



# Stochastic Primal–Dual Hybrid Gradient Algorithm with Adaptive Step Sizes

Antonin Chambolle<sup>1,2</sup> · Claire Delplancke<sup>3</sup> · Matthias J. Ehrhardt<sup>4</sup> · Carola-Bibiane Schönlieb<sup>5</sup> · Junqi Tang<sup>6</sup>

Received: 26 May 2023 / Accepted: 23 January 2024 / Published online: 16 March 2024  
© The Author(s) 2024

## Abstract

In this work, we propose a new primal–dual algorithm with adaptive step sizes. The stochastic primal–dual hybrid gradient (SPDHG) algorithm with constant step sizes has become widely applied in large-scale convex optimization across many scientific fields due to its scalability. While the product of the primal and dual step sizes is subject to an upper-bound in order to ensure convergence, the selection of the ratio of the step sizes is critical in applications. Up-to-now there is no systematic and successful way of selecting the primal and dual step sizes for SPDHG. In this work, we propose a general class of adaptive SPDHG (A-SPDHG) algorithms and prove their convergence under weak assumptions. We also propose concrete parameters-updating strategies which satisfy the assumptions of our theory and thereby lead to convergent algorithms. Numerical examples on computed tomography demonstrate the effectiveness of the proposed schemes.

## 1 Introduction

The stochastic primal–dual hybrid gradient (SPDHG) algorithm introduced in [8] is a stochastic version of the primal–dual hybrid gradient (PDHG) algorithm, also known as Chambolle–Pock algorithm [9]. SPDHG has proved more efficient than PDHG for a variety of problems in the framework of large-scale non-smooth convex inverse problems [13, 22, 24, 27]. Indeed, SPDHG only uses a subset of the

data at each iteration, hence reducing the computational cost of evaluating the forward operator and its adjoint; as a result, for the same computational burden, SPDHG attains convergence faster than PDHG. This is especially relevant in the context of medical imaging, where there is a need for algorithms whose convergence speed is compatible with clinical standards, and at the same time able to deal with convex, non-smooth priors like total variation (TV), which are well-suited to ill-posed imaging inverse problems, but preclude the recourse to scalable gradient-based methods.

Like PDHG, SPDHG is provably convergent under the assumption that the product of its primal and dual step sizes is bounded by a constant depending on the problem to solve. On the other hand, the ratio between the primal and dual step sizes is a free parameter, whose value needs to be chosen by the user. The value of this parameter, which can be interpreted as a control on balance between primal and dual convergence, can have a severe impact on the convergence speed of PDHG, and the same also holds true for SPDHG [12]. This leads to an important challenge in practice, as there is no known theoretical or empirical rule to guide the choice of the parameter. Manual tuning is computationally expensive, as it would require running and comparing the algorithm on a range of values, and there is no guarantee that a value leading to fast convergence for one dataset would keep being a good choice for another dataset. For PDHG, [14] have proposed an online primal–dual balancing strategy to solve the issue, where the values of the step sizes evolve

CD was at the Department of Mathematical Sciences, University of Bath, while the research presented in this article was undertaken.

✉ Claire Delplancke  
claire.delplancke@edf.fr

✉ Matthias J. Ehrhardt  
m.ehrhardt@bath.ac.uk

<sup>1</sup> CEREMADE, Université Paris-Dauphine, Place du Maréchal De Lattre De Tassigny, 75775 Paris, France

<sup>2</sup> MOKAPLAN, INRIA Paris, Paris, France

<sup>3</sup> EDF Lab Paris-Saclay, Route de Saclay, 91300 Palaiseau, France

<sup>4</sup> Department of Mathematical Sciences, University of Bath, Claverton Down, Bath BA2 7AY, UK

<sup>5</sup> Department of Applied Mathematics and Theoretical Physics, University of Cambridge, Wilberforce Road, Cambridge CB3 0WA, UK

<sup>6</sup> School of Mathematics, University of Birmingham, Edgbaston, Birmingham B15 2TT, UK

along the iterations. More generally, adaptive step sizes have been used for PDHG with backtracking in [14, 20], adapting to local smoothness in [25], and are widely used for a variety of other algorithms, namely gradient methods in [19], subgradient methods in [3] and splitting methods in [4–7, 18] to improve convergence speed and bypass the need for explicit model constants, like Lipschitz constants or operator norms. For SPDHG, an empirical adaptive scheme has been used for Magnetic Particle Imaging but without convergence proof [27].

On the theoretical side, a standard procedure to prove the convergence of proximal-based algorithms for convex optimization is to use the notion of F ej er monotonicity [2]. Constant step sizes lead to a fixed metric setting, while adaptive step sizes lead to a variable metric setting. Work [11] states the convergence of deterministic F ej er-monotone sequences in the variable metric setting, while work [10] is concerned by the convergence of random F ej er-monotone sequences in the fixed metric setting.

In this work, we introduce and study an adaptive version of SPDHG. More precisely:

- We introduce a broad class of strategies to adaptively choose the step sizes of SPDHG. This class includes, but is not limited to, the adaptive primal–dual balancing strategy, where the ratio of the step sizes, which controls the balance between convergence of the primal and dual variable, is tuned online.
- We prove the almost-sure convergence of SPDHG under the schemes of the class. In order to do that, we introduce the concept of  $C$ -stability, which generalizes the notion of F ej er monotonicity, and we prove the convergence of random  $C$ -stable sequences in a variable metric setting, hence generalizing results from [11] and [10]. We then show that our proposed algorithm falls within this novel theoretical framework by following similar strategies than in the almost-sure convergence proofs of [1, 16].
- We compare the performance of SPDHG for various adaptive schemes and the known fixed step-size scheme on large-scale imaging inverse tasks (sparse-view CT, limited-angle CT, low-dose CT). We observe that the primal–dual balancing adaptive strategy is always as fast or faster than all the other strategies. In particular, it consistently leads to substantial gains in convergence speed over the fixed strategy if the fixed step sizes, while in the theoretical convergence range, are badly chosen. This is especially relevant as it is impossible to know whether the fixed step sizes are well or badly chosen without running expensive comparative tests. Even in the cases where the SPDHG’s fixed step sizes are well tuned, meaning that they are in the range to which the adaptive step sizes are observed to converge, we observe that our adap-

tive scheme still provides convergence acceleration over the standard SPDHG after a certain number of iterations. Finally, we pay special attention to the hyperparameters used in the adaptive schemes. These hyperparameters are essentially controlling the degree of adaptivity for the algorithm and each of them has a clear interpretation and is easy to choose in practice. We observe in our extensive numerical tests that the convergence speed of our adaptive scheme is robust to the choices of these parameters within the empirical range we provide, hence can be applied directly to the problem at hand without fine-tuning, and solves the step-size choice challenge encountered by the user.

The rest of the paper is organized as follows. In Sect. 2, we introduce SPDHG with adaptive step sizes, state the convergence theorem, and carry the proof. In Sect. 3, we propose concrete schemes to implement the adaptiveness, followed by numerical tests on CT data in Sect. 4. We conclude in Sect. 5. Finally, Sect. 6 collects some useful lemmas and proofs.

## 2 Theory

### 2.1 Convergence Theorem

The variational problem to solve takes the form:

$$\min_{x \in X} \sum_{i=1}^n f_i(A_i x) + g(x),$$

where  $X$  and  $(Y_i)_{i \in \{1, \dots, n\}}$  are Hilbert spaces,  $A_i : X \rightarrow Y_i$  are bounded linear operators, and  $f_i : Y_i \rightarrow \mathbb{R} \cup \{+\infty\}$  and  $g : X \rightarrow \mathbb{R} \cup \{+\infty\}$  are convex functions. We define  $Y = Y_1 \times \dots \times Y_n$  with elements  $y = (y_1, \dots, y_n)$  and  $A : X \rightarrow Y$  such that  $Ax = (A_1 x, \dots, A_n x)$ . The associated saddle-point problem reads as

$$\min_{x \in X} \sup_{y \in Y} \sum_{i=1}^n \langle A_i x, y_i \rangle - f_i^*(y_i) + g(x), \tag{2.1}$$

where  $f_i^*$  stands for the Fenchel conjugate of  $f_i$ . The set of solution to (2.1) is denoted by  $\mathcal{C}$ , and the set of nonnegative integers by  $\mathbb{N}$  and  $\llbracket 1, n \rrbracket$  stands for  $\{1, \dots, n\}$ . Elements  $(x^*, y^*)$  of  $\mathcal{C}$  are called saddle points and characterized by

$$A_i x^* \in \partial f_i^*(y_i^*), \quad i \in \llbracket 1, n \rrbracket; \quad - \sum_{i=1}^n A_i^* y_i^* \in \partial g(x^*). \tag{2.2}$$

In order to solve the saddle-point problem, we introduce the adaptive stochastic primal–dual hybrid gradient

(A-SPDHG) algorithm in Algorithm 2.1. At each iteration  $k \in \mathbb{N}$ , A-SPDHG involves the following five steps:

---

**Algorithm 2.1:** A-SPDHG (variable step-sizes, serial sampling)

---

- 1: Input: dual step-sizes  $(\sigma_i^0)_{i \in \llbracket 1, n \rrbracket}$ , primal step-size  $\tau^0$ , update rule; probabilities  $(p_i)_{i \in \llbracket 1, n \rrbracket}$ ; primal variable  $x^0$ , dual variable  $y^0$
- 2: Initialize  $\bar{y}^0 = y^0$
- 3: **for**  $k \in \llbracket 0, K - 1 \rrbracket$  **do**
- 4: Determine  $(\sigma_i^{k+1})_{i \in \llbracket 1, n \rrbracket}$ ,  $\tau^{k+1}$  according to the update rule and the values of  $(\sigma_i^k)_{i \in \llbracket 1, n \rrbracket}$ ,  $\tau^k$ ,  $x^k$  and  $y^k$  for  $l \in \llbracket 0, k \rrbracket$ .
- 5:  $x^{k+1} = \text{prox}_{\tau^{k+1}g}(x^k - \tau^{k+1}A^*\bar{y}^k)$
- 6: Randomly pick  $i \in \llbracket 1, n \rrbracket$  with probability  $p_i$
- 7:  $y_j^{k+1} = \begin{cases} \text{prox}_{\sigma_i^{k+1}f_i^*}(y_i^k + \sigma_i^{k+1}A_i x^{k+1}) & \text{if } j = i \\ y_j^k & \text{if } j \neq i \end{cases}$
- 8:  $\bar{y}_j^{k+1} = \begin{cases} y_i^{k+1} + \frac{1}{p_i}(y_i^{k+1} - y_i^k) & \text{if } j = i \\ y_j^k & \text{if } j \neq i \end{cases}$
- 9: **end for**
- 10: **return**  $x^K$

---

- update the primal step size  $\tau^k$  and the dual step sizes  $(\sigma_i^k)_{i \in \llbracket 1, n \rrbracket}$  (line 4);
- update the primal variable  $x^k$  by a proximal step with step size  $\tau^{k+1}$  (line 5);
- randomly choose an index  $i$  with probability  $p_i$  (line 6);
- update the dual variable  $y_i^k$  by a proximal step with step size  $\sigma_i^{k+1}$  (line 7);
- compute the extrapolated dual variable (line 8).

A-SPDHG is *adaptive* in the sense that the step-size values are updated at each iteration according to an update rule which takes into account the value of the primal and dual iterates  $x^l$  and  $y^l$  up to the current iteration. As the iterates are stochastic, the step sizes are themselves stochastic, which must be carefully accounted for in the theory.

Before turning to the convergence of A-SPDHG, let us recall some facts about the state-of-the-art SPDHG. Each iteration of SPDHG involves the selection of a random subset of  $\llbracket 1, n \rrbracket$ . In the serial sampling case where the random subset is a singleton, SPDHG algorithm [8] is a special case of Algorithm 2.1 with the update rule

$$\begin{cases} \sigma_i^{k+1} = \sigma_i^k (= \sigma_i), & i \in \llbracket 1, n \rrbracket, \\ \tau^{k+1} = \tau^k (= \tau_i), & \end{cases} \quad k \in \mathbb{N}.$$

Under the condition

$$\tau \sigma_i < \frac{p_i}{\|A_i\|^2}, \quad i \in \llbracket 1, n \rrbracket, \tag{2.3}$$

SPDHG iterates converge almost surely to a solution of the saddle-point problem (2.1) [1, 16].

Let us now turn to the convergence of A-SPDHG. The main theorem, Theorem 2.1, gives conditions on the update rule under which A-SPDHG is provably convergent. Plainly speaking, these conditions are threefold:

- (i) the step sizes for step  $k + 1$ ,  $(\sigma_i^{k+1})_{i \in \llbracket 1, n \rrbracket}$  and  $\tau^{k+1}$ , depend only on the iterates up to step  $k$ ,
- (ii) the step sizes satisfy a uniform version of condition (2.3),
- (iii) the step-size sequences  $(\tau^k)_{k \geq 0}$  and  $(\sigma_i^k)_{k \geq 0}$  for  $i \in \llbracket 1, n \rrbracket$  do not decrease too fast. More precisely, they are uniformly almost surely quasi-increasing in the sense defined below.

In order to state the theorem rigorously, let us introduce some useful notation and definitions. For all  $k \in \mathbb{N}$ , the  $\sigma$ -algebra generated by the iterates up to point  $k$ ,  $\mathcal{F}((x^l, y^l), l \in \llbracket 0, k \rrbracket)$ , is denoted by  $\mathcal{F}^k$ . We say that a sequence  $(u^k)_{k \in \mathbb{N}}$  is  $(\mathcal{F}^k)_{k \in \mathbb{N}}$ -adapted if for all  $k \in \mathbb{N}$ ,  $u^k$  is measurable with respect to  $\mathcal{F}^k$ .

A positive real sequence  $(u^k)_{k \in \mathbb{N}}$  is said to be *quasi-increasing* if there exists a sequence  $(\eta^k)_{k \in \mathbb{N}}$  with values in  $[0, 1)$ , called the control on  $(u^k)_{k \in \mathbb{N}}$ , such that  $\sum_{k=1}^\infty \eta^k < \infty$  and:

$$u^{k+1} \geq (1 - \eta^k)u^k, \quad k \in \mathbb{N}. \tag{2.4}$$

By extension, we call a random positive real sequence  $(u^k)_{k \in \mathbb{N}}$  *uniformly almost surely quasi-increasing* if there exists a deterministic sequence  $(\eta^k)_{k \in \mathbb{N}}$  with values in  $[0, 1)$  such that  $\sum_{k=1}^\infty \eta^k < \infty$  and equation (2.4) above holds almost surely (a.s.).

