# Distributed Proximal Splitting Algorithms with Rates and Acceleration

Laurent Condat, Grigory Malinovsky, and Peter Richtárik

King Abdullah University of Science and Technology (KAUST), Thuwal, Saudi Arabia

Authors' final version. Published in Front. Signal Process., Jan. 2022. https://doi.org/10.3389/frsip.2021.776825

### Abstract

We analyze several generic proximal splitting algorithms well suited for large-scale convex nonsmooth optimization. We derive sublinear and linear convergence results with new rates on the function value suboptimality or distance to the solution, as well as new accelerated versions, using varying stepsizes. In addition, we propose distributed variants of these algorithms, which can be accelerated as well. While most existing results are ergodic, our nonergodic results significantly broaden our understanding of primal–dual optimization algorithms.

Keywords: convex nonsmooth optimization, proximal algorithm, splitting, convergence rate, distributed optimization

# 1 Introduction

We propose new algorithms for the generic convex optimization problem:

$$\underset{x \in \mathcal{X}}{\text{minimize}} \left\{ \Psi(x) \coloneqq \frac{1}{M} \sum_{m=1}^{M} \left( F_m(x) + H_m(K_m x) \right) + R(x) \right\}, \tag{1}$$

where  $M \geq 1$  is typically the number of parallel computing nodes in a distributed setting; the  $K_m : \mathcal{X} \to \mathcal{U}_m$  are linear operators;  $\mathcal{X}$  and  $\mathcal{U}_m$  are real Hilbert spaces (all spaces are supposed of finite dimension); R and  $H_m$  are proper, closed, convex functions with values in  $\mathbb{R} \cup \{+\infty\}$ , the proximity operators of which are easy to compute; and the  $F_m$  are convex  $L_{F_m}$ -smooth functions; that is  $\nabla F_m$  is  $L_{F_m}$ -Lipschitz continuous, for some  $L_{F_m} > 0$ .

This template problem covers most convex optimization problems met in signal and image processing, operations research, control, machine learning, and many other fields, and our goal is to propose new generic distributed algorithms able to deal with nonsmooth functions using their proximity operators, with acceleration in presence of strong convexity.

<sup>\*</sup>Corresponding author. Contact: see https://lcondat.github.io

# 1.1 Contributions

Our contributions are the following:

- 1. New algorithms: We propose the first distributed algorithms to solve (1) in whole generality, with proved convergence to an exact solution, and having the *full splitting*, or decoupling, property:  $\nabla F_m$ ,  $\operatorname{prox}_{H_m}$ ,  $K_m$  and  $K_m^*$  are applied at the *m*-th node, and the proximity operator of R is applied at the master node connected to all others. No other more complicated operation, like an inner loop or a linear system to solve, is involved.
- 2. Unified framework: The foundation of our distributed algorithms consists in two general principles, applied in a cascade, which are new contributions in themselves and could be used in other contexts:
  - (a) We show that problem (1) with M = 1, i.e. the minimization of  $F + R + H \circ K$ , can be reformulated as the minimization of  $\tilde{F} + \tilde{R} + \tilde{H}$  in a different space, with preserved smoothness and strong convexity properties. Hence, the linear operator disappears and the Davis–Yin algorithm [Davis and Yin, 2017] can be applied to this new problem. Through this lens, we recover many algorithms as particular cases of this unified framework, like the PD3O, Chambolle–Pock, Loris–Verhoeven algorithms.
  - (b) We design a non-straightforward lifting technique, so that the problem (1), with any M, is reformulated as the minimization of  $\hat{F} + \hat{R} + \hat{H} \circ \hat{K}$  in some product space.
- 3. New convergence analysis and acceleration: Even when M = 1, we improve upon the state of the art in two ways:
  - (a) For constant stepsizes, we recover existing algorithms, but we provide new, more precise, results about their convergence speed, see Theorems 1 and 5.
  - (b) With a particular strategy of varying stepsizes, we exhibit new algorithms, which are accelerated versions of them. We prove  $O(1/k^2)$  convergence rate on the last iterate, see Theorems 3 and 4, whereas current results in the literature are ergodic, e.g. Chambolle and Pock [2016b].

# 1.2 Related Work

Many estimation problems in a wide range of scientific fields can be formulated as large-scale convex optimization problems [Palomar and Eldar, 2009, Sra et al., 2011, Bach et al., 2012, Polson et al., 2015, Bubeck, 2015, Glowinski et al., 2016, Chambolle and Pock, 2016a, Stathopoulos et al., 2016, Condat, 2017b, Condat et al., 2019b]. Proximal splitting algorithms [Combettes and Pesquet, 2010, Boţ et al., 2014, Parikh and Boyd, 2014, Komodakis and Pesquet, 2015, Beck, 2017, Condat et al., 2019a] are particularly well suited to solve them; they consist of simple, easy to compute, steps that can deal with the terms in the objective function separately.

These algorithms are generally designed as sequential ones, for M = 1, and then they can be extended by lifting in product space to parallel versions, well suited to minimize  $F + R + \sum_m H_m \circ K_m$ , see for instance Condat et al., 2019a, Section 8. However, it is not straightforward to adapt lifting to the case of a finite-sum  $F = \frac{1}{M} \sum_m F_m$ , with each function  $F_m$  handled by a different node, which is of primary importance in machine learning. This generalization is one of our contributions.

There is a vast literature on distributed optimization to minimize  $\frac{1}{M} \sum_{m} F_m + R$ , with a focus on strategies based on (block-)coordinate or randomized activation, as well as replacing the gradients by cheaper stochastic estimates [Cevher et al., 2014, Richtárik and Takáč, 2014, Gorbunov et al., 2020, Salim et al., 2020]. Replacing the full gradient by a stochastic oracle in the accelerated algorithms

with varying stepsizes we propose is not straightforward; we leave this direction for future research. In any case, the generalized setting, with the smooth functions  $F_m$  at the nodes supplemented or replaced by nonsmooth functions  $H_m$ , possibly composed with linear operators, seems to have received little attention. We want to make up for that. Decentralized optimization over networks is an active research topic [Latafat et al., 2019, Alghunaim et al., 2021]. In this paper, we focus on the centralized client–server model, with one master node connected to several client nodes, working in parallel. We leave the study of decentralized algorithms for future work.

When M = 1 and K = I, where I denotes the identity, Davis and Yin [Davis and Yin, 2017] proposed an efficient algorithm, along with an extensive study of its convergence rates and possible accelerations. But the ability to handle a nontrivial K is behind the success of the Chambolle– Pock [Chambolle and Pock, 2011] or Condat–Vũ algorithms [Condat, 2013, Vũ, 2013]: they are well suited for regularized inverse problems in imaging [Chambolle and Pock, 2016a], for instance with the total variation and its variants [Condat, 2014, 2017a, Duran et al., 2016, Bredies et al., 2010]; other examples are computer vision problems [Cremers et al., 2011], overlapping group norms for sparse estimation in data science [Bach et al., 2012], and trend filtering on graphs [Wang et al., 2016]. Another prominent case is when H is an indicator function, so that the problem becomes: minimize F(x) + R(x) subject to Kx = b. If K is a gossip matrix like the minus graph Laplacian, decentralized optimization over a network can be tackled [Shi et al., 2015, Scaman et al., 2017, Salim et al., 2021].

When M = 1 and K is arbitrary, there exist algorithms to solve (1) in full generality, for example, the Combettes-Pesquet [Combettes and Pesquet, 2012], Condat-Vũ [Condat, 2013, Vũ, 2013], PD3O [Yan, 2018] and PDDY [Salim et al., 2020] algorithms. However, their convergence rates and possible accelerations are little understood. Our main contribution is to derive new convergence rates and accelerated versions of the PD3O and PDDY algorithms, and their particular cases, including Chambolle–Pock [Chambolle and Pock, 2011] and Loris–Verhoven [Loris and Verhoeven, 2011] algorithms. In order to do this, we show that these two algorithms can be viewed as instances of the Davis–Yin algorithm. This reformulation technique is inspired by the recent one of O'Connor and Vandenberghe [O'Connor and Vandenberghe, 2020]; it makes it possible to split the composition  $H \circ K$  and to derive algorithms, which call the operators  $\operatorname{prox}_H$ , K,  $K^*$  separately. This technique is fundamentally different from the one in Salim et al. [2020], showing that the PD3O and PDDY algorithms are primal-dual instances of the operator version of Davis-Yin splitting to solve monotone inclusions. Notably, we can derive convergence rates with respect to the objective function and accelerations, which is not possible with the primal-dual reformulation of Salim et al. [2020]. On the other hand, the latter encompasses the Condat–Vũ algorithm [Condat, 2013, Vũ, 2013], which is not the case of our approach. So, these are complementary interpretations.

### 1.3 Organization of the paper

In Section 2, we propose new nonstationary versions (i.e. with varying stepsizes) of several algorithms for optimization problems made of three terms, and we analyze their convergence rates. The derivation details are pushed to the end of the paper in Section 5 for ease of reading. In Section 3, we further propose distributed algorithms, which can minimize the sum of an arbitrary number of terms. Again, the derivation details are deferred to Section 6. Numerical experiments illustrating the good match between our theoretical results and practical performance are shown in Section 4.

# 2 Minimization of 3 Functions with a Linear Operator

Let us focus on the problem (1) when M = 1:

$$\underset{x \in \mathcal{X}}{\text{minimize }} \Psi(x) = F(x) + R(x) + H(Kx), \tag{2}$$

where  $K: \mathcal{X} \to \mathcal{U}$  is a linear operator,  $\mathcal{X}$  and  $\mathcal{U}$  are real Hilbert spaces, R and H are proper, closed, convex functions, and F is a convex and  $L_F$ -smooth function. We will see in Section 3 that using an adequate lifting technique, (2) can be extended to (1) and, accordingly, parallel or distributed versions of the sequential algorithms to solve (2) will be derived. That is why we first study the case M = 1. For any function G, we denote by  $\mu_G \geq 0$  some constant such that G is  $\mu_G$ -strongly convex; that is,  $G - (\mu_G/2) \| \cdot \|^2$  is convex.

The dual problem to (2) is

$$\underset{u \in \mathcal{U}}{\text{minimize}} (F+R)^* (-K^* u) + H^*(u), \tag{3}$$

where  $K^*$  is the adjoint operator of K and  $G^*$  is the convex conjugate of a function G [Bauschke and Combettes, 2017]; we recall the Moreau identity:  $\operatorname{prox}_{\tau G}(z) = z - \tau \operatorname{prox}_{G^*/\tau}(z/\tau)$  [Bauschke and Combettes, 2017]. We suppose that the following holds:

Assumption 1. There exists  $x^* \in \mathcal{X}$  such that  $0 \in \nabla F(x^*) + \partial R(x^*) + K^* \partial H(Kx^*)$ , which implies that  $x^*$  is a solution to (2); see for instance Combettes and Pesquet, 2012, Proposition 4.3 for sufficient conditions on the functions for this property to hold.

# 2.1 Deriving the Nonstationary PD3O and PDDY Algorithms

The main difficulty in (2) is the presence of the linear operator K. Indeed, if K = I, the Davis–Yin algorithm [Davis and Yin, 2017] is well suited to minimize F + R + H. Note that there is a minor mistake in the way Algorithm 3 in Davis and Yin [2017] is initialized. This is corrected here. Thus, the Davis–Yin algorithm is as follows:

Let  $(\gamma_k)_{k\in\mathbb{N}}$  be a sequence of stepsizes. Let  $x_H^0 \in \mathcal{X}$  and  $u^0 \in \mathcal{X}$ . For  $k = 0, 1, \ldots$  iterate

$$\begin{bmatrix} x^{k+1} = \operatorname{prox}_{\gamma_k R}(x_H^k + \gamma_k u^k) \\ u^{k+1} = u^k + \frac{1}{\gamma_k}(x_H^k - x^{k+1}) \\ x_H^{k+1} = \operatorname{prox}_{\gamma_{k+1} H}(x^{k+1} - \gamma_{k+1} u^{k+1} - \gamma_{k+1} \nabla F(x^{k+1})). \end{bmatrix}$$
(4)

To make this algorithm applicable to  $K \neq I$ , we reformulate the problem (2) as follows:

- 1. We choose a value  $\eta \ge ||K||^2$ ; we recommend to set  $\eta = ||K||^2$  in practice. Then there exists a real Hilbert space  $\mathcal{W}$  and a linear operator  $C: \mathcal{W} \to \mathcal{U}$  such that  $KK^* + CC^* = \eta I$ . C is not unique, for instance, we can set  $C = (\eta I KK^*)^{1/2}$ . We actually don't need to exhibit C, its existence is sufficient here and there will be no call to C in the algorithms.
- 2. Now, the problem (2) can be rewritten as:

$$\underset{x \in \mathcal{X}, w \in \mathcal{W}}{\text{minimize}} \quad \widetilde{F}(x, w) + \widetilde{R}(x, w) + \widetilde{H}(x, w), \tag{5}$$

where  $\widetilde{F}: (x, w) \mapsto F(x) + \frac{\mu_F}{2} ||w||^2$ ,  $\widetilde{R}: (x, w) \mapsto R(x) + \iota_0(w)$ , where  $\iota_0: w \mapsto \{0 \text{ if } w = 0, +\infty \text{ otherwise}\}$ , and  $\widetilde{H}: (x, w) = H(Kx + Cw)$ . Indeed, we introduce the variable w, but also the constraint that w = 0. Since  $\widetilde{F}(x, 0) = F(x)$ ,  $\widetilde{R}(x, 0) = R(x)$ ,  $\widetilde{H}(x, 0) = H(Kx)$ , the equivalence between (2) and (5) follows.

