

Optimization based machine learning for computational imaging

Optimization plays a central role in modern *machine learning* (ML). The most obvious application of optimization in ML is training of deep neural networks, where large-scale algorithms can be used to optimize over billions of parameters and data items. These lecture focuses on a complementary application, where the guiding idea is to adopt concepts from optimization to design ML methods for *computational imaging*.

Computational imaging is a rapidly growing field that seeks to overcome traditional limits of imaging instruments by viewing imaging as an *inverse problem*. Modern computational imaging algorithms achieve impressive results by exploiting complex mathematical models that characterize the physics of the imaging instrument (*forward model*) as well as prior knowledge on the class of desired images (*image prior*). Advances in computational imaging are having transformative effects across a wide range of scientific, engineering, and biomedical applications, including 3D live-cell imaging, structural analysis of complex materials, early diagnosis of neurodegenerative diseases, and improved patient comfort in radiology.

Our focus will be on optimization-based ML methods for solving inverse problems in computational imaging. This remains a very active area of research with many innovations emerging over the past few years. We will focus our exploration on a highly influential class of methods known as *plug-and-play priors* (PnP) (see [1] for a recent tutorial). What makes PnP interesting is that it is versatile, mathematically elegant, and state-of-the-art in many imaging applications. PnP is also easily relatable to other ML methods in computational imaging, including *deep unfolding* (DU), *deep equilibrium models* (DEQ), and diffusion models.

Inverse problems in computational imaging

Let us start by introducing the mathematical notation that we will use throughout these lectures. Our goal is to recover an *unknown* image $\mathbf{x} \in \mathbb{R}^n$ from its noisy measurements $\mathbf{y} \in \mathbb{R}^m$. We will denote the relationship between the image and its measurements as

$$\mathbf{y} = \mathbf{A}(\mathbf{x}) + \mathbf{e}, \quad (1)$$

where \mathbf{A} is the *measurement operator* that models the physics of the imaging system and \mathbf{e} is the *noise*. We will refer to (1) as the *forward model* or the *measurement model*. When the measurement operator is linear, it can be written as a measurement matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$.

Remark. Let us note that the measurement matrices \mathbf{A} are often too large to store or manipulate directly. To see this, consider an inverse problem of recovering a 1MP image (1024×1024 pixels) from its blurry observation. The corresponding matrix \mathbf{A} would have about 10^{12} entries. In practice, we do not need to store \mathbf{A} directly by instead implementing the effect of blur as a convolution.

Examples. Let us consider two examples.

- (a) *Image denoising* is an inverse problem that will play a central role in the next chapter. It corresponds to $\mathbf{A} = \mathbf{I}$, where \mathbf{I} is the identity matrix. Denoising seeks to recover two vectors \mathbf{x} and \mathbf{e} from one vector $\mathbf{s} = \mathbf{x} + \mathbf{e}$, which highlights the necessity of prior information on \mathbf{x} and \mathbf{e} .
- (b) In many applications, we have $m < n$, corresponding to an underdetermined system in (1). For example, in image super-resolution, the low-resolution observation \mathbf{y} corresponds to a subsampled version of \mathbf{x} . Since there are ∞ many high-resolution images consistent with \mathbf{y} , prior on \mathbf{x} is essential for the recovery of a meaningful solution.

Denoising networks as image priors

We would like to examine the following denoising problem

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

where the goal is to recover x from its noisy observation s . The problem (2) assumes that the image x is sampled from some probability distribution p_x and w is an *additive white Gaussian noise (AWGN)* vector. We are interested in denoising not for denoising, but due to the observation that a deep network pre-trained as a denoiser can give access to the prior p_x without requiring explicit knowledge of p_x (see Figure 1).

Proximal operator

A key mathematical concept that will enable us to easily deploy denoisers as priors is that of the *proximal operator*. The proximal operator, also known as *proximal mapping*, *proximity operator*, *Moreau proximal operator*, or simply *prox*, was originally introduced by Moreau [2].

Definition 1. The *proximal operator* of $h : \mathbb{R}^n \rightarrow (-\infty, \infty]$ is defined as

$$\text{prox}_h(z) = \arg \min_{x \in \mathbb{R}^n} \left\{ \frac{1}{2} \|x - z\|_2^2 + h(x) \right\}, \quad z \in \mathbb{R}^n.$$

The proximal operator thus maps a vector $z \in \mathbb{R}^n$ into a subset of \mathbb{R}^n . The example below show that the proximal operator is *not necessarily unique*. It can be *empty*, *multivalued*, or a *singleton*.

Example. Consider the function

$$h(x) = \begin{cases} 0 & \text{if } x \neq 0 \\ -2 & \text{if } x = 0. \end{cases}$$

The corresponding proximal operator returns the minimizers of the following function

$$\varphi(x) = \frac{1}{2}(x - z)^2 + h(x) = \begin{cases} \frac{1}{2}(x - z)^2 & \text{if } x \neq 0 \\ \frac{1}{2}z^2 - 2 & \text{if } x = 0, \end{cases} \Rightarrow \text{prox}_h(z) = \begin{cases} \{0\} & \text{if } |z| < 2 \\ \{z\} & \text{if } |z| > 2 \\ \{0, z\} & \text{if } |z| = 2. \end{cases}$$

To see this, note that for $z = 0$, the minimizer of φ is clearly $x = 0$. For $z \neq 0$, the minimum of $(x - z)^2/2$ over $\mathbb{R} \setminus \{0\}$ is attained at $x = z \neq 0$ with a minimum value 0. Thus, if $z^2/2 - 2 < 0$, then the unique minimizer of φ is $x = 0$. If $z^2/2 - 2 > 0$, then the unique minimizer of φ is $x = z$. Finally, when $z^2/2 - 2 = 0$, then 0 and z are the two minimizers of p .

The next result establishes conditions for the existence of prox .

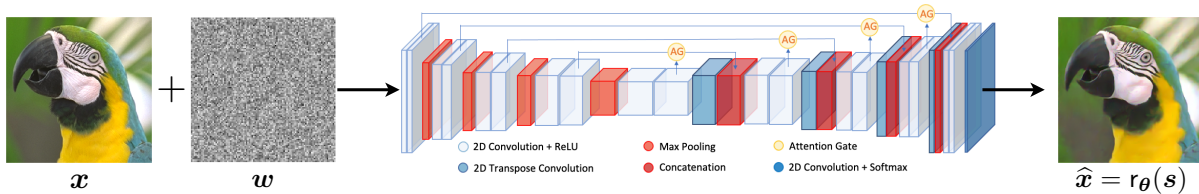


Figure 1: Priors in modern PnP methods are typically obtained by training a deep network as a denoiser.

Proposition 1. Let $h : \mathbb{R}^n \rightarrow (-\infty, \infty]$ be a proper and closed function. Assume that the function

$$\varphi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x} - \mathbf{z}\|_2^2 + h(\mathbf{x}),$$

is coercive for any $\mathbf{z} \in \mathbb{R}^n$. Then, $\text{prox}_h(\mathbf{z})$ is nonempty for any $\mathbf{z} \in \mathbb{R}^n$.

Remark. As a reminder, we say that φ is coercive if

$$\lim_{\|\mathbf{x}\|_2 \rightarrow \infty} \varphi(\mathbf{x}) = +\infty.$$

This means that

$$\forall \varphi_0 \in \mathbb{R}, \exists r > 0 \quad \text{s.t.} \quad \varphi(\mathbf{x}) > \varphi_0 \quad \text{whenever} \quad \|\mathbf{x}\| > r.$$

Thus, a coercive function φ becomes $+\infty$ for any path for which $\|\mathbf{x}\|_2$ becomes $+\infty$. A classical result in optimization states that proper, closed, and coercive functions always have a minimizer.

Proof. To prove Theorem 1, note that for any $\mathbf{z} \in \mathbb{R}^n$, φ is closed as a sum of two closed functions. Since φ is coercive, we know that $\text{prox}_f(\mathbf{z})$, which consists of minimizers of φ , is nonempty. ■

Proximal operators are commonly used for proper, closed, and convex functions $f \in \Gamma^0(\mathbb{R}^n)$. One reason is that the proximal operator of such functions always exists and is unique.

Proposition 2. Let $h \in \Gamma^0(\mathbb{R}^n)$, then prox_h is a singleton for any $\mathbf{z} \in \mathbb{R}^n$.

Remark. We will use the standard notation $\Gamma^0(\mathbb{R}^n)$ to denote the class of proper, closed, and convex functions. When the proximal operator is a singleton, we drop the set notation $\{\cdot\}$ from its output.

Proof. The function

$$\varphi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x} - \mathbf{z}\|_2^2 + h(\mathbf{x}).$$

is closed and strongly convex, since it is a sum of a closed strongly convex quadratic and closed convex h . Since h is proper, we know that φ is also proper. Since a strongly convex function always has a unique minimizer, we know that there exists a unique minimizer of φ . ■

Example: Quadratic and indicator functions

(a) The proximal operator of a quadratic function has a closed form expression

$$h(\mathbf{x}) = \frac{\tau}{2}\|\mathbf{x}\|_2^2 \quad \Rightarrow \quad \text{prox}_h(\mathbf{z}) = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2}\|\mathbf{x} - \mathbf{z}\|_2^2 + \frac{\tau}{2}\|\mathbf{x}\|_2^2 \right\} = \frac{\mathbf{z}}{1 + \tau}.$$

(b) The proximal operator of an indicator function of a nonempty, closed, and convex set $\mathcal{X} \subseteq \mathbb{R}^n$ is the *projection*

$$\begin{aligned} h(\mathbf{x}) = \mathbb{1}_{\mathcal{X}}(\mathbf{x}) \quad \Rightarrow \quad \text{prox}_h(\mathbf{z}) &= \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2}\|\mathbf{x} - \mathbf{z}\|_2^2 + \mathbb{1}_{\mathcal{X}}(\mathbf{x}) \right\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \frac{1}{2}\|\mathbf{x} - \mathbf{z}\|_2^2 \right\} = \text{proj}_{\mathcal{X}}(\mathbf{z}). \end{aligned}$$

The proximal operator has a number of important properties related to the fact that it solves an optimization problem. In particular, it can be fully characterized using a subdifferential.

Proposition 3. For $h \in \Gamma^0(\mathbb{R}^n)$, the following relationships are equivalent for any $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$

$$\mathbf{x} = \text{prox}_h(\mathbf{z}) \Leftrightarrow (\mathbf{z} - \mathbf{x}) \in \partial h(\mathbf{x}) \Leftrightarrow (\mathbf{z} - \mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \leq h(\mathbf{y}) - h(\mathbf{x}), \quad \forall \mathbf{y} \in \mathbb{R}^n.$$

Remark: Projection theorem. Classical *projection theorem* is the special case of Proposition 3. Let $h(\mathbf{x}) = \mathbb{1}_{\mathcal{X}}(\mathbf{x})$, where $\mathcal{X} \subseteq \mathbb{R}^n$ is nonempty, closed, and convex. Then, for any $\mathbf{z} \in \mathbb{R}^n$, we have

$$\mathbf{x} = \text{prox}_h(\mathbf{z}) = \text{proj}_{\mathcal{X}}(\mathbf{z}) \Leftrightarrow (\mathbf{z} - \text{proj}_{\mathcal{X}}(\mathbf{z}))^\top (\mathbf{y} - \text{proj}_{\mathcal{X}}(\mathbf{z})) \leq 0, \quad \forall \mathbf{y} \in \mathcal{X}.$$

Proof. We can apply the first-order optimality condition to the definition of the proximal, which leads to

$$\mathbf{0} \in \mathbf{x} - \mathbf{z} + \partial h(\mathbf{x}) \Leftrightarrow (\mathbf{z} - \mathbf{x}) \in \partial h(\mathbf{x}).$$

Additionally, from the definition of the subgradient, we have for any $\mathbf{g}(\mathbf{x}) \in \partial h(\mathbf{x})$

$$h(\mathbf{y}) \geq h(\mathbf{x}) + \mathbf{g}(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}),$$

which directly leads to the last equivalence. ■

An important consequence of Theorem 3 is that a *fixed point* of the proximal operator minimizes h .