**Theorem 2.1** (Convergence of A-SPDHG) *Let  $X$  and  $Y$  be separable Hilbert spaces,  $A_i : X \rightarrow Y_i$  bounded linear operators,  $f_i : Y_i \rightarrow \mathbb{R} \cup \{+\infty\}$  and  $g : X \rightarrow \mathbb{R} \cup \{+\infty\}$  proper, convex and lower semi-continuous functions for all  $i \in \llbracket 1, n \rrbracket$ . Assume that the set of saddle points  $\mathcal{C}$  is non-empty and the sampling is proper, that is to say  $p_i > 0$  for all  $i \in \llbracket 1, n \rrbracket$ . If the following conditions are met:*

- (i) the step-size sequences  $(\tau^{k+1})_{k \in \mathbb{N}}$ ,  $(\sigma_i^{k+1})_{k \in \mathbb{N}}$ ,  $i \in \llbracket 1, n \rrbracket$  are  $(\mathcal{F}^k)_{k \in \mathbb{N}}$ -adapted,
- (ii) there exists  $\beta \in (0, 1)$  such that for all indices  $i \in \llbracket 1, n \rrbracket$  and iterates  $k \in \mathbb{N}$ ,

$$\tau^k \sigma_i^k \frac{\|A_i\|^2}{p_i} \leq \beta < 1, \tag{2.5}$$

(iii) the initial step sizes  $\tau^0$  and  $\sigma_i^0$  for all indices  $i \in \llbracket 1, n \rrbracket$  are positive and the step-size sequences  $(\tau^k)_{k \in \mathbb{N}}$  and  $(\sigma_i^k)_{k \in \mathbb{N}}$  for all indices  $i \in \llbracket 1, n \rrbracket$  are uniformly almost surely quasi-increasing,

then the sequence of iterates  $(x^k, y^k)_{k \in \mathbb{N}}$  converges almost surely to an element of  $\mathcal{C}$ .

While the conditions (i)–(iii) are general enough to cover a large range of step-size update rules, we will focus in practice on the primal–dual balancing strategy, which consists in scaling the primal and the dual step sizes by an inverse factor at each iteration. In that case, the update rule depends on a random positive sequence  $(\gamma^k)_{k \in \mathbb{N}}$  and reads as:

$$\tau^{k+1} = \frac{\tau^k}{\gamma^k}, \quad \sigma_i^{k+1} = \gamma^k \sigma_i^k, \quad i \in \llbracket 1, n \rrbracket. \tag{2.6}$$

**Lemma 2.2** (Primal–dual balancing) *Let the step-size sequences satisfy equation (2.6) and assume in addition that  $(\gamma^k)_{k \in \mathbb{N}}$  is  $(\mathcal{F}^k)_{k \in \mathbb{N}}$ -adapted that the initial step sizes satisfy*

$$\tau^0 \sigma_i^0 \frac{\|A_i\|^2}{p_i} < 1, \quad i \in \llbracket 1, n \rrbracket,$$

and are positive, that there exists a deterministic sequence  $(\epsilon^k)_{k \in \mathbb{N}}$  with values in  $[0, 1)$  such that  $\sum \epsilon^k < \infty$  and for all  $k \in \mathbb{N}$  and  $i \in \llbracket 1, n \rrbracket$ ,

$$\min \left\{ \gamma^k, (\gamma^k)^{-1} \right\} \geq 1 - \epsilon^k. \tag{2.7}$$

Then, the step-size sequences satisfy assumptions (i)–(iii) of Theorem 2.1.

Lemma 2.2 is proved in Sect. 6.

Connection with the literature:

- The primal–dual balancing strategy has been introduced in [14] for PDHG and indeed for  $n = 1$  we recover with Lemma 2.2 the non-backtracking algorithm presented in [14]. As a consequence, our theorem also implies the pointwise convergence of this algorithm, whose convergence was established in the sense of vanishing residuals in [14].
- Still for PDHG, [20] proposes without proof an update rule where the ratio of the step sizes is either quasi-non-increasing or quasi-non-decreasing. This requirement is similar to but not directly connected with ours, where we ask the step sizes themselves to be quasi-non-increasing.
- For SPDHG, the angular constraint step-size rule proposed without convergence proof in [27] satisfies assumptions (i)–(iii).

*Outline of the proof:* Theorem 2.1 is proved in the following subsections. We first define in Sect. 2.2 metrics related to the algorithm step sizes on the primal–dual product space. As the step sizes are adaptive, we obtain a sequence of metrics. The proof of Theorem 2.1 is then similar in strategy to those of [1] and [16] but requires novel elements to deal with the metrics variability. In Theorem 2.5, we state convergence conditions for an abstract random sequence in a Hilbert space equipped with random variable metrics. In Sects. 2.4 and 2.5, we show that A-SPDHG falls within the scope of Theorem 2.5. We collect all elements and conclude the proof in Sect. 2.6.

### 2.2 Variable Metrics

For a Hilbert space  $H$ , we call  $\mathbb{S}(H)$  the set of bounded self-adjoint linear operators from  $H$  to  $H$ , and for all  $M \in \mathbb{S}(H)$  we introduce the notation:

$$\|u\|_M^2 = \langle Mu, u \rangle, \quad u \in H.$$

By an abuse of notation, we write  $\|\cdot\|_\alpha^2 = \|\cdot\|_{\alpha \text{Id}}^2$  for a scalar  $\alpha \in \mathbb{R}$ . Notice that  $\|\cdot\|_M$  is a norm on  $H$  if  $M$  is positive definite. Furthermore, we introduce the partial order  $\preceq$  on  $\mathbb{S}(H)$  such that for  $M, N \in \mathbb{S}(H)$ ,

$$N \preceq M \quad \text{if} \quad \forall u \in H, \|u\|_N \leq \|u\|_M.$$

We call  $\mathbb{S}_\alpha(H)$  the subset of  $\mathbb{S}(H)$  comprised of  $M$  such that  $\alpha \text{Id} \preceq M$ . Furthermore, a random sequence  $(M^k)_{k \in \mathbb{N}}$  in  $\mathbb{S}(H)$  is said to be *uniformly almost surely quasi-decreasing* if there exists a deterministic nonnegative sequence  $(\eta^k)_{k \in \mathbb{N}}$  such that  $\sum_{k=1}^\infty \eta^k < \infty$  and a.s.

$$M^{k+1} \preceq (1 + \eta^k)M^k, \quad k \in \mathbb{N}.$$

Coming back to A-SPDHG, let us define for every iteration  $k \in \mathbb{N}$  and every index  $i \in \llbracket 1, n \rrbracket$  two block operators of  $\mathbb{S}(X \times Y_i)$  as:

$$M_i^k = \begin{pmatrix} \frac{1}{\tau^k} \text{Id} & -\frac{1}{p_i} A_i^* \\ -\frac{1}{p_i} A_i & \frac{1}{p_i \sigma_i^k} \text{Id} \end{pmatrix}, \quad N_i^k = \begin{pmatrix} \frac{1}{\tau^k} \text{Id} & 0 \\ 0 & \frac{1}{p_i \sigma_i^k} \text{Id} \end{pmatrix},$$

and a block operator of  $\mathbb{S}(X \times Y)$  as:

$$N^k = \begin{pmatrix} \frac{1}{\tau^k} \text{Id} & & & & (0) \\ & \frac{1}{p_1 \sigma_1^k} \text{Id} & & & \\ & & \ddots & & \\ & & & \frac{1}{p_i \sigma_i^k} \text{Id} & \\ (0) & & & & \frac{1}{p_n \sigma_n^k} \text{Id} \end{pmatrix}. \tag{2.8}$$

The following lemma translates assumptions (i)–(iii) of Theorem 2.1 on properties on the variable metric sequences.

**Lemma 2.3** (Variable metric properties)

- (a) Assumption (i) of Theorem 2.1 implies that  $(M_i^{k+1})_{k \in \mathbb{N}}$ ,  $(N_i^{k+1})_{k \in \mathbb{N}}$ ,  $i \in \llbracket 1, n \rrbracket$  and  $(N^{k+1})_{k \in \mathbb{N}}$  are  $(\mathcal{F}^k)_{k \in \mathbb{N}}$ -adapted.
- (b) Assumption (ii) of Theorem 2.1 is equivalent to the existence of  $\beta \in (0, 1)$  such that for all indices  $i \in \llbracket 1, n \rrbracket$  and iterates  $k \in \mathbb{N}$ ,

$$(1 - \sqrt{\beta})N_i^k \preceq M_i^k.$$

- (c) Assumptions (ii) and (iii) of Theorem 2.1 imply that  $(M_i^k)_{k \in \mathbb{N}}$ ,  $(N_i^k)_{k \in \mathbb{N}}$ ,  $i \in \llbracket 1, n \rrbracket$  and  $(N^k)_{k \in \mathbb{N}}$  are uniformly a.s. quasi-decreasing.
- (d) Assumption (ii) and (iii) of Theorem 2.1 imply that the sequences  $(\tau^k)_{k \in \mathbb{N}}$  and  $(\sigma_i^k)_{k \in \mathbb{N}}$  for all  $i \in \llbracket 1, n \rrbracket$  are a.s. bounded from above and by below by positive constants. In particular, this implies that there exists  $\alpha > 0$  such that  $N_i^k \in \mathbb{S}_\alpha(X \times Y_i)$  for all  $i \in \llbracket 1, n \rrbracket$  and  $k \in \mathbb{N}$ , or equivalently that  $N^k \in \mathbb{S}_\alpha(X \times Y)$  for all  $k \in \mathbb{N}$ .

**Remark 2.4** (Step-size induced metrics on the primal–dual product space) The lemma implies that  $M_i^k$ ,  $N_i^k$  and  $N^k$  are positive definite and hence induce a metric on the corresponding spaces. If  $n = 1$  and for constant step sizes,  $M_i^k$  corresponds to the metric used in [17], where PDHG is reformulated as a proximal-point algorithm for a non-trivial metric on the primal–dual product space.

**Proof of Lemma 2.3** Assertion (a) of the lemma follows from the fact that for all iterate  $k \in \mathbb{N}$ , the operators  $M_i^{k+1}$ ,  $N_i^{k+1}$  and  $N^{k+1}$  are in the  $\sigma$ -algebra generated by  $\{\tau^{k+1}, \sigma_i^{k+1}, i \in \llbracket 1, n \rrbracket\}$ . Assertion (b) follows from equation (6.2) of Lemma 6.1 to be found in the complementary material. The proof of assertion (c) is a bit more involved. Let us assume that assumption (iii) of Theorem 2.1 holds and let  $(\eta_0^k)_{k \in \mathbb{N}}$  and  $(\eta_i^k)_{k \in \mathbb{N}}$  be the controls of  $(\tau^k)_{k \in \mathbb{N}}$  and  $(\sigma_i^k)_{k \in \mathbb{N}}$  for  $i \in \llbracket 1, n \rrbracket$ , respectively. We define the sequence  $(\eta^k)_{k \in \mathbb{N}}$

by:

$$\eta^k = \max \left\{ \eta_i^k, i \in \llbracket 0, n \rrbracket \right\}, \quad k \in \mathbb{N}, \tag{2.9}$$

which is a common control on  $(\tau^k)_{k \in \mathbb{N}}$  and  $(\sigma_i^k)_{k \in \mathbb{N}}$  for  $i \in \llbracket 1, n \rrbracket$  as the maximum of a finite number of controls. Let us fix  $k \in \mathbb{N}$  and  $i \in \llbracket 1, n \rrbracket$ . Because the intersection of a finite number of measurable events of probability one is again a measurable event of probability one, it holds almost surely that for all  $(x, y_i) \in X \times Y_i$ ,

$$\begin{aligned} \|(x, y_i)\|_{N_i^{k+1}}^2 &= \frac{1}{\tau^{k+1}} \|x\|^2 + \frac{1}{p_i \sigma_i^{k+1}} \|y_i\|^2 \\ &\leq \frac{1}{1 - \eta^k} \left( \frac{1}{\tau^k} \|x\|^2 + \frac{1}{p_i \sigma_i^k} \|y_i\|^2 \right) \\ &= \left( 1 + \frac{\eta^k}{1 - \eta^k} \right) \|(x, y_i)\|_{N_i^k}^2. \end{aligned}$$

Hence, the sequence  $(N_i^k)_{k \in \mathbb{N}}$  is uniformly quasi-decreasing with control  $(\eta^k(1 - \eta^k)^{-1})_{k \in \mathbb{N}}$ , which is indeed a positive sequence with bounded sum. (To see that  $(\eta^k(1 - \eta^k)^{-1})_{k \in \mathbb{N}}$  has a bounded sum, consider that  $(\eta^k)_{k \in \mathbb{N}}$  is summable, hence converges to 0, hence is smaller than 1/2 for all integers  $k$  bigger than a certain  $K$ ; in turn, for all integers  $k$  bigger than  $K$ , the term  $\eta^k(1 - \eta^k)^{-1}$  is bounded from below by 0 and from above by  $2\eta^k$ , hence is summable.) One can see by a similar proof that  $(N^k)_{k \in \mathbb{N}}$  is uniformly quasi-decreasing with the same control. To follow with the case of  $(M_i^k)_{k \in \mathbb{N}}$ , we have, as before:

$$\begin{aligned} M_i^{k+1} &= \begin{pmatrix} \frac{1}{\tau^{k+1}} \text{Id} & -\frac{1}{p_i} A_i^* \\ -\frac{1}{p_i} A_i & \frac{1}{p_i \sigma_i^{k+1}} \text{Id} \end{pmatrix} \preceq M_i^k + \frac{\eta^k}{1 - \eta^k} N_i^k \\ &\preceq \left( 1 + \frac{\eta^k}{1 - \eta^k} \frac{1}{1 - \sqrt{\beta}} \right) M_i^k \end{aligned}$$

thanks to (b).

Let us conclude with the proof of assertion (d). By assumption (iii), the sequences  $(\tau^k)_{k \in \mathbb{N}}$  and  $(\sigma_i^k)_{k \in \mathbb{N}}$  are uniformly a.s. quasi-increasing. We define a common control  $(\eta^k)_{k \in \mathbb{N}}$  as in (2.9). Then, the sequences  $(\tau^k)_{k \in \mathbb{N}}$  and  $(\sigma_i^k)_{k \in \mathbb{N}}$  are a.s. bounded from below by the same deterministic constant  $C = \min \{ \tau_0^0, \sigma_i^0, i \in \llbracket 1, n \rrbracket \} \prod_{j=0}^\infty (1 - \eta^j)$  which is positive as the initial step sizes are positive and  $(\eta^k)_{k \in \mathbb{N}}$  takes values in  $[0, 1)$  and has finite sum. Furthermore, by assumption (ii), the product of the sequences  $(\tau^k)_{k \in \mathbb{N}}$  and  $(\sigma_i^k)_{k \in \mathbb{N}}$  is almost surely bounded from above. As a consequence, each sequence  $(\tau^k)_{k \in \mathbb{N}}$  and  $(\sigma_i^k)_{k \in \mathbb{N}}$  is a.s. bounded from above. The equivalence with  $N_i^k \in \mathbb{S}_\alpha(X \times Y_i)$  for all  $i \in \llbracket 1, n \rrbracket$ , and with  $N^k \in \mathbb{S}_\alpha(X \times Y)$ , is straightforward.  $\square$

### 2.3 Convergence of Random C-stable Sequences in Random Variable Metrics

Let  $H$  be a Hilbert space and  $C \subset H$  a subset of  $H$ . Let  $(\Omega, \sigma(\Omega), \mathbb{P})$  be a probability space. All random variables in the following are assumed to be defined on  $\Omega$  and measurable with respect to  $\sigma(\Omega)$  unless stated otherwise. Let  $(Q^k)_{k \in \mathbb{N}}$  be a random sequence of  $\mathbb{S}(H)$ .

