

Accelerated Alternating Descent Methods for Dykstra-Like Problems

Antonin Chambolle¹  · Pauline Tan¹ · Samuel Vaiter²

Received: 12 July 2016 / Accepted: 4 March 2017
© Springer Science+Business Media New York 2017

Abstract This paper extends recent results by the first author and T. Pock (ICG, TU Graz, Austria) on the acceleration of alternating minimization techniques for quadratic plus nonsmooth objectives depending on two variables. We discuss here the strongly convex situation, and how ‘fast’ methods can be derived by adapting the overrelaxation strategy of Nesterov for projected gradient descent. We also investigate slightly more general alternating descent methods, where several descent steps in each variable are alternatively performed.

Keywords Alternating minimizations · Block descent algorithms · Accelerated methods · Total variation minimization

1 Introduction

This paper addresses the acceleration of alternating minimizations or descent methods for elementary problems which involve two variables coupled by a quadratic penalization. Such problems arise for instance in the computation of the proximity operators of sums of simple functions, for which in some setting (as we illustrate in an experimental section) it might be beneficial to perform such a splitting which decomposes the problem into tiny parallel subproblems, rather

than tackle the global problem by an accelerated descent or primal-dual algorithm such as [3, 9, 10, 19].

The present paper is a follow-up of [11] where this issue was already investigated, and a few contexts where acceleration was possible were investigated. In this paper, we extend these results in two directions: first, we consider strongly convex objectives and show how one can obtain nearly optimal linear convergence rates (in the sense of the lower bounds of [18, 19]) in the framework of alternating minimizations or descent. The case of minimizations is particular, as it can also be reduced into a forward–backward splitting method applied to auxiliary functions, and we could just refer to [12, 19] where a rate analysis is performed; however, it is roughly equivalent to perform an analysis adapted to the alternating minimizations algorithm.

As far as alternating descent methods are concerned, on the other hand, the problem does not boil down to a more standard structure, and the analysis is quite tedious (except if only one step of descent is performed at each step, as was studied in [11]). We perform this analysis in details; it leads however to algorithms which in theory would require to keep in memory a lot of intermediate states. We check experimentally that one can overlook these issues in implementations and still obtain good convergence properties. This is something we are unable to explain at this time.

Let us introduce now more precisely our problem and how it can be numerically tackled.

1.1 The Problem

We aim at solving convex minimization problems of the form

$$\min_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} \mathcal{E}(x, y) := f(x) + g(y) + \frac{1}{2} \|Ax + By\|^2 \quad (1)$$

✉ Antonin Chambolle
antonin.chambolle@cmmap.polytechnique.fr

Pauline Tan
pauline.tan@cmmap.polytechnique.fr

Samuel Vaiter
samuel.vaiter@u-bourgogne.fr

¹ CMAP, CNRS, Ecole Polytechnique, 91128 Palaiseau, France

² IMB, CNRS, Université de Bourgogne, 9 Ave Alain Savary, 21000 Dijon, France

where $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$ and $g : \mathcal{Y} \rightarrow \mathbb{R} \cup \{+\infty\}$ are two convex, proper, lower-semicontinuous (lsc) functions, and $A : \mathcal{X} \rightarrow \mathcal{Z}, B : \mathcal{Y} \rightarrow \mathcal{Z}$ two bounded linear operators. In the whole paper, $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$ should be thought as finite-dimensional Euclidean spaces, although the proofs carry on easily to the Hilbertian setting. These problems naturally arise in the computation of the ‘proximity operator’ of functions of the form $z \mapsto f(Kz)$:

$$\min_z f(Kz) + \frac{1}{2\tau} \|z - z^0\|^2$$

(given a point z^0 and $\tau > 0$), when f can be in turn split into two functions $f_1(K_1z) + f_2(K_2z)$. Examples will be provided in Sect. 7. The idea of Dykstra’s algorithm [5] (see also [13, Ex 10.11]) is to perform alternating minimizations on a dual problem: one writes (assuming formally that the min/max can be exchanged, which is generally true under quite mild assumptions)

$$\begin{aligned} & \min_z f_1(K_1z) + f_2(K_2z) + \frac{1}{2\tau} \|z - z^0\|^2 \\ &= \min_z \sup_{x_1, x_2} \langle x_1, K_1z \rangle + \langle x_2, K_2z \rangle \\ &\quad - f_1^*(x_1) - f_2^*(x_2) + \frac{1}{2\tau} \|z - z^0\|^2 \\ &= \max_{x_1, x_2} \inf_z \langle K_1^*x_1 + K_2^*x_2, z \rangle \\ &\quad + \frac{1}{2\tau} \|z - z^0\|^2 - f_1^*(x_1) - f_2^*(x_2) \\ &= \max_{x_1, x_2} \langle K_1^*x_1 + K_2^*x_2, z^0 \rangle \\ &\quad - f_1^*(x_1) - f_2^*(x_2) - \frac{\tau}{2} \|K_1^*x_1 + K_2^*x_2\|^2, \end{aligned}$$

which leads to a problem in form (1). Then, one can simply alternatively minimize the problem with respect to x_1 and then x_2 , provided the computations are tractable. Although it can be found in some cases, and for some geometries, that this method is efficient [16], in general its convergence rate can be quite poor [2,4]. In [11], it is observed that alternating minimization schemes on (1) can be accelerated using a FISTA-type overrelaxation [3] (see also [20] where a similar observation has been recently made). It is observed as well that this is still true that the exact minimizations are replaced with one step of a proximal-type descent for each variable. Of course, one should notice that such a problem is of the form smooth+nonsmooth minimization (in the variable (x, y)), so it is obvious that it can be tackled with an accelerated first-order method as in [3,19], and the techniques studied in [11] and in this note will certainly not improve the known rates of accelerated methods by an order of magnitude, but only the constants which appear in these rates. This was explained already in [11] and is confirmed experimentally: see in Table 1 the comparison between the ‘FISTA’ implementation and the alternating minimizations or descents. The

real speedup we can hope for lies mostly in the fact that, in some applications, we can split our problems into subproblems which are in turn split into many independent small dimensional problems, which one can hope to solve (almost) exactly, and in parallel.

