



**HAL**  
open science

## Acceleration of saddle-point methods in smooth cases

Pauline Tan

► **To cite this version:**

| Pauline Tan. Acceleration of saddle-point methods in smooth cases. 2017. hal-01415459v2

**HAL Id: hal-01415459**

**<https://hal.science/hal-01415459v2>**

Preprint submitted on 5 Jan 2017 (v2), last revised 16 Jan 2017 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# ACCELERATION OF SADDLE-POINT METHODS IN SMOOTH CASES

Pauline TAN

CMAP, École polytechnique, CNRS, Université Paris-Saclay, 91128, Palaiseau, France

## Abstract

In the present paper we propose a novel convergence analysis of the Alternating Direction Methods of Multipliers (ADMM), based on its equivalence with the overrelaxed Primal-Dual Hybrid Gradient (oPDHG) algorithm. We consider the smooth case, which correspond to the case where the objective function can be decomposed into one differentiable with Lipschitz continuous gradient part and one strongly convex part. An accelerated variant of the ADMM is also proposed, which is shown to converge linearly with same rate as the oPDHG.

## 1 Introduction

### 1.1 Context

The Alternating Direction Methods of Multipliers (ADMM) is a widely-used method aimed at minimizing constrained problems of form

$$\min_{\substack{(x,z) \in X \times Z \\ Ax+Bz=c}} g(x) + h(z). \quad (1)$$

The objective function is separable in  $(x, z)$  with  $g : X \rightarrow \mathbb{R} \cup \{+\infty\}$  and  $h : Z \rightarrow \mathbb{R} \cup \{+\infty\}$  two closed convex functions. The constraint involves two linear operators  $A : X \rightarrow Y$  and  $B : Z \rightarrow Y$  and a constant  $c \in Y$ . In this work,  $X$ ,  $Z$ , and  $Y$  are finite-dimensional real Hilbert spaces. The ADMM was initially introduced in the mid-70's by GABAY-MERCIER [9] and by GLOWINSKI-MARROCCO [10]. It considers the augmented Lagrangian associated to problem (1)

$$L_\tau(x, z; y) := g(x) + h(z) + \langle Ax + Bz - c, y \rangle + \frac{1}{2\tau} \|Ax + Bz - c\|^2 \quad (2)$$

for  $\tau > 0$  which leads to solve the saddle-point problem

$$\min_{(x,z) \in X \times Z} \sup_{y \in Y} L_\tau(x, z; y) \quad (3)$$

instead of the initial problem. One particular instance of these so-called augmented Lagrangian methods uses UZAWA's method to solve (3). Namely, the *method of multipliers* tackles this problem by alternating an exact minimization on the primal variable  $(x, z)$  and a gradient ascent step on the dual variable  $y$ . In such a method, the minimization step couples the primal variables. To decouple them, one may consider splitting this step into two partial minimizations, one over  $x$  and another over  $z$ . These two minimization can be done simultaneously, from the same initial points, or, in the case of the ADMM, one after the other, with an update in between. This leads to the following algorithm

$$\begin{cases} x_{n+1} = \arg \min_{x \in X} L_\tau(x, z_n; y_n) \\ z_{n+1} = \arg \min_{z \in Z} L_\tau(x_{n+1}, z; y_n) \\ y_{n+1} = y_n + \frac{1}{\tau} (Ax_{n+1} - Bz_{n+1}). \end{cases} \quad (4)$$

This method can be proved to be linked to another famous algorithm, which is known as the Primal-Dual Hybrid Gradient (PDHG) method [18]. The PDHG method tackles saddle-point problems by alternating gradient descent steps and gradient ascent steps. Such problems arise while considering a primal-dual formulation of a convex minimization problem, in a splitting strategy for instance. A noteworthy feature of the PDHG method is that it can be accelerated thanks to an overrelaxation step *à la* NESTEROV [13] on one of the variables [17, 3, 8, 4], which leads to the overrelaxed PDHG (oPDHG).

The ADMM has been intensively studied in the past years. One may see for instance a comprehensive review in [2]. The key point is the convergence of the algorithm and its convergence rate. Under assumptions on the matrix ranks and / or the regularity of the objective functions  $g$  and  $h$ , linear rates can be achieved [12]. Eventually, some accelerated variants of the ADMM have been proposed [7, 6].

As a recent developpement, we should mention [11], which also studied the convergence of the PDHG method and derived optimal step size choice, when only one function assumed to be strongly convex.

## 1.2 Contribution of this paper

In this paper, we provide a new analysis of the ADMM based on the equivalence between the ADMM and the oPDHG method. More specifically, we use the analysis to derive convergence rate for the ADMM in a case we refer to be *smooth*. We indeed made restrictive assumptions on the initial problem (1), which implies that we consider the following particular instance of (1):

$$\min_{\substack{(x,z) \in X \times Z \\ Ax=z}} g(x) + h(z). \quad (5)$$

which may be rewritten as the unconstrained composite problem

$$\min_{x \in X} g(x) + h(Ax) \quad (6)$$

with regularity assumptions on  $g$ , which is supposed to be strongly convex, and  $h$ , which has a Lipschitz gradient. We first establish new linear ergodic convergence rates of the oPDHG by generalizing the proofs of [3, 4]. This leads to a linear rate for the ADMM under these assumptions. Then, we introduce a slight variant of the ADMM which leads to a better rate, by relaxing the choice of the parameters in the convergence proof of the oPDHG method.

The reason why we only consider the case  $B = -\text{Id}$  and  $c = 0$  is that, otherwise, as the map  $y \mapsto h^*(B^*y)$  will be supposed to be strongly convex, this implies that  $\nabla f$  is Lipschitz continuous and that  $B$  is invertible. Such conditions are artificial when  $B$  is not  $-\text{Id}$ . However, the interested reader will easily extend our result to this case. Moreover, problems of standard form (6) often arise in many contexts, and thus can justify a special study by themselves.

## 1.3 Structure of the paper

This paper is organized as follows. In Section 2, we recall the equivalence between the ADMM and the oPDHG method. We also define what we call the *smooth case*, which is the case we will consider throughout this paper. In Section 3, we establish two linear convergence results for the oPDHG, and we provide the best parameter choice in the case where the overrelaxation parameter is fixed to be 1 or left unconstrained. In Section 4, we exploit the equivalence between the ADMM and the oPDHG to derive from the results

of Section 3 new linear convergence rate for the ADMM. We also propose a slight variant of the ADMM, which leads in the best case to the same convergence rate as the oPDHG method. In Section 5, we compare our results with some found in the literature for the classical ADMM or variants, in the case where the assumptions made on the problem yield a linear convergence rate. Those assumptions do not necessary include the smooth case studied here. Eventually, in Section 6, we applied our accelerated ADMM on two problems, and compared its convergent with the unaccelerated ADMM, the oPDHG and an adaptation of BECK and TEBoulLE's FISTA [1] for the strongly convex case [14, 5].

## 2 Equivalence between the ADMM and the oPDHG

### 2.1 Initial primal problem

Let  $X$  and  $Y$  be two finite-dimensional real Hilbert spaces. The inner product is denoted by  $\langle \cdot, \cdot \rangle$  and  $\|\cdot\|$  stands for the induced norm. We recall that we consider the minimization problem

$$\min_{x \in X} \left\{ f(x) := g(x) + h(Ax) \right\} \quad (7)$$

where  $g : X \rightarrow \mathbb{R} \cup \{+\infty\}$  and  $h : Y \rightarrow \mathbb{R} \cup \{+\infty\}$  are proper, convex, and lower semi-continuous (l.s.c.) functions. The map  $A : X \rightarrow Y$  is a continuous linear operator. Its adjoint is denoted by  $A^*$  and it is supposed to be bounded, of norm  $L_A$

$$L_A := \|A\| = \sup_{x \in X, \|x\| \leq 1} \|Ax\|. \quad (8)$$

### 2.2 Equivalence with oPDHG

Let us briefly recall how the ADMM is connected to the oPDHG algorithm, by rewriting the ADMM iterations (4) applied on Problem (7). Ignoring the constant terms in the minimization steps, we obtain

$$\begin{cases} x_{n+1} = \arg \min_{x \in X} \left\{ g(x) + \langle Ax, y_n \rangle + \frac{1}{2\tau} \|Ax - z_n\|^2 \right\} \\ z_{n+1} = \arg \min_{z \in X} \left\{ h(z) - \langle z, y_n \rangle + \frac{1}{2\tau} \|Ax_{n+1} - z\|^2 \right\} \\ y_{n+1} = y_n + \frac{1}{\tau} (Ax_{n+1} - z_{n+1}). \end{cases} \quad (9)$$

Defining  $\xi_{n+1} := Ax_{n+1}$  and introducing the map

$$g_A(\xi) := \inf_{x \in X, Ax=\xi} g(x) \quad (10)$$

we can make a change of variable in the  $x$ -update and rewrite the updates of  $x_{n+1}$  and  $y_{n+1}$  thanks to proximity operators. This yields

$$\begin{cases} \xi_{n+1} = \text{prox}_{\tau g_A}(\xi_n - \tau \bar{y}_n) \\ y_{n+1} = \text{prox}_{h^*/\tau}(y_n + \xi_{n+1}/\tau) \\ \bar{y}_{n+1} = y_{n+1} + (y_{n+1} - y_n) \end{cases} \quad (11)$$

and the  $z$ -update is given by  $z_{n+1} = \xi_{n+1} - \tau (y_{n+1} - y_n)$ . This primal-dual algorithm has been studied in [3]. It can be interpreted as an PDHG algorithm with an additional

overrelaxation step (of parameter 1) on the dual variable. It solves the saddle-point problem

$$\min_{\xi \in X} \sup_{y \in Y} \left\{ g_A(\xi) + \langle \xi, y \rangle - h^*(y) \right\} \quad (12)$$

which is of general form

$$\min_{\xi \in X} \sup_{y \in Y} \left\{ \mathcal{L}(\xi; y) := G(\xi) + \langle K\xi, y \rangle - H^*(y) \right\} \quad (13)$$

with  $K = \text{Id}$ ,  $G = g_A$  and  $H = h$ . Note that (13) is the primal-dual formulation of the minimization problem

$$\min_{\xi \in X} \left\{ G(\xi) + H(K\xi) \right\}. \quad (14)$$

### 2.3 Smooth case

From now on, we consider the *smooth case*. In the initial primal problem (7), the functions  $g$  and  $h^*$  are both supposed to be strongly convex, with respective parameter  $\gamma > 0$  and  $\delta > 0$ . We recall that a function  $f : X \rightarrow \mathbb{R} \cup \{+\infty\}$  is strongly convex of parameter  $\alpha > 0$  ( $f$  is also said to be  $\alpha$ -convex) if for any  $x_1, x_2 \in X$  and  $p \in \partial f(x_1)$

$$f(x_2) \geq f(x_1) + \langle p, x_2 - x_1 \rangle + \frac{\alpha}{2} \|x_2 - x_1\|^2 \quad (15)$$

where  $\partial f(x_1)$  denotes the subdifferential of  $f$  at point  $x_1$ . One can easily check that if  $f$  is  $\alpha$ -convex, then its convex conjugate  $f^*$  is differentiable, with a Lipschitz continuous gradient, of constant  $1/\alpha$ .

Let us study the regularity of Problem (13). The assumptions made above imply obviously that  $H^*$  is  $\delta$ -convex. Moreover, it is easy to show that  $G^*$  is differentiable and that  $\nabla G^*$  is Lipschitz continuous with constant  $L_A^2/\gamma$ , which follows from

$$g_A^*(y + t) = g^*(A^*(y + t)) = g^*(A^*y) + \langle \nabla g^*(A^*y), A^*t \rangle + o(\|A^*t\|) \quad (16)$$

since  $g$  is  $\gamma$ -convex. Hence,  $G$  is  $\gamma/L_A^2$ -convex. Let  $\tilde{\gamma} = \gamma/L_A^2$  and  $\tilde{\delta} = \delta$ .

We define  $\kappa_f := L_A^2/(\delta\gamma)$  the *condition number* of  $f$  as the ratio between  $L_A^2/\delta$  the Lipschitz constant of the smooth part  $h(K\cdot)$  and  $\gamma$  the strong convexity parameter of the non-smooth part  $g$ . In the case where  $f$  is both smooth with  $\nabla f$  Lipschitz continuous and strongly convex, this definition recovers the one usually used in such cases and the condition number is always larger than 1. In the general case, it can be less than 1. When  $\kappa_f$  is large, the function is said ill-conditioned.

### 2.4 Forward-backward splitting

If  $h$  is differentiable, it is possible to consider a forward-backward splitting (FBS) strategy to solve problem (7). The FBS applied on the sum  $f = g + h(A\cdot)$  gives updates of form

$$x_{n+1} = \text{prox}_{\tau g}(x_n - \tau A^* \nabla h(Ax_n)). \quad (17)$$

Hence, choosing to use the FBS instead of the ADMM or the oPDHG method suggests that  $\nabla h$  is supposed to be easier to compute than  $\text{prox}_g$ .

A variant of the FBS is FISTA [1], which adds an extra overrelaxation step. It can be adapted to solve for strongly convex problems following [14], see [5, Appendix B] for details. In other terms, the updates (17) are replaced by

$$\begin{cases} x_{n+1} = \text{prox}_{\tau g}(\bar{x}_n - \tau A^* \nabla h(A\bar{x}_n)) \\ \bar{x}_{n+1} = x_{n+1} + \theta_{n+1} (x_{n+1} - x_n) \end{cases} \quad (18)$$

where the variable overrelaxation parameter  $\theta_n$  is chosen in the strongly convex case by letting

$$t_{n+1} = \frac{1 - qt_n^2 + \sqrt{(1 - qt_n^2)^2 + 4t_n^2}}{2} \quad (19)$$

for  $q = \tau\gamma/(1 + \tau\gamma)$  for  $\tau \in (0, \delta/L_A^2]$ . Then,

$$\theta_n = (1 + \tau\gamma(1 - t_{n+1})) \frac{t_n - 1}{t_{n+1}}. \quad (20)$$

In the non-strongly convex case ( $\gamma = \delta = 0$ ), the quantity  $q$  is null, and the resulting updates of  $t_n$  and  $\theta_n$  are those in the original paper of BECK and TEBoulLE. When  $g$  is assumed to be strongly convex and  $h(A \cdot)$  has a  $L_A^2/\delta$ -LIPSCHITZ gradient, the convergence rate for the objective error of this algorithm has been proved to be linear. In the case where  $t$  (and thus,  $\theta$ ) is chosen to be constant

$$t_n = t = \frac{1}{\sqrt{q}} \quad \text{and} \quad \theta_n = \theta = (1 - \sqrt{q})^2 \frac{1 + \tau\gamma}{1 - \tau\gamma} \quad (21)$$

then the linear rate is of parameter  $1 - \sqrt{q}$  [5, Remark B.2]. This rate is minimal when  $\tau$  is maximal and equals

$$\omega = 1 - \sqrt{\frac{\delta\gamma/L_A^2}{1 + \delta\gamma/L_A^2}} = 1 - \sqrt{\frac{1}{\kappa_f + 1}}. \quad (22)$$

Note that, in the case where  $g$  is  $\gamma$ -convex then this is also the case for  $f$ . Hence, the optimality condition on  $x^*$  coupled with the strong convexity inequality recalled in (15) yields

$$f(x_n) - f(x^*) \geq \frac{\gamma}{2} \|x_n - x^*\|^2 \quad (23)$$

that is, a linear convergence for the objective error implies a linear convergence of at least same rate for the convergence of the primal iterate  $x_n$ .