Corollary 1. Let $h \in \Gamma^0(\mathbb{R}^n)$ and $\mathbf{x} \in \mathbb{R}^n$, then

$$\mathbf{x} \text{ is a minimizer of } h \Leftrightarrow \mathbf{x} = \text{prox}_h(\mathbf{x}).$$

Proof. We know that \mathbf{x} is a minimizer of h if and only if $\mathbf{0} \in \partial h(\mathbf{x})$. Thus, from theorem above we obtain the equivalence with $\mathbf{x} = \text{prox}_h(\mathbf{x})$. ■

Finally, one can use Theorem 3 to establish the firm nonexpansiveness of the proximal operator.

Proposition 4. Let $h \in \Gamma^0(\mathbb{R}^n)$, then the proximal operator is *firmly nonexpansive*, which means that

$$\|\text{prox}_h(\mathbf{x}) - \text{prox}_h(\mathbf{y})\|_2^2 \leq (\text{prox}_h(\mathbf{x}) - \text{prox}_h(\mathbf{y}))^\top (\mathbf{x} - \mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

By using Cauchy-Schwarz inequality, we can further show that the proximal operator is *nonexpansive*

$$\|\text{prox}_h(\mathbf{x}) - \text{prox}_h(\mathbf{y})\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

Proof. Let $\mathbf{u} = \text{prox}_h(\mathbf{x})$ and $\mathbf{v} = \text{prox}_h(\mathbf{y})$. Then, we have that

$$\begin{cases} (\mathbf{x} - \mathbf{u}) \in \partial h(\mathbf{u}) \\ (\mathbf{y} - \mathbf{v}) \in \partial h(\mathbf{v}) \end{cases} \Rightarrow (\mathbf{x} - \mathbf{u} - \mathbf{y} + \mathbf{v})^\top (\mathbf{u} - \mathbf{v}) \geq 0 \Rightarrow (\mathbf{x} - \mathbf{y})^\top (\mathbf{u} - \mathbf{v}) \geq \|\mathbf{u} - \mathbf{v}\|_2^2,$$

where in the first implication we used the monotonicity of the subdifferential. ■

Characterization of proximal operators

Given a function h , one can derive the value of the proximal operator $R = \text{prox}_h$. In this section, we are interested in the reverse question: given an operator R , what conditions would make it a proximal operator? We will start with the following definition.

Definition 2. An operator $R : \mathcal{X} \rightarrow \mathbb{R}^n$ is a proximal operator of h if $R(z) \in \text{prox}_h(x)$ for each $x \in \mathcal{X}$.

Proximal operators for functions in $\Gamma^0(\mathbb{R}^n)$ were characterized by Moreau in his original paper [2].

Proposition 5. An operator $R : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is the proximal operator of $h \in \Gamma^0(\mathbb{R}^n)$ if and only if the following two conditions are true:

- (a) there exists a convex and closed function ψ such that: $R(x) \in \partial\psi(x)$ for all $x \in \mathbb{R}^n$;
- (b) R is nonexpansive: $\|R(x) - R(z)\|_2 \leq \|x - z\|_2$ for all $x, z \in \mathbb{R}^n$.

Proof. See Corollary 10.c in [2]. ■

The original result by Moreau was extended to nonconvex functions by Gribonval and Nikolova [3].

Proposition 6. Consider $R : \mathcal{X} \rightarrow \mathbb{R}^n$, where $\mathcal{X} \subset \mathbb{R}^n$ is nonempty. We have that

There exists $h : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ such that $R(x) \in \text{prox}_h(x)$, $\forall x \in \mathcal{X}$
 \Leftrightarrow There exists convex and closed $\psi : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ such that $R(x) \in \partial\psi(x)$, $\forall x \in \mathcal{X}$.

Proof. See Theorem 3 in [3]. ■

Remark. The analysis in Proposition 6 further implies that h and ψ can be chosen such that

$$\psi(x) = x^T R(x) - \frac{1}{2} \|R(x)\|_2^2 - h(R(x)), \quad \forall x \in \mathcal{X}.$$

The results [3] can be used to recover classical characterization of certain proximal operators.

Proposition 7. Consider $R : \mathcal{X} \rightarrow \mathbb{R}$, where $\mathcal{X} \subset \mathbb{R}$ is nonempty. We have that

R is nondecreasing \Leftrightarrow There exists $h : \mathbb{R} \rightarrow (-\infty, +\infty]$ such that $R(x) \in \text{prox}_h(x)$, $\forall x \in \mathcal{X}$.

Proof. See Corollary 7 in [3]. ■

Example. Consider the classical soft-thresholding function

$$T_\mu(z) = \text{sign}(z) \cdot \max(|z| - \mu, 0), \quad z \in \mathbb{R}.$$

Since T_μ is nondecreasing, it is a proximal operator. In fact, it is known that it the soft-thresholding function is the proximal operator of $h(x) = \mu|x|$.

Moreau smoothing

We will next discuss the concept of Moreau envelope for proper, closed, and convex functions. This will give us an elegant interpretation of the proximal operator of h as a gradient-descent step of some function h_μ .

Definition 3. The Moreau envelope of $h \in \Gamma^0(\mathbb{R}^n)$ is defined as

$$h_\mu(\mathbf{z}) = \inf_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2\mu} \|\mathbf{x} - \mathbf{z}\|_2^2 + h(\mathbf{x}) \right\}, \quad \mathbf{z} \in \mathbb{R}^n.$$

Example (Huber function). Let $h(z) = |z|$. Then, its Moreau envelope is just the Huber function

$$h_\mu(z) = \inf_x \left\{ \frac{1}{2\mu} (x - z)^2 + |x| \right\} = \begin{cases} \frac{1}{2\mu} z^2 & \text{for } |z| \leq \mu \\ |z| - \frac{\mu}{2} & \text{for } |z| > \mu \end{cases}.$$

Proposition 8. The Moreau envelope h_μ of $h \in \Gamma^0(\mathbb{R}^n)$ is convex and has a $(1/\mu)$ -Lipschitz gradient

$$\nabla h_\mu(\mathbf{z}) = \frac{1}{\mu} (\mathbf{z} - \text{prox}_{\mu h}(\mathbf{z})), \quad \mathbf{z} \in \mathbb{R}^n.$$

Remark. Note that Proposition 8 implies that

$$\text{prox}_{\mu h}(\mathbf{x}) = \mathbf{x} - \mu \nabla h_\mu(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^n,$$

which allows us to conclude that the proximal operator is computing a gradient-descent step on h_μ .

Proof. We first show that h_μ is convex. Consider

$$\varphi(\mathbf{x}, \mathbf{z}) = \frac{1}{2\mu} \|\mathbf{x} - \mathbf{z}\|_2^2 + h(\mathbf{x}),$$

which is convex over (\mathbf{x}, \mathbf{z}) . Then, for any $0 \leq \theta \leq 1$ and any $(\mathbf{x}_1, \mathbf{z}_1), (\mathbf{x}_2, \mathbf{z}_2)$ in \mathbb{R}^{2n} , we have

$$\begin{aligned} h_\mu(\theta \mathbf{z}_1 + (1 - \theta) \mathbf{z}_2) &\leq \varphi(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2, \theta \mathbf{z}_1 + (1 - \theta) \mathbf{z}_2) \\ &\leq \theta \varphi(\mathbf{x}_1, \mathbf{z}_1) + (1 - \theta) \varphi(\mathbf{x}_2, \mathbf{z}_2), \end{aligned}$$

where we used the convexity of φ . Since this inequality holds everywhere, we have

$$h_\mu(\theta \mathbf{z}_1 + (1 - \theta) \mathbf{z}_2) \leq \theta h_\mu(\mathbf{z}_1) + (1 - \theta) h_\mu(\mathbf{z}_2),$$

with

$$h_\mu(\mathbf{z}_1) = \inf_{\mathbf{x}_1} \varphi(\mathbf{x}_1, \mathbf{z}_1) \quad \text{and} \quad h_\mu(\mathbf{z}_2) = \inf_{\mathbf{x}_2} \varphi(\mathbf{x}_2, \mathbf{z}_2).$$

To show the differentiability, note that

$$\begin{aligned} h_\mu(\mathbf{z}) &= \frac{1}{2\mu} \|\mathbf{z}\|_2^2 - \frac{1}{\mu} \sup_{\mathbf{x} \in \mathbb{R}^n} \left\{ \mathbf{z}^\top \mathbf{x} - \mu h(\mathbf{x}) - \frac{1}{2} \|\mathbf{x}\|_2^2 \right\} \\ &= \frac{1}{2\mu} \|\mathbf{z}\|_2^2 - \frac{1}{\mu} \phi^*(\mathbf{z}) \quad \text{with} \quad \phi^*(\mathbf{z}) := \frac{1}{2} \|\mathbf{z}\|_2^2 + \mu h(\mathbf{z}), \end{aligned}$$

where ϕ^* denotes the conjugate function of ϕ . The function ϕ is closed and 1-strongly convex. Hence, we know that ϕ^* is defined for all $z \in \mathbb{R}^n$ and that the gradient of the conjugate function $\phi^*(z)$ is equal to the optimal x^* at which $\phi^*(x)$ is achieved

$$\nabla \phi^*(z) = \arg \max_{x \in \mathbb{R}^n} \left\{ z^\top x - \mu h(x) - \frac{1}{2} \|x\|_2^2 \right\} = \text{prox}_{\mu h}(z).$$

Hence, we conclude that

$$\frac{1}{\mu} (z - \text{prox}_{\mu h}(z)) = \nabla h_\mu(z).$$

Note that since the proximal operator is firmly nonexpansive, $\mu \nabla h_\mu$ is also firmly nonexpansive, which implies that it is $(1/\mu)$ -Lipschitz continuous. ■

We next show that the Moreau envelope can serve as a smooth approximation to a nonsmooth function.

Proposition 9. Consider $h(z)$, its Moreau envelope $h_\mu(z)$ for $\mu > 0$, and any $g(z) \in \partial h(z)$. Then,

$$0 \leq h(z) - h_\mu(z) \leq \frac{\mu}{2} \|g(z)\|_2^2.$$

Proof. First note that we have for any $z \in \mathbb{R}^n$

$$h_\mu(z) = \inf_x \left\{ \frac{1}{2\mu} \|x - z\|_2^2 + h(x) \right\} \leq h(z),$$

which is due to the fact that $x = z$ is potentially suboptimal. We additionally have for any $g(z) \in \partial h(z)$

$$\begin{aligned} h_\mu(z) - h(z) &= \inf_x \left\{ h(x) - h(z) + \frac{1}{2\mu} \|x - z\|_2^2 \right\} \\ &\geq \inf_x \left\{ g(z)^\top (x - z) + \frac{1}{2\mu} \|x - z\|_2^2 \right\} \\ &= \inf_x \left\{ \frac{1}{2\mu} \|x - (z - \mu g(z))\|_2^2 - \frac{\mu}{2} \|g(z)\|_2^2 \right\} = -\frac{\mu}{2} \|g(z)\|_2^2. \end{aligned}$$

This directly leads to the conclusion. ■

Bayesian estimation

We are now equipped with tools to establish a relationship between an estimation problem (2) and the proximal operator. We will adopt the Bayesian perspective by designing an operator $R(s)$ for (2) that minimizes

$$\mathbb{E}_{s,x} [\mathcal{L}(x, R(s))], \quad (3)$$

where $\mathcal{L}(x, \hat{x})$ is some loss function and (x, s) is drawn from the joint distribution $p_{x,s}(x, s) = p_{s|x}(s|x) p_x(x)$. Note that the conditional distribution of s given x is simply the Gaussian probability density function

$$p_{s|x}(s|x) = \mathcal{G}_\sigma(s - x) = \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left(-\frac{\|s - x\|_2^2}{2\sigma^2}\right),$$

where we introduce the notation \mathcal{G}_σ for the Gaussian density with the standard deviation $\sigma > 0$.

