

Supplementary information for
**Iterative projection meets sparsity
regularization: towards practical
single-shot quantitative phase
imaging with in-line holography**

Yunhui Gao¹ and Liangcai Cao^{1*}

¹State Key Laboratory of Precision Measurement Technology and
Instruments, Department of Precision Instruments, Tsinghua
University, Beijing, 100084, China.

*Corresponding author: clc@tsinghua.edu.cn

1 Proof of duality

The CCTV-regularized denoising problem can be equivalently expressed as a constrained optimization problem as follows:

$$\min_{\mathbf{x}} \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|_2^2 + \lambda \|\mathbf{u}\|_1 + I_C(\mathbf{x}), \quad \text{subject to } \mathbf{u} = \mathbf{D}\mathbf{x}, \quad (\text{S1})$$

where $\mathbf{u} \in \mathbb{C}^{2n}$ is an auxiliary variable. Notice that the problem of Eq. (S1) is a convex optimization problem with an affine equality constraint, for which strong duality holds [1]. The Lagrangian is given by

$$L(\mathbf{x}, \mathbf{u}, \mathbf{w}) = \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|_2^2 + \lambda \|\mathbf{u}\|_1 + I_C(\mathbf{x}) + \text{Re}(\langle \mathbf{w}, \mathbf{D}\mathbf{x} - \mathbf{u} \rangle), \quad (\text{S2})$$

where $\mathbf{w} \in \mathbb{C}^{2n}$ is the dual variable, $\langle \cdot, \cdot \rangle$ denotes the inner product of two vectors, and $\text{Re}(\cdot)$ extracts the real part of a complex number. The Lagrange

dual function, by definition, is

$$\begin{aligned}
& \inf_{\mathbf{x}, \mathbf{u}} L(\mathbf{x}, \mathbf{u}, \mathbf{w}) \\
&= \inf_{\mathbf{x}, \mathbf{u}} \left\{ \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|_2^2 + \lambda \|\mathbf{u}\|_1 + I_C(\mathbf{x}) + \operatorname{Re}(\langle \mathbf{w}, \mathbf{D}\mathbf{x} - \mathbf{u} \rangle) \right\} \\
&= \inf_{\mathbf{x} \in C} \left\{ \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|_2^2 + \operatorname{Re}(\langle \mathbf{w}, \mathbf{D}\mathbf{x} \rangle) \right\} + \inf_{\mathbf{u}} \left\{ \lambda \|\mathbf{u}\|_1 - \operatorname{Re}(\langle \mathbf{w}, \mathbf{u} \rangle) \right\} \\
&\stackrel{(a)}{=} \inf_{\mathbf{x} \in C} \left\{ \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|_2^2 + \operatorname{Re}(\langle \mathbf{D}^H \mathbf{w}, \mathbf{x} \rangle) \right\} - I_S(\mathbf{w}) \\
&= \inf_{\mathbf{x} \in C} \left\{ \frac{1}{2\gamma} \|\mathbf{x} - (\mathbf{v} - \gamma \mathbf{D}^H \mathbf{w})\|_2^2 \right\} + \frac{1}{2\gamma} \|\mathbf{v}\|_2^2 - \frac{1}{2\gamma} \|\mathbf{v} - \gamma \mathbf{D}^H \mathbf{w}\|_2^2 - I_S(\mathbf{w}) \\
&= \|\mathcal{H}_C(\mathbf{v} - \gamma \mathbf{D}^H \mathbf{w})\|_2^2 + \frac{1}{2\gamma} \|\mathbf{v}\|_2^2 - \frac{1}{2\gamma} \|\mathbf{v} - \gamma \mathbf{D}^H \mathbf{w}\|_2^2 - I_S(\mathbf{w}), \tag{S3}
\end{aligned}$$

where (a) can be easily derived based on the fact that the convex conjugate of the ℓ_1 norm is the indicator function of $[-1, 1]^n$ [1]. The last equality in Eq. (S3), together with strong duality, suggests that the primal optimal solution \mathbf{x}^* is related to the dual optimal solution \mathbf{w}^* via $\mathbf{x}^* = \mathcal{P}_C(\mathbf{v} - \gamma \mathbf{D}^H \mathbf{w}^*)$. The dual problem is to maximize the Lagrange dual function with respect to \mathbf{w} , which is equivalent to Eq. (7) in the main text.