We have  $\nabla \widetilde{F}(x,w) = (\nabla F(x), \mu_F w)$ ,  $\operatorname{prox}_{\widetilde{R}}(x,w) = (\operatorname{prox}_R(x), 0)$ . Most importantly, for every  $\gamma > 0$ , we have [O'Connor and Vandenberghe, 2020]:

$$\operatorname{prox}_{\widetilde{H}^*/\gamma}(x,w) = (K^*u, C^*u), \text{ where } u = \operatorname{prox}_{H^*/(\gamma\eta)} \big( (Kx + Cw)/\eta \big).$$
(6)

# $\begin{array}{l} \label{eq:posterior} \begin{array}{c} \textbf{PD3O Algorithm} \quad (F+R+H\circ K) \\ \hline \textbf{input:} \quad (\gamma_k)_{k\in\mathbb{N}}, \ \eta \geq \|K\|^2, \ q^0 \in \mathcal{X}, \ u^0 \in \mathcal{U} \\ \textbf{for } k=0,1,\dots \ \textbf{do} \\ x^{k+1}\coloneqq \operatorname{prox}_{\gamma_k R} \big(\gamma_k(q^k-K^*u^k)\big) \\ q^{k+1}\coloneqq \frac{1}{\gamma_{k+1}}x^{k+1} - \nabla F(x^{k+1}) \\ u^{k+1}\coloneqq \operatorname{prox}_{H^*/(\gamma_{k+1}\eta)} \big(u^k \\ & + \frac{1}{\eta}K(\frac{1}{\gamma_k}x^{k+1}+q^{k+1}-q^k)\big) \\ \textbf{end for} \end{array}$

# **Davis–Yin Algorithm** (F + R + H)

 $\begin{array}{l} \mathbf{input:} \ (\gamma_k)_{k\in\mathbb{N}}, \ s^0 \in \mathcal{X} \\ \mathbf{for} \ k = 0, 1, \dots \mathbf{do} \\ x^{k+1} \coloneqq \mathrm{prox}_{\gamma_k R}(s^k) \\ x_H^{k+1} \coloneqq \mathrm{prox}_{\gamma_{k+1} H}\big((1 + \frac{\gamma_{k+1}}{\gamma_k})x^{k+1} \\ \quad - \frac{\gamma_{k+1}}{\gamma_k}s^k - \gamma_{k+1}\nabla F(x^{k+1})\big) \\ s^{k+1} \coloneqq x_H^{k+1} + \frac{\gamma_{k+1}}{\gamma_k}(s^k - x^{k+1}) \\ \mathbf{end} \ \mathbf{for} \end{array}$ 

**Chambolle–Pock Algorithm I**  $(R + H \circ K)$ 

 $\begin{array}{l} \textbf{input:} \ (\gamma_k)_{k \in \mathbb{N}}, \eta \geq \|K\|^2, x^{0} \in \mathcal{X}, u^0 \in \mathcal{U} \\ \textbf{for } k = 0, 1, \dots \textbf{do} \\ x^{k+1} \coloneqq \operatorname{prox}_{\gamma_k R} (x^k - \gamma_k K^* u^k) \\ u^{k+1} \coloneqq \operatorname{prox}_{H^*/(\gamma_{k+1}\eta)} (u^k + \frac{1}{\eta} K ((\frac{1}{\gamma_{k+1}} \\ + \frac{1}{\gamma_k}) x^{k+1} - \frac{1}{\gamma_k} x^k)) \\ \textbf{end for} \end{array}$ 

<b>Douglas–Rachford Algorithm</b> $(R+H)$
input: $(\gamma_k)_{k\in\mathbb{N}}, s^0 \in \mathcal{X}$
for $k = 0, 1,$ do
$x^{k+1} \coloneqq \operatorname{prox}_{\gamma_k R}(s^k)$
$x_{H}^{k+1} \coloneqq \operatorname{prox}_{\gamma_{k+1}H} \left( (1 + \frac{\gamma_{k+1}}{\gamma_{k}}) x^{k+1} - \frac{\gamma_{k+1}}{\gamma_{k}} s^{k} \right)$
$s^{k+1} \coloneqq x_H^{k+1} + \frac{\gamma_{k+1}}{\gamma_k} (s^k - x^{k+1})$ end for

**PDDY Algorithm**  $(F + R + H \circ K)$ 

 $\begin{array}{l} \textbf{input:} \ (\gamma_k)_{k \in \mathbb{N}}, \ \eta \geq \|K\|^2, \ x_R^0 \in \mathcal{X}, \ u^0 \in \mathcal{U} \\ \textbf{initialize:} \ p^0 \coloneqq K^* u^0 \\ \textbf{for } k = 0, 1, \dots \textbf{ do} \\ u^{k+1} \coloneqq \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( u^k + \frac{1}{\gamma_k \eta} K x_R^k \right) \\ p^{k+1} \coloneqq K^* u^{k+1} \\ x^{k+1} \coloneqq R^k_R - \gamma_k (p^{k+1} - p^k) \\ x_R^{k+1} \coloneqq \operatorname{prox}_{\gamma_{k+1}R} \left( x^{k+1} - \gamma_{k+1} \nabla F(x^{k+1}) - \gamma_{k+1} p^{k+1} \right) \\ - \gamma_{k+1} p^{k+1} \right) \\ \textbf{end for} \end{array}$ 

Loris–Verhoeven Algorithm  $(F + H \circ K)$ 

 $\begin{array}{l} \textbf{input:} \ (\gamma_k)_{k \in \mathbb{N}}, \ \eta \geq \|K\|^2, \ q^0 \in \mathcal{X}, \ u^0 \in \mathcal{U} \\ \textbf{for} \ k = 0, 1, \dots \textbf{do} \\ x^{k+1} \coloneqq \gamma_k (q^k - K^* u^k) \\ q^{k+1} \coloneqq \frac{1}{\gamma_{k+1}} x^{k+1} - \nabla F(x^{k+1}) \\ u^{k+1} \coloneqq \textbf{prox}_{H^*/(\gamma_{k+1}\eta)} (u^k \\ &\quad + \frac{1}{\eta} K(\frac{1}{\gamma_k} x^{k+1} + q^{k+1} - q^k)) \\ \textbf{end for} \end{array}$ 

 $\begin{array}{l} \hline \textbf{Chambolle-Pock Algorithm II} & (R+H\circ K) \\ \hline \textbf{input:} & (\gamma_k)_{k\in\mathbb{N}}, \eta \geq \|K\|^2, \, x_R^0 \in \mathcal{X}, \, u^0 \in \mathcal{U} \\ \textbf{for } k = 0, 1, \dots \, \textbf{do} \\ & u^{k+1} \coloneqq \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( u^k + \frac{1}{\gamma_k \eta} K x_R^k \right) \\ & x_R^{k+1} \coloneqq \operatorname{prox}_{\gamma_{k+1}R} \left( x_R^k - K^* \left( (\gamma_k \\ & + \gamma_{k+1}) u^{k+1} - \gamma_k u^k \right) \right) \\ \textbf{end for} \end{array}$ 

Forward–Backward Algorithm (F+R)input:  $(\gamma_k)_{k\in\mathbb{N}}, x_1 \in \mathcal{X},$ for k = 1, 2, ... do  $x^{k+1} \coloneqq \operatorname{prox}_{\gamma_k R} (x^k - \gamma_k \nabla F(x^k))$ end for

Note that in O'Connor and Vandenberghe [2020], the authors use  $\tilde{F}(x, w) = F(x)$ , whereas we add  $\frac{\mu_F}{2} ||w||^2$ . This difference is essential, so that  $\tilde{F}$  is  $L_F$ -smooth and  $\mu_F$ -strongly convex. Also,  $\tilde{R}$  is  $\mu_R$ -strongly convex.

Then, we can apply the Davis–Yin algorithm (4) to solve the problem (5). We set F, R, H in (4) as  $\tilde{F}$ ,  $\tilde{R}$ ,  $\tilde{H}$ , respectively. The details of the substitutions yielding the algorithms are deferred to Section 5 for the convenience of reading; most notably, whenever  $CC^*$  appears, it is replaced by  $\eta I - KK^*$ . The obtained algorithms turns out to be a nonstationary version of the PD3O

algorithm [Yan, 2018], shown above. On the other hand, if we exchange the two functions and set F, R, H in (4) as  $\tilde{F}$ ,  $\tilde{H}$ ,  $\tilde{R}$ , we obtain a different algorithm. It turns out to be a nonstationary version of the PDDY algorithm proposed recently [Salim et al., 2020], shown above too. With constant stepsizes  $\gamma_k \equiv \gamma \in (0, 2/L_F)$ , for both the PD3O and PDDY algorithms,  $x^k$  and  $u^k$  converge to some solutions  $x^*$  and  $u^*$  of (2) and (3), respectively; this result was known for  $\eta > ||K||^2$  [Yan, 2018, Salim et al., 2020] and shown for  $\eta = ||K||^2$  for the PD3O algorithm in O'Connor and Vandenberghe [2020], but convergence with  $\eta = ||K||^2$  for the PDDY algorithm, as stated in Theorem 2, is new.

Particular cases of the PD3O and PDDY algorithms, which are shown above, are the following:

- 1. If K = I and  $\eta = 1$ , the PD3O algorithm reverts to the Davis–Yin algorithm (4); the PDDY algorithm too, but with H and R exchanged in (4).
- 2. If F = 0, the PD3O and PDDY algorithms revert to the forms I and II [Condat et al., 2019a] of the Chambolle–Pock algorithm, a.k.a. Primal–Dual Hybrid Gradient algorithm [Chambolle and Pock, 2011], respectively.
- 3. If R = 0, the PD3O and PDDY algorithms revert to the Loris–Verhoeven algorithm [Loris and Verhoeven, 2011], also discovered independently as the PDFP2O [Chen et al., 2013] and PAPC [Drori et al., 2015] algorithms; see also Combettes et al. [2014], Condat et al. [2019a] for an analysis as a primal–dual forward–backward algorithm.
- 4. If F = 0 in the Davis–Yin algorithm or K = I and  $\eta = 1$  in the Chambolle–Pock algorithm, we obtain the Douglas–Rachford algorithm; it is equivalent to the ADMM, see the discussion in Condat et al. [2019a].
- 5. If H = 0, the PD3O and PDDY algorithms revert to the forward-backward algorithm, a.k.a. proximal gradient descent. The Loris-Verhoeven algorithm with K = I and  $\eta = 1$ , too.

# 2.2 Convergence Analysis

We first give convergence rates for the PD3O algorithm with constant stepsizes.

**Theorem 1** (convergence rate of the PD3O algorithm). In the PD3O algorithm, suppose that  $\gamma_k \equiv \gamma \in (0, 2/L_F)$  and  $\eta \geq ||K||^2$ . Then  $x^k$  and  $u^k$  converge to some solutions  $x^*$  and  $u^*$  of (2) and (3), respectively. In addition, suppose that H is continuous on an open ball centered at  $Kx^*$ . Then the following hold:

(i) 
$$\Psi(x^k) - \Psi(x^*) = o(1/\sqrt{k})$$

Define the weighted ergodic iterate  $\bar{x}^k = \frac{2}{k(k+1)} \sum_{i=1}^k ix^i$ , for every  $k \ge 1$ . Then

(ii) 
$$\Psi(\bar{x}^k) - \Psi(x^*) = O(1/k).$$

Furthermore, if H is L-smooth for some L > 0, we have a faster decay for the best iterate so far:

(iii) 
$$\min_{i=1,\dots,k} \Psi(x^i) - \Psi(x^*) = o(1/k).$$

Proof. The convergence of  $x^k$  follows from Davis and Yin, 2017, Theorem 2.1 and the convergence of  $u^k$  follows from the one of the variable  $u_B^k = (z^k - x_A^k)/\gamma$  in the notations of Davis and Yin [2017]. (i) follows from Davis and Yin, 2017, Theorem 3.1, using the following facts; first, in this theorem, the function corresponding to  $\tilde{H}$  is supposed to be Lipschitz-continuous on a certain ball, but since the rate is asymptotic and  $Kx^k \to Kx^*$ , it is sufficient to consider the property around  $Kx^*$ ; second, it is well known that if a convex real-valued function is continuous on a convex open set, it is Lipschitz-continuous on every compact subset of this set [Unknown author, 1972]; third, if H is continuous,  $\tilde{H}$  is continuous too. (ii) follows from Davis and Yin, 2017, Theorem 3.2 and (iii) follows from Theorem D.5 in the preprint of Davis and Yin [2017].

Theorem 1 applies to the particular cases of the PD3O algorithm, like the Loris–Verhoeven, Chambolle–Pock, Douglas–Rachford algorithms. Our results are new even for them.

**Remark 1.** We can note that the forward-backward algorithm  $x^{k+1} = \operatorname{prox}_{\gamma R}(x^k - \gamma \nabla F(x^k))$ , which is a particular case of the PD3O algorithm when H = 0, is monotonic. So, the best iterate so far is the last iterate. Hence, Theorem 1 (iii) yields  $\Psi(x^k) - \Psi(x^*) = o(1/k)$  for the forward-backward algorithm.

For the PDDY algorithm, we cannot derive a similar theorem, since  $\widetilde{R}$  is not continuous around  $(x^*, 0)$ . Still, we can establish convergence of the variables:

**Theorem 2** (convergence of the PDDY algorithm). In the PDDY algorithm, suppose that  $\gamma_k \equiv \gamma \in (0, 2/L_F)$  and  $\eta \geq ||K||^2$ . Then  $x^k$  and  $x_R^k$  both converge to some solution  $x^*$  of (2), and  $u^k$  converges to some solution  $u^*$  of (3).

Proof. The convergence of  $x^k$  and  $x^k_R$  to the same solution  $x^*$  of (2) follows from Davis and Yin, 2017, Theorem 2.1. The convergence of the variable  $u^k_B = (z^k - x^k_A)/\gamma$ , in the notations of Davis and Yin [2017], implies in our setting, according to (6), that  $K^*u^k$  and  $C^*u^k$  both converge to some elements. But since  $\eta u^k = KK^*u^k + CC^*u^k$ ,  $u^k$  converges to some element  $u^* \in \mathcal{U}$ . Finally, we have  $x^* = \operatorname{prox}_{\gamma R}(x^* - \gamma \nabla F(x^*) - \gamma K^*u^*)$ , so that  $0 \in \partial R(x^*) + \nabla F(x^*) + K^*u^*$ , and  $u^* = \operatorname{prox}_{H^*/(\gamma \eta)}(u^* + \frac{1}{\gamma \eta}Kx^*)$ , so that  $Kx^* \in (\partial H)^{-1}(u^*)$ . Hence,  $u^*$  is a solution to (3).

We now give accelerated convergence results using varying stepsizes, when F or R is strongly convex; that is,  $\mu_F + \mu_R > 0$ . In that case, we denote by  $x^*$  the unique solution to (2).

**Theorem 3** (convergence rate of the accelerated PD3O algorithm). Suppose that  $\mu_F + \mu_R > 0$ . Let  $\kappa \in (0,1)$  and  $\gamma_0 \in (0, 2(1-\kappa)/L_F)$ . Set  $\gamma_1 = \gamma_0$  and

$$\gamma_{k+1} = \frac{-\gamma_k^2 \mu_F \kappa + \gamma_k \sqrt{(\gamma_k \mu_F \kappa)^2 + 1 + 2\gamma_k \mu_R}}{1 + 2\gamma_k \mu_R}, \quad \text{for every } k \ge 1.$$
(7)

Suppose that  $\eta \ge ||K||^2$ . Then in the PD3O algorithm, there exists  $c_0 > 0$  (whose expression is given in Section 5) such that, for every  $k \ge 1$ ,

$$\|x^{k+1} - x^{\star}\|^2 \le \frac{\gamma_{k+1}^2}{1 - \gamma_{k+1}\mu_F\kappa}c_0 = O(1/k^2).$$

*Proof.* This result follows from Davis and Yin, 2017, Theorem 3.3, stated for convenience as Lemma 1 in Section 5.  $\Box$ 

Note that with the stepsize rule in (7), we have  $k \gamma_k \to 1/(\mu_F \kappa + \mu_R)$  as  $k \to +\infty$ , so that  $\gamma_k = O(1/k)$  and  $\gamma_{k+1}/\gamma_k \to 1$ . Also, when F = 0,  $L_F$  can be taken arbitrarily small, so that we can choose any  $\gamma_0 > 0$ .

Theorem 3 is new for the PD3O and Loris–Verhoeven algorithms, but has been derived in O'Connor and Vandenberghe [2020] for the Chambolle–Pock algorithm. For the forward–backward algorithm, strong convexity yields linear convergence with constant stepsizes, so this nonstationary version does not seem interesting.

