:<sup>5</sup> use best match to “score”  $\nabla \ln \tilde{p}_{\mathbf{x}}(\mathbf{x})$  among computationally friendly functions  $\psi(\mathbf{x}; \boldsymbol{\theta})$ :

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathbb{E}_{\tilde{p}_{\mathbf{x}}} \left\{ \left\| \psi(\mathbf{x}; \hat{\boldsymbol{\theta}}) - \nabla \ln \tilde{p}_{\mathbf{x}}(\mathbf{x}; \nu) \right\|^2 \right\}.$$

- Vincent<sup>6</sup> **connected to denoising**: if  $\psi(\mathbf{x}; \boldsymbol{\theta}) = [\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}) - \mathbf{x}]/\nu$ , then

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left\{ \left\| \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}_t + \mathcal{N}(0, \nu \mathbf{I})) - \mathbf{x}_t \right\|^2 \right\},$$

where  $\mathbf{f}_{\hat{\boldsymbol{\theta}}}(\cdot)$  is **MMSE optimal**  $\mathbf{f}_{\boldsymbol{\theta}} \in \mathcal{F}$ , where  $\mathcal{F} \triangleq \{\mathbf{f}_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \Theta\}$ .

<sup>5</sup>Hyvärinen'05, <sup>6</sup>Vincent'11

# Score-Matching by Denoising (SMD)

Key points:

- 1 SMD interpretation yields  $\nabla\rho(x) = x - f(x)$ , thus **explaining RED** algs.
- 2 SMD interpretation holds for **any**  $\hat{p}_x$ , any denoiser class  $\mathcal{F}$  (i.e.,  $f_\theta$  may be **non-JS and/or non-LH**), and any  $\theta$  (maybe **not MMSE**).
- 3 SMD arises naturally via **non-parametric estimation** and KDE. Matches construction of *learned* denoisers liked TNRD and DnCNN.

Related work:

Alain and Bengio<sup>7</sup> recently showed that learned auto-encoders can be explained by score-matching and *not* by minimization of an energy function.

---

<sup>7</sup>Alain/Bengio'14

# Fast RED Algorithms

Until now we focused on how to *explain* the RED algorithms, which solve

$$\text{RED fixed-point condition: } \frac{1}{\sigma^2} \mathbf{A}^\top (\mathbf{A} \hat{\mathbf{x}} - \mathbf{y}) + \lambda (\hat{\mathbf{x}} - \mathbf{f}(\hat{\mathbf{x}})) = \mathbf{0}$$

We now focus on interpretation/design of fast RED algorithms.

In the RED paper, three algorithms were described:

- 1 Steepest-Descent
- 2 ADMM with  $I$  inner iters (to solve  $\arg \min_{\mathbf{x}} \{\lambda \rho(\mathbf{x}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{r}_k\|^2\}$ )
- 3 A “fixed-point” method (we show equivalence to proximal gradient alg<sup>8</sup>)

We propose a couple more. . .

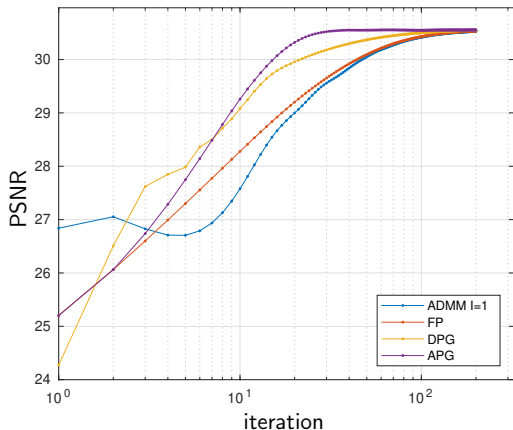
---

<sup>8</sup>Combettes/Pesquet'11

# Algorithm Comparison: Image Deblurring

New algorithms:

- **DPG**: “Dynamic” proximal gradient, which schedules the stepsize.
- **APG**: Accelerated proximal gradient, similar to FISTA.<sup>9</sup>



In this experiment, APG is about **3× faster** than the Fixed-Point method.

<sup>9</sup>Beck/Teboulle'09

# Conclusions

- The RED algorithms work very well in practice.
- But they do not minimize  $C_{\text{red}}(\mathbf{x}) = \ell(\mathbf{x}; \mathbf{y}) + \lambda \rho_{\text{red}}(\mathbf{x})$  for many  $\mathbf{f}(\cdot)$ .
  - Why? Practical denoisers  $\mathbf{f}(\cdot)$  are not sufficiently LH and JS.
  - Can construct examples of RED-SD *increasing*  $C_{\text{red}}(\mathbf{x})$  over the iterations.
- We explained RED algorithms as “score-matching by denoising”.
- We proposed new RED algorithms with faster convergence.

<http://arxiv.org/abs/1806.02296>