In practice, it can be observed experimentally that performing *several* steps of descent in each variable (before turning to the other variable) improves the performances (which is to be expected, as it gets closer to performing an exact minimization). One of the goals of this paper is to try to support this experimental observation, by showing that acceleration is still possible in this setting and leads to comparable rates. Another goal is to extend the analysis in [11] to strongly convex objectives, showing again that one can obtain in these cases a descent rate similar to the rate of a standard accelerated method [19]. As an application, we show how to implement fast parallel solvers for the proximity operator of the total variation, with or without regularization, for gray and color images.

This paper is divided as follows: in the next section we introduce the general type of updates we are considering, consisting in one or several minimization steps of the objective with respect to one variable, the other being frozen. We derive various sufficient descent rules for this technique, depending on the properties of the functions.

Then in Sect. 3 we discuss the acceleration of alternating minimizations in the strongly convex case. Since alternating minimizations are a variant of forward–backward splitting methods, it is clear that one can expect good convergence rates by adapting standard methods [3,12,19]. This is what we establish in Theorem 1, extending a result of [11] in the nonstrongly convex case.

In the following two sections, we extend this result to alternating descent with several descent steps. We first discuss the nonstrongly convex case (Sect. 4.1) and then the strongly convex situation (Sect. 5). However, except when only one step of descent is performed in each variable (case $K = L = 1$, already discussed in [11]), the algorithms which are found are not very practical (as they require the introduction of too many auxiliary variable, while experiments seem to show that this is not really needed to obtain good convergence properties, cf Sect. 7.4.1).

Eventually, in Sect. 7, we discuss our application of these results, as mentioned before, to the computation of the proximity operator of a smoothed version of the total variation, for both scalar and vectorial (color) images. This approximation is shown to be a correct approximation of the isotropic total variation in an ‘Appendix’ (Theorem 4).

1.2 Main Assumptions

We will assume here that the functions f and g in (1) are convex and possibly strongly convex, in, respectively, metrics

defined by positive semidefinite matrices F and G : for all x, x' and y, y' , and for all $p \in \partial f(x'), q \in \partial g(y')$

$$f(x) \geq f(x') + \langle p, x - x' \rangle + \frac{1}{2} \|x - x'\|_F^2,$$

$$g(y) \geq g(y') + \langle q, y - y' \rangle + \frac{1}{2} \|y - y'\|_G^2,$$

where here F and G are (possibly vanishing) symmetric positive semidefinite operators.

2 Sufficient Descent Rules

2.1 General Updates

We will first consider the following general updates for x and y : given $(\bar{x}, \bar{y}, \bar{y}') \in \mathcal{X} \times \mathcal{Y}^2$, the metrics M, N and integer numbers $K, L \geq 1$, we obtain $(\hat{x}, \hat{y}, \hat{x}', \hat{y}') = T_{K,L}(\bar{x}, \bar{y}, \bar{y}')$ by letting $\hat{x}_0 = \bar{x}, \hat{y}_0 = \bar{y}$, and solving:

$$\hat{x}_{k+1} \in \arg \min_{x \in \mathcal{X}} f(x) + \frac{1}{2} \|Ax + B\bar{y}'\|^2 + \frac{1}{2} \|x - \hat{x}_k\|_M^2, \quad k = 0, \dots, K - 1, \tag{2}$$

$$\hat{x} = \hat{x}_K, \quad \hat{x}' = \frac{1}{K} \sum_{k=1}^K \hat{x}_k, \tag{3}$$

$$\hat{y}_{l+1} \in \arg \min_{y \in \mathcal{Y}} g(y) + \frac{1}{2} \|A\hat{x}' + By\|^2 + \frac{1}{2} \|y - \hat{y}_l\|_N^2, \quad l = 0, \dots, L - 1, \tag{4}$$

$$\hat{y} = \hat{y}_L, \quad \hat{y}' = \frac{1}{L} \sum_{l=1}^L \hat{y}_l. \tag{5}$$

A basic observation is that, using the strong convexity of the norms and possibly of f , for all $x \in \mathcal{X}$,

$$f(x) + \frac{1}{2} \|Ax + B\bar{y}'\|^2 + \frac{1}{2} \|x - \hat{x}_k\|_M^2 \geq f(\hat{x}_{k+1}) + \frac{1}{2} \|A\hat{x}_{k+1} + B\bar{y}'\|^2 + \frac{1}{2} \|\hat{x}_{k+1} - \hat{x}_k\|_M^2 + \frac{1}{2} \|x - \hat{x}_{k+1}\|_{A^*A+M+F}^2 \tag{6}$$

and

$$g(y) + \frac{1}{2} \|A\hat{x}' + By\|^2 + \frac{1}{2} \|y - \hat{y}_l\|_N^2 \geq g(\hat{y}_{l+1}) + \frac{1}{2} \|A\hat{x}' + B\hat{y}_{l+1}\|^2 + \frac{1}{2} \|\hat{y}_{l+1} - \hat{y}_l\|_N^2 + \frac{1}{2} \|y - \hat{y}_{l+1}\|_{B^*B+N+G}^2 \tag{7}$$