Concerning the PDDY algorithm,  $\tilde{H}$  is not necessarily strongly convex, even if H is. So, we only consider the case where F is strongly convex. As a consequence of Lemma 1, we get:

**Theorem 4** (convergence rate of the accelerated PDDY algorithm). Suppose that  $\mu_F > 0$ . Let  $\kappa \in (0,1)$  and  $\gamma_0 \in (0, 2(1-\kappa)/L_F)$ . Set  $\gamma_1 = \gamma_0$  and

$$\gamma_{k+1} = -\gamma_k^2 \mu_F \kappa + \gamma_k \sqrt{(\gamma_k \mu_F \kappa)^2 + 1}, \quad \text{for every } k \ge 1.$$
(8)

Suppose that  $\eta \ge ||K||^2$ . Then in the PDDY algorithm, there exists  $c_0 > 0$  (whose expression is given in Section 5) such that, for every  $k \ge 1$ ,

$$\|x^{k+1} - x^{\star}\|^2 \le \frac{\gamma_{k+1}^2}{1 - \gamma_{k+1}\mu_F\kappa}c_0 = O(1/k^2).$$

Moreover, if  $\eta > \|K\|^2$ ,  $\|x_R^k - x^\star\|^2 = O(1/k^2)$  as well.

Finally, we consider the case where, in addition to strong convexity of F or R, H is smooth; in that case, the algorithms with constant stepsizes converge linearly; that is, as a consequence of Lemma 2, we have:

**Theorem 5** (linear convergence of the PD3O and PDDY algorithms). Suppose that  $\mu_F + \mu_R > 0$ and that H is  $L_H$ -smooth, for some  $L_H > 0$ . Let  $x^*$  and  $u^*$  be the unique solutions to (2) and (3), respectively. Suppose that  $\gamma_k \equiv \gamma \in (0, 2/L_F)$  and  $\eta \geq ||K||^2$ . Then the PD3O algorithm converges linearly: there exists  $\rho \in (0, 1]$  such that, for every  $k \in \mathbb{N}$ ,

$$||x^{k+1} - x^{\star}||^{2} \leq (1 - \rho)^{k} \Big( ||\gamma q^{0} - x^{\star} + \gamma \nabla F(x^{\star}) - \gamma K^{\star}(u^{0} - u^{\star})||^{2} \\ + \gamma^{2} \eta ||u^{0} - u^{\star}||^{2} - \gamma^{2} ||K^{\star}(u^{0} - u^{\star})||^{2} \Big).$$

The PDDY algorithm converges linearly too: there exists  $\rho \in (0,1]$  such that, for every  $k \in \mathbb{N}$ ,

$$\|x_R^{k+1} - x^{\star}\|^2 \le 4(1-\rho)^k \Big( \|x_R^0 - x^{\star} + \gamma K^*(u^0 - u^{\star})\|^2 + \gamma^2 \eta \|u^0 - u^{\star}\|^2 - \gamma^2 \|K^*(u^0 - u^{\star})\|^2 \Big).$$

Linear convergence of the other variables in the algorithms can be derived as well, see Proposition 1. Lower bounds for  $\rho$  can be derived from Theorem D.6 in the preprint version of Davis and Yin [2017]. We don't provide them, since they are not tight, as noticed in Remark D.2 of the same preprint. For instance, for the PDDY or Loris–Verhoeven algorithms with  $\mu_F > 0$ ,

$$\rho = \frac{\gamma \mu_F (2 - \gamma L_F)}{(1 + \gamma \eta L_H)^2}.$$

If H = 0, by setting  $L_H = 0$ , we get  $\rho = \gamma \mu_F (2 - \gamma L_F)$ . But then the PDDY algorithm reverts to the forward-backward algorithm, for which it is known that  $1 - \rho = (1 - \gamma \mu_F)^2$  whenever  $\gamma \leq 2/(L_F + \mu_F)$ , which corresponds to the larger value  $\rho = \gamma \mu_F (2 - \gamma \mu_F)$ .

We emphasize that linear convergence comes for free with the algorithms, if the conditions are met, without any modification. That is, there is no need to know  $\mu_F$ ,  $\mu_R$ ,  $L_H$ , since the conditions on the two parameters  $\gamma$  and  $\eta$  do not depend on these values. For the particular case of the Chambolle–Pock algorithm, as pointed out in O'Connor and Vandenberghe [2020], this is in contrast to existing linear convergence results [Chambolle and Pock, 2016a], derived for a modified version of the algorithm, which depends on these values.

# **3** Distributed Proximal Algorithms

# Distributed PD3O Algorithm

$$\begin{aligned} \text{input:} & (\gamma_k)_{k\in\mathbb{N}}, \eta \geq \|\widehat{K}\|^2, \ (\omega_m)_{m=1}^M, \\ & (q_m^0)_{m=1}^M \in \mathcal{X}^M, \ (u_m^0)_{m=1}^M \in \widehat{\mathcal{U}} \\ \text{initialize:} & a_m^0 \coloneqq q_m^0 - K_m^* u_m^0, \ m = 1...M \\ \text{for } k = 0, 1, \dots \text{ do} \\ & \text{at master, } \text{do} \\ & x^{k+1} \coloneqq \operatorname{prox}_{\gamma_k R} \left(\frac{\gamma_k}{M} \sum_{m=1}^M a_m^k\right) \\ & \text{broadcast } x^{k+1} \text{ to all nodes} \\ & \text{at all nodes, for } m = 1, \dots, M, \text{ do} \\ & q_m^{k+1} \coloneqq \frac{M\omega_m}{\gamma_{k+1}} x^{k+1} - \nabla F_m(x^{k+1}) \\ & u_m^{k+1} \coloneqq \operatorname{prox}_{M\omega_m H_m^*/(\gamma_{k+1}\eta)} \left(u_m^k \\ & + \frac{1}{\eta} K_m \left(\frac{M\omega_m}{\gamma_k} x^{k+1} + q_m^{k+1} - q_m^k\right)\right) \\ & a_m^{k+1} \coloneqq q_m^{k+1} - K_m^* u_m^{k+1} \\ & \text{transmit } a_m^{k+1} \text{ to master} \\ \text{end for} \end{aligned}$$

### Distributed Loris-Verhoeven Algorithm

 $\begin{aligned} & \text{input: } (\gamma_k)_{k\in\mathbb{N}}, \eta \geq \|\widehat{K}\|^2, \ (\omega_m)_{m=1}^M \in \widehat{\mathcal{U}} \\ & (q_m^0)_{m=1}^M \in \mathcal{X}^M, \ (u_m^0)_{m=1}^M \in \widehat{\mathcal{U}} \\ & \text{initialize: } a_m^0 \coloneqq q_m^0 - K_m^* u_m^0, \ m = 1...M \\ & \text{for } k = 0, 1, \dots \text{ do} \\ & \text{at master, } \text{do} \\ & x^{k+1} \coloneqq \frac{\gamma_k}{M} \sum_{m=1}^M a_m^k \\ & \text{broadcast } x^{k+1} \text{ to all nodes} \\ & \text{at all nodes, for } m = 1, \dots, M, \text{ do} \\ & q_m^{k+1} \coloneqq \frac{M\omega_m}{\gamma_{k+1}} x^{k+1} - \nabla F_m(x^{k+1}) \\ & u_m^{k+1} \coloneqq \operatorname{prox}_{M\omega_m H_m^*/(\gamma_{k+1}\eta)} \left( u_m^k \\ & + \frac{1}{\eta} K_m(\frac{M\omega_m}{\gamma_k} x^{k+1} + q_m^{k+1} - q_m^k) \right) \\ & a_m^{k+1} \coloneqq q_m^{k+1} - K_m^* u_m^{k+1} \\ & \operatorname{transmit} a_m^{k+1} \text{ to master} \end{aligned}$ 

### **Distributed PDDY Algorithm**

$$\begin{split} & \text{input: } (\gamma_k)_{k\in\mathbb{N}}, \eta \geq \|\widehat{K}\|^2, (\omega_m)_{m=1}^M, \\ & x_R^0 \in \mathcal{X}, (u_m^0)_{m=1}^M \in \widehat{\mathcal{U}} \\ & \text{initialize: } p_m^0 \coloneqq K_m^* u_m^0, m = 1, \dots, M \\ & \text{for } k = 0, 1, \dots \text{ do} \\ & \text{at all nodes, for } m = 1, \dots, M, \text{ do} \\ & u_m^{k+1} \coloneqq \operatorname{prox}_{M\omega_m} H_m^*/(\gamma_k \eta) \left( u_m^k \\ & + \frac{M\omega_m}{\gamma_k \eta} K_m x_R^k \right) \\ & p_m^{k+1} \coloneqq K_m^* u_m^{k+1} \\ & x_m^{k+1} \coloneqq x_R^k - \frac{\gamma_k}{M\omega_m} (p_m^{k+1} - p_m^k) \\ & a_m^k \coloneqq M\omega_m x_m^{k+1} - \gamma_{k+1} \nabla F_m(x_m^{k+1}) \\ & - \gamma_{k+1} p_m^{k+1} \\ & \text{transmit } a_m^k \text{ to master} \\ & \text{at master, } \text{ do} \\ & x_R^{k+1} \coloneqq \operatorname{prox}_{\gamma_{k+1}R} \left( \frac{1}{M} \sum_{m=1}^M a_m^k \right) \\ & \text{broadcast } x_R^{k+1} \text{ to all nodes} \\ & \text{end for} \end{split}$$



We now focus on the more general problem (1) and we derive distributed versions of the PD3O and PDDY algorithms to solve it. For this, we develop a lifting technique: we recast the minimization of  $R(x) + \frac{1}{M} \sum_{m=1}^{M} (F_m(x) + H_m(K_m x))$  as the minimization of

$$\widehat{R}(\widehat{x}) + \widehat{F}(\widehat{x}) + \widehat{H}(\widehat{K}\widehat{x}),$$

as follows. Let  $(\omega_m)_{m=1}^M$  be a sequence of positive weights, whose sum is 1; they can be used to mitigate different  $||K_m||$ , by setting  $\omega_m \propto 1/||K_m||^2$ , or different  $L_{F_m}$ , by setting  $\omega_m \propto L_{F_m}^2$ , as a rule of thumb.

We introduce the Hilbert space  $\widehat{\mathcal{X}} = \mathcal{X} \times \cdots \times \mathcal{X}$  (*M* times), endowed with the inner product

$$\langle \cdot , \cdot \rangle_{\widehat{\mathcal{X}}} : (\widehat{x}, \widehat{x}') \mapsto \sum_{m=1}^{M} \omega_m \langle x_m, x'_m \rangle,$$

# Distributed Chambolle–Pock Algorithm

$$\begin{aligned} & \text{input: } (\gamma_k)_{k \in \mathbb{N}}, \eta \geq \|\widehat{K}\|^2, (\omega_m)_{m=1}^M \\ & x_0 \in \mathcal{X}, (u_m^0)_{m=1}^M \in \widehat{\mathcal{U}} \\ & \text{initialize: } a_m^0 \coloneqq K_m^* u_m^0, m = 1, ..., M \\ & \text{for } k = 0, 1, \dots \text{ do} \\ & \text{at master, do} \\ & x^{k+1} \coloneqq \operatorname{prox}_{\gamma_k R} \left( x^k - \frac{\gamma_k}{M} \sum_{m=1}^M a_m^k \right) \\ & \text{broadcast } x^{k+1} \text{ to all nodes} \\ & \text{at all nodes, for } m = 1, \dots, M, \text{ do} \\ & u_m^{k+1} \coloneqq \operatorname{prox}_{M\omega_m H_m^*/(\gamma_{k+1}\eta)} \left( u_m^k \\ & + \frac{M\omega_m}{\eta} K_m \left( (\frac{1}{\gamma_k} + \frac{1}{\gamma_{k+1}}) x^{k+1} - \frac{1}{\gamma_k} x^k \right) \\ & a_m^{k+1} \coloneqq K_m^* u_m^{k+1} \\ & \text{transmit } a_m^{k+1} \text{ to master} \\ & \text{end for} \end{aligned}$$

# Distributed Chambolle–Pock Alg. Form II

$$\begin{array}{l} \text{input: } (\gamma_k)_{k\in\mathbb{N}}, \eta \geq \|\widehat{K}\|^2, \ (\omega_m)_{m=1}^M, \\ x_R^0 \in \mathcal{X}, \ (u_m^0)_{m=1}^M \in \widehat{\mathcal{U}} \\ \text{for } k = 0, 1, \dots \text{ do} \\ \text{at all nodes, for } m = 1, \dots, M, \text{ do} \\ u_m^{k+1} \coloneqq \operatorname{prox}_{M\omega_m H_m^*/(\gamma_k \eta)} \left( u_m^k \\ + \frac{M\omega_m}{\gamma_k \eta} K_m x_R^k \right) \\ a_m^k \coloneqq M\omega_m x_R^k - K_m^* \left( (\gamma_k + \gamma_{k+1}) u_m^{k+1} \\ - \gamma_k u_m^k \right) \\ \text{transmit } a_m^k \text{ to master} \\ \text{at master, do} \\ x_R^{k+1} \coloneqq \operatorname{prox}_{\gamma_{k+1}R} \left( \frac{1}{M} \sum_{m=1}^M a_m^k \right) \\ \text{broadcast } x_R^{k+1} \text{ to all nodes} \\ \text{end for} \end{array}$$

# Distributed Douglas–Rachford Algorithm

 $input: (\gamma_k)_{k \in \mathbb{N}}, (\omega_m)_{m=1}^M, (s_m^0)_{m=1}^M \in \mathcal{X}^M$  for  $k = 0, 1, \dots$  do at master, do  $x^{k+1} \coloneqq \operatorname{prox}_{\gamma_k R} \left( \sum_{m=1}^M \omega_m s_m^k \right)$ broadcast  $x^{k+1}$  to all nodes at all nodes, for  $m = 1, \dots, M$ , do  $x_m^{k+1} \coloneqq \operatorname{prox}_{\gamma_{k+1}H_m/(M\omega_m)}$   $((1 + \frac{\gamma_{k+1}}{\gamma_k})x^{k+1} - \frac{\gamma_{k+1}}{\gamma_k}s_m^k)$   $s_m^{k+1} \coloneqq x_m^{k+1} + \frac{\gamma_{k+1}}{\gamma_k}(s_m^k - x^{k+1})$ transmit  $s_m^{k+1}$  to master end for Distributed Forward–Backward Alg. input:  $(\gamma_k)_{k \in \mathbb{N}}, x_1 \in \mathcal{X}$ for k = 1, 2, ... do at all nodes, for m = 1, ..., M, do  $a_m^k \coloneqq \nabla F_m(x^k)$ transmit  $a_m^k$  to master at master, do  $x^{k+1} \coloneqq \operatorname{prox}_{\gamma_k R}(x^k - \frac{\gamma_k}{M} \sum_{m=1}^M a_m^k)$ broadcast  $x^{k+1}$  to all nodes end for

and the Hilbert space  $\widehat{\mathcal{U}} = \mathcal{U}_1 \times \cdots \times \mathcal{U}_M$ , endowed with the inner product

$$\langle \cdot, \cdot \rangle_{\widehat{\mathcal{U}}} : (\hat{u}, \hat{u}') \mapsto \sum_{m=1}^{M} \omega_m \langle u_m, u'_m \rangle.$$

Furthermore, we introduce  $\widehat{K} : \widehat{x} = (x_m)_{m=1}^M \in \widehat{\mathcal{X}} \mapsto (K_1 x_1, \dots, K_M x_M) \in \widehat{\mathcal{U}}$ , and the functions  $i_{=} : \widehat{x} \in \widehat{\mathcal{X}} \mapsto \{0 \text{ if } x_1 = \dots = x_M, +\infty \text{ otherwise}\}, \widehat{R} : \widehat{x} \in \widehat{\mathcal{X}} \mapsto R(x_1) + i_{=}(\widehat{x}), \widehat{H} : \widehat{u} \in \widehat{\mathcal{U}} \mapsto \frac{1}{M} \sum_{m=1}^M H_m(u_m), \text{ and } \widehat{F} : \widehat{x} \in \widehat{\mathcal{X}} \mapsto \frac{1}{M} \sum_{m=1}^M F_m(x_m).$  We have to be careful when defining the gradient and proximity operators, because of the weighted metrics; see in Section 6 for details.

