

On the Convergence of Learning-based Iterative Methods for Nonconvex Inverse Problems

Risheng Liu, *Member, IEEE*, Shichao Cheng, Yi He, Xin Fan, *Senior Member, IEEE*,
Zhouchen Lin, *Fellow, IEEE*, and Zhongxuan Luo

Abstract—Numerous tasks at the core of statistics, learning and vision areas are specific cases of ill-posed inverse problems. Recently, learning-based (e.g., deep) iterative methods have been empirically shown to be useful for these problems. Nevertheless, integrating learnable structures into iterations is still a laborious process, which can only be guided by intuitions or empirical insights. Moreover, there is a lack of rigorous analysis about the convergence behaviors of these reimplemented iterations, and thus the significance of such methods is a little bit vague. This paper moves beyond these limits and proposes Flexible Iterative Modularization Algorithm (FIMA), a generic and provable paradigm for nonconvex inverse problems. Our theoretical analysis reveals that FIMA allows us to generate globally convergent trajectories for learning-based iterative methods. Meanwhile, the devised scheduling policies on flexible modules should also be beneficial for classical numerical methods in the nonconvex scenario. Extensive experiments on real applications verify the superiority of FIMA.

Index Terms—Nonconvex optimization, Learning-based iteration, Global convergence, Computer vision.

1 INTRODUCTION

IN applications throughout statistics, machine learning and computer vision, one is often faced with the challenge of solving ill-posed inverse problems. In general, the basic inverse problem leads to a discrete linear system of the form $\mathcal{T}(\mathbf{x}) = \mathbf{y} + \mathbf{n}$, where $\mathbf{x} \in \mathbb{R}^D$ is the latent variable to be estimated, \mathcal{T} denotes some given linear operations on \mathbf{x} , and $\mathbf{y}, \mathbf{n} \in \mathbb{R}^D$ are the observation and an unknown error term, respectively. Typically, these inverse problems can be addressed by solving the composite minimization model:

$$\min_{\mathbf{x}} \Psi(\mathbf{x}) := f(\mathbf{x}; \mathcal{T}, \mathbf{y}) + g(\mathbf{x}), \quad (1)$$

where f is the fidelity that captures the loss of data fitting, and g refers to the prior that promotes desired distribution on the solution. Recent studies illustrate that many problems (e.g., image deconvolution, matrix factorization and dictionary learning) naturally require to be solved in the nonconvex scenario. This trend motivates us to investigate

Nonconvex Inverse Problems (NIPs) in the form of Eq. (1) and with the practical configuration that f is continuously differentiable, g is nonsmooth, and both f and g are possibly nonconvex.

Over the past decades, a broad class of first-order methods have been developed to solve special instances of Eq. (1). For example, by integrating Nesterov's acceleration [1] into the fundamental Proximal Gradient (PG) scheme, Accelerated Proximal Gradient (APG, a.k.a. FISTA [2]) method is initially proposed to solve convex models in the form of Eq. (1) for different applications, such as image restoration [2], image deblurring [3], and sparse/low-rank learning [4], etc. While these APGs generate a sequence of objectives that may oscillate [2], [5] developed a variant of APG that guarantees the monotonicity of the sequence. For nonconvex energies in Eq. (1), Li and Lin [6] investigated a monotone APG (mAPG) and proved the convergence under the Kurdyka-Łojasiewicz (KL) constraint [7]. The work in [8] developed another variation of APG (APGnc) for nonconvex problems, but their original analysis only characterized the fixed-point convergence. Recently, Li et al. [9] also proved the subsequence convergence of APGnc and estimated its convergence rates by further exploiting KL property.

Unfortunately, even with some theoretically proved convergence properties, these classical numerical solvers may still fail in real-world scenarios. This is mainly because that the abstractly designed and fixed updating schemes do not exploit the particular structure of the problem at hand nor the input data distribution [10].

In recent years, various learning-based strategies [11], [12], [13], [14], [15] have been proposed to address practical inverse problems in the form of Eq. (1). These methods first introduced hyperparameters into the classical numerical solvers and then performed discriminative learning on collected training data to obtain some data-specific (but possibly inconsistent) iteration schemes. Inspired by the

- R. Liu, Y. He, and X. Fan are with the DUT-RU International School of Information Science & Engineering, Dalian University of Technology, and also with the Key Laboratory for Ubiquitous Network and Service Software of Liaoning Province, Dalian 116024, China. E-mail: {rslu, xin.fan}@dlut.edu.cn, heyiking@outlook.com.
- S. Cheng is with the School of Mathematical Sciences, Dalian University of Technology, and also with the Key Laboratory for Ubiquitous Network and Service Software of Liaoning Province, Dalian 116024, China. E-mail: shichao.cheng@outlook.com.
- Z. Lin is with the Key Laboratory of Machine Perception (Ministry of Education), School of Electronics Engineering and Computer Science, Peking University, Beijing 100871, China, and also with the Cooperative Medianet Innovation Center, Shanghai Jiao Tong University, Shanghai 200240, China. E-mail: zlin@pku.edu.cn.
- Z. Luo is with the DUT-RU International School of Information Science & Engineering, Dalian University of Technology, the Key Laboratory for Ubiquitous Network and Service Software of Liaoning Province, and the School of Mathematical Sciences, Dalian University of Technology, Dalian 116024, China, and also with the Institute of Artificial Intelligence, Guilin University of Electronic Technology, Guilin 541004, China. E-mail: zxliao@dlut.edu.cn

Manuscript received April 19, 2018; revised August 26, 2015.

success of deep learning in different application fields, some preliminary studies considered the handcrafted network architectures as the implicit priors (a.k.a. deep priors) for inverse problems. Following this perspective, various deep priors are designed and nested into numerical iterations [16], [17], [18]. Alternately, the works in [19] and [20] addressed the iteration learning issues from the perspectives of deep reinforcement and recurrent learning, respectively.

Nevertheless, existing hyperparameters learning approaches can only build iterations based on the specific energy forms (e.g., ℓ_1 -penalty and MRFs), so that they are inapplicable for more generic inverse problems. Meanwhile, due to severe inconstancy of parameters during iterations, rigorous analysis on the resulted trajectories is also missing. Deep iterative methods have been executed in many learning and vision problems in practice. However, due to the complex network structure, little or even to no results have been proposed for the convergence behaviors of these methods. In summary, the lack of strict theoretical investigations is one of the most fundamental limits in prevalent learning-based iterative methods, especially in the challenging nonconvex scenario.

To break the limits of prevalent approaches, this paper explores Flexible Iterative Modularization Algorithm (FIMA), a *generic* and *convergent* algorithmic framework that combines together the learnable architecture (e.g., mainstream deep networks) with principled knowledges (formulated by mathematical models), to tackle challenging NIPs in Eq. (1). Specifically, derived from the fundamental forward-backward updating mechanism, FIMA replaces specific calculations corresponding to the fidelity and priors in Eq. (1) with two user-specified (learnable) computational modules. A series of theoretical investigations are established for FIMA. For example, we first prove the subsequence convergence of FIMA with explicit momentum policy (called eFIMA), which is as good as those mathematically designed nonconvex proximal methods with Nesterov's acceleration (e.g., various APGs in [6], [8], [9]). By introducing a carefully devised error-control policy (i.e., implicit momentum policy, called iFIMA), we further enhance the results and obtain a globally convergent Cauchy sequence for Eq. (1). We prove that this guarantee can also be preserved for FIMA with multiple blocks of unknown variables (called mFIMA). As a nontrivial byproduct, we finally show how to specify modules in FIMA for challenging inverse problems in low-level vision area (e.g., non-blind and blind image deconvolution). Our primary contributions are summarized as follows:

- 1) FIMA provides a generic framework that unifies almost all existing learning-based iterative methods, as well as a series of scheduling policies that make it possible to develop theoretically convergent learning-based iterations for challenging nonconvex inverse problems in the form of Eq. (1).
- 2) Even with highly flexible (learnable) iterations, the convergence guarantees obtained by FIMA is still as good as (eFIMA) or better (iFIMA) than prevalent mathematically designed nonconvex APGs. So it is worth noting that our devised scheduling policies together with the flexible algorithmic structures should also be beneficial for classical nonconvex

algorithms.

- 3) FIMA also provides us a practical and effective ensemble of domain knowledge and sophisticated learned data distributions for real applications. Thus we can bring the expressive power of knowledge-based and data-driven methodologies to yield state-of-the-art performance on challenging low-level vision tasks.

2 RELATED WORK

2.1 Classical First-order Numerical Solvers

We first briefly review a group of classical first-order algorithms, which have been widely used to solve inverse problems. The gradient descent (GD) scheme on a differentiable function f can be reformulated as minimizing the following quadratic approximation of f at given point \mathbf{v} with step size $\gamma > 0$, i.e., $Q_{\gamma f}(\mathbf{x}; \mathbf{v}) := f(\mathbf{v}) + \langle \nabla f(\mathbf{v}), \mathbf{x} - \mathbf{v} \rangle + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|^2$. As for the nonsmooth function g , its proximal mapping (PM) with parameter $\gamma > 0$ can be defined as $\text{prox}_{\gamma g}(\mathbf{v}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|^2$. So it is natural to consider PG as cascade of GD (on f) and PM (on g), or equivalently optimizing the quadratic approximation of Eq. (1), i.e., $\mathbf{x}^{k+1} \in \arg \min_{\mathbf{x}} g(\mathbf{x}) + Q_{\gamma^k f}(\mathbf{x}; \mathbf{v}^k)$, where \mathbf{v}^k is some calculated variable at k -th iteration. Thus most prevalent proximal schemes can be summarized as

$$\mathbf{v}^k = \begin{cases} \mathbf{x}^k, & \text{(A-1)} \\ \mathbf{x}^k + \beta^k(\mathbf{x}^k - \mathbf{x}^{k-1}), & \text{(A-2)} \end{cases}$$

$$\mathbf{x}^{k+1} \in \begin{cases} \text{prox}_{\gamma^k g}(\mathbf{v}^k - \gamma^k \nabla f(\mathbf{v}^k)), & \text{(B-1)} \\ \text{prox}_{\gamma^k g}(\mathbf{v}^k - \gamma^k \nabla f(\mathbf{v}^k + \mathbf{e}^k)), & \text{(B-2)} \end{cases}$$

where ε^k and \mathbf{e}^k in (B-2) denote the errors in PM and GD calculations, respectively [21]. Within this general scheme, we first obtain original PG by setting $\mathbf{v}^k = \mathbf{x}^k$ (i.e., (A-1)) and computing PM in (B-1) [2]. Using Nesterov's acceleration [1] (i.e., (A-2) with $\beta^k > 0$), we have the well-known APG method [2], [6], [9]. Moreover, by introducing ε^k and \mathbf{e}^k to respectively capture the inexactness of PM and GD (i.e., (B-2)), we actually consider inexact PG and APG for both convex [22] and nonconvex [21] problems. Notice that in the nonconvex scenario, most classical APGs can only guarantee the subsequence convergence to the critical points [6], [9].

2.2 Learning-based Iterative Methods

In [11], a trained version of FISTA (called LISTA) is introduced to approximate the solution of LASSO. [10], [23] extended LISTA for more generic sparse coding tasks and provided an adaptive acceleration. Unfortunately, LISTA is built on convex ℓ_1 regularization, thus may not be applicable for other complex nonconvex inverse problems (e.g., ℓ_0 prior). By introducing hyperparameters in MRF and solving the resulted variational model with different iteration schemes, various learning-based iterative methods are proposed for inverse problems in image domain (e.g., denoising, super-resolution, and MRI imaging). For example, [13], [14], [15], [24], [25] have considered half-quadratic splitting, gradient descent, Alternating Direction Method of Multiplier (ADMM) and primal-dual method, respectively. But their parameterizations are completely based on MRF priors. Even worse,

the original convergence properties are lost in these resulted iterations.

To better model complex image degradations, [16], [17], [18] considered Convolutional Neural Networks (CNNs) as implicit priors for image restoration. Since these methods discard the regularization term in Eq. (1), we may not enforce principled constraints on their solutions. It is also unclear when and where these iterative trajectories should stop. Another group of very recent works [19], [20] directly formulated the descent directions from reinforcement learning perspective or using recurrent networks. However, due to the high computational budgets, they can only be applied to relative simple tasks (e.g., linear regression). Besides, due to the complex topological network structure, it is extremely hard to provide strict theoretical analysis for these methods.

3 THE PROPOSED ALGORITHMS

This section develops Flexible Iterative Modularization Algorithm (FIMA) for nonconvex inverse problems in Eq. (1). The convergence behaviors are also investigated accordingly. Hereafter, some fairly loose assumptions are enforced on Eq. (1): f is proper and Lipschitz smooth (with modulus L) on a bounded set, g is proper, lower semi-continuous and proximable¹ and Ψ is coercive. Notice that the proofs and definitions are deferred until Supplementary Materials.

3.1 Abstract Iterative Modularization

As summarized in Sec. 2.1, a large amount of first-order methods can be summarized as forward-backward-type iterations. This motivates us to consider the following even more abstract updating principle:

$$\mathbf{x}^{k+1} = \mathcal{A}_g \circ \mathcal{A}_f(\mathbf{x}^k), \quad (2)$$

where \mathcal{A}_f and \mathcal{A}_g respectively stand for the user-specified modules for f and g , and \circ denotes operation composition. Building upon this formulation, it is easy to see that designing a learning-based iterative method reduces to the problem of iteratively specifying and learning \mathcal{A}_f and \mathcal{A}_g .

It is straightforward that most prevalent approaches [13], [14], [15], [16], [17], [18], [24] naturally fall into this general formulation. Nevertheless, currently it is still impossible to provide any strict theoretical results for practical trajectories of Eq. (2). This is mainly due to the lack of efficient mechanisms to control the propagations generated by these handcrafted operations. Fortunately, in the following, we will introduce different scheduling policies to automatically guide the iterations in Eq. (2), resulting in a series of *theoretically convergent learning-based iterative methods*.

3.2 Explicit Momentum: A Straightforward Strategy

The momentum of objective values is one of the most important properties for numerical iterations. This property is also necessary for analyzing the convergence of some classical algorithms. Inspired by these points, we present an explicit momentum FIMA (eFIMA) (i.e., Alg. 1), in which we explicitly compare $\Psi(\mathbf{u}^k)$ and $\Psi(\mathbf{x}^k)$ and choose the variable

with less objective value as our monitor (denoted as \mathbf{v}^k). Finally, a proximal refinement is performed to adjust the learning-based updating at each stage.

Algorithm 1 Explicit Momentum FIMA (eFIMA)

Require: \mathbf{x}^0 , $\mathcal{A} = \{\mathcal{A}_g, \mathcal{A}_f\}$, and $\{0 < \gamma^k < 1/L\}$.

```

1: while not converged do
2:    $\mathbf{u}^k = \mathcal{A}_g \circ \mathcal{A}_f(\mathbf{x}^k)$ .
3:   if  $\Psi(\mathbf{u}^k) \leq \Psi(\mathbf{x}^k)$  then
4:      $\mathbf{v}^k = \mathbf{u}^k$ .
5:   else
6:      $\mathbf{v}^k = \mathbf{x}^k$ .
7:   end if
8:    $\mathbf{x}^{k+1} \in \text{prox}_{\gamma^k g}(\mathbf{v}^k - \gamma^k \nabla f(\mathbf{v}^k))$ .
9: end while
```

The following theorem first verifies the sufficient descent of $\{\Psi(\mathbf{x}^k)\}_{k \in \mathbb{N}}$ and then proves the subsequence convergence of eFIMA. It is nice to observe that these results are not based on any specific choices of \mathcal{A}_f and \mathcal{A}_g .