2 Proof of convergence

2.1 Preliminaries

Since we are primarily dealing with real-valued functions over complex-valued variables, we adopt the CR-calculus as helpful mathematical tool for analysis. The CR-calculus extends the complex derivative to the general non-analytic functions. Readers may refer to Ref. [2] for a detailed introduction. The CR-calculus regards the complex variable \mathbf{x} and its conjugate $\bar{\mathbf{x}}$ as independent variables. Thus, the fidelity function $F(\mathbf{x})$ should be interpreted as a function over the pair of conjugate vectors $\hat{\mathbf{x}} = [\mathbf{x}^T, \bar{\mathbf{x}}^T]^T \in \mathbb{C}^{2n}$. Nevertheless, to keep notations consistent, we still denote the function as $F(\mathbf{x})$. The same applies to other functions as well.

The followings are some intermediate results from matrix analysis, which are helpful for proving the convergence theorems below.

Lemma 1 [3] *Given matrices $\mathbf{P} \in \mathbb{C}^{n \times n}$, $\mathbf{Q} \in \mathbb{C}^{n \times n}$, and $\mathbf{R} \in \mathbb{C}^{n \times n}$. The following holds:*

1. $\mathbf{P} \succ \mathbf{Q} \Rightarrow \mathbf{R}^H \mathbf{P} \mathbf{R} \succ \mathbf{R}^H \mathbf{Q} \mathbf{R}$,
2. $\mathbf{P} \succeq \mathbf{Q} \succ \mathbf{0} \Rightarrow \mathbf{Q}^{-1} \succeq \mathbf{P}^{-1} \succ \mathbf{0}$.

Lemma 2 (Schur Complement [3]) *Given a $2n \times 2n$ Hermitian matrix:*

$$\mathbf{P} = \begin{pmatrix} \mathbf{P}_{11} & \mathbf{P}_{12} \\ \mathbf{P}_{21} & \mathbf{P}_{22} \end{pmatrix}, \quad (\text{S4})$$

where each block is of size $n \times n$, and we have $\mathbf{P}_{11}^{\text{H}} = \mathbf{P}_{11}$, $\mathbf{P}_{22}^{\text{H}} = \mathbf{P}_{22}$, and $\mathbf{P}_{12}^{\text{H}} = \mathbf{P}_{21}$. Then

$$\mathbf{P} \succ \mathbf{0} \Leftrightarrow \mathbf{P}_{11} \succ \mathbf{0} \quad \text{and} \quad \mathbf{P}_{22} - \mathbf{P}_{21} \mathbf{P}_{11}^{-1} \mathbf{P}_{12} \succ \mathbf{0}. \quad (\text{S5})$$

Lemma 3 *Given a matrix $\mathbf{P} \in \mathbb{C}^{n \times n}$, and a scalar $\varepsilon > 0$,*

$$\mathbf{P} \left(\varepsilon \mathbf{I} + \mathbf{P}^{\text{H}} \mathbf{P} \right)^{-1} \mathbf{P}^{\text{H}} \prec \mathbf{I}. \quad (\text{S6})$$

Proof Suppose the singular value decomposition of \mathbf{P} is given by $\mathbf{P} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\text{H}}$, where $\mathbf{U} \in \mathbb{C}^{n \times n}$ and $\mathbf{V} \in \mathbb{C}^{n \times n}$ are unitary matrices, and $\mathbf{\Sigma} = \text{diag}(\boldsymbol{\sigma})$ is a real-valued diagonal matrix. Then, we have

$$\varepsilon \mathbf{I} + \mathbf{P}^{\text{H}} \mathbf{P} = \varepsilon \mathbf{I} + \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^{\text{H}} = \mathbf{V} \text{diag} \left(\varepsilon \mathbf{1} + \boldsymbol{\sigma}^2 \right) \mathbf{V}^{\text{H}}. \quad (\text{S7})$$

That is, $\mathbf{P}^{\text{H}} \mathbf{P}$ is diagonalizable with real-valued non-negative eigenvalues $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$. $\varepsilon \mathbf{I} + \mathbf{P}^{\text{H}} \mathbf{P}$ is nonsingular and its inverse is given by

$$\left(\varepsilon \mathbf{I} + \mathbf{P}^{\text{H}} \mathbf{P} \right)^{-1} = \mathbf{V} \text{diag} \left(\frac{\mathbf{1}}{\varepsilon \mathbf{1} + \boldsymbol{\sigma}^2} \right) \mathbf{V}^{\text{H}}. \quad (\text{S8})$$

Thus, we arrive at the result:

$$\mathbf{P} \left(\varepsilon \mathbf{I} + \mathbf{P}^{\text{H}} \mathbf{P} \right)^{-1} \mathbf{P}^{\text{H}} = \mathbf{U} \text{diag} \left(\frac{\boldsymbol{\sigma}^2}{\varepsilon \mathbf{1} + \boldsymbol{\sigma}^2} \right) \mathbf{U}^{\text{H}} \prec \mathbf{U} \mathbf{U}^{\text{H}} = \mathbf{I}. \quad (\text{S9})$$