If we sum (6) from $k = 0$ to $K - 1$, we obtain (remember $\bar{x} = \hat{x}_0$)

$$Kf(x) + \frac{K}{2} \|Ax + B\bar{y}'\|^2 + \frac{1}{2} \|x - \bar{x}\|_M^2 \geq \sum_{k=1}^K \left(f(\hat{x}_k) + \frac{1}{2} \|A\hat{x}_k + B\bar{y}'\|^2 + \frac{1}{2} \|\hat{x}_k - \hat{x}_{k-1}\|_M^2 \right)$$

$$+ \frac{1}{2} \|x - \hat{x}_k\|_{A^*A+F}^2 \Big) + \frac{1}{2} \|x - \hat{x}\|_M^2.$$

Dividing by K and using the convexity of f and the norms, it follows:

$$f(x) + \frac{1}{2} \|Ax + B\bar{y}'\|^2 + \frac{1}{2K} \|x - \bar{x}\|_M^2 \geq f(\hat{x}') + \frac{1}{2} \|A\hat{x}' + B\bar{y}'\|^2 + \frac{1}{2} \|x - \hat{x}'\|_{A^*A+F}^2 + \frac{1}{2K} \|x - \hat{x}\|_M^2. \tag{8}$$

Remark It is maybe suboptimal to do so, as we do not exploit the fact that letting $x = \hat{x}^k$ in (6) yields

$$f(\hat{x}_k) + \frac{1}{2} \|A\hat{x}_k + B\bar{y}'\|^2 \geq f(\hat{x}_{k+1}) + \frac{1}{2} \|A\hat{x}_{k+1} + B\bar{y}'\|^2 + \frac{1}{2} \|\hat{x}_{k+1} - \hat{x}_k\|_{A^*A+2M+F}^2,$$

which would allow to evaluate the first two terms in the right-hand side of (8) at \hat{x} rather than \hat{x}' , yielding a smaller right-hand side. However, we need later on to exploit cancellations between the terms involving the norm $\|\cdot\|_{A^*A}$, which we cannot do anymore if we improve the inequality in this way.

Similarly, one finds

$$g(y) + \frac{1}{2} \|A\hat{x}' + By\|^2 + \frac{1}{2L} \|y - \bar{y}\|_N^2 \geq g(\hat{y}') + \frac{1}{2} \|A\hat{x}' + B\hat{y}'\|^2 + \frac{1}{2} \|y - \hat{y}'\|_{B^*B+G}^2 + \frac{1}{2L} \|y - \hat{y}\|_N^2. \tag{9}$$

We observe again that one also has, for $l = 0, \dots, L - 1$,

$$g(\hat{y}_l) + \frac{1}{2} \|A\hat{x}' + B\hat{y}_l\|^2 \geq g(\hat{y}_{l+1}) + \frac{1}{2} \|A\hat{x}' + B\hat{y}_{l+1}\|^2 + \frac{1}{2} \|\hat{y}_{l+1} - \hat{y}_l\|_{B^*B+2N+G}^2,$$

so that in particular, recalling $\hat{y} = \hat{y}^L$,

$$\frac{1}{L} \sum_{l=1}^L g(\hat{y}_l) + \frac{1}{2} \|A\hat{x}' + B\hat{y}_l\|^2 \geq g(\hat{y}) + \frac{1}{2} \|A\hat{x}' + B\hat{y}\|^2$$

(plus a term controlling the differences, which is hard to exploit), so that one also can write

$$g(y) + \frac{1}{2} \|A\hat{x}' + By\|^2 + \frac{1}{2L} \|y - \bar{y}\|_N^2 \geq g(\hat{y}) + \frac{1}{2} \|A\hat{x}' + B\hat{y}\|^2 + \frac{1}{2} \|y - \hat{y}'\|_{B^*B+G}^2 + \frac{1}{2L} \|y - \hat{y}\|_N^2. \tag{10}$$

Summing (8) and (9), we obtain:

$$\mathcal{E}(x, y) + \frac{1}{2K} \|x - \bar{x}\|_M^2 + \frac{1}{2L} \|y - \bar{y}\|_N^2 \geq \mathcal{E}(\hat{x}', \hat{y}') + \frac{1}{2K} \|x - \hat{x}\|_M^2 + \frac{1}{2L} \|y - \hat{y}\|_N^2 + \frac{1}{2} \|x - \hat{x}'\|_{A^*A+F}^2 + \frac{1}{2} \|y - \hat{y}'\|_{B^*B+G}^2$$

$$\begin{aligned}
 &+ \frac{1}{2} \|Ax + By\|^2 - \frac{1}{2} \|Ax + B\bar{y}'\|^2 \\
 &+ \frac{1}{2} \|A\hat{x}' + B\bar{y}'\|^2 - \frac{1}{2} \|A\hat{x}' + By\|^2.
 \end{aligned}$$

Now we observe that

$$\begin{aligned}
 &\frac{1}{2} \|Ax + By\|^2 - \frac{1}{2} \|Ax + B\bar{y}'\|^2 \\
 &+ \frac{1}{2} \|A\hat{x}' + B\bar{y}'\|^2 - \frac{1}{2} \|A\hat{x}' + By\|^2 \\
 &= \langle A(x - \hat{x}'), B(y - \bar{y}') \rangle \geq -\frac{1}{2} \|A(x - \hat{x}')\|^2 \\
 &- \frac{1}{2} \|B(y - \bar{y}')\|^2,
 \end{aligned}$$

and we deduce

$$\begin{aligned}
 \mathcal{E}(x, y) &+ \frac{1}{2K} \|x - \bar{x}\|_M^2 + \frac{1}{2L} \|y - \bar{y}\|_N^2 \\
 &+ \frac{1}{2} \|y - \bar{y}'\|_{B^*B}^2 \\
 &\geq \mathcal{E}(\hat{x}', \hat{y}') + \frac{1}{2K} \|x - \hat{x}\|_M^2 + \frac{1}{2L} \|y - \hat{y}\|_N^2 \\
 &+ \frac{1}{2} \|x - \hat{x}'\|_F^2 + \frac{1}{2} \|y - \hat{y}'\|_{B^*B+G}^2.
 \end{aligned}$$