Example. The *maximum a posteriori (MAP)* estimator is often been used to justify various inverse problem solvers. MAP can be seen as minimizing (3) with the pseudo-cost $\mathcal{L}(\mathbf{x}, \hat{\mathbf{x}}) = \delta(\mathbf{x} - \hat{\mathbf{x}})$. A MAP denoiser simply selects a vector that maximizes $p(\mathbf{x}|\mathbf{s}) \propto p(\mathbf{x}|\mathbf{s})p_{\mathbf{x}}(\mathbf{x})$

$$\mathbf{R}_{\text{map}}(\mathbf{s}) \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} p(\mathbf{x}|\mathbf{s}) = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{s}\|_2^2 - \log p_{\mathbf{x}}(\mathbf{x}) \right\}.$$

Therefore, a MAP denoiser can be directly expressed as the proximal operator

$$\mathbf{R}_{\text{map}}(\mathbf{s}) \in \text{prox}_{\mu h_{\text{map}}}(\mathbf{s}) \quad \text{with } \mu := \sigma^2 \text{ and } h_{\text{map}}(\mathbf{x}) := -\log p_{\mathbf{x}}(\mathbf{x}).$$

Example. The *minimum mean squared error (MMSE)* estimator seeks to minimize (3) with the cost $\mathcal{L}(\mathbf{x}, \hat{\mathbf{x}}) = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$ and the solution is the conditional expectation

$$\mathbf{R}_{\text{mmse}}(\mathbf{s}) = \mathbb{E}[\mathbf{x}|\mathbf{s}] = \int_{\mathbb{R}^n} \mathbf{x} p(\mathbf{x}|\mathbf{s}) d\mathbf{x} = \frac{1}{p(\mathbf{s})} \int_{\mathbb{R}^n} \mathbf{x} \mathcal{G}_{\sigma}(\mathbf{x} - \mathbf{s}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}.$$

Proposition 10 (Tweedie formula). Let $\mathbf{R}_{\text{mmse}}(\mathbf{s}) = \mathbb{E}[\mathbf{x}|\mathbf{s}]$ for problem (2), then we have

$$\mathbf{R}_{\text{mmse}}(\mathbf{s}) = \mathbf{s} + \sigma^2 \nabla \log p(\mathbf{s}), \quad \mathbf{s} \in \mathbb{R}^n.$$

Proof. First note that gradient of the Gaussian density is given by

$$\nabla_{\mathbf{s}} \mathcal{G}_{\sigma}(\mathbf{s} - \mathbf{x}) = \frac{1}{\sigma^2} (\mathbf{x} - \mathbf{s}) \mathcal{G}_{\sigma}(\mathbf{s} - \mathbf{x}).$$

Now consider the marginal distribution and its gradient

$$\begin{aligned} p(\mathbf{s}) &= \int p(\mathbf{x}, \mathbf{s}) d\mathbf{x} = \int \mathcal{G}_{\sigma}(\mathbf{s} - \mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \\ \Rightarrow \quad \nabla p(\mathbf{s}) &= \frac{1}{\sigma^2} \int (\mathbf{x} - \mathbf{s}) \mathcal{G}_{\sigma}(\mathbf{s} - \mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} = \frac{1}{\sigma^2} \left[\int \mathbf{x} p(\mathbf{x}, \mathbf{s}) d\mathbf{x} - \mathbf{s} p(\mathbf{s}) \right], \end{aligned}$$

where we used the gradient of the Gaussian density and $p(\mathbf{x}, \mathbf{s}) = \mathcal{G}_{\sigma}(\mathbf{s} - \mathbf{x}) p_{\mathbf{x}}(\mathbf{x})$. Dividing both sides by $p(\mathbf{s})$ and using the definition of conditional mean, we obtain

$$\nabla \log p(\mathbf{s}) = \frac{\nabla p(\mathbf{s})}{p(\mathbf{s})} = \frac{1}{\sigma^2} (\mathbb{E}[\mathbf{x}|\mathbf{s}] - \mathbf{s}),$$

which establishes the desired result. ■

MAP and MMSE denoisers are generally not equivalent. However, they are when the prior is Gaussian.

Example. Consider the Gaussian prior $p_{\mathbf{x}}(\mathbf{x}) = \mathcal{G}_{\nu}(\mathbf{x})$ where $\nu > 0$. Then, we have that

$$\mathbf{R}_{\text{map}}(\mathbf{s}) = \mathbf{R}_{\text{mmse}}(\mathbf{s}) = \left(\frac{\nu^2}{\nu^2 + \sigma^2} \right) \mathbf{s}.$$

Note how in the Gaussian case, both estimators correspond to the same *linear* function.

Proof. Let $p_{\mathbf{x}}(\mathbf{x}) = \mathcal{G}_{\nu}(\mathbf{x})$ where $\nu > 0$. For MAP, we have that

$$\mathbf{R}_{\text{map}}(\mathbf{s}) = \arg \min_{\mathbf{x}} \left\{ \frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{s}\|_2^2 + \frac{1}{2\nu^2} \|\mathbf{x}\|_2^2 \right\} = \left(\frac{\nu^2}{\nu^2 + \sigma^2} \right) \mathbf{s}.$$

For MMSE, we have that

$$p(\mathbf{s}) = \int \mathcal{G}_{\sigma}(\mathbf{s} - \mathbf{x}) \mathcal{G}_{\nu}(\mathbf{x}) d\mathbf{x} = \mathcal{G}_{\alpha}(\mathbf{s}) \quad \Rightarrow \quad \log p(\mathbf{s}) = -\frac{1}{\alpha^2} \|\mathbf{s}\|_2^2,$$

where $\alpha^2 := \sigma^2 + \nu^2$. Using the Tweedie formula, we obtain

$$\mathbf{R}_{\text{mmse}}(\mathbf{s}) = \mathbf{s} - \left(\frac{\sigma^2}{\nu^2 + \sigma^2} \right) \mathbf{s} = \left(\frac{\nu^2}{\nu^2 + \sigma^2} \right) \mathbf{s}.$$

■

The following result is due to Gribonval [4] and shows that the MMSE denoiser is a proximal operator.

Proposition 11. For a non-degenerate prior $p_{\mathbf{x}}$, there exists a function h_{mmse} such that

$$\mathbf{R}_{\text{mmse}}(\mathbf{s}) = \text{prox}_{h_{\text{mmse}}}(\mathbf{s}), \quad \forall \mathbf{s} \in \mathbb{R}^n.$$

The function h_{mmse} is given by

$$h_{\text{mmse}}(\mathbf{x}) = \begin{cases} -\frac{1}{2} \|\mathbf{x} - \mathbf{R}_{\text{mmse}}^{-1}(\mathbf{x})\|_2^2 - \sigma^2 \log p_{\mathbf{s}}(\mathbf{R}_{\text{mmse}}^{-1}(\mathbf{x})) & \text{for } \mathbf{x} \in \mathcal{X} \\ +\infty & \text{for } \mathbf{x} \notin \mathcal{X} \end{cases},$$

where $\mathcal{X} := \text{Im}(\mathbf{R}_{\text{mmse}})$, $\mathbf{R}_{\text{mmse}}^{-1} : \mathcal{X} \rightarrow \mathbb{R}^n$ is the inverse mapping, which is well defined and smooth over \mathcal{X} , and $p_{\mathbf{s}}$ is the probability distribution of the AWGN corrupted observation in (2).

Proof. See Theorem 1 in [4].

■

Remark. Proposition 11 states that there exists a function h_{mmse} such that the MMSE denoiser is its proximal operator. The underlying analysis additionally shows that h_{mmse} is C^∞ over $\text{Im}(\mathbf{R}_{\text{mmse}})$. One interesting conclusion from this results is that one can in principle misinterpret a MMSE denoiser as a MAP denoiser corresponding to a prior probability $q(\mathbf{x}) \propto \exp(-h_{\text{mmse}}(\mathbf{x}))$.

The next result provides necessary and sufficient conditions on the convexity and separability of the implicit regularizer h specified via the MMSE estimator.

Proposition 12. Let p_x be a non-degenerate prior and h be the regularizer associated with R_{mmse} via Proposition 11. The following statements are true:

- (a) h is convex if and only if $p_s = (\mathcal{G}_\sigma * p_x)$ is log-concave.
- (b) h is additively separable if and only if p_s is multiplicatively separable.

Proof. See Lemma 2 in [5]. ■

Remark. Proposition 12 uses the following notions of separability.

- (a) $R : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is separable if there exists a set of functions $R_1, \dots, R_n : \mathbb{R} \rightarrow \mathbb{R}$ such that $R(\mathbf{x}) = (R_1(x_1), \dots, R_n(x_n))$, for all $\mathbf{x} \in \mathbb{R}^n$.
- (b) $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is additively separable if there exists a set of functions $h_1, \dots, h_n : \mathbb{R} \rightarrow \mathbb{R}$ such that $h(\mathbf{x}) = h_1(x_1) + \dots + h_n(x_n)$ for all $\mathbf{x} \in \mathbb{R}^n$.
- (c) $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is multiplicatively separable if there exists a set of functions $h_1, \dots, h_n : \mathbb{R} \rightarrow \mathbb{R}$ such that $h(\mathbf{x}) = h_1(x_1) \cdots h_n(x_n)$ for all $\mathbf{x} \in \mathbb{R}^n$.

The following result shows that log-concavity of p_x is sufficient for the convexity of h associated with R_{mmse} .

Proposition 13. The log-concavity of the prior p_x implies the log-concavity of p_s .

Proof. It is well-known that convolution preserves log-concavity. Since both \mathcal{G}_σ and p_x are log-concave, we have that $p_s = (\mathcal{G}_\sigma * p_x)$ is also log-concave. ■

Remark. From Propositions 12 and 13, we can conclude that

$$p_x \text{ is log-concave} \Rightarrow p_s \text{ is log-concave} \Leftrightarrow h_{\text{mmse}} \text{ is convex.}$$

Note how log-concavity of p_x is not a necessary condition for that of p_s . This means that there are some priors p_x that are not log-concave such that h_{mmse} associated with R_{mmse} is convex.

Deep denoising networks

In practice, denoisers are obtained by training a deep network R_θ to minimize an empirical loss function

$$\frac{1}{k} \sum_{i=1}^k \mathcal{L}(\mathbf{x}_i, R_\theta(\mathbf{s}_i)), \quad (4)$$

where $\mathbf{x}_1, \dots, \mathbf{x}_k$ are clean images assumed to be sampled from an unknown image prior p_x and each \mathbf{s}_i is the noisy observation of type (2). Since the loss (4) is an approximation of the true expectation (3), for $\mathcal{L}(\mathbf{x}, \hat{\mathbf{x}}) = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$, we can interpret the solution R_{θ^*} of (4) to be an approximation the true MMSE denoiser

$$R_{\theta^*}(\mathbf{s}) = R_{\text{mmse}}(\mathbf{s}) + \varepsilon(\mathbf{s}) = \mathbb{E}[\mathbf{x}|\mathbf{s}] + \varepsilon(\mathbf{s}), \quad \mathbf{s} \in \mathbb{R}^n, \quad (5)$$

where we use $\varepsilon(\mathbf{s})$ denote the error between the true MMSE and the learned denoisers. While the design and training of state-of-the-art denoisers is beyond the scope of these lectures, interested readers can find plenty of material in the literature [6, 7].

