

Regularization by Denoising: Clarifications and New Interpretations

Phil Schniter and Ted Reehorst



THE OHIO STATE UNIVERSITY

With support from NSF CCF-1716388

Allerton Conference (Monticello, IL) — Oct. 4, 2018

Outline

- Introduction to RED
- Clarifications on RED
- New Interpretations of RED
- Fast and Convergent RED Algorithms

Inverse Problems in Imaging

- Inverse problems in imaging:

Recover \mathbf{x}^0 from measurements $\mathbf{y} = \text{corrupted}(\mathbf{A}\mathbf{x}^0)$,
where \mathbf{A} is a known linear operator.

- Corruptions include noise, quantization, loss of phase, Poisson. . .
- Operator \mathbf{A} depends on the application:
 - deblurring
 - super-resolution
 - compressive imaging
 - inpainting
 - etc

Optimization-Based Recovery and MAP Estimation

- A common approach to recovering image \mathbf{x} is through **optimization**:

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

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

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

- The loss function $\ell(\cdot; \mathbf{y})$ is usually straightforward to choose.
But how do we **choose the regularization** $\rho(\cdot)$?

Plug-and-Play ADMM

- A common approach to convex optimization is **ADMM**: For $k = 1, 2, \dots$

$$\mathbf{x}_k = \arg \min_{\mathbf{x}} \left\{ \ell(\mathbf{x}; \mathbf{y}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{v}_{k-1} + \mathbf{u}_{k-1}\|^2 \right\}$$

$$\mathbf{v}_k = \arg \min_{\mathbf{v}} \left\{ \rho(\mathbf{v}) + \frac{\beta}{2} \|\mathbf{v} - \mathbf{x}_k + \mathbf{u}_{k-1}\|^2 \right\} \triangleq \text{prox}_{\rho/\beta}(\mathbf{x}_k - \mathbf{u}_{k-1})$$

$$\mathbf{u}_k = \mathbf{u}_{k-1} + \mathbf{x}_k - \mathbf{v}_k$$

- The prox performs **denoising** (eg, soft-thresholding when $\rho(\mathbf{x}) = \|\mathbf{x}\|_1$).
- Bouman et al. proposed **plug-and-play (PnP) ADMM**,¹ where the prox is replaced by a sophisticated image denoiser $\mathbf{f}(\cdot)$ like BM3D.

¹ Venkatakrisnan, Bouman, Wolkberg'13

Regularization by Denoising (RED)

- Recently, Romano, Elad and Milanfar² proposed a new family of PnP algorithms that find the image estimate $\hat{\mathbf{x}}$ that obeys

$$\nabla \ell(\hat{\mathbf{x}}; \mathbf{y}) + \lambda(\hat{\mathbf{x}} - \mathbf{f}(\hat{\mathbf{x}})) = \mathbf{0}$$

- They claimed these algs result from optimization under the regularizer

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

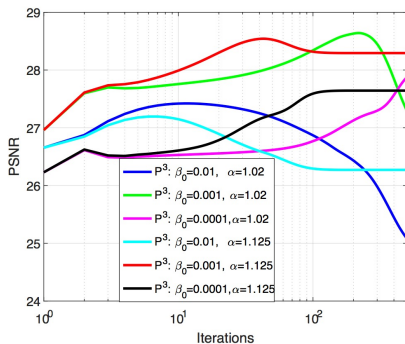
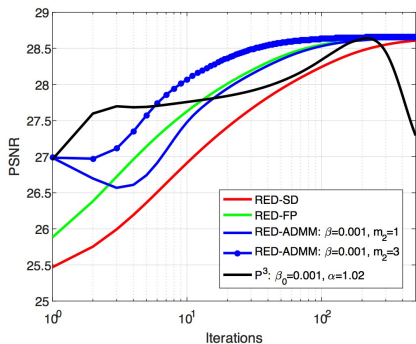
and thus coined the approach **Regularization by Denoising (RED)**.

- They furthermore claimed that $\rho_{\text{red}}(\cdot)$ was convex in practice.

²Romano, Elad, Milanfar'17

RED versus PnP-ADMM

Experiments in the RED paper² suggest advantages over PnP-ADMM:



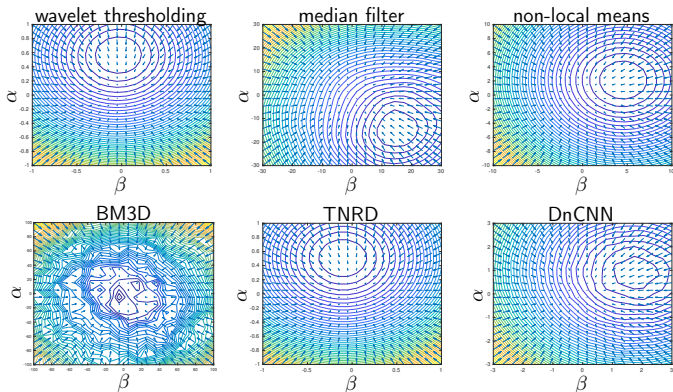
Super-resolution recovery, averaged over 10 test images.