A random sequence  $(u^k)_{k \in \mathbb{N}}$  with values in  $H$  is said to be *stable with respect to the target  $C$  relative to  $(Q^k)_{k \in \mathbb{N}}$*  if for all  $u \in C$ , the sequence  $(\|u^k - u\|_{Q^k})_{k \in \mathbb{N}}$  converges almost surely. The following theorem then states sufficient conditions for the convergence of such sequences.

**Theorem 2.5** (Convergence of  $C$ -stable sequences) *Let  $H$  be a separable Hilbert space,  $C$  a closed non-empty subset of  $H$ ,  $(Q^k)_{k \in \mathbb{N}}$  a random sequence of  $\mathbb{S}(H)$ , and  $(u^k)_{k \in \mathbb{N}}$  a random sequence of  $H$ . If the following conditions are met:*

- (i)  $(Q^k)_{k \in \mathbb{N}}$  takes values in  $\mathbb{S}_\alpha(H)$  for a given  $\alpha > 0$  and is uniformly a.s. quasi-decreasing,
- (ii)  $(u^k)_{k \in \mathbb{N}}$  is stable with respect to the target  $C$  relative to  $(Q^k)_{k \in \mathbb{N}}$ ,
- (iii) every weak sequential cluster point of  $(u^k)_{k \in \mathbb{N}}$  is almost surely in  $C$ , meaning that there exists  $\Omega_{(iii)}$  a measurable subset of  $\Omega$  of probability one such that for all  $\omega \in \Omega$ , every weak sequential cluster point of  $(u^k(\omega))_{k \in \mathbb{N}}$  is in  $C$ .

then  $(u^k)_{k \in \mathbb{N}}$  converges almost surely weakly to a random variable in  $C$ .

Stability with respect to a target set  $C$  is implied by Féjer and quasi-Féjer monotonicity with respect to  $C$ , which have been studied either for random sequences [10] or in the framework of variable metrics [11], but to the best of our knowledge not both at the same time. The proof of Theorem 2.5 follows the same lines than [10, Proposition 2.3 (iii)] and uses two results from [11].

**Proof** The set  $C$  is a subset of the separable Hilbert space  $H$ , hence is separable. As  $C$  is a closed and separable, there exists  $\{c^n, n \in \mathbb{N}\}$  a countable subset of  $C$  whose closure is equal to  $C$ . Thanks to assumption (ii), there exists for all  $n \in \mathbb{N}$  a measurable subset  $\Omega_{(ii)}^n$  of  $\Omega$  with probability one such that the sequence  $(\|u^k(\omega) - c^n\|_{Q^k(\omega)})_{k \in \mathbb{N}}$  converges for all  $\omega \in \Omega_{(ii)}^n$ . Furthermore, let  $\Omega_{(i)}$  be a measurable subset of  $\Omega$  of probability one corresponding to the almost-sure property for assumption (i). Let

$$\tilde{\Omega} = \left( \bigcap_{n \geq 0} \Omega_{(ii)}^n \right) \cap \Omega_{(i)} \cap \Omega_{(iii)}.$$

As the intersection of a countable number of measurable subsets of probability one,  $\tilde{\Omega}$  is itself a measurable set of  $\Omega$  with  $\mathbb{P}(\tilde{\Omega}) = 1$ . Fix  $\omega \in \tilde{\Omega}$  for the rest of the proof.

The sequence  $(Q^k(\omega))_{k \in \mathbb{N}}$  takes values in  $\mathbb{S}_\alpha(H)$  for  $\alpha > 0$  and is quasi-decreasing with control  $(\eta^k(\omega))_{k \in \mathbb{N}}$ . Furthermore, for all  $k \in \mathbb{N}$ ,

$$\begin{aligned} \|Q^k(\omega)\| &\leq \left( \prod_{j=0}^{k-1} (1 + \eta^j) \right) \|Q^0(\omega)\| \\ &\leq \left( \prod_{j=0}^{\infty} (1 + \eta^j) \right) \|Q^0(\omega)\|, \end{aligned}$$

where the product  $\prod_{j=0}^{\infty} (1 + \eta^j)$  is finite because  $(\eta^k)_{k \in \mathbb{N}}$  is positive and summable. By [11, Lemma 2.3],  $(Q^k(\omega))_{k \in \mathbb{N}}$  converges pointwise strongly to some  $Q(\omega) \in \mathbb{S}_\alpha(H)$ .

Furthermore, for all  $x \in C$ , there exists a sequence  $(x^n)_{n \in \mathbb{N}}$  with values in  $\{c^n, n \in \mathbb{N}\}$  converging strongly to  $x$ . By assumption, for all  $n \in \mathbb{N}$ , the sequence  $(\|u^k(\omega) - x^n\|_{Q^k(\omega)})_{k \in \mathbb{N}}$  converges to a limit which shall be called  $l^n(\omega)$ . For all  $n \in \mathbb{N}$  and  $k \in \mathbb{N}$ , we can write thanks to the triangular inequality:

$$\begin{aligned} -\|x^n - x\|_{Q^k(\omega)} &\leq \|u^k(\omega) - x\|_{Q^k(\omega)} - \|u^k(\omega) - x^n\|_{Q^k(\omega)} \\ &\leq \|x^n - x\|_{Q^k(\omega)}. \end{aligned}$$

By taking the limit  $k \rightarrow +\infty$ , it follows that:

$$\begin{aligned} -\|x^n - x\|_{Q(\omega)} &\leq \liminf_{k \rightarrow \infty} \|u^k(\omega) - x\|_{Q^k(\omega)} - l^n(\omega) \\ &\leq \limsup_{k \rightarrow \infty} \|u^k(\omega) - x\|_{Q^k(\omega)} - l^n(\omega) \\ &\leq \|x^n - x\|_{Q(\omega)}. \end{aligned}$$

Taking now the limit  $n \rightarrow +\infty$  shows that the sequence  $(\|u^k(\omega) - x\|_{Q^k(\omega)})_{k \in \mathbb{N}}$  converges for all  $x \in C$ . On the other hand, because  $\omega \in \Omega_{(iii)}$ , the weak cluster points of  $(u^k(\omega))_{k \in \mathbb{N}}$  lie in  $C$ . Hence, by [11, Theorem 3.3], the sequence  $(u^k(\omega))_{k \in \mathbb{N}}$  converges almost surely to a point  $u(\omega) \in C$ . □

We are now equipped to prove Theorem 2.1. We show in Sects. 2.4 and 2.5 that A-SPDHG satisfies points (ii) and (iii) of Theorem 2.5, respectively, and conclude the proof in Sect. 2.6. Interestingly, the proofs of point (ii) and of point (iii) rely on two different ways of apprehending A-SPDHG. Point (ii) relies on a convex optimization argument: By taking advantage of the measurability of the primal variable at step  $k + 1$  with respect to  $\mathcal{F}^k$ , one can write a contraction-type inequality relating the conditional expectation of the

iterates' norm at step  $k + 1$  to the iterates' norm at step  $k$ . Point (iii) relies on monotone operator theory: We use the fact that the update from the half-shifted iterations  $(y^k, x^{k+1})$  to  $(y^{k+1}, x^{k+2})$  can be interpreted as a step of a proximal-point algorithm on  $X \times Y_i$  conditionally to  $i$  being the index randomly selected at step  $k$ .

### 2.4 A-SPDHG is Stable with Respect to the Set of Saddle Points

In this section, we show that  $(x^k, y^k)_{k \in \mathbb{N}}$  is stable with respect to  $\mathcal{C}$  relative to the variable metrics sequence  $(N^k)_{k \in \mathbb{N}}$  defined in equation (2.8) above. We introduce the operators  $P \in \mathbb{S}(Y)$  and  $\Sigma^k \in \mathbb{S}(Y)$  defined, respectively, by

$$(Py)_i = p_i y_i, \quad (\Sigma^k y)_i = \sigma_i^k y_i, \quad i \in \llbracket 1, n \rrbracket,$$

and the functionals  $(U^k)_{k \in \mathbb{N}}, (V^k)_{k \in \mathbb{N}}$  defined for all  $(x, y) \in X \times Y$  as:

$$U^k(y) = \|y\|_{(P\Sigma^k)^{-1}}^2, \\ V^k(x, y) = \|x\|_{(\tau^k)^{-1}}^2 - 2\langle P^{-1}Ax, y \rangle + \|y\|_{(P\Sigma^k)^{-1}}^2.$$

We begin by recalling the cornerstone inequality satisfied by the iterates of SPDHG stated first in [8] and reformulated in [1].

**Lemma 2.6** ([1], Lemma 4.1) *For every saddle-point  $(x^*, y^*)$ , it a.s. stands that for all  $k \in \mathbb{N} \setminus \{0\}$ ,*

$$\mathbb{E} \left[ V^{k+1}(x^{k+1} - x^*, y^{k+1} - y^k) + U^{k+1}(y^{k+1} - y^*) | \mathcal{F}^k \right] \\ \leq V^{k+1}(x^k - x^*, y^k - y^{k-1}) + U^{k+1}(y^k - y^*) \\ - V^{k+1}(x^{k+1} - x^k, y^k - y^{k-1}). \tag{2.10}$$

The second step is to relate the assumptions of Theorem 2.1 to properties of the functionals appearing in (2.10). Let us introduce  $Y_{\text{sparse}} \subset Y$  the set of elements  $(y_1, \dots, y_n)$  having at most one non-vanishing component.

**Lemma 2.7** (Properties of functionals of interest) *Under the assumptions of Theorem 2.1, there exists a nonnegative, summable sequence  $(\eta^k)_{k \in \mathbb{N}}$  such that a.s. for every iterate  $k \in \mathbb{N}$  and  $x \in X, y \in Y, z \in Y_{\text{sparse}}$ :*

$$U^{k+1}(y) \leq (1 + \eta^k)U^k(y), \tag{2.11a}$$

$$V^{k+1}(x, z) \leq (1 + \eta^k)V^k(x, z), \tag{2.11b}$$

$$\|(x, z)\|_{N^k}^2 \geq \alpha \|(x, z)\|^2, \tag{2.11c}$$

$$V^k(x, z) \geq (1 - \beta)\|(x, z)\|_{N^k}^2, \tag{2.11d}$$

$$\left| \langle P^{-1}Ax, z \rangle \right| \leq \sqrt{\beta} \|x\|_{(\tau^k)^{-1}} \|z\|_{(P\Sigma^k)^{-1}}. \tag{2.11e}$$

**Proof** Let  $(\eta_i^k)_{k \in \mathbb{N}}$  and  $(\tilde{\eta}_i^k)_{k \in \mathbb{N}}$  be the controls of  $(M_i^k)_{k \in \mathbb{N}}$  and  $(N_i^k)_{k \in \mathbb{N}}$ , respectively, for all  $i \in \llbracket 1, n \rrbracket$ . We define the common control  $(\eta^k)_{k \in \mathbb{N}}$  by:

$$\eta^k = \max \left\{ \max \left\{ \eta_i^k, \tilde{\eta}_i^k \right\}, i \in \llbracket 1, n \rrbracket \right\}, \quad k \in \mathbb{N}. \tag{2.12}$$

For all  $y \in Y$ , we can write

$$U^{k+1}(y) = \sum_{i=1}^n \|(0, y_i)\|_{N_i^{k+1}}^2 \leq (1 + \eta^k) \sum_{i=1}^n \|(0, y_i)\|_{N_i^k}^2 \\ = (1 + \eta^k)U^k(y),$$

which proves (2.11a). Let us now fix  $x \in X, z \in Y_{\text{sparse}}$  and  $k \in \mathbb{N}$ . By definition, there exists  $i \in \llbracket 1, n \rrbracket$  such that  $z_j = 0$  for all  $j \neq i$ . We obtain the inequalities (2.11b)–(2.11d) by writing:

$$V^{k+1}(x, z) = \|(x, z_i)\|_{M_i^{k+1}}^2 \leq (1 + \eta^k)\|(x, z_i)\|_{M_i^k}^2 \\ = (1 + \eta^k)V^k(x, z), \\ \|(x, z)\|_{N^k}^2 = \|(x, z_i)\|_{N_i^k}^2 \geq \alpha \|(x, z_i)\|^2 = \alpha \|(x, z)\|^2, \\ V^k(x, z) = \|(x, z_i)\|_{M_i^k}^2 \geq (1 - \beta)\|(x, z_i)\|_{N_i^k}^2 \\ = (1 - \beta)\|(x, z)\|_{N^k}^2.$$

Finally, we obtain inequality (2.11e) by writing:

$$\left| \langle P^{-1}Ax, z \rangle \right| = \frac{1}{p_i} |\langle A_i x, z_i \rangle| \\ \leq \frac{\|A_i\|}{p_i} \|x\| \|z_i\| \\ = \frac{\|A_i\|}{p_i} * (\tau^k \sigma_i^k p_i)^{1/2} \|x\|_{(\tau^k)^{-1}} \|z\|_{(P\Sigma^k)^{-1}} \\ \leq \sqrt{\beta} \|x\|_{(\tau^k)^{-1}} \|z\|_{(P\Sigma^k)^{-1}},$$

where the last inequality is a consequence of (2.5). □

**Lemma 2.8** (A-SPDHG is  $\mathcal{C}$ -stable) *Under the assumptions of Theorem 2.1,*

- (i) *The sequence  $(x^k, y^k)_{k \in \mathbb{N}}$  of Algorithm 2.1 is stable with respect to  $\mathcal{C}$  relative to  $(N^k)_{k \in \mathbb{N}}$ ,*
- (ii) *the following results hold:*

$$\mathbb{E} \left[ \sum_{k=1}^{\infty} \left\| (x^{k+1} - x^k, y^k - y^{k-1}) \right\|^2 \right] \\ < \infty \quad \text{and a.s.} \quad \|x^{k+1} - x^k\| \rightarrow 0.$$

**Proof** Let us begin with the proof of point (i). By definition of A-SPDHG with serial sampling, the difference between

two consecutive dual iterates is almost surely sparse:

$$\text{a.s. } \forall k \in \mathbb{N} \setminus \{0\}, y^k - y^{k-1} \in Y_{\text{sparse}}.$$