Had we used (10) instead of (9), we would have rather obtained, in the same way:

$$\begin{aligned}
 \mathcal{E}(x, y) &+ \frac{1}{2K} \|x - \bar{x}\|_M^2 \\
 &+ \frac{1}{2L} \|y - \bar{y}\|_N^2 + \frac{1}{2} \|y - \bar{y}'\|_{B^*B}^2 \\
 &\geq \mathcal{E}(\hat{x}', \hat{y}') + \frac{1}{2K} \|x - \hat{x}\|_M^2 + \frac{1}{2L} \|y - \hat{y}\|_N^2 \\
 &+ \frac{1}{2} \|x - \hat{x}'\|_F^2 + \frac{1}{2} \|y - \hat{y}'\|_{B^*B+G}^2.
 \end{aligned}$$

In general, we will consider a new point (\tilde{x}, \tilde{y}) such that $\mathcal{E}(\tilde{x}, \tilde{y}) \leq \mathcal{E}(\hat{x}', \hat{y}')$ (so one could have $\tilde{x} = \hat{x}', \tilde{y} = \hat{y}'$), and use the general sufficient descent rule:

$$\begin{aligned}
 \mathcal{E}(x, y) &+ \frac{1}{2K} \|x - \bar{x}\|_M^2 + \frac{1}{2L} \|y - \bar{y}\|_N^2 \\
 &+ \frac{1}{2} \|y - \bar{y}'\|_{B^*B}^2 \\
 &\geq \mathcal{E}(\tilde{x}, \tilde{y}) + \frac{1}{2K} \|x - \hat{x}\|_M^2 + \frac{1}{2L} \|y - \hat{y}\|_N^2 \\
 &+ \frac{1}{2} \|x - \hat{x}'\|_F^2 + \frac{1}{2} \|y - \hat{y}'\|_{B^*B+G}^2. \tag{11}
 \end{aligned}$$

2.2 The Case of Alternating Minimizations

The case of alternating minimizations, discussed in [11], is substantially simpler. It corresponds to having $M = N = 0$, $K = L = 1$; in particular the points \bar{x}, \bar{y} are not used, and $(\hat{x}', \hat{y}') = (\hat{x}, \hat{y})$. To simplify, in this case, we drop the prime and denote \bar{y} the initial point \bar{y}' . Equation (11) becomes

$$\begin{aligned}
 \mathcal{E}(x, y) &+ \frac{1}{2} \|y - \bar{y}\|_{B^*B}^2 \geq \mathcal{E}(\tilde{x}, \tilde{y}) \\
 &+ \frac{1}{2} \|x - \hat{x}\|_F^2 + \frac{1}{2} \|y - \hat{y}\|_{B^*B+G}^2. \tag{12}
 \end{aligned}$$

In general, in that case, the most natural choice for (\tilde{x}, \tilde{y}) is of course the point (\hat{x}, \hat{y}) , as by construction it has the lowest energy encountered so far.

2.3 A Further Improvement

Now, we use the ‘FISTA’ trick which consists, given (x^k, y^k) a current iterate, in replacing (x, y) in (11) with points of the form $(x + (t - 1)x^k)/t, (y + (t - 1)y^k)/t$, for $t \geq 1$. We let also $(x^{k+1}, y^{k+1}) = (\tilde{x}, \tilde{y}), (\hat{x}^{k+1}, \hat{y}^{k+1}) = (\hat{x}', \hat{y}')$, $(\hat{x}^{k+1}, \hat{y}^{k+1}) = (\hat{x}, \hat{y})$, and, as well, $(\bar{x}^k, \bar{y}^k) = (\bar{x}, \bar{y}), \bar{y}^k = \bar{y}'$. It follows, after a multiplication by t^2

$$\begin{aligned}
 &t(t - 1) \left(\mathcal{E}(x^k, y^k) - \mathcal{E}(x, y) \right) - \frac{t-1}{2} \left(\|x - x^k\|_F^2 + \|y - y^k\|_G^2 \right. \\
 &\quad \left. + \|A(x - x^k) + B(y - y^k)\|^2 \right) + \frac{1}{2K} \|x + (t - 1)x^k - t\bar{x}^k\|_M^2 \\
 &\quad + \frac{1}{2L} \|y + (t - 1)y^k - t\bar{y}^k\|_N^2 + \frac{1}{2} \|y + (t - 1)y^k - t\bar{y}^k\|_{B^*B}^2 \\
 &\geq t^2 \left(\mathcal{E}(x^{k+1}, y^{k+1}) - \mathcal{E}(x, y) \right) + \frac{1}{2K} \|x + (t - 1)x^k - t\hat{x}^{k+1}\|_M^2 \\
 &\quad + \frac{1}{2L} \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_N^2 + \frac{1}{2} \|x + (t - 1)x^k - t\hat{x}^{k+1}\|_F^2 \\
 &\quad + \frac{1}{2} \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_{B^*B+G}^2. \tag{13}
 \end{aligned}$$