Theorem 1. *Let $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ be the sequence generated by eFIMA. Then at the k -th iteration, there exists a sequence $\{\alpha^k | \alpha^k > 0\}_{k \in \mathbb{N}}$, such that*

$$\Psi(\mathbf{x}^{k+1}) \leq \Psi(\mathbf{v}^k) - \alpha^k \|\mathbf{x}^{k+1} - \mathbf{v}^k\|^2, \quad (3)$$

where \mathbf{v}^k is the monitor in Alg. 1. Furthermore, $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ is bounded and any of its accumulation points are the critical points of $\Psi(\mathbf{x})$ in Eq. (1).

Based on Theorem 1 and considering Ψ as a semi-algebraic function², the convergence rate of eFIMA can be straightforwardly estimated as follows.

Corollary 1. *Let $\phi(s) = \frac{t}{\theta} s^\theta$ be a desingularizing function with a constant $t > 0$ and a parameter $\theta \in (0, 1]$ [27]. Then $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ generated by eFIMA converges after finite iterations if $\theta = 1$. The linear and sub-linear rates can be obtained if choosing $\theta \in [1/2, 1)$ and $\theta \in (0, 1/2)$, respectively.*

3.3 Implicit Momentum via Error Control

Indeed, even with the explicit momentum schedule, we may still not obtain a globally convergent iteration. This is mainly because that there is no policy to efficiently control the inexactness of the user-specified modules (i.e., \mathcal{A}). In this subsection, we show how to address this issue by controlling the first-order optimality error during iterations.

Specifically, we consider the auxiliary of Ψ at \mathbf{x}^k (denoted as Ψ^k) and denote its sub-differential (denoted as $\mathbf{d}_{\Psi^k}^{\mathbf{x}}$)³ as

$$\begin{aligned} \Psi^k(\mathbf{x}) &= f(\mathbf{x}) + g(\mathbf{x}) + \frac{\mu^k}{2} \|\mathbf{x} - \mathbf{x}^k\|^2, \\ \mathbf{d}_{\Psi^k}^{\mathbf{x}} &= \mathbf{d}_g^{\mathbf{x}} + \nabla f(\mathbf{x}) + \mu^k(\mathbf{x} - \mathbf{x}^k) \in \partial \Psi^k(\mathbf{x}), \end{aligned} \quad (4)$$

where $\mu^k > 0$ is the penalty parameter and $\mathbf{d}_g^{\mathbf{x}} \in \partial g(\mathbf{x})$.

As shown in Alg. 2, at stage k , a variable $\tilde{\mathbf{u}}^k$ is obtained by proximally minimizing Ψ^k at \mathbf{u}^k (i.e., Step 3 of Alg. 2).

2. Indeed, a variety of functions (e.g., the indicator function of polyhedral set, ℓ_0 and rational ℓ_p penalties) satisfy the semi-algebraic property [26].

3. Strictly speaking, $\partial \Psi^k(\mathbf{x})$ is the so-called limiting Frechét sub-differential. We state its formal definition and propose a practical computation scheme for $\mathbf{d}_{\Psi^k}^{\mathbf{x}}$ in Appendix.

1. The function g is proximable if $\min_{\mathbf{x}} g(\mathbf{x}) + \frac{\gamma}{2} \|\mathbf{x} - \mathbf{y}\|^2$ can be easily solved by the given \mathbf{y} and $\gamma > 0$.

Algorithm 2 Implicit Momentum FIMA (iFIMA)

Require: $\mathbf{x}^0, \mathcal{A} = \{\mathcal{A}_g, \mathcal{A}_f\}, \{0 < 2C^k < \mu^k < \infty\},$ and $\{0 < \gamma^k < 1/L\}.$

- 1: **while** not converged **do**
- 2: $\mathbf{u}^k = \mathcal{A}_g \circ \mathcal{A}_f(\mathbf{x}^k).$
- 3: $\tilde{\mathbf{u}}^k \in \text{prox}_{\gamma^k g}(\mathbf{u}^k - \gamma^k(\nabla f(\mathbf{u}^k) + \mu^k(\mathbf{u}^k - \mathbf{x}^k))).$
- 4: **if** $\|\mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}\| \leq C^k \|\tilde{\mathbf{u}}^k - \mathbf{x}^k\|$ **then**
- 5: $\mathbf{v}^k = \tilde{\mathbf{u}}^k.$
- 6: **else**
- 7: $\mathbf{v}^k = \mathbf{x}^k.$
- 8: **end if**
- 9: $\mathbf{x}^{k+1} \in \text{prox}_{\gamma^k g}(\mathbf{v}^k - \gamma^k \nabla f(\mathbf{v}^k)).$
- 10: **end while**

Roughly, this new variable is just an ensemble of the last updated \mathbf{x}^k and the output \mathbf{u}^k of user-specified \mathcal{A} following the specific proximal structure in Eq. (1). Then the monitor is obtained by checking the boundedness of $\mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}$. Notice that the constant C^k actually reveals our tolerance to the inexactness of \mathcal{A} at k -th iteration.

Proposition 1. *Let $\{\mathbf{x}^k, \tilde{\mathbf{u}}^k, \mathbf{v}^k\}_{k \in \mathbb{N}}$ be the sequences generated by Alg. 2. Then there exist two sequences $\{\alpha^k | \alpha^k > 0\}_{k \in \mathbb{N}}$ and $\{\beta^k | \beta^k > 0\}_{k \in \mathbb{N}}$, such that the inequality (3) in Theorem 1 and $\Psi(\tilde{\mathbf{u}}^k) \leq \Psi(\mathbf{x}^k) - \beta^k \|\tilde{\mathbf{u}}^k - \mathbf{x}^k\|^2$ are respectively satisfied.*

Equipped with Proposition 1, it will be straightforward to guarantee that the objective values generated by Alg. 2 (i.e., $\{\Psi(\mathbf{x}^k)\}_{k \in \mathbb{N}}$) also has sufficient descent. So we call this version of FIMA as implicit momentum FIMA (iFIMA). Then the global convergence of iFIMA is proved as follows.

Theorem 2. *Let $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ be the sequence generated by iFIMA. Then $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ is bounded and any of its accumulation points are the critical points of Ψ . If Ψ is semi-algebraic, we further have that $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ is a Cauchy sequence, thus globally converges to a critical point of $\Psi(\mathbf{x})$ in Eq. (1).*

Indeed, based on Theorem 2, it is also easy to obtain the same convergence rate as that in Corollary 1 for iFIMA.

3.3.1 Practical Calculation of $\mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}$ in iFIMA

Here we propose a practical calculation scheme for $\mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k} \in \partial\Psi^k(\tilde{\mathbf{u}}^k)$ defined in Eq. (4) and used in Alg. 2. In fact, it is challenging to directly calculate $\mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}$ since the sub-differential $\mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}$ is often intractable in the non-convex scenario [28]. Fortunately, our following analysis provides an efficient practical calculation scheme for $\mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}$ within FIMA framework. Specifically, by considering $\tilde{\mathbf{u}}^k$ as the solution for Step 3 of Alg. 2, we have

$$\begin{aligned} \tilde{\mathbf{u}}^k &\in \text{prox}_{\gamma^k g}(\mathbf{u}^k - \gamma^k(\nabla f(\mathbf{u}^k) + \mu^k(\mathbf{u}^k - \mathbf{x}^k))) \\ &= \arg \min_{\mathbf{u}} \frac{1}{2} \|\mathbf{u} - (\mathbf{u}^k - \gamma^k(\nabla f(\mathbf{u}^k) + \mu^k(\mathbf{u}^k - \mathbf{x}^k)))\|^2 \\ &\quad + g(\mathbf{u}). \end{aligned}$$

By the first-order optimality condition, we further have

$$\begin{aligned} 0 &\in \partial g(\tilde{\mathbf{u}}^k) + \frac{1}{\gamma^k}(\tilde{\mathbf{u}}^k - (\mathbf{u}^k - \gamma^k(\nabla f(\mathbf{u}^k) + \mu^k(\mathbf{u}^k - \mathbf{x}^k)))) \\ &\Rightarrow -(\nabla f(\mathbf{u}^k) + \mu^k(\mathbf{u}^k - \mathbf{x}^k)) + \frac{1}{\gamma^k}(\mathbf{u}^k - \tilde{\mathbf{u}}^k) \in \partial g(\tilde{\mathbf{u}}^k). \end{aligned}$$

Recalling the definition of $\mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}$ in Eq. (4) and $\mathbf{d}_g^{\tilde{\mathbf{u}}^k} \in \partial g(\tilde{\mathbf{u}}^k)$, we finally obtain the practical calculation scheme for $\mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}$:

$$\begin{aligned} \mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k} &= \mathbf{d}_g^{\tilde{\mathbf{u}}^k} + \nabla f(\tilde{\mathbf{u}}^k) + \mu^k(\tilde{\mathbf{u}}^k - \mathbf{x}^k) \\ &= (\mu^k - 1/\gamma^k)(\tilde{\mathbf{u}}^k - \mathbf{u}^k) - (\nabla f(\mathbf{u}^k) - \nabla f(\tilde{\mathbf{u}}^k)). \end{aligned}$$

3.4 Multi-block Extension

Algorithm 3 Multi-block FIMA (mFIMA)

Require: $\mathfrak{X}^0, \mathcal{A} = \{\mathcal{A}_{g_1}, \dots, \mathcal{A}_{g_N}, \mathcal{A}_f\}, \{0 < 2C_n^k < \mu_n^k < \infty\},$ and $\{0 < \gamma_n^k < 1/L_n\}.$

- 1: **while** not converged **do**
- 2: **for** $n = 1 : N$ **do**
- 3: $\mathbf{u}_n^k = \mathcal{A}_{g_n} \circ \mathcal{A}_f(\mathfrak{X}_{[<n]}^{k+1}, \mathfrak{X}_{[\geq n]}^k).$
- 4: $\tilde{\mathbf{u}}_n^k \in \text{prox}_{\gamma_n^k g_n}(\mathbf{u}_n^k - \gamma_n^k(\nabla_n f(\mathfrak{X}_{[<n]}^{k+1}, \mathbf{u}_n^k, \mathfrak{X}_{[\geq n]}^k) + \mu_n^k(\mathbf{u}_n^k - \mathbf{x}_n^k))).$
- 5: **if** $\|\mathbf{d}_{\Psi_n^k}^{\tilde{\mathbf{u}}_n^k}\| \leq C_n^k \|\tilde{\mathbf{u}}_n^k - \mathbf{x}_n^k\|$ **then**
- 6: $\mathbf{v}_n^k = \tilde{\mathbf{u}}_n^k.$
- 7: **else**
- 8: $\mathbf{v}_n^k = \mathbf{x}_n^k.$
- 9: **end if**
- 10: $\mathbf{x}_n^{k+1} \in \text{prox}_{\gamma_n^k g_n}(\mathbf{v}_n^k - \gamma_n^k \nabla_n f(\mathfrak{X}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathfrak{X}_{[\geq n]}^k)).$
- 11: **end for**
- 12: **end while**

In order to tackle the inverse problems with blocks of unknown variables (e.g., blind deconvolution and dictionary learning), we now discuss how to extend FIMA for multi-block NIPs, which is formulated as $\mathcal{T}(\mathfrak{X}) = \mathbf{y} + \mathbf{n}$, where $\mathfrak{X} = \{\mathbf{x}_n\}_{n=1}^N \in \mathbb{R}^{D_1} \times \dots \times \mathbb{R}^{D_N}$ is a set of $N \geq 2$ unknown variables to be estimated. Notice that here \mathcal{T} should be some given linear operations on \mathfrak{X} . The inference of such problem can be addressed by solving

$$\min_{\mathfrak{X}} \Psi(\mathfrak{X}) := f(\mathfrak{X}; \mathcal{T}, \mathbf{y}) + \sum_{n=1}^N g_n(\mathbf{x}_n), \quad (5)$$

where $f(\mathfrak{X}) : \mathbb{R}^{D_1} \times \dots \times \mathbb{R}^{D_N} \rightarrow (-\infty, +\infty]$ is still differentiable and each $g_n(\mathbf{x}_n) : \mathbb{R}^{D_n} \rightarrow (-\infty, +\infty]$ may also nonsmooth and possibly nonconvex. Here both f and block-wise $g_n(\mathbf{x}_n)$ follow the same assumptions as that in Eq. (1) and f should also satisfy the generalized Lipschitz smooth property on bounded subsets of $\mathbb{R}^{D_1} \times \dots \times \mathbb{R}^{D_N}$. For ease of presentation, we denote $\mathfrak{X}_{[<n]} = \{\mathbf{x}_i\}_{i=1}^{n-1}$, $\mathfrak{X}_{[\leq n]} = \{\mathbf{x}_i\}_{i=1}^n$ and the subscripts $[> n]$ and $[\geq n]$ are defined in the same manner. Then we summarize the main iterations of multivariable FIMA (mFIMA) as follows⁴:

$$\begin{aligned} \mathbf{u}_n^k &= \mathcal{A}_{g_n} \circ \mathcal{A}_f(\mathfrak{X}_{[<n]}^{k+1}, \mathfrak{X}_{[\geq n]}^k), \\ \mathbf{x}_n^{k+1} &\in \text{prox}_{\gamma_n^k g_n}(\mathbf{v}_n^k - \gamma_n^k \nabla_n f(\mathfrak{X}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathfrak{X}_{[\geq n]}^k)). \end{aligned}$$

Here \mathbf{v}_n^k is the monitor of \mathbf{x}_n^k , obtained by the same error control strategy as that in iFIMA. Then we summarize our multi-block FIMA in Alg. 3 and prove the convergence of mFIMA in Corollary 2.

4. Due to space limit, the details of mFIMA are presented in Supplemental Material.

Corollary 2. Let $\{\mathcal{X}^k\}_{k \in \mathbb{N}}$ be the sequence generated by mFIMA. Then we have the same convergence properties as that in Theorem 2 and Corollary 1 for $\{\mathcal{X}^k\}_{k \in \mathbb{N}}$.

Then we summarize our multi-block FIMA in Alg. 3. Notice that here we adopt the error-control policy in iFIMA to guide the iterations of mFIMA.

4 DISCUSSIONS

Here we would like to discuss some important aspects of our FIMA, including the difference with traditional first order methods, learning-based iterative methods, and existing meta-learning techniques. We also discuss the relation between our eFIMA and iFIMA. For convenience, we specify the components in FIMA as three categories, including the flexible modules $\mathcal{A}_f \circ \mathcal{A}_g$, the various criteria (Steps 3-7 in Alg. 1 and Steps 3-8 in Alg. 2) and PG operator.

4.1 FIMA vs. Traditional First Order Methods

We first point out that by specifying the user-specified modules \mathcal{A}_f and \mathcal{A}_g as numerical calculations, our FIMA can reduce to some traditional first-order numerical algorithms (e.g., PG/APG/PALM) when ignoring the criteria and PG operator. Specifically, by respectively specifying \mathcal{A}_f and \mathcal{A}_g as gradient descent and proximal operator, our e/iFIMA will be the standard PG [29]. Our e/iFIMA also can reduce to the standard APG [2] when adopting Nesterov's acceleration and proximal gradient operator as \mathcal{A}_f and \mathcal{A}_g , respectively. Considering the two-block case in mFIMA, our method becomes PALM [26] when configuring the same strategy with PG. In these degradations, we can obtain the same convergence results with existing PG/APG.