Let us define the sequences

$$\begin{aligned} a^k &= V^k(x^k - x^*, y^k - y^{k-1}) + U^k(y^k - y^*), \quad b^k \\ &= V^{k+1}(x^{k+1} - x^k, y^k - y^{k-1}), \end{aligned}$$

which are a.s. nonnegative thanks to (2.11c) and (2.11d). Notice that the primal iterates  $x^l$  from  $l = 0$  up to  $l = k + 1$  are measurable with respect to  $\mathcal{F}^k$ , whereas the dual iterates  $y^l$  from  $l = 0$  up to  $l = k$  are measurable with respect to  $\mathcal{F}^k$ . Hence,  $a^k$  and  $b^k$  are measurable with respect to  $\mathcal{F}^k$ . Furthermore, inequalities (2.10), (2.11a) and (2.11b) imply that almost surely for all  $k \in \mathbb{N} \setminus \{0\}$ ,

$$\mathbb{E} [a^{k+1} | \mathcal{F}^k] \leq (1 + \eta^k)a^k - b^k.$$

By Robbins–Siegmund lemma [23],  $(a^k)$  converges almost surely,  $\sup_k \mathbb{E} [a^k] < \infty$  and  $\sum_{k=1}^\infty \mathbb{E} [b^k] < \infty$ . From the last point in particular, we can write thanks to (2.11d) and the monotone convergence theorem:

$$\begin{aligned} &\mathbb{E} \left[ \sum_{k=1}^\infty \|y^k - y^{k-1}\|_{(P^{\Sigma^{k+1}})^{-1}}^2 \right] \\ &\leq \mathbb{E} \left[ \sum_{k=1}^\infty \|(x^{k+1} - x^k, y^k - y^{k-1})\|_{N^{k+1}}^2 \right] \\ &\leq (1 - \beta)^{-1} \mathbb{E} \left[ \sum_{k=1}^\infty b^k \right] = (1 - \beta)^{-1} \sum_{k=1}^\infty \mathbb{E} [b^k] < \infty, \end{aligned}$$

hence  $\sum_{k=1}^\infty \|y^k - y^{k-1}\|_{(P^{\Sigma^{k+1}})^{-1}}^2$  is almost surely finite, thus  $(\|y^k - y^{k-1}\|_{(P^{\Sigma^{k+1}})^{-1}})_{k \in \mathbb{N} \setminus \{0\}}$ , and in turn  $(\|y^k - y^{k-1}\|_{(P^{\Sigma^{k+1}})^{-1}})_{k \in \mathbb{N} \setminus \{0\}}$ , converge almost surely to 0. Furthermore,  $\sup_k \mathbb{E} [a^k] < \infty$  hence  $\sup_k \|x^k - x^*\|_{(\tau^k)^{-1}}$ , and in turn  $\sup_k \|x^k - x^*\|_{(\tau^k)^{-1}}$ , are finite, and by (2.11e), one can write that for  $k \in \mathbb{N} \setminus \{0\}$ ,

$$\begin{aligned} &\left| \langle P^{-1}A(x^k - x^*), y^k - y^{k-1} \rangle \right| \\ &\leq \sqrt{\beta} \|x^k - x^*\|_{(\tau^{k+1})^{-1}} \|y^k - y^{k-1}\|_{(P^{\Sigma^{k+1}})^{-1}} \\ &\leq \sqrt{\beta(1 + \eta^k)} \|x^k - x^*\|_{(\tau^k)^{-1}} \|y^k - y^{k-1}\|_{(P^{\Sigma^{k+1}})^{-1}}. \end{aligned}$$

We know that  $(\eta^k)_{k \in \mathbb{N}}$  is summable hence converges to 0. As a consequence,

$$|\langle P^{-1}A(x^k - x^*), y^k - y^{k-1} \rangle| \rightarrow 0 \quad \text{almost surely.}$$

To conclude with, thanks to the identity

$$\begin{aligned} a^k &= \|(x^k - x^*, y^k - y^*)\|_{N^k}^2 \\ &\quad + \langle P^{-1}A(x^k - x^*), y^k - y^{k-1} \rangle, \quad k \in \mathbb{N} \setminus \{0\}, \end{aligned}$$

the almost-sure convergence of  $(a^k)_{k \in \mathbb{N}}$  implies in turn that of  $(\|(x^k - x^*, y^k - y^*)\|_{N^k}^2)_{k \in \mathbb{N}}$ .

Let us now turn to point (ii). The first assertion is a straightforward consequence of

$$\mathbb{E} \left[ \sum_{k=1}^\infty b^k \right] = \sum_{k=1}^\infty \mathbb{E} [b^k] < \infty$$

and bounds (2.11c) and (2.11d). Furthermore, it implies that  $\sum_{k=1}^\infty \|(x^{k+1} - x^k, y^k - y^{k-1})\|^2$  is a.s. finite, hence  $(\|(x^{k+1} - x^k, y^k - y^{k-1})\|)$  a.s. converges to 0, and so does  $(\|x^{k+1} - x^k\|)$ .  $\square$

### 2.5 Weak Cluster Points of A-SPDHG are Saddle Points

The goal of this section is to prove that A-SPDHG satisfies point (iii) of Theorem 2.5. On the event  $\{I^k = i\}$ , A-SPDHG update procedure can be rewritten as

$$\begin{aligned} y_i^{k+1} &= \text{prox}_{\sigma_i^k f_i^*}(y_i^k + \sigma_i^{k+1} A_i x^{k+1}), \quad \bar{y}_i^{k+1} = y_i^{k+1} \\ &\quad + \frac{1}{p_i} (y_i^{k+1} - y_i^k), \quad \bar{y}_j^{k+1} = y_j^k, \quad j \neq i \\ x^{k+2} &= \text{prox}_{\tau^{k+2} g}(x^{k+1} - \tau^{k+2} A^* \bar{y}^{k+1}). \end{aligned}$$

We define  $T_i^{\sigma, \tau} : (x, y) \mapsto (\hat{x}, \hat{y}_i)$  by:

$$\begin{aligned} \hat{y}_i &= \text{prox}_{\sigma_i f_i^*}(y_i + \sigma_i A_i x), \quad \hat{x} \\ &= \text{prox}_{\tau g} \left( x - \tau A^* y - \tau \frac{1 + p_i}{p_i} A_i^*(\hat{y}_i - y_i) \right), \end{aligned}$$

so that  $(x^{k+2}, y_i^{k+1}) = T_i^{\sigma_i^{k+1}, \tau^{k+2}}(x^{k+1}, y^k)$  on the event  $\{I^k = i\}$  (and  $y_j^{k+1} = y_j^k$  for  $j \neq i$ ).

**Lemma 2.9** (Cluster points of A-SPDHG are saddle points) *Let  $(\bar{x}, \bar{y})$  a.s. be a weak cluster point of  $(x^k, y^k)_{k \in \mathbb{N}}$  (meaning that there exists a measurable subset  $\bar{\Omega}$  of  $\Omega$  of probability one such that for all  $\omega \in \bar{\Omega}$ ,  $(\bar{x}(\omega), \bar{y}(\omega))$  is a weak sequential cluster point of  $(x^k(\omega), y^k(\omega))_{k \in \mathbb{N}}$ ) and assume that the assumptions of Theorem 2.1 hold. Then,  $(\bar{x}, \bar{y})$  is a.s. in  $\mathcal{C}$ .*

**Proof** Thanks to Lemma 2.8-(ii) and the monotone convergence theorem,

$$\sum_{k=1}^\infty \mathbb{E} \left[ \|(x^{k+1} - x^k, y^k - y^{k-1})\|^2 \right]$$



$$= \mathbb{E} \left[ \sum_{k=1}^{\infty} \left\| (x^{k+1} - x^k, y^k - y^{k-1}) \right\|^2 \right] < \infty.$$

Now,

$$\begin{aligned} & \sum_{k=1}^{\infty} \mathbb{E} \left[ \left\| (x^{k+1} - x^k, y^k - y^{k-1}) \right\|^2 \right] \\ &= \sum_{k=1}^{\infty} \mathbb{E} \left[ \mathbb{E} \left[ \left\| (x^{k+1} - x^k, y^k - y^{k-1}) \right\|^2 \middle| I^{k-1} \right] \right] \\ &= \sum_{k=1}^{\infty} \sum_{i=1}^n \mathbb{P}(I^{k-1} = i) \mathbb{E} \left[ \left\| T_i^{\sigma_i^k, \tau^{k+1}}(x^k, y_i^{k-1}) - (x^k, y_i^{k-1}) \right\|^2 \right] \\ &= \mathbb{E} \left[ \sum_{i=1}^n p_i \sum_{k=1}^{\infty} \left\| T_i^{\sigma_i^k, \tau^{k+1}}(x^k, y_i^{k-1}) - (x^k, y_i^{k-1}) \right\|^2 \right]. \end{aligned}$$

Hence, we can deduce that

$$\mathbb{E} \left[ \sum_{k=1}^{\infty} \sum_{i=1}^n p_i \left\| T_i^{\sigma_i^k, \tau^{k+1}}(x^k, y_i^{k-1}) - (x^k, y_i^{k-1}) \right\|^2 \right] < \infty.$$

It follows that the series in the expectation is a.s. finite, and since  $p_i > 0$  we deduce that almost surely,

$$\left\| T_i^{\sigma_i^k, \tau^{k+1}}(x^k, y_i^{k-1}) - (x^k, y_i^{k-1}) \right\| \xrightarrow{k \rightarrow \infty} 0 \tag{2.13}$$

for all  $i = 1, \dots, n$ . We consider a sample  $(x^k, y^k)$  which is bounded and such that (2.13) holds. We let for each  $i$ ,  $(\hat{x}^{i,k+1}, \hat{y}_i^{i,k}) = T_i^{\sigma_i^k, \tau^{k+1}}(x^k, y_i^{k-1})$ , so that  $\|(\hat{x}^{i,k+1}, \hat{y}_i^{i,k}) - (x^k, y_i^{k-1})\| \rightarrow 0$  for  $i = 1, \dots, n$ . Then, one has

$$\begin{aligned} \partial f_i^*(\hat{y}_i^{i,k}) &\ni \frac{y_i^{k-1} - \hat{y}_i^{i,k}}{\sigma_i^k} + A_i x^k =: A_i x^k + \delta_y^{i,k} \\ \partial g(\hat{x}^{i,k+1}) &\ni \frac{x^k - \hat{x}^{i,k+1}}{\tau^{k+1}} - A^* y^{k-1} \\ &\quad - \frac{1 + p_i}{p_i} A_i^*(\hat{y}_i^{i,k} - y_i^{k-1}) \\ &=: -A^* y^{k-1} + \delta_x^{i,k} \end{aligned}$$

where  $\delta_{x,y}^{i,k} \rightarrow 0$  as  $k \rightarrow \infty$ . Given a test point  $(x, y)$ , one may write for any  $k$ :

$$\begin{aligned} f_i^*(y_i) &\geq f_i^*(\hat{y}_i^{i,k}) + \langle A_i x^k, y_i - y_i^{k-1} \rangle + \langle A_i x^k, y_i^{k-1} - \hat{y}_i^{i,k} \rangle \\ &\quad + \langle \delta_y^{i,k}, y_i - \hat{y}_i^{i,k} \rangle, \quad i = 1, \dots, n \\ g(x) &\geq g(\hat{x}^{1,k+1}) - \langle A^* y^{k-1}, x - x^k \rangle \\ &\quad - \langle A^* y^{k-1}, x^k - \hat{x}^{1,k+1} \rangle + \langle \delta_x^{i,k}, x - \hat{x}^{1,k+1} \rangle \end{aligned}$$

and summing all these inequalities, we obtain:

$$\begin{aligned} g(x) + \sum_{i=1}^n f_i^*(y_i) &\geq g(\hat{x}^{1,k+1}) + \sum_{i=1}^n \left( f_i^*(\hat{y}_i^{i,k}) + \langle A_i x^k, y_i \rangle \right) \\ &\quad - \langle A^* y^{k-1}, x \rangle + \delta^k \end{aligned}$$

where  $\delta^k \rightarrow 0$  as  $k \rightarrow \infty$ . We deduce that if  $(\bar{x}, \bar{y})$  is the weak limit of a subsequence  $(x^{k_l}, y^{k_l-1})$  (as well as, of course,  $(x^{k_l}, y^{k_l})$ ), then:

$$\begin{aligned} g(x) + \sum_{i=1}^n f_i^*(y_i) &\geq g(\bar{x}) \\ &\quad + \sum_{i=1}^n \left( f_i^*(\bar{y}_i) + \langle A_i \bar{x}, y_i \rangle \right) - \langle A^* \bar{y}, x \rangle. \end{aligned}$$

Since  $(x, y)$  is arbitrary, we find that (2.2) holds for  $(\bar{x}, \bar{y})$ . □

### 2.6 Proof of Theorem 2.1

Under the assumptions of Theorem 2.1, the set  $\mathcal{C}$  of saddle points is closed and non-empty and  $X \times Y$  is a separable Hilbert space. By Lemma 2.3, the variable metrics sequence  $(N^k)_{k \in \mathbb{N}}$  defined in (2.8) satisfies condition (i) of Theorem 2.5. Furthermore, the iterates of Algorithm 2.1 comply with condition (ii) and (iii) of Theorem 2.5 by Lemma 2.8 and Lemma 2.9, respectively, and hence converge almost surely to a point in  $\mathcal{C}$ .

## 3 Algorithmic Design and Practical Implementations

In this section, we present practical instances of our A-SPDHG algorithm, where we specify a step-size adjustment rule which satisfies our assumptions in convergence proof. We extend the adaptive step-size balancing rule for deterministic PDHG, which is proposed by [14], into our stochastic setting, with minibatch approximation to minimize the computational overhead.

### 3.1 A-SPDHG Rule (a)—Tracking and Balancing the Primal–Dual Progress

Let’s first briefly introduce the foundation of our first numerical scheme, which is built upon the deterministic adaptive PDHG algorithm proposed by Goldstein et al [14], with the iterates:

$$\begin{aligned} x^{k+1} &= \text{prox}_{\tau^{k+1}g}(x^k - \tau^{k+1}A^*y^k), \quad y^{k+1} \\ &= \text{prox}_{\sigma^{k+1}f^*}(y^k + \sigma^{k+1}A(2x^{k+1} - x^k)) \end{aligned}$$

In this foundational work of Goldstein et al [14], they proposed to evaluate two sequences in order to track and balance the progresses of the primal and dual iterates of deterministic PDHG (denoted here as  $v_k^*$  and  $d_k^*$ ):

$$\begin{aligned} v_k^* &:= \|(x^k - x^{k+1})/\tau^{k+1} - A^*(y^k - y^{k+1})\|_1, \quad d_k^* \\ &:= \|(y^k - y^{k+1})/\sigma^{k+1} - A(x^k - x^{k+1})\|_1. \end{aligned} \tag{3.1}$$

These two sequences measure the lengths of the primal and dual subgradients for the objective  $\min_{x \in X} \max_{y \in Y} g(x) + \langle Ax, y \rangle - f^*(y)$ , which can be demonstrated by the definition of proximal operators. The primal update of deterministic PDHG can be written as:

$$x^{k+1} = \arg \min_x \frac{1}{2} \|x - (x^k - \tau^{k+1} A^* y^k)\|_2^2 + \tau^{k+1} g(x). \tag{3.2}$$

The optimality condition of the above objective declares:

$$0 \in \partial g(x^{k+1}) + A^* y^k + \frac{1}{\tau^{k+1}} (x^{k+1} - x^k). \tag{3.3}$$

By adding  $-A^* y^{k+1}$  on both sides and rearranging the terms, one can derive:

$$(x^k - x^{k+1})/\tau^{k+1} - A^*(y^k - y^{k+1}) \in \partial g(x^{k+1}) + A^* y^{k+1} \tag{3.4}$$

and similarly for the dual update one can also derive:

$$(y^k - y^{k+1})/\sigma^{k+1} - A(x^k - x^{k+1}) \in \partial f^*(y^{k+1}) - Ax^{k+1}, \tag{3.5}$$

which indicates that the sequences  $v_k^*$  and  $d_k^*$  given by (3.1) should effectively track the primal progress and dual progress of deterministic PDHG, and hence, Goldstein et al [14] propose to utilize these as the basis of balancing the primal and dual step sizes for PDHG.