Doing these substitutions in the PD3O and PDDY algorithms, we obtain the new Distributed PD3O and Distributed PDDY algorithms, shown above. Their particular cases, also shown above, are the distributed Davis–Yin algorithm when  $K_m \equiv I$  and  $\eta = 1$ , the distributed Loris–Verhoeven algorithm when R = 0, the distributed Chambolle–Pock algorithm when  $F_m \equiv 0$ , the distributed Douglas–Rachford algorithm when  $F_m \equiv 0$ ,  $K_m \equiv I$  and  $\eta = 1$ , the (classical) distributed forward–backward algorithm when  $H_m \equiv 0$ .

We can easily translate Theorems 1–5 to these distributed algorithms; the corresponding theorems are given in Section 6. In a nutshell, we obtain the same convergence results and rates with any



Figure 1: Convergence error, in log-log scale, for the experiment of image deblurring regularized with the total variation, see Section 4.1 for details.

number of nodes  $M \ge 1$  as in the non-distributed setting, for any  $\gamma_0 \in (0, 2/L_{\widehat{F}})$  and  $\eta \ge \|\widehat{K}\|^2$ , where  $L_{\widehat{F}}$  and  $\widehat{K}$  are detailed in Section 6. Hence, to our knowledge, we are the first to propose distributed proximal splitting methods with guaranteed, possibly accelerated, convergence, to minimize an arbitrary sum of smooth or nonsmooth functions, possibly composed with linear operators.

# 4 Experiments

### 4.1 Image Deblurring Regularized with Total Variation

We first consider the non-distributed problem (2), for the imaging inverse problem of deblurring, which consists in restoring an image y corrupted by blur and noise [Chambolle and Pock, 2016a]. So, we set

$$F: x \mapsto \frac{1}{2} \|Ax - y\|^2,$$

where the linear operator A corresponds to a 2-D convolution with a lowpass filter, with  $L_F = 1$ . The filter is approximately Gaussian and chosen so that F is  $\mu_F$ -strongly convex with  $\mu_F = 0.01$ . y is obtained by applying A to the classical  $256 \times 256$  Shepp–Logan phantom image, with additive Gaussian noise.  $R = i_0$  enforces nonnegativity of the pixel values.  $H \circ K$  corresponds to the classical 'isotropic' total variation (TV) [Chambolle and Pock, 2016a, Condat, 2017a], with H = 0.6 times the  $l_{1,2}$  norm and K the concatenation of vertical and horizontal finite differences.

We compare the nonaccelerated, i.e. with constant  $\gamma_k$ , and accelerated versions, with decaying  $\gamma_k$ , of the PD3O, PDDY and Condat–Vũ algorithms. We initialize the dual variables at zero and the estimate of the solution as y. We set  $\gamma_0 = 1.7$ ,  $\kappa = 0.15$ ,  $\eta = 8 \ge ||K||^2$  (except for the accelerated Condat–Vũ algorithm proposed in Chambolle and Pock [2016b], for which  $\eta = 16$  and  $\gamma = 0.5$ ).

The results are illustrated in Figure 1 (implementation in Matlab). We observe that the PD3O and PDDY algorithms have almost identical variables: the pink, red, blue curves are superimposed; we know that both algorithms are identical and revert to the Loris–Verhoeven algorithm when R = 0. Here  $R \neq 0$  but the nonnegativity constraint does not change the solution significantly, which explains the similarity of the two algorithms.

Note that  $x^k$  in the PDDY algorithm is not feasible with respect to nonnegativity, and the red curve actually shows  $F(x^k) + H(Kx^k) - \Psi(x^*)$ . In the nonaccelerated case,  $\Psi(x^k)$  decays faster



Figure 2: Convergence error, in log-log scale, for the experiment of image deblurring regularized with the smooth Huber-total-variation, so that linear convergence occurs, see Section 4.2 for details.

than O(1/k) but slower than  $O(1/k^2)$ , which is coherent with Theorem 1. The same holds for  $||x^k - x^*||^2 \leq \frac{2}{\mu_F} (\Psi(x^k) - \Psi(x^*)).$ 

The accelerated versions improve the convergence speed significantly:  $\Psi(x^k)$  and  $||x^k - x^*||^2$  decay even faster than  $O(1/k^2)$ , in line with Theorems 3 and 4. In all cases, the Condat–Vũ algorithm is outperformed. Also, there is no interest in considering the ergodic iterate instead of the last iterate, since the former converges at the same asymptotic rate as the latter, but slower.

# 4.2 Image Deblurring Regularized with Huber-TV

We consider the same deblurring experiment as before, but we make H smooth by taking the Huber function instead of the  $l_1$  norm in the total variation; that is,  $\lambda | \cdot |$  in the latter is replaced by

$$h: t \in \mathbb{R} \mapsto \begin{cases} \frac{\lambda}{2\nu} t^2 & \text{if } |t| \le \nu, \\ \lambda\left(|t| - \frac{\nu}{2}\right) & \text{otherwise,} \end{cases}$$

for some  $\nu > 0$  and  $\lambda > 0$  (set here as 0.1 and 0.6, respectively). We can also write h without branching as  $h(t) = \frac{\lambda}{2\nu} \max(\nu - |t|, 0)^2 + \lambda(|t| - \frac{\nu}{2})$ . It is known that h is  $L_h$ -smooth with  $L_h = \lambda/\nu$ . For any  $\gamma > 0$  and  $t \in \mathbb{R}$ , we have  $\operatorname{prox}_{h^*/\gamma}(t) = t/\max(|t|/\lambda, 1 + \frac{\nu}{\lambda\gamma})$ . Except for H, everything is unchanged.

The results are illustrated in Figure 2. Again, the PD3O and PDDY algorithms behave very similarly; they converge linearly, as proved in Theorem 5, and achieve machine precision in finite time.  $x^k$  in the PDDY algorithm is not feasible and  $F(x^k) + H(Kx^k) - \Psi(x^*)$  (red curve) takes negative values (not shown in log scale); so,  $x_R^k$  is the variable to study in this setting. We tested the 'accelerated' versions of the algorithms with decaying  $\gamma_k$ , but in this scenario, they are much slower and not suitable. Again, the Condat–Vũ algorithm is outperformed and the ergodic sequences converge much slower. Interestingly, the image  $x^*$  is visually the same with TV and with Huber-TV.

# 4.3 SVM with Hinge Loss

Here we consider Problem (1) in the special case with  $\mathcal{X} = \mathbb{R}^d$ , for some  $d \ge 1$ ,  $F_m \equiv 0$ , and  $K_m \equiv I$ ; that is, the problem of minimizing

$$\Psi(x) = \frac{1}{M} \sum_{m=1}^{M} H_m(x) + R(x).$$
(9)



Figure 3: Convergence error, in log-log scale, for the SVM binary classification experiment with hinge loss, see Section 4.3 for details.

In particular, to train a binary classifier, we consider the classical SVM problem with hinge loss, which has the form (9) with  $R(x) = \frac{\alpha}{2} ||x||^2$ , for some  $\alpha > 0$ , and  $H_m(x) = \max(1 - b_m a_m^T x, 0)$ , with data samples  $a_m \in \mathbb{R}^d$  and  $b_m \in \{-1, 1\}$ .

For any  $\gamma > 0$  we have  $\operatorname{prox}_{\gamma R}(x) = x/(1 + \gamma \alpha)$ . We could view the dot product  $x \mapsto b_m a_m^{\mathrm{T}} x$  as a linear operator  $K_m$ , but it is more interesting to integrate it in the function  $H_m$ . Indeed, as is perhaps not well known, the proximity operator of  $H_m$  has a closed form: for any  $\gamma > 0$ ,

$$\operatorname{prox}_{\gamma H_m} : x \in \mathbb{R}^d \mapsto x - \frac{b_m}{\eta_m} \max\left(\min(b_m a_m^{\mathrm{T}} x - 1, 0), -\eta_m \gamma\right) a_m,$$

where  $\eta_m = a_m^{\mathrm{T}} a_m = ||a_m||^2$ . Thus, we use the Distributed Douglas–Rachford algorithm, a particular case of the distributed PD3O and PDDY algorithms. Since R is  $\alpha$ -strongly convex, we also use the accelerated version of the algorithm with varying stepsizes, like in Theorem 3. We can note that in the context of Federated learning [Konečný et al., 2016, Malinovsky et al., 2020], where each mcorresponds to the smart phone or computer of a different user with its own data  $(a_m, b_m)$  stored locally, the problem is solved in a collaborative way but with preserved privacy, without the users sharing their data.

The method was implemented in Python on a single machine and tested on the dataset 'australian' from the LibSVM base [Chang and Lin, 2011], with d = 15 and M = 680. We set  $\omega_m \equiv 1/M$ ,  $\alpha = 0.1$ ,  $\gamma_0 = 0.1$ , and we used zero vectors for the initialization. The results are shown in Figure 3. Despite the oscillations, we observe that both the objective suboptimality and the squared distance to the solution converge sublinearly, with rates looking like  $o(1/\sqrt{k})$  and  $O(1/k^2)$  for the nonaccelerated and accelerated algorithms, respectively, as guaranteed by Theorems 1 and 3. The proposed accelerated version of the distributed Douglas–Rachford algorithm yields a significant speedup.

# 5 Derivation of the Algorithms

In this section, we give the details of the derivation of the PD3O and PPDY algorithms, and their particular cases, to solve:

$$\underset{x \in \mathcal{X}}{\text{minimize}} F(x) + R(x) + H(Kx),$$

with same notations and assumptions as above. Let  $\eta \geq ||K||^2$ , let  $\mathcal{W}$  be a real Hilbert space and  $C: \mathcal{W} \to \mathcal{U}$  be a linear operator, such that  $KK^* + CC^* = \eta I$ . We set  $Q: (x, w) \mapsto Kx + Cw$ . We have  $QQ^* = \eta I$ . Let  $(\gamma_k)_{k \in \mathbb{N}}$  be a sequence of positive stepsizes.

# 5.1 The Davis–Yin Algorithm

In this section, we state the results on the Davis–Yin algorithm, which we be needed to analyze the PD3O and PPDY algorithms.

The Davis–Yin algorithm to minimize the sum of 3 convex functions  $\widetilde{F} + G + J$  over a real Hilbert space  $\mathcal{Z}$  (assuming that there exists a solution  $z^*$  such that  $0 \in \nabla \widetilde{F}(z^*) + \partial G(z^*) + \partial J(z^*)$ ) is [Davis and Yin, 2017]:

Let  $z_J^0 \in \mathcal{Z}, u_G^0 \in \mathcal{Z}$ . For  $k = 0, 1, \ldots$  iterate:

$$\begin{aligned} z_{G}^{k+1} &= \operatorname{prox}_{\gamma_{k}G}(z_{J}^{k} + \gamma_{k}u_{G}^{k}) \\ u_{G}^{k+1} &= u_{G}^{k} + \frac{1}{\gamma_{k}}(z_{J}^{k} - z_{G}^{k+1}) \\ z_{J}^{k+1} &= \operatorname{prox}_{\gamma_{k+1}J}(z_{G}^{k+1} - \gamma_{k+1}u_{G}^{k+1} - \gamma_{k+1}\nabla\widetilde{F}(z_{G}^{k+1})). \end{aligned}$$
(10)

Equivalently, introducing the variable  $r^k \coloneqq z_J^k + \gamma_k u_G^k$ : let  $r^0 \in \mathcal{Z}$ . For  $k = 0, 1, \ldots$  iterate:

$$\begin{bmatrix} z_G^{k+1} = \operatorname{prox}_{\gamma_k G}(r^k) \\ z_J^{k+1} = \operatorname{prox}_{\gamma_{k+1} J} \left( (1 + \frac{\gamma_{k+1}}{\gamma_k}) z_G^{k+1} - \frac{\gamma_{k+1}}{\gamma_k} r^k - \gamma_{k+1} \nabla \widetilde{F}(z_G^{k+1}) \right) \\ r^{k+1} = z_J^{k+1} + \frac{\gamma_{k+1}}{\gamma_k} (r^k - z_G^{k+1}). \tag{11}$$

Equivalently: let  $r^0 \in \mathcal{Z}$ . For k = 0, 1, ... iterate:

$$\begin{bmatrix} z_G^{k+1} = \operatorname{prox}_{\gamma_k G}(r^k) \\ u_J^{k+1} = \operatorname{prox}_{J^*/\gamma_{k+1}} \left( \left( \frac{1}{\gamma_{k+1}} + \frac{1}{\gamma_k} \right) z_G^{k+1} - \frac{1}{\gamma_k} r^k - \nabla \widetilde{F}(z_G^{k+1}) \right) \\ r^{k+1} = z_G^{k+1} - \gamma_{k+1} \nabla \widetilde{F}(z_G^{k+1}) - \gamma_{k+1} u_J^{k+1}. \tag{12}$$

In our notations, Theorem 3.3 of Davis and Yin [2017] translates into Lemma 1 as follows; we assume that  $\widetilde{F}$  is  $L_{\widetilde{F}}$ -smooth and  $\mu_{\widetilde{F}}$ -strongly convex and that G is  $\mu_G$ -strongly convex, for some  $L_{\widetilde{F}} > 0, \ \mu_{\widetilde{F}} \ge 0, \ \mu_G \ge 0.$ 

**Lemma 1** (accelerated Davis–Yin algorithm). Suppose that  $\mu_{\widetilde{F}} + \mu_G > 0$ . Let  $z^*$  be the unique minimizer of  $\widetilde{F} + G + J$ ; that is,  $0 \in \nabla \widetilde{F}(z^*) + \partial G(z^*) + \partial J(z^*)$ . Let  $u_G^*$  be such that  $u_G^* \in \partial G(z^*)$  and  $0 \in \nabla \widetilde{F}(z^*) + \partial J(z^*) + u_G^*$ . Let  $\kappa \in (0, 1)$  and  $\gamma_0 \in (0, 2(1 - \kappa)/L_{\widetilde{F}})$ . Set  $\gamma_1 = \gamma_0$  and

$$\gamma_{k+1} = \frac{-\gamma_k^2 \mu_{\widetilde{F}} \kappa + \gamma_k \sqrt{(\gamma_k \mu_{\widetilde{F}} \kappa)^2 + 1 + 2\gamma_k \mu_G}}{1 + 2\gamma_k \mu_G}, \quad \text{for every } k \ge 1.$$

Then, for every  $k \geq 1$ ,

$$\|z_G^{k+1} - z^{\star}\|^2 \le \frac{\gamma_{k+1}^2}{1 - \gamma_{k+1}\mu_{\widetilde{F}}\kappa}c_0 = O(1/k^2)$$

where

$$c_0 = \frac{1 - \gamma_0 \mu_{\widetilde{F}} \kappa}{\gamma_0^2} \|z_G^1 - z^\star\|^2 + \|u_G^1 - u_G^\star\|^2.$$

Note that  $u_G^1 = (r^0 - z_G^1) / \gamma_0$ .

