

Plug-and-Play: Algorithms, Parameters Tuning and Interpretation

Zhonglin Xie

Peking University

February 28, 2022

Linear Inverse Problem

- Formulation:

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

where \mathbf{A} is a known block-diagonal matrix, and \mathbf{w} denotes Gaussian random vector with mean $\mathbf{0}$ and covariance $\sigma^2 \mathbf{I}$.

- Aim: Recovery \mathbf{x} from \mathbf{y} .
- Application: Magnetic Resonance Imaging (MRI).

Signal Recovery and Denoising

- The maximum likelihood (ML) estimate:

$$\hat{\mathbf{x}}_{\text{ml}} \triangleq \underset{x}{\operatorname{argmax}} p(\mathbf{y} | \mathbf{x}),$$

where $p(\mathbf{y} | \mathbf{x})$, the probability density of \mathbf{y} conditioned on \mathbf{x} , is known as the likelihood function.

- The equivalent form:

$$\hat{\mathbf{x}}_{\text{ml}} = \underset{x}{\operatorname{argmin}} \{-\ln p(\mathbf{y} | \mathbf{x})\}.$$

- In the case of additive white Gaussian noise (AWGN) $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, we have

$$-\ln p(\mathbf{y} | \mathbf{x}) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \text{const.}$$

Maximum A Posteriori (MAP)

- Since \mathbf{A} is fat, we can not perform least-squares estimation.
- Use a regularization term $\phi(\mathbf{x})$ to encodes priori knowledge:

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \left\{ \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \phi(\mathbf{x}) \right\}.$$

- Bayes rule implies

$$\ln p(\mathbf{x} | \mathbf{y}) = \ln p(\mathbf{y} | \mathbf{x}) + \ln p(\mathbf{x}) - \ln p(\mathbf{y}).$$

- The maximum a posteriori (MAP) estimate:

$$\hat{\mathbf{x}}_{\text{map}} = \operatorname{argmin}_{\mathbf{x}} \{-\ln p(\mathbf{y} | \mathbf{x}) - \ln p(\mathbf{x})\}.$$

- $\hat{\mathbf{x}}$ can be recognized as $\hat{\mathbf{x}}_{\text{map}}$ with $p(\mathbf{x}) \propto \exp(-\phi(\mathbf{x}))$.

More on $\phi(\mathbf{x})$

- $\phi(\mathbf{x})$ should mimic the negative log of the prior density.
- $\phi(\mathbf{x})$ must enable tractable optimization.
- Common choice: $\phi(\mathbf{x}) = \lambda \|\Psi \mathbf{x}\|_1$, where $\Psi^T \Psi = \mathbf{I}$, $\lambda > 0$.
- Advantages:
 - The problem remains convex.
 - The transform output $\Psi \mathbf{x}$ is sparse.

Denoising

- When $\mathbf{A} = \mathbf{I}$, the linear inverse problem reduces to

$$\mathbf{z} = \mathbf{x} + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

Recovering \mathbf{x} from noisy \mathbf{z} is known as **denoising**.

- State-of-the-art approaches are either algorithmic or neural.
- Can these state-of-the-art denoisers be leveraged for MRI?

- $\hat{\mathbf{x}}$ can be derived from an equivalent optimization

$$(\hat{\mathbf{x}}, \hat{\mathbf{v}}) = \arg \min_{\mathbf{x}, \mathbf{v} \in \mathbb{R}^n} \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \phi(\mathbf{v}), \quad \text{s.t. } \mathbf{x} = \mathbf{v}.$$

- The augmented Lagrangian:

$$\begin{aligned} L(\mathbf{x}, \mathbf{v}; \boldsymbol{\lambda}) &= \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \phi(\mathbf{v}) - \boldsymbol{\lambda}^\top (\mathbf{x} - \mathbf{v}) + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{v}\|_2^2 \\ &= \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \phi(\mathbf{v}) + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{v} + \mathbf{u}\|_2^2, \end{aligned}$$

where $\mathbf{u} = \eta\boldsymbol{\lambda}$.