□

2.2 Convergence of the proximal gradient algorithm

The Wirtinger derivatives of $F(\mathbf{x})$ with respect to \mathbf{x} and $\bar{\mathbf{x}}$ are given by [4]

$$\frac{\partial F(\mathbf{x})}{\partial \mathbf{x}} = \frac{1}{2} (|\mathbf{A}\mathbf{x}| - \mathbf{y})^{\text{H}} \text{diag} \left(\frac{\overline{\mathbf{A}\mathbf{x}}}{|\mathbf{A}\mathbf{x}|} \right) \mathbf{A}, \quad (\text{S10})$$

$$\frac{\partial F(\mathbf{x})}{\partial \bar{\mathbf{x}}} = \frac{1}{2} (|\mathbf{A}\mathbf{x}| - \mathbf{y})^{\text{T}} \text{diag} \left(\frac{\mathbf{A}\mathbf{x}}{|\mathbf{A}\mathbf{x}|} \right) \bar{\mathbf{A}}. \quad (\text{S11})$$

Let $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m]^{\text{H}}$ where $\mathbf{a}_i \in \mathbb{C}^n$ denotes the i -th sampling vector. It should be noted that, the Wirtinger derivatives are not well-defined for $\mathbf{x} \in Z$ where Z is defined as

$$Z \stackrel{\text{def}}{=} \{ \mathbf{x} \in \mathbb{C}^n : \exists 1 \leq i \leq M, \text{ s.t. } \mathbf{a}_i^{\text{H}} \mathbf{x} = 0, \mathbf{a}_i \neq \mathbf{0} \}. \quad (\text{S12})$$

For any $\mathbf{x} \in \mathbb{C}^n \setminus Z$, the complex Hessian is defined as

$$\nabla^2 F(\mathbf{x}) = \mathbf{H}_{\hat{\mathbf{x}}\hat{\mathbf{x}}} = \begin{pmatrix} \mathbf{H}_{\mathbf{x}\mathbf{x}} & \mathbf{H}_{\bar{\mathbf{x}}\mathbf{x}} \\ \mathbf{H}_{\mathbf{x}\bar{\mathbf{x}}} & \mathbf{H}_{\bar{\mathbf{x}}\bar{\mathbf{x}}} \end{pmatrix}, \quad (\text{S13})$$

where the four second-order partial derivatives are calculated as follows:

$$\begin{aligned}
\mathbf{H}_{\mathbf{x}\mathbf{x}} &= \frac{\partial}{\partial \mathbf{x}} \left(\frac{\partial F(\mathbf{x})}{\partial \mathbf{x}} \right)^{\text{H}} \\
&= \frac{1}{2} \frac{\partial}{\partial \mathbf{x}} \left(\mathbf{A}^{\text{H}} \text{diag} \left(\frac{\mathbf{A}\mathbf{x}}{|\mathbf{A}\mathbf{x}|} \right) (|\mathbf{A}\mathbf{x}| - \mathbf{y}) \right) \\
&= \frac{1}{2} \frac{\partial}{\partial \mathbf{x}} \left(\mathbf{A}^{\text{H}} \mathbf{A}\mathbf{x} - \mathbf{A}^{\text{H}} \text{diag}(\mathbf{y}) \frac{\mathbf{A}\mathbf{x}}{|\mathbf{A}\mathbf{x}|} \right) \\
&= \frac{1}{2} \mathbf{A}^{\text{H}} \mathbf{A} - \frac{1}{2} \mathbf{A}^{\text{H}} \text{diag}(\mathbf{y}) \frac{\partial}{\partial \mathbf{x}} \left(\frac{\mathbf{A}\mathbf{x}}{|\mathbf{A}\mathbf{x}|} \right) \\
&= \frac{1}{2} \mathbf{A}^{\text{H}} \mathbf{A} - \frac{1}{4} \mathbf{A}^{\text{H}} \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right) \mathbf{A}, \tag{S14}
\end{aligned}$$