In light of this, we propose our first practical implementation of A-SPDHG in Algorithm 3.1 as our rule-(a), where we use a unique dual step-size  $\sigma^k = \sigma_j^k$  for all iterates  $k$  and indices  $j$  and where we estimate the progress of achieving optimality on the primal and dual variables via the two sequences  $v^k$  and  $d^k$  defined at each iteration  $k$  with  $I^k = i$  as:

$$\begin{aligned} v_{k+1} &:= \|(x^k - x^{k+1})/\tau^{k+1} - \frac{1}{p_i} A_i^*(y_i^k - y_i^{k+1})\|_1, \quad d_{k+1} \\ &:= \frac{1}{p_i} \|(y_i^k - y_i^{k+1})/\sigma^{k+1} - A_i(x^k - x^{k+1})\|_1, \end{aligned} \tag{3.6}$$

which are minibatch extension of (3.1) tailored for our stochastic setting. By making them balanced on the fly via adjusting the primal–dual step-size ratio when appropriate, we can enforce the algorithm to achieve similar progress in both primal and dual steps and hence improve the convergence. To be more specific, as shown in Algorithm 3.1, in each iteration the values of  $v_k$  and  $d_k$  are evaluated and compared. If the value of  $v_k$  (which tracks the primal subgradients) is significantly larger than  $d_k$  (which tracks the dual subgradients), then we know that the primal progress is slower than the dual progress, and hence, the algorithm would boost the primal step size while shrinking the dual step size. If  $v_k$  is noticeably smaller than  $d_k$ , then the algorithm would do the opposite.

Note that here we adopt the choice of  $\ell_1$ -norm as the length measure for  $v^k$  and  $d^k$  as done by Goldstein et al [14, 15], since we also observe numerically the benefit over the more intuitive choice of  $\ell_2$ -norm.

For full-batch case ( $n = 1$ ), it reduces to the adaptive PDHG proposed by [14, 15]. We adjust the ratio between primal and dual step sizes according to the ratio between  $v^k$  and  $d^k$ , and whenever the step-size change, we shrink  $\alpha$  (which controls the amplitude of the changes) by a factor  $\eta \in (0, 1)$ —we typically choose  $\eta = 0.995$  in our experiments. For the choice of  $s$ , we choose  $s = \|A\|$  as our default.<sup>1</sup>

### 3.1.1 Reducing the Overhead with Subsampling

Noting that unlike the deterministic case which does not have the need of extra matrix–vector multiplication since  $A^* y^k$  and  $Ax^k$  can be memorized, our stochastic extension

<sup>1</sup> The choice of  $s$  is crucial for the convergence behavior of rule (a), and we found numerically that it is better to scale with the operator norm  $\|A\|$  instead of depending on the range of pixel values as suggested in [15].

---

**Algorithm 3.1:** A-SPDHG, rule (a)

---

Input: dual step-size  $\sigma^0$ , primal step-size  $\tau^0$ ,  $\alpha^0 \in (0, 1)$ ,  $\eta \in (0, 1)$ ,  $\delta > 1$ , probabilities  $(p_i)_{1 \leq i \leq n}$ ; primal variable  $x^0$ , dual variable  $y^0$   
 Initialize  $\bar{y}^0 = y^0$ ,  $v^0 = d^0 = 0$ ,  $s = \|A\|$   
**for**  $k \in \llbracket 0, K - 1 \rrbracket$  **do**  
   **if**  $v^k > s d^k \delta$  **then**  $\tau^{k+1} = \frac{\tau^k}{1 - \alpha^k}$ ,  $\sigma^{k+1} = \sigma^k (1 - \alpha^k)$ ,  
    $\alpha^{k+1} = \alpha^k \eta$   
   **if**  $v^k < s d^k / \delta$  **then**  $\tau^{k+1} = \tau^k (1 - \alpha^k)$ ,  $\sigma^{k+1} = \frac{\sigma^k}{1 - \alpha^k}$ ,  
    $\alpha^{k+1} = \alpha^k \eta$   
   **if**  $s d^k / \delta \leq v^k \leq s d^k \delta$  **then**  $\tau^{k+1} = \tau^k$ ,  $\sigma^{k+1} = \sigma^k$ ,  $\alpha^{k+1} = \alpha^k$   
    $x^{k+1} = \text{prox}_{\tau^{k+1} g}(x^k - \tau^{k+1} A^* \bar{y}^k)$   
   Randomly pick  $i \in \llbracket 1, n \rrbracket$  with probability  $p_i$   
    $y_j^{k+1} = \begin{cases} \text{prox}_{\sigma^{k+1} f_i^*}(y_i^k + \sigma^{k+1} A_i x^{k+1}) & \text{if } j = i \\ y_j^k & \text{if } j \neq i \end{cases}$   
    $\bar{y}_j^{k+1} = \begin{cases} y_i^{k+1} + \frac{1}{p_i} (y_i^{k+1} - y_i^k) & \text{if } j = i \\ y_j^k & \text{if } j \neq i \end{cases}$   
    $v^{k+1} = \|(x^k - x^{k+1})/\tau^{k+1} - \frac{1}{p_i} A_i^*(y_i^k - y_i^{k+1})\|_1$   
    $d^{k+1} = \frac{1}{p_i} \|(y_i^k - y_i^{k+1})/\sigma^{k+1} - A_i(x^k - x^{k+1})\|_1$  – or  
   approximate this step by (3.7)  
**end for**  
**return**  $x^K$

---

will require the computation of  $A_i x^k$  since we will sample different subsets between back-to-back iterations with high probability. When using this strategy, we will only have a maximum 50% overhead in terms of FLOP counts, which is numerically negligible compared to the significant acceleration it will bring toward SPDHG especially when the primal–dual step-size ratio is suboptimal, as we will demonstrate later in the experiments. Moreover, we found numerically that we can significantly reduce this overhead by approximation tricks such as subsampling:

$$d^{k+1} \approx \frac{\rho}{p_i} \|S^k (y_i^k - y_i^{k+1}) / \sigma^{k+1} - S^k A_i (x^k - x^{k+1})\|_1 \tag{3.7}$$

with  $S^k$  being a random subsampling operator such that  $\mathbb{E}[(S^k)^T S^k] = \frac{1}{\rho} \text{Id}$ . In our experiments, we choose 10% subsampling for this approximation and hence the overhead is reduced from 50% to only 5% which is negligible, without compromising the convergence rates in practice.

### 3.2 A-SPDHG Rule (b)—Exploiting Angle Alignments

More recently, Yokota and Hontani [26] propose a variant of adaptive step-size balancing scheme for PDHG, utilizing the angles between the subgradients  $\partial g(x^{k+1}) + A^* y^{k+1}$  and the difference of the updates  $x^k - x^{k+1}$ .

If these two directions are highly aligned, then the primal step size can be increased for bigger step. If these two directions have a large angle, then the primal step size should be shrunken. By extending this scheme to stochastic setting, we obtain another choice of adaptive scheme for SPDHG.

---

#### Algorithm 3.2: A-SPDHG, rule (b)

---

Input: dual step-size  $\sigma^0$ , primal step-size  $\tau^0$ ,  $\eta \in (0, 1)$ , probabilities  $(p_i)_{1 \leq i \leq n}$ ; primal variable  $x^0$ , dual variable  $y^0$   
 Initialize  $\bar{y}^0 = y^0$ ,  $w^0 = 0$ ,  $\alpha^0 = 1$   
**for**  $k \in \llbracket 0, K - 1 \rrbracket$  **do**  
   **If**  $w^k < 0$  **then**  $\tau^{k+1} = \frac{\tau^k}{1 + \alpha^k}$ ,  $\sigma^{k+1} = \sigma^k (1 + \alpha^k)$ ,  $\alpha^{k+1} = \alpha^k \eta$   
   **If**  $w^k \geq c$  **then**  $\tau^{k+1} = \tau^k (1 + \alpha^k)$ ,  $\sigma^{k+1} = \frac{\sigma^k}{1 + \alpha^k}$ ,  $\alpha^{k+1} = \alpha^k \eta$   
   **If**  $0 \leq w^k < c$  **then**  $\tau^{k+1} = \tau^k$ ,  $\sigma^{k+1} = \sigma^k$ ,  $\alpha^{k+1} = \alpha^k$   
    $x^{k+1} = \text{prox}_{\tau^{k+1} g}(x^k - \tau^{k+1} A^* \bar{y}^k)$   
   Randomly pick  $i \in \llbracket 1, n \rrbracket$  with probability  $p_i$   
    $y_j^{k+1} = \begin{cases} \text{prox}_{\sigma^{k+1} f_i^*}(y_i^k + \sigma^{k+1} A_i x^{k+1}) & \text{if } j = i \\ y_j^k & \text{if } j \neq i \end{cases}$   
    $\bar{y}_j^{k+1} = \begin{cases} y_i^{k+1} + \frac{1}{p_i} (y_i^{k+1} - y_i^k) & \text{if } j = i \\ y_j^k & \text{if } j \neq i \end{cases}$   
    $q^{k+1} = (x^k - x^{k+1}) / \tau^{k+1} - \frac{1}{p_i} A_i^* (y_i^k - y_i^{k+1})$   
    $w^{k+1} = \langle x^k - x^{k+1}, q^{k+1} \rangle / (\|x^k - x^{k+1}\|_2 \|q^{k+1}\|_2)$   
**end for**  
**return**  $x^K$

---

We present this scheme in Algorithm 3.2 as our rule (b). At iteration  $k$  with  $I^k = i$ , compute:

$$q^{k+1} = (x^k - x^{k+1}) / \tau^{k+1} - \frac{1}{p_i} A_i^* (y_i^k - y_i^{k+1}), \tag{3.8}$$

as an estimate of  $\partial g(x^{k+1}) + A^* y^{k+1}$ , then measure the cosine of the angle between this and  $x^k - x^{k+1}$ :

$$w^{k+1} = \frac{\langle x^k - x^{k+1}, q^{k+1} \rangle}{(\|x^k - x^{k+1}\|_2 \|q^{k+1}\|_2)}. \tag{3.9}$$

The threshold  $c$  for the cosine value (which triggers the increase of the primal step size) typically needs to be very close to 1 (we use  $c = 0.999$ ) due to the fact that we mostly apply these type of algorithms in high-dimensional problems, following the choice in [26] which was for deterministic PDHG.

Recently, Zdun et al [27] proposed a heuristic similar to our rule (b), but they choose  $q^{k+1}$  to be the approximation for an element of  $\partial g(x^{k+1})$  instead of  $\partial g(x^{k+1}) + A^* y^{k+1}$ . Our choice follows more closely to the original scheme of Yokota and Hontani [26]. We numerically found that their scheme is not competitive in our settings.

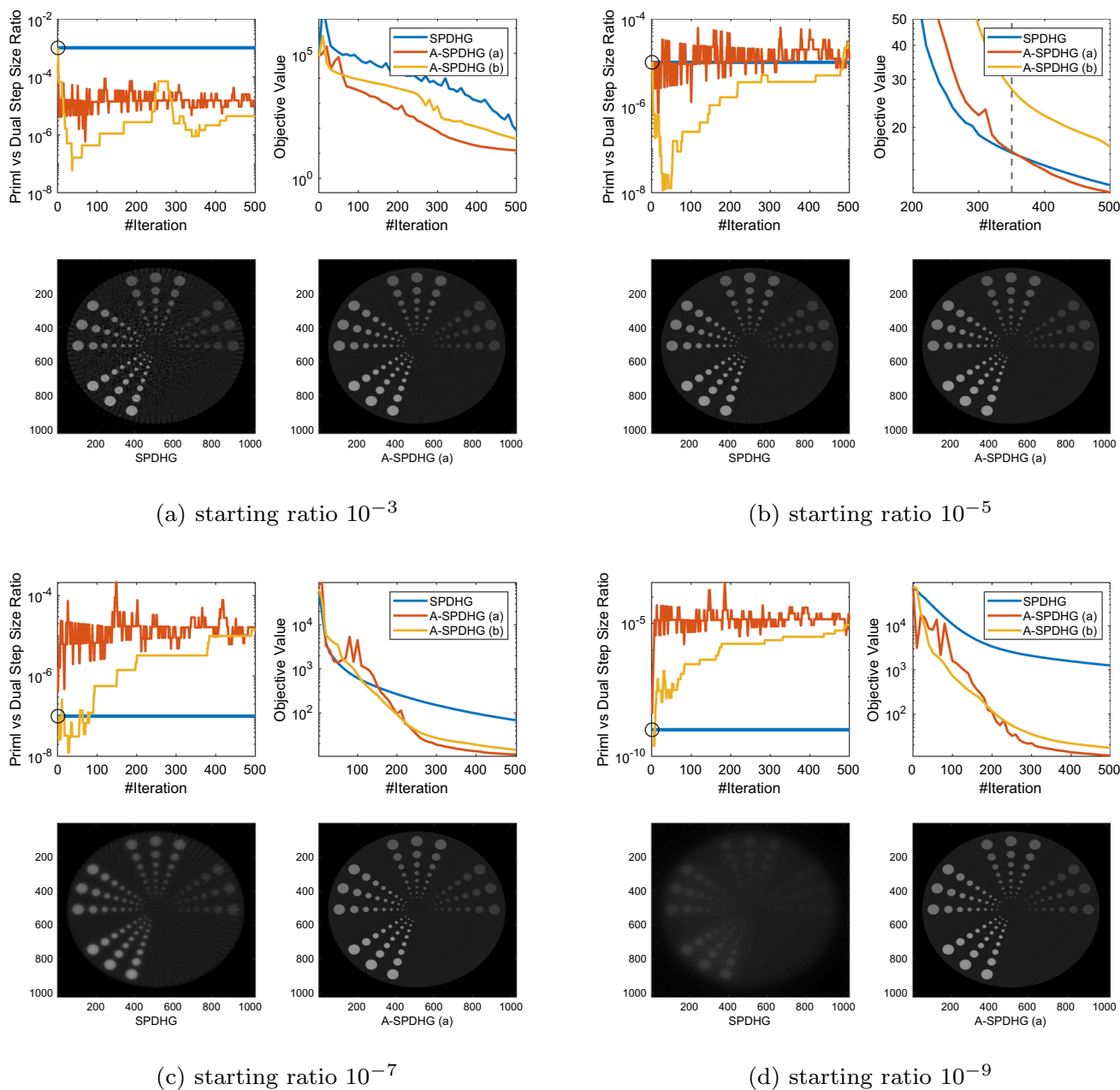
## 4 Numerical Experiments

In this section, we present numerical studies of the proposed scheme in solving one of the most typical imaging inverse problems, the computed tomography (CT). We compare A-SPDHG algorithm with the original SPDHG, on different choices of starting ratio of the primal and dual step sizes.