ADMM

Alternating the optimization of \mathbf{x} , \mathbf{v} with gradient ascent of \mathbf{u} :

$$\mathbf{x}_k = \mathbf{h}(\mathbf{v}_{k-1} - \mathbf{u}_{k-1}; \sigma^2/\eta)$$

$$\mathbf{v}_k = \text{prox}_\phi(\mathbf{x}_k + \mathbf{u}_{k-1}; \eta)$$

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

where

$$\text{prox}_\phi(\mathbf{z}; \eta) \triangleq \arg \min_{\mathbf{x} \in \mathbb{R}^n} \phi(\mathbf{x}) + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{z}\|^2,$$

$$\begin{aligned} \mathbf{h}(\mathbf{z}; \sigma^2/\eta) &\triangleq \arg \min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{z}\|^2 \\ &= \left(\mathbf{A}^\top \mathbf{A} + \frac{\sigma^2}{\eta} \mathbf{I} \right)^{-1} \left(\mathbf{A}^\top \mathbf{y} + \frac{\sigma^2}{\eta} \mathbf{z} \right). \end{aligned}$$

PnP-ADMM

- $\text{prox}_\phi(\mathbf{z}; \eta)$ can be **recognized** as the MAP **denoiser** of \mathbf{z} .
- PnP plug in a image denoiser in place of the $\text{prox}_\phi(\mathbf{z}; \eta)$.
- Denoting the denoiser as $\mathbf{f}(\cdot; \eta)$, PnP-ADMM writes:

$$\mathbf{x}_k = \mathbf{h}(\mathbf{v}_{k-1} - \mathbf{u}_{k-1}; \sigma^2/\eta)$$

$$\mathbf{v}_k = \mathbf{f}(\mathbf{x}_k + \mathbf{u}_{k-1}; \eta)$$

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

- The fixed point of the ADMM is independent of η , while η affects the fixed-point of the PnP-ADMM.
- To promote the PnP-ADMM, we untie the parameters:

$$\mathbf{x}_k = \mathbf{h}(\mathbf{v}_{k-1} - \mathbf{u}_{k-1}; \mu_k)$$

$$\mathbf{v}_k = \mathbf{f}(\mathbf{x}_k + \mathbf{u}_{k-1}; \eta_k)$$

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

Discussion: The Effects of the Denoising Strength η

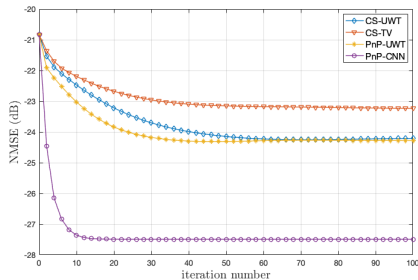
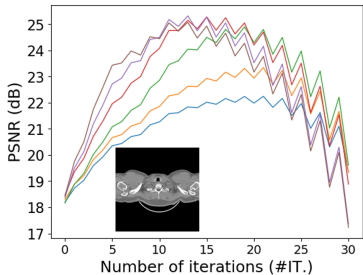
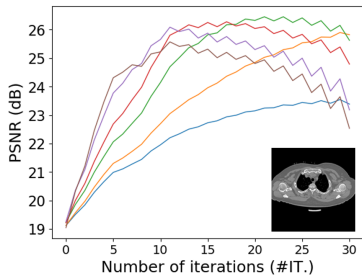
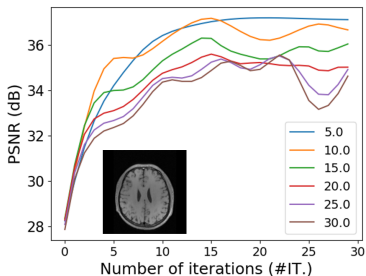
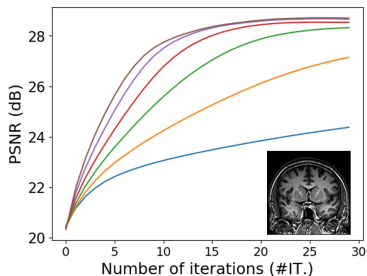


Figure: The normalized mean-squared error (NMSE) versus iteration.

- ADMM
- PnP-ADMM

The Effects of the Denoising Strength η_k



RL Formulation for Automated Parameter Selection

- Motivation: η, μ affect the result of the PnP-ADMM.
- Manually tuned parameters are time-cost.
- Aim: automatically select

$$\tau \text{ and } (\eta_0, \mu_0, \eta_1, \mu_1, \dots, \eta_{\tau-1}, \mu_{\tau-1})$$

to recover \mathbf{x}_τ that close to \mathbf{x} .

- Tool: Reinforcement Learning (RL).

Markov decision process (MDP) $(\mathcal{S}, \mathcal{A}, p, r)$

- State space \mathcal{S} : any feasible value of $(\mathbf{x}_k, \mathbf{v}_k, \mathbf{u}_k)$.
- Action space \mathcal{A} : any feasible value of τ and (μ_k, η_k) .
- Transition function $p: \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$

several iterations of the PnP-ADMM.

- Reward function $r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$

$$r(s_t, a_t) = [\zeta(p(s_t, a_t)) - \zeta(s_t)] - \eta.$$

$\zeta(s_t)$: the PSNR of the recovered image at step t .