$$\begin{aligned}
\mathbf{H}_{\bar{\mathbf{x}}\mathbf{x}} &= \frac{\partial}{\partial \bar{\mathbf{x}}} \left(\frac{\partial F(\mathbf{x})}{\partial \mathbf{x}} \right)^{\text{H}} \\
&= \frac{1}{2} \frac{\partial}{\partial \bar{\mathbf{x}}} \left(\mathbf{A}^{\text{H}} \mathbf{A}\mathbf{x} - \mathbf{A}^{\text{H}} \text{diag}(\mathbf{y}) \frac{\mathbf{A}\mathbf{x}}{|\mathbf{A}\mathbf{x}|} \right) \\
&= \mathbf{0} - \frac{1}{2} \mathbf{A}^{\text{H}} \text{diag}(\mathbf{y}) \frac{\partial}{\partial \bar{\mathbf{x}}} \left(\frac{\mathbf{A}\mathbf{x}}{|\mathbf{A}\mathbf{x}|} \right) \\
&= \frac{1}{4} \mathbf{A}^{\text{H}} \text{diag}(\mathbf{y}) \text{diag} \left(\frac{(\mathbf{A}\mathbf{x})^2}{|\mathbf{A}\mathbf{x}|^3} \right) \bar{\mathbf{A}}, \tag{S15}
\end{aligned}$$

$$\begin{aligned}
\mathbf{H}_{\mathbf{x}\bar{\mathbf{x}}} &= \frac{\partial}{\partial \mathbf{x}} \left(\frac{\partial F(\mathbf{x})}{\partial \bar{\mathbf{x}}} \right)^{\text{H}} = \mathbf{H}_{\bar{\mathbf{x}}\mathbf{x}}^{\text{H}} \\
&= \frac{1}{4} \mathbf{A}^{\text{T}} \text{diag}(\mathbf{y}) \text{diag} \left(\frac{(\overline{\mathbf{A}\mathbf{x}})^2}{|\mathbf{A}\mathbf{x}|^3} \right) \mathbf{A}, \tag{S16}
\end{aligned}$$

$$\begin{aligned}
\mathbf{H}_{\bar{\mathbf{x}}\bar{\mathbf{x}}} &= \frac{\partial}{\partial \bar{\mathbf{x}}} \left(\frac{\partial F(\mathbf{x})}{\partial \bar{\mathbf{x}}} \right)^{\text{H}} = \mathbf{H}_{\mathbf{x}\mathbf{x}}^{\text{T}} \\
&= \frac{1}{2} \mathbf{A}^{\text{T}} \bar{\mathbf{A}} - \frac{1}{4} \mathbf{A}^{\text{T}} \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right) \bar{\mathbf{A}}. \tag{S17}
\end{aligned}$$

We now prove that the gradient of $F(\mathbf{x})$ is upper Lipschitz bounded by a constant. This is a particularly useful property of the amplitude-based fidelity term, enabling us to use prespecified algorithm parameters while ensuring convergence.

Lemma 4 *For any $\mathbf{x} \in \mathbb{C}^n \setminus \mathcal{Z}$, the Lipschitz constant for the gradient of the data-fidelity function $\nabla F(\mathbf{x})$ is bounded above by $(1/2)\rho(\mathbf{A}^{\text{H}}\mathbf{A})$.*

Proof We only need to prove that for any $\tau > (1/2)\rho(\mathbf{A}^{\text{H}}\mathbf{A})$, we have

$$\mathbf{G} \equiv \tau \mathbf{I} - \mathbf{H}_{\hat{\mathbf{x}}\hat{\mathbf{x}}} = \begin{pmatrix} \tau \mathbf{I} - \mathbf{H}_{\mathbf{x}\mathbf{x}} & -\mathbf{H}_{\bar{\mathbf{x}}\mathbf{x}} \\ -\mathbf{H}_{\mathbf{x}\bar{\mathbf{x}}} & \tau \mathbf{I} - \mathbf{H}_{\bar{\mathbf{x}}\bar{\mathbf{x}}} \end{pmatrix} \succ \mathbf{0}. \tag{S18}$$

Let $\varepsilon = \tau - (1/2)\rho(\mathbf{A}^H \mathbf{A}) > 0$ and denote $\mathbf{G}_{11}, \mathbf{G}_{12}, \mathbf{G}_{21}, \mathbf{G}_{22} \in \mathbb{C}^{n \times n}$ as the four block matrices of \mathbf{G} , we have

$$\begin{aligned} \mathbf{G}_{11} &= \left(\tau \mathbf{I} - \frac{1}{2} \mathbf{A}^H \mathbf{A} \right) + \frac{1}{4} \mathbf{A}^H \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right) \mathbf{A} \\ &\succ \varepsilon \mathbf{I} + \frac{1}{4} \mathbf{A}^H \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right) \mathbf{A}. \end{aligned} \quad (\text{S19})$$