When remaining the criteria and PG operator in FIMA, e/iFIMA will be the inexact APGs by setting $\mathcal{A}_g \circ \mathcal{A}_f$ similar with PG. For example, eFIMA will be the monotone APG [6] when arranging \mathcal{A}_f and \mathcal{A}_g as Nesterov's acceleration and proximal gradient operator, respectively. Furthermore, when transforming these arrangements into our iFIMA framework, it generates a new inexact APG method. However, the convergence performance based on iFIMA is even better than that for the prevalent nonconvex APGs [6], [8]. This actually suggests that our devised error-control policy together with the flexible algorithmic structures should also be beneficial for inexact nonconvex algorithms.

Moreover, our experimental results have demonstrated that thanks to the plug-and-play architectures, FIMA can achieve much better numerical performance and final results than these traditional first order numerical methods, especially in real-world applications. Indeed, FIMA provides a generic, flexible and theoretically guaranteed way to extend standard numerical methods using plug-and-play architectures.

4.2 FIMA vs. Learning-based Iterative Methods

As discussed above, most existing learning-based iterative methods only replace their numerical computations by trained architectures, which directly lead to the missing of necessary theoretical guarantees. Fortunately, within FIMA framework, we can prove in Corollary 1 and Theorem 2 that

our newly proposed learning-based scheme does not depend on the particular choices of \mathcal{A}_f and \mathcal{A}_g in general. It actually provides us a unified methodology to analyze and improve the convergence issues for learning-based methods. Within FIMA, we can provide an easily-implemented and strictly convergent way to extend almost all the learning-based methods reviewed in Sec. 2.2. For example, we can design the gradient descent operator and the encoder architecture in LISTA [11] as \mathcal{A}_f and \mathcal{A}_g , respectively. Then we cascade the explicit or implicit momentum in our algorithms to build the convergence guaranteed iterations. Similarly, we can also regard the learnable priors in MRF [24] as \mathcal{A}_g and the rest part as \mathcal{A}_f to generate our eFIMA or iFIMA. When designing the solution of the subproblem about fidelity as \mathcal{A}_f and exploring the data distribution by denoise CNN [17] as \mathcal{A}_g , we also provide the convergence guarantee for these plug-and-play learnable iterations. The criterion in our algorithms actually provides the guidance to judge whether the output of each iteration in learning based methods is reasonable. Thus, almost all the learning-based methods can strictly converge with tiny assistance under our algorithm framework.

4.3 FIMA vs. Existing Meta-learning Techniques

Meta-learning (a.k.a. "learning to learn") aims to design methods that can learn how to learn new tasks by reusing previous experience, rather than considering each new task in isolation [30], [31], [32]. It should be noticed that FIMA can also be categorized as a specific meta-learning technique from the perspective of "learning to optimize". However, compared with existing approaches [20], [33], which just learn all the hyperparameters in their optimization processes in heuristic manners and thus miss solid theoretical investigations, the main advantages of our FIMA is that we provide a *theoretically guaranteed* framework to learn *strict convergence* optimization schemes for meta-learning. But please notice that to obtain these convergence results, we have to set some algorithmic parameters following the theoretical guidance.

4.4 eFIMA vs. iFIMA

We first clarify that the main difference between eFIMA and iFIMA is the error-control condition. It can be seen that the optimally-based condition in iFIMA is a little bit stricter than the loss-based eFIMA condition. Thus we can obtain better convergence properties but additional computations are needed at each iterations.

As for the plug-and-play computational modules in e/iFIMA, it will be shown in Sec. 6 that the choices of \mathcal{A}_f and \mathcal{A}_g do affect our speed and accuracy in practice. Fortunately, the proposed scheduling of learnable and numerical modules are automatically and adaptively adjusted by error-control conditions of both eFIMA and iFIMA. So in the specific task, FIMA actually can automatically reject improper \mathcal{A}_f and \mathcal{A}_g during iterations.

5 APPLICATIONS

As a nontrivial byproduct, this section illustrates how to apply FIMA to tackle practical inverse problems in low-level

vision area, such as image deconvolution in the standard non-blind and even more challenging blind scenarios.

Non-blind Deconvolution (Uni-block) aims to restore the latent image \mathbf{z} from corrupted observation \mathbf{y} with known blur kernel \mathbf{b} . In this part, we utilize the well-known sparse coding formulation [2]: $\mathbf{y} = \mathbf{D}\mathbf{x} + \mathbf{n}$, where \mathbf{x} , \mathbf{D} and \mathbf{n} are the sparse code, given dictionary and unknown noises, respectively. Indeed, the form of \mathbf{D} is given as $\mathbf{D} = \mathbf{B}\mathbf{W}^\top$, where \mathbf{B} is the matrix form of \mathbf{b} , \mathbf{W}^\top denotes the inverse of the wavelet transform \mathbf{W} (i.e., $\mathbf{x} = \mathbf{W}\mathbf{z}$ and $\mathbf{z} = \mathbf{W}^\top\mathbf{x}$). So by defining $f(\mathbf{x}; \mathbf{D}, \mathbf{y}) = \|\mathbf{y} - \mathbf{D}\mathbf{x}\|^2$ and $g(\mathbf{x}) = \lambda\|\mathbf{x}\|_p$ ($0 \leq p < 1$), we obtain a special case of Eq. (1) as follows

$$\min_{\mathbf{x}} f(\mathbf{x}; \mathbf{D}, \mathbf{y}) + g(\mathbf{x}). \quad (6)$$

Now we are ready to design iterative modules (i.e., \mathcal{A}_f and \mathcal{A}_g) to optimize the SC model in Eq. (6). With the well-known imaging formulation $\mathbf{y} = \mathbf{b} \otimes \mathbf{z} + \mathbf{n}$ (\otimes denotes the convolution operator), we actually update \mathbf{z} by solving $\mathcal{A}_f(\mathbf{z}^k) := \arg \min_{\mathbf{z}} \|\mathbf{y} - \mathbf{b} \otimes \mathbf{z}\|^2 + \tau\|\mathbf{z} - \mathbf{z}^k\|^2$ to aggregate principles of the task and information from last updated variable, where $\mathbf{z}^k = \mathbf{W}^\top\mathbf{x}^k$ and τ is a positive constant. Then \mathcal{A}_f on \mathbf{x} can be defined as $\mathcal{A}_f(\mathbf{x}^k) = \mathbf{W}\mathcal{A}_f(\mathbf{z}^k)$, i.e.,

$$\mathcal{A}_f(\mathbf{x}^k) = \mathbf{W}(\mathbf{B}^\top\mathbf{B} + \tau\mathbf{I})^{-1} (\mathbf{B}^\top\mathbf{y} + \tau\mathbf{W}^\top\mathbf{x}^k), \quad (7)$$

where \mathbf{I} is the identity matrix. It is easy to check that \mathcal{A}_f can be efficiently calculated by FFT [24]. As for \mathcal{A}_g , we consider solving $\mathcal{A}_g(\mathcal{A}_f(\mathbf{z}^k)) := \arg \min_{\mathbf{z}} g(\mathbf{z}) + \tau\|\mathbf{z} - \mathcal{A}_f(\mathbf{z}^k)\|^2$ by a network to describe the distribution of latent image.

Blind Deconvolution (Multi-block) involves the joint estimation of both the latent image \mathbf{z} and blur kernel \mathbf{b} , given only an observed \mathbf{y} . Here we formulate this problem on image gradient domain and solve the following special case of Eq. (5) with two unknown variables (\mathbf{x}, \mathbf{b}) ⁵:

$$\min_{\mathbf{x}, \mathbf{b}} f(\mathbf{x}, \mathbf{b}; \nabla\mathbf{y}) + g_{\mathbf{x}}(\mathbf{x}) + g_{\mathbf{b}}(\mathbf{b}), \quad (8)$$

where $f(\mathbf{x}, \mathbf{b}; \nabla\mathbf{y}) = \|\nabla\mathbf{y} - \mathbf{b} \otimes \mathbf{x}\|^2$, $g_{\mathbf{x}}(\mathbf{x}) = \lambda_{\mathbf{x}}\|\mathbf{x}\|_0$, and $g_{\mathbf{b}}(\mathbf{b}) = \chi_{\Omega_{\mathbf{b}}}(\mathbf{b})$. Here $\chi_{\Omega_{\mathbf{b}}}$ is the indicator function of the set $\Omega_{\mathbf{b}} := \{\mathbf{b} \in \mathbb{R}^{D_{\mathbf{b}}} : [\mathbf{b}]_i \geq 0, \sum_{i=1}^{D_{\mathbf{b}}} [\mathbf{b}]_i = 1\}$, where $[\cdot]_i$ denotes the i -th element. So the proximal updating in mFIMA corresponding to $g_{\mathbf{x}}$ and $g_{\mathbf{b}}$ can be respectively calculated by hard-thresholding [3] and simplex projection [34]. Here we need to specify three modules (i.e., \mathcal{A}_f , $\mathcal{A}_{g_{\mathbf{x}}}$ and $\mathcal{A}_{g_{\mathbf{b}}}$) for mFPGA. We first follow similar idea in the non-blind case to define $\mathcal{A}_f(\mathbf{x}^k, \mathbf{b}^k)$ using the aggregated deconvolution energy

$$\mathcal{A}_f(\mathbf{x}^k, \mathbf{b}^k) := \arg \min_{\mathbf{x}, \mathbf{b}} \|\nabla\mathbf{y} - \mathbf{b} \otimes \mathbf{x}\|^2 + \tau_{\mathbf{x}}\|\mathbf{x} - \mathbf{x}^k\|^2 + \tau_{\mathbf{b}}\|\mathbf{b} - \mathbf{b}^k\|^2, \quad (9)$$

where $\tau_{\mathbf{b}}$ and $\tau_{\mathbf{x}}$ are positive constants. We then train CNNs on image gradient domain and solve $\min_{\mathbf{b}} \|\nabla\mathbf{y} - \mathbf{b} \otimes \mathbf{x}\|^2 + \lambda_{\mathbf{b}}\|\mathbf{b}\|^2$ using conjugate gradient method [35] to formulate $\mathcal{A}_{g_{\mathbf{x}}}$ and $\mathcal{A}_{g_{\mathbf{b}}}$, respectively.

Rain Streaks Removal (Multi-block) is another challenging task which focuses on removing the sporadic rain streak \mathbf{r} and restoring rain free background scenes \mathbf{z} from several types of visibility distorted observations \mathbf{o} . The observation

5. Notice that in this section, \mathbf{x} is defined with different meanings, i.e., image gradient in Eq. (8), while sparse code in Eq. (6).

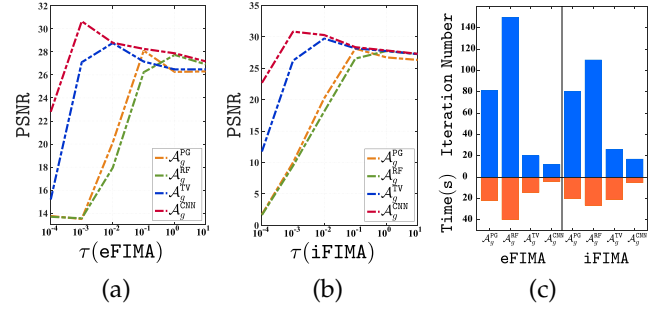


Fig. 1. Comparisons of FIMA with different \mathcal{A}_f ($\tau \in [10^{-4}, 10^{-1}]$) and $\mathcal{A}_g \in \{\mathcal{A}_g^{\text{PG}}, \mathcal{A}_g^{\text{RF}}, \mathcal{A}_g^{\text{TV}}, \mathcal{A}_g^{\text{CNN}}\}$. The bar charts in the rightmost subfigure compares the overall iteration number and running time (in seconds, “Time(s)” for short).

\mathbf{o} can be generated by $\mathbf{o} = \mathbf{z} + \mathbf{r}$. Considering the different sparsity of rain streak and background image, we formulate the sparse coding model as to minimize the following energy function:

$$\min_{\mathbf{x}, \mathbf{c}} f(\mathbf{x}, \mathbf{c}; \mathbf{o}) + g_{\mathbf{x}}(\mathbf{x}) + g_{\mathbf{c}}(\mathbf{c}), \quad (10)$$

where $f(\mathbf{x}, \mathbf{c}; \mathbf{o}) = \|\mathbf{o} - \mathbf{W}^\top\mathbf{x} - \mathbf{W}^\top\mathbf{c}\|$, $g_{\mathbf{x}}(\mathbf{x}) = \lambda_{\mathbf{x}}\|\mathbf{x}\|_{0.8}$ and $g_{\mathbf{c}}(\mathbf{c}) = \lambda_{\mathbf{c}}\|\mathbf{c}\|_0$. Here, \mathbf{W} is the wavelet transform which has explained in non-blind deconvolution. \mathbf{x} , \mathbf{c} are the sparse codes of \mathbf{z} and \mathbf{r} , respectively (i.e., $\mathbf{z} = \mathbf{W}^\top\mathbf{x}$, $\mathbf{r} = \mathbf{W}^\top\mathbf{c}$). By applying our multi-block FIMA to solve Eq. (10), we design \mathcal{A}_f by the similar strategy in Eq. (9). As for $\mathcal{A}_{g_{\mathbf{x}}}$ and $\mathcal{A}_{g_{\mathbf{c}}}$, we achieve them by same network architecture but different training data.

6 EXPERIMENTAL RESULTS

This section conducts experiments to verify our theoretical results and compares the performance of FIMA with other state-of-the-art learning-based iterative methods on real-world inverse problems. All experiments are performed on a PC with Intel Core i7 CPU at 3.4 GHz, 32 GB RAM and a NVIDIA GeForce GTX 1050 Ti GPU. More results can also be found in Supplemental Materials.

6.1 Non-blind Image Deconvolution

We first evaluate FIMA on solving Eq. (6) for image restoration. The test images are collected by [24], [36] and different levels of Gaussian noise are further added to generate our corrupted observations.