η : penalizing the policy as it does not terminate at step t .

Workflow of TFPnP

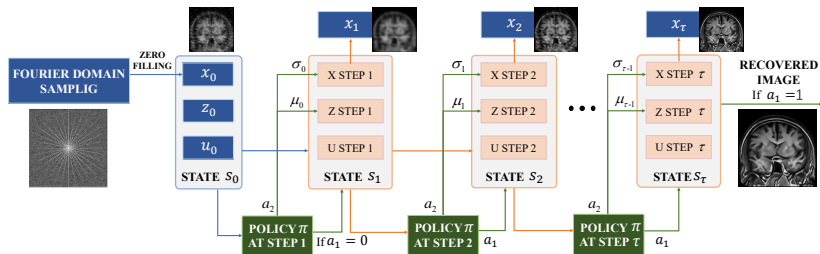


Figure: Workflow of the TFPnP instantiated by the PnP-ADMM.

Formal Definition of the Goal

- $s_k = (\mathbf{x}_k, \mathbf{v}_k, \mathbf{u}_k)$, $a_k = (a_{k,1}, a_{k,2})$, $r_k = r(s_k, a_k)$.
- $a_{k,2} = (\mu_k, \eta_k)$. $a_{k,1}$ is a boolean that terminates the iterate at step k when $a_{k,1} = 0$ and versus verse.
- Trajectory: $T = \{s_0, a_0, r_0, \dots, s_N, a_N, r_N\}$.
- Given T and $\rho \in [0, 1]$, define the return as

$$R_t = \sum_{t'=0}^{N-t} \rho^{t'} r(s_{t+t'}, a_{t+t'}).$$

- Goal: learn a policy π to maximize

$$J(\pi) = \mathbb{E}_{s_0, \pi} [R_0], \quad \pi(a | s) : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1].$$

RL-based Policy Learning

- State-value function:

$$V^\pi(s) = \mathbb{E}_\pi [R_0 \mid s_0 = s]$$

- Action-value function:

$$Q^\pi(s, a) = \mathbb{E}_\pi [R_0 \mid s_0 = s, a_0 = a]$$

- Policy: $\pi = (\pi_1, \pi_2)$.

π_1 : a **stochastic** policy that generate $a_{t,1}$ to decide whether to terminate.

π_2 : a **deterministic** policy that generate $a_{t,2}$.

Actor-critic Framework

- Policy network (actor): $\pi_\theta = (\pi_1, \pi_2)$ with $\theta = (\theta_1, \theta_2)$.
 $\pi_1(\cdot | s): \mathcal{S} \times \{0, 1\} \rightarrow [0, 1]$, controlled by θ_1 .
 $\pi_2(s): \mathcal{S} \rightarrow \mathcal{A}$, controlled by θ_2 .
- Value network (critic): $V_\phi^\pi(s_t)$.
- Train the value network:

$$L_\phi = \mathbb{E}_{s \sim B, a \sim \pi_\theta(s)} \left[\frac{1}{2} \left(r(s, a) + \gamma V_{\hat{\phi}}^\pi(p(s, a)) - V_\phi^\pi(s) \right)^2 \right]$$

- Model-free training of π_1 :

$$\nabla_{\theta_1} J(\pi_\theta) = \mathbb{E}_{s \sim B, a \sim \pi_\theta(s)} [\nabla_{\theta_1} \log \pi_1(a_1 | s) A^\pi(s, a)]$$

- Model-based training of π_2 :

$$\nabla_{\theta_2} J(\pi_\theta) = \mathbb{E}_{s \sim B, a \sim \pi_\theta(s)} [\nabla_{a_2} Q^\pi(s, a) \nabla_{\theta_2} \pi_2(s)]$$

Training Scheme

Algorithm 1 Training Scheme

Require: Image dataset D , degradation operator $g(\cdot)$, learning rates l_θ, l_ϕ , weight parameter β .

1: Initialize network parameters $\theta, \phi, \hat{\phi}$ and state buffer B .