According to Lemma 1(a) and 1(b), we have

$$\mathbf{G}_{21} \mathbf{G}_{11}^{-1} \mathbf{G}_{12} \prec \mathbf{G}_{21} \left(\varepsilon \mathbf{I} + \frac{1}{4} \mathbf{A}^H \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right) \mathbf{A} \right)^{-1} \mathbf{G}_{12} \quad (\text{S20})$$

Let $\mathbf{P} = (1/2) \text{diag} \left((\mathbf{y}/|\mathbf{A}\mathbf{x}|)^{1/2} \right) \mathbf{A}$ and use Lemma 1(b), we have

$$\begin{aligned} \mathbf{G}_{21} \mathbf{G}_{11}^{-1} \mathbf{G}_{12} &\prec \frac{1}{4} \mathbf{A}^T \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right)^{\frac{1}{2}} \text{diag} \left(\frac{(\overline{\mathbf{A}\mathbf{x}})^2}{|\mathbf{A}\mathbf{x}|^2} \right) \\ &\quad \times \mathbf{P} \left(\varepsilon \mathbf{I} + \mathbf{P}^H \mathbf{P} \right)^{-1} \mathbf{P}^H \text{diag} \left(\frac{(\mathbf{A}\mathbf{x})^2}{|\mathbf{A}\mathbf{x}|^2} \right) \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right)^{\frac{1}{2}} \bar{\mathbf{A}} \\ &\prec \frac{1}{4} \mathbf{A}^T \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right)^{\frac{1}{2}} \text{diag} \left(\frac{(\overline{\mathbf{A}\mathbf{x}})^2}{|\mathbf{A}\mathbf{x}|^2} \right) \text{diag} \left(\frac{(\mathbf{A}\mathbf{x})^2}{|\mathbf{A}\mathbf{x}|^2} \right) \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right)^{\frac{1}{2}} \bar{\mathbf{A}} \\ &\prec \frac{1}{4} \mathbf{A}^T \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right) \bar{\mathbf{A}} \\ &\prec \frac{1}{4} \text{diag} \left(\frac{\mathbf{y}}{|\mathbf{A}\mathbf{x}|} \right) \bar{\mathbf{A}} + \left(\lambda \mathbf{I} - \frac{1}{2} \mathbf{A}^T \bar{\mathbf{A}} \right) \\ &= \mathbf{G}_{22}. \end{aligned} \quad (\text{S21})$$

Therefore, according to Lemma 2, \mathbf{G} is positive-definite. This implies that for $\mathbf{x} \in \mathbb{C}^n \setminus Z$ the Lipschitz constant of ∇F is upper-bounded by $(1/2)\rho(\mathbf{A}^H \mathbf{A})$. \square

The above Lemma implies that the fidelity function is upper-bounded by a quadratic function for all $\mathbf{x} \in \mathbb{C}^n \setminus Z$. The following theorem states that $F(\mathbf{x})$ is in fact globally upper-bounded by the same quadratic function for all $\mathbf{x} \in \mathbb{C}^n$.

Lemma 5 *Given any $\mathbf{z} \in \mathbb{C}^n$, the fidelity function $F(\mathbf{x})$ is upper-bounded by a quadratic function $Q(\mathbf{x})$:*

$$F(\mathbf{x}) \leq Q(\mathbf{x}) \stackrel{\text{def}}{=} F(\mathbf{z}) + \langle \nabla F(\mathbf{z}), \hat{\mathbf{x}} - \hat{\mathbf{z}} \rangle + \frac{L}{2} \|\hat{\mathbf{x}} - \hat{\mathbf{z}}\|_2^2, \quad (\text{S22})$$

where $L = (1/2)\rho(\mathbf{A}^H \mathbf{A})$.

Proof Let $\Delta \mathbf{x} = \mathbf{x} - \mathbf{z}$, then either of the two following cases occurs:

1) The line between \mathbf{x} and \mathbf{z} does not pass through any nonsmooth points, i.e., $\mathbf{z} + \alpha \Delta \mathbf{x} \in \mathbb{C} \setminus Z, \forall \alpha \in [0, 1]$, or \mathbf{x} and \mathbf{z} lie in the subspace, i.e., $\mathbf{z} + \alpha \Delta \mathbf{x} \in Z, \forall \alpha \in [0, 1]$, the result is obtained directly according to the multivariate Taylor expansion of F :

$$F(\mathbf{x}) = F(\mathbf{z}) + \langle \nabla F(\mathbf{z}), \hat{\mathbf{x}} - \hat{\mathbf{z}} \rangle + \frac{1}{2} (\hat{\mathbf{x}} - \hat{\mathbf{z}})^H \nabla^2 F(\mathbf{u}) (\hat{\mathbf{x}} - \hat{\mathbf{z}})$$

$$\begin{aligned} &\leq F(\mathbf{z}) + \langle \nabla F(\mathbf{z}), \hat{\mathbf{x}} - \hat{\mathbf{z}} \rangle + \frac{L}{2} \|\hat{\mathbf{x}} - \hat{\mathbf{z}}\|_2^2 \\ &= Q(\mathbf{x}), \end{aligned} \tag{S23}$$

where \mathbf{u} is a convex combination of \mathbf{x} and \mathbf{z} .