Modules Evaluation: Firstly, the influences of different choices of \mathcal{A} in FIMA is studied. Following Eq. (7), we adopt \mathcal{A}_f^τ with varying τ . As for \mathcal{A}_g , different choices are also considered: classical PG ($\mathcal{A}_g^{\text{PG}}$), Recursive Filter [37] ($\mathcal{A}_g^{\text{RF}}$), Total Variation [38] ($\mathcal{A}_g^{\text{TV}}$) and CNNs ($\mathcal{A}_g^{\text{CNN}}$). For $\mathcal{A}_g^{\text{CNN}}$, we introduce a residual structure $\mathbf{x} = \mathbf{x} + \mathcal{R}(\mathbf{x})$ [39] and define \mathcal{R} as a cascade of 7 dilated convolution layers (with filter size 3×3). ReLUs are added between each two linear layers and batch normalizations are used for the 2-nd to 6-th linear layers. We collect 800 images, in which 400 have been used in [24] and the other 400 are randomly sampled from ImageNet [40] as training data. We set the patch size as 160×160 , and feed 10 patches in each epoch. Finally, Adam is adopted and executed 60000 epochs for learning rate ranging from $1e-4$ to $1e-6$. Here we just adopt similar strategies in [17] to

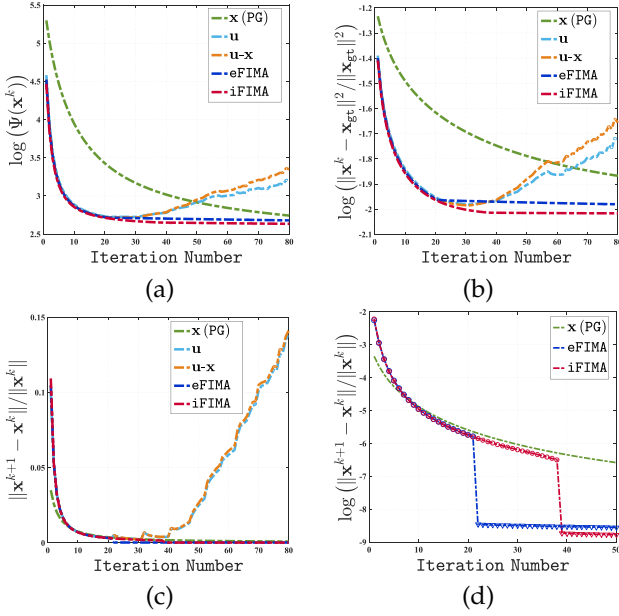


Fig. 2. The iteration curves of FIMA with different settings. The first three subfigures express the function values, constructive errors, and iteration errors, respectively. Subfigure (d) only plots the first 50 iterations for illustrate the scheduling policies of FIMA.

train $\mathcal{A}_g^{\text{CNN}}$ with different noise levels. Thus, we only training $\mathcal{A}_g^{\text{CNN}}$ independently in our iterations. As for the algorithmic parameters, we set $\gamma^k = 0.1$, $\mu^k = 0.2$, and $C^k = 0.09$. Fig. 1 analyzes the contributions of \mathcal{A}_f^r ($r \in [10^{-4}, 10^1]$) and $\mathcal{A}_g \in \{\mathcal{A}_g^{\text{PG}}, \mathcal{A}_g^{\text{RF}}, \mathcal{A}_g^{\text{TV}}, \mathcal{A}_g^{\text{CNN}}\}$. We observe that $\mathcal{A}_g^{\text{TV}}$ is relatively better than $\mathcal{A}_g^{\text{PG}}$ and $\mathcal{A}_g^{\text{RF}}$, while $\mathcal{A}_g^{\text{CNN}}$ performs consistently better and faster than other strategies. So hereafter we always utilize $\mathcal{A}_g^{\text{CNN}}$ in eFIMA and iFIMA. We also observe that even with different \mathcal{A}_g , relatively large τ in \mathcal{A}_f^r will result in analogous quantitative results. Thus we experimentally set $\tau = 10^{-3}$ for \mathcal{A}_f^r in eFIMA and iFIMA for all the experiments.

Convergence Behaviors: We then verify the convergence properties of FIMA. The convergence behaviors of both each module in our algorithms and other nonconvex APGs are considered. To be fair and comprehensive, we adopt specific iteration numbers ($K = 80$) and iteration errors ($\|x^{k+1} - x^k\| / \|x^k\| \leq 10^{-4}$) as stopping criterion in Figs. 2 and 3, respectively.

In Fig. 2(a), (b), and (c), we plot the curves of objective values ($\log(\Psi(x^k))$), reconstruction errors ($\log(\|x^{k+1} - x^k\|^2 / \|x^k\|^2)$) and iteration errors for FIMA with different settings. The legends “x”, “u”, and “u-x” respectively denote that at each iteration, we only perform classical PG (i.e., only the last step in Algs. 1 and 2), task-driven modules \mathcal{A} (i.e., only perform Eq. (2)), and their naive combination (without any scheduling policies). It can be seen that the function values and reconstruction errors of PG decrease slower than our FIMA strategies, while both “u”-curve (i.e., naive $\mathcal{A}_g \circ \mathcal{A}_f$) and “u-x”-curve (i.e., \mathcal{A} with PG refinement but no “explicit momentum” or “error-control” policy) have oscillations and could not converge after only 30 iterations. Moreover, we observe that adding PG to “u” (i.e., “u-x”) make the curve worse rather than correct it to the descent direction. It illustrates that the pure adding strategies

TABLE 1

The number of iterations (including plug-and-play modules) in FIMA. “No. Iter.” reports the number of whole iterations and “No. \mathcal{A} ” denotes the times the plug-and-play modules $\mathcal{A}_g \circ \mathcal{A}_f$ has been performed by FIMA during iterations. We also report the number of iterations for standard PG in the rightmost column.

Image	eFIMA		iFIMA		PG
	No. Iter.	No. \mathcal{A}	No. Iter.	No. \mathcal{A}	No. Iter.
Fig. 4	13	11	22	21	542
Fig. 5	19	15	26	25	577

indeed break the convergence guarantee. In contrast, since of the choice mechanism in our algorithms, both eFIMA and iFIMA can provide a reliable variable (v^k) in the current iteration to satisfy the convergence condition. We further explore the choice mechanism of FIMA in Fig. 2(d). The “circles” in each curve represent the “explicit momentum” or “error-control” policy is satisfied, while the “triangles” denote only perform PG in the current stage. It can be seen that the eFIMA strategy is more strict than iFIMA, the judgment policy fails only 20 iterations in eFIMA while remains almost 40 iterations in iFIMA. Both eFIMA and iFIMA have better performance than other compared schemes, thus verifies the efficiency of our proposed scheduling policies in Sec. 3.

We also compare the iteration behaviors of FIMA to classical nonconvex APGs, including mAPG [6], APGnc [9] and inexact niAPG [8] on the dataset collected by [24], which consists of 68 images corrupted by different blur kernels of the size ranging from 17×17 to 37×37 . We add 1% and 1% Gaussian noise to generate our corrupted observations, respectively. In Fig. 3, the left four subfigures compare curves of iteration errors and reconstruction errors on an example image and the rightmost one illustrate the averaged iteration numbers and run time on the whole dataset. It can be seen that our eFIMA and iFIMA are faster and better than these abstractly designed classical solvers under the same iteration error ($\leq 1e-4$). Moreover, we observe that the performance of these nonconvex APGs is not satisfied when the noise level is bigger. The reconstruction errors of them (Fig. 3(d)) ascend after dozens of steps, while our eFIMA and iFIMA remain lower reconstruction error and fewer iterations. It illustrates that our strategy is more stable than traditional nonconvex APGs in image restoration because of the flexible modules and effective choice mechanisms.

In Fig. 4, we illustrate the visual results of eFIMA and iFIMA with comparisons to both convex image restoration approaches, including FISTA [2] (APG) and FTVd [41] (ADMM), and nonconvex mAPG, APGnc, and niAPG on an example image with 1% noise level but large kernel size (i.e., 75×75) [36]. Here FISTA and FTVd solve their original convex models, while mAPG, APGnc, and niAPG are based on the nonconvex model in Eq. (6). We have that APGs outperformed the original PG. The inexact niAPG is better than exact mAPG and APGnc. Since FTVd is specifically designed for this task, it is the best among all classical solvers, but worse than our FIMA. Overall, iFIMA obtain higher PSNR than eFIMA since the error-control mechanism actually tend to perform more accurate refinements.

We also analyze the iteration behaviors of FIMA in Tab. 1. We report the number of whole iterations and the times

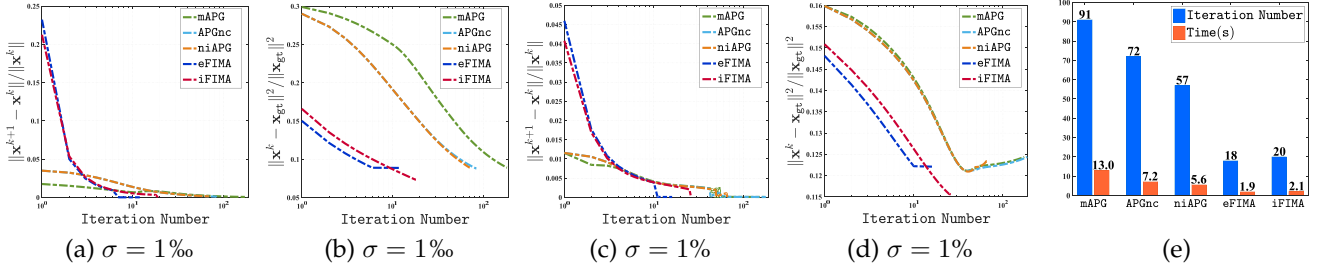


Fig. 3. Comparing iteration behaviors of FIMA to classical nonconvex APGs, including exact ones (mAPG, and APGnc) and inexact niAPG. The left four subfigures compare curves of iteration errors and reconstruction errors with different noise level (1% and 1%), respectively. The rightmost subfigure plot bar charts of the averaged iteration number and “Time(s)” on the dataset [24].

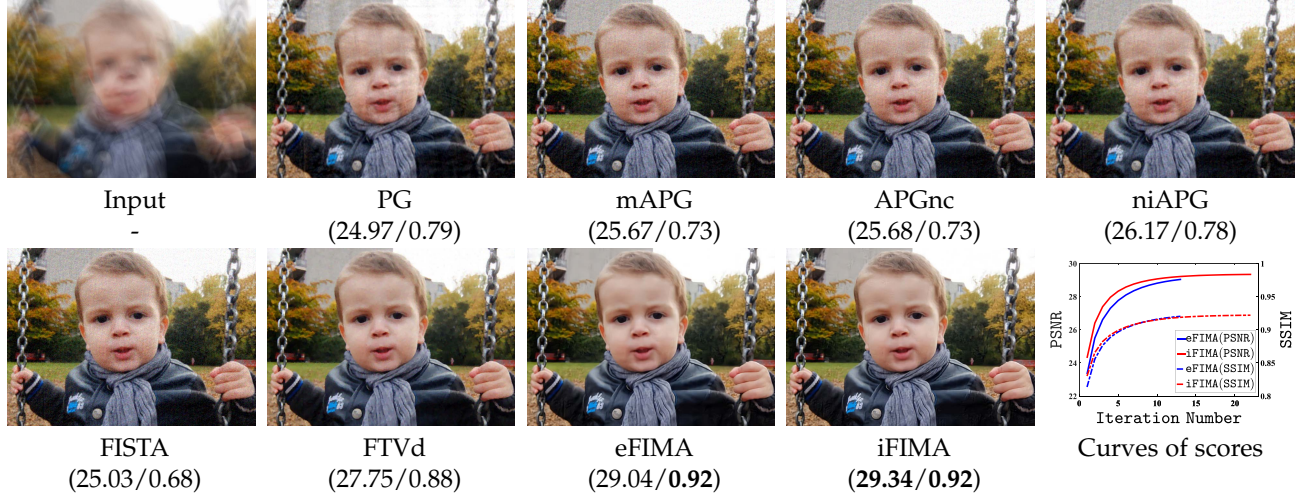


Fig. 4. The non-blind deconvolution performances (1% noise level) of eFIMA and iFIMA with comparisons to convex optimization based algorithms (i.e., FISTA and FTVd), and non-convex solvers (i.e., APGnc, mAPG, and niAPG). The quantitative scores (PSNR/SSIM) are reported below each image. The rightmost subfigure on the bottom row plots the curves of PSNR and SSIM of our methods.



Fig. 5. The non-blind image deconvolution performance (5% noise level) of FIMA with comparisons to existing plug-and-play type methods (i.e., PPADMM and IRCNN). The quantitative scores (PSNR/SSIM) are reported below each image.

the plug-and-play modules $\mathcal{A}_g \circ \mathcal{A}_f$ has been performed by FIMA during iterations. It can be seen that $\mathcal{A}_g \circ \mathcal{A}_f$ are performed in most of the iterations. Moreover, thanks to these user-specified modules, FIMA only needs a dozen or twenty iterations to obtain our desired solutions. In contrast, there are more than 500 iterations in standard PG method. But the practical performances of PG are still worse than our FIMA (see Figs. 4 and 5 for comparisons). These results verify the efficiency and effectiveness of the mechanism of FIMA in real-world applications.

State-of-the-art Comparisons: We compare FIMA with state-of-the-art image restoration approaches, such as ID-DBM3D [42], EPLL [43], PPADMM [25], RTF [44] and

IRCNN [17]. Fig. 5 first compares our FIMA with two prevalent learning-based iterative approaches (i.e., PPADMM and IRCNN) on an example image with 5% noise. Tab. 2 then reports the averaged quantitative results of all the compared methods on the image set (collected by [24]) with different levels of Gaussian noise (i.e., 1%, 2%, 3% and 4%). We have that eFIMA and iFIMA not only outperform classical numerical solvers by a large margin in terms of speed and accuracy, but also achieve better performance than other state-of-the-art approaches. Within FIMA, it can be seen that the speed of eFIMA is faster, while PSNR and SSIM of iFIMA are relatively higher. This is mainly because the “error control” strategy tends to perform more refinements than the

TABLE 2
Averaged PSNR, SSIM and Time(s) on the benchmark image set [24]. Here σ denotes the noise levels.

σ	Metric	State-of-the-art Image Restoration Methods					Classical Nonconvex Methods				Ours	
		IDDBM3D	EPLL	PPADMM	RTF	IRCNN	PG	mAPG	APGnc	niAPG	eFIMA	iFIMA
1%	PSNR	28.83	28.67	28.01	29.12	29.78	27.32	26.68	26.69	27.24	29.81	29.85
	SSIM	0.81	0.81	0.78	0.83	0.84	0.71	0.67	0.67	0.73	0.85	0.85
	Time(s)	193.13	112.03	293.99	249.83	2.67	20.36	13.02	7.16	5.29	1.89	2.06
2%	PSNR	27.60	26.79	26.54	25.58	27.90	25.61	25.20	25.28	25.63	28.02	28.06
	SSIM	0.76	0.74	0.72	0.66	0.78	0.63	0.60	0.61	0.64	0.79	0.79
	Time(s)	198.66	100.52	270.45	254.26	2.68	15.43	7.70	4.66	3.30	1.90	2.07
3%	PSNR	26.72	25.68	25.78	21.18	26.81	24.63	24.39	24.48	24.76	27.05	27.07
	SSIM	0.72	0.69	0.68	0.42	0.73	0.57	0.55	0.56	0.61	0.74	0.75
	Time(s)	191.25	96.32	257.94	252.47	2.68	13.89	6.44	5.37	2.63	1.89	2.07
4%	PSNR	26.06	24.88	25.27	17.95	26.10	24.05	23.88	23.95	24.14	26.20	26.37
	SSIM	0.69	0.65	0.66	0.28	0.70	0.54	0.53	0.53	0.59	0.70	0.72
	Time(s)	183.44	93.82	258.45	255.84	2.67	11.99	6.01	7.82	2.35	1.89	2.07

TABLE 3
Averaged quantitative scores on Levin et al's benchmark.

Method	PSNR	SSIM	ER	KS	Time(s)
Perrone et al.	29.27	0.88	1.35	0.80	113.70
Levin et al.	29.03	0.89	1.40	0.81	41.77
Sun et al.	29.71	0.90	1.32	0.82	209.47
Zhang et al.	28.01	0.86	1.25	0.58	37.45
Pan et al.	29.78	0.89	1.33	0.80	102.60
Ours	30.37	0.91	1.20	0.83	5.65

“explicit momentum” rule during iterations.

6.2 Blind Image Deconvolution

Blind deconvolution is known as one of the most challenging low-level vision tasks. Here we evaluate mFIMA on solving Eq. (8) to address this fundamentally ill-posed multi-variables inverse problem. We adopt the same CNN module $\mathcal{A}_g^{\text{CNN}}$ as that in Sec. 6.1 but train it on image gradient domain to enhance its ability for sharp edge detection.