### 3 Convergence of oPDHG in the smooth case

In this section, we establish the general convergence proof of the following algorithm

$$\begin{cases} y_{n+1} = \text{prox}_{\sigma H^*}(y_n + \sigma K \bar{\xi}_n) \\ \xi_{n+1} = \text{prox}_{\tau G}(\xi_n - \tau K^* y_{n+1}) \\ \bar{\xi}_{n+1} = \xi_{n+1} + \theta(\xi_{n+1} - \xi_n). \end{cases} \quad (24)$$

which aims at solving problem (13), in the general case where  $K : Z \rightarrow Y$  is bounded of norm  $L_K$ ,  $G : Z \rightarrow \mathbb{R} \cup \{+\infty\}$  is  $\tilde{\gamma}$ -convex and  $H^* : Y \rightarrow \mathbb{R} \cup \{+\infty\}$  is  $\tilde{\delta}$ -convex. The step sizes  $\tau, \sigma > 0$  and the relaxation parameter  $0 < \theta \leq 1$  are to be specified.

When  $\theta = 0$ , this algorithm is known as the PDHG method [18]. It consists in a proximal gradient ascent step for the dual variable, followed by a proximal gradient descent step for the primal variable. The overrelaxation step has been added in [17] for minimizing the MUMFORD-SHAH functional, and studied in a wider framework in [8] and more recently in [4]. The case  $\theta = 1$  and  $\tau\sigma = 1$  corresponds to the equivalence with the ADMM, as recalled in the previous section. When  $\theta = 1$  and  $\tau\sigma \neq 1$ , the iterations are equivalent to the ADMM with an additional proximal term [3], which leads to a preconditioned version of the ADMM [8].

Now we can formulate our main result.

**Theorem 1** Assume problem (13) has a solution, which is a saddle-point of  $\mathcal{L}$ , denoted by  $(\xi^*, y^*)$ . Choose  $\tau > 0$ ,  $\sigma > 0$  and  $0 < \theta \leq 1$  such that

$$\max \left\{ \frac{1}{\tau\tilde{\gamma} + 1}, \frac{1}{\sigma\tilde{\delta} + 1} \right\} \leq \theta \leq \frac{1}{L_K^2 \tau \sigma}. \quad (25)$$

Then, for any  $\omega$  such that

$$\max \left\{ \frac{1}{\tau\tilde{\gamma} + 1}, \frac{\theta + 1}{\sigma\tilde{\delta} + 2} \right\} \leq \omega \leq \theta \quad (26)$$

we have the following majoration for any  $N \in \mathbb{N}$  and any  $(\xi, y) \in Z \times Y$ :

$$\begin{aligned} \frac{1}{2\tau} \|\xi_N - \xi\|^2 + (1 - \omega L_K^2 \tau \sigma) \frac{1}{2\sigma} \|y_N - y\|^2 + \sum_{n=1}^N \frac{\omega^n}{\omega^{n-1}} (\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n)) \\ \leq \frac{\omega^N}{2\tau} \|\xi_0 - \xi\|^2 + \frac{\omega^N}{2\sigma} \|y_0 - y\|^2 \end{aligned} \quad (27)$$

where  $(\xi_n, y_n)_n$  are generated by Algorithm (24). Hence, if we define

$$T_N := \sum_{n=1}^N \frac{1}{\omega^{n-1}} = \frac{1 - \omega^N}{\omega^{N-1}(1 - \omega)} \quad (28)$$

and let

$$\Xi_N := \frac{1}{T_N} \sum_{n=1}^N \frac{1}{\omega^{n-1}} \xi_n \quad \text{and} \quad Y_N := \frac{1}{T_N} \sum_{n=1}^N \frac{1}{\omega^{n-1}} y_n. \quad (29)$$

Then we have the following bound for any  $(\xi, y) \in Z \times Y$ :

$$\begin{aligned} \frac{1 - \omega}{\omega(1 - \omega^N)} \frac{1}{2\tau} \|\xi - \xi_N\|^2 + \frac{1 - \omega}{\omega(1 - \omega^N)} (1 - \omega L_K^2 \tau \sigma) \frac{1}{2\sigma} \|y - y_N\|^2 \\ + \mathcal{L}(\Xi_N; y) - \mathcal{L}(\xi; Y_N) \\ \leq \frac{1}{T_N} \frac{1}{2\tau} \|\xi - \xi_0\|^2 + \frac{1}{T_N} \frac{1}{2\sigma} \|y - y_0\|^2. \end{aligned} \quad (30)$$

This theorem provides a linear ergodic convergence rate, namely for the sequences  $(\Xi_N)$  and  $(Y_N)$ . This rate can be compared with [14], and will be proved to be better with optimal parameters. Also note that no assumption is made about the rank of the linear operator  $K$ . Equation (30) can be applied to  $\xi = \xi^*$  and  $y = y^*$ , which yields a non-ergodic linear convergence rate for the variable convergence (see subsection 3.1.3). The objective error can be measured thanks to the partial primal-dual gap:

$$\mathcal{G}_{\mathcal{B}_Z \times \mathcal{B}_Y}(\Xi; Y) := \sup_{y \in \mathcal{B}_Y} \mathcal{L}(\Xi; y) - \inf_{\xi \in \mathcal{B}_Z} \mathcal{L}(\xi; Y) \quad (31)$$

for  $\mathcal{B}_Z \times \mathcal{B}_Y$  an bounded open subset of  $Z \times Y$  which contains  $(\xi^*, y^*)$ . The partial gap is nonnegative and expected to be zero iff  $(\Xi, Y)$  is a saddle-point of  $\mathcal{L}$ . Of course, this implies that  $\mathcal{B}_Z$  and  $\mathcal{B}_Y$  can be estimated.

A similar result may be found in [3], but the rate we provide here is better, since no restrictive assumptions are made on the parameters values, unless necessary.

### 3.1 Proof of convergence

We proceed analogously to the proof in [3], but we do not specify any parameter unless needed. This proof is also inspired by the one found in [4], which does not allow  $\theta \neq 1$ . For now, we only assume that  $0 < \theta \leq 1$ .

#### 3.1.1 Preliminaries

Let us consider the general updates  $(\hat{y}, \hat{\xi})$  by setting for any  $(\bar{\xi}, \tilde{\xi}) \in Z^2$  and  $(\bar{y}, \tilde{y}) \in Y^2$

$$\begin{cases} \hat{y} := \text{prox}_{\sigma H^*}(\bar{y} + \sigma K \tilde{\xi}) \\ \hat{\xi} := \text{prox}_{\tau G}(\bar{\xi} - \tau K^* \tilde{y}). \end{cases} \quad (32)$$

In other terms,  $\hat{y}$  and  $\hat{\xi}$  are the output of an iteration, and are respectively computed from initial points  $(\bar{y}, \bar{\xi})$  and  $(\tilde{y}, \tilde{\xi})$ . These points are related by first-order optimality conditions. For instance, the point  $\hat{\xi}$  is defined as the solution of a minimization problem

$$\hat{\xi} = \arg \min_{\xi \in Z} \left\{ \frac{1}{2\tau} \|\bar{\xi} - \tau K^* \tilde{y} - \xi\|^2 + G(\xi) \right\} \quad (33)$$

so, by optimality, we obtain

$$-\frac{1}{\tau} (\hat{\xi} - \bar{\xi}) - K^* \tilde{y} \in \partial G(\hat{\xi}). \quad (34)$$

Similarly, the definition of  $\hat{y}$  yields

$$-\frac{1}{\sigma} (\hat{y} - \bar{y}) + K \tilde{\xi} \in \partial F^*(\hat{y}). \quad (35)$$

Using the definition of strong convexity recalled in (15), we get (after expanding the scalar products)

$$G(\xi) + \frac{1}{2\tau} \|\xi - \bar{\xi}\|^2 \geq G(\hat{\xi}) + \langle K(\hat{\xi} - \xi), \tilde{y} \rangle + \frac{1}{2\tau} \|\hat{\xi} - \bar{\xi}\|^2 + \frac{1}{2\tau} \|\xi - \hat{\xi}\|^2 + \frac{\tilde{\gamma}}{2} \|\xi - \hat{\xi}\|^2 \quad (36)$$

$$H^*(y) + \frac{1}{2\sigma} \|y - \bar{y}\|^2 \geq H^*(\hat{y}) - \langle K \tilde{\xi}, \hat{y} - y \rangle + \frac{1}{2\sigma} \|\hat{y} - \bar{y}\|^2 + \frac{1}{2\sigma} \|y - \hat{y}\|^2 + \frac{\tilde{\delta}}{2} \|y - \hat{y}\|^2. \quad (37)$$

Now, summing (36) and (37), we have after rearrangement

$$\begin{aligned} \mathcal{L}(\hat{\xi}; y) - \mathcal{L}(\xi; \hat{y}) &\leq \frac{1}{2\tau} \|\xi - \bar{\xi}\|^2 - \frac{1 + \tau\tilde{\gamma}}{2\tau} \|\xi - \hat{\xi}\|^2 - \frac{1}{2\tau} \|\bar{\xi} - \hat{\xi}\|^2 \\ &\quad + \frac{1}{2\sigma} \|y - \bar{y}\|^2 - \frac{1 + \sigma\tilde{\delta}}{2\sigma} \|y - \hat{y}\|^2 - \frac{1}{2\sigma} \|\bar{y} - \hat{y}\|^2 \\ &\quad + \langle K(\hat{\xi} - \xi), \hat{y} - \tilde{y} \rangle - \langle K(\hat{\xi} - \bar{\xi}), \hat{y} - y \rangle. \end{aligned} \quad (38)$$

#### 3.1.2 First inequality

Let us now prove the following lemma:



**Lemma 1** *Let  $(\xi_n, y_n)_n$  be generated by Algorithm (24). Then, for any  $n \in \mathbb{N}$ ,  $\tau, \sigma > 0$  and  $0 < \omega \leq \theta$ , we have*

$$\begin{aligned}
\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n) &\leq \frac{1}{2\tau} \|\xi - \xi_n\|^2 + \frac{1}{2\sigma} \|y - y_n\|^2 \\
&\quad - \frac{1}{\omega} \left( \frac{1}{2\tau} \|\xi - \xi_{n+1}\|^2 + \frac{1}{2\sigma} \|y - y_{n+1}\|^2 \right) \\
&\quad + \omega \frac{1}{2\tau} \|\xi_{n-1} - \xi_n\|^2 - \frac{1}{2\tau} \|\xi_n - \xi_{n+1}\|^2 \\
&\quad + \omega \langle K(\xi_{n-1} - \xi_n), y - y_n \rangle - \langle K(\xi_n - \xi_{n+1}), y - y_{n+1} \rangle.
\end{aligned} \tag{39}$$

**Proof.** We specify the six variables in (38), by choosing on one hand

$$\hat{\xi} = \xi_{n+1}, \quad \bar{\xi} = \xi_n, \quad \text{and} \quad \tilde{\xi} = \xi_n + \theta(\xi_n - \xi_{n-1}) \tag{40}$$

for  $1 \geq \theta > 0$  not specified yet, and

$$\hat{y} = y_{n+1}, \quad \bar{y} = y_n, \quad \text{and} \quad \tilde{y} = y_{n+1} \tag{41}$$

on the other hand, which leads to the iterations in (24). After a simplification, we get

$$\begin{aligned}
\mathcal{L}(\xi_{n+1}; y) - \mathcal{L}(\xi; y_{n+1}) &\leq \frac{1}{2\tau} \|\xi - \xi_n\|^2 + \frac{1}{2\sigma} \|y - y_n\|^2 \\
&\quad - \frac{1 + \tau\tilde{\gamma}}{2\tau} \|\xi - \xi_{n+1}\|^2 - \frac{1 + \sigma\tilde{\delta}}{2\sigma} \|y - y_{n+1}\|^2 \\
&\quad - \frac{1}{2\tau} \|\xi_n - \xi_{n+1}\|^2 - \frac{1}{2\sigma} \|y_n - y_{n+1}\|^2 \\
&\quad + \theta \langle K(\xi_{n-1} - \xi_n), y - y_{n+1} \rangle - \langle K(\xi_n - \xi_{n+1}), y - y_{n+1} \rangle.
\end{aligned} \tag{42}$$

Now, we define  $\tau\tilde{\gamma} = \mu > 0$  and  $\sigma\tilde{\delta} = \mu' > 0$ . For any  $n \in \mathbb{N}$ , we set

$$\Delta_n = \frac{1}{2\tau} \|\xi - \xi_n\|^2 + \frac{1}{2\sigma} \|y - y_n\|^2. \tag{43}$$

Hence, we can rewrite (42) with  $\Delta_n$ , which yields

$$\begin{aligned}
\mathcal{L}(\xi_{n+1}; y) - \mathcal{L}(\xi; y_{n+1}) &\leq \Delta_n - (1 + \mu) \Delta_{n+1} - \frac{1}{2\tau} \|\xi_n - \xi_{n+1}\|^2 - \frac{1}{2\sigma} \|y_n - y_{n+1}\|^2 \\
&\quad + \theta \langle K(\xi_{n-1} - \xi_n), y - y_{n+1} \rangle - \langle K(\xi_n - \xi_{n+1}), y - y_{n+1} \rangle \\
&\quad + \frac{\mu - \mu'}{2\sigma} \|y - y_{n+1}\|^2.
\end{aligned} \tag{44}$$

Let us bound the scalar products in (44). For any  $0 < \omega \leq \theta$ , we have the decomposition

$$\begin{aligned}
\theta \langle K(\xi_{n-1} - \xi_n), y - y_{n+1} \rangle &= \omega \langle K(\xi_{n-1} - \xi_n), y - y_n \rangle \\
&\quad + \omega \langle K(\xi_{n-1} - \xi_n), y_n - y_{n+1} \rangle \\
&\quad + (\theta - \omega) \langle K(\xi_{n-1} - \xi_n), y - y_{n+1} \rangle.
\end{aligned} \tag{45}$$

Let us have a closer look at the last two terms. Let  $\alpha > 0$ . Since  $\omega \geq 0$ , we have

$$\begin{aligned} \omega \langle K(\xi_{n-1} - \xi_n), y_n - y_{n+1} \rangle &\leq \omega L_K \|\xi_{n-1} - \xi_n\| \cdot \|y_n - y_{n+1}\| \\ &\leq \omega L_K \left( \frac{\alpha}{2} \|\xi_{n-1} - \xi_n\|^2 + \frac{1}{2\alpha} \|y_n - y_{n+1}\|^2 \right). \end{aligned} \quad (46)$$

Similarly, since  $\theta - \omega \geq 0$ ,

$$(\theta - \omega) \langle K(\xi_{n-1} - \xi_n), y - y_{n+1} \rangle \leq (\theta - \omega) L_K \left( \frac{\alpha}{2} \|\xi_{n-1} - \xi_n\|^2 + \frac{1}{2\alpha} \|y - y_{n+1}\|^2 \right). \quad (47)$$

After simplification, the majoration (44) becomes, thanks to inequalities (46) and (47),

$$\begin{aligned} \mathcal{L}(\xi_{n+1}; y) - \mathcal{L}(\xi; y_{n+1}) &\leq \Delta_n - (1 + \mu) \Delta_{n+1} \\ &\quad + \theta L_K \frac{\alpha}{2} \|\xi_{n-1} - \xi_n\|^2 - \frac{1}{2\tau} \|\xi_n - \xi_{n+1}\|^2 \\ &\quad + \left( \frac{\omega L_K}{2\alpha} - \frac{1}{2\sigma} \right) \|y_n - y_{n+1}\|^2 \\ &\quad + \omega \langle K(\xi_{n-1} - \xi_n), y - y_n \rangle - \langle K(\xi_n - \xi_{n+1}), y - y_{n+1} \rangle \\ &\quad + \left( \frac{(\theta - \omega) L_K}{2\alpha} + \frac{\mu - \mu'}{2\sigma} \right) \|y - y_{n+1}\|^2. \end{aligned} \quad (48)$$