Linear convergence occurs in the following conditions, according to Theorem D.6 in the preprint version of Davis and Yin [2017], which translates into Lemma 2 as follows. We assume that  $\tilde{F}$  is  $L_{\tilde{F}}$ -smooth and  $\mu_{\tilde{F}}$ -strongly convex, G is  $\mu_G$ -strongly convex, and J is  $\mu_J$ -strongly convex, for some  $L_{\tilde{F}} > 0, \ \mu_{\tilde{F}} \ge 0, \ \mu_G \ge 0, \ \mu_J \ge 0$ . We consider constant stepsizes  $\gamma_k \equiv \gamma$ , for some  $\gamma \in (0, 2/L_{\tilde{F}})$ .

**Lemma 2** (linear convergence of the Davis–Yin algorithm). Suppose that  $\mu_{\widetilde{F}} + \mu_G + \mu_J > 0$  and that G is  $L_G$ -smooth, for some  $L_G > 0$ , or J is  $L_J$ -smooth, for some  $L_J > 0$ . Let  $z^*$  be the unique minimizer of  $\widetilde{F} + G + J$ ; that is,  $0 \in \nabla \widetilde{F}(z^*) + \partial G(z^*) + \partial J(z^*)$ . The dual problem of minimizing  $(\widetilde{F} + J)^*(-u) + G^*(u)$  over  $u \in \mathbb{Z}$  is strongly convex too; let  $u_G^*$  be its unique solution. We have  $u_G^* \in \partial G(z^*)$  and  $0 \in \nabla \widetilde{F}(z^*) + \partial J(z^*) + u_G^*$ . Set  $r^* = z^* + \gamma u_G^*$ . Then, the Davis–Yin algorithm (11) converges linearly: there exists  $\rho \in (0, 1]$  such that, for every  $k \in \mathbb{N}$ ,

$$\|r^{k} - r^{\star}\|^{2} \le (1 - \rho)^{k} \|r^{0} - r^{\star}\|^{2}.$$
(13)

Loose lower bounds for  $\rho$  are given in Davis and Yin, 2017, Theorem D.6.

We have the following corollary of Lemma 2:

**Proposition 1** (linear convergence of the other variables in the Davis–Yin algorithm). In the same conditions and notations as in Lemma 2, we have, for every  $k \in \mathbb{N}$ ,

$$\begin{aligned} \|z_G^{k+1} - z^{\star}\|^2 &\leq (1-\rho)^k \|r^0 - r^{\star}\|^2 \\ \|z_J^{k+1} - z^{\star}\|^2 &\leq 4(1-\rho)^k \|r^0 - r^{\star}\|^2. \end{aligned}$$
(14)

Also, in the form (12) of the algorithm,

$$\|u_J^{k+1} + u_G^{\star} + \nabla \widetilde{F}(z^{\star})\|^2 \le \frac{4}{\gamma^2} (1-\rho)^k \|r^0 - r^{\star}\|^2$$

and, in the form (10) of the algorithm,

$$||u_G^{k+1} - u_G^{\star}||^2 \le \frac{1}{\gamma^2} (1 - \rho)^k ||r^0 - r^{\star}||^2.$$

*Proof.* Let  $k \in \mathbb{N}$ . By nonexpansiveness of the proximity operator, in view of the first line in (11), we have  $||z_G^{k+1} - z^*|| \le ||r^k - r^*||$ , so that (14) follows from (13). In addition, in view of the second line in (11), we have

$$\begin{aligned} \|z_J^{k+1} - z^{\star}\|^2 &\leq \|2(z_G^{k+1} - z^{\star}) - (r^k - r^{\star}) - \gamma(\nabla \widetilde{F}(z_G^{k+1}) - \nabla \widetilde{F}(z^{\star}))\|^2 \\ &= \|(z_G^{k+1} - z^{\star}) - (r^k - r^{\star}) + (I - \gamma \nabla \widetilde{F})(z_G^{k+1}) - (I - \gamma \nabla \widetilde{F})(z^{\star})\|^2 \\ &= \|(I - \operatorname{prox}_{\gamma G})(r^k) - (I - \operatorname{prox}_{\gamma G})(r^{\star}) + (I - \gamma \nabla \widetilde{F})(z_G^{k+1}) - (I - \gamma \nabla \widetilde{F})(z^{\star})\|^2 \end{aligned}$$

and, by nonexpansiveness of  $I - \operatorname{prox}_{\gamma G}$  and  $I - \gamma \nabla \widetilde{F}$ ,

$$\begin{aligned} \|z_J^{k+1} - z^{\star}\|^2 &\leq \left(\|r^k - r^{\star}\| + \|z_G^{k+1} - z^{\star}\|\right)^2 \\ &\leq 4\|r^k - r^{\star}\|^2. \end{aligned}$$

Using the same arguments, in view of the second line in (12),

$$\begin{aligned} \|u_J^{k+1} + u_G^{\star} + \nabla \widetilde{F}(z^{\star})\|^2 &\leq \frac{1}{\gamma^2} \big( \|r^k - r^{\star}\| + \|z_G^{k+1} - z^{\star}\| \big)^2 \\ &\leq \frac{4}{\gamma^2} \|r^k - r^{\star}\|^2. \end{aligned}$$

Finally, as visible in the first line of (16), since  $r^k = z_J^k + \gamma_k u_G^k$ , and using the Moreau identity, we have  $u_G^{k+1} = \operatorname{prox}_{G^*/\gamma}(\frac{1}{\gamma}z_J^k + u_G^k) = \operatorname{prox}_{G^*/\gamma}(\frac{1}{\gamma}r^k)$ , so that

$$||u_G^{k+1} - u_G^{\star}||^2 \le \frac{1}{\gamma^2} ||r^k - r^{\star}||^2.$$

# 5.2 The PD3O Algorithm

We set  $\mathcal{Z} = \mathcal{X} \times \mathcal{W}$ ,  $\widetilde{F}$ ,  $G = \widetilde{R}$ ,  $J = \widetilde{H}$ , as defined in Section 2. Doing the substitutions in (12), we get the algorithm:

Let  $s^0 \in \mathcal{X}$  and  $r_w^0 \in \mathcal{W}$ . For k = 0, 1, ... iterate:

$$\begin{aligned} x^{k+1} &= \operatorname{prox}_{\gamma_k R}(s^k) \\ u^{k+1} &= \operatorname{prox}_{H^*/(\gamma_{k+1}\eta)} \left( K \left( \left( \frac{1}{\gamma_{k+1}} + \frac{1}{\gamma_k} \right) x^{k+1} - \frac{1}{\gamma_k} s^k - \nabla F(x^{k+1}) \right) / \eta - C r_w^k / (\gamma_k \eta) \right) \\ s^{k+1} &= x^{k+1} - \gamma_{k+1} \nabla F(x^{k+1}) - \gamma_{k+1} K^* u^{k+1} \\ r_w^{k+1} &= -\gamma_{k+1} C^* u^{k+1}. \end{aligned}$$

We can remove the variable  $r_w$  and the algorithm becomes: Let  $s^0 \in \mathcal{X}$  and  $u^0 \in \mathcal{U}$ . For k = 0, 1, ... iterate:

$$x^{k+1} = \operatorname{prox}_{\gamma_k R}(s^k)$$
  

$$u^{k+1} = \operatorname{prox}_{H^*/(\gamma_{k+1}\eta)} \left( \frac{1}{\eta} K \left( \left( \frac{1}{\gamma_{k+1}} + \frac{1}{\gamma_k} \right) x^{k+1} - \frac{1}{\gamma_k} s^k - \nabla F(x^{k+1}) \right) + \frac{1}{\eta} C C^* u^k \right)$$
  

$$s^{k+1} = x^{k+1} - \gamma_{k+1} \nabla F(x^{k+1}) - \gamma_{k+1} K^* u^{k+1}.$$

After replacing  $CC^*$  by  $\eta I - KK^*$ , the iteration becomes:

$$\left[ \begin{array}{l} x^{k+1} = \operatorname{prox}_{\gamma_k R}(s^k) \\ u^{k+1} = \operatorname{prox}_{H^*/(\gamma_{k+1}\eta)} \left( u^k + \frac{1}{\eta} K \left( (\frac{1}{\gamma_{k+1}} + \frac{1}{\gamma_k}) x^{k+1} - \frac{1}{\gamma_k} s^k - \nabla F(x^{k+1}) - K^* u^k \right) \right) \\ s^{k+1} = x^{k+1} - \gamma_{k+1} \nabla F(x^{k+1}) - \gamma_{k+1} K^* u^{k+1}. \end{array} \right]$$

We can change the variables, so that only one call to  $\nabla F$  and  $K^*$  appears, which yields the algorithm: Let  $q^0 \in \mathcal{X}$  and  $u^0 \in \mathcal{U}$ . For k = 0, 1, ... iterate:

$$\begin{bmatrix} x^{k+1} = \operatorname{prox}_{\gamma_k R} \left( \gamma_k (q^k - K^* u^k) \right) \\ q^{k+1} = \frac{1}{\gamma_{k+1}} x^{k+1} - \nabla F(x^{k+1}) \\ u^{k+1} = \operatorname{prox}_{H^*/(\gamma_{k+1}\eta)} \left( u^k + \frac{1}{\eta} K(\frac{1}{\gamma_k} x^{k+1} + q^{k+1} - q^k) \right). \end{bmatrix}$$

When  $\gamma_k \equiv \gamma$  is constant, we recover the PD3O algorithm [Yan, 2018].

To derive Theorem 3 from Lemma 1, we simply have to notice that the variable  $z_G^{k+1}$  in the latter corresponds to the pair  $(x^{k+1}, 0)$ . Also, in the conditions of Theorem 3, let  $u^*$  be any solution of (3); that is,  $u^* \in \partial H(Kx^*)$  and  $0 \in \partial R(x^*) + \nabla F(x^*) + K^*u^*$ . Then the constant  $c_0$  is

$$c_0 = \frac{1 - \gamma_0 \mu_F \kappa}{\gamma_0^2} \|x^1 - x^\star\|^2 + \|q^0 - \frac{1}{\gamma_0} x^1 - K^* (u^0 - u^\star) + \nabla F(x^\star)\|^2 + \eta \|u^0 - u^\star\|^2 - \|K^* (u^0 - u^\star)\|^2.$$

If K = I and  $\eta = 1$ , the PD3O algorithm reverts to the Davis–Yin algorithm, as given in (4). In the conditions of Theorem 3, let  $u^*$  be any solution of (3); that is,  $u^* \in \partial H(x^*)$  and  $0 \in \partial R(x^*) + \nabla F(x^*) + u^*$ . Then the constant  $c_0$  is

$$c_0 = \frac{1 - \gamma_0 \mu_F \kappa}{\gamma_0^2} \|x^1 - x^\star\|^2 + \|\frac{1}{\gamma_0} (s^0 - x^1) + u^\star + \nabla F(x^\star)\|^2.$$
(15)

### 5.3 The PDDY Algorithm

The PDDY algorithm is obtained like the PD3O algorithm from the David–Yin algorithm, but after swapping the roles of  $\tilde{H}$  and  $\tilde{R}$ .

To obtain the PDDY algorithm, starting from (10), let us first write the Davis–Yin algorithm as: Let  $z_J^0 \in \mathcal{Z}$  and  $u_G^0 \in \mathcal{Z}$ . For k = 0, 1, ... iterate:

$$\begin{aligned} u_{G}^{k+1} &= \operatorname{prox}_{G^{*}/\gamma_{k}}(\frac{1}{\gamma_{k}}z_{J}^{k} + u_{G}^{k}) \\ z_{G}^{k+1} &= z_{J}^{k} - \gamma_{k}(u_{G}^{k+1} - u_{G}^{k}) \\ z_{J}^{k+1} &= \operatorname{prox}_{\gamma_{k+1}J}(z_{G}^{k+1} - \gamma_{k+1}\nabla\widetilde{F}(z_{G}^{k+1}) - \gamma_{k+1}u_{G}^{k+1}). \end{aligned}$$

Equivalently: Let  $r^0 \in \mathcal{Z}$ . For k = 0, 1, ... iterate:

$$\begin{bmatrix}
 u_{G}^{k+1} = \operatorname{prox}_{G^{*}/\gamma_{k}}(r^{k}/\gamma_{k}) \\
 z_{G}^{k+1} = r^{k} - \gamma_{k}u_{G}^{k+1} \\
 z_{J}^{k+1} = \operatorname{prox}_{\gamma_{k+1}J}(z_{G}^{k+1} - \gamma_{k+1}\nabla\widetilde{F}(z_{G}^{k+1}) - \gamma_{k+1}u_{G}^{k+1}) \\
 r^{k+1} = z_{J}^{k+1} + \gamma_{k+1}u_{G}^{k+1}.
\end{cases}$$
(16)

We set  $\mathcal{Z} = \mathcal{X} \times \mathcal{W}$ ,  $\tilde{F}$ ,  $G = \tilde{H}$ ,  $J = \tilde{R}$ , as defined in Section 2. Doing the substitutions in (16), we get the algorithm: Let  $r_x^0 \in \mathcal{X}$ ,  $r_w^0 \in \mathcal{W}$ . For  $k = 0, 1, \ldots$  iterate:

$$u^{k+1} = \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( (Kr_x^k + Cr_w^k)/(\gamma_k \eta) \right)$$
  

$$x^{k+1} = r_x^k - \gamma_k K^* u^{k+1}$$
  

$$x_R^{k+1} = \operatorname{prox}_{\gamma_{k+1}R} \left( x^{k+1} - \gamma_{k+1} \nabla F(x^{k+1}) - \gamma_{k+1} K^* u^{k+1} \right)$$
  

$$r_x^{k+1} = x_R^{k+1} + \gamma_{k+1} K^* u^{k+1}$$
  

$$r_w^{k+1} = \gamma_{k+1} C^* u^{k+1}.$$

We can remove the variable  $r_w$  and rename  $r_x$  as s:

$$u^{k+1} = \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( Ks^k/(\gamma_k \eta) + CC^* u^k/\eta \right)$$
  

$$x^{k+1} = s^k - \gamma_k K^* u^{k+1}$$
  

$$x^{k+1}_R = \operatorname{prox}_{\gamma_{k+1}R} \left( x^{k+1} - \gamma_{k+1} \nabla F(x^{k+1}) - \gamma_{k+1} K^* u^{k+1} \right)$$
  

$$s^{k+1} = x^{k+1}_R + \gamma_{k+1} K^* u^{k+1}.$$

The algorithm becomes: Let  $s^0 \in \mathcal{X}$ ,  $u^0 \in \mathcal{U}$ . For k = 0, 1, ... iterate:

$$u^{k+1} = \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( u^k + K(s^k/\gamma_k - K^*u^k)/\eta \right)$$
  

$$x^{k+1} = s^k - \gamma_k K^* u^{k+1}$$
  

$$x^{k+1}_R = \operatorname{prox}_{\gamma_{k+1}R} \left( x^{k+1} - \gamma_{k+1} \nabla F(x^{k+1}) - \gamma_{k+1} K^* u^{k+1} \right)$$
  

$$s^{k+1} = x^{k+1}_R + \gamma_{k+1} K^* u^{k+1}.$$

$$\begin{bmatrix} u_{R}^{k} - prox_{\gamma_{k+1}R}(u^{k} - \gamma_{k+1}\nabla T(u^{k})) & \gamma_{k+1}R^{k}u^{k} \\ s^{k+1} = x_{R}^{k+1} + \gamma_{k+1}K^{*}u^{k+1}. \end{bmatrix}$$
  
Equivalently: Let  $x_{R}^{0} \in \mathcal{X}, u^{0} \in \mathcal{U}$ . For  $k = 0, 1, ...$  iterate:  
$$\begin{bmatrix} u^{k+1} = prox_{H^{*}/(\gamma_{k}\eta)} (u^{k} + Kx_{R}^{k}/(\gamma_{k}\eta)) \\ x^{k+1} = x_{R}^{k} - \gamma_{k}K^{*}(u^{k+1} - u^{k}) \\ x_{R}^{k+1} = prox_{\gamma_{k+1}R} (x^{k+1} - \gamma_{k+1}\nabla F(x^{k+1}) - \gamma_{k+1}K^{*}u^{k+1}). \end{bmatrix}$$

We can write the algorithm with only one call of  $K^*$  per iteration by introducing an additional variable p: Let  $x_R^0 \in \mathcal{X}$ ,  $u^0 \in \mathcal{U}$ . Set  $p^0 = K^* u^0$ . For k = 0, 1, ... iterate:

$$\begin{cases} u^{k+1} = \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( u^k + \frac{1}{\gamma_k \eta} K x_R^k \right) \\ p^{k+1} = K^* u^{k+1} \\ x^{k+1} = x_R^k - \gamma_k (p^{k+1} - p^k) \\ x_R^{k+1} = \operatorname{prox}_{\gamma_{k+1}R} \left( x^{k+1} - \gamma_{k+1} \nabla F(x^{k+1}) - \gamma_{k+1} p^{k+1} \right) \end{cases}$$

When  $\gamma_k \equiv \gamma$  is constant, we recover the PDDY algorithm [Salim et al., 2020].