In our CT imaging example, we seek to reconstruct the tomography images from fanbeam X-ray measurement data, by solving the following TV-regularized objective:

$$x^* \in \arg \min_{x \in \mathbb{R}^d} \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|Dx\|_1 \tag{4.1}$$

where  $D$  denotes the 2D differential operator,  $A \in \mathbb{R}^{m \times d}$  and  $x \in \mathbb{R}^d$ . We consider three fanbeam CT imaging modalities: sparse-view CT, low-dose CT and limited-angle CT. We test the A-SPDHG and SPDHG on two images of different sizes (Example 1 on a phantom image sized  $1024 \times 1024$ , while Example 2 being an image from the Mayo Clinic Dataset [21] sized  $512 \times 512$ .), on 4 different starting ratios ( $10^{-3}$ ,  $10^{-5}$ ,  $10^{-7}$  and  $10^{-9}$ ). We interleave partitioned the measurement data and operator into  $n = 10$  minibatches for both algorithms. To be more specific, we first collect all the X-ray measurement data and list them consecutively from 0 degree to 360 degree to form the full  $A$  and  $b$ , and then interleavingly group every 10-th of the measurements into one minibatch, to form the partition  $\{A_i\}_{i=1}^{10}$  and  $\{b_i\}_{i=1}^{10}$ .



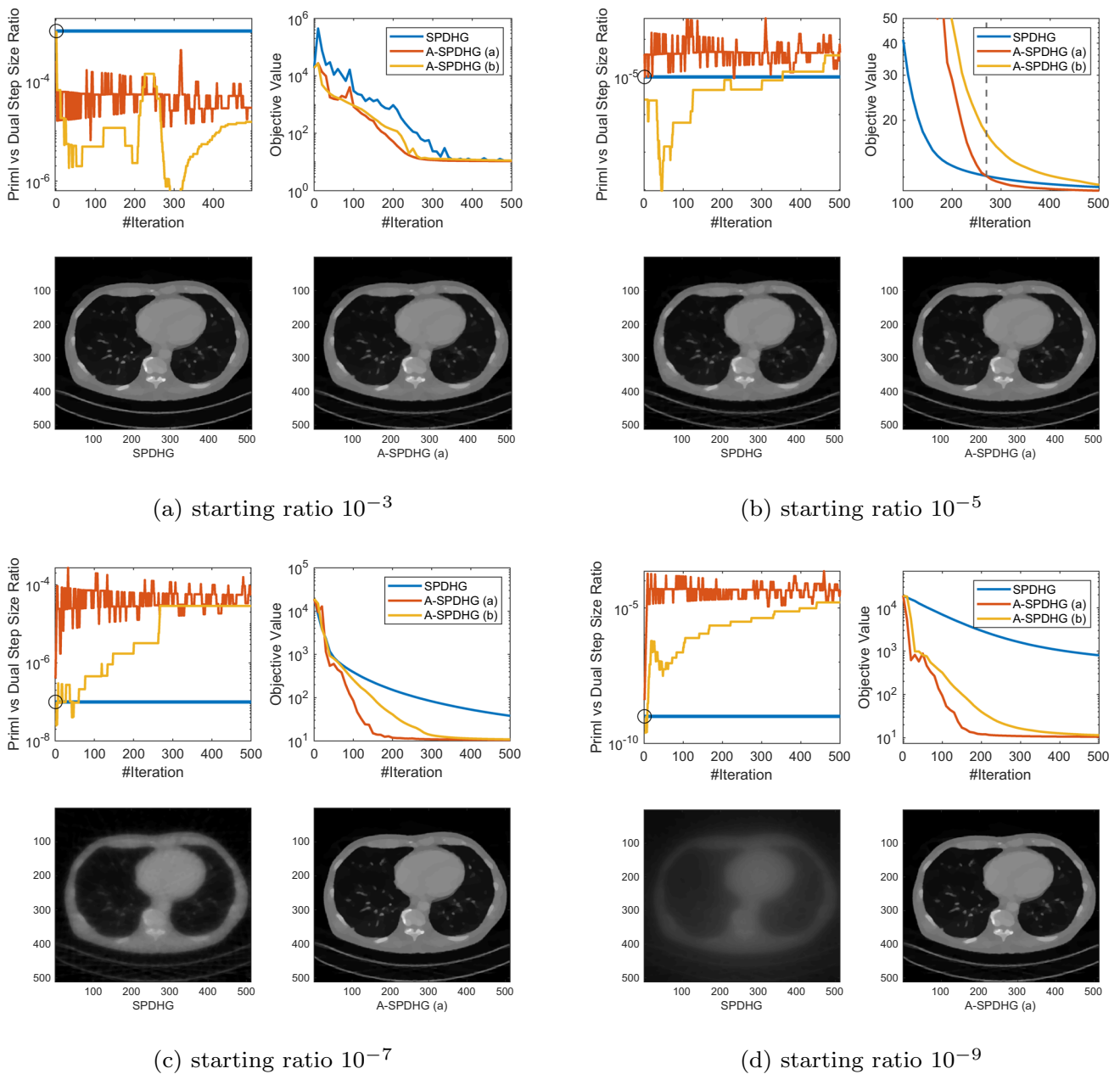
**Fig. 1** Comparison between SPDHG and A-SPDHG on sparse-view CT (Example 1), with a variety of starting primal–dual step size ratios. Here, the forward operator is  $A \in \mathbb{R}^{m \times d}$  with dimensions  $m = 368640$ ,  $d = 1048576$ . We include the images reconstructed by the algorithms

at termination (50th epoch). In the first plot of each subfigure, the black circle indicates the starting step-size ratio for all the algorithms, same for the following figures

For A-SPDHG, we choose to use the approximation step for  $d^k$  presented in (3.7) with 10% subsampling and hence the computational overhead is negligible in this experiment. We initialize all algorithms from a zero image.

We present our numerical results in Figs. 1, 2, 3 and 6. In these plots, we compare the convergence rates of the algorithms in terms of number of iterations (the execution time per iteration for the algorithms are almost the same, as the

overhead of A-SPDHG is trivial numerically). Among these, Figs. 1 and 2 report the results for large-scale sparse-view CT experiments on a phantom image and a lung CT image from Mayo Clinic dataset [21], while Fig. 3 reports the results for low-dose CT experiments where we simulate a large number of measurements corrupted with a significant amount Poisson noise, and then, in Fig. 6 we report the results for limited-angle CT which only a range of 0-degree to 150-degree of

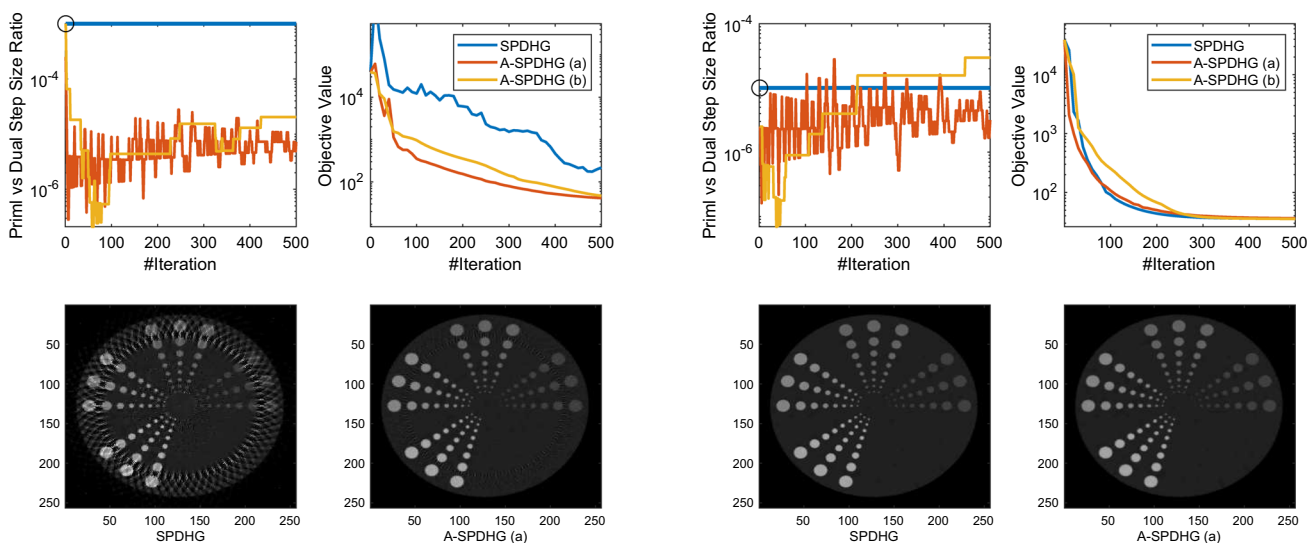


**Fig. 2** Comparison between SPDHG and A-SPDHG on sparse-view CT (Example 2), with a variety of starting primal–dual step-size ratios. Here, the forward operator is  $A \in \mathbb{R}^{m \times d}$  with dimensions  $m = 92160, d = 262144$ . We include the images reconstructed by the algorithms at termination (50th epoch)

measurement angles are present, while the measurements from the rest [150, 360] degrees of angles are all missing. In all these examples, we can consistently observe that no matter how we initialize the primal–dual step-size ratio, A-SPDHG can automatically and consistently adjust the step-size ratio to the optimal choice which is around either  $10^{-5}$  or  $10^{-7}$  for these four different CT problems and significantly outperform the vanilla SPDHG for the cases where the starting ratio is away from the optimal range. Meanwhile, even for the cases where the starting ratio of SPDHG algorithm is

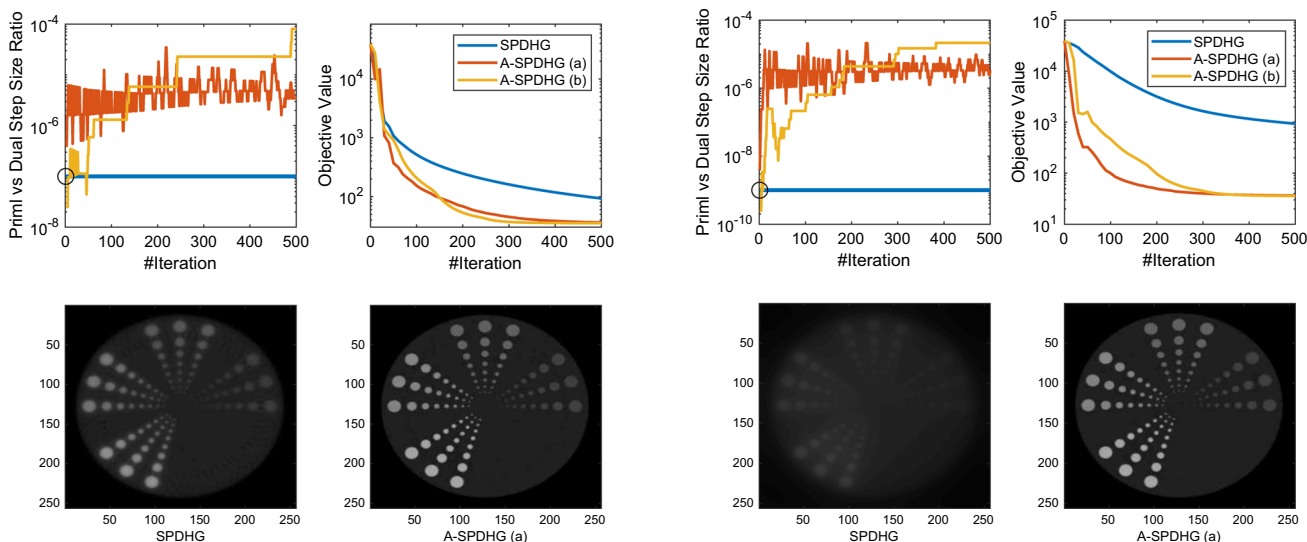
near-optimal, we can observe consistently from most of these examples that our scheme outperforms the vanilla SPDHG algorithm locally after a certain number of iterations (highlighted by the vertical dash lines in relevant subfigures), which further indicates the benefit of adaptivity for this class of algorithms<sup>2</sup>. Note that throughout all these different exam-

<sup>2</sup> The most typical example here would be Fig. 1b where the optimal step-size ratio selected by the adaptive scheme at convergence is almost exactly  $10^{-5}$ , where we have set SPDHG to run with this ratio. We



(a) starting ratio  $10^{-3}$

(b) starting ratio  $10^{-5}$



(c) starting ratio  $10^{-7}$

(d) starting ratio  $10^{-9}$

**Fig. 3** Comparison between SPDHG and A-SPDHG on low-dose CT (where we use a large number of highly-noisy X-ray measurements), with a variety of starting primal–dual step-size ratios. Here, the for-

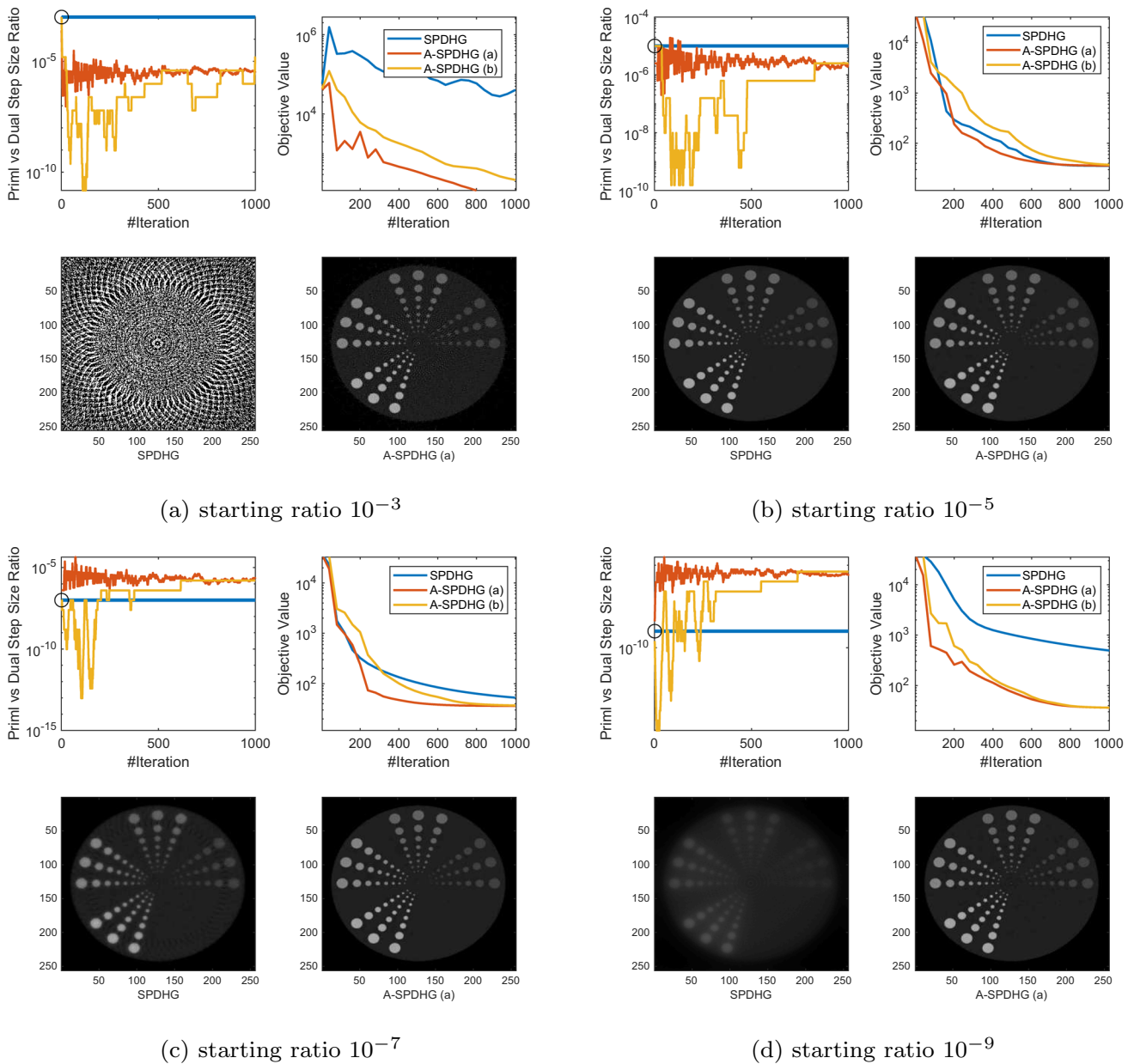
ward operator is  $A \in \mathbb{R}^{m \times d}$  with dimensions  $m = 184320, d = 65536$ . We resized the phantom image to 256 by 256. We include the images reconstructed by the algorithms at termination (50th epoch)

ples, we use only one fixed set of parameters for A-SPDHG suggested in the previous section, which again indicates the strong practicality of our scheme.