In Fig. 6, we show the visual performances of mFIMA in different settings (i.e., with and without \mathcal{A}) on an example blurry image from [45]. We observe that mFIMA without \mathcal{A} almost failed on this experiment. This is not surprising since [45], [46] have proved that standard optimization strategy is likely to lead to degenerate global solutions like the delta kernel (frequently called the no-blur solution), or many suboptimal local minima. In contrast, the CNN-based modules successfully avoid trivial results and significantly improve the deconvolution performance. We also plot the curves of quantitative scores (i.e., PSNR for the latent image and Kernel Similarity (KS) for the blur kernel) on the bottom row for these two strategies on the bottom row. As these scores are stable after 20 iterations, here we only plot curves of the first 20 iterations. We then compare mFIMA with state-of-the-art deblurring methods⁶, such as Perrone et al. [47], Levin et al. [45], Sun et al. [46], Zhang et al. [48] and Pan et al. [49] on the most widely-used Levin et al's benchmark [45], which consists of 32 blurred images generated by 4 clean

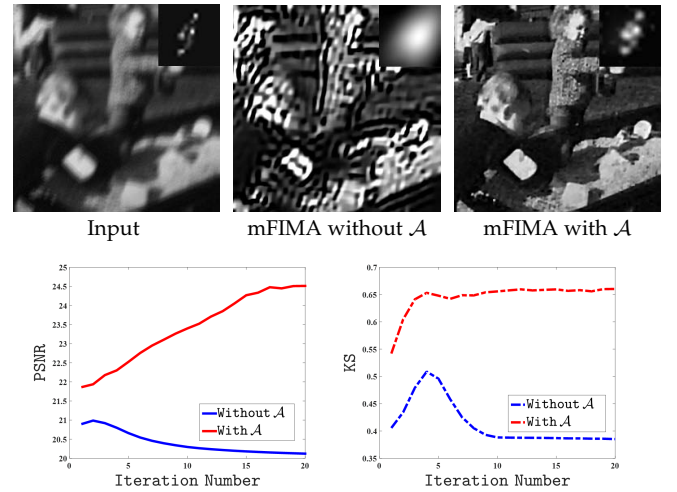


Fig. 6. The comparisons of mFIMA with and without the module \mathcal{A} . The top row compares the visual results of these different strategies. The bottom row plots the curves of PSNR and KS scores during iterations.

images and 8 blur kernels. Tab. 3 reports the averaged quantitative scores, including PSNR, SSIM, and Error Rate (ER) for the latent image, Kernel Similarity (KS) for the blur kernel and the overall run time. Fig. 7 further compares the visual performance of mFIMA to Perrone et al., Sun et al. and Pan et al. (i.e., top 3 in Tab. 3) on a real-world challenging blurry image collected by [36]. It can be seen that mFIMA consistently outperforms all the compared methods both quantitatively and qualitatively, which verifies the efficiency of our proposed learning-based iteration methodology.

In Figs. 8 and 9, we further compare the blind image deconvolution performance of mFIMA with Perrone et al. [47], Sun et al. [46] and Pan et al. [49] (top 3 among all the compared methods in Tab. 2) on example images corrupted by not only unknown blur kernels, but also different levels of Gaussian noises (1% and 3% in Figs. 8 and 9, respectively). It can be seen that mFIMA is robust to these corruptions and outperforms all the compared state-of-the-art deblurring methods.

6. In this and the following experiments, the widely used multi-scale techniques are adopted for all the compared methods.

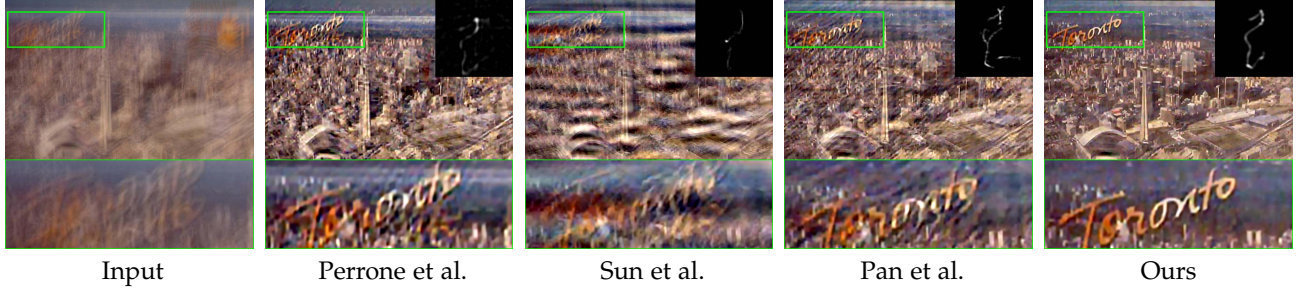


Fig. 7. Visual comparisons between mFIMA and other competitive methods (top 3 in Tab. 3) on a real blurry image.

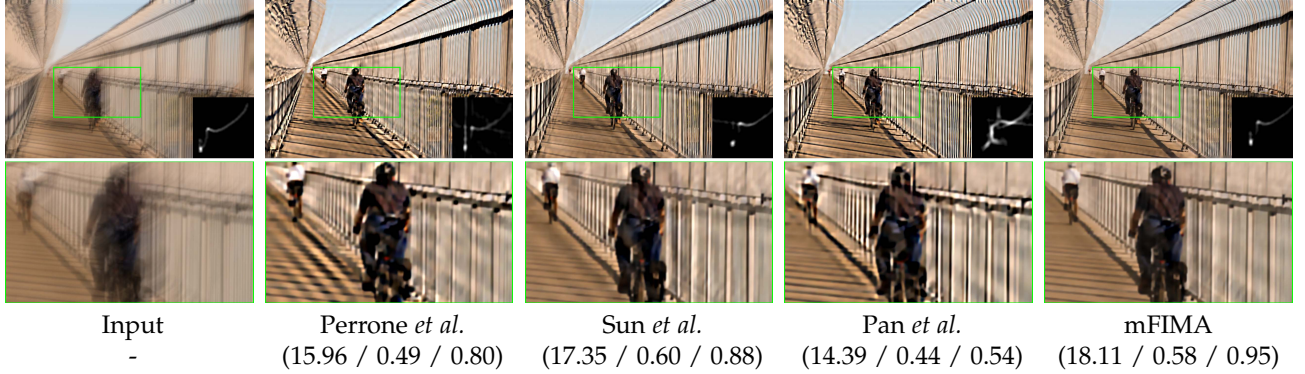


Fig. 8. The blind image deconvolution results of mFIMA with comparisons to state-of-the-art approaches on blurry image with 1% Gaussian noise. The quantitative scores (i.e., PSNR / SSIM / KS) are reported below each image.

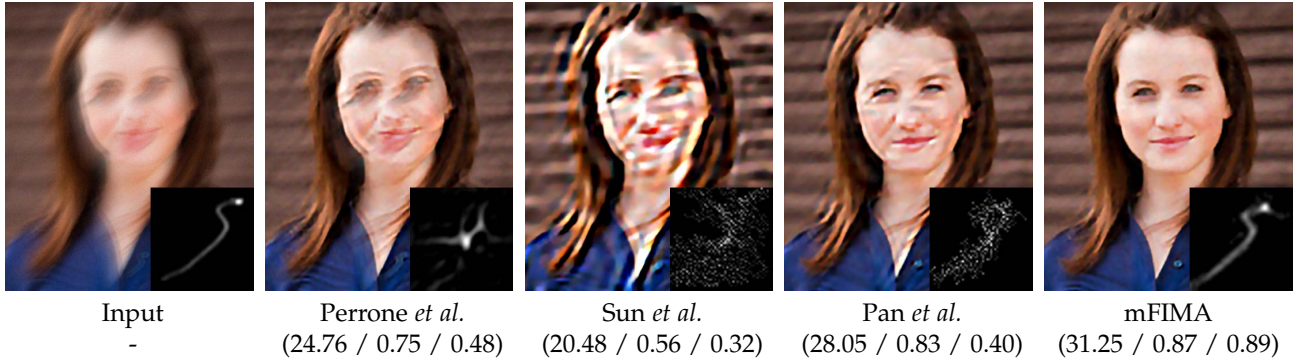


Fig. 9. The blind image deconvolution results of mFIMA with comparisons to state-of-the-art approaches on blurry facial image with 3% Gaussian noise. The quantitative scores (i.e., PSNR / SSIM / KS) are reported below each image.

6.3 Rain Streaks Removal

To further verify our method can deal with various vision tasks, we provide the performance of our mFIMA on rain streaks removal. As we claimed in Sec. 5, we adopt the same CNN architecture to train the learnable \mathcal{A}_{g_x} and \mathcal{A}_{g_c} . It should be noticed that we feed rainy observations into network and output the rain streak images to train the \mathcal{A}_{g_c} , while adopting the similar strategy with deconvolution to train \mathcal{A}_{g_x} .

Firstly, we reported the quantitative scores (PSNR / SSIM) on a widely used synthetic Rain12 dataset [50] and compared with state-of-the-art deraining methods including SR [51], DSC [52], LP [50], DerainNet [53], and DetailNet [54], and UGSM [55]. As Tab. 4 shown, our mFIMA is easily superior to all of the competitive methods. Moreover, we also provide the visual results on the challenging real-world rainy image from [52] in Fig. 10. As can be observed, our proposed

method can remove more rain streaks and preserve the more detail textures than others.

7 CONCLUSION

This paper provided FIMA, a framework to analyze the convergence behaviors of learning-based iterative methods for nonconvex inverse problems. We proposed two novel mechanisms to adaptively guide the trajectories of learning-based iterations and proved their strict convergence. We also showed how to apply FIMA for real-world applications, such as non-blind deconvolution, blind image deconvolution, and rain streaks removal.

APPENDIX A PROOFS

We first give some preliminaries on variational analysis and nonconvex optimization in Sec. A.1. Secs. A.2-A.4 then prove

TABLE 4
Averaged PSNR and SSIM on the benchmark image set [50] for rain streaks removal.

Methods	Metric	Image Index												Avg.
		1	2	3	4	5	6	7	8	9	10	11	12	
SR	PSNR	28.17	29.46	25.55	31.40	26.21	29.05	30.66	27.35	28.71	27.87	26.40	27.97	28.23
	SSIM	0.83	0.88	0.92	0.93	0.87	0.94	0.94	0.94	0.88	0.89	0.88	0.90	0.90
DSC	PSNR	28.07	27.81	22.72	32.59	23.26	26.32	29.85	24.95	29.16	26.21	26.94	27.70	27.13
	SSIM	0.85	0.91	0.84	0.97	0.88	0.96	0.96	0.94	0.90	0.85	0.90	0.87	0.90
LP	PSNR	32.01	32.81	28.99	33.31	28.10	30.33	32.32	30.26	31.15	30.58	29.08	31.21	30.85
	SSIM	0.90	0.94	0.95	0.97	0.93	0.98	0.97	0.97	0.94	0.93	0.93	0.95	0.95
DerainNet	PSNR	28.62	29.58	24.86	33.96	27.22	30.67	33.66	26.75	30.05	27.72	26.22	28.00	28.94
	SSIM	0.90	0.93	0.92	0.98	0.94	0.97	0.98	0.95	0.94	0.91	0.91	0.93	0.94
DetailNet	PSNR	32.55	32.97	29.00	36.27	30.03	31.67	36.13	30.82	33.13	31.56	29.20	31.85	32.10
	SSIM	0.93	0.96	0.95	0.99	0.96	0.98	0.99	0.98	0.97	0.94	0.94	0.96	0.96
UGSM	PSNR	32.42	34.03	27.85	38.03	29.40	34.92	37.94	28.40	30.28	30.33	29.31	30.80	31.98
	SSIM	0.93	0.96	0.94	0.99	0.96	0.99	0.99	0.97	0.92	0.93	0.94	0.94	0.95
mFIMA	PSNR	34.10	36.24	29.46	40.76	31.94	35.10	39.51	35.37	36.66	33.02	31.17	33.41	34.73
	SSIM	0.95	0.98	0.96	0.99	0.97	0.99	0.99	0.99	0.98	0.96	0.95	0.97	0.97



Fig. 10. Rain streaks removal results of mFIMA with comparisons to the state-of-the-art approaches on real-world rainy images.

the main results in our manuscript.

A.1 Preliminaries

Definition 1. [56] The necessary function properties, including proper, lower semi-continuous, Lipschitz smooth, and coercive are summarized as follows. Let $f : \mathbb{R}^D \rightarrow (-\infty, +\infty]$. Then we have

- Proper and lower semi-continuous: f is proper if $\text{dom} f := \{\mathbf{x} \in \mathbb{R}^D : f(\mathbf{x}) < +\infty\}$ is nonempty and $f(\mathbf{x}) > -\infty$. f is lower semi-continuous if $\liminf_{\mathbf{x} \rightarrow \mathbf{y}} f(\mathbf{x}) \geq f(\mathbf{y})$ at any point $\mathbf{y} \in \text{dom} f$.
- Coercive: f is said to be coercive, if f is bounded from below and $f \rightarrow \infty$ if $\|\mathbf{x}\| \rightarrow \infty$, where $\|\cdot\|$ is the ℓ_2 norm.
- L -Lipschitz smooth (i.e., $C_L^{1,1}$): f is L -Lipschitz smooth if f is differentiable and there exists $L > 0$ such that

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^D.$$

If f is L -Lipschitz smooth, we have the following inequality

$$f(\mathbf{x}) \leq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^D.$$

Definition 2. [7], [56] Let $g : \mathbb{R}^D \rightarrow (-\infty, +\infty]$ be a proper and lower semi-continuous function. Then we have

- Sub-differential: The Frecht sub-differential (denoted as $\hat{\partial}g$) of g at point $\mathbf{x} \in \text{dom}(g)$ is the set of all vectors \mathbf{z} which satisfies

$$\liminf_{\mathbf{y} \neq \mathbf{x}, \mathbf{y} \rightarrow \mathbf{x}} \frac{g(\mathbf{y}) - g(\mathbf{x}) - \langle \mathbf{z}, \mathbf{y} - \mathbf{x} \rangle}{\|\mathbf{y} - \mathbf{x}\|} \geq 0,$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product. Then the limiting Frecht sub-differential (denoted as ∂g) at $\mathbf{x} \in \text{dom} g$ is the following closure of $\hat{\partial}g$:

$$\{\mathbf{z} \in \mathbb{R}^n : \exists (\mathbf{x}^k, g(\mathbf{x}^k)) \rightarrow (\mathbf{x}, g(\mathbf{x}))\},$$

where $\mathbf{z}^k \in \hat{\partial}g(\mathbf{x}^k) \rightarrow \mathbf{z}$ when $k \rightarrow \infty$.

- Kurdyka-Łojasiewicz property: g is said to have the Kurdyka-Łojasiewicz property at $\bar{\mathbf{x}} \in \text{dom} \partial g := \{\mathbf{x} \in \mathbb{R}^D : \partial g(\mathbf{x}) \neq \emptyset\}$ if there exist $\eta \in (0, \infty]$, a neighborhood $\mathcal{U}_{\bar{\mathbf{x}}}$ of $\bar{\mathbf{x}}$ and a desingularizing function $\phi : [0, \eta) \rightarrow \mathbb{R}_+$ which satisfies (1) ϕ is continuous at 0 and $\phi(0) = 0$; (2) ϕ is concave and C^1 on $(0, \eta)$; (3) for all $s \in (0, \eta) : \phi'(s) > 0$, such that for all

$$\mathbf{x} \in \mathcal{U}_{\bar{\mathbf{x}}} \cap [g(\bar{\mathbf{x}}) < g(\mathbf{x}) < g(\bar{\mathbf{x}}) + \eta],$$

the following inequality holds

$$\phi'(g(\mathbf{x}) - g(\bar{\mathbf{x}})) \text{dist}(0, \partial g(\mathbf{x})) \geq 1.$$

Moreover, if g satisfies the KL property at each point of $\text{dom} g$ then g is called a KL function.