Let us now derive Theorem 4 from Lemma 1. The variable  $z_G^{k+1}$  in the latter corresponds to the pair  $(x^{k+1}, \gamma_k C^*(u^k - u^{k+1}))$ , so that  $||z_G^{k+1} - z^*||^2$  becomes

$$\begin{aligned} \|x^{k+1} - x^{\star}\|^{2} + \|\gamma_{k}C^{*}(u^{k} - u^{k+1})\|^{2} &= \|x^{k+1} - x^{\star}\|^{2} + \gamma_{k}^{2}\langle CC^{*}(u^{k} - u^{k+1}), u^{k} - u^{k+1}\rangle \\ &= \|x^{k+1} - x^{\star}\|^{2} + \gamma_{k}^{2}\langle (\eta I - KK^{*})(u^{k} - u^{k+1}), u^{k} - u^{k+1}\rangle \\ &= \|x^{k+1} - x^{\star}\|^{2} + \gamma_{k}^{2}\eta\|u^{k} - u^{k+1}\|^{2} - \gamma_{k}^{2}\|K^{*}(u^{k} - u^{k+1})\|^{2}. \end{aligned}$$

$$(17)$$

Therefore, in the conditions of Theorem 4, let  $u^*$  be any solution of (3); that is,  $u^* \in \partial H(Kx^*)$  and  $0 \in \partial R(x^*) + \nabla F(x^*) + K^*u^*$ . Then the constant  $c_0$  is

$$c_0 = \frac{1 - \gamma_0 \mu_F \kappa}{\gamma_0^2} \Big( \|x^1 - x^\star\|^2 + \gamma_0^2 \eta \|u^1 - u^0\|^2 - \gamma_0^2 \|K^*(u^1 - u^0)\|^2 \Big) + \eta \|u^1 - u^\star\|^2.$$

The last statement in Theorem 4 is obtained as follows. First, for every  $k \ge 1$ ,  $x_R^k = x^{k+1} - \gamma_k K^*(u^k - u^{k+1})$ , so that  $\|x_R^k - x^*\|^2 \le 2\|x^{k+1} - x^*\|^2 + 2\|K\|^2\|\gamma_k(u^k - u^{k+1})\|^2$ . Second, from (17),  $\|x^{k+1} - x^*\|^2 = O(1/k^2)$  and  $(\eta - \|K\|^2)\|\gamma_k(u^k - u^{k+1})\|^2 \le \gamma_k^2 \langle (\eta I - KK^*)(u^k - u^{k+1}), u^k - u^{k+1} \rangle = O(1/k^2)$ . So, assuming that  $\eta > \|K\|^2$ ,  $\|\gamma_k(u^k - u^{k+1})\|^2 = O(1/k^2)$ . Hence,  $\|x_R^k - x^*\|^2 = O(1/k^2)$ .

If K = I and  $\eta = 1$ , the PDDY algorithm reverts to the Davis–Yin algorithm, as given in (4), but with R and H exchanged. In the conditions of Theorem 4, let  $u^*$  be any solution of (3); that is,  $u^* \in \partial H(x^*)$  and  $0 \in \partial R(x^*) + \nabla F(x^*) + u^*$ . Then the constant  $c_0$  is

$$c_0 = \frac{1 - \gamma_0 \mu_F \kappa}{\gamma_0^2} \|x^1 - x^\star\|^2 + \|\frac{1}{\gamma_0}(s^0 - x^1) - u^\star\|^2.$$

This is the same value as in (15), corresponding to the Davis–Yin algorithm, viewed as the PD3O algorithm, with R and H exchanged. Indeed,  $u^*$  is defined differently in both cases; that is, with the exchange,  $u^* \in \partial R(x^*)$  in (15).

# 5.4 R = 0: The Loris–Verhoeven Algorithm

If R = 0, the PD3O algorithm becomes: Let  $q^0 \in \mathcal{X}$  and  $u^0 \in \mathcal{U}$ . For k = 0, 1, ... iterate:

$$\begin{bmatrix}
x^{k+1} = \gamma_k (q^k - K^* u^k) \\
q^{k+1} = \frac{1}{\gamma_{k+1}} x^{k+1} - \nabla F(x^{k+1}) \\
u^{k+1} = \operatorname{prox}_{H^*/(\gamma_{k+1}\eta)} \left( u^k + \frac{1}{\eta} K(\frac{1}{\gamma_k} x^{k+1} + q^{k+1} - q^k) \right),$$
(18)

whereas the PDDY algorithm becomes: Let  $x_R^0 \in \mathcal{X}$ ,  $u^0 \in \mathcal{U}$ . Set  $p^0 = K^* u^0$ . For k = 0, 1, ... iterate:

$$\begin{bmatrix} u^{k+1} = \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( u^k + \frac{1}{\gamma_k \eta} K x_R^k \right) \\ p^{k+1} = K^* u^{k+1} \\ x^{k+1} = x_R^k - \gamma_k (p^{k+1} - p^k) \\ x_R^{k+1} = x^{k+1} - \gamma_{k+1} \nabla F(x^{k+1}) - \gamma_{k+1} p^{k+1}. \end{bmatrix}$$

Equivalently,

$$u^{k+1} = \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( u^k + \frac{1}{\gamma_k \eta} K(x^k - \gamma_k \nabla F(x^k) - \gamma_k K^* u^k) \right)$$
$$x^{k+1} = x^k - \gamma_k \nabla F(x^k) - \gamma_k K^* u^{k+1},$$

or:

$$\begin{bmatrix} q^{k+1} = \frac{1}{\gamma_k} x^k - \nabla F(x^k) \\ u^{k+1} = \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( u^k + \frac{1}{\gamma_k \eta} K(\gamma_k q^{k+1} - \gamma_k K^* u^k) \right) \\ x^{k+1} = \gamma_k q^{k+1} - \gamma_k K^* u^{k+1}, \end{bmatrix}$$

which is equivalent to (18). So, when R = 0, both the PD3O and PPDY revert to an algorithm which, for  $\gamma_k \equiv \gamma$ , is the Loris–Verhoeven algorithm [Loris and Verhoeven, 2011, Combettes et al., 2014, Condat et al., 2019a].

Let  $u^*$  be any solution of (3); that is,  $u^* \in \partial H(Kx^*)$  and  $0 \in \nabla F(x^*) + K^*u^*$ . In the conditions of Theorem 3,  $c_0$  is:

$$c_0 = \frac{1 - \gamma_0 \mu_F \kappa}{\gamma_0^2} \|x^1 - x^\star\|^2 + \|q^0 - \frac{1}{\gamma_0} x^1 - K^* (u^0 - u^\star) + \nabla F(x^\star)\|^2 + \eta \|u^0 - u^\star\|^2 - \|K^* (u^0 - u^\star)\|^2.$$

On the other hand, in Theorem 4,

$$c_0 = \frac{1 - \gamma_0 \mu_F \kappa}{\gamma_0^2} \Big( \|x^1 - x^\star\|^2 + \gamma_0^2 \eta \|u^1 - u^0\|^2 - \gamma_0^2 \|K^*(u^1 - u^0)\|^2 \Big) + \eta \|u^1 - u^\star\|^2.$$

It is not clear how these two values compare to each other. They are both valid, in any case.

# 5.5 F = 0: The Chambolle–Pock and Douglas–Rachford Algorithms

If F = 0, the PD3O algorithms reverts to: Let  $x^0 \in \mathcal{X}$  and  $u^0 \in \mathcal{U}$ . For k = 0, 1, ... iterate:

For  $\gamma_k \equiv \gamma$ , this is the form I [Condat et al., 2019a] of the Chambolle–Pock algorithm [Chambolle and Pock, 2011].

In the conditions of Theorem 3, let  $u^*$  be any solution of (3); that is,  $u^* \in \partial H(Kx^*)$  and  $0 \in \partial R(x^*) + K^*u^*$ . Then the constant  $c_0$  is

$$c_0 = \frac{1}{\gamma_0^2} \|x^1 - x^\star\|^2 + \|\frac{1}{\gamma_0}(x^0 - x^1) - K^*(u^0 - u^\star)\|^2 + \eta \|u^0 - u^\star\|^2 - \|K^*(u^0 - u^\star)\|^2.$$

On the other hand, if F = 0, the PDDY algorithm reverts to: Let  $x_R^0 \in \mathcal{X}$ ,  $u^0 \in \mathcal{U}$ . Set  $p^0 = K^* u^0$ . For  $k = 0, 1, \ldots$  iterate:

$$\begin{bmatrix} u^{k+1} = \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( u^k + \frac{1}{\gamma_k \eta} K x_R^k \right) \\ p^{k+1} = K^* u^{k+1} \\ x^{k+1} = x_R^k - \gamma_k (p^{k+1} - p^k) \\ x_R^{k+1} = \operatorname{prox}_{\gamma_{k+1}R} \left( x^{k+1} - \gamma_{k+1} p^{k+1} \right), \end{bmatrix}$$

which can be simplified as: Let  $x_R^0 \in \mathcal{X}$ ,  $u^0 \in \mathcal{U}$ . For k = 0, 1, ... iterate:

$$\begin{bmatrix} u^{k+1} = \operatorname{prox}_{H^*/(\gamma_k \eta)} \left( u^k + \frac{1}{\gamma_k \eta} K x_R^k \right) \\ x_R^{k+1} = \operatorname{prox}_{\gamma_{k+1}R} \left( x_R^k - K^* \left( (\gamma_k + \gamma_{k+1}) u^{k+1} - \gamma_k u^k \right) \right),$$

knowing that we can retrieve the variable  $x^k$  as  $x^{k+1} = x_R^k - \gamma_k K^* (u^{k+1} - u^k)$ .

For  $\gamma_k \equiv \gamma$ , this is the form II [Condat et al., 2019a] of the Chambolle–Pock algorithm [Chambolle and Pock, 2011].

Note that with constant stepsizes, the Chambolle–Pock form II can be viewed as the form I applied to the dual problem. This interpretation does not hold with varying stepsizes as in Theorem 3: the stepsize playing the role of  $\gamma_k$  would be  $1/(\gamma_k \eta)$ , which tends to  $+\infty$  instead of 0, so that the theorem does not apply.

Note, also, that Theorem 4 does not apply, since F = 0 is not strongly convex. Finally, if the accelerated Chambolle–Pock algorithm form I is applied to the dual problem, our results do not

guarantee convergence of the primal variable  $x^k$  to a solution. So, we cannot derive an accelerated Chambolle–Pock algorithm form II.

If K = I,  $\mathcal{U} = \mathcal{X}$  and  $\eta = 1$ , the Chambolle-Pock algorithm form I becomes the Douglas–Rachford algorithm: Let  $x^0 \in \mathcal{X}$  and  $u^0 \in \mathcal{X}$ . For  $k = 0, 1, \ldots$  iterate:

$$x^{k+1} = \operatorname{prox}_{\gamma_k R} (x^k - \gamma_k u^k) u^{k+1} = \operatorname{prox}_{H^*/\gamma_{k+1}} (u^k + (\frac{1}{\gamma_{k+1}} + \frac{1}{\gamma_k}) x^{k+1} - \frac{1}{\gamma_k} x^k).$$

We can rewrite the algorithm using only the meta-variable  $s^k = x^k - \gamma_k u^k$ : Let  $s^0 \in \mathcal{X}$ . For  $k = 0, 1, \ldots$  iterate:

$$\begin{aligned} x^{k+1} &= \operatorname{prox}_{\gamma_k R}(s^k) \\ u^{k+1} &= \operatorname{prox}_{H^*/\gamma_{k+1}} \left( \left( \frac{1}{\gamma_{k+1}} + \frac{1}{\gamma_k} \right) x^{k+1} - \frac{1}{\gamma_k} s^k \right) \\ s^{k+1} &= x^{k+1} - \gamma_{k+1} u^{k+1}. \end{aligned}$$

Using the Moreau identity, we obtain: Let  $s^0 \in \mathcal{X}$ . For k = 0, 1, ... iterate:

$$\begin{bmatrix} x^{k+1} = \operatorname{prox}_{\gamma_k R}(s^k) \\ x_H^{k+1} = \operatorname{prox}_{\gamma_{k+1} H} \left( (1 + \frac{\gamma_{k+1}}{\gamma_k}) x^{k+1} - \frac{\gamma_{k+1}}{\gamma_k} s^k \right) \\ s^{k+1} = x_H^{k+1} + \frac{\gamma_{k+1}}{\gamma_k} (s^k - x^{k+1}), \tag{19}$$

and for  $\gamma_k \equiv \gamma$ , we recognize the classical form of the Douglas–Rachford algorithm [Combettes and Pesquet, 2010].