Model-based optimization

It is common to formulate the solutions of inverse problems (1) as an optimization problem of form

$$\hat{\mathbf{x}} \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad \text{with} \quad f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x}), \quad (6)$$

where g is a data fidelity term that quantifies consistency with the observed measurements \mathbf{y} and h is a regularizer that enforces prior knowledge on \mathbf{x} . The formulation (6) is known as the *variational optimization*, *regularized inversion*, and *model-based reconstruction*.

Remark. The formulation in (6) can be interpreted as the MAP estimator when

$$g(\mathbf{x}) = -\log(p_{\mathbf{y}|\mathbf{x}}(\mathbf{x})) \quad \text{and} \quad h(\mathbf{x}) = -\log(p_{\mathbf{x}}(\mathbf{x})), \quad (7)$$

where $p_{\mathbf{y}|\mathbf{x}}$ is the likelihood relating \mathbf{x} to measurements \mathbf{y} and $p_{\mathbf{x}}$ is the prior distribution.

Example: least squares. One of the most popular functions used in practice is the quadratic function

$$g(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 = \frac{1}{2} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x} - y_i)^2 = \frac{1}{2} \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} - \mathbf{y}^\top \mathbf{A} \mathbf{x} + \frac{1}{2} \mathbf{y}^\top \mathbf{y},$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a matrix with m rows $\mathbf{a}_i \in \mathbb{R}^n$, $\mathbf{x} \in \mathbb{R}^n$ is the optimization variable, and $\mathbf{y} \in \mathbb{R}^m$ is the observation vector. Quadratic function above corresponds to the negative log-likelihood for the AWGN. The gradient and the Hessian of the least-squares is given by

$$\nabla g(\mathbf{x}) = \mathbf{A}^\top (\mathbf{A} \mathbf{x} - \mathbf{y}) \quad \text{and} \quad \text{H}g(\mathbf{x}) = \mathbf{A}^\top \mathbf{A} \succeq 0.$$

Example: indicator function. An *indicator function* of a closed set $\mathcal{X} \subseteq \mathbb{R}^n$ corresponds to a regularizer defined as follows

$$\mathbb{1}_{\mathcal{X}}(\mathbf{x}) := \begin{cases} 0 & \text{if } \mathbf{x} \in \mathcal{X} \\ +\infty & \text{if } \mathbf{x} \notin \mathcal{X} \end{cases}.$$

Indicator function is often used for imposing hard constraints on the solution, such as nonnegativity. Note that when $\mathcal{X} \neq \emptyset$ is a convex set, the indicator function is a proper, closed, and convex function.

Example: total variation. Another classical regularizer in computational imaging is the *total variation* (TV) function. It has two variants. The *isotropic* TV is defined as follows

$$h_{\text{iv}}(\mathbf{x}) = \lambda \sum_{i=1}^n \|[D\mathbf{x}]_i\|_2 = \lambda \sum_{i=1}^n \sqrt{\sum_{j=1}^d ([D_j\mathbf{x}]_i)^2},$$

while the *anisotropic* TV is defined as

$$r_{\text{av}}(\mathbf{x}) = \lambda \sum_{i=1}^n \|[D\mathbf{x}]_i\|_1 = \lambda \sum_{i=1}^n \sum_{j=1}^d |[D_j\mathbf{x}]_i|,$$

where $\lambda > 0$ controls the strength of the regularization and $\mathbf{D} = (\mathbf{D}_1, \dots, \mathbf{D}_d)$ denotes the d -dimensional discrete gradient. Both variants of TV are nonsmooth convex functions.

Plug-and-play priors

There have been many algorithms developed over the years for solving composite optimization functions of form (6). *Proximal algorithms* have received significant attention due to their ability to solve (6) when g or h is nonsmooth [8]. One important property of proximal algorithms is that they do not explicitly require knowledge of h , relying instead on the prox_h . We know that both MAP and MMSE estimators correspond to proximal operators of some h , motivating us to simply replace prox_h in proximal algorithms by R_θ .

Alternating direction method of multipliers

One of the most widely used and effective proximal algorithms is the *alternating direction method of multipliers* (ADMM), which uses an augmented Lagrangian formulation to allow for alternating minimization of each function in turn (see [9] for an overview of ADMM). To develop the algorithm, we consider the following optimization problem over (z, x) equivalent to the original problem

$$\text{minimize } g(z) + h(x) \quad \text{subject to } z = x. \quad (8)$$

This process of introducing additional convenient variables is known as *variable splitting*. We will denote the optimal value of the problem by

$$f^* = \inf\{g(z) + h(x) : z - x = \mathbf{0}\}.$$

We form the *augmented Lagrangian*

$$\begin{aligned} L_\gamma(z, x, \mu) &= g(z) + h(x) + \mu^\top(z - x) + \frac{1}{2\gamma}\|z - x\|_2^2 \\ &= g(z) + h(x) + \frac{1}{2\gamma}\|z - x + \gamma\mu\|_2^2 - \frac{\gamma}{2}\|\mu\|_2^2, \end{aligned}$$

where $\gamma > 0$ is the quadratic parameter and $\mu \in \mathbb{R}^n$ is the Lagrangian. The augmented Lagrangian can be viewed as a traditional Lagrangian associated with the problem

$$\text{minimize } g(z) + h(x) + \frac{1}{2\gamma}\|z - x\|_2^2 \quad \text{subject to } z = x.$$

This problem is equivalent to the original problem, since for any feasible (z, x) , the term added to the objective is zero. We can re-write the augmented Lagrangian by introducing a *scaled dual variable* $s := \gamma\mu$, which leads to

$$L_\gamma(z, x, s) = g(z) + h(x) + \frac{1}{2\gamma}\|z - x + s\|_2^2 - \frac{1}{2\gamma}\|s\|_2^2.$$

We next solve this problem using the *method of multipliers* that has the following form

$$\begin{aligned} (z^t, x^t) &= \arg \min_{x, z} L_\gamma(x, z, s^{t-1}) \\ s^t &= s^{t-1} + (z^t - x^t). \end{aligned}$$

The difficulty of running this algorithm in practice lies in the first step, where there is a need to jointly minimize over both z and x . ADMM splits this step into two, substantially simplifying the optimization.

$$\begin{aligned} z^t &= \arg \min_{z \in \mathbb{R}^n} L_\gamma(z, x^{t-1}, s^{t-1}) = \text{prox}_{\gamma g}(x^{t-1} - s^{t-1}) \\ x^t &= \arg \min_{x \in \mathbb{R}^n} L_\gamma(z^t, x, s^{t-1}) = \text{prox}_{\gamma h}(z^t + s^{t-1}) \\ s^t &= s^{t-1} + (z^t - x^t). \end{aligned}$$

PnP-ADMM [10, 11] enables the substitution of prox_h by a deep denoiser R_θ .

Plug-and-play ADMM (PnP-ADMM)

Input: An initial values $\mathbf{x}^0 \in \mathbb{R}^n$ and $\mathbf{s}^0 = \mathbf{0}$, and a parameter $\gamma > 0$.

Iterate: For $t = 1, 2, 3, \dots$, do

$$\mathbf{z}^t \leftarrow \text{prox}_{\gamma g}(\mathbf{x}^{t-1} - \mathbf{s}^{t-1})$$

$$\mathbf{x}^t \leftarrow \mathbf{R}_\theta(\mathbf{z}^t + \mathbf{s}^{t-1})$$

$$\mathbf{s}^t \leftarrow \mathbf{s}^{t-1} + (\mathbf{z}^t - \mathbf{x}^t).$$

Remark. PnP-ADMM has an important feature of modularity; it explicitly separate the application of the physical models (the data fidelity update) from that of the learned models (the image denoising). This observation reveals a key strength of PnP methods: they can be easily customized for different measurement operators by changing the data fidelity term, thus enabling the use of the same learned network over a wide range of applications without retraining.

It is possible to establish the convergence of ADMM under some basic assumptions.

Proposition 14. Let $g \in \Gamma^0(\mathbb{R}^n)$ and $h \in \Gamma^0(\mathbb{R}^m)$ and assume that the unaugmented Lagrangian L_0 has a saddle point. Then, ADMM iterates satisfy the following:

- *Residual convergence:* The iterates approach feasibility

$$\lim_{t \rightarrow \infty} (\mathbf{z}^t - \mathbf{x}^t) = \mathbf{0}.$$

- *Objective convergence:* The objective function of the iterates approaches the optimum

$$\lim_{t \rightarrow \infty} (g(\mathbf{z}^t) + h(\mathbf{x}^t)) = f^*.$$

- *Dual variable convergence:* Lagrangian converges to the dual optimal point

$$\lim_{t \rightarrow \infty} \boldsymbol{\mu}^t = \boldsymbol{\mu}^*.$$

Remark. The unaugmented Lagrangian L_0 has a saddle point if there exist $(\mathbf{z}^*, \mathbf{x}^*, \mathbf{s}^*)$ for which

$$L_0(\mathbf{z}^*, \mathbf{x}^*, \mathbf{s}) \leq L_0(\mathbf{z}^*, \mathbf{x}^*, \mathbf{s}^*) \leq L_0(\mathbf{z}, \mathbf{x}, \mathbf{s}^*), \quad \forall (\mathbf{z}, \mathbf{x}, \mathbf{s}).$$

Together, the convexity and the saddle point condition imply that the strong duality will hold for the constrained optimization problem. This will imply that $(\mathbf{z}^*, \mathbf{x}^*)$ is a solution to the problem such that $\mathbf{z} - \mathbf{x} = \mathbf{0}$, which means that $g(\mathbf{x}^*) < \infty$ and $h(\mathbf{x}^*) < \infty$.

Proof. See Section 3.2 in [9]. ■

Although, the analysis of ADMM above assumes that both g and h are proper, closed, and convex, it can be extend in numerous ways to be more compatible with PnP. For example, one can get the convergence specifically for the MMSE denoisers without any assumptions on the convexity of g and h .

Proposition 15. Assume that p_x is non-degenerate over \mathbb{R}^n , g is continuously differentiable, h associated with $\mathbf{R}_\theta = \mathbf{R}_{\text{mmse}}$ has a M -Lipschitz continuous gradient over $\mathcal{X} := \text{Im}(\mathbf{R}_\theta)$. Then, PnP-ADMM iterates for $0 < \gamma \leq 1/(2M)$ satisfy $\|\nabla f(\mathbf{x}^t)\|_2 \rightarrow 0$ as $t \rightarrow \infty$.

Proof. See Theorem 1 in [12]. ■

Half-quadratic splitting

Deep plug-and-play image restoration (DPIR) is one of the state-of-the-art implementations of PnP based on the DRUNet denoiser [7]. DPIR iterations are not based on ADMM. They instead correspond to the well-known *half-quadratic splitting (HQS)* method that seeks to solve (8) using a quadratic penalty method

$$L_\gamma(\mathbf{z}, \mathbf{x}) = g(\mathbf{z}) + h(\mathbf{x}) + \frac{1}{2\gamma} \|\mathbf{z} - \mathbf{x}\|_2^2.$$

DPIR minimizes L_γ by changing the noise level σ parameter of DRUNet at every iteration in a way that corresponds to decreasing the value of γ . In addition to $\gamma > 0$ (denoted as $\mu = 1/\gamma$) The presentation of the original DPIR explicitly includes the noise level $\sigma > 0$ and a parameter $\lambda > 0$ for the regularizer. However, one can conceptually simplify the algorithm as follows by absorbing the parameters into functions g and h .