- **Semi-algebraic set and function:** A subset Ω of \mathbb{R}^D is a real semi-algebraic set if there exist a finite number of real polynomial functions $r_{ij}, h_{ij} : \mathbb{R}^D \rightarrow \mathbb{R}$ such that

$$\Omega = \bigcup_{j=1}^p \bigcap_{i=1}^q \left\{ \mathbf{x} \in \mathbb{R}^D : r_{ij}(\mathbf{x}) = 0 \text{ and } h_{ij}(\mathbf{x}) < 0 \right\}. \quad (11)$$

g is called semi-algebraic if its graph $\{(\mathbf{x}, z) \in \mathbb{R}^{D+1} : g(\mathbf{x}) = z\}$ is a semi-algebraic subset of \mathbb{R}^{D+1} . It is verified in [7] that all semi-algebraic functions satisfy the KL property.

A.2 Explicit Momentum FIMA (eFIMA)

A.2.1 Proof of Theorem 1

Proof. We first prove the inequality relationship of $\Psi(\mathbf{x}^{k+1})$ and $\Psi(\mathbf{v}^k)$. According to the update rule of \mathbf{x}^{k+1} (Step 8 in Alg. 1): $\mathbf{x}^{k+1} \in \text{prox}_{\gamma^k g}(\mathbf{v}^k - \gamma^k \nabla f(\mathbf{v}^k))$, we have

$$\mathbf{x}^{k+1} \in \arg \min_{\mathbf{x}} g(\mathbf{x}) + \langle \nabla f(\mathbf{v}^k), \mathbf{x} - \mathbf{v}^k \rangle + \frac{1}{2\gamma^k} \|\mathbf{x} - \mathbf{v}^k\|^2, \quad (12)$$

thus

$$g(\mathbf{x}^{k+1}) + \langle \nabla f(\mathbf{v}^k), \mathbf{x}^{k+1} - \mathbf{v}^k \rangle + \frac{1}{2\gamma^k} \|\mathbf{x}^{k+1} - \mathbf{v}^k\|^2 \leq g(\mathbf{v}^k). \quad (13)$$

Since f is $C_L^{1,1}$, we have

$$f(\mathbf{x}^{k+1}) \leq f(\mathbf{v}^k) + \langle \nabla f(\mathbf{v}^k), \mathbf{x}^{k+1} - \mathbf{v}^k \rangle + \frac{L}{2} \|\mathbf{x}^{k+1} - \mathbf{v}^k\|^2, \quad (14)$$

where L is the Lipschitz moduli of ∇f . Combining this with Eqs. (13) and (14), we have

$$\Psi(\mathbf{x}^{k+1}) \leq \Psi(\mathbf{v}^k) - \left(\frac{1}{2\gamma^k} - \frac{L}{2} \right) \|\mathbf{x}^{k+1} - \mathbf{v}^k\|^2. \quad (15)$$

Set $\gamma^k < 1/L$ and define $\alpha^k = \frac{1}{2\gamma^k} - \frac{L}{2}$, we have $\alpha^k > 0$ and $\Psi(\mathbf{x}^{k+1}) \leq \Psi(\mathbf{v}^k) - \alpha^k \|\mathbf{x}^{k+1} - \mathbf{v}^k\|^2$.

Then we prove the boundness and convergence of $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$. Based on the momentum scheduling policy in Alg. 1, we obviously have $\Psi(\mathbf{v}^k) \leq \Psi(\mathbf{x}^k)$. This together with the result in Eq. (15) (i.e., $\Psi(\mathbf{x}^{k+1}) \leq \Psi(\mathbf{v}^k)$ with $\gamma^k < 1/L$) concludes that for any $k \in \mathbb{N}_+$,

$$\Psi(\mathbf{x}^{k+1}) \leq \Psi(\mathbf{v}^k) \leq \Psi(\mathbf{x}^k) \leq \Psi(\mathbf{v}^{k-1}) \leq \Psi(\mathbf{x}^0). \quad (16)$$

Since both f and g are proper, we also have $\Psi(\mathbf{v}^k) \geq \inf \Psi > -\infty$. Thus both sequences $\{\Psi(\mathbf{x}^k)\}_{k \in \mathbb{N}}$ and $\{\Psi(\mathbf{v}^k)\}_{k \in \mathbb{N}}$ are non-increasing and bounded. This together with the coercive of Ψ concludes that both $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ and $\{\mathbf{v}^k\}_{k \in \mathbb{N}}$ are bounded and thus have accumulation points.

Then we prove that all accumulation points are the critical points of Ψ . From Eq. (16), we actually have that the objective sequences $\{\Psi(\mathbf{x}^k)\}_{k \in \mathbb{N}}$ and $\{\Psi(\mathbf{v}^k)\}_{k \in \mathbb{N}}$ converge to the same value Ψ^* , i.e.,

$$\lim_{k \rightarrow \infty} \Psi(\mathbf{x}^k) = \lim_{k \rightarrow \infty} \Psi(\mathbf{v}^k) = \Psi^*. \quad (17)$$

From Eqs. (15) and (16), we have

$$\left(\frac{1}{2\gamma^k} - \frac{L}{2} \right) \|\mathbf{x}^{k+1} - \mathbf{v}^k\|^2 \leq \Psi(\mathbf{v}^k) - \Psi(\mathbf{x}^{k+1}) \leq \Psi(\mathbf{x}^k) - \Psi(\mathbf{x}^{k+1}). \quad (18)$$

Summing over k , we further have

$$\min_k \left\{ \frac{1}{2\gamma^k} - \frac{L}{2} \right\} \sum_{k=0}^{\infty} \|\mathbf{x}^{k+1} - \mathbf{v}^k\|^2 \leq \Psi(\mathbf{x}^0) - \Psi^* < \infty. \quad (19)$$

The above inequality implies that $\|\mathbf{x}^{k+1} - \mathbf{v}^k\| \rightarrow 0$ and hence $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ and $\{\mathbf{v}^k\}_{k \in \mathbb{N}}$ share the same set of accumulation points (denoted as Ω). Consider that $\mathbf{x}^* \in \Omega$ is any accumulation point of $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$, i.e., $\mathbf{x}^{k_j} \rightarrow \mathbf{x}^*$ if $j \rightarrow \infty$. Then by Eq. (12), we have

$$g(\mathbf{x}^{k+1}) + \langle \nabla f(\mathbf{v}^k), \mathbf{x}^{k+1} - \mathbf{v}^k \rangle + \frac{1}{2\gamma^k} \|\mathbf{x}^{k+1} - \mathbf{v}^k\|^2 \leq g(\mathbf{x}^*) + \langle \nabla f(\mathbf{v}^k), \mathbf{x}^* - \mathbf{v}^k \rangle + \frac{1}{2\gamma^k} \|\mathbf{x}^* - \mathbf{v}^k\|^2. \quad (20)$$

Let $k_j = k + 1$ in Eq. (20) and $j \rightarrow \infty$, by taking \limsup on both sides of Eq. (20), we have $\limsup_{j \rightarrow \infty} g(\mathbf{x}^{k_j}) \leq g(\mathbf{x}^*)$.

On the other hand, since g is lower semi-continuous and $\mathbf{x}^{k_j} \rightarrow \mathbf{x}^*$, it follows that $\liminf_{j \rightarrow \infty} g(\mathbf{x}^{k_j}) \geq g(\mathbf{x}^*)$. So we have $\lim_{j \rightarrow \infty} g(\mathbf{x}^{k_j}) = g(\mathbf{x}^*)$. Note that the continuity of f yields $\lim_{j \rightarrow \infty} f(\mathbf{x}^{k_j}) = f(\mathbf{x}^*)$, so we conclude

$$\lim_{j \rightarrow \infty} \Psi(\mathbf{x}^{k_j}) = \Psi(\mathbf{x}^*). \quad (21)$$

Recall that $\lim_{k \rightarrow \infty} \Psi(\mathbf{x}^{k+1}) = \Psi^*$ in Eq. (17), we have $\lim_{j \rightarrow \infty} \Psi(\mathbf{x}^{k_j}) = \Psi^*$, so

$$\Psi(\mathbf{x}^*) = \Psi^*, \forall \mathbf{x}^* \in \Omega. \quad (22)$$

By first-order optimality condition of Eq. (12) and $k_j = k + 1$, we have

$$\mathbf{0} \in \partial g(\mathbf{x}^{k_j}) + \nabla f(\mathbf{v}^k) + \frac{1}{\gamma^k} (\mathbf{x}^{k_j} - \mathbf{v}^k). \quad (23)$$

Thus, we have

$$\begin{aligned} \nabla f(\mathbf{x}^{k_j}) - \nabla f(\mathbf{v}^k) - \frac{1}{\gamma^k} (\mathbf{x}^{k_j} - \mathbf{v}^k) &\in \partial \Psi(\mathbf{x}^{k_j}) \\ \Rightarrow \|\nabla f(\mathbf{x}^{k_j}) - \nabla f(\mathbf{v}^k) - \frac{1}{\gamma^k} (\mathbf{x}^{k_j} - \mathbf{v}^k)\| &\leq \left(L + \frac{1}{\gamma^k} \right) \|\mathbf{x}^{k_j} - \mathbf{v}^k\| \rightarrow 0, \text{ as } j \rightarrow \infty. \end{aligned} \quad (24)$$

Then from the definition of sub-differential and Eqs. (21), (23), and (24), we conclude that

$$\mathbf{0} \in \partial \Psi(\mathbf{x}^*), \forall \mathbf{x}^* \in \Omega. \quad (25)$$

Therefore, we have that all accumulation points \mathbf{x}^* are the critical points of Ψ . \square

A.2.2 Proof of Corollary 1

Proof. Considering the semi-algebraic (thus KL) property of $\Psi(\mathbf{x})$ and defining a desingularizing function with the form $\phi(s) = \frac{t}{\theta} s^\theta$, we can prove **Corollary 1** by Eqs. (15), (16), and (23) using similar methodology as that in [7], [27]. Since these derivations are quite standard, we omit details of this proof in our Supplemental Materials. \square

A.3 Implicit Momentum FIMA (iFIMA)

A.3.1 Proof of Proposition 1

Proof. First, by using the same derivations as that in Eq. (15), we can directly obtain the inequality in Theorem 1 for iFIMA. Then we show how to build the relationship between $\Psi(\tilde{\mathbf{u}}^k)$

and $\Psi(\mathbf{x}^k)$. It is known that $\tilde{\mathbf{u}}^k$ is actually an inexact minimizer of Ψ^k . But by defining its sub-differential $\mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k} \in \partial\Psi^k(\tilde{\mathbf{u}}^k)$ as that in Eq. (4), we can also consider it as the exact solution to the following problem

$$\tilde{\mathbf{u}}^k \in \arg \min_{\mathbf{x}} \Psi^k(\mathbf{x}) - \langle \mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}, \mathbf{x} \rangle. \quad (26)$$

Thus, we have

$$\begin{aligned} & \Psi(\tilde{\mathbf{u}}^k) + \frac{\mu^k}{2} \|\tilde{\mathbf{u}}^k - \mathbf{x}^k\|^2 - \langle \mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}, \tilde{\mathbf{u}}^k \rangle \\ & \leq \Psi(\mathbf{x}^k) - \langle \mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}, \mathbf{x}^k \rangle \\ & \Rightarrow \Psi(\tilde{\mathbf{u}}^k) \leq \Psi(\mathbf{x}^k) - \frac{\mu^k}{2} \|\tilde{\mathbf{u}}^k - \mathbf{x}^k\|^2 + \langle \mathbf{d}_{\Psi^k}^{\tilde{\mathbf{u}}^k}, \tilde{\mathbf{u}}^k - \mathbf{x}^k \rangle \\ & \leq \Psi(\mathbf{x}^k) - \frac{\mu^k}{2} \|\tilde{\mathbf{u}}^k - \mathbf{x}^k\|^2 + C^k \|\tilde{\mathbf{u}}^k - \mathbf{x}^k\|^2 \\ & = \Psi(\mathbf{x}^k) - \left(\frac{\mu^k}{2} - C^k \right) \|\tilde{\mathbf{u}}^k - \mathbf{x}^k\|^2, \end{aligned} \quad (27)$$

in which the second inequality holds under Cauchy-Schwarz inequality and our error-control-based scheduling policy in Alg. 2. Set $C^k < \frac{\mu^k}{2}$ and define $\beta^k = \frac{\mu^k}{2} - C^k$, we have $\beta^k > 0$ and $\Psi(\tilde{\mathbf{u}}^k) \leq \Psi(\mathbf{x}^k) - \beta^k \|\tilde{\mathbf{u}}^k - \mathbf{x}^k\|^2$, which concludes the proof. \square

A.3.2 Proof of Theorem 2

Proof. We first prove the boundedness of $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$. According to **Proposition 1** we have $\Psi(\tilde{\mathbf{u}}^k) \leq \Psi(\mathbf{x}^k)$ when $\mu^k/2 > C^k$. So if the error-control criteria in Alg. 2 is satisfied, we have $\mathbf{v}^k = \tilde{\mathbf{u}}^k$, $\Psi(\mathbf{v}^k) = \Psi(\tilde{\mathbf{u}}^k) \leq \Psi(\mathbf{x}^k)$, otherwise, we have $\mathbf{v}^k = \mathbf{x}^k$, $\Psi(\mathbf{v}^k) = \Psi(\mathbf{x}^k)$. This together with the results in Theorem 1 (i.e., $\Psi(\mathbf{x}^{k+1}) \leq \Psi(\mathbf{v}^k)$ with $\gamma^k < 1/L$) concludes that for any $k \in \mathbb{N}$,

$$\Psi(\mathbf{x}^{k+1}) \leq \Psi(\mathbf{v}^k) \leq \Psi(\mathbf{x}^k) \leq \Psi(\mathbf{v}^{k-1}) \leq \Psi(\mathbf{x}^0). \quad (28)$$

Then by using similar derivations as that in **Theorem 1**, we have that all accumulation points \mathbf{x}^* are the critical points of Ψ .