In the case of alternating minimizations, this simplifies a lot. Assuming that $(x^{k+1}, y^{k+1}) = (\tilde{x}, \tilde{y}) = (\hat{x}, \hat{y})$, one deduces from (12) that

$$\begin{aligned}
 &t(t - 1) \left(\mathcal{E}(x^k, y^k) - \mathcal{E}(x, y) \right) \\
 &\quad - \frac{t-1}{2} \left(\|x - x^k\|_F^2 \right. \\
 &\quad \left. + \|y - y^k\|_G^2 + \|A(x - x^k) + B(y - y^k)\|^2 \right) \\
 &\quad + \frac{1}{2} \|y + (t - 1)y^k - t\bar{y}^k\|_{B^*B}^2 \\
 &\geq t^2 \left(\mathcal{E}(x^{k+1}, y^{k+1}) - \mathcal{E}(x, y) \right) \\
 &\quad + \frac{1}{2} \|x + (t - 1)x^k - t\hat{x}^{k+1}\|_F^2 \\
 &\quad + \frac{1}{2} \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_{B^*B+G}^2. \tag{14}
 \end{aligned}$$

The convergence rates will be derived from these main inequalities.

3 Accelerated Alternating Minimization

In [11], it is shown (in case $F = G = 0$) how one can derive an accelerated algorithm, in the spirit of the ‘FISTA’ [3] method, from inequality (14). In this section, we extend these results to the strongly convex case, yielding better (linear) rates. As alternating minimizations for two variables are essentially equivalent to a forward–backward algorithm [11, 13], it is clear that an acceleration in the spirit of [19, Thm. 2.2.2] will provide efficient convergence rates. A derivation from an equality similar to (14) is provided in [12, Appendix B]. We provide here an adaption of that proof to our particular situation (only the parameters are slightly differing, so that we will sketch most of the arguments). We make the assumption that for some nonnegative parameters γ, δ ,

$$G \geq \gamma B^* B, \quad F \geq \delta A^* A, \tag{15}$$

and we assume $\gamma + \delta > 0$. This assumption, which may look strange at first glance, expresses that g and/or f are, respectively, strongly convex with respect to the variables $B y$ and $A x$, which appear in the quadratic term. They are obviously satisfied if g, f are strongly convex in the classical sense. Observe that

$$\begin{aligned} & \|x - x^k\|_F^2 + \|y - y^k\|_G^2 + \|A(x - x^k) + B(y - y^k)\|^2 \\ & \geq (\gamma + 1)\|y - y^k\|_{B^* B}^2 \\ & \quad + (\delta + 1)\|x - x^k\|_{A^* A}^2 + 2\langle A(x - x^k), B(y - y^k) \rangle \\ & \geq \left(\gamma + \frac{\delta}{1 + \delta}\right) \|y - y^k\|_{B^* B}^2, \end{aligned}$$

hence, denoting $\gamma' = \gamma + \delta/(1 + \delta) > 0$, it follows from (14) that

$$\begin{aligned} & t(t - 1) \left(\mathcal{E}(x^k, y^k) - \mathcal{E}(x, y) \right) - \gamma'^{\frac{t-1}{2}} \|y - y^k\|_{B^* B}^2 \\ & \quad + \frac{1}{2} \|y + (t - 1)y^k - t\bar{y}^k\|_{B^* B}^2 \\ & \geq t^2 \left(\mathcal{E}(x^{k+1}, y^{k+1}) - \mathcal{E}(x, y) \right) \\ & \quad + \frac{1+\gamma}{2} \|y + (t - 1)y^k - t\bar{y}^{k+1}\|_{B^* B}^2. \tag{16} \end{aligned}$$

As in [12, Appendix B], we first collapse the two quadratic terms in the left-hand side as follows (assuming $(t - 1)\gamma' \neq 1$):

$$\begin{aligned} & -\gamma'^{\frac{t-1}{2}} \|y - y^k\|_{B^* B}^2 \\ & \quad + \frac{1}{2} \|y - y^k + t(y^k - \bar{y}^k)\|_{B^* B}^2 \\ & = \frac{1-(t-1)\gamma'}{2} \|y - y^k\|_{B^* B}^2 + t \langle y - y^k, y^k - \bar{y}^k \rangle_{B^* B} \\ & \quad + \frac{t^2}{2} \|y^k - \bar{y}^k\|_{B^* B}^2 \\ & = \frac{1-(t-1)\gamma'}{2} \|y - y^k + \frac{t}{1-(t-1)\gamma'}(y^k - \bar{y}^k)\|_{B^* B}^2 \\ & \quad + \left(\frac{t^2}{2} - \frac{t^2}{2(1-(t-1)\gamma')} \right) \|y^k - \bar{y}^k\|_{B^* B}^2 \\ & \leq \frac{1-(t-1)\gamma'}{2} \left\| y - y^k + \frac{t}{1-(t-1)\gamma'}(y^k - \bar{y}^k) \right\|_{B^* B}^2 \end{aligned}$$

provided $0 < 1 - (t - 1)\gamma' \leq 1$, which we now assume (that is, $1 \leq t < 1 + 1/\gamma'$). Equation (16) becomes, assuming $t = t_{k+1}$ is now a variable parameter,

$$\begin{aligned} & t_{k+1}(t_{k+1} - 1) \left(\mathcal{E}(x^k, y^k) - \mathcal{E}(x, y) \right) \\ & \quad + \frac{1-(t_{k+1}-1)\gamma'}{2} \left\| y - y^k + \frac{t_{k+1}}{1-(t_{k+1}-1)\gamma'}(y^k - \bar{y}^k) \right\|_{B^* B}^2 \\ & \geq t_{k+1}^2 \left(\mathcal{E}(x^{k+1}, y^{k+1}) - \mathcal{E}(x, y) \right) \\ & \quad + \frac{1+\gamma}{2} \|y + (t_{k+1} - 1)y^k - t_{k+1}\bar{y}^{k+1}\|_{B^* B}^2. \end{aligned}$$

