

Regularization by Denoising: Clarifications and New Interpretations

Phil Schniter and Ted Reehorst

Supported by NSF grants IIP-0968910 and CCF-1716388 and NIH grant R01HL135489



THE OHIO STATE UNIVERSITY

BASP Frontiers Workshop 2019

Inverse Problems in Imaging

- Consider the basic **inverse problem** 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 photons**, etc.

- The 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 the image \mathbf{x} is through posing and solving an **optimization** problem:

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

- This 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 some $\beta > 0$ and $k = 1, 2, 3, \dots$

$$\begin{aligned} \mathbf{x}_k &= \arg \min_{\mathbf{x}} \{ \ell(\mathbf{x}; \mathbf{y}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{v}_{k-1} + \mathbf{u}_{k-1}\|^2 \} \\ \mathbf{v}_k &= \arg \min_{\mathbf{v}} \{ \rho(\mathbf{v}) + \frac{\beta}{2} \|\mathbf{v} - \mathbf{x}_k + \mathbf{u}_{k-1}\|^2 \} \triangleq \text{prox}_{\rho/\beta}(\mathbf{x}_k - \mathbf{u}_{k-1}) \\ \mathbf{u}_k &= \mathbf{u}_{k-1} + \mathbf{x}_k - \mathbf{v}_k \end{aligned}$$

- The prox operation performs **denoising** (eg, soft-thresholding when $\rho(\mathbf{x}) = \|\mathbf{x}\|_1$).

- In 2013, Bouman et al. proposed **plug-and-play (PnP) ADMM**, where the prox is replaced by a sophisticated image denoiser $\mathbf{f}(\cdot)$, such as BM3D.

Regularization by Denoising (RED)

- In 2017, Romano, Elad, and Milanfar proposed a new family of PnP algorithms that find the image estimate $\hat{\mathbf{x}}$ that obeys

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

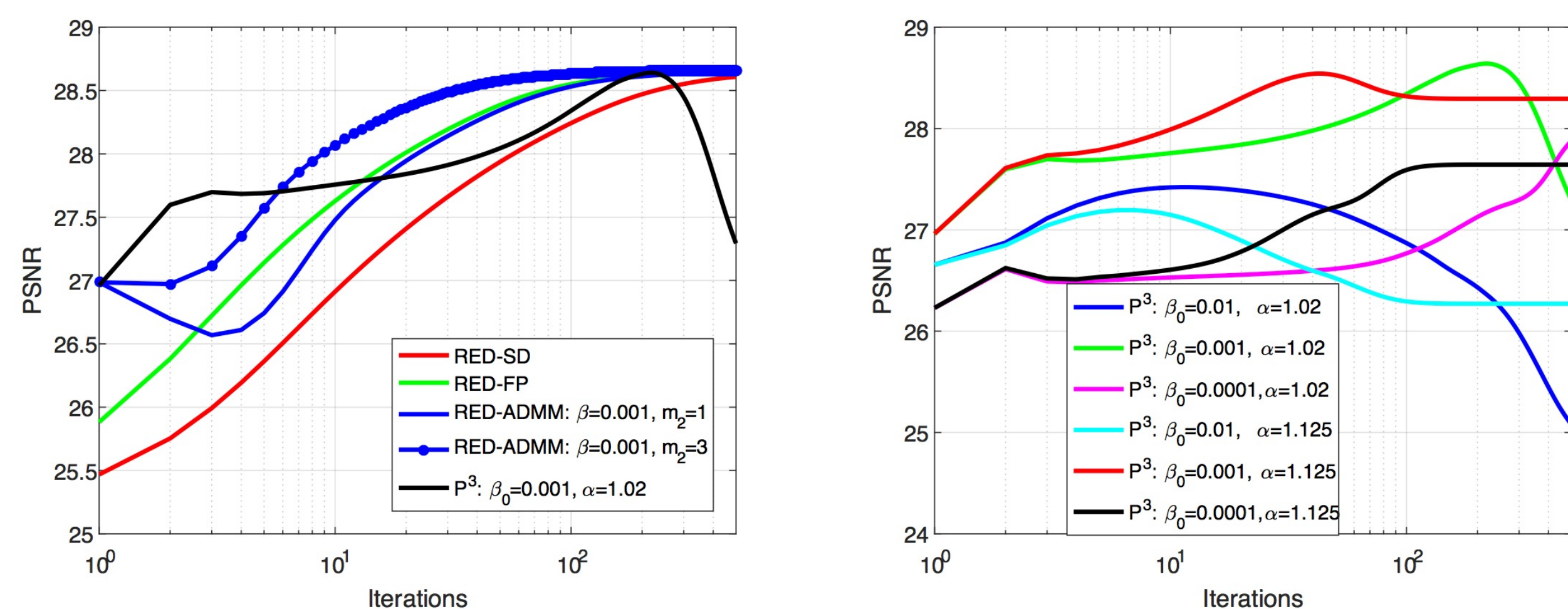
- They claimed these algorithms result from optimization under the regularizer