For the low-dose CT example, we run two extra sets of experiments, regarding a larger number of partitioning of minibatches (40) in Fig. 4, and warm-start from a better initialization image obtained via filter backprojection in

can still observe benefit of local convergence acceleration given by our adaptive scheme.

Fig. 5. We found that in all these extra examples we consistently observe superior performances of A-SPDHG over the vanilla SPDHG especially when the primal–dual step-size ratios are suboptimal. Interestingly, we found that the warm-start’s effect does not have noticeable impact of the comparative performances between SPDHG and A-SPDHG. This is mainly due to the fact that the SPDHG with suboptimal primal–dual step-size ratio will converge very slowly in



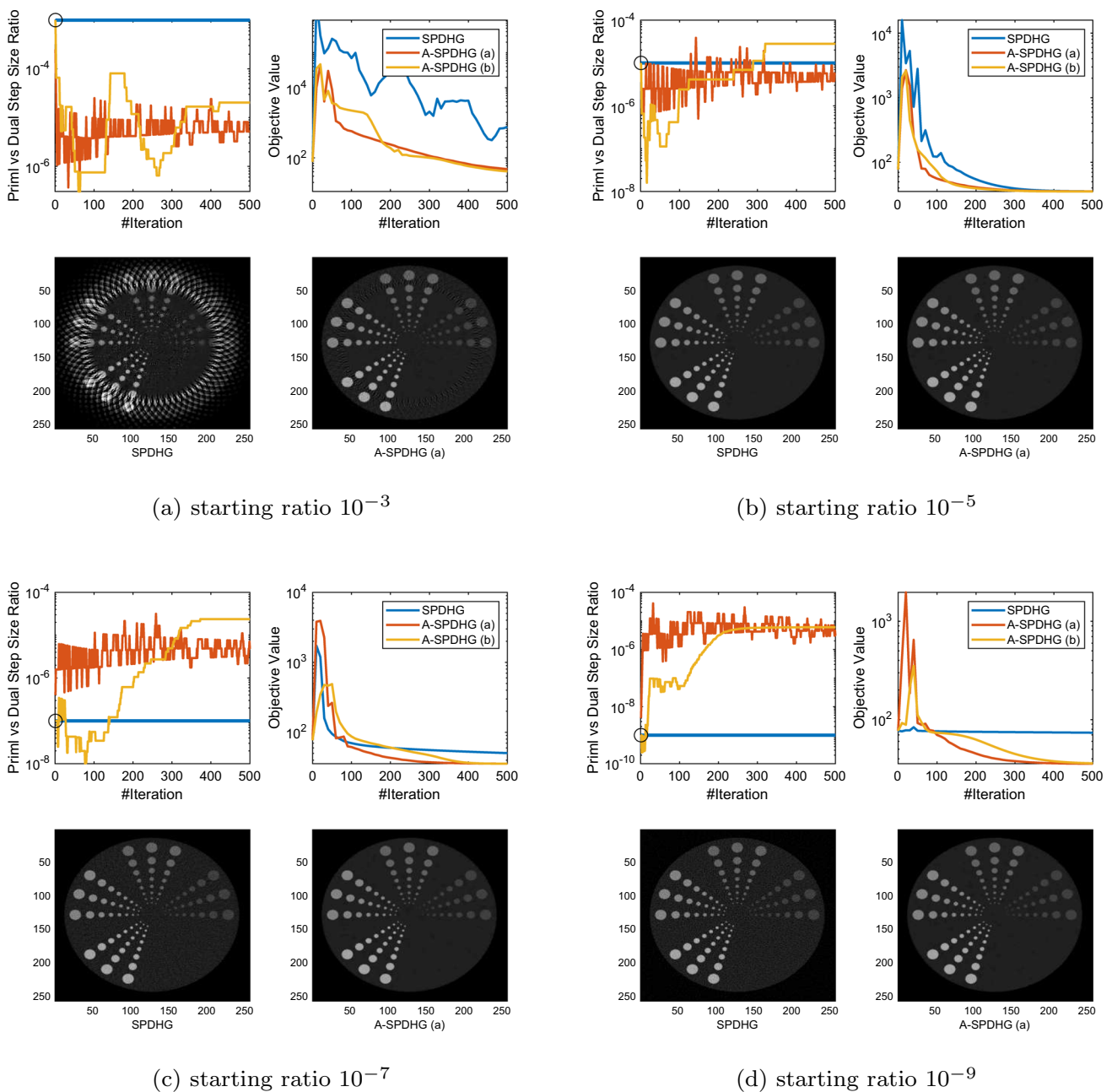
**Fig. 4** Comparison between SPDHG and A-SPDHG with the data being split to 40 minibatches on low-dose CT. Comparing to the results presented in Fig. 3 which used 10 minibatches, we obtain similar results and our A-SPDHG continues to perform more favorably comparing to SPDHG.

high accuracy regimes (see Fig. 5d for example) in practice hence the warm-start won't help much here.

We should also note that conceptually all the hyperparameters in our adaptive schemes are basically the controllers of the adaptivity of the algorithm (while for extreme choices we recover the vanilla SPDHG). In Figs. 7 and 9, we present some numerical studies on the choices of hyperparameters of rule (a) and rule (b) of A-SPDHG algorithm. We choose the fixed starting ratio of  $10^{-7}$  for primal–dual step sizes in these experiments. For rule (a), we found that it is robust to the choice of the starting shrinking rate  $\alpha_0$ , shrinking speed  $\eta$

and the gap  $\delta$ . Overall, we found that these parameters have weak impact of the convergence performance of our rule (a) and easy to choose.

For rule (b), we found that the performance is more sensitive to the choice of parameter  $c$  and  $\eta$  comparing to rule (a), although the dependence is still weak. Our numerical studies suggest that rule (a) is a better-performing choice than rule (b), but each of them have certain mild weaknesses (the first rule has a slight computational overhead which can be partially addressed with subsampling scheme, while the second rule seems often being slower than the first rule), which



**Fig. 5** Comparison between SPDHG and A-SPDHG with warm-start using a FBP (filtered backprojection) on low-dose CT. Comparing to the results shown in Fig. 3 which are without warm-start, actually our methods seem to compare even more favorably with warm-start. Please

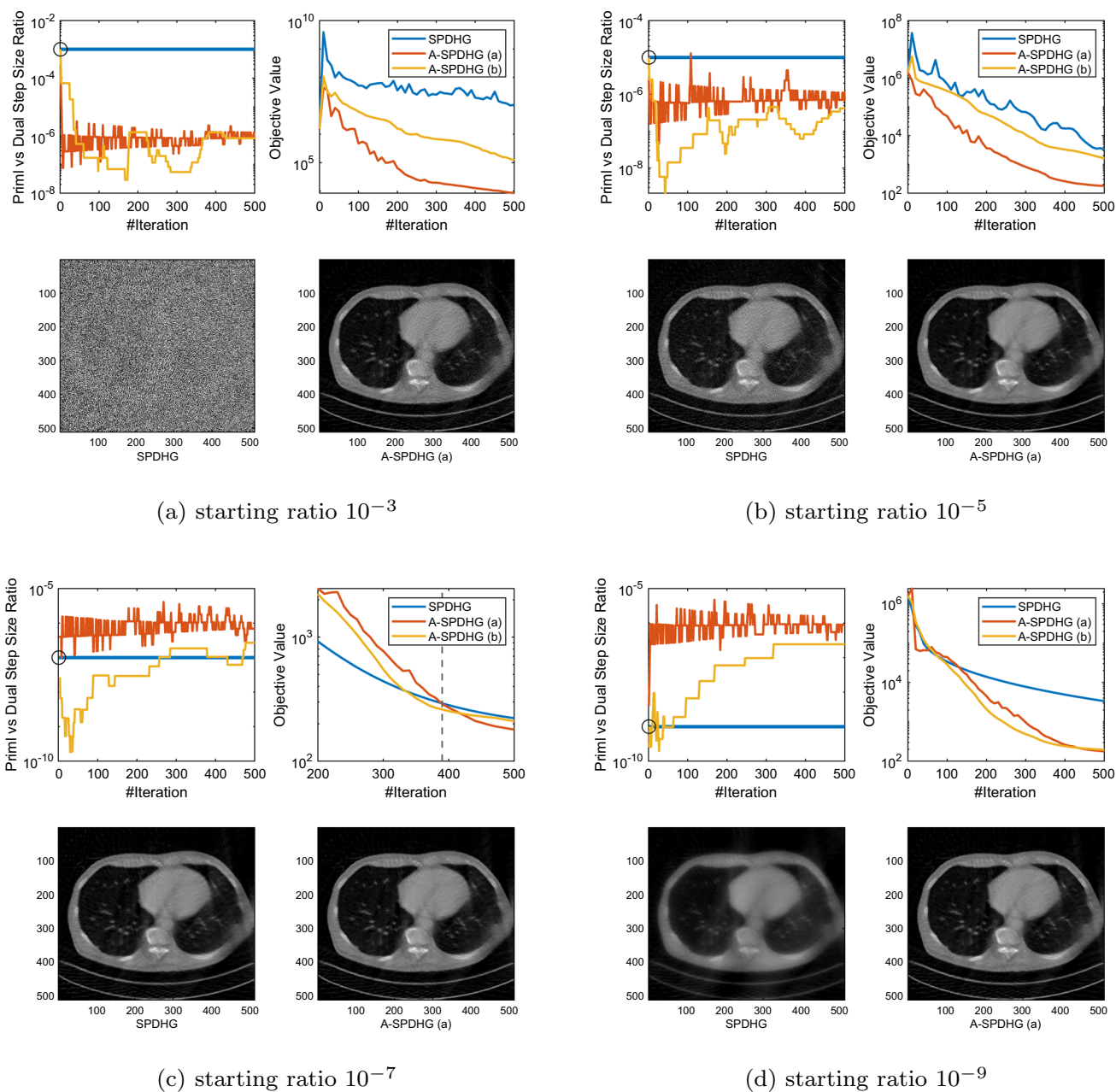
also note that the early jump in terms of function value is within our expectation due to the stochasticity of the algorithms. We include the images reconstructed by the algorithms at termination (50th epoch)

require further studies and improvements. Nevertheless, we need to emphasize that all these parameters are essentially controlling the degree of adaptivity of the algorithms and fairly easy to choose, noting that for all these CT experiments with varying sizes/dimensions and modalities we only use one fixed set of the hyperparameters in A-SPDHG, and we are already able to consistently observe numerical improvements over vanilla SPDHG.

### 5 Conclusion

In this work, we propose a new framework (A-SPDHG) for adaptive step-size balancing in stochastic primal–dual hybrid gradient methods. We first derive theoretically sufficient conditions on the adaptive primal and dual step sizes for ensuring convergence in the stochastic setting. We then propose a number of practical schemes which satisfy the condition for





**Fig. 6** Comparison between SPDHG and A-SPDHG on limited-angle CT (Example 2), with a variety of starting primal–dual step-size ratios. Here, the forward operator is  $A \in \mathbb{R}^{m \times d}$  with dimensions  $m = 92160, d = 262144$ . We include the images reconstructed by the algorithms at termination (50th epoch)

convergence, and our numerical results on imaging inverse problems support the effectiveness of the proposed approach.

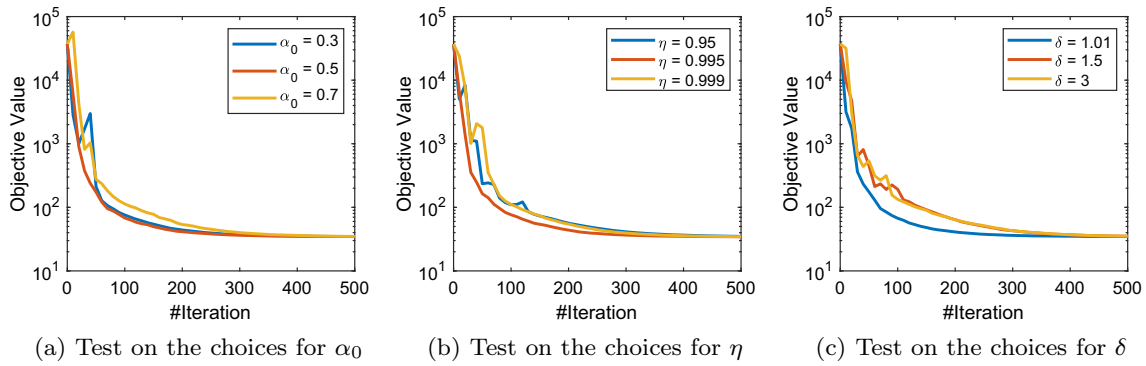
To our knowledge, this work constitutes the first theoretical analysis of adaptive step sizes for a stochastic primal–dual algorithm. Our ongoing work includes the theoretical analysis and algorithmic design of further accelerated stochastic primal–dual methods with line-search schemes for even faster convergence rates.

### 6 Complementary Material for Sect. 2

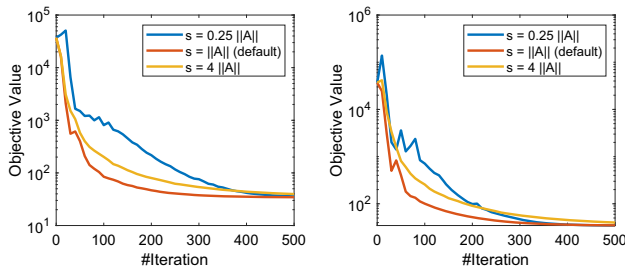
We begin by a useful lemma.