Choose  $\alpha = \omega L_K \sigma$ . Hence, we have  $\omega L_K / \alpha = 1/\sigma$ , so that the  $\|y_n - y_{n+1}\|^2$  term cancels. This leads to:

$$\begin{aligned} \mathcal{L}(\xi_{n+1}; y) - \mathcal{L}(\xi; y_{n+1}) &\leq \Delta_n - (1 + \mu) \Delta_{n+1} \\ &\quad + \omega \frac{\theta L_K^2 \tau \sigma}{2\tau} \|\xi_{n-1} - \xi_n\|^2 - \frac{1}{2\tau} \|\xi_n - \xi_{n+1}\|^2 \\ &\quad + \omega \langle K(\xi_{n-1} - \xi_n), y - y_n \rangle - \langle K(\xi_n - \xi_{n+1}), y - y_{n+1} \rangle \\ &\quad + \left( \frac{\theta - \omega}{\omega} + \mu - \mu' \right) \frac{1}{2\sigma} \|y - y_{n+1}\|^2. \end{aligned} \quad (49)$$

Since  $1 + \mu = 1/\omega + 1 + \mu - 1/\omega$ , we have

$$-(1 + \mu) \Delta_{n+1} = -\frac{1}{\omega} \Delta_{n+1} + \left( \frac{1}{\omega} - \mu - 1 \right) \left( \frac{1}{2\tau} \|\xi - \xi_{n+1}\|^2 + \frac{1}{2\sigma} \|y - y_{n+1}\|^2 \right) \quad (50)$$

so the right-hand side of (49) becomes

$$\begin{aligned} \Delta_n - \frac{1}{\omega} \Delta_{n+1} &+ \omega \frac{\theta L_K^2 \tau \sigma}{2\tau} \|\xi_n - \xi_{n-1}\|^2 - \frac{1}{2\tau} \|\xi_n - \xi_{n+1}\|^2 \\ &\quad + \omega \langle K(\xi_{n-1} - \xi_n), y - y_n \rangle - \langle K(\xi_n - \xi_{n+1}), y - y_{n+1} \rangle \\ &\quad + \left( \frac{1}{\omega} - \mu - 1 \right) \frac{1}{2\tau} \|\xi - \xi_{n+1}\|^2 \\ &\quad + \left( \frac{\theta - \omega}{\omega} + \frac{1}{\omega} - \mu' - 1 \right) \frac{1}{2\sigma} \|y - y_{n+1}\|^2. \end{aligned} \quad (51)$$

It is now time to set conditions on  $\omega$ ,  $\theta$ ,  $\tau$  and  $\sigma$ . First, choose  $\theta$ ,  $\tau$  and  $\sigma$  so that  $\theta L_K^2 \tau \sigma \leq 1$ . Then, choose  $\theta$  so that both  $1/\omega - \mu - 1$  and  $(\theta - \omega)/\omega + 1/\omega - \mu' - 1$

are nonpositive, which implies that

$$\frac{1}{\mu+1} \leq \omega \leq \theta \quad \text{and} \quad \frac{\theta+1}{\mu'+2} \leq \omega \leq \theta. \quad (52)$$

Then we can bound (51) by

$$\begin{aligned} \Delta_n - \frac{1}{\omega} \Delta_{n+1} + \omega \frac{1}{2\tau} \|\xi_{n-1} - \xi_n\|^2 - \frac{1}{2\tau} \|\xi_n - \xi_{n+1}\|^2 \\ + \omega \langle K(\xi_{n-1} - \xi_n), y - y_n \rangle - \langle K(\xi_n - \xi_{n+1}), y - y_{n+1} \rangle. \end{aligned} \quad (53)$$

Eventually, back to (49) we get the wanted inequality

$$\begin{aligned} \mathcal{L}(\xi_{n+1}; y) - \mathcal{L}(\xi; y_{n+1}) \leq \Delta_n - \frac{1}{\omega} \Delta_{n+1} \\ + \omega \frac{1}{2\tau} \|\xi_{n-1} - \xi_n\|^2 - \frac{1}{2\tau} \|\xi_n - \xi_{n+1}\|^2 \\ + \omega \langle K(\xi_{n-1} - \xi_n), y - y_n \rangle - \langle K(\xi_n - \xi_{n+1}), y - y_{n+1} \rangle. \end{aligned} \quad (54)$$

### 3.1.3 Linear convergence of the iterates

Multiplying (54) by  $1/\omega^n$  and summing between  $n = 0$  and  $n = N-1$  (choose  $\xi^{-1} = \xi^0$ ) cancels most of the terms:

$$\begin{aligned} \sum_{n=1}^N \frac{1}{\omega^{n-1}} (\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n)) \leq \Delta_0 - \frac{1}{\omega^N} \Delta_N - \frac{1}{2\tau\omega^{N-1}} \|\xi_{N-1} - \xi_N\|^2 \\ - \frac{1}{\omega^{N-1}} \langle K(\xi_{N-1} - \xi_N), y - y_N \rangle. \end{aligned} \quad (55)$$

Once again, we bound the scalar product: let  $\beta > 0$ ,

$$- \frac{1}{\omega^{N-1}} \langle K(\xi_{N-1} - \xi_N), y - y_N \rangle \leq \frac{L_K}{\omega^{N-1}} \left( \frac{\beta}{2} \|\xi_{N-1} - \xi_N\|^2 + \frac{1}{2\beta} \|y - y_N\|^2 \right) \quad (56)$$

and inequality (55) becomes

$$\begin{aligned} \sum_{n=1}^N \frac{1}{\omega^{n-1}} (\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n)) \leq \Delta_0 - \frac{1}{\omega^N} \Delta_N \\ + \left( \frac{L_K\beta}{2\omega^{N-1}} - \frac{1}{2\tau\omega^{N-1}} \right) \|\xi_{N-1} - \xi_N\|^2 \\ + \frac{L_K}{\omega^{N-1}} \frac{1}{2\beta} \|y - y_N\|^2. \end{aligned} \quad (57)$$

Now choose  $\beta = 1/(L_K\tau)$ , which cancels the  $\|\xi_{N-1} - \xi_N\|^2$  term, and we get

$$\sum_{n=1}^N \frac{1}{\omega^{n-1}} (\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n)) \leq \Delta_0 - \frac{1}{\omega^N} \Delta_N + \frac{L_K^2\tau\sigma}{\omega^{N-1}} \frac{1}{2\sigma} \|y - y_N\|^2. \quad (58)$$

Replacing  $\Delta_0$  and  $\Delta_n$  by their respective definition, we obtain

$$\begin{aligned} \sum_{n=1}^N \frac{1}{\omega^{n-1}} (\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n)) \leq \frac{1}{2\tau} \|\xi - \xi_0\|^2 + \frac{1}{2\sigma} \|y - y_0\|^2 \\ - \frac{1}{\omega^N} \frac{1}{2\tau} \|\xi - \xi_N\|^2 - \frac{1}{\omega^N} (1 - \omega L_K^2\tau\sigma) \frac{1}{2\sigma} \|y - y_N\|^2. \end{aligned} \quad (59)$$

Since  $\omega L_K^2 \tau \sigma \leq \theta L_K^2 \tau \sigma \leq 1$  and  $\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n) \geq 0$  for any  $n \in \mathbb{N}$ , we have

$$\begin{aligned} 0 &\leq \frac{1}{2\tau} \|\xi - \xi_N\|^2 + (1 - \omega L_K^2 \tau \sigma) \frac{1}{2\sigma} \|y - y_N\|^2 + \sum_{n=1}^N \frac{\omega^N}{\omega^{n-1}} (\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n)) \\ &\leq \frac{\omega^N}{2\tau} \|\xi - \xi_0\|^2 + \frac{\omega^N}{2\sigma} \|y - y_0\|^2. \end{aligned} \quad (60)$$

The latter inequality proves the linear convergence of the iterates:

**Corollary 1** *Assume problem (13) has a solution, which is a saddle-point of  $\mathcal{L}$ , denoted by  $(\xi^*, y^*)$ . Let  $(\xi_n, y_n)_n$  be generated by Algorithm (24). Suppose there exist  $\tau$ ,  $\sigma$ ,  $\theta$  and  $\omega$  satisfying both conditions (25) and (26). Then, for any  $N \in \mathbb{N}$ , we have*

$$\|\xi^* - \xi_N\|^2 \leq \omega^N \left( \|\xi^* - \xi_0\|^2 + \frac{\tau}{\sigma} \|y^* - y_0\|^2 \right). \quad (61)$$

Moreover, if  $\omega L_K^2 \tau \sigma \neq 1$ , then we also have

$$\|y^* - y_N\|^2 \leq \frac{\omega^N}{1 - \omega L_K^2 \tau \sigma} \left( \frac{\sigma}{\tau} \|\xi^* - \xi_0\|^2 + \|y^* - y_0\|^2 \right). \quad (62)$$

REMARK: The convergence rate in Corollary 1 can be improved if we use the fact that, by definition,

$$\begin{aligned} \mathcal{L}(\xi_{n+1}; y^*) - \mathcal{L}(\xi^*; y_{n+1}) &= G(\xi_{n+1}) - G(\xi^*) + H^*(y_{n+1}) - H^*(y^*) \\ &\quad + \langle K \xi_{n+1}, y^* \rangle - \langle K \xi^*, y_{n+1} \rangle. \end{aligned} \quad (63)$$

The strong convexity of  $G$  et  $H^*$  and the optimality of  $\xi^*$  and  $y^*$  yield the following inequalities:

$$G(\xi_{n+1}) - G(\xi^*) \geq \langle -K^* y^*, \xi_{n+1} - \xi^* \rangle + \frac{\tilde{\gamma}}{2} \|\xi_{n+1} - \xi^*\|^2 \quad (64)$$

$$H^*(y_{n+1}) - H^*(y^*) \geq \langle K \xi^*, y_{n+1} - y^* \rangle + \frac{\tilde{\delta}}{2} \|y_{n+1} - y^*\|^2. \quad (65)$$

This implies that

$$\frac{\tilde{\gamma}}{2} \|\xi_{n+1} - \xi^*\|^2 + \frac{\tilde{\delta}}{2} \|y_{n+1} - y^*\|^2 \leq \mathcal{L}(\xi_{n+1}; y^*) - \mathcal{L}(\xi^*; y_{n+1}) \quad (66)$$

since the sum of the scalar products cancels. Hence, if we choose not to control the primal-dual gap, choosing  $(\xi, y) = (\xi^*, y^*)$ , in (42) becomes

$$\begin{aligned} 0 &\leq \frac{1}{2\tau} \|\xi^* - \xi_n\|^2 + \frac{1}{2\sigma} \|y^* - y_n\|^2 \\ &\quad - \frac{1 + 2\tau\tilde{\gamma}}{2\tau} \|\xi^* - \xi_{n+1}\|^2 - \frac{1 + 2\sigma\tilde{\delta}}{2\sigma} \|y^* - y_{n+1}\|^2 \\ &\quad - \frac{1}{2\tau} \|\xi_n - \xi_{n+1}\|^2 - \frac{1}{2\sigma} \|y_n - y_{n+1}\|^2 \\ &\quad + \theta \langle K(\xi_{n-1} - \xi_n), y^* - y_{n+1} \rangle - \langle K(\xi_n - \xi_{n+1}), y^* - y_{n+1} \rangle \end{aligned} \quad (67)$$

which means that all the computations from (42) to (60) hold, with  $\mu$  and  $\mu'$  replaced by  $\tilde{\mu} = 2\mu$  and  $\tilde{\mu}' = 2\mu'$  and without  $\mathcal{L}$ -terms, as well as the constraints on the parameters. In others terms, the same computations prove that

**Corollary 2** Assume problem (13) has a solution, which is a saddle-point of  $\mathcal{L}$ , denoted by  $(\xi^*, y^*)$ . Let  $(\xi_n, y_n)_n$  be generated by Algorithm (24). Suppose there exist  $\tau$ ,  $\sigma$ ,  $\theta$  and  $\omega$  satisfying both conditions

$$\max \left\{ \frac{1}{2\tau\tilde{\gamma} + 1}, \frac{1}{2\sigma\tilde{\delta} + 1} \right\} \leq \theta \leq \frac{1}{L_K^2 \tau \sigma}. \quad (68)$$

Then, for any  $\tilde{\omega}$  such that

$$\max \left\{ \frac{1}{2\tau\tilde{\gamma} + 1}, \frac{\theta + 1}{2\sigma\tilde{\delta} + 2} \right\} \leq \tilde{\omega} \leq \theta. \quad (69)$$

Then, for any  $N \in \mathbb{N}$ , we have

$$\|\xi^* - \xi_N\|^2 \leq \tilde{\omega}^N \left( \|\xi^* - \xi_0\|^2 + \frac{\tau}{\sigma} \|y^* - y_0\|^2 \right). \quad (70)$$

Moreover, if  $\tilde{\omega} L_K^2 \tau \sigma \neq 1$ , then we also have

$$\|y^* - y_N\|^2 \leq \frac{\tilde{\omega}^N}{1 - \tilde{\omega} L_K^2 \tau \sigma} \left( \frac{\sigma}{\tau} \|\xi^* - \xi_0\|^2 + \|y^* - y_0\|^2 \right). \quad (71)$$

For given  $\tau$ ,  $\sigma$  and  $\theta$ , the lower bounds  $1/(2\tau\tilde{\gamma} + 1)$  and  $(\theta + 1)/(2\sigma\tilde{\delta} + 2)$  for  $\tilde{\omega}$  are smaller than those for  $\omega$ . Thus, the new rate  $\tilde{\omega}$  can be expected to be better than the global one  $\omega$  (which is called *global* since it also holds for the objective error, as shown in the next paragraph). This will be checked in Subsection 3.2.

### 3.1.4 Ergodic convergence of the objective error

We can now complete the proof of Theorem 1. Dividing (60) by  $\omega^N \neq 0$  and by  $T_N \neq 0$ , we get

$$\begin{aligned} & \frac{1 - \omega}{\omega(1 - \omega^N)} \frac{1}{2\tau} \|\xi - \xi_N\|^2 + \frac{1 - \omega}{\omega(1 - \omega^N)} (1 - \omega L_K^2 \tau \sigma) \frac{1}{2\sigma} \|y - y_N\|^2 \\ & \quad + \frac{1}{T_N} \sum_{n=1}^N \frac{1}{\omega^{n-1}} (\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n)) \\ & \leq \frac{1}{T_N} \frac{1}{2\tau} \|\xi - \xi_0\|^2 + \frac{1}{T_N} \frac{1}{2\sigma} \|y - y_0\|^2. \end{aligned} \quad (72)$$

But, by convexity,

$$\mathcal{L}(\Xi_N; y) - \mathcal{L}(\xi; Y_N) \leq \frac{1}{T_N} \sum_{n=1}^N \frac{1}{\omega^{n-1}} (\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n)) \quad (73)$$

Therefore, (72) becomes

$$\begin{aligned} & \frac{1 - \omega}{\omega(1 - \omega^N)} \frac{1}{2\tau} \|\xi - \xi_N\|^2 + \frac{1 - \omega}{\omega(1 - \omega^N)} (1 - \omega L_K^2 \tau \sigma) \frac{1}{2\sigma} \|y - y_N\|^2 \\ & \quad + \mathcal{L}(\Xi_N; y) - \mathcal{L}(\xi; Y_N) \\ & \leq \frac{1}{T_N} \frac{1}{2\tau} \|\xi - \xi_0\|^2 + \frac{1}{T_N} \frac{1}{2\sigma} \|y - y_0\|^2 \end{aligned} \quad (74)$$

which completes the proof of Theorem 1.