Denoting $\omega_k = (1 - (t_{k+1} - 1)\gamma')/(1 + \gamma) \leq 1$, we find that provided $t_{k+1}^2 - t_{k+1} = \omega_k t_k^2$, this inequality becomes

$$\begin{aligned} & t_{k+1}^2 \left(\mathcal{E}(x^{k+1}, y^{k+1}) - \mathcal{E}(x, y) \right) \\ & \quad + \frac{1+\gamma}{2} \|y + (t_{k+1} - 1)y^k - t_{k+1}\bar{y}^{k+1}\|_{B^* B}^2 \\ & \leq \omega_k \left(t_k^2 \left(\mathcal{E}(x^k, y^k) - \mathcal{E}(x, y) \right) \right. \\ & \quad \left. + \frac{1+\gamma}{2} \|y - y^k + \frac{t_{k+1}}{1-(t_{k+1}-1)\gamma'}(y^k - \bar{y}^k)\|_{B^* B}^2 \right). \end{aligned}$$

Hence, provided one ensures

$$y - y^k + \frac{t_{k+1}}{1-(t_{k+1}-1)\gamma'}(y^k - \bar{y}^k) = y + (t_k - 1)y^{k-1} - t_k y^k, \tag{17}$$

it will follow

$$\begin{aligned} & \mathcal{E}(x^k, y^k) - \mathcal{E}(x, y) \\ & \leq \frac{1}{t_k^2} \left(\prod_{i=0}^{k-1} \omega_i \right) \left(t_0^2 \left(\mathcal{E}(x^0, y^0) - \mathcal{E}(x, y) \right) \right. \\ & \quad \left. + \frac{1+\gamma}{2} \|y - y^0\|_{B^* B}^2 \right). \tag{18} \end{aligned}$$

The update rules for t_k, ω_k, \bar{y}^k , ensuring in particular (17), should be as follows:

$$t_{k+1} = \frac{1}{2} \left(1 - \frac{\gamma'}{1+\gamma} t_k^2 + \sqrt{\left(1 - \frac{\gamma'}{1+\gamma} t_k^2 \right)^2 + 4 \frac{1+\gamma'}{1+\gamma} t_k^2} \right), \tag{19}$$

$$\beta_k = (1 - (t_{k+1} - 1)\gamma') \frac{t_k - 1}{t_{k+1}}, \tag{20}$$

$$\bar{y}^k = y^k + \beta_k (y^k - y^{k-1}), \tag{21}$$

$$\begin{aligned} \omega_k & = \frac{1 + \gamma' - t_{k+1}\gamma'}{1 + \gamma} = 1 - t_{k+1} \frac{\gamma}{1 + \gamma} \\ & \quad - (t_{k+1} - 1) \frac{\delta}{(1 + \delta)(1 + \gamma)}. \tag{22} \end{aligned}$$

It remains to estimate the rates which these updates yield. This is done in a similar way as in [19, Chap. 2] and [12, Appendix B]. A starting point is the update equation for t_{k+1} , which is chosen as the nonnegative root of:

$$t_{k+1}^2 - t_{k+1} = \omega_k t_k^2 = \frac{1 - (t_{k+1} - 1)\gamma'}{1 + \gamma} t_k^2 \tag{23}$$

which, letting $q' := \gamma'/(1 + \gamma')$, also reads

$$t_{k+1}^2 = t_{k+1} + \frac{1 + \gamma'}{1 + \gamma} (1 - q' t_{k+1}) t_k^2.$$

A first fact is that $1 \leq t_{k+1} < 1/q'$. Indeed if $q' t_{k+1} \geq 1$ one obtains that $t_{k+1}^2 \leq t_{k+1}$, hence $t_{k+1} \leq 1$ and $q' t_{k+1} < 1$,

a contradiction. Hence, $q't_{k+1} < 1$, and $t_{k+1}^2 \geq t_{k+1}$ so that $t_{k+1} \geq 1$ (moreover, if $t_k = 0$, which is possible only for $k = 0$ one has $t_{k+1} = 1$, otherwise, $t_{k+1} > 1$). As a consequence for any $k \geq 0$,

$$0 < \omega_k \leq \frac{1}{1 + \gamma} < 1. \tag{24}$$

But we can prove better and actually show that the convergence rate is linear.

For this, we observe first that thanks to (22), for any $q'' > 0$ one has

$$q''t_{k+1}^2 = q''t_{k+1} + \left(1 - t_{k+1} \frac{\gamma}{1+\gamma} - (t_{k+1} - 1) \frac{\delta}{(1+\delta)(1+\gamma)}\right) q''t_k^2$$

and we deduce that $q''t_{k+1}^2$ is less than (or equal to) a convex combination of 1 and $q''t_k^2$ as soon as

$$q''t_{k+1} \leq t_{k+1} \frac{\gamma}{1+\gamma} + (t_{k+1} - 1) \frac{\delta}{(1+\delta)(1+\gamma)},$$

which is as soon as

$$\mathcal{E}(x^k, y^k) - \mathcal{E}(x, y) \leq \begin{cases} (1 - \sqrt{q''})^{k-1} (t_0^2 (\mathcal{E}(x^0, y^0) - \mathcal{E}(x, y)) + \frac{1+\gamma}{2} \|y - y^0\|_{B^*B}^2) \\ (1 - \sqrt{q''})^k \left(\mathcal{E}(x^0, y^0) - \mathcal{E}(x, y) + \frac{1+\gamma}{2t_0^2} \|y - y^0\|_{B^*B}^2 \right) \end{cases} \quad (\text{if } t_0 \neq 0).$$

$$\left(\frac{\gamma'}{1+\gamma} - q''\right) t_{k+1} \geq \frac{\delta}{(1+\delta)(1+\gamma)}. \tag{25}$$