Deep plug-and-play image restoration (DPIR)

Input: An initial values $\mathbf{x}^0 \in \mathbb{R}^n$ and a decreasing set of positive parameters $\{\gamma_t\}$.

Iterate: For $t = 1, 2, 3, \dots$, do

$$\begin{aligned} \mathbf{z}^t &\leftarrow \text{prox}_{\gamma_t g}(\mathbf{x}^{t-1}) \\ \mathbf{x}^t &\leftarrow \mathbf{R}_\theta(\mathbf{z}^t, \gamma_t) \end{aligned}$$

Proximal gradient method

Another widely-used proximal algorithm in the context of PnP is *proximal gradient method (PGM)*, which is also known as *forward backward splitting* and *iterative shrinkage/thresholding algorithm*.

Plug-and-play proximal gradient method (PnP-PGM)

Input: An initial value $\mathbf{x}^0 \in \mathbb{R}^n$ and a step-size $\gamma > 0$.

For $t = 1, 2, 3, \dots$, do

$$\mathbf{x}^t \leftarrow \mathbf{R}_\theta(\mathbf{x}^{t-1} - \gamma \nabla g(\mathbf{x}^{t-1})).$$

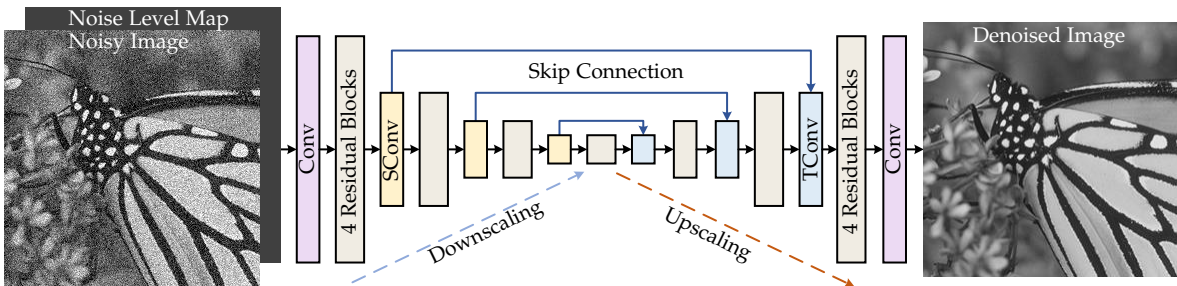


Figure 2: DRUNet is a deep denoiser for DPIR that takes a noisy image s and the noise level σ as inputs [7].

Plug-and-play accelerated PGM (PnP-APGM)

Input: An initial value $\mathbf{x}^0 = \mathbf{s}^0 \in \mathbb{R}^n$, a step-size $\gamma > 0$, and a set of positive parameters $\{\theta_t\}$.

For $t = 1, 2, 3, \dots$, do

$$\begin{aligned}\mathbf{x}^t &\leftarrow \mathbf{R}_\theta(\mathbf{s}^{t-1} - \gamma \nabla g(\mathbf{s}^{t-1})) \\ \mathbf{s}^t &\leftarrow (1 - \theta_t)\mathbf{x}^t + \theta_t \mathbf{x}^{t-1}\end{aligned}$$

Remark. Note the difference in the treatment of data fidelity term g between different PnP algorithm. While PnP-PGM uses the standard (explicit) gradient step $\mathbf{z} = \mathbf{x} - \gamma \nabla g(\mathbf{x})$, PnP-ADMM uses the proximal operator $\text{prox}_{\gamma g}$, which can be written as an implicit gradient step $\mathbf{z} = \mathbf{x} - \gamma \nabla g(\mathbf{z})$ with ∇g evaluated at \mathbf{z} . For the quadratic loss $g(\mathbf{x}) = (1/2)\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$, we obtain the following updates

$$\begin{aligned}\mathbf{x} - \gamma g(\mathbf{x}) &= \mathbf{A}^\top (\mathbf{A}\mathbf{x} - \mathbf{y}) \\ \text{prox}_{\gamma g}(\mathbf{x}) &= (\mathbf{I} + \gamma \mathbf{A}^\top \mathbf{A})^{-1}(\mathbf{x} + \gamma \mathbf{A}^\top \mathbf{y}).\end{aligned}$$

The convergence of both PGM and APGM have been extensively analyzed in the literature.

Theorem 1. Let $f = g + h$ have a finite minimum f^* , g have a L -Lipschitz continuous gradient, and h be sub-differentiable. Run PGM using a fixed step-size $0 < \gamma < 1/L$. Then, the sequence $\{f(\mathbf{x}^t)\}$ is non-increasing and $\|\mathbf{w}(\mathbf{x}^t)\|_2 \rightarrow 0$ as $t \rightarrow \infty$ with $\mathbf{w}(\mathbf{x}^t) \in \partial f(\mathbf{x}^t)$.

Remark. Note that since h associated with \mathbf{R}_{mmse} is differentiable on the iterates of PnP-PGM, we have that $\mathbf{w}(\mathbf{x}^t) = \nabla f(\mathbf{x}^t) = \nabla g(\mathbf{x}^t) + \nabla h(\mathbf{x}^t)$.

Proof. Without loss of generality set $\gamma = 1/(\alpha L)$ with $\alpha > 1$. Consider a single iteration of PGM

$$\mathbf{x}^+ \in \text{prox}_{\gamma h}(\mathbf{x} - \gamma \nabla g(\mathbf{x})) = \arg \min_{\mathbf{z}} \left\{ \frac{1}{2\gamma} \|\mathbf{z} - \mathbf{x}\|_2^2 + \nabla g(\mathbf{x})^\top (\mathbf{z} - \mathbf{x}) + h(\mathbf{z}) \right\}.$$

- From the Lipschitz continuity of ∇g , we have

$$g(\mathbf{x}^+) \leq g(\mathbf{x}) + \nabla g(\mathbf{x})^\top (\mathbf{x}^+ - \mathbf{x}) + \frac{L}{2} \|\mathbf{x}^+ - \mathbf{x}\|_2^2.$$

- From the first order optimality of prox, we have

$$h(\mathbf{x}^+) \leq h(\mathbf{x}) - \nabla g(\mathbf{x})^\top (\mathbf{x}^+ - \mathbf{x}) - \frac{1}{2\gamma} \|\mathbf{x}^+ - \mathbf{x}\|_2^2.$$

By combining these two inequalities, we obtain

$$f(\mathbf{x}^+) \leq f(\mathbf{x}) - (\alpha - 1) \frac{L}{2} \|\mathbf{x}^+ - \mathbf{x}\|_2^2 \quad \Leftrightarrow \quad \|\mathbf{x}^+ - \mathbf{x}\|_2^2 \leq \frac{2}{L(\alpha - 1)} (f(\mathbf{x}) - f(\mathbf{x}^+)).$$

Note that this inequality directly implies that $f(\mathbf{x}^+) \leq f(\mathbf{x})$.

On the other hand, from the first-order optimality conditions for the prox, we know that

$$\mathbf{0} \in \frac{1}{\gamma}(\mathbf{x}^+ - \mathbf{x}) + \nabla g(\mathbf{x}) + \partial h(\mathbf{x}^+) \quad \Rightarrow \quad \frac{1}{\gamma}(\mathbf{x} - \mathbf{x}^+) \in \nabla g(\mathbf{x}) + \partial h(\mathbf{x}^+).$$

This directly implies that the following vector is a subgradient of $f = g + h$ at \mathbf{x}^+

$$\mathbf{w}(\mathbf{x}^+) = \frac{1}{\gamma}(\mathbf{x} - \mathbf{x}^+) + \nabla g(\mathbf{x}^+) - \nabla g(\mathbf{x}) \in \partial f(\mathbf{x}^+).$$

This implies that

$$\|\mathbf{w}(\mathbf{x}^+)\|_2 \leq \frac{1}{\gamma} \|\mathbf{x}^+ - \mathbf{x}\|_2 + \|\nabla g(\mathbf{x}^+) - \nabla g(\mathbf{x})\|_2 \leq L(\alpha + 1) \|\mathbf{x}^+ - \mathbf{x}\|_2,$$

where we used the triangular inequality and L -Lipschitz continuity of ∇g .

We can thus obtain this bound

$$\|\mathbf{w}^+\|_2^2 \leq L^2(\alpha + 1)^2 \|\mathbf{x}^+ - \mathbf{x}\|^2 \leq L \frac{(\alpha + 1)^2}{(\alpha - 1)} (f(\mathbf{x}) - f(\mathbf{x}^+)),$$

which we can iterate to obtain

$$\sum_{i=0}^t \|\mathbf{w}(\mathbf{x}^i)\|_2^2 \leq L \frac{(\alpha + 1)^2}{(\alpha - 1)} (f(\mathbf{x}^0) - f(\mathbf{x}^{t+1})) \leq L \frac{(\alpha + 1)^2}{(\alpha - 1)} (f(\mathbf{x}^0) - f^*),$$

which implies that $\mathbf{w}(\mathbf{x}^t) \rightarrow \mathbf{0}$ as $t \rightarrow \infty$. ■

Remark. We will skip other analyses of PGM and APGM, since they can be readily found elsewhere. The key point to remember is that APGM enables faster $O(1/t^2)$ convergence in terms of f when the function is convex (see for example Nesterov [13]).

Regularization by denoising

Regularization by denoising (RED) is a variant of PnP that seeks to use the gradient of the regularizer associated with a given denoiser rather than its proximal operator. It was originally introduced [14] and its theory has been significantly clarified in subsequent papers [15, 16].

Gradient method regularization by denoising (GM-RED)

Input: An initial value $\mathbf{x}^0 \in \mathbb{R}^n$, a step-size $\gamma > 0$, and regularization parameter $\tau > 0$.

For $t = 1, 2, 3, \dots$, do

$$\mathbf{x}^t \leftarrow \mathbf{x}^{t-1} - \gamma(\nabla g(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathbf{R}_\theta(\mathbf{x}^{t-1}))).$$

PGM-RED

Input: An initial value $\mathbf{x}^0 \in \mathbb{R}^n$, a step-size $\gamma > 0$, and regularization parameter $\tau > 0$.

For $t = 1, 2, 3, \dots$, do

$$\mathbf{x}^t \leftarrow \text{prox}_{\gamma g}(\mathbf{x}^{t-1} - \gamma\tau(\mathbf{x}^{t-1} - \mathbf{R}_\theta(\mathbf{x}^{t-1}))).$$

Propositions 8 and 10 provide theoretical background necessary to understand the convergence of the RED algorithms for MAP and MMSE denoisers.

Remark. RED seeks to minimize a function $f = g + h_\mu$ where h_μ depends on the type of the denoiser. For simplicity, we will let $\tau = 1/\mu$, where $\mu := \sigma^2$.