2) The line between \mathbf{x} and \mathbf{z} passes through a finite number of nonsmooth points. For simplicity, we consider the case of passing through a single nonsmooth point indexed by j , that is, we have

$$|\mathbf{a}_j^H(\mathbf{z} + \alpha^* \Delta \mathbf{x})| = 0, \tag{S24}$$

for some $0 < \alpha^* < 1$. According to 1), for any $0 \leq \alpha \leq \alpha^*$, $F(\mathbf{x})$ is upper-bounded by $Q(\mathbf{x})$. We now prove that this also holds for any $\alpha^* < \alpha < 1$. The fidelity function can be written as a function over α for any point that lies on the line between \mathbf{x} and \mathbf{z} :

$$\begin{aligned} g(\alpha) &= F(\mathbf{z} + \alpha \Delta \mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{z} + \alpha \Delta \mathbf{x}) \\ &= \sum_{i=1, i \neq j}^m f_i(\mathbf{z} + \alpha \Delta \mathbf{x}) + f_j(\mathbf{z} + \alpha \Delta \mathbf{x}) \\ &= \sum_{i=1, i \neq j}^m f_i(\mathbf{z} + \alpha \Delta \mathbf{x}) + \left(|\mathbf{a}_j^H(\mathbf{z} + \alpha \Delta \mathbf{x})| - y_j \right)^2 \\ &= \sum_{i=1, i \neq j}^m f_i(\mathbf{z} + \alpha \Delta \mathbf{x}) + \left(|\alpha - \alpha^*| |\mathbf{a}_j^H \Delta \mathbf{x}| - y_j \right)^2 \\ &\leq \sum_{i=1, i \neq j}^m f_i(\mathbf{z} + \alpha \Delta \mathbf{x}) + \left((\alpha - \alpha^*) |\mathbf{a}_j^H \Delta \mathbf{x}| + y_j \right)^2 \\ &\leq h(\alpha). \end{aligned} \tag{S25}$$

where $f_i(\mathbf{x}) \stackrel{\text{def}}{=} (1/2)(|\mathbf{a}_i^H \mathbf{x}| - y_i)^2$ and $h(\alpha) \stackrel{\text{def}}{=} Q(\mathbf{z} + \alpha \Delta \mathbf{x})$. As a result, we have

$$F(\mathbf{x}) = F(\mathbf{z} + \Delta \mathbf{x}) = g(1) \leq h(1) = Q(\mathbf{x}). \tag{S26}$$

The above derivation can be easily extended to the case of multiple nonsmooth points.

With this, we conclude that for any $\mathbf{x} \in \mathbb{C}^n$, we have

$$F(\mathbf{x}) \leq Q(\mathbf{x}), \tag{S27}$$

which completes the proof. \square

We are now ready to prove the main theorem, which establishes the convergence of the basic proximal gradient method.

Theorem 1 *The basic proximal gradient algorithm (with $\beta_t \equiv 0$ in Algorithm 1) for the problem of Eq. (3) converges to a stationary point using a fixed step size γ that satisfies*

$$\gamma \leq \frac{2}{\rho(\mathbf{A}^H \mathbf{A})}. \tag{S28}$$

Proof The proof is adapted from Ref. [5]. Recall that the proximal update is given by

$$\mathbf{x}^{(t+1)} = \text{prox}_{\gamma R}(\mathbf{x}^{(t)} - \gamma \nabla_{\mathbf{x}} F(\mathbf{x}^{(t)})). \quad (\text{S29})$$

According to Lemma 5, we have that

$$\begin{aligned} F(\mathbf{x}^{(t+1)}) &\leq Q(\mathbf{x}^{(t+1)}) = F(\mathbf{x}^{(t)}) \\ &\quad + \langle \nabla F(\mathbf{x}^{(t)}), \hat{\mathbf{x}}^{(t+1)} - \hat{\mathbf{x}}^{(t)} \rangle + \frac{L}{2} \|\hat{\mathbf{x}}^{(t+1)} - \hat{\mathbf{x}}^{(t)}\|_2^2. \end{aligned} \quad (\text{S30})$$