We choose q'' given by the equation

$$\left(\frac{\gamma'}{1+\gamma} - q''\right) \frac{1}{\sqrt{q''}} = \frac{\delta}{(1+\delta)(1+\gamma)}. \tag{26}$$

Then, if $q''t_k^2 < 1$, we find that in case $t_{k+1} \geq 1/\sqrt{q''}$, inequality (25) holds (thanks to (26)), so that $q''t_{k+1}^2$ is less or equal to a convex combination of 1 and $q''t_k^2$, and it follows $q''t_{k+1}^2 < 1$ which is contradictory. In consequence, $q''t_{k+1}^2 < 1$. Hence by induction, if we assume that

$$t_0 \in [0, 1/\sqrt{q''}], \tag{27}$$

we find that $t_k < 1/\sqrt{q''}$ for all $k \geq 1$. The value of $\sqrt{q''}$, obtained by solving equation (26), is

$$\sqrt{q''} = \sqrt{\frac{\gamma}{1+\gamma} + \frac{\delta}{(1+\delta)(1+\gamma)} + \frac{\delta^2}{4(1+\delta)^2(1+\gamma)^2} - \frac{\delta}{2(1+\delta)(1+\gamma)}}. \tag{28}$$

If $\delta = 0, q'' = \gamma/(1+\gamma) > 0$. If not one can check that it is larger (as the function $t \mapsto \sqrt{a + 2t + t^2} - t$, for $0 \leq a < 1$, is concave and increasing). Using that $\sqrt{2a + 2b} \geq \sqrt{a} + \sqrt{b}$ and $\delta/(1 + \delta)(1 + \gamma) \geq (\delta/(1 + \delta)(1 + \gamma))^2$, one can also check that

$$q'' \geq \frac{1}{2} \frac{\gamma}{1 + \gamma} + \frac{3}{8} \frac{\delta}{(1 + \delta)(1 + \gamma)}.$$

In particular, it follows from (23) that $t_k^2 \omega_k / t_{k+1}^2 = 1 - 1/t_{k+1} \leq 1 - \sqrt{q''}$ so that

$$\begin{aligned} \theta_k &:= \frac{1}{t_k^2} \left(\prod_{i=0}^{k-1} \omega_i \right) = \frac{\omega_0}{t_1^2} \prod_{i=1}^{k-1} \frac{t_i^2 \omega_i}{t_{i+1}^2} \\ &\leq \frac{\omega_0}{t_1^2} \left(1 - \sqrt{q''}\right)^{k-1} \leq \left(1 - \sqrt{q''}\right)^{k-1}, \end{aligned}$$

where we have used (24) and $t_1 \geq 1$. In addition, if $t_0 > 0$, one also finds similarly the bound $\theta_k \leq (1 - \sqrt{q''})^k / t_0^2$.

One deduces from (18) that if (27), (19), (20) and (21) hold, then

Eventually, it is straightforward to check that also [12, (B.10)] holds, that is, $\theta_k \leq 4/(k + 1)^2$. It follows the result:

Algorithm 1 Accelerated alternating minimizations

Input: Metrics F, G , parameters γ, δ satisfying (15).
 Let then $q := \gamma/(1 + \gamma) \in [0, 1], \gamma' := \gamma + \delta/(1 + \delta)$.
 Choose $(x^0, y^0), t_0 \in [0, 1/\sqrt{q''}]$ and let $\bar{y}^0 = y^0$.
for all $k \geq 1$ **do**
 $x^k = \arg \min_{x \in \mathcal{X}} \mathcal{E}(x, \bar{y}^{k-1})$,
 $y^k = \arg \min_{y \in \mathcal{Y}} \mathcal{E}(x^k, y)$,
 then compute $t_{k+1}, \beta_k, \bar{y}^k$ according to (19), (20), (21).
end for

Theorem 1 Let (x^k, y^k) be computed according to Algorithm 1. Then for any $k \geq 1$, one has

$$\begin{aligned} \mathcal{E}(x^k, y^k) - \min_{x,y} \mathcal{E} &\leq \theta_k (t_0^2 (\mathcal{E}(x^0, y^0) - \mathcal{E}(x, y)) \\ &\quad + \frac{1+\gamma}{2} \|y - y^0\|_{B^*B}^2) \end{aligned} \tag{29}$$

where

$$\theta_k \leq \min \left\{ \frac{4}{(k + 1)^2}, (1 - \sqrt{q''})^{k-1}, \frac{(1 - \sqrt{q''})^k}{t_0^2} \right\}.$$

4 Accelerated Alternating Descent

Now, we show that this analysis can be adapted, in theory, to yield also accelerated algorithms for the alternating descent method. The idea is to adapt the previous proof to the more complex inequality (13). We will do this in a slightly sub-optimal way, considering only, in the strong convex case, a constant overrelaxation, in order to make the paper more readable.