- **MAP denoiser:** Consider a MAP denoiser for a log-concave prior p_x

$$R_\theta(x) = \text{prox}_{\mu h_{\text{map}}}(x) \quad \text{with} \quad h_{\text{map}}(x) := -\log p_x(x) \quad \Rightarrow \quad h_\mu \text{ is the Moreau envelope of } h_{\text{map}}$$

- **MMSE denoiser:** Consider a MMSE denoiser

$$R_\theta(x) = R_{\text{mmse}}(x) \quad \Rightarrow \quad h_\mu(x) = \nabla \log p_s(x) = \nabla \log (\mathcal{G}_\sigma * p_x)(x).$$

Gradient-step PnP

The traditional PnP strategy is to use a deep network $R_\theta : \mathbb{R}^n \rightarrow \mathbb{R}^n$ to parametrize a MAP or MMSE denoisers, thus giving us an implicit access to a regularizer h that depends on the prior p_x . An alternative approach is to use a deep network $h_\theta : \mathbb{R}^n \rightarrow \mathbb{R}$ to parametrize the regularizer directly. An innovative approach is to define the following denoiser for over $x \in \mathbb{R}^n$ [18, 17]

$$h_\theta(x) := \frac{1}{2} \|x - D_\theta(x)\|_2^2 \quad \Rightarrow \quad R_\theta(x) := x - \gamma \nabla h_\theta(x),$$

where $D_\theta : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a deep network (such as DRUNet or DnCNN). The operator R_θ is the *gradient-step* (GS) denoiser that can be trained and deployed as any traditional PnP denoiser. Two key benefit of the GS approach is that (a) the denoiser exactly represents a conservative vector field and (b) there is a direct access to the regularizer h_θ .

Online PnP methods

The traditional PnP methods are batch algorithms in the sense that they compute the gradient ∇g or the proximal operator $\text{prox}_{\gamma g}$ of the data fidelity term g by using the whole measurement vector $y \in \mathbb{R}^m$. The per-iteration computational and memory complexity of batch PnP algorithms depends on the total number of measurements. For example, in tomography with b projections, the complexity of evaluating ∇g scales linearly with b , making it computationally expensive for a large b . This has motivated interest in *online*, *stochastic*, and/or *incremental* PnP algorithms that approximate the batch ∇g with an approximation $\hat{\nabla} g$ based on a single element or a small subset of the measurements [19, 20, 21].

Consider the decomposition of \mathbb{R}^m into $b \geq 1$ blocks

$$\mathbb{R}^m = \mathbb{R}^{m_1} \times \mathbb{R}^{m_2} \times \dots \times \mathbb{R}^{m_b} \quad \text{with} \quad m = m_1 + m_2 + \dots + m_b.$$

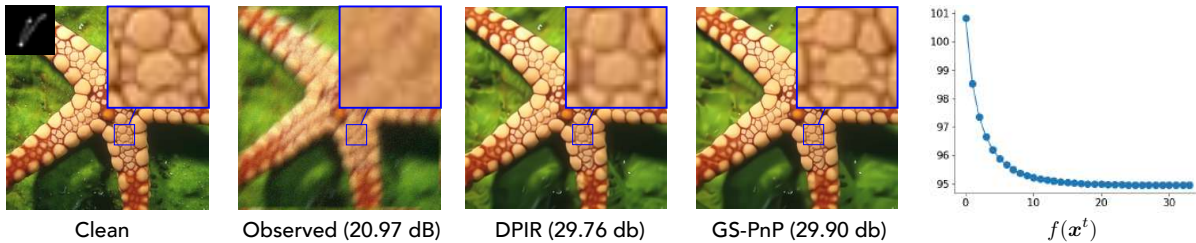


Figure 3: The strategy of explicitly parametrizing the regularizer h_θ in GS-PnP provides a competitive performance relative to DPIR, which is a state-of-the-art PnP method [17].

In this setting, the data fidelity term is given by

$$g(\mathbf{x}) = \frac{1}{b} \sum_{i=1}^b g_i(\mathbf{x})$$

where each g_i is evaluated only on the subset $\mathbf{y}_i \in \mathbb{R}^{m_i}$ of the full measurement vector $\mathbf{y} \in \mathbb{R}^m$. For example, each individual term can be set to the quadratic function $g_i(\mathbf{x}) = (1/2)\|\mathbf{y}_i - \mathbf{A}_i\mathbf{x}\|_2^2$, where \mathbf{A}_i is the operator corresponding to the measurement block \mathbf{y}_i . Online PnP algorithms improve the scalability to large-scale measurements by using only a single component gradient $\nabla g_i(\mathbf{x})$ or a single component proximal operator $\text{prox}_{\gamma g_i}(\mathbf{x})$ with $i \in \{1, \dots, b\}$, making their per-iteration complexity independent of b .

Online PnP algorithms can be implemented using different block selection rules. The strategy commonly adopted for the theoretical analysis focuses on selecting indices i_t as independent identically distributed random variables distributed uniformly over $\{1, \dots, b\}$. An alternative would be to proceed in epochs of b consecutive iterations, where the set $\{1, \dots, b\}$ is reshuffled at the start of each epoch and the index i_t is selected from this ordered set.

Online PnP-PGM

Input: An initial value $\mathbf{x}^0 \in \mathbb{R}^n$ and a step-size $\gamma > 0$.

For $t = 1, 2, 3, \dots$, do

Choose an index $i_t \in \{1, \dots, b\}$

$$\mathbf{x}^t \leftarrow \mathbf{R}_\theta(\mathbf{x}^{t-1} - \gamma \nabla g_{i_t}(\mathbf{x}^{t-1}))$$

Incremental plug-and-play ADMM (IPA)

Input: An initial values $\mathbf{x}^0 \in \mathbb{R}^n$ and $\mathbf{s}^0 = \mathbf{0}$, and a parameter $\gamma > 0$.

Iterate: For $t = 1, 2, 3, \dots$, do

Choose an index $i_t \in \{1, \dots, b\}$

$$\mathbf{z}^t \leftarrow \text{prox}_{\gamma g_{i_t}}(\mathbf{x}^{t-1} - \mathbf{s}^{t-1})$$

$$\mathbf{x}^t \leftarrow \mathbf{R}_\theta(\mathbf{z}^t + \mathbf{s}^{t-1})$$

$$\mathbf{s}^t \leftarrow \mathbf{s}^{t-1} + (\mathbf{z}^t - \mathbf{x}^t).$$

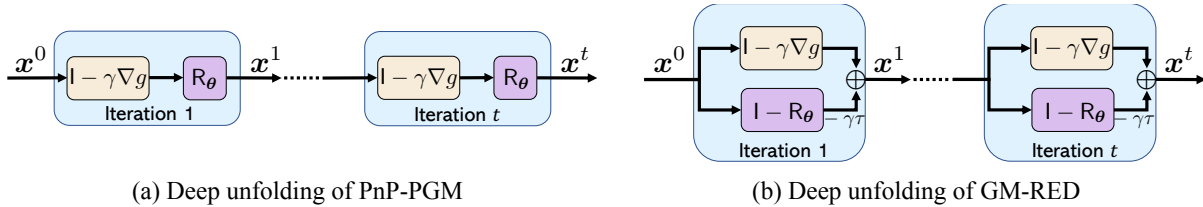


Figure 4: PnP is related to two other popular computational-imaging paradigms, *deep unfolding* (DU) and *deep equilibrium models* (DEQ). PnP-PGM or GM-RED can be turned into a DU architecture by truncating the algorithm to $t \geq 1$ iterations and training the weights θ of the network \mathbf{R}_θ end-to-end. Similarly, a DEQ architecture can be obtained by running the PnP algorithm until convergence and using the implicit differentiation at the fixed point to train θ . The operator \mathbf{R}_θ in DU/DEQ is not necessarily an AWGN denoiser, instead it is an *artifact-removal* (AR) operator trained to remove artifacts specific to the PnP iterations.

Deep unfolding and deep equilibrium learning

Deep unfolding (DU) (also known as *deep unrolling* or *algorithm unrolling*) is a DL paradigm that has gained popularity in computational imaging due to its ability to provide a systematic connection between iterative algorithms and deep neural network architectures. PnP algorithms can be easily turned into DU architectures by parameterizing the operator R_θ as a CNN with weights θ , truncating the PnP algorithm to a fixed number of iterations $t \geq 1$, and training the corresponding architecture end-to-end in a supervised fashion. For example, Figure 4 illustrates the representation of t iterations of PnP-ISTA and RED-SD as DU architectures.

Consider a set of paired data (x_i, y_i) , where x_i is the desired “ground truth” image and $y_i = Ax_i + e_i$ is its noisy observation. Consider also the iterate $x_i^t(\theta)$ of a PnP algorithm truncated to $t \geq 1$ iterations, where we made explicit the dependence of the PnP output on the weights θ of the deep network parameterizing R_θ . DU interprets the steps required for mapping the input vector y_i and the initialization x_i^0 to the output $x_i^t(\theta)$ as layers of a deep neural network architecture. The DU training is performed by solving

$$\hat{\theta} \in \arg \min_{\theta} \sum_i \mathcal{L}(x_i, x_i^t(\theta)), \quad (9)$$

where \mathcal{L} is a loss function that quantifies the discrepancy between the true and predicted solutions. Once trained using (9), the truncated PnP algorithm can be directly used for imaging [22].

Deep equilibrium models (DEQ) is an extension of DU to an arbitrary number of iterations [23]. DEQ can be implemented by replacing $x_i^t(\theta)$ in (9) by a fixed-point $\bar{x}_i(\theta)$ of a given PnP algorithm and using implicit differentiation for updating the weights θ . The benefit of DEQ over DU is that it doesn’t require the storage of the intermediate variables for solving (9), thus reducing the memory complexity of training. However, DEQ requires the computation of the fixed-point $\bar{x}_i(\theta)$, which can increase the computational complexity.

There are some important differences between traditional PnP and DU/DEQ. Traditional PnP relies on an AWGN denoiser as an image prior. On the other hand, the operator R_θ in DU/DEQ is not an AWGN denoiser; instead, it is an *artifact-removal (AR)* operator trained to remove artifacts specific to the PnP iterations. As seen in Fig. 5, which shows the relative performance of PnP using an AWGN denoiser and using a pre-trained AR operator, this problem-specific training can yield significantly improved results. However, this performance comes at a cost; while the prior in traditional PnP is *fully decoupled* from the measurement operator, that of DU/DEQ is trained by accounting for the measurement operator A . Hence the DU/DEQ approach has reduced generality and higher computational/memory complexity of training, since the AR prior is trained for the specific task of reconstruction from random projections rather than for AWGN denoising.

Monotone operator theory

Monotone operator theory [24] provides a set of mathematical tools for analyzing convergence of PnP, DU with weight sharing, and DEQ in a unified fashion. The approach is to view the algorithm as a fixed-point

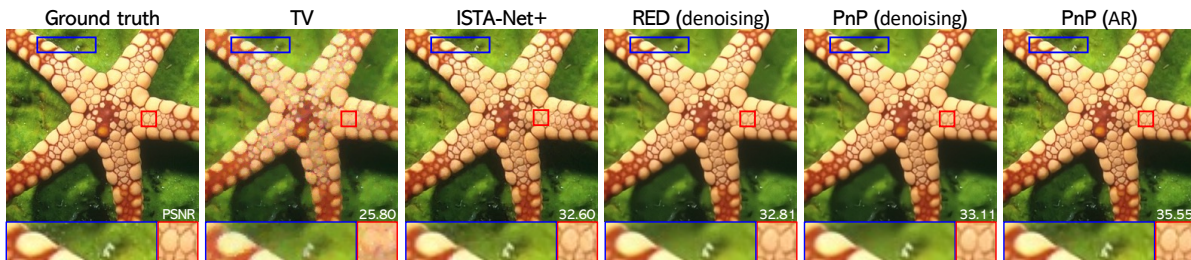


Figure 5: Visual evaluation of color image recovery in compressive sensing from random projections with 20% subsampling. The results of TV and a well-known DU architecture ISTA-Net+ are provided for reference. The methods PnP (denoising) and RED (denoising) use a pre-trained AWGN denoiser as an image prior. The method PnP (AR) uses a problem-dependent AR operator pre-trained using DU. Note that the choice of denoiser affects the reconstruction significantly (PSNR shown in white).

iteration of a contractive or nonexpansive operator.