Are the RED algs explained by the RED regularization?

Visualize by probing in two random directions: $\mathbf{x}_{\alpha,\beta} = \hat{\mathbf{x}} + \alpha\mathbf{r}_1 + \beta\mathbf{r}_2$.

Contours show cost: $C_{\text{red}}(\mathbf{x}_{\alpha,\beta}) \triangleq \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{x}_{\alpha,\beta}\|^2 + \rho_{\text{red}}(\mathbf{x}_{\alpha,\beta})$.

Arrows show gradient: $\nabla_{\alpha,\beta} C_{\text{red}}(\mathbf{x}_{\alpha,\beta})$.



Zero of gradient field is not at cost minimizer!

And cost is not convex!

Clarifications on RED Gradient

It can be shown³ that...

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

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

- adding **local-homogeneity** (LH), i.e., $\mathbf{f}((1 + \epsilon)\mathbf{x}) = (1 + \epsilon)\mathbf{f}(\mathbf{x})$, we get

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

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

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

But practical denoisers are **not LH and JS!**

And there exists no regularizer ρ_{red} for a non-JS denoiser \mathbf{f} !

³Reehorst & Schniter, 2018.

How To Explain the RED Algorithms?

The RED algorithms solve $\nabla \ell(\hat{\mathbf{x}}; \mathbf{y}) + \lambda(\hat{\mathbf{x}} - \mathbf{f}(\hat{\mathbf{x}})) = \mathbf{0}$ and work well.

Can we justify this approach?

Even when $\mathbf{f}(\cdot)$ is not locally homogeneous or Jacobian symmetric?

Yes! Using **score matching**.⁴ We explain this in 3 steps:

- 1 kernel density estimation,
- 2 Tweedie's formula,
- 3 score matching.

⁴Hyvärinen'05.

Kernel Density Estimation (KDE)

- Given training data $\{\mathbf{x}_t\}_{t=1}^T$, consider forming 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 **Parzen windowing** or **KDE**:

$$\begin{aligned} \tilde{p}_{\mathbf{x}}(\mathbf{x}; \nu) &= \frac{1}{T} \sum_{t=1}^T \mathcal{N}(\mathbf{x}; \mathbf{x}_t, \nu \mathbf{I}) && \text{“smoothed prior”} \\ &= \int_{\mathbb{R}^N} \mathcal{N}(\mathbf{r}; \mathbf{x}, \nu \mathbf{I}) \hat{p}_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

- Using the smoothed prior $\tilde{p}_{\mathbf{x}}$ for MAP image recovery, we get

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \{ \ell(\mathbf{x}; \mathbf{y}) - \ln \tilde{p}_{\mathbf{x}}(\mathbf{x}; \nu) \}.$$

Tweedie's Formula

- Assuming differentiability, the MAP estimation problem is solved by

$$\mathbf{0} = \nabla \ell(\mathbf{x}; \mathbf{y}) - \nabla \ln \tilde{p}_{\mathbf{x}}(\mathbf{x}; \nu).$$

- Tweedie's formula⁵ says that

$$\nabla \ln \tilde{p}_{\mathbf{x}}(\mathbf{x}; \nu) = \frac{1}{\nu} (\mathbf{f}_{\text{mmse}, \nu}(\mathbf{x}) - \mathbf{x}),$$

with $\mathbf{f}_{\text{mmse}, \nu}(\mathbf{r})$ the MMSE denoiser of $\mathbf{x} \sim \hat{p}_{\mathbf{x}}$ from $\mathbf{r} = \mathbf{x} + \mathcal{N}(\mathbf{0}, \nu \mathbf{I})$.

- Together, these results match the RED fixed-point equation

$$\mathbf{0} = \nabla \ell(\mathbf{x}; \mathbf{y}) + \lambda (\mathbf{x} - \mathbf{f}_{\text{mmse}, \nu}(\mathbf{x})) \quad \text{with} \quad \lambda = \frac{1}{\nu}$$

for the specific denoiser $\mathbf{f}_{\text{mmse}, \nu}$. What about other \mathbf{f} ?