Now we are ready to prove that $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ is a Cauchy sequence. Since Ψ is a KLfunction, we have

$$\varphi'(\Psi(\mathbf{x}^{k+1}) - \Psi(\mathbf{x}^*)) \text{dist}(0, \partial\Psi(\mathbf{x}^{k+1})) \geq 1.$$

From Eq. (24) we get that

$$\varphi'(\Psi(\mathbf{x}^{k+1}) - \Psi(\mathbf{x}^*)) \geq \frac{1}{L + \frac{1}{\gamma^k}} \|\mathbf{x}^{k+1} - \mathbf{v}^k\|^{-1}$$

On the other hand, from the concavity of φ and Eqs. (15) and (24) we have that

$$\begin{aligned} & \varphi(\Psi(\mathbf{x}^{k+1}) - \Psi(\mathbf{x}^*)) - \varphi(\Psi(\mathbf{x}^{k+2}) - \Psi(\mathbf{x}^*)) \\ & \geq \varphi'(\Psi(\mathbf{x}^{k+1}) - \Psi(\mathbf{x}^*))(\varphi(\Psi(\mathbf{x}^{k+1}) - \Psi(\mathbf{x}^{k+2}))) \\ & \geq \frac{1}{L + \frac{1}{\gamma^k}} \|\mathbf{x}^{k+1} - \mathbf{v}^k\|^{-1} \left(\frac{1}{2\gamma^k} - \frac{L}{2} \right) \|\mathbf{x}^{k+2} - \mathbf{v}^{k+1}\|^2 \end{aligned}$$

For convenience, we define for all $m, n \in \mathbb{N}$ and \mathbf{x}^* the following quantities

$$\Delta_{m,n} := \varphi(\Psi(\mathbf{x}^m) - \Psi(\mathbf{x}^*)) - \varphi(\Psi(\mathbf{x}^n) - \Psi(\mathbf{x}^*)),$$

and

$$B := \min_k \frac{2L\gamma^k + 2}{1 - L\gamma^k} \in (0, \infty).$$

These deduce that

$$\Delta_{k+1,k+2} \geq \frac{\|\mathbf{x}^{k+2} - \mathbf{v}^{k+1}\|^2}{B \|\mathbf{x}^{k+1} - \mathbf{v}^k\|}.$$

and hence

$$\begin{aligned} & \|\mathbf{x}^{k+2} - \mathbf{v}^{k+1}\|^2 \leq E\Delta_{k+1,k+2} \|\mathbf{x}^{k+1} - \mathbf{v}^k\| \\ & \Rightarrow 2\|\mathbf{x}^{k+2} - \mathbf{v}^{k+1}\| \leq E\Delta_{k+1,k+2} + \|\mathbf{x}^{k+1} - \mathbf{v}^k\|. \end{aligned}$$

Summing up the above inequality for $k = p, \dots, q \in \mathbb{N}$ and $p < q$ yields

$$\begin{aligned} & 2 \sum_{k=p}^q \|\mathbf{x}^{k+2} - \mathbf{v}^{k+1}\| \leq \sum_{k=p}^q \|\mathbf{x}^{k+1} - \mathbf{v}^k\| + E \sum_{k=p}^q \Delta_{k+1,k+2} \\ & \leq \sum_{k=p}^q \|\mathbf{x}^{k+2} - \mathbf{v}^{k+1}\| + \|\mathbf{x}^{p+1} - \mathbf{v}^p\| + E\Delta_{p+1,q+2}. \end{aligned}$$

where the last inequality holds under the fact $\Delta_{m,n} + \Delta_{n,r} = \Delta_{m,r}$ for all $m, n, r \in \mathbb{N}$. Since $\varphi > 0$, we thus have for any $k > l$ that

$$\sum_{k=p}^q \|\mathbf{x}^{k+2} - \mathbf{v}^{k+1}\| \leq \|\mathbf{x}^{p+1} - \mathbf{v}^p\| + E\varphi(\Psi(\mathbf{x}^{p+1}) - \Psi(\mathbf{x}^*)). \quad (29)$$

Moreover, recalled the Eq. (27), we also has

$$\begin{aligned} & \min_k \left\{ \frac{\mu^k}{2} - C^k \right\} \sum_{k=p}^q \|\mathbf{v}^{k+1} - \mathbf{x}^{k+1}\| \\ & \leq \sum_{k=p}^q \Psi(\mathbf{x}^{k+1}) - \Psi(\mathbf{x}^{k+2}) = \Psi(\mathbf{x}^{p+1}) - \Psi(\mathbf{x}^{q+2}) \end{aligned} \quad (30)$$

Combing with Eqs. (29) and (30), we easily deduce

$$\begin{aligned} & \sum_{k=1}^{\infty} \|\mathbf{x}^{k+1} - \mathbf{x}^k\| \\ & \leq \sum_{k=1}^{\infty} \|\mathbf{x}^{k+1} - \mathbf{v}^k\| + \sum_{k=1}^{\infty} \|\mathbf{v}^k - \mathbf{x}^k\| < \infty \end{aligned} \quad (31)$$

The above inequality implies that $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ is a Cauchy sequence. Thus the sequence globally converges to a critical point of $\Psi(\mathbf{x})$ in Eq. (1). \square

A.4 Multi-block FIMA

A.4.1 Definition Extension

As for the generalized Lipschitz smooth property of f , we actually need f satisfy that

- For each \mathbf{x}_n with other variables fixed, there exists $L_n > 0$ such that

$$\begin{aligned} & \|\nabla_n f(\mathfrak{X}_{[<n]}, \mathbf{x}_n, \mathfrak{X}_{[>n]}) - \nabla_n f(\mathfrak{X}_{[<n]}, \mathbf{y}_n, \mathfrak{X}_{[>n]})\| \\ & \leq L_n (\mathfrak{X}_{[<n]}, \mathfrak{X}_{[>n]}) \|\mathbf{x}_n - \mathbf{y}_n\|, \forall \mathbf{x}_n, \mathbf{y}_n \in \mathbb{R}^{D_n}, \end{aligned} \quad (32)$$

where ∇_n denotes the gradient with respect to \mathbf{x}_n .

- For each bounded subset $\Omega_1 \times \dots \times \Omega_N \subseteq \mathbb{R}^{D_1} \times \dots \times \mathbb{R}^{D_N}$, there exists $M > 0$ such that

$$\begin{aligned} & \|(\nabla_1 f(\mathfrak{X}) - \nabla_1 f(\mathfrak{Y}), \dots, \nabla_N f(\mathfrak{X}) - \nabla_N f(\mathfrak{Y}))\| \\ & \leq M \|\mathfrak{X} - \mathfrak{Y}\|, \forall \mathfrak{X}, \mathfrak{Y} \in \Omega_1 \times \dots \times \Omega_N. \end{aligned} \quad (33)$$

A.4.2 Proof of Corollary 2

Proof. We first prove the boundedness of $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$. Using the inequality in Theorem 1 and Step 10 in Alg. 3, we have

$$\begin{aligned} & \mathbf{x}_n^{k+1} \in \arg \min_{\mathbf{x}_n} g_n(\mathbf{x}_n) + \frac{1}{2\gamma_n^k} \|\mathbf{x}_n - \mathbf{v}_n^k\|^2 \\ & \quad + \left\langle \nabla_n f(\mathfrak{X}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathfrak{X}_{[>n]}^k), \mathbf{x}_n - \mathbf{v}_n^k \right\rangle. \end{aligned} \quad (34)$$

This together with Eq. (32) concludes that

$$\begin{aligned} & \Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k) - \Psi(\mathbf{x}_{[\leq n]}^{k+1}, \mathbf{x}_{[>n]}^k) \\ & \geq \left(\frac{1}{2\gamma_n^k} - \frac{L_n(\mathbf{x}_{[<n]}^{k+1}, \mathbf{x}_{[>n]}^k)}{2} \right) \|\mathbf{x}_n^{k+1} - \mathbf{v}_n^k\|^2. \end{aligned} \quad (35)$$

Define the auxiliary function $\Psi_n^k(\mathbf{x}_n) = f(\mathbf{x}_{[<n]}^{k+1}, \mathbf{x}_n, \mathbf{x}_{[>n]}^k) + g_n(\mathbf{x}_n) + \frac{\mu_n^k}{2} \|\mathbf{x}_n - \mathbf{x}_n^k\|^2$. Then by considering that $\tilde{\mathbf{u}}_n^k$ is an inexact solution of the auxiliary function $\Psi_n^k(\mathbf{x}_n)$, and applying **Proposition 1**, we have

$$\begin{aligned} & \Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{x}_n^k, \mathbf{x}_{[>n]}^k) - \Psi(\mathbf{x}_{[<n]}^{k+1}, \tilde{\mathbf{u}}_n^k, \mathbf{x}_{[>n]}^k) \\ & \geq \left(\frac{\mu_n^k}{2} - C_n^k \right) \|\tilde{\mathbf{u}}_n^k - \mathbf{x}_n^k\|^2. \end{aligned} \quad (36)$$

Let $\lambda_n^k = L_n(\mathbf{x}_{[<n]}^{k+1}, \mathbf{x}_{[>n]}^k)$, then consider Eq. (35) with $\gamma_n^k < 1/\lambda_n^k$, Eq. (36) with $\mu_n^k > 2C_n^k$, and our error-control updating rule, we have

$$\begin{aligned} & \Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{x}_{[>n]}^k) \leq \Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k) \\ & \leq \Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{x}_n^k, \mathbf{x}_{[>n]}^k). \end{aligned}$$

It concludes that for any $k \in \mathbb{N}_+$ and $n \in \{1, \dots, N\}$,

$$\begin{aligned} & \Psi(\mathbf{x}^{k+1}) = \Psi(\mathbf{x}_{[\leq N]}^{k+1}, \mathbf{x}_{[>N]}^k) \\ & \leq \Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k) \leq \Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{x}_n^k, \mathbf{x}_{[>n]}^k) \\ & \leq \Psi(\mathbf{x}_{[<1]}^{k+1}, \mathbf{x}_1^k, \mathbf{x}_{[>1]}^k) = \Psi(\mathbf{x}^k) \leq \dots \leq \Psi(\mathbf{x}^0). \end{aligned} \quad (37)$$

Since f, g_n are proper, we also have $-\infty < \inf \Psi \leq \Psi(\mathbf{x}^{k+1})$. Thus $\{\Psi(\mathbf{x}^k)\}_{k \in \mathbb{N}}$ and $\{\Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k)\}_{k \in \mathbb{N}}$ are all non-increasing and bounded. This together with the coercive of Ψ concludes that the sequences $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ and $\{\mathbf{v}_n^k\}_{k \in \mathbb{N}} (1 \leq n \leq N)$ are bounded and thus have accumulation points.

Then we prove that all accumulation points are the critical points of Ψ . From Eq. (37), we have that the function value sequences $\{\Psi(\mathbf{x}^k)\}_{k \in \mathbb{N}}$ and $\{\Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k)\}_{k \in \mathbb{N}}$ converge to the same value Ψ^* , i.e.,

$$\lim_{k \rightarrow \infty} \Psi(\mathbf{x}^k) = \lim_{k \rightarrow \infty} \Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k) = \Psi^*. \quad (38)$$

From Eqs. (35) and (37), summing over k and n we have

$$\begin{aligned} & \min_{k,n} \left\{ \frac{1}{2\gamma_n^k} - \frac{\lambda_n^k}{2} \right\} \sum_{k=0}^{\infty} \sum_{n=1}^N \|\mathbf{x}_n^{k+1} - \mathbf{v}_n^k\|^2 \\ & \leq \sum_{k=0}^{\infty} \sum_{n=1}^N \left(\Psi(\mathbf{x}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k) - \Psi(\mathbf{x}_{[\leq n]}^{k+1}, \mathbf{x}_{[>n]}^k) \right) \\ & \leq \sum_{k=0}^{\infty} (\Psi(\mathbf{x}^k) - \Psi(\mathbf{x}^{k+1})) = \Psi(\mathbf{x}^0) - \Psi^* < \infty. \end{aligned} \quad (39)$$

The above inequality implies that $\|\mathbf{x}_n^{k+1} - \mathbf{v}_n^k\| \rightarrow 0$, hence $\{\mathbf{x}_n^k\}_{k \in \mathbb{N}}$ and $\{\mathbf{v}_n^k\}_{k \in \mathbb{N}}$ share the same set of accumulation points. Consider that $\mathbf{x}^* = \{\mathbf{x}_1^*, \dots, \mathbf{x}_N^*\}$ is any accumulation point of $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$, there exists a subsequence $\{\mathbf{x}^{k_j}\}_{j \in \mathbb{N}}$ such that

$$\mathbf{x}_n^{k_j} \rightarrow \mathbf{x}_n^*, \text{ as } j \rightarrow \infty. \quad (40)$$

From Eq. (34), we have

$$\begin{aligned} & g_n(\mathbf{x}_n^{k+1}) + \left\langle \nabla_n f(\mathbf{x}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k), \mathbf{x}_n^{k+1} - \mathbf{v}_n^k \right\rangle \\ & + \frac{1}{2\gamma_n^k} \|\mathbf{x}_n^{k+1} - \mathbf{v}_n^k\|^2 \\ & \leq g_n(\mathbf{x}_n^*) + \left\langle \nabla_n f(\mathbf{x}_{[<n]}^{k+1}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k), \mathbf{x}_n^* - \mathbf{v}_n^k \right\rangle \\ & + \frac{1}{2\gamma_n^k} \|\mathbf{x}_n^* - \mathbf{v}_n^k\|^2. \end{aligned} \quad (41)$$

Let $k_j = k + 1$ in Eq. (41) and $j \rightarrow \infty$, by taking \limsup on both sides, we have $\limsup_{j \rightarrow \infty} g_n(\mathbf{x}_n^{k_j}) \leq g_n(\mathbf{x}_n^*)$. On the

other hand, since g_n is lower semi-continuous and $\mathbf{x}_n^{k_j} \rightarrow \mathbf{x}_n^*$, it follows that $\liminf_{j \rightarrow \infty} g_n(\mathbf{x}_n^{k_j}) \geq g_n(\mathbf{x}_n^*)$. Thus we have

$\lim_{j \rightarrow \infty} g_n(\mathbf{x}_n^{k_j}) = g_n(\mathbf{x}_n^*)$. Note that the continuity of f yields $\lim_{j \rightarrow \infty} f(\mathbf{x}^{k_j}) = f(\mathbf{x}^*)$. Therefore, we conclude

$$\begin{aligned} & \lim_{j \rightarrow \infty} f(\mathbf{x}^{k_j}) + \sum_{n=1}^N g_n(\mathbf{x}_n^{k_j}) = f(\mathbf{x}^*) + \sum_{n=1}^N g_n(\mathbf{x}_n^*) \\ & \Rightarrow \lim_{j \rightarrow \infty} \Psi(\mathbf{x}^{k_j}) = \Psi(\mathbf{x}^*). \end{aligned} \quad (42)$$

Considering $\lim_{k \rightarrow \infty} \Psi(\mathbf{x}^{k+1}) = \Psi^*$ in Eq. (38), we have $\lim_{j \rightarrow \infty} \Psi(\mathbf{x}^{k_j}) = \Psi^*$, and thus

$$\Psi(\mathbf{x}^*) = \Psi^*. \quad (43)$$

Similar with the deduction in **Theorem 2**, we conclude that $\|\mathbf{x}^{k+1} - \mathbf{x}^k\| \rightarrow 0, \|\mathbf{x}_n^{k+1} - \mathbf{v}_n^k\| \rightarrow 0$ when $k \rightarrow \infty$.

By considering the first-order optimality condition of Eq. (34) and setting $k_j = k + 1$, we have

$$\begin{aligned} & 0 \in \partial_n g_n(\mathbf{x}_n^{k_j}) + \nabla_n f(\mathbf{x}_{[<n]}^{k_j}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k) + \frac{1}{\gamma_n^k} (\mathbf{x}_n^{k_j} - \mathbf{v}_n^k) \\ & \Leftrightarrow \nabla_n f(\mathbf{x}^{k_j}) - \nabla_n f(\mathbf{x}_{[<n]}^{k_j}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k) - \frac{1}{\gamma_n^k} (\mathbf{x}_n^{k_j} - \mathbf{v}_n^k) \\ & \in \partial_n \Psi(\mathbf{x}^{k_j}) \\ & \Rightarrow \left\| \nabla_n f(\mathbf{x}^{k_j}) - \nabla_n f(\mathbf{x}_{[<n]}^{k_j}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k) - \frac{1}{\gamma_n^k} (\mathbf{x}_n^{k_j} - \mathbf{v}_n^k) \right\| \\ & \leq M \left\| \mathbf{x}^{k_j} - (\mathbf{x}_{[<n]}^{k_j}, \mathbf{v}_n^k, \mathbf{x}_{[>n]}^k) \right\| + \frac{1}{\gamma_n^k} \|\mathbf{x}_n^{k_j} - \mathbf{v}_n^k\| \\ & \leq M \left(\sum_{i=1}^N \|\mathbf{x}_i^{k_j} - \mathbf{x}_i^k\| + \|\mathbf{x}_n^{k_j} - \mathbf{v}_n^k\| \right) + \frac{1}{\gamma_n^k} \|\mathbf{x}_n^{k_j} - \mathbf{v}_n^k\| \\ & \leq M \|\mathbf{x}^{k_j} - \mathbf{x}^k\| + (M + \frac{1}{\gamma_n^k}) \|\mathbf{x}_n^{k_j} - \mathbf{v}_n^k\| \rightarrow 0 \end{aligned} \quad (44)$$

when $j \rightarrow \infty$. Here ∂_n denotes the partial sub-differential with respect to \mathbf{x}_n , $\gamma^- = \inf \{\gamma_n^k : k \in \mathbb{N}, n = 1, \dots, N\}$ and M is defined in Eq. (33). Therefore, combining Eqs. (38), (43) and (44) with the definition of sub-differential, we finally concludes

$$\begin{aligned} & \|\partial \Psi(\mathbf{x}^{k_j})\| \leq \sum_{n=1}^N \|\partial_n \Psi(\mathbf{x}^{k_j})\| \\ & \leq NM \|\mathbf{x}^{k_j} - \mathbf{x}^k\| + (M + \frac{1}{\gamma^-}) \sum_{n=1}^N \|\mathbf{x}_n^{k_j} - \mathbf{v}_n^k\| \rightarrow 0 \end{aligned} \quad (45)$$

when $j \rightarrow \infty$, i.e.,

$$0 \in \partial \Psi(\mathbf{x}^*), \forall \mathbf{x}^* \in \Omega, \quad (46)$$

where Ω denotes the set including all accumulation points of $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$. Therefore, we have that all accumulation points \mathbf{x}^* are the critical points of Ψ .