By the second prox theorem (Theorem 6.39) in Ref. [5], we have

$$\begin{aligned} \langle \hat{\mathbf{x}}^{(t)} - \gamma \nabla F(\mathbf{x}^{(t)}) - \hat{\mathbf{x}}^{(t+1)}, \hat{\mathbf{x}}^{(t)} - \hat{\mathbf{x}}^{(t+1)} \rangle \\ \leq \gamma R(\mathbf{x}^{(t)}) - \gamma R(\mathbf{x}^{(t+1)}), \end{aligned} \quad (\text{S31})$$

from which it follows that

$$\begin{aligned} \langle \nabla F(\mathbf{x}^{(t)}), \hat{\mathbf{x}}^{(t+1)} - \hat{\mathbf{x}}^{(t)} \rangle \\ \leq R(\mathbf{x}^{(t)}) - R(\mathbf{x}^{(t+1)}) - \frac{1}{\gamma} \|\hat{\mathbf{x}}^{(t)} - \hat{\mathbf{x}}^{(t+1)}\|_2^2. \end{aligned} \quad (\text{S32})$$

Let $J(\mathbf{x}) = F(\mathbf{x}) + R(\mathbf{x})$. Combining Eqs. (S30) and (S32), we arrive at

$$\begin{aligned} J(\mathbf{x}^{(t+1)}) &\leq J(\mathbf{x}^{(t)}) + \left(\frac{L}{2} - \frac{1}{\gamma} \right) \|\hat{\mathbf{x}}^{(t)} - \hat{\mathbf{x}}^{(t+1)}\|_2^2 \\ &\leq J(\mathbf{x}^{(t)}) - \frac{L}{2} \|\hat{\mathbf{x}}^{(t+1)} - \hat{\mathbf{x}}^{(t)}\|_2^2. \end{aligned} \quad (\text{S33})$$

The second inequality holds because $\gamma \leq 1/L$. Thus, the updating step for each iteration is upper-bounded:

$$\|\hat{\mathbf{x}}^{(t+1)} - \hat{\mathbf{x}}^{(t)}\|_2^2 \leq \frac{2}{L} \left(J(\mathbf{x}^{(t)}) - J(\mathbf{x}^{(t+1)}) \right). \quad (\text{S34})$$

By summing up T iterations, we arrive at

$$\begin{aligned} \sum_{t=0}^T \|\hat{\mathbf{x}}^{(t+1)} - \hat{\mathbf{x}}^{(t)}\|_2^2 &\leq \frac{2}{L} \sum_{t=0}^T \left(J(\mathbf{x}^{(t)}) - J(\mathbf{x}^{(t+1)}) \right) \\ &\leq \frac{2}{L} \left(J(\mathbf{x}^{(0)}) - J^* \right), \end{aligned} \quad (\text{S35})$$

where $J^* \geq 0$ denotes the global minimum value of the objective function. This implies that

$$\lim_{t \rightarrow \infty} \|\hat{\mathbf{x}}^{(t+1)} - \hat{\mathbf{x}}^{(t)}\|_2 = 0. \quad (\text{S36})$$

That is, the algorithm converges to a stationary point. \square

A similar result has been reported in Ref. [6] regarding the Wirtinger gradient descent algorithm for ptychographic phase retrieval. We consider the more general proximal gradient algorithm and present above an alternative proof.

2.3 Convergence of the accelerated gradient projection algorithm

In order to prove the convergence of the denoising algorithm, we first derive an upper Lipschitz bound for the gradient of $G(\mathbf{w})$.

Lemma 6 *The Lipschitz constant of the gradient of $G(\mathbf{w})$ is upper-bounded by $\gamma^2 \rho(\mathbf{D}^H \mathbf{D})$.*

Proof The proof is adapted from Ref. [7]. Given any $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{C}^{2n}$, we have

$$\begin{aligned}
& \|\nabla_{\mathbf{w}} G(\mathbf{w}_1) - \nabla_{\mathbf{w}} G(\mathbf{w}_2)\|_2 \\
&= \gamma \|\mathbf{D} \mathcal{P}_C(\mathbf{v} - \gamma \mathbf{D}^H \mathbf{w}_1) - \mathbf{D} \mathcal{P}_C(\mathbf{v} - \gamma \mathbf{D}^H \mathbf{w}_2)\|_2 \\
&\leq \gamma \|\mathbf{D}\|_2 \|\mathcal{P}_C(\mathbf{v} - \gamma \mathbf{D}^H \mathbf{w}_1) - \mathcal{P}_C(\mathbf{v} - \gamma \mathbf{D}^H \mathbf{w}_2)\|_2 \\
&\stackrel{(b)}{\leq} \gamma^2 \|\mathbf{D}\|_2 \|\mathbf{D}^H \mathbf{w}_1 - \mathbf{D}^H \mathbf{w}_2\|_2 \\
&\leq \gamma^2 \|\mathbf{D}\|_2^2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \\
&= \gamma^2 \rho(\mathbf{D}^H \mathbf{D}) \|\mathbf{w}_1 - \mathbf{w}_2\|_2,
\end{aligned} \tag{S37}$$