2: **for** each training iteration **do**

3: sample initial state s_0 from D via $g(\cdot)$

4: **for** environment step $t \in [0, N)$ **do**

5: $a_t \sim \pi_\theta(a_t|s_t)$

6: $s_{t+1} \sim p(s_{t+1}|s_t, a_t)$

7: $B \leftarrow B \cup \{s_{t+1}\}$

8: break if the boolean outcome of a_t equals to 1

9: **end for**

10: **for** each gradient step **do**

11: sample states from the state buffer B

12: $\theta_1 \leftarrow \theta_1 + l_\theta \nabla_{\theta_1} J(\pi_\theta)$

13: $\theta_2 \leftarrow \theta_2 + l_\theta \nabla_{\theta_2} J(\pi_\theta)$

14: $\phi \leftarrow \phi - l_\phi \nabla_\phi L_\phi$

15: $\hat{\phi} \leftarrow \beta \phi + (1 - \beta) \hat{\phi}$

16: **end for**

17: **end for**

Ensure: Learned policy network π_θ

Experiment Results: CS-MRI

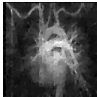
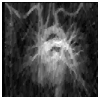


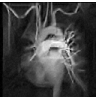
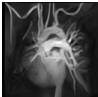

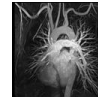
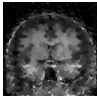
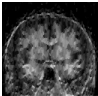
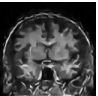
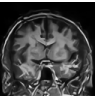
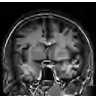
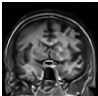
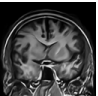
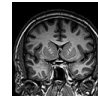

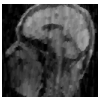
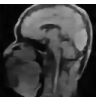
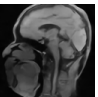
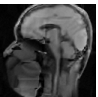
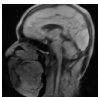
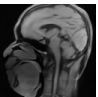
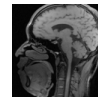
RecPF	FCSA	ADMMNet	ISTANet	BM3D-MRI	IRCNN	Ours	Ground Truth
							
22.57	22.27	24.15	24.61	23.64	24.16	25.28	PSNR
							
18.74	19.23	20.48	21.37	20.62	20.91	22.02	PSNR
							
24.89	24.47	26.85	27.90	26.72	27.74	28.65	PSNR

Figure: Visual and numerical CS-MRI reconstruction comparison against the state-of-the-art techniques on medical images. The numerical values denote the PSNR obtained by each technique.

Rethinking of the PnP-ADMM: On Derivation

- \mathbf{f} is not the proximal map of any regularizer ϕ .
- \mathbf{f} coincides with $\text{prox}_\phi(\mathbf{z}; \eta)$ **only when**

$$p(\mathbf{x}) \propto \exp(-\phi(\mathbf{x})), \quad \mathbf{z} - \mathbf{z}_{\text{true}} \sim \mathcal{N}(\mathbf{0}, \eta^2 \mathbf{I}).$$

- However, $p(\mathbf{x})$ may **not** prompt to $\exp(-\phi(\mathbf{x}))$ and the distribution of

$$(\mathbf{x}_k + \mathbf{u}_{k-1}) - (\mathbf{x}_k + \mathbf{u}_{k-1})_{\text{true}}$$

is **unknown!**

- PnP-ADMM is a result of the similarity of the formulation.

Rethinking of the PnP-ADMM: On Convergence

- PnP-ADMM may not be an implementation of ADMM.
- If the PnP-ADMM converges?
- If it does converge, what it converges to?

PnP FISTA

$$\min_{\mathbf{x}} \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \phi(\mathbf{x})$$

$$\mathbf{z}_k = \mathbf{s}_{k-1} - \frac{\eta}{\sigma^2} \mathbf{A}^\top (\mathbf{A}\mathbf{s}_{k-1} - \mathbf{y}) \quad \mathbf{z}_k = \mathbf{s}_{k-1} - \frac{\eta}{\sigma^2} \mathbf{A}^\top (\mathbf{A}\mathbf{s}_{k-1} - \mathbf{y})$$

$$\mathbf{x}_k = \text{prox}_\phi(\mathbf{z}_k; \eta) \quad \mathbf{x}_k = \mathbf{f}(\mathbf{z}_k)$$

$$\mathbf{s}_k = \mathbf{x}_k + \frac{q_{k-1} - 1}{q_k} (\mathbf{x}_k - \mathbf{x}_{k-1}) \quad \mathbf{s}_k = \mathbf{x}_k + \frac{q_{k-1} - 1}{q_k} (\mathbf{x}_k - \mathbf{x}_{k-1})$$

FISTA

PnP FISTA

where it is typical to use $q_k = \left(1 + \sqrt{1 + 4q_{k-1}^2}\right) / 2$ and $q_0 = 1$ with step-size $\eta \in (0, \sigma^2 \|A\|_2^{-2})$.

Regularization by Denoising (RED)

- Recover \mathbf{x} from measurements \mathbf{y} by solving

$$\mathbf{0} = \frac{1}{\sigma^2} \mathbf{A}^\top (\mathbf{A} \hat{\mathbf{x}} - \mathbf{y}) + \frac{1}{\eta} (\hat{\mathbf{x}} - \mathbf{f}(\hat{\mathbf{x}})).$$

- \mathbf{f} is an arbitrary image denoiser.
- When \mathbf{f} is a sophisticated denoiser and η is well tuned, the solutions $\hat{\mathbf{x}}$ are state-of-the-art.