Finally, based on the KL property of $\Psi(\mathfrak{x})$ and using similar derivations as that in the proof of **Theorem 2**, we also have

$$\sum_{k=0}^{\infty} \|\mathfrak{x}^{k+1} - \mathfrak{x}^k\| < \infty. \quad (47)$$

It is clear that Eq. (47) implies that the sequence $\{\mathfrak{x}^k\}_{k \in \mathbb{N}}$ is a Cauchy sequence, thus is globally converged to the critical points of $\Psi(\mathfrak{x})$ in Eq. (5).

Considering $\Psi(\mathfrak{x}^k)$ is semi-algebraic and choosing $\phi(s) = \frac{t}{\theta} s^\theta$ as the desingularizing function, it is also easy to conclude that mFIMA still shares the same convergence rates stated in **Corollary 1**. \square

ACKNOWLEDGMENTS

This work is partially supported by the National Natural Science Foundation of China (Nos. 61672125, 61733002, 61572096, 61432003 and 61632019), and the Fundamental Research Funds for the Central Universities.

REFERENCES

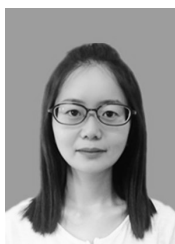
- [1] Y. Nesterov, "A method of solving a convex programming problem with convergence rate $o(1/k^2)$," in *Soviet Mathematics Doklady*, vol. 27, no. 2, 1983, pp. 372–376.
- [2] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM Journal on Imaging Sciences*, vol. 2, no. 1, pp. 183–202, 2009.
- [3] L. Xu, C. Lu, Y. Xu, and J. Jia, "Image smoothing via ℓ_0 gradient minimization," *ACM Trans. on Graphics*, vol. 30, no. 6, p. 174, 2011.
- [4] E. J. Candes, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis," *Journal of the ACM*, vol. 58, no. 3, p. 11, 2011.
- [5] A. Beck and M. Teboulle, "Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems," *IEEE Trans. Image Process.*, vol. 18, no. 11, pp. 2419–2434, 2009.
- [6] H. Li and Z. Lin, "Accelerated proximal gradient methods for nonconvex programming," in *Proc. Advances in Neural Inf. Process. Systems*, 2015, pp. 379–387.
- [7] H. Attouch and J. Bolte, "On the convergence of the proximal algorithm for nonsmooth functions involving analytic features," *Mathematical Programming*, vol. 116, no. 1, pp. 5–16, 2009.
- [8] Q. Yao, J. T. Kwok, F. Gao, W. Chen, and T.-Y. Liu, "Efficient inexact proximal gradient algorithm for nonconvex problems," in *Proc. Int. Joint Conf. Artif. Intell.*, 2017, pp. 3308–3314.
- [9] Q. Li, Y. Zhou, Y. Liang, and P. K. Varshney, "Convergence analysis of proximal gradient with momentum for nonconvex optimization," in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 2111–2119.
- [10] T. Moreau and J. Bruna, "Understanding trainable sparse coding via matrix factorization," *arXiv preprint arXiv:1609.00285*, 2016.
- [11] K. Gregor and Y. LeCun, "Learning fast approximations of sparse coding," in *Proc. Int. Conf. Mach. Learn.*, 2010, pp. 399–406.
- [12] R. Liu, G. Zhong, J. Cao, Z. Lin, S. Shan, and Z. Luo, "Learning to diffuse: A new perspective to design pdes for visual analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 12, pp. 2457–2471, 2016.
- [13] Y. Chen and T. Pock, "Trainable nonlinear reaction diffusion: A flexible framework for fast and effective image restoration," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1256–1272, 2017.
- [14] Y. Yang, J. Sun, H. Li, and Z. Xu, "Admm-net: A deep learning approach for compressive sensing mri," in *Proc. Advances in Neural Inf. Process. Systems*, 2016, pp. 10–18.
- [15] S. Wang, S. Fidler, and R. Urtasun, "Proximal deep structured models," in *Proc. Advances in Neural Inf. Process. Systems*, 2016, pp. 865–873.
- [16] S. Diamond, V. Sitzmann, F. Heide, and G. Wetzstein, "Unrolled optimization with deep priors," *arXiv preprint arXiv:1705.08041*, 2017.
- [17] K. Zhang, W. Zuo, S. Gu, and L. Zhang, "Learning deep cnn denoiser prior for image restoration," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 3929–3938.
- [18] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Deep image prior," *arXiv preprint arXiv:1711.10925*, 2017.
- [19] K. Li and J. Malik, "Learning to optimize," *arXiv preprint arXiv:1606.01885*, 2016.
- [20] M. Andrychowicz, M. Denil, S. Gomez, M. W. Hoffman, D. Pfau, T. Schaul, and N. de Freitas, "Learning to learn by gradient descent by gradient descent," in *Proc. Advances in Neural Inf. Process. Systems*, 2016, pp. 3981–3989.
- [21] B. Gu, Z. Huo, and H. Huang, "Inexact proximal gradient methods for non-convex and non-smooth optimization," *arXiv preprint arXiv:1612.06003*, 2016.
- [22] M. Schmidt, N. L. Roux, and F. R. Bach, "Convergence rates of inexact proximal-gradient methods for convex optimization," in *Proc. Advances in Neural Inf. Process. Systems*, 2011, pp. 1458–1466.
- [23] A. Bronstein, P. Sprechmann, and G. Sapiro, "Learning efficient structured sparse models," *arXiv preprint arXiv:1206.4649*, 2012.
- [24] U. Schmidt and S. Roth, "Shrinkage fields for effective image restoration," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2014, pp. 2774–2781.
- [25] S. H. Chan, X. Wang, and O. A. Elgendy, "Plug-and-play admm for image restoration: Fixed-point convergence and applications," *IEEE Trans. Comput. Imaging*, vol. 3, no. 1, pp. 84–98, 2017.
- [26] J. Bolte, S. Sabach, and M. Teboulle, "Proximal alternating linearized minimization for nonconvex and nonsmooth problems," *Mathematical Programming*, vol. 146, no. 1–2, pp. 459–494, 2014.
- [27] E. Chouzenoux, J.-C. Pesquet, and A. Repetti, "A block coordinate variable metric forward-backward algorithm," *Journal of Global Optimization*, vol. 66, no. 3, pp. 457–485, 2016.
- [28] H. Attouch, J. Bolte, and B. F. Svaiter, "Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forwardbackward splitting, and regularized gaussseidel methods," *Mathematical Programming*, vol. 137, pp. 91–129, 2013.
- [29] I. Daubechies, M. Defrise, and C. De Mol, "An iterative thresholding algorithm for linear inverse problems with a sparsity constraint," *Communications on Pure and Applied Mathematics*, vol. 57, no. 11, pp. 1413–1457, 2004.
- [30] L.-Y. Gui, Y.-X. Wang, D. Ramanan, and J. M. Moura, "Few-shot human motion prediction via meta-learning," in *Proc. Eur. Conf. Comput. Vis.*, 2018.
- [31] Y.-X. Wang and M. Hebert, "Learning from small sample sets by combining unsupervised meta-training with cnns," in *Proc. Advances in Neural Inf. Process. Systems*, 2016, pp. 244–252.
- [32] C. Zhang, Y. Yu, and Z.-H. Zhou, "Learning environmental calibration actions for policy self-evolution," in *Proc. Int. Joint Conf. Artif. Intell.*, 2018, pp. 3061–3067.
- [33] S. Thrun and L. Pratt, *Learning to learn*. Springer Science & Business Media, 2012.
- [34] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra, "Efficient projections onto the ℓ_1 -ball for learning in high dimensions," in *Proc. Int. Conf. Mach. Learn.*, 2008, pp. 272–279.
- [35] S. Cho, J. Wang, and S. Lee, "Handling outliers in non-blind image deconvolution," in *Proc. IEEE Conf. Int. Conf. Comput. Vis.*, 2011, pp. 495–502.
- [36] W. S. Lai, J. B. Huang, Z. Hu, N. Ahuja, and M. H. Yang, "A comparative study for single image blind deblurring," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 1701–1709.
- [37] M. Unser, A. Aldroubi, and M. Eden, "Recursive regularization filters: design, properties, and applications," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 13, no. 3, pp. 272–277, 1991.
- [38] Y. Wang, J. Yang, W. Yin, and Y. Zhang, "A new alternating minimization algorithm for total variation image reconstruction," *SIAM Journal on Imaging Sciences*, vol. 1, no. 3, pp. 248–272, 2008.
- [39] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [40] J. Deng, W. Dong, R. Socher, L. J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009, pp. 248–255.
- [41] Y. Wang, J. Yang, W. Yin, and Y. Zhang, "A new alternating minimization algorithm for total variation image reconstruction," *SIAM Journal on Imaging Sciences*, vol. 1, no. 3, pp. 248–272, 2008.
- [42] A. Danielyan, V. Katkovnik, and K. Egiazarian, "Bm3d frames and variational image deblurring," *IEEE Trans. Image Process.*, vol. 21, no. 4, pp. 1715–1728, 2012.
- [43] D. Zoran and Y. Weiss, "From learning models of natural image patches to whole image restoration," in *Proc. IEEE Conf. Int. Conf. Comput. Vis.*, 2011, pp. 479–486.

- [44] U. Schmidt, J. Jancsary, S. Nowozin, S. Roth, and C. Rother, "Cascades of regression tree fields for image restoration," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 4, pp. 677–689, 2016.
- [45] A. Levin, Y. Weiss, F. Durand, and W. T. Freeman, "Understanding and evaluating blind deconvolution algorithms," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009, pp. 1964–1971.
- [46] L. Sun, S. Cho, J. Wang, and J. Hays, "Edge-based blur kernel estimation using patch priors," in *Proc. IEEE Int. Conf. Comput. Photography*, 2013, pp. 1–8.
- [47] D. Perrone and P. Favaro, "Total variation blind deconvolution: The devil is in the details," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2014, pp. 2909–2916.
- [48] H. Zhang, D. Wipf, and Y. Zhang, "Multi-image blind deblurring using a coupled adaptive sparse prior," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2013, pp. 1051–1058.
- [49] J. Pan, D. Sun, H. Pfister, and M.-H. Yang, "Deblurring images via dark channel prior," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 1, no. 99, pp. 1–1, 2017.
- [50] Y. Li, R. T. Tan, X. Guo, J. Lu, and M. S. Brown, "Rain streak removal using layer priors," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 2736–2744.
- [51] L.-W. Kang, C.-W. Lin, and Y.-H. Fu, "Automatic single-image-based rain streaks removal via image decomposition," *IEEE Trans. Image Process.*, vol. 21, no. 4, pp. 1742–1755, Apr. 2012.
- [52] Y. Luo, Y. Xu, and H. Ji, "Removing rain from a single image via discriminative sparse coding," in *Proc. IEEE Conf. Int. Conf. Comput. Vis.*, Dec. 2015, pp. 3397–3405.
- [53] X. Fu, J. Huang, X. Ding, Y. Liao, and J. Paisley, "Clearing the skies: A deep network architecture for single-image rain removal," *IEEE Trans. Image Process.*, vol. 26, no. 6, pp. 2944–2956, Apr. 2017.
- [54] X. Fu, J. Huang, Y. Huang, Delu Zeng, X. Ding, and J. Paisley, "Removing rain from single images via a deep detail network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jul. 2017, pp. 1715–1723.
- [55] L.-J. Deng, T.-Z. Huang, X.-L. Zhao, and T.-X. Jiang, "A directional global sparse model for single image rain removal," *Applied Mathematical Modelling*, vol. 59, pp. 662–679, 2018.
- [56] R. T. Rockafellar and R. J.-B. Wets, *Variational analysis*. Springer Science & Business Media, 2009, vol. 317.



Risheng Liu (M'12-) received the BSc and PhD degrees both in mathematics from the Dalian University of Technology in 2007 and 2012, respectively. He was a visiting scholar in the Robotics Institute of Carnegie Mellon University from 2010 to 2012. He served as Hong Kong Scholar Research Fellow at the Hong Kong Polytechnic University from 2016 to 2017. He is currently an associate professor with the Key Laboratory for Ubiquitous Network and Service Software of Liaoning Province, Internal School of Information

and Software Technology, Dalian University of Technology. His research interests include machine learning, optimization, computer vision and multimedia. He was a co-recipient of the IEEE ICME Best Student Paper Award in both 2014 and 2015. Two papers were also selected as Finalist of the Best Paper Award in ICME 2017. He is a member of the IEEE and ACM.



Shichao Cheng received the B.E. degree in Mathematics and Applied Mathematics from Henan Normal University, Xinxiang, China, in 2013. She is currently pursuing the PhD degree in Computational Mathematics at Dalian University of Technology, Dalian, China. Her research interests include computer vision, machine learning and optimization.



Yi He received the B.E. degree in Computer Science from Shihezi University, Shihezi, China, in 2017. He is currently pursuing the master degree in software engineering at Dalian University of Technology, Dalian, China. His research interests include computer vision, deep learning.



DTI-MR image analysis.

Xin Fan was born in 1977. He received the B.E. and Ph.D. degrees in information and communication engineering from Xian Jiaotong University, Xian, China, in 1998 and 2004, respectively. He was with Oklahoma State University, Stillwater, from 2006 to 2007, as a post-doctoral research Fellow. He joined the School of Software, Dalian University of Technology, Dalian, China, in 2009. His current research interests include computational geometry and machine learning, and their applications to lowlevel image processing and



member of AAAI 2016/2017/2018 and IJCAI 2016/2018. He is an Associate Editor of the IEEE Transactions on Pattern Analysis and Machine Intelligence and the International Journal of Computer Vision. He is an IAPR/IEEE fellow.

Zhouchen LIN received the Ph.D. degree in applied mathematics from Peking University in 2000. He is currently a Professor with the Key Laboratory of Machine Perception, School of Electronics Engineering and Computer Science, Peking University. His research interests include computer vision, image processing, machine learning, pattern recognition, and numerical optimization. He is an area chair of ACCV 2009/2018, CVPR 2014/2016, ICCV 2015, NIPS 2015/2018 and AAAI 2019, and senior program committee



Zhongxuan Luo received the B.S. degree in Computational Mathematics from Jilin University, China, in 1985, the M.S. degree in Computational Mathematics from Jilin University in 1988, and the PhD degree in Computational Mathematics from Dalian University of Technology, China, in 1991. He has been a full professor of the School of Mathematical Sciences at Dalian University of Technology since 1997. His research interests include computational geometry and computer vision.