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

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

- They furthermore claimed that $\rho_{\text{red}}(\cdot)$ was convex with practical image denoisers $\mathbf{f}(\cdot)$.

- Experiments in the RED paper suggest advantages for RED over PnP-ADMM:



Super-resolution recovery, averaged over 10 test images.

The RED algorithms are not explained by the RED regularization!

- Visualize by probing in two random directions:

$$\mathbf{x}_{\alpha,\beta} = \hat{\mathbf{x}} + \alpha \mathbf{r}_1 + \beta \mathbf{r}_2.$$

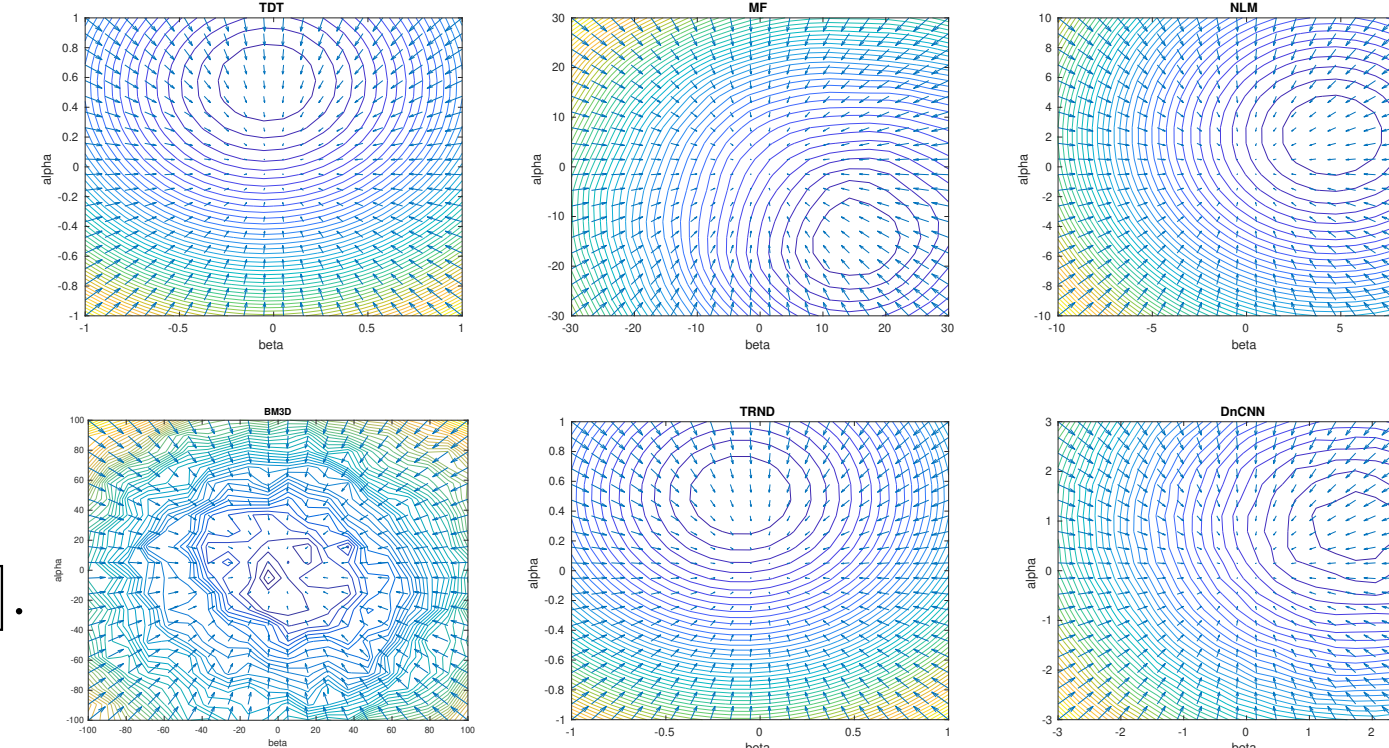
- Contours show cost:

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

- Arrows show claimed RED gradient:

$$\frac{1}{\sigma^2} (\mathbf{x}_{\alpha,\beta} - \mathbf{y})^\top [\mathbf{r}_1 \ \mathbf{r}_2] + \lambda (\mathbf{x}_{\alpha,\beta} - \mathbf{f}(\mathbf{x}_{\alpha,\beta}))^\top [\mathbf{r}_1 \ \mathbf{r}_2].$$

- Figures show that 1) zero of gradient field is not at cost minimizer, and 2) cost may not be convex!



Clarifications on the RED Gradient

In the full paper, we established that...

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

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

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

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

But practical denoisers are **not LH and JS!** And there exists no regularizer ρ_{red} for a non-JS denoiser $\mathbf{f}!$

How can we 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**, a framework first described by Hyvärinen in 2005. We explain this in 3 steps:

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

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**:

$$\tilde{p}_{\mathbf{x}}(\mathbf{x}; \nu) = \frac{1}{T} \sum_{t=1}^T \mathcal{N}(\mathbf{x}; \mathbf{x}_t, \nu \mathbf{I}) = \int_{\mathbb{R}^N} \mathcal{N}(\mathbf{r}; \mathbf{x}, \nu \mathbf{I}) \hat{p}_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \quad \text{"smoothed prior"}$$

- 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 differentiable $\ell(\cdot; \mathbf{y})$, 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 formula** (see [Robbins'56]) 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 $\mathbf{x} \sim \tilde{p}_{\mathbf{x}}$ from $\mathbf{r} = \mathbf{x} + \mathcal{N}(\mathbf{0}, \nu \mathbf{I})$.

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

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

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

Score-Matching by Denoising

- Recall $\mathbf{f}_{\text{mmse},\nu} = \arg \min_{\mathbf{f}} \mathbb{E} \{ \|\mathbf{x} - \mathbf{f}(\mathbf{r})\|^2 \}$ for $\begin{cases} \mathbf{r} = \mathbf{x} + \mathcal{N}(\mathbf{0}, \nu \mathbf{I}) \\ \mathbf{x} \sim \tilde{p}_{\mathbf{x}} \end{cases}$

- Since $\mathbf{f}_{\text{mmse},\nu}$ is expensive to implement, we typically use some approximation $\mathbf{f}_{\hat{\theta}}$ with

$$\hat{\theta} = \arg \min_{\theta} \mathbb{E} \{ \|\mathbf{x} - \mathbf{f}_{\theta}(\mathbf{r})\|^2 \} \quad \text{e.g., deep network}$$

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

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

$$= \arg \min_{\theta} \mathbb{E} \left\{ \left\| \underbrace{\nabla \ln \tilde{p}_{\mathbf{x}}(\mathbf{r}; \nu)}_{\text{"score"}} - \underbrace{\frac{1}{\nu} (\mathbf{r} - \mathbf{f}_{\theta}(\mathbf{r}))}_{\text{RED log-prior}} \right\|^2 \right\} \quad \text{via Tweedie's formula}$$

- Thus RED algorithm with general denoiser \mathbf{f}_{θ} can be interpreted as **"score matching."**

- Key points:

- RED algs 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 $-\nabla \ln \tilde{p}_{\mathbf{x}}(\mathbf{x}; \nu)$.
- This SMD interpretation holds for any $\tilde{p}_{\mathbf{x}}$, any denoiser class \mathbf{f}_{θ} (i.e., may be **non-JS** and/or **non-LH**), and any θ .
- SMD arises naturally via non-parametric estimation (i.e., KDE). Matches construction of *learned* denoisers like TNRD and DnCNN.
- Related work:
 - In 2014, Alain and Bengio showed that learned auto-encoders are explained by score-matching and *not* by minimization of an energy function.
 - In 2017, Bigdeli and Zwicker used Tweedie's formula to interpret autoencoding-based image priors.

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.