RED: Assumptions

Define

$$\rho_{\text{RED}}(\mathbf{x}) \triangleq \frac{1}{2} \langle \mathbf{x}, \mathbf{x} - \mathbf{f}(\mathbf{x}) \rangle, \quad \ell(\mathbf{x}; \mathbf{y}) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2.$$

We get

$$\hat{\mathbf{x}}_{\text{RED}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \ell(\mathbf{x}; \mathbf{y}) + \rho_{\text{RED}}(\mathbf{x}).$$

The denoiser $\mathbf{f}(x)$ obeys the following assumption:

1 Local Homogeneity:

$$\mathbf{f}((1 + \varepsilon)\mathbf{x}) = (1 + \varepsilon)\mathbf{f}(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^n, 0 < \varepsilon \ll 1.$$

2 $\mathbf{f}(\cdot)$ is differentiable where $J\mathbf{f} \in \mathbb{R}^{n \times n}$ denotes its Jacobian.

3 **Jacobian Symmetry:** $J\mathbf{f}(\mathbf{x})^\top = J\mathbf{f}(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^n.$

4 The spectral radius the Jacobian satisfies $\eta(J\mathbf{f}(\mathbf{x})) \leq 1.$

RED: Proof

- From the multivariate calculus:

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

- Local homogeneity implies $[J\mathbf{f}(\mathbf{x})]\mathbf{x} = \mathbf{f}(\mathbf{x})$:

$$\begin{aligned} 0 &= \lim_{\varepsilon \rightarrow 0} \frac{\|\mathbf{f}(\mathbf{x} + \varepsilon \mathbf{x}) - \mathbf{f}(\mathbf{x}) - [J\mathbf{f}(\mathbf{x})]\mathbf{x}\varepsilon\|}{\|\varepsilon \mathbf{x}\|} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{\|(1 + \varepsilon)\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}) - [J\mathbf{f}(\mathbf{x})]\mathbf{x}\varepsilon\|}{\|\varepsilon \mathbf{x}\|} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{\|\mathbf{f}(\mathbf{x}) - [J\mathbf{f}(\mathbf{x})]\mathbf{x}\|}{\|\mathbf{x}\|}. \end{aligned}$$

- Jacobian symmetry gives $\nabla \rho_{\text{RED}}(\mathbf{x}) = \mathbf{x} - \mathbf{f}(\mathbf{x})$.
- $\eta(J\mathbf{f}(\mathbf{x})) \leq 1$ guarantees the convexity.

RED: Remark

- When the denoiser $\mathbf{f}(\cdot)$ is locally homogeneous, then

$$\nabla \rho_{\text{RED}}(\mathbf{x}) = \mathbf{x} - \mathbf{f}(\mathbf{x}) \quad \Leftrightarrow \quad J\mathbf{f}(\mathbf{x}) = [J\mathbf{f}(\mathbf{x})]^\top.$$

- When $J\mathbf{f}(\cdot) \neq J\mathbf{f}(\cdot)^\top$, there exists no regularizer $\rho(\cdot)$ for which $\nabla \rho(\mathbf{x}) = \mathbf{x} - \mathbf{f}(\mathbf{x})$.
- Many popular denoisers lack symmetric Jacobian, making the gradient expression invalid.

Proximal-based PnP v.s. RED: a Toy Example

- $f(z) = \mathbf{W}z$ with $\mathbf{W} = \mathbf{W}^\top$.
- f is the proximal map of $\phi(\mathbf{x}) = (1/2\eta)\mathbf{x}^\top (\mathbf{W}^{-1} - \mathbf{I}) \mathbf{x}$.
- Proximal-based PnP:

$$\hat{\mathbf{x}}_{\text{pnp}} = \underset{\mathbf{x}}{\operatorname{argmin}} \left\{ \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \frac{1}{2\eta} \mathbf{x}^\top (\mathbf{W}^{-1} - \mathbf{I}) \mathbf{x} \right\}$$

- RED:

$$\hat{\mathbf{x}}_{\text{red}} = \underset{\mathbf{x}}{\operatorname{argmin}} \left\{ \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \frac{1}{2\eta} \mathbf{x}^\top (\mathbf{I} - \mathbf{W}) \mathbf{x} \right\}$$

Algorithms for RED

- GD, inexact ADMM, and a “fixed-point” heuristic that was later recognized as a special case of the proximal gradient (PG) algorithm.
- Accelerated proximal gradient (fastest):

$$\begin{aligned}\mathbf{x}_k &= \mathbf{h}(\mathbf{v}_{k-1}; \eta/L) \\ \mathbf{z}_k &= \mathbf{x}_k + \frac{q_{k-1} - 1}{q_k} (\mathbf{x}_k - \mathbf{x}_{k-1}) \\ \mathbf{v}_k &= \frac{1}{L} \mathbf{f}(\mathbf{z}_k) + \left(1 - \frac{1}{L}\right) \mathbf{z}_k\end{aligned}$$

where $L > 0$ is a design parameter that can be related to the Lipschitz constant of $\phi_{\text{red}}(\cdot)$.