We can now establish the ergodic linear convergence of the objective function. Let  $\mathcal{B}_Z \times \mathcal{B}_Y$  be an open subset of  $Z \times Y$  which contains  $(\xi^*, y^*)$  and  $(\xi_0, y_0)$ . The partial primal-dual gap is bounded by

$$\begin{aligned} \mathcal{G}_{\mathcal{B}_Z \times \mathcal{B}_Y}(\Xi_N; Y_N) &\leq \frac{1}{T_N} \frac{1}{2\tau} \sup_{\xi \in \mathcal{B}_Z} \|\xi - \xi_0\|^2 + \frac{1}{T_N} \frac{1}{2\sigma} \sup_{y \in \mathcal{B}_Y} \|y - y_0\|^2 \\ &\leq \frac{1}{T_N} \left( \frac{(\text{diam}(\mathcal{B}_Z))^2}{2\tau} + \frac{(\text{diam}(\mathcal{B}_Y))^2}{2\sigma} \right) \end{aligned} \quad (75)$$

with  $(T_N)^{-1} \sim \omega^{N-1}/(1-\omega)$ . Hence, we prove the linear convergence of the primal-dual gap applied to the sequences  $(\Xi_N, Y_N)$ .

### 3.2 Choice of parameters

Theorem 1 holds provided one can properly choose the steps  $\tau$  and  $\sigma$  and the relaxation parameter  $\theta$ . We study some particular choices for those parameters and the convergence rate they yield. Since a smaller  $\omega$  leads to a faster convergence, we tune the algorithm parameters to minimize the lower bound of  $\omega$ . Here is how we proceed:

1. Fix  $\tau > 0$ .
2. Find conditions on  $\sigma$  so that inequalities (25) hold.
3. Minimize  $(\theta + 1)/(\sigma\tilde{\delta} + 2)$  with respect to (w.r.t.)  $\theta$  satisfying (25) and w.r.t  $\sigma$  given by the previous step.
4. Compare this minimum to  $1/(\tau\tilde{\gamma} + 1)$  and deduce the lower bound  $\omega^*(\tau)$  for  $\omega$ .
5. Minimize  $\omega^*(\tau)$  and derive the optimal rate  $\omega^*$ .

Since the resulting parameters are compatible with conditions (68) and (69), the left-hand member in (70) yields a better theoretical rate for the convergence of the variables.

Besides, the same computations (with  $\tilde{\gamma}$  and  $\tilde{\delta}$  doubled) may be used to choose the parameters so that the rate  $\tilde{\omega}$  (Corollary 2) is minimal.

#### 3.2.1 Case $\theta = 1$

We first fix  $\theta = 1$ . As shown in [3], this choice is equivalent to the ADMM with an additional proximal term.

Fix  $\tau > 0$ . Replacing  $\theta = 1$  in (25), we obtain that the steps  $\tau$  and  $\sigma$  are constrained as following

$$1 \leq \frac{1}{L_K^2 \tau \sigma} \quad (76)$$

which implies that  $\sigma \leq 1/(L_K^2 \tau)$ . Then, (26) in Theorem 1 states that the convergence rate  $\omega$  satisfies

$$\max \left\{ \frac{1}{\tau\tilde{\gamma} + 1}, \frac{1}{\sigma\tilde{\delta}/2 + 1} \right\} \leq \omega \leq 1 \quad (77)$$

Let us minimize  $1/(\sigma\tilde{\delta}/2 + 1)$  w.r.t.  $\sigma$  satisfying (76). Since the map  $\sigma \mapsto 1/(\sigma\tilde{\delta}/2 + 1)$  is nondecreasing, its minimum is reached when  $\sigma$  is maximal, which leads to

$$\min_{\sigma \text{ subject to (76)}} \left\{ \frac{1}{\sigma\tilde{\delta}/2 + 1} \right\} = \frac{1}{\tilde{\delta}/(2L_K^2 \tau) + 1}. \quad (78)$$

Now, compare it to  $1/(\tau\tilde{\gamma}+1)$ . It is clear that  $1/(\tilde{\delta}/(2L_K^2\tau)+1)$  is greater than  $1/(\tau\tilde{\gamma}+1)$  as soon as  $\tau^2 \geq \tilde{\delta}/(2\tilde{\gamma}L_K^2)$ . Hence, the lower bound  $\omega^*(\tau)$  is given by

$$\omega^*(\tau) = \max \left\{ \frac{1}{\tau\tilde{\gamma}+1}, \frac{1}{\tilde{\delta}/(2L_K^2\tau)+1} \right\} = \begin{cases} \frac{1}{\tau\tilde{\gamma}+1} & \text{if } 0 < \tau < \sqrt{\tilde{\delta}/(2\tilde{\gamma}L_K^2)} \\ \frac{1}{\tilde{\delta}/(2L_K^2\tau)+1} & \text{if } \tau \geq \sqrt{\tilde{\delta}/(2\tilde{\gamma}L_K^2)} \end{cases} \quad (79)$$

which is minimal for  $\tau^* = \sqrt{\tilde{\delta}/(2\tilde{\gamma}L_K^2)}$  and leads to the optimal rate

$$\omega^* = \omega^*(\tau^*) = \frac{1}{\sqrt{(\tilde{\gamma}\tilde{\delta})/(2L_K^2)+1}} = \frac{1}{\sqrt{1/(2\kappa_F)+1}}. \quad (80)$$

This rate is reached for

$$\tau = \tau^* = \sqrt{\frac{\tilde{\delta}}{2\tilde{\gamma}L_K^2}} \quad \text{and} \quad \sigma = \frac{1}{L_K^2\tau^*} = \sqrt{\frac{2\tilde{\gamma}}{\tilde{\delta}L_K^2}}. \quad (81)$$

One can check that the same choice for  $\tau$  and  $\sigma$  yield the minimal value for the solution error rate  $\tilde{\omega}$ , which is

$$\tilde{\omega}^* = \frac{1}{2\sqrt{(\tilde{\gamma}\tilde{\delta})/(2L_K^2)+1}} = \frac{1}{\sqrt{2/\kappa_F+1}}. \quad (82)$$

In other terms, in the case where  $\theta = 1$ , the best choice for the global rate  $\omega$  and for the solution error rate  $\tilde{\omega}$  coincide.

### 3.2.2 The best convergence rate ( $\theta < 1$ )

In this section, we want to derive the best convergence rate given the constraints in (1).

**Theorem 2** *The best convergence rate in Theorem 1 is obtained when choosing*

$$\tau = \frac{\tilde{\delta}}{2L_K^2} \left( 1 + \sqrt{1 + \frac{4L_K^2}{\tilde{\gamma}\tilde{\delta}}} \right) \quad \text{and} \quad \sigma = \frac{\tilde{\gamma}}{2L_K^2} \left( 1 + \sqrt{1 + \frac{4L_K^2}{\tilde{\gamma}\tilde{\delta}}} \right) \quad (83)$$

and, if  $\kappa_F = L_K^2/(\tilde{\gamma}\tilde{\delta})$ ,

$$\theta = \frac{\sqrt{1 + (4L_K^2)/(\tilde{\gamma}\tilde{\delta})} - 1}{\sqrt{1 + (4L_K^2)/(\tilde{\gamma}\tilde{\delta})} + 1} = \frac{\sqrt{1 + 4\kappa_F} - 1}{\sqrt{1 + 4\kappa_F} + 1} < 1 \quad (84)$$

which satisfy  $\tau\tilde{\gamma} = \sigma\tilde{\delta}$ . The resulting rate is  $\omega^* = \theta$ .

**Proof** Fix  $\tau > 0$  and find out which conditions  $\sigma$  must satisfy to ensure the existence of  $\theta$  satisfying (25). There exists  $\theta$  satisfying (25) if

$$\frac{1}{\tau\tilde{\gamma}+1} \leq \frac{1}{L_K^2\tau\sigma} \quad \text{and} \quad \frac{1}{\sigma\tilde{\delta}+1} \leq \frac{1}{L_K^2\tau\sigma}. \quad (85)$$

which also reads

$$\sigma \leq \frac{1}{L_K^2 \tau} + \frac{\tilde{\gamma}}{L_K^2} \quad \text{and} \quad (L_K^2 \tau - \tilde{\delta}) \sigma \leq 1. \quad (86)$$

Let us determine conditions on  $\sigma$  so that these inequalities hold. If  $L_K^2 \tau - \tilde{\delta} \leq 0$ , i.e.  $\tau \leq \tilde{\delta}/L_K^2$ , the second inequality is always true. Hence, let us study the case  $L_K^2 \tau - \tilde{\delta} > 0$ , i.e.  $\tau > \tilde{\delta}/L_K^2$ . It implies that  $\sigma$  must satisfy both majorations

$$\sigma \leq \frac{1}{L_K^2 \tau} + \frac{\tilde{\gamma}}{L_K^2} \quad \text{and} \quad \sigma \leq \frac{1}{L_K^2 \tau - \tilde{\delta}}. \quad (87)$$

Let us compare these two bounds. Since

$$\frac{1}{L_K^2 \tau} + \frac{\tilde{\gamma}}{L_K^2} - \frac{1}{L_K^2 \tau - \tilde{\delta}} = \frac{\tilde{\gamma} L_K^2 \tau^2 - \tilde{\gamma} \tilde{\delta} \tau - \tilde{\delta}}{L_K^2 \tau (L_K^2 \tau - \tilde{\delta})} \quad (88)$$

with  $L_K^2 \tau (L_K^2 \tau - \tilde{\delta})$  positive,  $1/(L_K^2 \tau) + \tilde{\gamma}/L_K^2$  is greater than  $1/(L_K^2 \tau - \tilde{\delta})$  iff  $\tilde{\gamma} L_K^2 \tau^2 - \tilde{\gamma} \tilde{\delta} \tau - \tilde{\delta} \geq 0$ , i.e. iff  $\tau \geq \tau^*$ , given by

$$\tau^* = \frac{\tilde{\delta}}{2L_K^2} \left( 1 + \sqrt{1 + \frac{4L_K^2}{\tilde{\gamma}\tilde{\delta}}} \right) > \frac{\tilde{\delta}}{L_K^2}. \quad (89)$$

Therefore, for any  $\tilde{\delta}/L_K^2 < \tau \leq \tau^*$ , (87) becomes  $\sigma \leq 1/(L_K^2 \tau) + \tilde{\gamma}/L_K^2$ . If  $\tau > \tau^*$ , (87) reads  $\sigma \leq 1/(L_K^2 \tau - \tilde{\delta})$ . As a conclusion, we have the following upper bounds for  $\sigma$ :

$$\sigma \leq \begin{cases} \frac{1}{L_K^2 \tau} + \frac{\tilde{\gamma}}{L_K^2} & \text{if } 0 < \tau \leq \tau^* \\ \frac{1}{L_K^2 \tau - \tilde{\delta}} & \text{if } \tau^* < \tau. \end{cases} \quad (90)$$

Now, fix  $\sigma$  satisfying (90) and let us minimize  $(\theta + 1)/(\sigma \tilde{\delta} + 2)$  subject to (25). The map  $\theta \mapsto (\theta + 1)/(\sigma \tilde{\delta} + 2)$  is minimal when  $\theta$  is minimal. Hence, let us determine the lower bound of  $\theta$ , which is given by

$$\max \left\{ \frac{1}{\tau \tilde{\gamma} + 1}, \frac{1}{\sigma \tilde{\delta} + 1} \right\}. \quad (91)$$

First, remark that, if  $\tau > \tilde{\delta}/L_K^2$ , then

$$\frac{\tilde{\delta}}{L_K^2 \tau - \tilde{\delta}} \leq \tau \tilde{\gamma} \quad \iff \quad \tilde{\gamma} L_K^2 \tau^2 - \tilde{\gamma} \tilde{\delta} \tau - \tilde{\delta} \geq 0 \quad \iff \quad \tau \geq \tau^*. \quad (92)$$

Suppose that  $\tau > \tau^*$ , which implies that  $\tau > \tilde{\delta}/L_K^2$ . Since  $\sigma$  is bounded from above by  $1/(L_K^2 \tau - \tilde{\delta})$ , we deduce that  $\sigma \tilde{\delta} \leq \tau \tilde{\gamma}$ , which yields

$$\max \left\{ \frac{1}{\tau \tilde{\gamma} + 1}, \frac{1}{\sigma \tilde{\delta} + 1} \right\} = \frac{1}{\sigma \tilde{\delta} + 1} \quad \text{if} \quad 0 < \sigma \leq \frac{1}{L_K^2 \tau - \tilde{\delta}}. \quad (93)$$

Now, let us consider the case  $\tau \leq \tau^*$ . Since

$$\frac{1}{L_K^2 \tau} + \frac{\tilde{\gamma}}{L_K^2} \geq \frac{\tau \tilde{\gamma}}{\tilde{\delta}} \quad \iff \quad \tilde{\gamma} L_K^2 \tau^2 - \tilde{\gamma} \tilde{\delta} \tau - \tilde{\delta} \leq 0 \quad \iff \quad \tau \leq \tau^*, \quad (94)$$

we deduce that

$$\max \left\{ \frac{1}{\tau \tilde{\gamma} + 1}, \frac{1}{\sigma \tilde{\delta} + 1} \right\} = \begin{cases} \frac{1}{\sigma \tilde{\delta} + 1} & \text{if } 0 < \sigma \leq \frac{\tau \tilde{\gamma}}{\tilde{\delta}} \\ \frac{1}{\tau \tilde{\gamma} + 1} & \text{if } \frac{\tau \tilde{\gamma}}{\tilde{\delta}} < \sigma \leq \frac{1}{L_K^2 \tau} + \frac{\tilde{\gamma}}{L_K^2}. \end{cases} \quad (95)$$



Let us minimize  $(\theta+1)/(\sigma\tilde{\delta}+2)$  w.r.t. to  $\sigma$ , when  $\theta$  is equal to its lower bound  $\theta^*(\sigma)$ , given by (93) and (95). This leads to minimize the following quantity w.r.t.  $\sigma$ :

$$\frac{\theta^*(\sigma)+1}{\sigma\tilde{\delta}+2} = \begin{cases} \frac{1}{\sigma\tilde{\delta}+1} & \text{if } \tau > \tau^* \text{ or } \left( \tau \leq \tau^* \text{ and } 0 < \sigma \leq \frac{\tau\tilde{\gamma}}{\tilde{\delta}} \right) \\ \frac{1}{\tau\tilde{\gamma}+1} \frac{\tau\tilde{\gamma}+2}{\sigma\tilde{\delta}+2} & \text{if } \left( \tau \leq \tau^* \text{ and } \frac{\tau\tilde{\gamma}}{\tilde{\delta}} < \sigma \leq \frac{1}{L_K^2\tau} + \frac{\tilde{\gamma}}{L_K^2} \right). \end{cases} \quad (96)$$

In both cases, the minimum is reached when  $\sigma$  is maximal, equal to its upper bound given by (90). Hence,

$$\min_{\substack{\sigma \text{ subject to (90)} \\ \theta \text{ subject to (25)}}} \frac{\theta+2}{\sigma+1} = \begin{cases} 1 - \frac{\tilde{\delta}}{L_K^2\tau} & \text{if } \tau > \tau^* \\ \min \left\{ \frac{1}{\tau\tilde{\gamma}+1}, \frac{1}{\tau\tilde{\gamma}+1} \frac{\tau\tilde{\gamma}+2}{\tilde{\delta}/(L_K^2\tau) + \tilde{\delta}\tilde{\gamma}/L_K^2 + 2} \right\} & \text{if } \tau \leq \tau^*. \end{cases} \quad (97)$$