In the conditions of Theorem 3, let  $u^*$  be any solution of (3); that is,  $u^* \in \partial H(x^*)$  and  $0 \in \partial R(x^*) + u^*$ . Then the constant  $c_0$  is

$$c_0 = \frac{1}{\gamma_0^2} \|x^1 - x^\star\|^2 + \|\frac{1}{\gamma_0}(s^0 - x^1) + u^\star\|^2.$$

On the other hand, if K = I,  $\mathcal{U} = \mathcal{X}$  and  $\eta = 1$ , the Chambolle-Pock algorithm form II becomes: Let  $x_R^0 \in \mathcal{X}$ ,  $u^0 \in \mathcal{U}$ . For k = 0, 1, ... iterate:

$$\begin{bmatrix} u^{k+1} = \operatorname{prox}_{H^*/\gamma_k} \left( u^k + \frac{1}{\gamma_k} x_R^k \right) \\ x^{k+1} = x_R^k - \gamma_k (u^{k+1} - u^k) \\ x_R^{k+1} = \operatorname{prox}_{\gamma_{k+1}R} \left( x^{k+1} - \gamma_{k+1} u^{k+1} \right) \end{bmatrix}$$

Using the Moreau identity, we obtain: Let  $x_R^0 \in \mathcal{X}$ ,  $u^0 \in \mathcal{U}$ . For k = 0, 1, ... iterate:

$$\begin{bmatrix} x^{k+1} = \operatorname{prox}_{\gamma_k H}(x_R^k + \gamma_k u^k) \\ u^{k+1} = u^k + (x_R^k - x^{k+1})/\gamma_k \\ x_R^{k+1} = \operatorname{prox}_{\gamma_{k+1} R}(x^{k+1} - \gamma_{k+1} u^{k+1}). \end{bmatrix}$$

Introducing the meta-variable  $s^k = x_R^k + \gamma_k u^k$ , we obtain: Let  $s^0 \in \mathcal{X}$ . For k = 0, 1, ... iterate:

$$\begin{bmatrix} x^{k+1} = \operatorname{prox}_{\gamma_k H}(s^k) \\ x_R^{k+1} = \operatorname{prox}_{\gamma_{k+1}R} \left( (1 + \frac{\gamma_{k+1}}{\gamma_k}) x^{k+1} - \frac{\gamma_{k+1}}{\gamma_k} s^k \right) \\ s^{k+1} = x_R^{k+1} + \frac{\gamma_{k+1}}{\gamma_k} (s^k - x^{k+1}). \end{bmatrix}$$

Thus, we recover exactly the Douglas–Rachford algorithm (19), with R and H exchanged.

### **Derivation of the Distributed Algorithms** 6

### 6.1 The Distributed PD3O Algorithm and its Particular Cases

Let us adopt the notations of Section 3 and precise the different operators. The gradient of  $\widehat{F}$  in  $\widehat{\mathcal{X}}$  is

$$\nabla \widehat{F}(\widehat{x}) = \left(\frac{1}{M\omega_1} \nabla F_1(x_1), \dots, \frac{1}{M\omega_M} \nabla F_M(x_M)\right), \quad \forall \widehat{x} \in \widehat{\mathcal{X}}$$

We define the linear subspace  $\mathcal{S} = \{ \hat{x} \in \hat{\mathcal{X}} : x_1 = \cdots = x_M \}$ .  $\hat{F}$  is  $L_{\hat{F}}$ -smooth, with  $L_{\hat{F}} = C_{\hat{F}}$  $\max_{m} \frac{L_{F_m}}{M\omega_m}.$  But since  $\nabla \widehat{F}$  is applied to an element of  $\mathcal{S}$  in the algorithms, we can weaken the condition on  $L_{\widehat{F}} > 0$  to be: for every  $\widehat{x} = (x)_{m=1}^M \in \mathcal{S}$  and  $\widehat{x}' = (x')_{m=1}^M \in \mathcal{S}$ ,

$$\begin{aligned} \|\nabla\widehat{F}(\hat{x}) - \nabla\widehat{F}(\hat{x}')\|_{\widehat{\mathcal{X}}}^2 &= \sum_{m=1}^M \omega_m \left\| \frac{1}{M\omega_m} \nabla F_m(x) - \frac{1}{M\omega_m} \nabla F_m(x') \right\|^2 \\ &\leq L_{\widehat{F}}^2 \|\widehat{x} - \widehat{x}'\|_{\widehat{\mathcal{X}}}^2 = L_{\widehat{F}}^2 \|x - x'\|^2. \end{aligned}$$

That is,  $L_{\widehat{F}}$  is such that, for every  $(x, x') \in \mathcal{X}^2$ ,

$$\frac{1}{M^2} \sum_{m=1}^M \frac{1}{\omega_m} \|\nabla F_m(x) - \nabla F_m(x')\|^2 \le L_{\widehat{F}}^2 \|x - x'\|^2.$$
(20)

Notably,

$$L_{\hat{F}}^{2} = \frac{1}{M^{2}} \sum_{m=1}^{M} \frac{L_{F_{m}}^{2}}{\omega_{m}}$$

satisfies the condition.

The adjoint operator of K is

$$\widehat{K}^*: \widehat{u} \in \widehat{\mathcal{U}} \mapsto \left(K_1^* u_1, \dots, K_M^* u_M\right) \in \widehat{\mathcal{X}}.$$

Thus,

$$\|\widehat{K}\|^{2} = \|\widehat{K}^{*}\widehat{K}\| = \max_{m} \|K_{m}\|^{2}.$$
(21)

But if  $F_1 = \cdots = F_M$ , we can restrict the norm to  $\mathcal{S}$  and

$$\begin{aligned} \|\widehat{K}\|^{2} &= \sup_{\widehat{x}\in\mathcal{S}} \langle \widehat{x}, \widehat{K}^{*}\widehat{K}\widehat{x} \rangle_{\widehat{\mathcal{X}}} / \|\widehat{x}\|_{\widehat{\mathcal{X}}}^{2} \\ &= \sup_{x\in\mathcal{X}} \langle x, \sum_{m=1}^{M} \omega_{m}K_{m}^{*}K_{m}x \rangle / \|x\|^{2} \\ &= \left\| \sum_{m=1}^{M} \omega_{m}K_{m}^{*}K_{m} \right\|, \end{aligned}$$

$$(22)$$

which is  $\leq \sum_{m=1}^{M} \omega_m \|K_m\|^2$ . For any  $\zeta > 0$ , we have  $\operatorname{prox}_{\zeta \widehat{R}} : \widehat{x} \mapsto (x', \dots, x')$ , where  $x' = \operatorname{prox}_{\zeta R} \left( \sum_{m=1}^{M} \omega_m x_m \right)$  and  $\operatorname{prox}_{\zeta \widehat{H}} :$  $\hat{u} \mapsto \left( \operatorname{prox}_{\zeta H_1/(M\omega_1)}(u_1), \dots, \operatorname{prox}_{\zeta H_M/(M\omega_M)}(u_M) \right). \text{ We also have } \partial \widehat{H} : \hat{u} \mapsto \frac{1}{M\omega_1} \partial H_1(u_1) \times \dots \times \frac{1}{M\omega_M} \partial H_M(u_M), \ \hat{H}^* : \hat{u} \mapsto \frac{1}{M} \sum_{m=1}^M H_m^*(M\omega_m u_m), \text{ and } \operatorname{prox}_{\zeta \widehat{H}^*} : \hat{u} \mapsto \left( \frac{1}{M\omega_1} \operatorname{prox}_{\zeta M\omega_1 H_1^*}(M\omega_1 u_1), \dots, \frac{1}{M\omega_M} \operatorname{prox}_{\zeta M\omega_M H_M^*}(M\omega_M u_M) \right).$ By doing all these substitutions in the PD3O algorithm, we obtain the distributed PD3O algorithm,

and all its particular cases, shown above. Theorem 1 becomes Theorem 6 as follows. The objective function is  $\Psi: x \in \mathcal{X} \mapsto R(x) + \frac{1}{M} \sum_{m=1}^{M} (F_m(x) + H_m(K_m x)).$ 

**Theorem 6** (convergence rate of the Distributed PD3O Algorithm). In the Distributed PD3O Algorithm, suppose that  $\gamma_k \equiv \gamma \in (0, 2/L_{\widehat{F}})$ , where  $\widehat{F}$  satisfies (20); if  $F_m \equiv 0$ , we can choose any  $\gamma > 0$ . Also, suppose that  $\eta \geq \|\widehat{K}\|^2$ , where  $\|\widehat{K}\|^2$  is defined in (21) or (22). Then  $x^k$  converges to some solution  $x^*$  of (1). Also,  $u_m^k$  converges to some element  $u_m^* \in \mathcal{U}_m$ , for every  $m = 1, \ldots, M$ . In addition, suppose that every  $H_m$  is continuous on an open ball centered at  $K_m x^*$ . Then the following hold:

(i) 
$$\Psi(x^k) - \Psi(x^*) = o(1/\sqrt{k}).$$

Define the weighted ergodic iterate  $\bar{x}^k = \frac{2}{k(k+1)} \sum_{i=1}^k ix^i$ , for every  $k \ge 1$ . Then

(ii) 
$$\Psi(\bar{x}^k) - \Psi(x^*) = O(1/k).$$

Furthermore, if every  $H_m$  is  $L_m$ -smooth for some  $L_m > 0$ , we have a faster decay for the best iterate so far:

(iii) 
$$\min_{i=1,\dots,k} \Psi(x^i) - \Psi(x^*) = o(1/k).$$

The theorem applies to the particular cases of the Distributed PD3O Algorithm, like the distributed Loris–Verhoeven, Chambolle–Pock, Douglas–Rachford algorithms. We can note that the distributed forward–backward algorithm is monotonic, so Theorem 6 (iii) (with  $H_m \equiv 0$ ) yields  $\Psi(x^k) - \Psi(x^*) = o(1/k)$  for this algorithm.

We now give accelerated convergence results using varying stepsizes, in presence of strong convexity. For this, we have to define the strong convexity constants  $\mu_{\widehat{F}}$  and  $\mu_{\widehat{R}}$ . Like for the smoothness constant, we can restrict their definition to  $\mathcal{S}$ . So,  $\mu_{\widehat{F}}$  becomes the strong convexity constant of the average function  $\frac{1}{M} \sum_{m=1}^{M} F_m$ . That is,  $\mu_{\widehat{F}} \geq 0$  is such that the function

$$x \in \mathcal{X} \mapsto \frac{1}{M} \sum_{m=1}^{M} F_m(x) - \frac{\mu_{\widehat{F}}}{2} \|x\|^2$$

is convex. It is much weaker to require  $\mu_{\hat{F}} > 0$  than to ask all  $F_m$  to be strongly convex. Similarly, we have  $\mu_{\hat{R}} = \mu_R$ , the strong convexity constant of R. Thus, since the Accelerated Distributed PD3O Algorithm can be viewed as the accelerated PD3O algorithm applied to the minimization of  $\hat{F}(\hat{x}) + \hat{R}(\hat{x}) + \hat{H}(\hat{K}\hat{x})$ , we have all the ingredients to invoke Theorem 3, which is transposed as:

**Theorem 7** (Accelerated Distributed PD3O Algorithm). Suppose that  $\mu_{\widehat{F}} + \mu_R > 0$ . Let  $x^*$  be the unique solution to (1). Let  $\kappa \in (0, 1)$  and  $\gamma_0 \in (0, 2(1 - \kappa)/L_{\widehat{F}})$ . Set  $\gamma_1 = \gamma_0$  and

$$\gamma_{k+1} = \frac{-\gamma_k^2 \mu_{\widehat{F}} \kappa + \gamma_k \sqrt{(\gamma_k \mu_{\widehat{F}} \kappa)^2 + 1 + 2\gamma_k \mu_R}}{1 + 2\gamma_k \mu_R}, \quad \text{for every } k \ge 1$$

Suppose that  $\eta \geq \|\widehat{K}\|^2$ , where  $\|\widehat{K}\|^2$  is defined in (21) or (22). Then in the Distributed PD3O Algorithm, there exists  $\widehat{c}_0 > 0$  such that, for every  $k \geq 1$ ,

$$\|x^{k+1} - x^{\star}\|^2 \le \frac{\gamma_{k+1}^2}{1 - \gamma_{k+1}\mu_{\widehat{F}}\kappa}\hat{c}_0 = O(1/k^2).$$

As for Theorem 5, its counterpart in the distributed setting is:

**Theorem 8** (linear convergence of the Distributed PD3O Algorithm). Suppose that  $\mu_{\widehat{F}} + \mu_R > 0$ and that every  $H_m$  is  $L_m$ -smooth, for some  $L_m > 0$ . Let  $x^*$  be the unique solution to (1). We suppose that  $\gamma_k \equiv \gamma \in (0, 2/L_{\widehat{F}})$  and  $\eta \geq \|\widehat{K}\|^2$ , where  $\|\widehat{K}\|^2$  is defined in (21) or (22). Then the Distributed PD3O Algorithm converges linearly: there exists  $\rho \in (0, 1]$  and  $\hat{c}_0 > 0$  such that, for every  $k \in \mathbb{N}$ ,

$$\|x^{k+1} - x^{\star}\|^2 \le (1 - \rho)^k \hat{c}_0.$$

We can remark that the Distributed Davis–Yin algorithm (with  $\omega_m = 1/M$  and  $\gamma_k \equiv \gamma$ ) has been proposed in an unpublished paper by Ryu and Yin [Ryu and Yin, 2017], where it is named Proximal-Proximal-Gradient Method. Their results are similar to ours in Theorems 6 and 8 for this algorithm, but their condition  $\gamma < 3/(2L)$ , with  $L = \max_m L_{F_m}$ , is worse than ours. Also, our accelerated version with varying stepsizes in Theorem 7 is new.

# 6.2 The Distributed PDDY Algorithm

The Distributed PDDY Algorithm, shown above, is derived the same way as the Distributed PD3O Algorithm. However, the smoothness constant cannot be defined only on S, so that we have

$$L_{\widehat{F}} = \max_{m=1,\dots,M} \frac{L_{F_m}}{M\omega_m}$$

and

$$\mu_{\widehat{F}} = \min_{m=1,\dots,M} \frac{\mu_{F_m}}{M\omega_m}.$$

Moreover,

$$\|\widehat{K}\|^2 = \max_{m=1,\dots,M} \|K_m\|^2,$$
(23)

except if  $F_m \equiv 0$ , in which case the Distributed PDDY Algorithm becomes the Distributed Chambolle– Pock Algorithm Form II, for which we can set

$$\|\widehat{K}\|^{2} = \left\|\sum_{m=1}^{M} \omega_{m} K_{m}^{*} K_{m}\right\|.$$
(24)

We can note that when  $K_m \equiv I$ , the Distributed PDDY Algorithm reverts to a form of distributed Davis–Yin algorithm, which is different from the Distributed Davis–Yin Algorithm obtained from the PD3O algorithm, shown above. Similarly, when R = 0, we obtain a different algorithm than the Distributed Loris–Verhoeven Algorithm shown above. When  $F_m \equiv 0$ , the Distributed PDDY Algorithm reverts to the Distributed Chambolle–Pock Algorithm Form II, which is still different from the Distributed Douglas–Rachford Algorithm when  $K_m \equiv I$ .

The counterpart of Theorem 2 is:

**Theorem 9** (convergence of the Distributed PDDY Algorithm). In the Distributed PDDY Algorithm, suppose that  $\gamma_k \equiv \gamma \in (0, 2/L_F)$  and  $\eta \geq \|\widehat{K}\|^2$ , where  $\|\widehat{K}\|^2$  is defined in (23) or (24). Then all  $x_m^k$ as well as  $x_R^k$  converge to the same solution  $x^*$  of (1), and every  $u_m^k$  converges to some element  $u_m^*$ .