RED as Score Matching

- Given a training set $\{\mathbf{x}_t\}_{t=1}^T$, the empirical prior model is

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

- Build a prior model using kernel density estimation (KDE):

$$\tilde{p}(\mathbf{x}; \eta) \triangleq \frac{1}{T} \sum_{t=1}^T \mathcal{N}(\mathbf{x}; \mathbf{x}_t, \eta \mathbf{I})$$

- Adopting \tilde{p} as the prior, MAP becomes

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 - \ln \tilde{p}(\mathbf{x}; \eta)$$

RED as Score Matching

- Because $\ln \tilde{p}$ is differentiable, $\hat{\mathbf{x}}$ must obey

$$\mathbf{0} = \frac{1}{\sigma^2} \mathbf{A}^\top (\mathbf{A}\hat{\mathbf{x}} - \mathbf{y}) - \nabla \ln \tilde{p}(\hat{\mathbf{x}}; \eta)$$

- $\mathbf{f}_{\text{mmse}}(\mathbf{z}; \eta) = \mathbb{E}[\mathbf{x}|\mathbf{z}]$, where $\mathbf{z} = \mathbf{x} + \mathcal{N}(\mathbf{0}, \eta \mathbf{I})$, $\mathbf{x} \sim \hat{p}$
- Tweedie's formula says that

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

- The MAP estimate $\hat{\mathbf{x}}$ under the KDE prior \tilde{p} obeys

$$\mathbf{0} = \frac{1}{\sigma^2} \mathbf{A}^\top (\mathbf{A}\hat{\mathbf{x}} - \mathbf{y}) + \frac{1}{\eta} (\hat{\mathbf{x}} - \mathbf{f}_{\text{mmse}}(\hat{\mathbf{x}}; \eta))$$

which matches the RED condition when $\mathbf{f} = \mathbf{f}_{\text{mmse}}(\cdot; \eta)$

RED as Score Matching: $f \neq f_{\text{mmse}}(\cdot; \eta)$

- f_{θ} : neural denoiser parameterized by θ
- Training strategy:

$$\min_{\theta} \mathbb{E} \|\mathbf{x} - f_{\theta}(\mathbf{z})\|^2, \quad \text{where } \mathbf{x} \sim \hat{p}, \quad \mathbf{z} = \mathbf{x} + \mathcal{N}(\mathbf{0}, \eta \mathbf{I})$$

- MMSE orthogonality principle:

$$\begin{aligned} \mathbb{E} \|\mathbf{x} - f_{\theta}(\mathbf{z})\|^2 &= \mathbb{E} \|\mathbf{x} - f_{\text{mmse}}(\mathbf{z}; \eta)\|^2 \\ &\quad + \mathbb{E} \|f_{\text{mmse}}(\mathbf{z}; \eta) - f_{\theta}(\mathbf{z})\|^2 \end{aligned}$$

- Using Tweedie's formula, we get

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmin}} \mathbb{E} \|\mathbf{x} - f_{\theta}(\mathbf{z})\|^2 \\ &= \underset{\theta}{\operatorname{argmin}} \mathbb{E} \|f_{\text{mmse}}(\mathbf{z}; \eta) - f_{\theta}(\mathbf{z})\|^2 \\ &= \underset{\theta}{\operatorname{argmin}} \mathbb{E} \left\| \nabla \ln \tilde{p}(\mathbf{z}; \eta) - \frac{1}{\eta} (f_{\theta}(\mathbf{z}) - \mathbf{z}) \right\|^2 \end{aligned}$$

- Choose θ so that $(f_{\theta}(\mathbf{z}) - \mathbf{z})/\eta$ matches the “score” $\nabla \ln \tilde{p}$

CE for Prox-based PnP

- View Prox-based PnP as seeking a solution to

$$\hat{\mathbf{x}}_{\text{pnp}} = \mathbf{h}(\hat{\mathbf{x}}_{\text{pnp}} - \hat{\mathbf{u}}_{\text{pnp}}; \eta)$$

$$\hat{\mathbf{x}}_{\text{pnp}} = \mathbf{f}(\hat{\mathbf{x}}_{\text{pnp}} + \hat{\mathbf{u}}_{\text{pnp}})$$