Compare it to  $1/(\tau\tilde{\gamma}+1)$ , and deduce the lower bound  $\omega^*(\tau)$ :

$$\omega^*(\tau) = \max \left\{ \frac{1}{\tau\tilde{\gamma}+1}, \min_{\substack{\sigma \text{ subject to (90)} \\ \theta \text{ subject to (25)}}} \left\{ \frac{\theta+2}{\sigma+1} \right\} \right\}. \quad (98)$$

Thanks to (92), it follows that

$$\omega^*(\tau) = \begin{cases} 1 - \frac{\tilde{\delta}}{L_K^2\tau} & \text{if } \tau > \tau^* \\ \frac{1}{\tau\tilde{\gamma}+1} & \text{if } \tau \leq \tau^*. \end{cases} \quad (99)$$

In the second case,  $\sigma\tilde{\delta}$  is supposed to be greater than  $\tau\tilde{\gamma}$ , so  $(\theta^*+1)/(\sigma\tilde{\delta}+2)$  is always smaller than  $1/(\tau\tilde{\gamma}+1)$ . Therefore, the best rate is bounded from below by  $1/(\tau\tilde{\gamma}+1)$ . Eventually, we get the following best rate:

$$\omega^*(\tau) = \begin{cases} 1 - \frac{\tilde{\delta}}{L_K^2\tau} & \text{if } \tau > \tau^* \\ \frac{1}{\tau\tilde{\gamma}+1} & \text{if } \tau \leq \tau^* \end{cases} \quad (100)$$

which is minimal for  $\tau = \tau^*$ . This eventually leads to the best rate

$$\omega^* = 1 - \frac{\tilde{\delta}}{L_K^2\tau^*} = \frac{1}{\tau^*\tilde{\gamma}+1} = \frac{\sqrt{1+(4L_K^2)/(\tilde{\gamma}\tilde{\delta})} - 1}{\sqrt{1+(4L_K^2)/(\tilde{\gamma}\tilde{\delta})} + 1} = \frac{\sqrt{1+4\kappa_F} - 1}{\sqrt{1+4\kappa_F} + 1} \quad (101)$$

obtained when  $\tau = \tau^*$  and  $\sigma = \tau^*\tilde{\gamma}/\tilde{\delta}$ .  $\square$

This choice leads to the following value for the solution error rate  $\tilde{\omega}$ :

$$\tilde{\omega} = \frac{1}{2\tau^*\tilde{\gamma}+1} = \frac{\sqrt{1+(4L_K^2)/(\tilde{\gamma}\tilde{\delta})} - 1}{\sqrt{1+(4L_K^2)/(\tilde{\gamma}\tilde{\delta})} + 3} = \frac{\sqrt{1+4\kappa_F} - 1}{\sqrt{1+4\kappa_F} + 3}. \quad (102)$$

Once again, the same computations prove that the best solution error rate  $\tilde{\omega}$  is reached when

$$\tilde{\tau} = \frac{\tilde{\delta}}{L_K^2} \left( 1 + \sqrt{1 + \frac{L_K^2}{\tilde{\gamma}\tilde{\delta}}} \right) \quad \text{and} \quad \tilde{\sigma} = \frac{\tilde{\gamma}}{L_K^2} \left( 1 + \sqrt{1 + \frac{L_K^2}{\tilde{\gamma}\tilde{\delta}}} \right) \quad (103)$$

and leads to

$$\tilde{\omega}^* = \theta = \frac{\sqrt{1 + L_K^2/(\tilde{\gamma}\tilde{\delta})} - 1}{\sqrt{1 + L_K^2/(\tilde{\gamma}\tilde{\delta})} + 1} = \frac{\sqrt{1 + \kappa_F} - 1}{\sqrt{1 + \kappa_F} + 1} < \omega^*. \quad (104)$$

### 3.3 Overrelaxation on the dual variable

Thanks to the symmetry of Problem (13), similar results still hold if the relaxation is done on the dual variable  $y$  instead of the primal variable  $\xi$ , namely if the updates are replaced by

$$\begin{cases} \xi_{n+1} = \text{prox}_{\tau G}(\xi_n - \tau K^* \bar{y}_n) \\ y_{n+1} = \text{prox}_{\sigma H^*}(y_n + \sigma K \xi_{n+1}) \\ \bar{y}_{n+1} = y_{n+1} + \theta(y_{n+1} - y_n). \end{cases} \quad (105)$$

As seen in (11), such an overrelaxation will be useful for the analysis of the ADMM. It is equivalent to inverting the role of the dual and the primal variables. Indeed, Problem (13) can be rewritten

$$\min_{y \in Y} \sup_{\xi \in X} \left\{ H^*(y) - \langle K^* y, \xi \rangle - G(\xi) \right\} \quad (106)$$

which shares the same regularity assumptions as Problem (13). Hence, applying Theorem 1 yields the following result:

**Theorem 3** *Assume problem (13) has a solution, which is a saddle-point of  $\mathcal{L}$ , denoted by  $(\xi^*, y^*)$ . Choose  $\tau > 0$ ,  $\sigma > 0$  and  $0 < \theta \leq 1$  such that*

$$\max \left\{ \frac{1}{\tau \tilde{\gamma} + 1}, \frac{1}{\sigma \tilde{\delta} + 1} \right\} \leq \theta \leq \frac{1}{L_K^2 \tau \sigma}. \quad (107)$$

*Then, for any  $\omega$  such that*

$$\max \left\{ \frac{\theta + 1}{\tau \tilde{\gamma} + 2}, \frac{1}{\sigma \tilde{\delta} + 1} \right\} \leq \omega \leq \theta \quad (108)$$

*we have the following majoration for any  $N \in \mathbb{N}$  and any  $(\xi, y) \in \Xi \times Y$ :*

$$\begin{aligned} (1 - \omega L_K^2 \tau \sigma) \frac{1}{2\tau} \|\xi_N - \xi\|^2 + \frac{1}{2\sigma} \|y_N - y\|^2 \\ + \sum_{n=1}^N \frac{\omega^n}{\omega^{n-1}} (\mathcal{L}(\xi_n; y) - \mathcal{L}(\xi; y_n)) \\ \leq \frac{\omega^N}{2\tau} \|\xi_0 - \xi\|^2 + \frac{\omega^N}{2\sigma} \|y_0 - y\|^2. \end{aligned} \quad (109)$$

*Now, define*

$$T_N := \sum_{n=1}^N \frac{1}{\omega^{n-1}} = \frac{1 - \omega^N}{\omega^{N-1}(1 - \omega)} \quad (110)$$

*and let*

$$\Xi_N := \frac{1}{T_N} \sum_{n=1}^N \frac{1}{\omega^{n-1}} \xi_n \quad \text{and} \quad Y_N := \frac{1}{T_N} \sum_{n=1}^N \frac{1}{\omega^{n-1}} y_n. \quad (111)$$

Then we have the following bound for any  $(\xi, y) \in \Xi \times Y$ :

$$\begin{aligned} & \frac{1-\omega}{\omega(1-\omega^N)} (1-\omega L_K^2 \tau \sigma) \frac{1}{2\tau} \|\xi - \xi_N\|^2 + \frac{1-\omega}{\omega(1-\omega^N)} \frac{1}{2\sigma} \|y - y_N\|^2 \\ & \quad + \mathcal{L}(\Xi_N; y) - \mathcal{L}(\xi; Y_N) \\ & \leq \frac{1}{T_N} \frac{1}{2\tau} \|\xi - \xi_0\|^2 + \frac{1}{T_N} \frac{1}{2\sigma} \|y - y_0\|^2. \end{aligned} \quad (112)$$

Note that the conditions on the parameters now slightly differ from the previous case. A variant can be found in [5, Appendix C2].

As in the previous case where the overrelaxation is done over the primal variable, we can prove the following result for the linear convergence of the solution errors:

**Corollary 3** *Assume problem (13) has a solution, which is a saddle-point of  $\mathcal{L}$ , denoted by  $(\xi^*, y^*)$ . Let  $(\xi_n, y_n)_n$  be generated by Algorithm (24). Suppose there exist  $\tau$ ,  $\sigma$ ,  $\theta$  and  $\omega$  satisfying both conditions*

$$\max \left\{ \frac{1}{2\tau\tilde{\gamma} + 1}, \frac{1}{2\sigma\tilde{\delta} + 1} \right\} \leq \theta \leq \frac{1}{L_K^2 \tau \sigma}. \quad (113)$$

Then, for any  $\tilde{\omega}$  such that

$$\max \left\{ \frac{\theta + 1}{2\tau\tilde{\gamma} + 1}, \frac{1}{2\sigma\tilde{\delta} + 2} \right\} \leq \tilde{\omega} \leq \theta. \quad (114)$$

Then, for any  $N \in \mathbb{N}$ , we have

$$\|y^* - y_N\|^2 \leq \tilde{\omega}^N \left( \|y^* - y_0\|^2 + \frac{\tau}{\sigma} \|\xi^* - \xi_0\|^2 \right). \quad (115)$$

Moreover, if  $\tilde{\omega} L_K^2 \tau \sigma \neq 1$ , then we also have

$$\|\xi^* - \xi_N\|^2 \leq \frac{\tilde{\omega}^N}{1 - \tilde{\omega} L_K^2 \tau \sigma} \left( \frac{\sigma}{\tau} \|y^* - y_0\|^2 + \|\xi^* - \xi_0\|^2 \right). \quad (116)$$

Similar computations as in the previous section show that the best rate  $\omega^*$  is achieved when choosing the following parameters:

$$\tau = \frac{\tilde{\delta}}{2L_K^2} \left( 1 + \sqrt{1 + \frac{4L_K^2}{\tilde{\gamma}\tilde{\delta}}} \right) \quad \text{and} \quad \sigma = \frac{\tilde{\gamma}}{2L_K^2} \left( 1 + \sqrt{1 + \frac{4L_K^2}{\tilde{\gamma}\tilde{\delta}}} \right) \quad (117)$$

and, with  $\kappa_F = L_K^2 / (\tilde{\gamma}\tilde{\delta})$ ,

$$\theta = \frac{\sqrt{1 + (4L_K^2)/(\tilde{\gamma}\tilde{\delta})} - 1}{\sqrt{1 + (4L_K^2)/(\tilde{\gamma}\tilde{\delta})} + 1} = \frac{\sqrt{1 + 4\kappa_F} - 1}{\sqrt{1 + 4\kappa_F} + 1} < 1 \quad (118)$$

which leads to  $\omega^* = \theta$ .

## 4 Application : convergence rate for the ADMM in the smooth case

As the ADMM is nothing but a particular instance of the oPDHG method with additional constraints on the parameter choice, its convergence rate is expected to be worse than that of the latter. In subsection 4.1, it will indeed be derived from the computations of the previous section and shown to be greater than that of the oPDHG method.

However, as we will show it in the subsection 4.2, it is possible to recover the same convergence rate as in the oPDHG method by introducing a slight modification in the ADMM iterations.

### 4.1 Unaccelerated ADMM

As recalled in Section 2.2, the ADMM iterations (9) are equivalent to the oPDHG iterations (11) applied to the functions  $G = g_A$  and  $H = h$ , and to the identity operator  $K = \text{Id}$ , of norm  $L_K = 1$ . The functions  $G$  and  $H$  are respectively  $\tilde{\gamma} = \gamma/L_A^2$ -convex and  $\tilde{\delta} = \delta$ -convex. Note that the relaxation is done on the dual variable, of parameter  $\theta = 1$ . The stepsize for the primal (resp. dual) proximal ascent is  $\tau > 0$  (resp.  $\sigma = 1/\tau$ ).

#### 4.1.1 Ergodic linear convergence

Apply Theorem 3. Parameters  $\theta$  and  $\sigma$  being constrained as stated above, Theorem 3 ensures that, provided one can find  $\tau > 0$  such that

$$\max \left\{ \frac{1}{\tau\gamma/L_A^2 + 1}, \frac{1}{\delta/\tau + 1} \right\} \leq 1 \quad (119)$$

for any  $\omega$  such that

$$\max \left\{ \frac{2}{\tau\gamma/L_A^2 + 2}, \frac{1}{\delta/\tau + 1} \right\} \leq \omega \leq 1 \quad (120)$$

we have the following bound for any  $(x, y) \in X \times Y$ :

$$\begin{aligned} & \frac{(1-\omega)^2}{\omega(1-\omega^N)} \frac{1}{2\tau} \|Ax - Ax_N\|^2 + \frac{1-\omega}{\omega(1-\omega^N)} \frac{1}{2/\tau} \|y - y_N\|^2 \\ & \quad + \mathcal{L}(Ax_N; y) - \mathcal{L}(Ax; Y_N) \\ & \leq \frac{1}{T_N} \frac{1}{2\tau} \|Ax - Ax_0\|^2 + \frac{1}{T_N} \frac{1}{2/\tau} \|y - y_0\|^2. \end{aligned} \quad (121)$$

First note that, if we apply this inequality to  $(x, y) = (x^*, y^*)$ , then its left-hand side is nonnegative. Hence, the linear convergence of the dual iterates comes naturally. However, though the strong convexity ensures the convergence of the primal iterates  $x_N$ , their convergence speed is not clear. We can solely estimate the convergence of  $Ax_N$ , which is linear. Thanks to

$$z_{N+1} - z^* = Ax_{N+1} - Ax^* + \tau(y^* - y_{N+1}) + \tau(y_N - y^*) \quad (122)$$

we can nevertheless deduce the linear convergence of the primal iterates  $z_N$ . This also implies the linear convergence for the feasibility error  $Ax_N - z_N$ . If we now apply (121) to  $x = x^*$ , using

$$f(x^*) = \mathcal{L}(Ax^*; y^*) = \sup_{y \in Y} \mathcal{L}(Ax^*; y) \geq \mathcal{L}(Ax^*; Y_N) \quad (123)$$

we get for any  $y \in Y$

$$\begin{aligned} & \frac{(1-\omega)^2}{\omega(1-\omega^N)} \frac{1}{2\tau} \|Ax^* - Ax_N\|^2 + \frac{1-\omega}{\omega(1-\omega^N)} \frac{1}{2/\tau} \|y - y_N\|^2 \\ & \quad + \mathcal{L}(AX_N; y) - f(x^*) \\ & \leq \frac{1}{T_N} \frac{1}{2\tau} \|Ax^* - Ax_0\|^2 + \frac{1}{T_N} \frac{1}{2/\tau} \|y - y_0\|^2. \end{aligned} \quad (124)$$

Let define  $y^*(X_N) \in Y$  as

$$y^*(X_N) = \arg \max_{y \in Y} \mathcal{L}(AX_N; y) \quad (125)$$

so that  $\mathcal{L}(AX_N; y^*(X_N)) = f(X_N)$ . The left-hand side of Inequality (124) applied to  $y = y^*(X_N)$  is then nonnegative and implies in particular that

$$\begin{aligned} & \frac{1-\omega}{\omega(1-\omega^N)} \frac{1}{2/\tau} \|y^*(X_N) - y_N\|^2 \\ & \leq \frac{1}{T_N} \frac{1}{2\tau} \|Ax^* - Ax_0\|^2 + \frac{1}{T_N} \frac{1}{2/\tau} \|y^*(X_N) - y_0\|^2 \\ & \leq \frac{1}{T_N} \frac{1}{2\tau} \|Ax^* - Ax_0\|^2 + \frac{1}{T_N} \frac{1}{1/\tau} (\|y^*(X_N) - y_N\|^2 + \|y_N - y_0\|^2) \end{aligned} \quad (126)$$

which yields

$$\begin{aligned} & \left( \frac{1-\omega}{\omega(1-\omega^N)} - \frac{2}{T_N} \right) \frac{1}{2/\tau} \|y^*(X_N) - y_N\|^2 \\ & \leq \frac{1}{T_N} \frac{1}{2\tau} \|Ax^* - Ax_0\|^2 + \frac{1}{T_N} \frac{1}{\tau} \|y_N - y_0\|^2. \end{aligned} \quad (127)$$

This inequality is nonnegative for any sufficiently large  $N$ . Thus, we deduce that the quantity  $\|y^*(X_N) - y_N\|$  goes to zero, since  $(y_N)_N$  has been proved to converge. Hence, the quantity  $\|y^*(X_N) - y_0\|$  is bounded and we have proved the ergodic linear convergence of the ADMM in terms of objective error.