Example. PnP-ADMM is equivalent to running the following fixed-point iteration

$$\mathbf{v}^t \leftarrow \mathbf{T}_\theta(\mathbf{v}^{t-1}) \quad \text{with} \quad \mathbf{T}_\theta := \frac{1}{2}\mathbf{I} + \frac{1}{2}(2\text{prox}_{\gamma g} - 1)(2\mathbf{R}_\theta - 1),$$

where \mathbf{I} is the identity and $\text{prox}_{\gamma g}$ is the proximal operator of g . Similarly, PnP-PGM is equivalent to running the following fixed-point iteration

$$\mathbf{v}^t \leftarrow \mathbf{F}_\theta(\mathbf{v}^{t-1}) \quad \text{with} \quad \mathbf{F}_\theta := \mathbf{R}_\theta(\mathbf{I} - \gamma \nabla g),$$

where ∇g is the gradient of g .

Definition 4. An operator \mathbf{T} is *Lipschitz continuous* with constant $\lambda > 0$ if

$$\|\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y}\|_2 \leq \lambda \|\mathbf{x} - \mathbf{y}\|_2, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

When $\lambda = 1$, we say that \mathbf{T} is *nonexpansive*. When $\lambda < 1$, we say that \mathbf{T} is a *contraction*.

Definition 5. \mathbf{T} is *monotone* if

$$(\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \geq 0, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

We say that it is *strongly monotone* or *coercive* with parameter $\mu > 0$ if

$$(\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \geq \mu \|\mathbf{x} - \mathbf{y}\|^2, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

Definition 6. \mathbf{T} is *cocoercive* with constant $\beta > 0$ if

$$(\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \geq \beta \|\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y}\|^2, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

When $\beta = 1$, we say that \mathbf{T} is *firmly nonexpansive*.

Definition 7. For a constant $\alpha \in (0, 1)$, we say that \mathbf{T} is α -averaged, if there exists a nonexpansive operator \mathbf{N} such that $\mathbf{T} = (1 - \alpha)\mathbf{I} + \alpha\mathbf{N}$.

Using the definitions above, one can show the following

Proposition 16. Let $g : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. Then, we have that

$$\nabla g \text{ is Lipschitz continuous with constant } L \quad \Leftrightarrow \quad \nabla g \text{ is cocoercive with constant } (1/L).$$

Proof. This proof will be left to the students as part of the exercises. ■

Example. Let $g : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function with a L -Lipschitz continuous gradient and \mathbf{R}_θ be a firmly-nonexpansive operator. What can we say about the convergence of PnP-PGM and PnP-ADMM?

Appendix: Monotone operators

The following results are derived from the definition above.

Proposition 17. Consider $R = I - T$ where $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$.

T is nonexpansive $\Leftrightarrow R$ is $(1/2)$ -cocoercive.

Proof. First suppose that R is $1/2$ cocoercive. Let $\mathbf{h} := \mathbf{x} - \mathbf{y}$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. We then have

$$\frac{1}{2} \|\mathbf{R}\mathbf{x} - \mathbf{R}\mathbf{y}\|^2 \leq (\mathbf{R}\mathbf{x} - \mathbf{R}\mathbf{y})^\top \mathbf{h} = \|\mathbf{h}\|^2 - (\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y})^\top \mathbf{h}.$$

We also have that

$$\frac{1}{2} \|\mathbf{R}\mathbf{x} - \mathbf{R}\mathbf{y}\|^2 = \frac{1}{2} \|\mathbf{h}\|^2 - (\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y})^\top \mathbf{h} + \frac{1}{2} \|\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y}\|^2.$$

By combining these two and simplifying the expression

$$\|\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y}\| \leq \|\mathbf{h}\|.$$

The converse can be proved by following this logic in reverse. ■

Proposition 18. Consider $R = I - T$ where $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$.

T is Lipschitz continuous with constant $\lambda < 1 \Rightarrow R$ is $(1 - \lambda)$ -strongly monotone.

Proof. By using the Cauchy-Schwarz inequality, we have for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

$$\begin{aligned} (\mathbf{R}\mathbf{x} - \mathbf{R}\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) &= \|\mathbf{x} - \mathbf{y}\|^2 - (\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \\ &\geq \|\mathbf{x} - \mathbf{y}\|^2 - \|\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y}\| \|\mathbf{x} - \mathbf{y}\| \geq \|\mathbf{x} - \mathbf{y}\|^2 - \lambda \|\mathbf{x} - \mathbf{y}\|^2 \geq (1 - \lambda) \|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$
■

The following characterization is often convenient.

Proposition 19. For a nonexpansive operator T , a constant $\alpha \in (0, 1)$, and the operator $R := I - T$, the following are equivalent

- (a) T is α -averaged
- (b) $(1 - 1/\alpha)I + (1/\alpha)T$ is nonexpansive
- (c) $\|\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y}\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2 - \left(\frac{1-\alpha}{\alpha}\right) \|\mathbf{R}\mathbf{x} - \mathbf{R}\mathbf{y}\|^2, \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$

Proof. See Proposition 4.35 in [25]. ■

Proposition 20. Consider $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\beta > 0$. Then, the following are equivalent

- (a) T is β -cocoercive
- (b) βT is firmly nonexpansive
- (c) $I - \beta T$ is firmly nonexpansive.
- (d) βT is $(1/2)$ -averaged.
- (e) $I - 2\beta T$ is nonexpansive.

Proof. For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, let $\mathbf{h} := \mathbf{x} - \mathbf{y}$. The equivalence between (a) and (b) is readily observed by defining $\mathbf{P} := \beta \mathbf{T}$ and noting that

$$\begin{aligned} (\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y})^\top \mathbf{h} &= \beta(\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y})^\top \mathbf{h} \quad \text{and} \\ \|\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y}\|^2 &= \beta^2 \|\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y}\|^2. \end{aligned}$$

Define $\mathbf{R} := \mathbf{I} - \mathbf{P}$ and suppose (b) is true, then

$$\begin{aligned} &(\mathbf{R}\mathbf{x} - \mathbf{R}\mathbf{y})^\top \mathbf{h} \\ &= \|\mathbf{h}\|^2 - (\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y})^\top \mathbf{h} \\ &= \|\mathbf{R}\mathbf{x} - \mathbf{R}\mathbf{y}\|^2 + (\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y})^\top \mathbf{h} - \|\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y}\|^2 \\ &\geq \|\mathbf{R}\mathbf{x} - \mathbf{R}\mathbf{y}\|^2. \end{aligned}$$

By repeating the same argument for $\mathbf{P} = \mathbf{I} - \mathbf{R}$, we establish the full equivalence between (b) and (c).

The equivalence of (b) and (d) can be seen by noting that

$$\begin{aligned} &2\|\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y}\|^2 \leq 2(\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y})^\top \mathbf{h} \\ \Leftrightarrow &\|\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y}\|^2 \leq 2(\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y})^\top \mathbf{h} - \|\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y}\|^2 \\ &= \|\mathbf{h}\|^2 - (\|\mathbf{h}\|^2 - 2(\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y})^\top \mathbf{h} + \|\mathbf{P}\mathbf{x} - \mathbf{P}\mathbf{y}\|^2) \\ &= \|\mathbf{h}\|^2 - \|\mathbf{R}\mathbf{x} - \mathbf{R}\mathbf{y}\|^2. \end{aligned}$$

To show the equivalence with (e), first suppose that $\mathbf{N} := \mathbf{I} - 2\mathbf{P}$ is nonexpansive, then $\mathbf{P} = \frac{1}{2}(\mathbf{I} + (-\mathbf{N}))$ is 1/2-averaged, which means that it is firmly nonexpansive. On the other hand, if \mathbf{P} is firmly nonexpansive, then it is 1/2-averaged, which means that from Proposition 19(b) we have that $(1 - 2)\mathbf{I} + 2\mathbf{P} = 2\mathbf{P} - \mathbf{I} = -\mathbf{N}$ is nonexpansive. This directly means that \mathbf{N} is nonexpansive. ■

Proposition 21. Let f be a proper, closed, and convex function. Then for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $\mathbf{g} \in \partial f(\mathbf{x})$, and $\mathbf{h} \in \partial f(\mathbf{y})$, ∂f is a monotone operator

$$(\mathbf{g} - \mathbf{h})^\top (\mathbf{x} - \mathbf{y}) \geq 0.$$

Additionally if f is strongly convex with constant $\mu > 0$, then ∂f is strongly monotone with the same constant.

$$(\mathbf{g} - \mathbf{h})^\top (\mathbf{x} - \mathbf{y}) \geq \mu \|\mathbf{x} - \mathbf{y}\|^2.$$

Proof. Consider a strongly convex function f with a constant $\mu \geq 0$. Then, we have that

$$\begin{aligned} &\begin{cases} f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2 \\ f(\mathbf{x}) \geq f(\mathbf{y}) + \mathbf{h}^\top (\mathbf{x} - \mathbf{y}) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2 \end{cases} \\ \Rightarrow &(\mathbf{g} - \mathbf{h})^\top (\mathbf{x} - \mathbf{y}) \geq \mu \|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

The proof for a weakly convex f is obtained by considering $\mu = 0$ in the inequalities above. ■

Appendix: Subgradient and subdifferential

This chapter reviews the concept of *subgradient* and *subdifferential* that generalize that of the gradient.

Definition and technicalities

We start by considering extended real-valued functions and defining some useful properties.

Definition 8. We define the *domain* of function f as the set

$$\text{dom}(f) = \{x \in \mathbb{R}^n : f(x) < \infty\}.$$

We always assume that $\text{dom}(f) \neq \emptyset$.

Definition 9. We say that a function $f : \mathbb{R}^n \rightarrow [-\infty, \infty]$ is *proper* if

- $\exists x \in \mathbb{R}^n$ such that $f(x) < \infty$ (*non-empty domain*)
- $\forall x \in \mathbb{R}^n$ we have $f(x) > -\infty$. (*never attains $-\infty$*)

Remark. When we use the term *proper*, we imply that the function takes values in the extended real-line $(-\infty, \infty]$ with at least one $x \in \mathbb{R}^n$ such that $f(x) < \infty$.

Definition 10. We say that a function f is *closed*, if for each $\theta \in \mathbb{R}$ the *sublevel set*

$$L_f(\theta) := \{x \in \mathbb{R}^n : f(x) \leq \theta\}$$

is closed.

Remark. We will denote the class of all convex, proper, and closed functions with $\Gamma^0(\mathbb{R}^n)$.

Definition 11. Let f be a proper function. The vector $g \in \mathbb{R}^n$ is a *subgradient* of f at $x \in \mathbb{R}^n$ if

$$f(y) \geq f(x) + g^\top(y - x), \quad \forall y \in \text{dom}(f).$$

We also define the *subdifferential* of f at $x \in \mathbb{R}^n$ as a set

$$\partial f(x) := \{g \in \mathbb{R}^n : f(y) \geq f(x) + g^\top(y - x), \forall y \in \text{dom}(f)\}.$$

We will use the notation $g(x)$ for highlighting the dependence of the subgradient on x .