⁵Robbins'56

Score-Matching by Denoising

- Recall $\mathbf{f}_{\text{mmse},\nu} = \arg \min_{\mathbf{f}} \mathbb{E}\{\|\mathbf{x} - \mathbf{f}(\mathbf{r})\|^2\}$ for $\begin{cases} \mathbf{r} = \mathbf{x} + \mathcal{N}(\mathbf{0}, \nu \mathbf{I}) \\ \mathbf{x} \sim \hat{p}_{\mathbf{x}}. \end{cases}$
- Since $\mathbf{f}_{\text{mmse},\nu}$ is expensive to implement, use approximation $\mathbf{f}_{\hat{\theta}}$ with $\hat{\theta} = \arg \min_{\theta} \mathbb{E}\{\|\mathbf{x} - \mathbf{f}_{\theta}(\mathbf{r})\|^2\}$ e.g., deep network

$$= \arg \min_{\theta} \mathbb{E}\{\|\mathbf{x} - \mathbf{f}_{\text{mmse},\nu}(\mathbf{r})\|^2\} + \mathbb{E}\{\|\mathbf{f}_{\text{mmse},\nu}(\mathbf{r}) - \mathbf{f}_{\theta}(\mathbf{r})\|^2\}$$
 via orthog principle

$$= \arg \min_{\theta} \mathbb{E}\{\|\mathbf{f}_{\text{mmse},\nu}(\mathbf{r}) - \mathbf{f}_{\theta}(\mathbf{r})\|^2\}$$

$$= \arg \min_{\theta} \mathbb{E}\left\{\underbrace{\|\nabla \ln \tilde{p}_{\mathbf{x}}(\mathbf{r}; \nu)\|}_{\text{"score"}} + \underbrace{\frac{1}{\nu}(\mathbf{f}_{\theta}(\mathbf{r}) - \mathbf{r})}_{\text{RED with } \mathbf{f}_{\theta}}\right\}^2$$
 via Tweedie.
- Thus RED with general \mathbf{f}_{θ} can be interpreted as “score matching.”

Score-Matching by Denoising (SMD)

Key points:

- 1 RED alg's solve $\mathbf{0} = \nabla \ell(\mathbf{x}; \mathbf{y}) + \lambda(\mathbf{x} - \mathbf{f}_\theta(\mathbf{x}))$ where $\lambda(\mathbf{x} - \mathbf{f}_\theta(\mathbf{x}))$ approximates the score $-\nabla \ln \tilde{p}_x(\mathbf{x}; \nu)$.
- 2 This SMD interpretation holds for any \hat{p}_x , any denoiser class \mathbf{f}_θ (i.e., may be non-JS and/or non-LH), and any θ .
- 3 SMD arises naturally via non-parametric estimation (i.e., KDE).
Matches construction of *learned* denoisers liked TNRD and DnCNN.

Related work:

Alain and Bengio⁶ showed that learned auto-encoders are explained by score-matching and *not* by minimization of an energy function.

⁶Alain/Bengio'14

Fast RED Algorithms

Until now we focused on how to explain the RED method, which solves

$$\mathbf{0} = \nabla \ell(\hat{\mathbf{x}}; \mathbf{y}) + \lambda(\hat{\mathbf{x}} - \mathbf{f}(\hat{\mathbf{x}})).$$

Now we focus on algorithms that try to solve this equation.

In the RED paper, three algorithms were described:

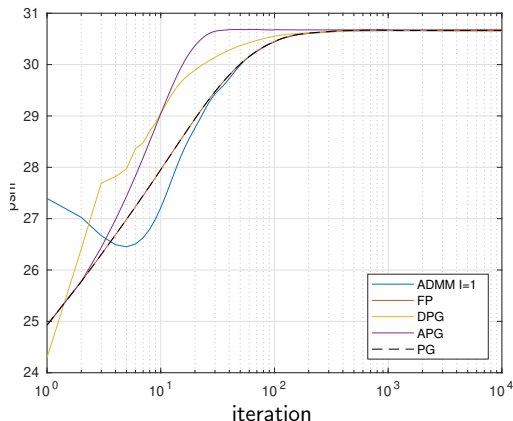
- 1 steepest-descent
- 2 ADMM with I inner iters (to solve $\arg \min_{\mathbf{x}} \{\lambda \rho_{\text{red}}(\mathbf{x}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{r}_t\|^2\}$)
- 3 a heuristic “fixed-point” method.

We propose a several others. . .

Algorithm Comparison: Image Deblurring

New algorithms:

- **PG**: Proximal gradient with stepsize $L > 0$.
- **DPG**: “Dynamic” proximal gradient, which schedules L_t .
- **APG**: Accelerated proximal gradient, similar to FISTA.⁷



In this experiment, APG is about $3\times$ faster than the Fixed-Point method.

⁷Beck/Teboulle'09

Convergence to a Fixed Point

Theorem

If $\ell(\cdot)$ is proper, convex, and continuous; $f(\cdot)$ is non-expansive; $L > 1$; and RED-PG has at least one fixed point, then RED-PG converges to a fixed point.

Proof.

Uses α -averaged operators and Mann iteration. □

Conclusions

- RED algorithms seem to work well in practice.
- But, in practice, they are *not* minimizing any cost function.
 - Practical denoisers $f(\cdot)$ are not LH and JS.
 - Non-JS $f \Rightarrow$ that there exists no regularizer ρ s.t. $\nabla\rho(\mathbf{x}) = \mathbf{x} - f(\mathbf{x})$.
- The RED methodology can be explained as “score-matching by denoising”.
- We proposed new RED algorithms with i) faster recovery and ii) guaranteed convergence to a fixed point.

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