where (b) is based on the fact that projection onto convex sets is non-expansive. \square

Theorem 2 *Assuming that the constraint set C is closed and convex, the accelerated gradient projection algorithm for the problem of Eq. (7) converges to the global optimum using a fixed step size η that satisfies*

$$\eta \leq \frac{1}{\gamma^2 \rho(\mathbf{D}^H \mathbf{D})}. \tag{S38}$$

Proof The gradient projection algorithm can be viewed as a special case of the proximal gradient algorithm with the nonsmooth term being an indicator function. Because Eq. (6) is a convex optimization problem with a closed and convex C , it is sufficient to prove that the objective function is Lipschitz continuous with a constant no greater than $\gamma^2 \rho(\mathbf{D}^H \mathbf{D})$, which is accomplished by Lemma 6. Based on the convergence results of the accelerated proximal gradient algorithm for convex functions [8], the proof is completed. \square

For the particular case of \mathbf{D} being the finite difference operator, the Lipschitz bound (and thus the step size η) can be explicitly given according to the following observation.

Lemma 7 *If \mathbf{D} represents the finite-difference operator defined by Eq. (4), we have*

$$\rho(\mathbf{D}^H \mathbf{D}) \leq 8. \tag{S39}$$

Proof Given any $\mathbf{x} \in \mathbb{C}^n$, we have

$$\|\mathbf{D}\mathbf{x}\|_2^2 = \sum_{i=1}^{n_\xi-1} \sum_{j=1}^{n_\nu} |X_{i+1,j} - X_{i,j}|^2 + \sum_{i=1}^{n_\xi} \sum_{j=1}^{n_\nu-1} |X_{i,j+1} - X_{i,j}|^2$$

$$\begin{aligned}
&\leq \sum_{i=1}^{n_\xi} \sum_{j=1}^{n_\nu} |X_{i+1,j} - X_{i,j}|^2 + \sum_{i=1}^{n_\xi} \sum_{j=1}^{n_\nu} |X_{i,j+1} - X_{i,j}|^2 \\
&\leq 2 \sum_{i=1}^{n_\xi} \sum_{j=1}^{n_\nu} (|X_{i+1,j}|^2 + |X_{i,j}|^2) + 2 \sum_{i=1}^{n_\xi} \sum_{j=1}^{n_\nu} (|X_{i,j+1}|^2 + |X_{i,j}|^2) \\
&\leq 8 \sum_{i=1}^{n_\xi} \sum_{j=1}^{n_\nu} |X_{i,j}|^2 \\
&= 8 \|\mathbf{x}\|_2^2, \tag{S40}
\end{aligned}$$

which implies that

$$\rho(\mathbf{D}^H \mathbf{D}) = \|\mathbf{D}\|_2^2 \leq 8. \tag{S41}$$

□

References

- [1] S. Boyd, S.P. Boyd, L. Vandenberghe, *Convex Optimization* (Cambridge University Press, 2004)
- [2] K. Kreutz-Delgado, The complex gradient operator and the CR-calculus. arXiv preprint arXiv:0906.4835 (2009)
- [3] R.A. Horn, C.R. Johnson, *Matrix Analysis* (Cambridge University Press, 2012)
- [4] Y. Gao, L. Cao, Generalized optimization framework for pixel super-resolution imaging in digital holography. *Optics Express* **29**(18), 28,805–28,823 (2021)
- [5] A. Beck, *First-Order Methods in Optimization* (SIAM, 2017)
- [6] R. Xu, M. Soltanolkotabi, J.P. Haldar, W. Unglaub, J. Zusman, A.F. Levi, R.M. Leahy, Accelerated Wirtinger flow: a fast algorithm for ptychography. arXiv preprint arXiv:1806.05546 (2018)
- [7] A. Beck, M. Teboulle, Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems. *IEEE Transactions on Image Processing* **18**(11), 2419–2434 (2009)
- [8] A. Beck, M. Teboulle, A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences* **2**(1), 183–202 (2009)