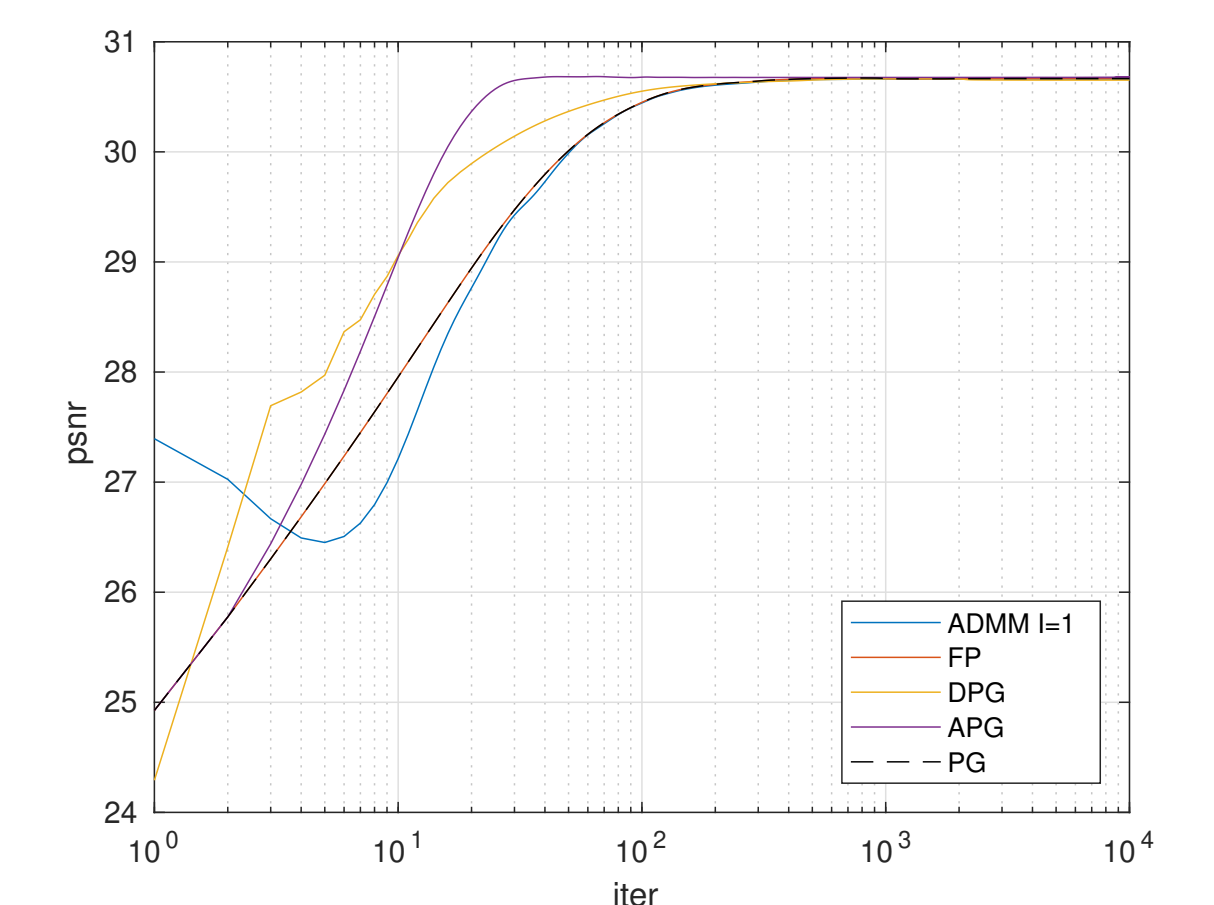
In the RED paper by Romano, Elad, and Milanfar, three algorithms were described:

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

We proposed three new algorithms:

- PG**: Proximal gradient with stepsize $L > 0$.
- DPG**: "Dynamic" proximal gradient that schedules L with the iterations.
- APG**: Accelerated proximal gradient, similar in spirit to FISTA.

In the de-blurring experiment on the right, APG is about **3x faster** than the "fixed-point" method.



Convergence to a Fixed Point

Theorem:

If $\ell(\cdot)$ is proper, convex, and continuous; $\mathbf{f}(\cdot)$ 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 the Mann iteration.

Conclusions

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

- For more details (e.g., an equilibrium analysis), please see:

E.T. Reehorst and P. Schniter, "Regularization by denoising: Clarifications and new interpretations," *IEEE Trans. Computational Imaging*, vol. 5, no. 1, pp. 52-67, Mar 2019. <http://arxiv.org/abs/1806.02296>