#### 4.1.2 Convergence rate

Let us estimate the best convergence rate which can be achieved by the ADMM. Condition (119) is always true. Hence, for any  $\tau > 0$ , the convergence rate satisfies

$$\max \left\{ \frac{1}{(\tau\gamma)/(2L_A^2) + 1}, \frac{1}{\delta/\tau + 1} \right\} \leq \omega \leq 1. \quad (128)$$

The lower bound is equal to  $1/((\tau\gamma)/(2L_A^2) + 1)$  when  $\tau \leq \sqrt{(2\delta L_A^2)/\gamma}$  and is equal to  $1/(\delta/\tau + 1)$  otherwise. This leads to the best rate

$$\omega^* = \frac{1}{\sqrt{(\gamma\delta)/(2L_A^2)} + 1} = \frac{1}{\sqrt{1/(2\kappa_f)} + 1} \quad \text{reached when} \quad \tau = \sqrt{\frac{2\delta L_A^2}{\gamma}}. \quad (129)$$

We call this parameter the *optimal parameter* for the ADMM. Using this parameter also yields the following theoretical rate for the dual variable and  $Ax_n$ , given by Corollary 3:

$$\tilde{\omega} = \max \left\{ \frac{1}{(\tau\gamma)/L_A^2 + 1}, \frac{1}{2\delta/\tau + 1} \right\} = \frac{1}{\sqrt{(2\gamma\delta)/L_A^2} + 1} = \frac{1}{\sqrt{2/\kappa_f} + 1}. \quad (130)$$

This value can be easily proved to be the optimal one for  $\tilde{\omega}$ .

## 4.2 Accelerated ADMM

We propose to relax the choice of step  $\tau$  in the updates of  $z$  and of  $y$  in the ADMM. Replacing  $\tau$  by  $\tau' \leq \tau$  in these two updates leads to the following algorithm:

$$\begin{cases} x_{n+1} = \arg \min_{x \in X} \left\{ g(x) + \langle Ax, y_n \rangle + \frac{1}{2\tau} \|Ax - z_n\|^2 \right\} \\ z_{n+1} = \arg \min_{z \in X} \left\{ h(z) - \langle z, y_n \rangle + \frac{1}{2\tau'} \|Ax_{n+1} - z\|^2 \right\} \\ y_{n+1} = y_n + \frac{1}{\tau'} (Kx_{n+1} - z_{n+1}). \end{cases} \quad (131)$$

### 4.2.1 Equivalent oPDHG

Following the same computations as in 2.2, we show that iterations in Algorithm (131) are equivalent to those of the following oPDHG algorithm

$$\begin{cases} \xi^{n+1} = \text{prox}_{\tau g_A}(\xi^n - \tau \bar{y}^n) \\ y^{n+1} = \text{prox}_{h^*/\tau'}(y^n + \xi^{n+1}/\tau') \\ \bar{y}^{n+1} = y^{n+1} + \frac{\tau'}{\tau} (y^{n+1} - y^n) \end{cases} \quad (132)$$

where the relaxation parameter  $\theta = \tau'/\tau$  is linked to the ascent steps  $\tau$  and  $\sigma = 1/\tau'$ .

Once again, Theorem 3 reads for any suitable  $\omega$ ,  $\tau$  and  $\tau'$ :

$$\begin{aligned} 0 \leq \frac{1-\omega}{\omega(1-\omega^N)} (1-\omega\tau/\tau') \frac{1}{2\tau} \|Ax - Ax_N\|^2 + \frac{1-\omega}{\omega(1-\omega^N)} \frac{1}{2/\tau'} \|y - y_N\|^2 \\ + \mathcal{L}(AX_N; y) - \mathcal{L}(Ax; Y_N) \\ \leq \frac{1}{T_N} \frac{1}{2\tau} \|Ax - Ax_0\|^2 + \frac{1}{T_N} \frac{1}{2/\tau'} \|y - y_0\|^2 \end{aligned} \quad (133)$$

which yields a linear convergence in terms of objective error (in an ergodic sense). However, the best convergence rate achieved by the algorithm is expected to be better than that of the unaccelerated ADMM. Indeed, introducing the relaxed step  $\tau'$  add a degree of freedom in the constraints over the value of  $\omega$ . Hence, it is minimized over a larger set and its minimal value is thus smaller.

Similarly to the unaccelerated case, (133) ensures the linear convergence of the dual iterates. If  $1 - \omega\tau/\tau'$  does not cancel, it also implies the linear convergence of the primal iterates ( $z_N$ ). Otherwise, we lose the control on the convergence of ( $Ax_N$ ), thus on that of ( $z_N$ ).

### 4.2.2 Convergence rate

Let us derive the best convergence rate for Algorithm (131). We may use Theorem 3, which ensures that steps  $\tau$  and  $\tau'$  are constrained by the relations

$$\max \left\{ \frac{1}{\tau\gamma/L_A^2 + 1}, \frac{1}{\delta/\tau' + 1} \right\} \leq \frac{\tau'}{\tau} \leq 1 \quad (134)$$

and that the convergence rate is constrained by

$$\max \left\{ \frac{\tau'/\tau + 1}{\tau\gamma/L_A^2 + 2}, \frac{1}{\delta/\tau' + 1} \right\} \leq \omega \leq \frac{\tau'}{\tau}. \quad (135)$$

Hence, it is sufficient to find  $(\tau, \tau')$  satisfying both (134) and (135) which minimize the left-hand member in the latter.

One can also first use the remark made after Corollary 3. If no constraint on  $\theta$  is made, then the best rate is achieved when

$$\tau = \frac{\delta}{2} \left( 1 + \sqrt{1 + \frac{4L_A^2}{\gamma\delta}} \right) \quad \text{and} \quad \sigma = \frac{1}{\tau'} = \frac{\gamma}{2L_A^2} \left( 1 + \sqrt{1 + \frac{4L_A^2}{\gamma\delta}} \right) \quad (136)$$

and

$$\theta = \frac{\sqrt{1 + 4L_A^2/(\gamma\delta)} - 1}{\sqrt{1 + 4L_A^2/(\gamma\delta)} + 1}. \quad (137)$$

Let us check that such a choice satisfy  $\theta = \tau'/\tau$ . First, we have

$$\tau' = \frac{2L_A^2/\gamma}{1 + \sqrt{1 + 4L_A^2/(\gamma\delta)}} = \frac{(\delta/2)(\sqrt{1 + 4L_A^2/(\gamma\delta)} + 1)(\sqrt{1 + 4L_A^2/(\gamma\delta)} - 1)}{1 + \sqrt{1 + 4L_A^2/(\gamma\delta)}} \quad (138)$$

which implies that

$$\frac{\tau'}{\tau} = \frac{\sqrt{1 + 4L_A^2/(\gamma\delta)} - 1}{\sqrt{1 + 4L_A^2/(\gamma\delta)} + 1}. \quad (139)$$

Hence, these parameters can be chosen for the accelerated ADMM, and yields to the best rate. Thus, they are called *optimal parameters* for the accelerated ADMM. With this parameter choice, we have  $\omega^* = \theta$ . Note that the resulting rate is the same as the best one expected when applying the oPDHG on Problem (6). However, unlike in the oPDHG algorithm, this choice implies a loss of control on both  $x$ -iterates and  $z$ -iterates. Moreover, this choice leads to the following rate  $\tilde{\omega}$ :

$$\tilde{\omega} = \frac{1}{2\delta/(\tau')^* + 1} = \frac{\sqrt{1 + (4L_A^2)/(\gamma\delta)} - 1}{\sqrt{1 + (4L_A^2)/(\gamma\delta)} + 3} = \frac{\sqrt{1 + \kappa_f} - 1}{\sqrt{1 + \kappa_f} + 3} \quad (140)$$

To minimize the latter rate, we use the previous computations with  $\gamma$  and  $\delta$  doubled, which leads to the parameter choice

$$\tau' = \frac{\delta}{2} \left( \sqrt{1 + \frac{L_A^2}{\gamma\delta}} - 1 \right) \quad \text{and} \quad \tau = \frac{\delta}{2} \left( \sqrt{1 + \frac{L_A^2}{\gamma\delta}} + 1 \right) \quad (141)$$

and the resulting rate:

$$\tilde{\omega}^* = \frac{\sqrt{1 + L_A^2/(\gamma\delta)} - 1}{\sqrt{1 + L_A^2/(\gamma\delta)} + 1} = \frac{\sqrt{1 + \kappa_f} - 1}{\sqrt{1 + \kappa_f} + 1}. \quad (142)$$

### 4.3 Theoretical rate comparison

Figure 1 compares the theoretical rates of the unaccelerated ADMM, the accelerated ADMM, the oPDHG method and FISTA with constant step, by plotting for each algorithm the best rate with respect to the condition number  $\kappa_f$ . The rate achieved by FISTA is the best one, but remains comparable with the accelerated ADMM and the oPDHG method. As expected, the unaccelerated ADMM yield larger rate values.

## 5 Relations of other methods

In this section, we make a quick review on other linear convergence results for variant of the ADMM found in the literature. Generally, their differ from our result on the hypotheses made on the problem (both on the regularity of the objective function and on the operators).

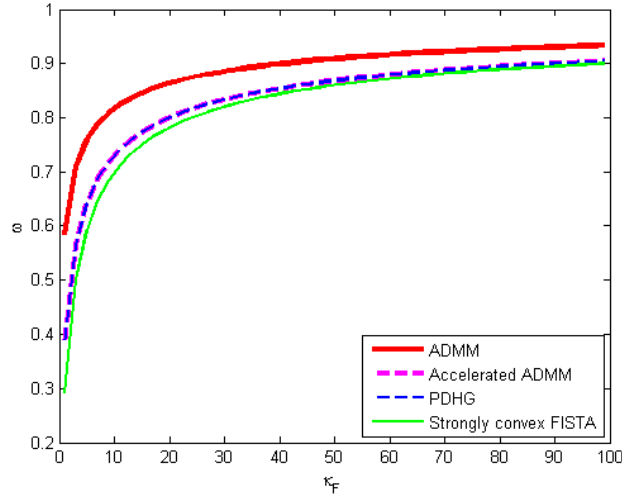


Figure 1: theoretical rate comparison. In red/thick the unaccelerated ADMM, in magenta/thick dotted the accelerated ADMM, in blue/dotted the oPDHG and in green strongly convex FISTA with constant step.

### 5.1 Overrelaxed ADMM

In [15], the authors propose to add an overrelaxation step in the spirit of NESTEROV’s acceleration. They showed linear convergence rate when  $h$  is assumed to be strongly convex and with LIPSCHITZ-continuous gradient, while  $B$  is invertible and  $A$  is full column rank.

### 5.2 Generalized ADMM

In [7], the authors studied the ADMM in a wider framework, by allowing in each partial minimization to add an extra proximal term, which leads to a generalized ADMM. Linear convergence rates are proved for four scenarios in which at least one of the functions  $g$  or  $h$  is strongly convex *and* has a LIPSCHITZ gradient, which is not assumed here. The case we treated is considered, but with extra assumptions (in particular,  $h$  is supposed to be strongly convex). They provided an explicit convergence rate for only one scenario [7, Corollary 3.6].

### 5.3 Relaxed ADMM

It can be shown that the ADMM iterations are also equivalent to applying the DOUGLAS-RACHFORD splitting (DRS) to the dual of (6). A relaxed version of the DRS, called PEACEMAN-RACHFORD splitting (PRS), can be obtained by introducing a relaxed parameter in the DRS iterations. Applying the PRS on the dual of (6) hence leads to a so-called relaxed ADMM [6]. In [6, Theorem 6.3], the authors proved the linear convergence rate of the relaxed ADMM in various cases (including the one we studied here), which depend on the assumptions made on the operators  $A$  and  $B$  (which is not supposed to be the negative identity) and / or on the regularity of the functions  $g$  and  $h$ . However, the study is theoretical and does not provide explicit optimal rates.



## 5.4 $K$ -block ADMM

In [12], the authors proved a linear convergence rate in the case where one can make assumptions on  $g$  and  $h$  which are supposed to be decomposable into a strictly convex term and a polyhedral one. This includes for instance the strongly convex case, but do not recover the smooth case studied in this paper. Furthermore, hypothesis on the rank of operators  $A$  and  $B$  (not necessary the negative identity) are made. Moreover, their proof still holds when the objective function is a sum of  $K$  separable convex functions (with an according number of variables).

## 6 Applications

### 6.1 A toy example

#### 6.1.1 Problem

Let  $N$  be a integer. We consider the following constrained problem:

$$\min_{\substack{x=(x_i)_{i=0,\dots,N-1} \in \mathbb{R}^N \\ x_0=1}} \left\{ f(x) := \frac{M-m}{2} \|K_N x\|_2^2 + \frac{m}{2} \|x\|_2^2 \right\} \quad (143)$$

where the linear operator  $K_N : \mathbb{R}^N \rightarrow \mathbb{R}^{N-1}$  is defined by  $(K_N x)_i = (x_{i+1} - x_i)/2$  for any  $i = 0, \dots, N-2$ , of norm  $\|K_N\| \leq 1$ . The condition number of this problem is  $M/m$ . Hence, if  $m$  is negligible compared to  $M$ , then the problem is ill-conditioned. Let  $h(z) := (M-m)\|z\|_2^2/2$  for any  $z \in \mathbb{R}^{N-1}$  and  $g(x) := m\|x\|_2^2/2 + \chi_{\{1\}}(x_0)$  for any  $x = (x_i)_{i=0,\dots,N-1} \in \mathbb{R}^N$ . The function  $g$  is  $m$ -convex and the convex conjugate  $h^* : y \mapsto (M-m)^{-1}\|y\|_2^2/2$  is  $(M-m)^{-1}$ -convex.

#### 6.1.2 Solution

The minimizer of problem (143) may be explicitly computed, by introducing the sub-vector  $\hat{x}$  given by:

$$\forall i = 0, \dots, N-2, \quad \hat{x}_i = x_{i+1}. \quad (144)$$

such that  $x = (1, \hat{x})$ . The constrained problem (143) can thus be rewritten in the unconstrained form

$$\min_{\hat{x}=(\hat{x}_i)_{i=0,\dots,N-2} \in \mathbb{R}^{N-1}} \left\{ \frac{M-m}{2} \left( \|K_{N-1} \hat{x}\|_2^2 + \frac{(\hat{x}_0 - 1)^2}{4} \right) + \frac{m}{2} (\|\hat{x}\|_2^2 + 1) \right\} \quad (145)$$

The minimizer  $\hat{x}^*$  is then given by the EULER equation, namely  $\hat{x}^* = A^{-1}b$  with

$$A = m I_{N-1} + (M-m) K_{N-1}^* K_{N-1} + \frac{M-m}{4} e_{0,0} \quad (146)$$

where  $e_{0,0}$  denote the matrix of size  $N-1$  with null coefficients except the one at index  $(0,0)$  equal to 1. The vector  $b$  is given by  $b := (M-m)e_0/4$ , with  $e_0$  the first vector of the canonical basis of  $\mathbb{R}^{N-1}$ . Hence, the minimizer of the initial problem (143) is  $x^* = (1, \hat{x}^*)$ . For  $N = 15$ ,  $M = 1000$ , and  $m = 1$ , Figure 2 plots  $x^*$ .