The counterpart of Theorem 4 is:

**Theorem 10** (Accelerated Distributed PDDY Algorithm). Suppose that  $\mu_{\widehat{F}} > 0$ . Let  $x^*$  be the unique solution to (1). Let  $\kappa \in (0, 1)$  and  $\gamma_0 \in (0, 2(1 - \kappa)/L_{\widehat{F}})$ . Set  $\gamma_1 = \gamma_0$  and

$$\gamma_{k+1} = -\gamma_k^2 \mu_{\widehat{F}} \kappa + \gamma_k \sqrt{(\gamma_k \mu_{\widehat{F}} \kappa)^2 + 1}, \quad \text{for every } k \ge 1.$$

Suppose that  $\eta \geq \|\widehat{K}\|^2$ , where  $\|\widehat{K}\|^2$  is defined in (23) or (24). Then in the Distributed PDDY Algorithm, there exists  $\widehat{c}_0 > 0$  such that, for every  $k \geq 1$ ,

$$\sum_{m=1}^{M} \omega_m \|x_m^{k+1} - x^\star\|^2 \le \frac{\gamma_{k+1}^2}{1 - \gamma_{k+1} \mu_F \kappa} c_0 = O(1/k^2).$$

Consequently, for every  $m = 1, \ldots, M$ ,

$$||x_m^k - x^*||^2 = O(1/k^2).$$

Moreover, if  $\eta > \|\widehat{K}\|^2$ ,  $\|x_R^k - x^\star\|^2 = O(1/k^2)$  as well.

Distributed Condat–Vũ Alg. Form I	Distributed Condat-Vỹ Alg. Form II
$\begin{aligned} & \text{input: } \gamma > 0,  \sigma > 0,  (\omega_m)_{m=1}^M \\ & x_0 \in \mathcal{X},  (u_m^0)_{m=1}^M \in \widehat{\mathcal{U}} \\ & \text{initialize: } a_m^0 \coloneqq K_m^* u_m^0 + \nabla F_m(x^0),  \forall m \\ & \text{for } k = 0, 1, \dots \text{ do} \\ & \text{at master, } \text{do} \\ & x^{k+1} \coloneqq \operatorname{prox}_{\gamma R} \left( x^k - \frac{\gamma}{M} \sum_{m=1}^M a_m^k \right) \\ & \text{broadcast } x^{k+1} \text{ to all nodes} \\ & \text{at all nodes, for } m = 1, \dots, M,  \text{do} \\ & u_m^{k+1} \coloneqq \operatorname{prox}_{M\omega_m \sigma H_m^*} \left( u_m^k \\ & + M\omega_m \sigma K_m (2x^{k+1} - x^k) \right) \\ & a_m^{k+1} \coloneqq K_m^* u_m^{k+1} + \nabla F_m(x^{k+1}) \\ & \text{transmit } a_m^{k+1} \text{ to master} \\ & \text{end for} \end{aligned}$	Distributed Condat–Vu Alg. Form II input: $\gamma > 0, \sigma > 0, (\omega_m)_{m=1}^M$ $x_0 \in \mathcal{X}, (u_m^0)_{m=1}^M \in \widehat{\mathcal{U}}$ for $k = 0, 1, \dots$ do at all nodes, for $m = 1, \dots, M$ , do $u_m^{k+1} \coloneqq \operatorname{prox}_{M\omega_m \sigma H_m^*}(u_m^k)$ $+ M\omega_m \sigma K_m x^k)$ $a_m^k \coloneqq K_m^*(2u_m^{k+1} - u_m^k) + \nabla F_m(x^k)$ transmit $a_m^k$ to master at master, do $x^{k+1} \coloneqq \operatorname{prox}_{\gamma R}(x^k - \frac{\gamma}{M} \sum_{m=1}^M a_m^k)$ broadcast $x^{k+1}$ to all nodes end for

The counterpart of Theorem 5 is:

**Theorem 11** (linear convergence of the Distributed PDDY Algorithm). Suppose that  $\mu_{\widehat{F}} + \mu_R > 0$ and that every  $H_m$  is  $L_m$ -smooth, for some  $L_m > 0$ . Let  $x^*$  be the unique solution to (1). Suppose that  $\gamma_k \equiv \gamma \in (0, 2/L_{\widehat{F}})$  and  $\eta \geq \|\widehat{K}\|^2$ , where  $\|\widehat{K}\|^2$  is defined in (23) or (24). Then the Distributed PDDY Algorithm converges linearly: there exists  $\rho \in (0, 1]$  and  $\hat{c}_0 > 0$  such that, for every  $k \in \mathbb{N}$ ,

$$||x_R^{k+1} - x^\star||^2 \le (1-\rho)^k \hat{c}_0$$

# 6.3 The Distributed Condat–Vũ Algorithm

We can apply our product-space technique to other algorithms; in particular, we can derive distributed versions, shown below, of the Condat–Vũ algorithm [Condat, 2013, Vũ, 2013, Condat et al., 2019a], which is a well known algorithm for the problem (2).

The smoothness constant  $L_{\widehat{F}}^2$  is the same as for the Distributed PD3O Algorithm; we can set  $L_{\widehat{F}}^2 = \frac{1}{M^2} \sum_{m=1}^M L_{F_m}^2 / \omega_m$ .

Moreover, the norm of  $\hat{K}$  is smaller for the Condat–Vũ algorithm: we have  $\|\hat{K}\|^2 = \|\sum_{m=1}^{M} \omega_m K_m^* K_m\|$ , whatever the functions  $F_m$ . This is because the gradient descent step is completely decoupled from the dual variables in the Condat–Vũ algorithm.

The price to pay is a stronger condition on the parameters for convergence:

**Theorem 12** (convergence of the Distributed Condat–Vũ Algorithm). Suppose that the parameters  $\gamma > 0$  and  $\sigma > 0$  are such that

$$\gamma\left(\sigma \|\sum_{m=1}^{M} \omega_m K_m^* K_m \| + \frac{L_{\widehat{F}}}{2}\right) < 1.$$

Then  $x^k$  converges to a solution  $x^*$  of (1). Also,  $u_m^k$  converges to some element  $u_m^* \in \mathcal{U}_m$ , for every  $m = 1, \ldots, M$ .

When  $F_m \equiv 0$ , the two forms of the Distributed Condat–Vũ Algorithm revert to the two forms of the Distributed Chambolle–Pock Algorithm, respectively. In that case, with constant stepsizes  $\gamma_k \equiv \gamma$ , the convergence condition is  $\gamma \sigma \| \sum_{m=1}^M \omega_m K_m^* K_m \| \leq 1$ , which is the same as above with  $\sigma = 1/(\eta \gamma)$ .

# Author Contributions

Grigory Malinovsky wrote the code and generated the results for the SVM experiment in Section 4.3. Peter Richtárik contributed to the paper writing and to the project management. Laurent Condat did all the rest.

# References

- S. A. Alghunaim, E. K. Ryu, K. Yuan, and A. H. Sayed. Decentralized proximal gradient algorithms with linear convergence rates. *IEEE Trans. Autom. Control*, 66(6):2787–2794, June 2021.
- F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with sparsity-inducing penalties. Found. Trends Mach. Learn., 4(1):1–106, 2012.
- H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces.* Springer, New York, 2nd edition, 2017.
- A. Beck. First-Order Methods in Optimization. MOS-SIAM Series on Optimization. SIAM, 2017.
- R. I. Boţ, E. R. Csetnek, and C. Hendrich. Recent developments on primal-dual splitting methods with applications to convex minimization. In P. M. Pardalos and T. M. Rassias, editors, *Mathematics Without Boundaries: Surveys in Interdisciplinary Research*, pages 57–99. Springer New York, 2014.
- K. Bredies, K. Kunisch, and T. Pock. Total generalized variation. *SIAM J. Imaging Sci.*, 3(3): 492–526, 2010.
- S. Bubeck. Convex optimization: Algorithms and complexity. *Found. Trends Mach. Learn.*, 8(3–4): 231–357, 2015.
- V. Cevher, S. Becker, and M. Schmidt. Convex optimization for big data: Scalable, randomized, and parallel algorithms for big data analytics. *IEEE Signal Process. Mag.*, 31(5):32–43, 2014.
- A. Chambolle and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. J. Math. Imaging Vision, 40(1):120–145, May 2011.
- A. Chambolle and T. Pock. An introduction to continuous optimization for imaging. Acta Numerica, 25:161–319, 2016a.
- A. Chambolle and T. Pock. On the ergodic convergence rates of a first-order primal-dual algorithm. Math. Program., 159(1-2):253-287, Sept. 2016b.
- C.-C. Chang and C.-J. Lin. LibSVM: A library for support vector machines. ACM Transactions on Intelligent Systems and Technology (TIST), 2(3):27, 2011.
- P. Chen, J. Huang, and X. Zhang. A primal-dual fixed point algorithm for convex separable minimization with applications to image restoration. *Inverse Problems*, 29(2), 2013.
- P. L. Combettes and J.-C. Pesquet. Proximal splitting methods in signal processing. In H. H. Bauschke, R. Burachik, P. L. Combettes, V. Elser, D. R. Luke, and H. Wolkowicz, editors, *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, pages 185–212. Springer-Verlag, New York, 2010.
- P. L. Combettes and J.-C. Pesquet. Primal-dual splitting algorithm for solving inclusions with mixtures of composite, Lipschitzian, and parallel-sum type monotone operators. *Set-Val. Var. Anal.*, 20(2):307–330, 2012.

- P. L. Combettes, L. Condat, J.-C. Pesquet, and B. C. Vũ. A forward–backward view of some primal–dual optimization methods in image recovery. In *Proc. of IEEE ICIP*, pages 4141–4145, Paris, France, Oct. 2014.
- L. Condat. A primal-dual splitting method for convex optimization involving Lipschitzian, proximable and linear composite terms. J. Optim. Theory Appl., 158(2):460–479, 2013.
- L. Condat. A generic proximal algorithm for convex optimization—Application to total variation minimization. *IEEE Signal Process. Lett.*, 21(8):985–989, Aug. 2014.
- L. Condat. Discrete total variation: New definition and minimization. SIAM J. Imaging Sci., 10(3): 1258–1290, 2017a.
- L. Condat. A convex approach to K-means clustering and image segmentation. In Proc. of EMMCVPR. In: M. Pelillo and E. Hancock eds., Lecture Notes in Computer Science vol. 10746, Springer, 2018, pages 220–234, Venice, Italy, Oct. 2017b.
- L. Condat, D. Kitahara, A. Contreras, and A. Hirabayashi. Proximal splitting algorithms for convex optimization: A tour of recent advances, with new twists. preprint arXiv:1912.00137, 2019a.
- L. Condat, D. Kitahara, and A. Hirabayashi. A convex lifting approach to image phase unwrapping. In Proc. of IEEE ICASSP, pages 1852–1856, Brighton, UK, 2019b.
- D. Cremers, T. Pock, K. Kolev, and A. Chambolle. Convex relaxation techniques for segmentation, stereo and multiview reconstruction. In *Markov Random Fields for Vision and Image Processing*. MIT Press, 2011.
- D. Davis and W. Yin. A three-operator splitting scheme and its optimization applications. Set-Val. Var. Anal., 25:829–858, 2017.
- Y. Drori, S. Sabach, and M. Teboulle. A simple algorithm for a class of nonsmooth convex concave saddle-point problems. Oper. Res. Lett., 43(2):209–214, 2015.
- J. Duran, M. Moeller, C. Sbert, and D. Cremers. Collaborative total variation: A general framework for vectorial TV models. *SIAM J. Imaging Sci.*, 9(1):116–151, 2016.
- R. Glowinski, S. J. Osher, and W. Yin, editors. *Splitting Methods in Communication, Imaging, Science, and Engineering.* Springer International Publishing, 2016.
- E. Gorbunov, F. Hanzely, and P. Richtárik. A unified theory of SGD: Variance reduction, sampling, quantization and coordinate descent. In Proc. of Int. Conf. Artif. Intell. Stat. (AISTATS), PMLR 108, pages 680–690, Palermo, Sicily, Italy, June 2020.
- N. Komodakis and J.-C. Pesquet. Playing with duality: An overview of recent primal-dual approaches for solving large-scale optimization problems. *IEEE Signal Process. Mag.*, 32(6):31–54, Nov. 2015.
- J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon. Federated learning: Strategies for improving communication efficiency. In NIPS Private Multi-Party Machine Learning Workshop, 2016. paper arXiv:1610.05492.
- P. Latafat, N. M. Freris, and P. Patrinos. A new randomized block-coordinate primal-dual proximal algorithm for distributed optimization. *IEEE Trans. Autom. Control*, 64(10):4050–4065, Oct. 2019.
- I. Loris and C. Verhoeven. On a generalization of the iterative soft-thresholding algorithm for the case of non-separable penalty. *Inverse Problems*, 27(12), 2011.

- G. Malinovsky, D. Kovalev, E. Gasanov, L. Condat, and P. Richtárik. From local SGD to local fixed point methods for federated learning. In *Proc. of 37th Int. Conf. Machine Learning (ICML)*, *PMLR 119*, pages 6692–6701, 2020.
- D. O'Connor and L. Vandenberghe. On the equivalence of the primal-dual hybrid gradient method and Douglas–Rachford splitting. *Math. Program.*, 179:85–108, 2020.
- D. P. Palomar and Y. C. Eldar, editors. Convex Optimization in Signal Processing and Communications. Cambridge University Press, 2009.
- N. Parikh and S. Boyd. Proximal algorithms. Foundations and Trends in Optimization, 3(1):127–239, 2014.
- N. G. Polson, J. G. Scott, and B. T. Willard. Proximal algorithms in statistics and machine learning. *Statist. Sci.*, 30(4):559–581, 2015.
- P. Richtárik and M. Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Math. Program.*, 144(1–2):1–38, Apr. 2014.
- E. K. Ryu and W. Yin. Proximal-proximal-gradient method. preprint arXiv:1708.06908, 2017.
- A. Salim, L. Condat, K. Mishchenko, and P. Richtárik. Dualize, split, randomize: Fast nonsmooth optimization algorithms. preprint arXiv:2004.02635, 2020.
- A. Salim, L. Condat, D. Kovalev, and P. Richtárik. An optimal algorithm for strongly convex minimization under affine constraints. preprint arXiv:2102.11079. Accepted at AISTATS 2022, 2021.
- K. Scaman, F. Bach, S. Bubeck, Y. T. Lee, and L. Massoulié. Optimal algorithms for smooth and strongly convex distributed optimization in networks. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, volume 70, pages 3027–3036, 2017.
- W. Shi, Q. Ling, G. Wu, and W. Yin. EXTRA: An exact first-order algorithm for decentralized consensus optimization. SIAM J. Optim., 25(2):944–966, 2015.
- S. Sra, S. Nowozin, and S. J. Wright. Optimization for Machine Learning. The MIT Press, 2011.
- G. Stathopoulos, H. Shukla, A. Szucs, Y. Pu, and C. N. Jones. Operator splitting methods in control. Foundations and Trends in Systems and Control, 3(3):249–362, 2016.
- Unknown author. Every convex function is locally Lipschitz. The American Mathematical Monthly, 79(10):1121–1124, Dec. 1972.
- B. C. Vũ. A splitting algorithm for dual monotone inclusions involving cocoercive operators. Adv. Comput. Math., 38(3):667–681, Apr. 2013.
- Y.-X. Wang, J. Sharpnack, A. Smola, and R. Tibshirani. Trend filtering on graphs. Journal of Machine Learning Research, 17(105):1–41, 2016.
- M. Yan. A new primal-dual algorithm for minimizing the sum of three functions with a linear operator. J. Sci. Comput., 76(3):1698–1717, Sept. 2018.