Remark. The inequality in Definition 11 is known as the *subgradient inequality*. Thus, if f has a subgradient at x , then it has a linear lower bound at that location. Note that the inequality is trivial for $y \notin \text{dom}(f)$. When $x \notin \text{dom}(f)$, we define $\partial f(x) = \emptyset$, since the subgradient inequality does not hold for any $y \in \text{dom}(f)$. We will use the notation $g(x) \in \partial f(x)$ when we want to make it explicit that the subgradient is evaluated at x .

The next results establishes that the set $\partial f(x)$ is closed and convex.

Theorem 2. Let f be a proper function. Then, the set $\partial f(x)$ is closed and convex for any $x \in \mathbb{R}^n$.

Proof. For any $x \in \mathbb{R}^n$, the subdifferential can be represented as

$$\partial f(x) = \bigcap_{y \in \mathbb{R}^n} H_y \quad \text{with} \quad H_y = \{g \in \mathbb{R}^n : f(y) \geq f(x) + g^\top(y - x)\}.$$

Since each H_y is a half-space, which means that it is closed and convex, the set $\partial f(x)$ is closed and convex since it is an intersection of closed and convex sets. ■

However, $\partial f(x)$ may be empty at x , in which case we say that f is *not* subdifferentiable there.

Definition 12. A proper function f is *subdifferentiable* at $x \in \mathbb{R}^n$ if $\partial f(x) \neq \emptyset$. We will denote the collection of points of subdifferentiability by

$$\text{dom}(\partial f) := \{x \in \mathbb{R}^n : \partial f(x) \neq \emptyset\}.$$

Convex functions are not necessarily subdifferentiable in their domain; however, they are subdifferentiable at any point in the *interior* of their domain.

Theorem 3. If f be proper and convex, then $\partial f(x)$ is nonempty and bounded for any $x \in \text{int}(\text{dom}(f))$.

Proof. This proof is highly technical, but can be found in standard textbooks on convex optimization. ■

Remark. Consider an open ball centered at $x \in \mathbb{R}^n$ of radius $\epsilon > 0$

$$\mathcal{B}(x, \epsilon) := \{y \in \mathbb{R}^n : \|x - y\|_2 < \epsilon\}$$

A point $x \in \mathbb{R}^n$ is an *interior point* of $\text{dom}(f)$ if

$$\exists \epsilon > 0 \quad \text{such that} \quad \mathcal{B}(x, \epsilon) \subset \text{dom}(f),$$

which basically means that there exists an open ball centered at x fully contained in $\text{dom}(f)$.

The direct consequence of Theorem 3 is that real-valued convex functions are subdifferentiable everywhere.

Corollary 2.

$$f : \mathbb{R}^n \rightarrow \mathbb{R} \text{ is a convex function} \quad \Rightarrow \quad f \text{ is subdifferentiable over } \mathbb{R}^n.$$

When f is differentiable, the subdifferential and gradient are equivalent.

Theorem 4. Let f be a proper convex function and let $x \in \text{int}(\text{dom}(f))$. Then

$$f \text{ is differentiable at } x \quad \Leftrightarrow \quad \partial f(x) = \{\nabla f(x)\}.$$

Proof. We prove the result only for real-valued functions and in the \Rightarrow direction. Let us consider an arbitrary $g \in \partial f(x)$. Then, for all $u \in \mathbb{R}^n$ and $\epsilon > 0$, we have from the definition of the subgradient

$$f(x + \epsilon u) \geq f(x) + \epsilon g^\top u \quad \Rightarrow \quad \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon u) - f(x)}{\epsilon} := \nabla f(x)^\top u \geq g^\top u \quad \Rightarrow \quad (\nabla f(x) - g)^\top u \geq 0,$$

where we used the definition of the directional derivative. By choosing $u = g - \nabla f(x)$, we obtain

$$-\|\nabla f(x) - g\|_2^2 \geq 0 \quad \Rightarrow \quad \|\nabla f(x) - g\|_2^2 \leq 0 \quad \Rightarrow \quad g = \nabla f(x).$$

■

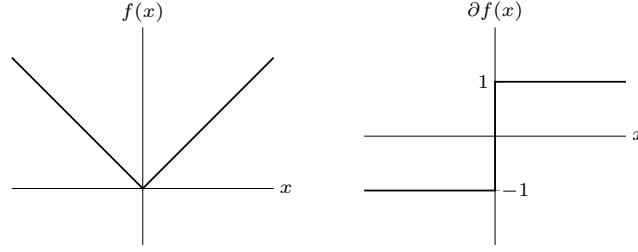


Figure 6: The only point at which $f(x) = |x|$ is not differentiable is $x = 0$, where we have $\partial f(x) = [-1, 1]$.

Example 1. Consider the function $f(x) = |x|$ with $x \in \mathbb{R}$. The only point at which f is not differentiable is $x = 0$. At this point, the subgradients of f are characterized by the inequality

$$f(x) - f(0) \geq g \cdot x \Rightarrow |x| \geq gx \Rightarrow g \in [-1, 1].$$

For $x \neq 0$, the subderivative coincides with the derivative, which means that

$$\partial f(x) = \begin{cases} 1 & \text{if } x > 0 \\ [-1, 1] & \text{if } x = 0 \\ -1 & \text{if } x < 0 \end{cases}.$$

We can also generalize the notion of the *directional derivative* of f .

Theorem 5. Let f be a proper convex function. Then for any $x \in \text{int}(\text{dom}(f))$ and $d \in \mathbb{R}^n$

$$\nabla_d f(x) = \lim_{\gamma \rightarrow 0^+} \frac{f(x + \gamma d) - f(x)}{\gamma} = \max\{g^\top d : g \in \partial f(x)\}.$$

Proof. We will show the proof for “ \geq ” and omit the proof for “ \leq ”. Let $x \in \text{int}(\text{dom}(f))$ and $d \in \mathbb{R}^n$. From the subgradient inequality we have that for any $g \in \partial f(x)$

$$\nabla_d f(x) = \lim_{\gamma \rightarrow 0^+} \frac{f(x + \gamma d) - f(x)}{\gamma} \geq \lim_{\gamma \rightarrow 0^+} g^\top d = g^\top d \Rightarrow \nabla_d f(x) \geq \max\{g^\top d : g \in \partial f(x)\}.$$

■

References

- [1] U. S. Kamilov, C. A. Bouman, G. T. Buzzard, and B. Wohlberg, “Plug-and-play methods for integrating physical and learned models in computational imaging,” *IEEE Signal Process. Mag.*, vol. 40, no. 1, pp. 85–97, Jan. 2023.
- [2] J. J. Moreau, “Proximité et dualité dans un espace hilbertien,” *Bull. Soc. Math. France*, vol. 93, pp. 273–299, 1965.
- [3] R. Gribonval and M. Nikolova, “A characterization of proximity operators,” *J. Math. Imaging Vis.*, vol. 62, pp. 773–789, 2020.
- [4] R. Gribonval, “Should penalized least squares regression be interpreted as maximum a posteriori estimation?” *IEEE Trans. Signal Process.*, vol. 59, no. 5, pp. 2405–2410, May 2011.
- [5] R. Gribonval and P. Machart, “Reconciling “priors” & “priors” without prejudice?” in *Proc. Advances in Neural Information Processing Systems 26*, Lake Tahoe, NV, USA, December 5–10, 2013, pp. 2193–2201.

- [6] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, “Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising,” *IEEE Trans. Image Process.*, vol. 26, no. 7, pp. 3142–3155, July 2017.
- [7] K. Zhang, Y. Li, W. Zuo, L. Zhang, L. Van Gool, and R. Timofte, “Plug-and-play image restoration with deep denoiser prior,” *IEEE Trans. Patt. Anal. and Machine Intell.*, vol. 44, no. 10, pp. 6360–6376, Oct. 2022, doi: 10.1109/tpami.2021.3088914.
- [8] N. Parikh and S. Boyd, “Proximal algorithms,” *Foundations and Trends in Optimization*, vol. 1, no. 3, pp. 123–231, 2014.
- [9] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, “Distributed optimization and statistical learning via the alternating direction method of multipliers,” *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [10] S. V. Venkatakrishnan, C. A. Bouman, and B. Wohlberg, “Plug-and-play priors for model based reconstruction,” in *Proc. IEEE Global Conf. Signal Process. and Inf. Process. (GlobalSIP)*, Austin, TX, USA, December 3-5, 2013, pp. 945–948.
- [11] S. H. Chan, X. Wang, and O. A. Elgendy, “Plug-and-play ADMM for image restoration: Fixed-point convergence and applications,” *IEEE Trans. Comp. Imag.*, vol. 3, no. 1, pp. 84–98, March 2017.
- [12] C. Park, S. Shoushtari, W. Gan, and U. S. Kamilov, “Convergence of nonconvex pnp-admm with mmse denoisers,” in *Proc. Int. Workshop on Computational Advances in Multi-Sensor Adaptive Process. (CAMSAP 2023)*, Los Suenos, Costa Rica, Dec. 2023, pp. 511–515.
- [13] Y. Nesterov, *Introductory Lectures on Convex Optimization: A Basic Course*. Kluwer Academic Publishers, 2004.
- [14] Y. Romano, M. Elad, and P. Milanfar, “The little engine that could: Regularization by denoising (RED),” *SIAM J. Imaging Sci.*, vol. 10, no. 4, pp. 1804–1844, 2017.
- [15] E. T. Reehorst and P. Schniter, “Regularization by denoising: Clarifications and new interpretations,” *IEEE Trans. Comput. Imag.*, vol. 5, no. 1, pp. 52–67, Mar. 2019.
- [16] Y. Sun, J. Liu, and U. S. Kamilov, “Block coordinate regularization by denoising,” in *Proc. Advances in Neural Information Processing Systems 32*, Vancouver, BC, Canada, Dec. 2019, pp. 382–392.
- [17] S. Hurault, A. Leclaire, and N. Papadakis, “Gradient step denoiser for convergent plug-and-play,” in *International Conference on Learning Representations (ICLR)*, Kigali, Rwanda, May 1-5, 2022.
- [18] R. Cohen, Y. Blau, D. Freedman, and E. Rivlin, “It has potential: Gradient-driven denoisers for convergent solutions to inverse problems,” in *Proc. Advances in Neural Information Processing Systems 34*, December 6-14, 2021, pp. 18 152–18 164.
- [19] Y. Sun, B. Wohlberg, and U. S. Kamilov, “An online plug-and-play algorithm for regularized image reconstruction,” *IEEE Trans. Comput. Imag.*, vol. 5, no. 3, pp. 395–408, Sep. 2019.
- [20] Z. Wu, Y. Sun, A. Matlock, J. Liu, L. Tian, and U. S. Kamilov, “SIMBA: Scalable inversion in optical tomography using deep denoising priors,” *IEEE J. Sel. Topics Signal Process.*, vol. 14, no. 6, pp. 1163–1175, 2020.
- [21] Y. Sun, Z. Wu, X. Xu, B. Wohlberg, and U. S. Kamilov, “Scalable plug-and-play ADMM with convergence guarantees,” *IEEE Trans. Comput. Imag.*, vol. 7, pp. 849–863, Jul. 2021.
- [22] J. Liu, S. Asif, B. Wohlberg, and U. S. Kamilov, “Recovery analysis for plug-and-play priors using the restricted eigenvalue condition,” in *Proc. Advances in Neural Information Processing Systems 34*, December 6-14, 2021, pp. 5921–5933.
- [23] D. Gilton, G. Ongie, and R. Willett, “Deep equilibrium architectures for inverse problems in imaging,” *IEEE Trans. Comput. Imag.*, vol. 7, pp. 1123–1133, Oct. 2021.

- [24] H. H. Bauschke and P. L. Combettes, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, 2010.
- [25] —, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, 2nd ed. Springer, 2017.