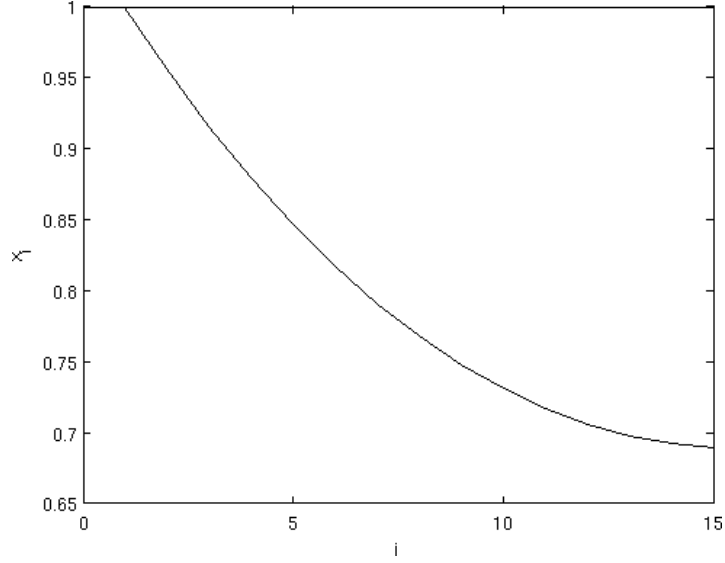


Figure 2: Minimizer of (143).

### 6.1.3 ADMM

We apply the accelerated ADMM, which yields

$$\begin{cases} x_{n+1} = \arg \min_{\substack{x=(x_i)_{i \in \mathbb{R}^N} \\ x_0=1}} \left\{ \frac{m}{2} \|x\|_2^2 + \langle K_N x, y_n \rangle + \frac{1}{2\tau} \|K_N x - z_n\|_2^2 \right\} \\ z_{n+1} = \arg \min_{z \in \mathbb{R}^{N-1}} \left\{ \frac{M-m}{2} \|z\|_2^2 - \langle z, y_n \rangle + \frac{1}{2\tau'} \|K_N x_{n+1} - z\|_2^2 \right\} \\ y_{n+1} = y_n + \frac{1}{\tau'} (K_N x_{n+1} - z_{n+1}). \end{cases}$$

The  $z$ -update is computed thanks to the EULER equation:

$$z_{n+1} = \frac{y_n + K_N x_{n+1} / \tau'}{M - m + 1 / \tau'}. \quad (147)$$

The  $x$ -update is computed thanks to the subvectors we introduced above and is equivalent to minimizing

$$\frac{m}{2} \|\hat{x}\|_2^2 + \langle \hat{x}, K_{N-1}^* \hat{y}_n \rangle + \frac{1}{2\tau} \left( \|K_{N-1} \hat{x} - \hat{z}_n\|_2^2 + \left( \frac{\hat{x}_0 - 1}{2} - (z_n)_0 \right)^2 \right). \quad (148)$$

The EULER equation ensures that  $\hat{x}_{n+1} = A_n^{-1} b_n$  with

$$A_n = m I_{N-1} + \frac{1}{\tau} K_{N-1}^* K_{N-1} + \frac{1}{4\tau} e_{0,0} \quad (149)$$

$$b_n = -K_{N-1}^* \hat{y}_n + \frac{1}{\tau} K_{N-1}^* \hat{z}_n + \left( -\frac{(y_n)_0}{2} + \frac{1}{2\tau} (z_n)_0 + \frac{1}{4\tau} \right) e_0. \quad (150)$$

We eventually have  $x_{n+1} = (1, \hat{x}_{n+1})$ .

#### 6.1.4 Parameters

We tested two sets of parameters:

1. optimal parameter for the unaccelerated ADMM:

$$\tau = \tau' = \sqrt{\frac{2}{m(M-m)}} \quad (151)$$

(we assume that  $L = 1$ ).

2. optimal parameters for the accelerated ADMM:

$$\tau = \frac{1}{2(M-m)} \left( \sqrt{1 + \frac{4(M-m)}{m}} + 1 \right) \quad (152)$$

$$\tau' = \tau - \frac{1}{M-m} = \frac{1}{2(M-m)} \left( \sqrt{1 + \frac{4(M-m)}{m}} - 1 \right). \quad (153)$$

The convergence rates achieved in each case are respectively  $1/(\sqrt{1/(M/m-1)}/2 + 1)$  and  $(\sqrt{4M/m-3}-1)/(\sqrt{4M/m-3}+1)$ .

#### 6.1.5 Comparison with oPDHG and FISTA

To solve problem (145), we can use the oPDHG iterations, by considering its primal-dual formulation

$$\min_{\substack{x=(x_i)_{i=0,\dots,N-1} \in \mathbb{R}^N \\ x_0=1}} \sup_{z' \in \mathbb{R}^{N-1}} \left\{ \frac{m}{2} \|x\|_2^2 + \langle Kx, z' \rangle - \frac{1}{2(M-m)} \|z'\|_2^2 \right\}. \quad (154)$$

Hence, we are considering the following algorithm:

$$\begin{cases} z'_{n+1} = \text{prox}_{\sigma h^*}(z'_n + \sigma K_N \bar{x}_n) \\ x_{n+1} = \text{prox}_{\tau g}(x_n - \tau K_N^* z'_{n+1}) \\ \bar{x}_{n+1} = x_{n+1} + \theta (x_{n+1} - x_n) \end{cases} \quad (155)$$

for which the best theoretical convergence rate is achieved when choosing

$$\tau = \frac{1}{2(M-m)} \left( 1 + \sqrt{1 + \frac{4(M-m)}{m}} \right) \quad (156)$$

$$\sigma = \frac{m}{2} \left( 1 + \sqrt{1 + \frac{4(M-m)}{m}} \right) \quad (157)$$

$$\theta = \frac{\sqrt{1 + \frac{4(M-m)}{m}} - 1}{\sqrt{1 + \frac{4(M-m)}{m}} + 1} < 1. \quad (158)$$

The  $z'$ -iterates are explicitly given by

$$z'_{n+1} = \frac{M-m}{M-m+\sigma} (z'_n + \sigma K_N \bar{x}_n) \quad (159)$$

$$\hat{x}_{n+1} = \frac{\hat{x}_n/\tau - K_{N-1}^* \hat{z}'_{n+1}}{1/\tau + m} \quad \text{and} \quad x_{n+1} = (1, \hat{x}_{n+1}). \quad (160)$$

Note that, unlike in the ADMM iterations, there is no operator to invert.

We can also use the FISTA algorithm, which solves problem (143) by an accelerated FBS which can be written

$$\begin{cases} x_{n+1} = \text{prox}_{\tau g}(\bar{x}_n - \tau \nabla h(\bar{x}_n)) \\ \bar{x}_{n+1} = x_{n+1} + \theta_{n+1} (x_{n+1} - x_n). \end{cases} \quad (161)$$

where  $\theta_n$  is given by (21), with  $\tau = 1/(M - m)$ . The  $x$ -iterates are explicitly given by

$$\hat{x}_{n+1} = \frac{\widehat{\bar{x}}_n/\tau - (M - m)(K_N^* \widehat{K}_N \widehat{\bar{x}}_n)}{1/\tau + m}. \quad (162)$$

### 6.1.6 Results

To compare the convergence of each set of parameters, we used two tools:

1. the solution error  $\|x_n - x^*\|_2^2$ ;
2. the objective error  $f(x_n) - f(x^*)$ .

Figure 3 displays the evolution of both measures, as well as the theoretical convergence decays expected in each case ( $\tilde{\omega}$  and  $\omega$ ). We chose  $m = 0.1$  and  $M = 10$ , so that  $\kappa_f = 100$ .

We first observe that, as expected, the accelerated ADMM has a better convergence than the unaccelerated ADMM. The empirical rates are better than the theoretical ones, which can be explained by the over-smoothness of the quadratic problem, compared to the assumptions required by the smooth case.

We also observe oscillations for both the oPDHG and FISTA. Rippling for FISTA has been already observed for quadratic problems of this kind [16]. This phenomena occurs when the overrelaxation parameter  $\theta$  is chosen too large compared to the eigenvalues of  $m I_N + (M - m) K_N^* K_N$ . Similar cause may explain the oscillations in the oPDHG, namely using overrelaxation steps can introduce oscillations when the according parameter are unproperly chosen. Hence, we do not expect to observe such oscillations for ADMM-like schemes.

## 6.2 Denoising with TV-Huber

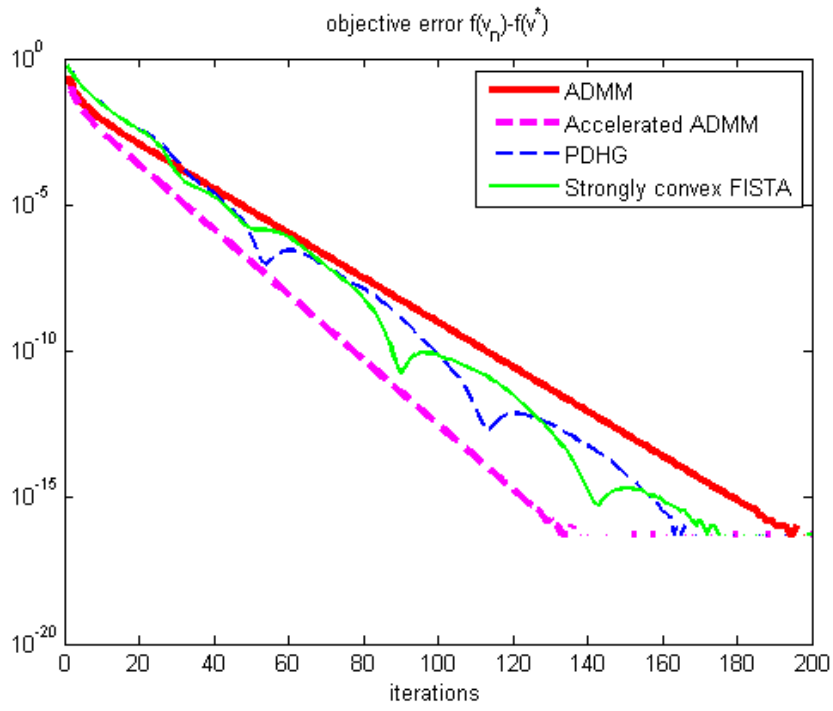
### 6.2.1 Problem

We now apply the accelerated ADMM to a denoising problem, which is less smooth and more realistic than the toy example. Let  $g \in \mathbb{R}^{3N_x N_y}$  be a RGB-color (noisy) image. We want to solve the following problem:

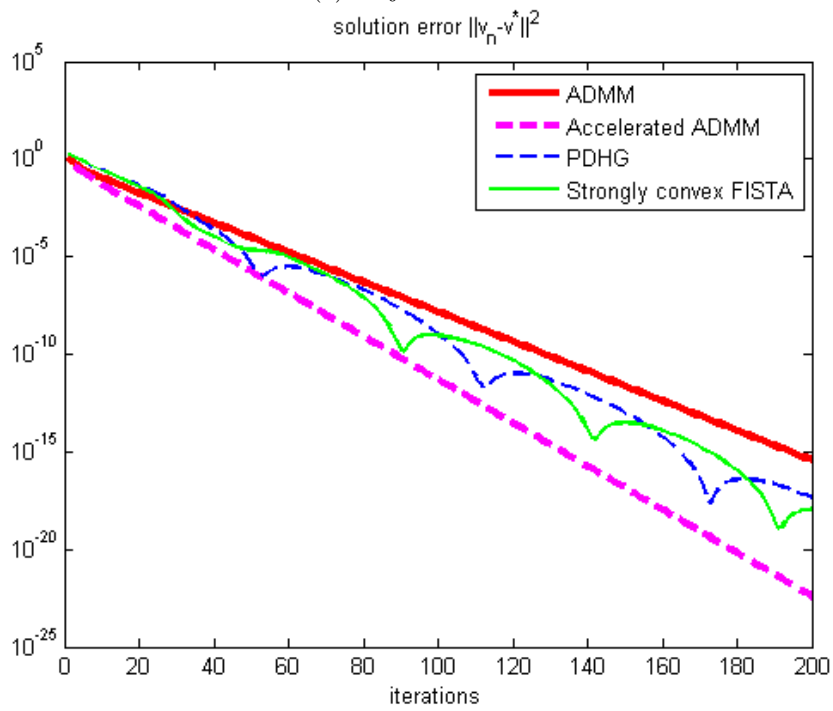
$$\min_{v \in \mathbb{R}^{3N_x N_y}} \left\{ f(v) := \frac{\mu}{2} \|v - u\|_2^2 + h(\nabla v) \right\} \quad (163)$$

where the gradient linear operator  $\nabla : \mathbb{R}^{3N_x N_y} \rightarrow \mathbb{R}^{3N_x N_y} \times \mathbb{R}^{3N_x N_y}$  is defined for any color image  $v$  by a pair of color images  $\nabla v = (\delta_x v, \delta_y v)^T$ . The finite differences are given at any index  $(i, j) \in [0, N_x - 1] \times [0, N_y - 1]$  by

$$(\delta_x v)_{i,j} = \begin{cases} v_{i+1,j} - v_{i,j} & \text{if } i < N_x - 1 \\ 0 & \text{otherwise} \end{cases} \quad (164)$$



(a) Objective error



(b) Solution error

Figure 3: Empirical convergence for the toy example.

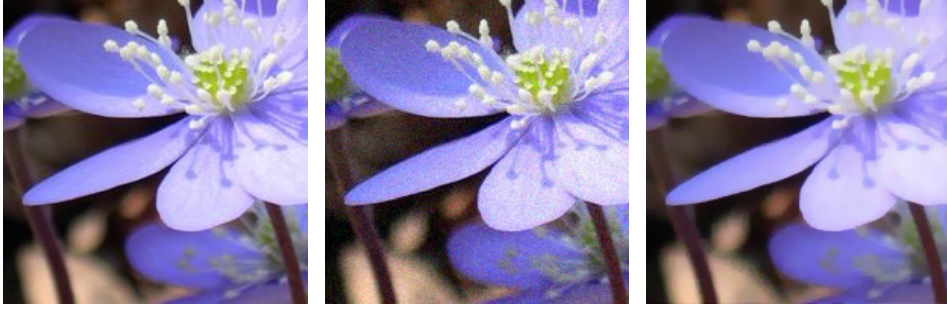


Figure 4: We added a white Gaussian noise to an ideal image (left) to get noisy image (middle). The noise is of standard variation 10 (the image values are between 0 and 255). The denoising is made by solving (163) (right). Source : *Hepatica nobilis* flowers, by Archenzo (detail).

and

$$(\delta_y v)_{i,j} = \begin{cases} v_{i,j+1} - v_{i,j} & \text{if } j < N_y - 1 \\ 0 & \text{otherwise.} \end{cases} \quad (165)$$

The TV-HUBER regularization term is defined by

$$h(\nabla v) = \sum_{i=0}^{N_x-1} \sum_{j=0}^{N_y-1} h_0(\|(\nabla v)_{i,j}\|) \quad (166)$$

with

$$h_0(z) = \begin{cases} |z|^2/2 & \text{if } |z| \leq 1 \\ |z| - 1/2 & \text{if } |z| > 1 \end{cases} \quad \text{and} \quad h'_0(z) = \begin{cases} z & \text{if } |z| \leq 1 \\ z/|z| & \text{if } |z| > 1. \end{cases} \quad (167)$$

Hence, this term acts like a quadratic regularization when the image variations are small and like a TV regularization when they are larger (see Figure 4). The quantity  $\mu > 0$  is a weight parameter.