**Lemma 6.1** *Let  $a, b$  be positive scalars,  $\beta \in (0, 1)$ , and  $P$  a bounded linear operator from a Hilbert space  $X$  to a Hilbert space  $Y$ . Then,*

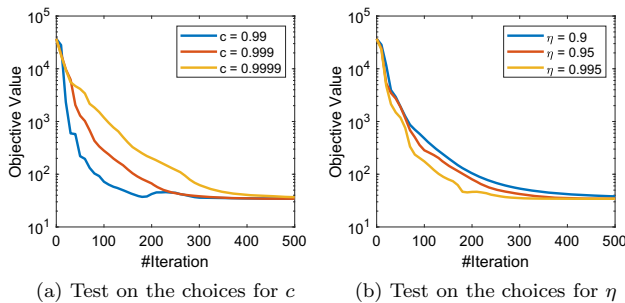
$$(ab)^{-1/2} \|P\| \leq 1 \iff \begin{pmatrix} a Id & P^* \\ P & b Id \end{pmatrix} \succcurlyeq 0. \tag{6.1}$$



**Fig. 7** Test on different choices of parameters of A-SPDHG (rule-a) on X-ray low-dose fanbeam CT example, starting ratio of primal–dual step sizes:  $10^{-7}$ . We can observe that the performance of A-SPDHG has only minor dependence on these parameter choices



**Fig. 8** Test on the default choice  $s = \|A\|$  of A-SPDHG (rule-a) on X-ray low-dose fanbeam CT example. Left figure: starting ratio of primal–dual step sizes:  $10^{-7}$ . Right figure: starting ratio of primal–dual step sizes:  $10^{-5}$ . We can observe that our default choice of  $s$  is indeed a reasonable choice (at least near-optimal) in practice, and when deviating from it may lead to slower convergence



**Fig. 9** Test on different choices of parameters of A-SPDHG (rule-b) on X-ray low-dose fanbeam CT example, starting ratio of primal–dual step sizes:  $10^{-7}$

$$(ab)^{-1/2}\|P\| \leq \beta \Leftrightarrow \begin{pmatrix} a \text{Id} & P^* \\ P & b \text{Id} \end{pmatrix} \succcurlyeq (1 - \beta) \begin{pmatrix} a \text{Id} & 0 \\ 0 & b \text{Id} \end{pmatrix}. \tag{6.2}$$

**Proof** Let us call

$$M = \begin{pmatrix} a \text{Id} & P^* \\ P & b \text{Id} \end{pmatrix}.$$

For all  $(x, y) \in X \times Y$ ,

$$\begin{aligned} \|(x, y)\|_M^2 &\geq a\|x\|^2 + b\|y\|^2 - 2\|P\|\|x\|\|y\| \\ &= \|x\|_a^2 + \|y\|_b^2 - 2(ab)^{-1/2}\|P\|\|x\|_a\|y\|_b, \end{aligned}$$

which proves the direct implication of (6.1). For the converse implication, consider  $x \in X \setminus \{0\}$  such that  $\|Px\| = \|P\|\|x\|$  and  $y = -\lambda Px$  for a scalar  $\lambda$ . Then, the nonnegativity of the polynomial

$$\frac{\|(x, y)\|_M^2}{\|x\|^2} = b\|P\|^2\lambda^2 - 2\|P\|^2\lambda + a$$

for all  $\lambda \in \mathbb{R}$  implies that  $\|P\|^4 - ab\|P\|^2 \leq 0$ , which is equivalent to the desired conclusion  $(ab)^{-1/2}\|P\| \leq 1$ .

Equivalence (6.2) is straightforward by noticing that

$$\begin{pmatrix} a \text{Id} & P^* \\ P & b \text{Id} \end{pmatrix} \succcurlyeq (1 - \beta) \begin{pmatrix} a \text{Id} & 0 \\ 0 & b \text{Id} \end{pmatrix} \Leftrightarrow \begin{pmatrix} \beta a \text{Id} & P^* \\ P & \beta b \text{Id} \end{pmatrix} \succcurlyeq 0. \quad \square$$

Let us now turn to the proof of Lemma 2.2.

**Proof of Lemma 2.2** Let us assume that the step sizes satisfy the assumptions of the lemma. Then, Assumption (i) of Theorem 2.1 is straightforwardly satisfied. Moreover, for  $i \in \llbracket 1, n \rrbracket$ , the product sequence  $(\tau^k \sigma_i^k)_{k \in \mathbb{N}}$  is constant along the iterations by equation (2.6) and satisfies equation (2.5) for iterate  $k = 0$  and thus satisfies (2.5) for all  $k \in \mathbb{N}$  for  $\beta = \max_i \{\tau^0 \sigma_i^0 \|A_i\|^2 / p_i\}$ , which proves Assumption (ii). Finally, equation (2.7) implies that Assumption (iii) is satisfied.  $\square$

**Author Contributions** CD, MJE and AC elaborated the proof strategy, and CD wrote parts 1, 2 and 6. JT worked on the algorithmic design, performed the numerical experiments and wrote parts 3-5. All authors reviewed the manuscript.

**Funding** CD acknowledges support from the EPSRC (EP/S026045/1). MJE acknowledges support from the EPSRC (EP/S026045/1, EP/T026693/1, EP/V026259/1) and the Leverhulme Trust (ECF-2019-478). CBS acknowledges support from the Philip Leverhulme Prize, the Royal Society Wolfson Fellowship, the EPSRC advanced career fellowship EP/V029428/1, EPSRC grants EP/S026045/1 and EP/T003553/1, EP/N014588/1, EP/T017961/1, the Wellcome Innovator Awards 215733/Z/19/Z and 221633/Z/20/Z, the European Union Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No. 777826 NoMADS, the Cantab Capital Institute for the Mathematics of Information and the Alan Turing Institute.

**Availability of data and materials** The related implementation of the algorithms and the image data used in the experiment will be made available on the website <https://junqitang.com>. For the phantom image example, we use the one in the experimental section of [8], while for the lung CT image example we use an image from the Mayo Clinic Dataset [21] which is publicly available.

## Declarations

**Conflict of interest** There are no competing interests to declare.

**Ethical approval** This declaration is not applicable.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Alacaoglu, A., Fercoq, O., Cevher, V.: On the convergence of stochastic primal-dual hybrid gradient. *SIAM J. Optim.* **32**(2), 1288–1318 (2022)
- Bauschke, H.H., Combettes, P.L.: *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, vol. 408. Springer (2011)
- Bonettini, S., Benfenati, A., Ruggiero, V.: Scaling techniques for epsilon-subgradient methods. *SIAM J. Optim.* **26**(3), 1741–1772 (2016)
- Bonettini, S., Porta, F., Ruggiero, V., Zanni, L.: Variable metric techniques for forward-backward methods in imaging. *J. Comput. Appl. Math.* **385**, 113192 (2021)
- Bonettini, S., Prato, M., Rebegoldi, S.: A block coordinate variable metric linesearch based proximal gradient method. *Comput. Optim. Appl.* **71**(1), 5–52 (2018)
- Bonettini, S., Rebegoldi, S., Ruggiero, V.: Inertial variable metric techniques for the inexact forward-backward algorithm. *SIAM J. Sci. Comput.* **40**(5), A3180–A3210 (2018)
- Bonettini, S., Ruggiero, V.: On the convergence of primal-dual hybrid gradient algorithms for total variation image restoration. *J. Math. Imaging Vis.* **44**(3), 236–253 (2012)
- Chambolle, A., Ehrhardt, M.J., Richtárik, P., Schönlieb, C.-B.: Stochastic primal-dual hybrid gradient algorithm with arbitrary sampling and imaging applications. *SIAM J. Optim.* **28**(4), 2783–2808 (2018)
- Chambolle, A., Pock, T.: A first-order primal-dual algorithm for convex problems with applications to imaging. *J. Math. Imaging Vis.* **40**(1), 120–145 (2011)
- Combettes, P.L., Pesquet, J.-C.: Stochastic quasi-Fejér block-coordinate fixed point iterations with random sweeping. *SIAM J. Optim.* **25**(2), 1221–1248 (2015)
- Combettes, P.L., Vũ, B.C.: Variable metric quasi-Fejér monotonicity. *Nonlinear Anal.: Theory Methods Appl.* **78**, 17–31 (2013)
- Delplancke, C., Gurnell, M., Latz, J., Markiewicz, P.J., Schönlieb, C.-B., Ehrhardt, M.J.: Improving a stochastic algorithm for regularized PET image reconstruction. In: 2020 IEEE Nuclear Science Symposium and Medical Imaging Conference (NSS/MIC), pp. 1–3. IEEE (2020)
- Ehrhardt, M.J., Markiewicz, P., Schönlieb, C.-B.: Faster PET reconstruction with non-smooth priors by randomization and preconditioning. *Phys. Med. Biol.* **64**(22), 225019 (2019)
- Goldstein, T., Li, M., Yuan, X.: Adaptive primal-dual splitting methods for statistical learning and image processing. *Adv. Neural. Inf. Process. Syst.* **28**, 2089–2097 (2015)
- Goldstein, T., Li, M., Yuan, X., Esser, E., Baraniuk, R.: Adaptive primal-dual hybrid gradient methods for saddle-point problems. *arXiv preprint arXiv:1305.0546* (2013)
- Gutiérrez, E.B., Delplancke, C., Ehrhardt, M.J.: On the convergence and sampling of randomized primal-dual algorithms and their application to parallel MRI reconstruction. *arXiv preprint arXiv:2207.12291* (2022)
- He, B., Yuan, X.: *Convergence Analysis of Primal-dual Algorithms for Total Variation Image Restoration*. Rapport technique, Citeseer (2010)
- Malitsky, Y.: Golden ratio algorithms for variational inequalities. *Math. Program.* **184**(1), 383–410 (2020)
- Malitsky, Y., Mishchenko, K.: Adaptive gradient descent without descent. In: Daumé III, H. Singh, A., (eds) *Proceedings of the 37th International Conference on Machine Learning*, vol. 119 of *Proceedings of Machine Learning Research*, pp. 6702–6712 (2020)
- Malitsky, Y., Pock, T.: A first-order primal-dual algorithm with linesearch. *SIAM J. Optim.* **28**(1), 411–432 (2018)
- McCollough, C.: TU-FG-207A-04: overview of the low dose CT grand challenge. *Med. Phys.* **43**(6Part35), 3759–3760 (2016)
- Papoutsellis, E., Ametova, E., Delplancke, C., Fardell, G., Jørgensen, J.S., Pasca, E., Turner, M., Warr, R., Lionheart, W.R.B., Withers, P.J.: Core imaging library-part II: multichannel reconstruction for dynamic and spectral tomography. *Philos. Trans. R. Soc. A* **379**(2204), 20200193 (2021)
- Robbins, H., Siegmund, D.: A convergence theorem for non negative almost supermartingales and some applications. In: *Optimizing Methods in Statistics*, pp. 233–257. Elsevier (1971)
- Schramm, G., Holler, M.: Fast and memory-efficient reconstruction of sparse poisson data in listmode with non-smooth priors with application to time-of-flight PET. *Phys. Med. Biol.* (2022)
- Vladarean, M.-L., Malitsky, Y., Cevher, V.: A first-order primal-dual method with adaptivity to local smoothness. *Adv. Neural. Inf. Process. Syst.* **34**, 6171–6182 (2021)
- Yokota, T., Hontani, H.: An efficient method for adapting step-size parameters of primal-dual hybrid gradient method in application to total variation regularization. In: 2017 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC), pp. 973–979. IEEE (2017)
- Zdun, L., Brandt, C.: Fast MPI reconstruction with non-smooth priors by stochastic optimization and data-driven splitting. *Phys. Med. Biol.* **66**(17), 175004 (2021)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Antonin Chambolle** is a CNRS senior scientist at CEREMADE, CNRS and Paris-Dauphine University (PSL), France. He received a Ph.D. from U. Paris-Dauphine in 1993, in Mathematics applied to image analysis, and a Habilitation in 2002. He worked in SISSA, Trieste, CEREMADE, CMAP (CNRS and École Polytechnique), and now back to CEREMADE. His main research topics are related to the calculus of variations, the theoretical and numerical analysis of variational

and evolution problems involving discontinuities and boundaries, the numerical optimization, especially for non-smooth convex problems. He is part also of the INRIA team “Mokaplan” which studies numerical methods for optimal transportation problems.



**Claire Delplancke** is a researcher and engineer at EDF Research & Development (EDF Lab Paris-Saclay). After studying at ENS Cachan, she received her Ph.D. in Applied Mathematics from the University of Toulouse in 2017. She held two postdoctoral positions, first at the University of Chile, then at the University of Bath, before joining EDF in 2022, as a researcher, engineer and now project manager. Her research interests lie with stochastic algorithms and optimization for a variety of

applications: inverse problems, medical imaging, and more recently energy management.



**Matthias J. Ehrhardt** received the Diploma degree (Hons.) in industrial mathematics from the University of Bremen, Germany, in 2011, and the Ph.D. degree in medical imaging from University College London, U.K., in 2015. He held a postdoctoral position with the Cambridge Image Analysis group, Department for Applied Mathematics and Theoretical Physics, University of Cambridge, U.K., from 2016 to 2018. He moved to the University of Bath, U.K. as a Prize Fellow in

2018, where since 2021 he is a Reader at the Department of Mathematical Sciences. He is heading the Bath Imaging Group, a co-director of the Bath Centre for Mathematics and Algorithms for Data and the deputy director of the EPSRC Programme Grant on the Mathematics

of Deep Learning. His research interests include optimisation, inverse problems, computational imaging, and machine learning.



**Carola-Bibiane Schönlieb** graduated from the Institute for Mathematics, University of Salzburg (Austria) in 2004. From 2004 to 2005 she held a teaching position in Salzburg. She received her Ph.D. degree from the University of Cambridge (UK) in 2009. After one year of postdoctoral activity at the University of Göttingen (Germany), she became a Lecturer at Cambridge in 2010, promoted to Reader in 2015 and promoted to Professor in 2018. Since 2011 she is a fellow of Jesus College Cambridge.

She currently is Professor of Applied Mathematics at the University of Cambridge, where she is head of the Cambridge Image Analysis group and co-Director of the EPSRC Cambridge Mathematics of Information in Healthcare Hub. Her current research interests focus on variational methods, partial differential equations and machine learning for image analysis, image processing and inverse imaging problems.



**Junqi Tang** is an Assistant Professor in the School of Mathematics, University of Birmingham. He received the M.Sc. and Ph.D. from the Institute for Digital Communications, University of Edinburgh, U.K., in 2015 and 2019, respectively. He worked as a postdoctoral research associate with the Department of Applied Mathematics and Theoretical Physics (DAMTP), University of Cambridge before joining the University of Birmingham in 2023. His research interests include machine

learning, large-scale optimization and multi-agent systems, with applications in computational imaging and computational social science.