- It equals to find a fixed point of

$$\underline{z} = (2\mathbf{G} - \mathbf{I})(2\mathcal{F} - \mathbf{I})\underline{z}$$

$$\underline{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}, \quad \mathcal{F}(\underline{z}) = \begin{bmatrix} \mathbf{h}(z_1; \eta) \\ \mathbf{f}(z_2) \end{bmatrix}, \quad \mathcal{G}(\underline{z}) = \begin{bmatrix} (z_1 + z_2) / 2 \\ (z_1 + z_2) / 2 \end{bmatrix}$$

- Mann iteration writes:

$$\underline{z}^{(k+1)} = (1 - \gamma)\underline{z}^{(k)} + \gamma(2\mathbf{G} - \mathbf{I})(2\mathcal{F} - \mathbf{I})\underline{z}^{(k)}$$

CE for RED

- CE for ADMM-based RED:

$$\hat{\mathbf{x}}_{\text{red}} = \mathbf{h}(\hat{\mathbf{x}}_{\text{red}} - \hat{\mathbf{u}}_{\text{red}}; \eta)$$

$$\hat{\mathbf{x}}_{\text{red}} = \left(\left(1 + \frac{1}{L} \right) \mathbf{I} - \frac{1}{L} \mathbf{f} \right)^{-1} (\hat{\mathbf{x}}_{\text{red}} + \hat{\mathbf{u}}_{\text{red}})$$

- A more intuitive form:

$$\hat{\mathbf{x}}_{\text{red}} = \mathbf{h}(\hat{\mathbf{x}}_{\text{red}} - \hat{\mathbf{u}}_{\text{red}}; \eta)$$

$$\hat{\mathbf{x}}_{\text{red}} = \mathbf{f}(\hat{\mathbf{x}}_{\text{red}}) + L\hat{\mathbf{u}}_{\text{red}}$$

- Solving the first equation gives:

$$\hat{\mathbf{u}}_{\text{red}} = \frac{\eta}{\sigma^2} \mathbf{A}^\top (\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}_{\text{red}})$$

- Plugging $\hat{\mathbf{u}}_{\text{red}}$ back:

$$\frac{L\eta}{\sigma^2} \mathbf{A}^\top (\mathbf{A}\hat{\mathbf{x}}_{\text{red}} - \mathbf{y}) = \mathbf{f}(\hat{\mathbf{x}}_{\text{red}}) - \hat{\mathbf{x}}_{\text{red}}$$

RED via Fixed-point Projection (RED-PRO)

RED-PRO problem writes:

$$\hat{\mathbf{x}}_{\text{RED-PRO}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \ell(\mathbf{x}; \mathbf{y}), \quad \text{s.t. } \mathbf{x} \in \text{Fix}(\mathbf{f}).$$

- Interpretation: searching for a minimizer of $\ell(\mathbf{x}; \mathbf{y})$ over the set of “clean” images.
- The manifold of natural images \mathcal{M} is generally not well-defined, it is not easy accessible and it is not convex, making the search within this domain difficult. Therefore, as an alternative, we propose to use $\text{Fix}(f)$ which is well-behaved for demicontractive denoisers and should satisfy $\mathcal{M} \subset \text{Fix}(f)$ for a “perfect” denoiser.
- Common denoisers are far from being ideal, hence, the solution is sensitive to the choice of the denoiser and it may vary considerably for different choices.

d -demicontractive Mapping

A mapping T is d -demicontractive ($d \in [0, 1)$) if for any $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{z} \in \text{Fix}(T)$ it holds that

$$\|T(\mathbf{x}) - \mathbf{z}\|^2 \leq \|\mathbf{x} - \mathbf{z}\|^2 + d\|T(\mathbf{x}) - \mathbf{x}\|^2$$

or equivalently

$$\frac{1-d}{2}\|\mathbf{x} - T(\mathbf{x})\|^2 \leq \langle \mathbf{x} - T(\mathbf{x}), \mathbf{x} - \mathbf{z} \rangle$$

RED-PRO

- Assume the denoiser $\mathbf{f}(\cdot)$ is a d -demicontractive mapping. Then, RED-PRO defines a convex minimization problem.
- Consider a demicontractive denoiser $\mathbf{f}(\cdot)$ and assume $\mathbf{f}(0) = 0$. Then,

$$\rho_{\text{RED}}(\mathbf{x}) = \frac{1}{2} \langle \mathbf{x}, \mathbf{x} - \mathbf{f}(\mathbf{x}) \rangle = 0 \text{ iff } \mathbf{x} \in \text{Fix}(\mathbf{f}).$$