The convex conjugate  $h^*$  of the regularization function  $h$  can be proved to be

$$h^*(y) = \sum_{i=0}^{N_x-1} \sum_{j=0}^{N_y-1} \left( \frac{1}{2} |y_{i,j}|^2 + \chi_{[0,1]}(|y_{i,j}|) \right) \quad (168)$$

where  $\chi_{[0,1]}(t) = 0$  if  $t \in [0, 1]$  and  $+\infty$  otherwise. This implies that  $h^*$  is 1-convex.

### 6.2.2 ADMM

Let  $g := \mu \|\cdot - u\|_2^2/2$ . We apply the accelerated ADMM to problem (163), which leads to the following iterations:

$$\begin{cases} v_{n+1} = \arg \min_{v \in \mathbb{R}^{3N_x N_y}} \left\{ \frac{\mu}{2} \|v - u\|_2^2 + \langle \nabla v, \xi_n \rangle + \frac{1}{2\tau} \|\nabla v - \phi_n\|_2^2 \right\} \\ \phi_{n+1} = \arg \min_{\phi \in (\mathbb{R}^{3N_x N_y})^2} \left\{ h(\phi) - \langle \phi, \xi_n \rangle + \frac{1}{2\tau'} \|\nabla v_{n+1} - \phi\|_2^2 \right\} \\ \xi_{n+1} = \xi_n + \frac{1}{\tau'} (\nabla v_{n+1} - \phi_{n+1}). \end{cases}$$

Each minimization is solved thanks to the EULER equation: the  $v$ -update reads

$$v_{n+1} = \left( \mu \mathbf{I} + \frac{1}{\tau} \nabla^* \nabla \right)^{-1} \left( \mu u + \frac{1}{\tau} \nabla^* \phi_n - \nabla^* \xi_n \right) \quad (169)$$

whereas the  $\phi$ -update is given by

$$(\phi_{n+1})_{i,j} = \frac{\tau'(\xi_n)_{i,j} + (\nabla v_{n+1})_{i,j}}{|\tau'(\xi_n)_{i,j} + (\nabla v_{n+1})_{i,j}|} |(\phi_{n+1})_{i,j}| \quad (170)$$

with

$$|(\phi_{n+1})_{i,j}| = \begin{cases} \frac{\tau' |(\xi_n)_{i,j} + (\nabla v_{n+1})_{i,j}|}{\tau' + 1} & \text{if } |\tau'(\xi_n)_{i,j} + (\nabla v_{n+1})_{i,j}| \leq \tau' + 1 \\ |\tau'(\xi_n)_{i,j} + (\nabla v_{n+1})_{i,j}| - \tau' & \text{if } |\tau'(\xi_n)_{i,j} + (\nabla v_{n+1})_{i,j}| > \tau' + 1. \end{cases} \quad (171)$$

### 6.2.3 Parameters

Before choosing the parameters, we recall the regularity of the problem. Functions  $h^*$  and  $g$  are respectively 1-convex and  $\mu$ -convex. The gradient operator is bounded, of norm  $L \leq 2\sqrt{2}$  (this bound being tight when  $N_x$  or  $N_y$  go to  $+\infty$ ). Thus, we set  $L = 2\sqrt{2}$ . We tested two sets of parameters:

1. optimal parameter for the unaccelerated ADMM:  $\tau = \tau' = 4/\sqrt{\mu}$ ;
2. optimal paramters for the accelerated ADMM:

$$(\tau, \tau') = \left( \frac{1}{2} \left( \sqrt{1 + \frac{32}{\mu}} + 1 \right), \frac{1}{2} \left( \sqrt{1 + \frac{32}{\mu}} - 1 \right) \right). \quad (172)$$

These choices lead to the convergence rates  $1/(\sqrt{\mu/8} + 1)$  for the unaccelerated ADMM and  $(\sqrt{1 + 32/\mu} - 1)/(\sqrt{1 + 32/\mu} + 1)$  for the accelerated one.

### 6.2.4 oPDHG and FISTA

The primal-dual formulation of problem (163) is given by

$$\min_{v \in \mathbb{R}^{3N_x N_y}} \sup_{\phi \in (\mathbb{R}^{3N_x N_y})^2} \left\{ \frac{\mu}{2} \|v - u\|_2^2 + \langle \nabla v, \phi \rangle - h^*(\phi) \right\}. \quad (173)$$

Hence using the oPDHG algorithm to solve it leads to the following iterations:

$$\begin{cases} \phi'_{n+1} = \arg \min_{\phi' \in (\mathbb{R}^{3N_x N_y})^2} \left\{ h^*(\phi') + \frac{1}{2\sigma} \|\phi' - \phi'_n - \sigma \nabla \bar{v}_n\|^2 \right\} \\ v_{n+1} = \arg \min_{v \in \mathbb{R}^{3N_x N_y}} \left\{ \frac{\mu}{2} \|v - u\|^2 + \frac{1}{2\tau} \|v - v_n + \tau \nabla^* \phi'_{n+1}\|^2 \right\} \\ \bar{v}_{n+1} = v_{n+1} + \theta (v_{n+1} - v_n) \end{cases} \quad (174)$$

which are computed thanks to the EULER equation:

$$(\phi'_{n+1})_{i,j} = \text{proj}_{[-1,1]} \left( \frac{(\phi'_n)_{i,j} + \sigma (\nabla \bar{v}_n)_{i,j}}{1 + \sigma} \right) \quad \text{and} \quad v_{n+1} = \frac{v_n/\tau + \mu u - \nabla^* \phi'_{n+1}}{1/\tau + \mu}. \quad (175)$$

The best choice of parameters for this algorithm is (Theorem 2):

$$\tau = \frac{1 + \sqrt{1 + 32/\mu}}{16}, \quad \sigma = \frac{1 + \sqrt{1 + 32/\mu}}{16/\mu} \quad \text{and} \quad \theta = \frac{\sqrt{1 + 32/\mu} - 1}{\sqrt{1 + 32/\mu} + 1}. \quad (176)$$

If we apply FISTA to this problem, this leads to the following updates:

$$\begin{cases} v_{n+1} = \arg \min_{v \in \mathbb{R}^{3N_x N_y}} \left\{ \frac{\mu}{2} \|v - u\|^2 + \frac{1}{2\tau} \|v - \bar{v}_n + \tau \nabla^* \nabla (\nabla h(\bar{v}_n))\|^2 \right\} \\ \bar{v}_{n+1} = v_{n+1} + \theta_n (v_{n+1} - v_n) \end{cases} \quad (177)$$

which leads to the explicit update

$$v_{n+1} = \frac{\bar{v}_n/\tau + \mu u - \nabla^* \nabla (\nabla h(\bar{v}_n))}{1/\tau + \mu}. \quad (178)$$

The variable relaxation parameter follows the update rule (21) with  $\tau = 1/8$ .

### 6.2.5 Results

To measure the convergence of the algorithm, we used the same two tools as in the previous case: the solution error and the objective error.

Figure 5(a) displays the evolution of the objective error, while Figure 5(b) shows the decay of the solution error, for the accelerated ADMM and the oPDHG method. In the latter, the theoretical linear rate  $\tilde{\omega}$  is also plotted for comparison. We chose  $\mu = 10$ . The solution error decreases as expected for all methods except FISTA, for which we did not estimate a finer theoretical rate for the solution error. In practice, it seems that it converges with same rate as the oPDHG. Hence, in terms of solution error convergence, the accelerated ADMM provides the best empirical decay. For the objective error, the accelerated ADMM, the oPDHG method and FISTA yield comparable decay rate. However, one should keep in mind that both the unaccelerated ADMM and the accelerated ADMM require an operator inversion, unlike the oPDHG method and FISTA. Hence, even if comparable number of iterations are needed to achieve convergence, the ADMMs iterations are more time consuming than the other methods and should be used only when the inversion of the operator can be implemented efficiently.

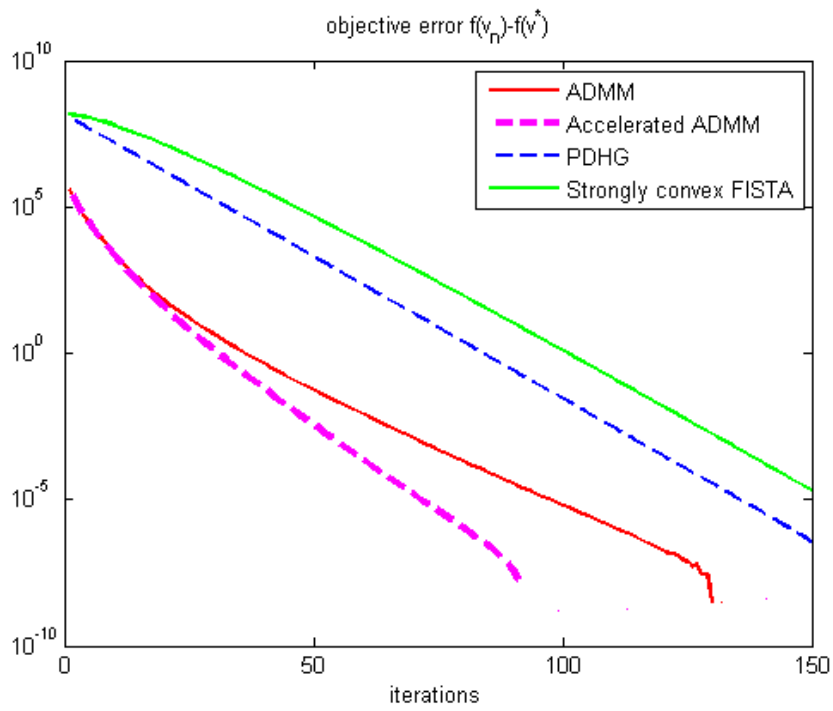
## 7 Conclusion

In this work, we studied the convergence of the oPDHG scheme in the case where the composite problem has a strongly convex part and a differentiable with a Lipschitz continuous gradient part. Using the equivalence between this algorithm and the ADMM, we provided a new convergence analysis of the latter. This analysis allowed us to introduce an accelerated variant of the ADMM, which is proved to have same convergence rate as the oPDHG method. Experimental results confirmed this theoretical analysis. In particular, it has been observed that the accelerated ADMM does not introduce oscillations in some cases, unlike the oPDHG algorithm and FISTA, which are known to be in practice more efficient than the ADMM-like scheme, since they require no operator inversion.

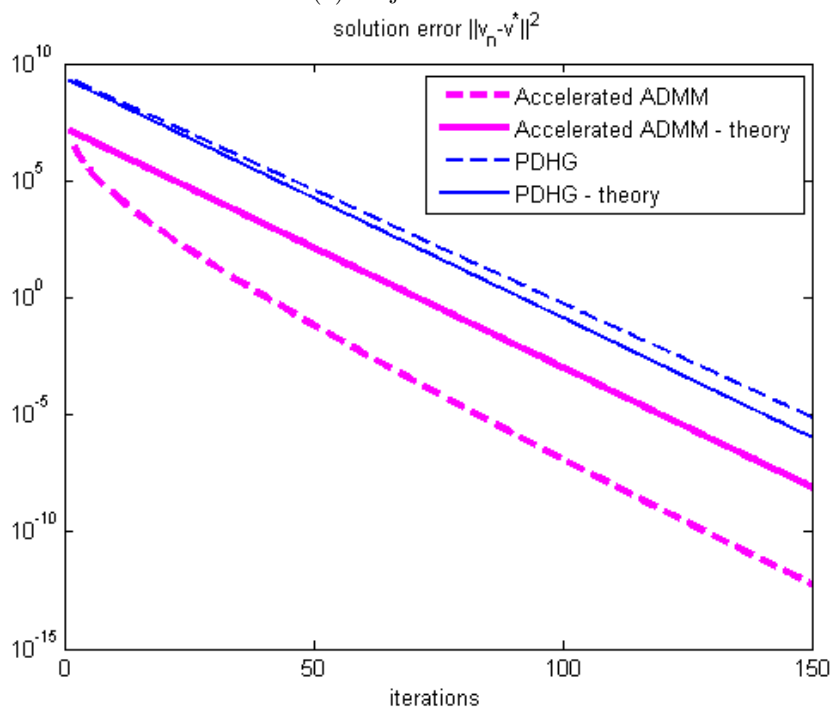
## References

- [1] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.





(a) Objective error



(b) Solution error - Comparison with the theoretical rate

Figure 5: Empirical convergence for TV-Huber denoising.

- [2] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.
- [3] Antonin Chambolle and Thomas Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of Mathematical Imaging and Vision*, 40(1):120–145, 2011.
- [4] Antonin Chambolle and Thomas Pock. On the ergodic convergence rates of a first-order primal-dual algorithm. *Mathematical Programming*, pages 1–35, 2015.
- [5] Antonin Chambolle and Thomas Pock. An introduction to continuous optimization for imaging. *Acta Numerica*, 25:161–319, 2016.
- [6] Damek Davis and Wotao Yin. Faster convergence rates of relaxed peaceman-rachford and admm under regularity assumptions. *arXiv preprint arXiv:1407.5210*, 2014.
- [7] Wei Deng and Wotao Yin. On the global and linear convergence of the generalized alternating direction method of multipliers. *Journal of Scientific Computing*, 66(3):889–916, 2016.
- [8] Ernie Esser, Xiaoqun Zhang, and Tony F Chan. A general framework for a class of first order primal-dual algorithms for convex optimization in imaging science. *SIAM Journal on Imaging Sciences*, 3(4):1015–1046, 2010.
- [9] Daniel Gabay and Bertrand Mercier. A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers & Mathematics with Applications*, 2(1):17–40, 1976.
- [10] Roland Glowinski and A Marrocco. Sur l’approximation, par éléments finis d’ordre un, et la résolution, par pénalisation-dualité d’une classe de problèmes de dirichlet non linéaires. *Revue française d’automatique, informatique, recherche opérationnelle. Analyse numérique*, 9(2):41–76, 1975.
- [11] Bingsheng He, Yanfei You, and Xiaoming Yuan. On the convergence of primal-dual hybrid gradient algorithm. *SIAM Journal on Imaging Sciences*, 7(4):2526–2537, 2014.
- [12] Mingyi Hong and Zhi-Quan Luo. On the linear convergence of the alternating direction method of multipliers. *arXiv preprint arXiv:1208.3922*, 2012.
- [13] Yurii Nesterov. A method of solving a convex programming problem with convergence rate  $o(1/k^2)$ . In *Soviet Mathematics Doklady*, volume 27, pages 372–376, 1983.
- [14] Yurii Nesterov. Introductory lectures on convex optimization. applied optimization, vol. 87, 2004.
- [15] Robert Nishihara, Laurent Lessard, Benjamin Recht, Andrew Packard, and Michael I Jordan. A general analysis of the convergence of admm. *arXiv preprint*, 2015.
- [16] Brendan O’Donoghue and Emmanuel Candes. Adaptive restart for accelerated gradient schemes. *Foundations of computational mathematics*, 15(3):715–732, 2015.

- [17] Thomas Pock, Daniel Cremers, Horst Bischof, and Antonin Chambolle. An algorithm for minimizing the mumford-shah functional. In *Computer Vision, 2009 IEEE 12th International Conference on*, pages 1133–1140. IEEE, 2009.
- [18] Mingqiang Zhu and Tony Chan. An efficient primal-dual hybrid gradient algorithm for total variation image restoration. *UCLA CAM Report*, pages 08–34, 2008.