- Hybrid steepest descent method for RED-PRO:

$$\begin{aligned}\mathbf{v}_{k+1} &= \mathbf{x}_k - \mu_k \nabla \ell(\mathbf{x}_k; \mathbf{y}), \\ \mathbf{z}_{k+1} &= \mathbf{f}(\mathbf{v}_{k+1}), \\ \mathbf{x}_{k+1} &= (1 - \alpha)\mathbf{v}_{k+1} + \alpha\mathbf{z}_{k+1},\end{aligned}$$

- which is equivalent to

$$\mathbf{x}_{k+1} = f_\alpha(\mathbf{x}_k - \mu_k \nabla \ell(\mathbf{x}_k; \mathbf{y})), \text{ where } f_\alpha = (1 - \alpha)\text{Id} + \alpha f.$$

Uniform Algorithm Framework

- accelerated-PG (proximal gradient) RED algorithm, which uses the iterative update:

$$\mathbf{v}_{k+1} = \mathbf{x}_k - \mu_k \nabla \ell(\mathbf{x}_k; \mathbf{y}),$$

$$\mathbf{z}_{k+1} = \mathbf{v}_{k+1} + q_k (\mathbf{v}_{k+1} - \mathbf{v}_k), \quad (\text{FISTA-like acceleration})$$

$$\mathbf{x}_{k+1} = (1 - \alpha) \mathbf{z}_{k+1} + \alpha f(\mathbf{z}_{k+1}), \quad (\text{SOR-like acceleration})$$

- Thus, when we set $q_k \equiv 0$, i.e. when we skip the acceleration step, the above RED variant reduces to the iterative update of the Hybrid steepest for RED-PRO.
- When we continue and set $\alpha = 1$, we obtain the PnP-PGD method (Proximal-based).

Projection Gradient Descent

- Projected Gradient Descent writes

$$\mathbf{x}_{k+1} = P_{\text{Fix}(f)}(\mathbf{x}_k - \mu_k \nabla \ell(\mathbf{x}; y))$$

- Replacing the projection operator $P_{\text{Fix}(f)}(\cdot)$ with denoiser (Plug and Play) $f(\cdot)$ we get PnP-PGD:

$$\mathbf{x}_{k+1} = f(\mathbf{x}_k - \mu_k \nabla \ell(\mathbf{x}; y))$$

Convergence Theorem

Let $f(\cdot)$ be a continuous d -demicontractive denoiser and $\ell(\cdot; \mathbf{y})$ be a proper convex lower semicontinuous differentiable function with L -Lipschitz gradient $\nabla\ell(\cdot; \mathbf{y})$. Assume the following:

$$(A1) \quad \alpha \in (0, \frac{1-d}{2}).$$

$$(A2) \quad \{\mu_k\}_{k \in \mathbb{N}} \subset [0, \infty) \text{ where } \mu_k \xrightarrow[k \rightarrow \infty]{} 0 \text{ and } \sum_{k \in \mathbb{N}} \mu_k = \infty.$$

Then, the sequence $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ generated by

$$\mathbf{x}_{k+1} = f_\alpha(\mathbf{x}_k - \mu_k \nabla\ell(\mathbf{x}_k; \mathbf{y})), \text{ where } f_\alpha = (1 - \alpha)\text{Id} + \alpha f,$$

converges to an optimal solution of the RED-PRO problem:

$$\hat{\mathbf{x}}_{\text{RED-PRO}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \ell(\mathbf{x}; \mathbf{y}), \quad \text{s.t. } \mathbf{x} \in \text{Fix}(f).$$

Conclusion

There are various ways to model the denoising problem:

- 1 PnP: Inspired by ADMM, Proximal gradient, while lacking objective function.
- 2 RED: Regularization by Denoising, while many denoisers do not satisfy the assumptions.
- 3 RED-PRO: require the denoisers to be demicontractive.

However, as pointed by, when applying practical algorithms (e.g. PnP-ADMM and PnP primal-dual hybrid gradient method (PnP-PDHG) , satisfy the same fixed-point equation as PnP-PGM (Proximal Gradient Method)) to solve these models, different models have the same aim:

$$\mathbf{x}_* = f_\alpha(\mathbf{x}_* - \mu_k \nabla \ell(\mathbf{x}_*; \mathbf{y})), \text{ where } f_\alpha = (1 - \alpha)\text{Id} + \alpha f.$$

Thus, we only need to guarantee the convergence of the above formulation.

Future Directions

- RL for general parameters tuning
- The convergence theory of the PnP with weaker assumptions
- PnP for general ADMM-based algorithms