4.1 The Nonstrongly Convex Case

We first consider the case $F, G = 0$, for which the computations are substantially simpler to read. In this case, (13) boils down to

$$\begin{aligned}
 & t(t-1)(\mathcal{E}(x^k, y^k) - \mathcal{E}(x, y)) + \frac{1}{2K} \|x + (t-1)x^k - t\bar{x}^k\|_M^2 \\
 & + \frac{1}{2L} \|y + (t-1)y^k - t\bar{y}^k\|_N^2 + \frac{1}{2} \|y + (t-1)y^k - t\bar{y}^k\|_{B^*B}^2 \\
 & \geq t^2 (\mathcal{E}(x^{k+1}, y^{k+1}) - \mathcal{E}(x, y)) + \frac{1}{2K} \|x + (t-1)x^k - t\hat{x}^{k+1}\|_M^2 \\
 & + \frac{1}{2L} \|y + (t-1)y^k - t\hat{y}^{k+1}\|_N^2 + \frac{1}{2} \|y + (t-1)y^k - t\hat{y}^{k+1}\|_{B^*B}^2.
 \end{aligned} \tag{30}$$

The standard proof of ‘FISTA’ [3] consists then in letting $t_0 = 0$, and for $k \geq 0$, $t_{k+1} = (1 + \sqrt{1 + 4t_k^2})/2$ (so that $t_{k+1}(t_{k+1} - 1) = t_k^2$) (one can more generally choose, for $k \geq 1$, $t_k = 1 + (k-1)/a$, $a \geq 2$, so that $t_{k+1}(t_{k+1} - 1) \leq t_k^2$ for all k [8]; then the following inequalities will continue to hold as long as we also assume that (x, y) is a minimizer of the energy). It follows

$$\begin{aligned}
 & t_{k+1}^2 (\mathcal{E}(x^{k+1}, y^{k+1}) - \mathcal{E}(x, y)) \\
 & + \frac{1}{2K} \|x + (t_{k+1}-1)x^k - t_{k+1}\bar{x}^{k+1}\|_M^2 \\
 & + \frac{1}{2L} \|y + (t_{k+1}-1)y^k - t_{k+1}\bar{y}^{k+1}\|_N^2 + \frac{1}{2} \|y + (t_{k+1}-1)y^k - t_{k+1}\bar{y}^{k+1}\|_{B^*B}^2 \\
 & \leq t_k^2 (\mathcal{E}(x^k, y^k) - \mathcal{E}(x, y)) + \frac{1}{2K} \|x + (t_{k+1}-1)x^k - t_{k+1}\bar{x}^k\|_M^2 \\
 & + \frac{1}{2L} \|y + (t_{k+1}-1)y^k - t_{k+1}\bar{y}^k\|_N^2 + \frac{1}{2} \|y + (t_{k+1}-1)y^k - t_{k+1}\bar{y}^k\|_{B^*B}^2.
 \end{aligned} \tag{31}$$

It remains to choose $\bar{x}^k, \bar{y}^k, \hat{y}^k$, to ensure the following equalities:

$$\begin{aligned}
 x + (t_{k+1} - 1)x^k - t_{k+1}\bar{x}^k &= x + (t_k - 1)x^{k-1} - t_k\hat{x}^k \\
 y + (t_{k+1} - 1)y^k - t_{k+1}\bar{y}^k &= y + (t_k - 1)y^{k-1} - t_k\hat{y}^k \\
 y + (t_{k+1} - 1)y^k - t_{k+1}\bar{y}^k &= y + (t_k - 1)y^{k-1} - t_k\hat{y}^k.
 \end{aligned}$$

This is obtained by letting

$$\bar{x}^k = x^k + \frac{t_k-1}{t_{k+1}} (x^k - x^{k-1}) + \frac{t_k}{t_{k+1}} (\hat{x}^k - x^k) \tag{32}$$

$$\bar{y}^k = y^k + \frac{t_k-1}{t_{k+1}} (y^k - y^{k-1}) + \frac{t_k}{t_{k+1}} (\hat{y}^k - y^k) \tag{33}$$

$$\bar{y}^k = y^k + \frac{t_k-1}{t_{k+1}} (y^k - y^{k-1}) + \frac{t_k}{t_{k+1}} (\hat{y}^k - y^k) \tag{34}$$

With these choices, one can eventually sum (31) from $n = 0$ to $k - 1$ and it follows

$$\mathcal{E}(x^n, y^n) - \mathcal{E}(x, y) \leq \frac{\|x - x^0\|_{M/K}^2 + \|y - y^0\|_{N/L+B^*B}^2}{2t_n^2}.$$

Using the fact that by construction, $t_{k+1} \geq t_k + 1/2$ and $t_1 \geq 1$, and choosing for (x, y) a minimizer, we deduce the following theorem:

Algorithm 2 Accelerated alternating descent method (general case)

Input: Metrics M, N , number of inner loops $K, L \geq 1$.
 Choose $(x^0, y^0) \in \mathcal{X} \times \mathcal{Y}$, $t_0 \geq 0$, let $\bar{x}^0 = x^0, \bar{y}^0 = \bar{y}^0 = y^0$.
for all $k \geq 1$ **do**
 Find $(\hat{x}^k, \hat{y}^k, \hat{x}^k, \hat{y}^k) = T_{K,L}(\bar{x}^{k-1}, \bar{y}^{k-1}, \bar{y}^{k-1})$ (cf Eq. (2-5)).
 Choose a point (x^k, y^k) such that $\mathcal{E}(x^k, y^k) \leq \mathcal{E}(\hat{x}^k, \hat{y}^k)$, for instance $(x^k, y^k) = (\hat{x}^k, \hat{y}^k)$,
 then compute $t_{k+1} = (1 + \sqrt{1 + 4t_k^2})/2$, $(\bar{x}^k, \bar{y}^k, \bar{y}^k)$ according to (32), (33), (34).
end for

Theorem 2 Let (x^k, y^k) be computed using Algorithm 2, starting from initial points (x^0, y^0) , and let (x^*, y^*) be a minimizer of \mathcal{E} . Then one has the global rate:

$$\begin{aligned}
 & \mathcal{E}(x^k, y^k) - \mathcal{E}(x^*, y^*) \\
 & \leq 2 \frac{\|x^* - x^0\|_{M/K}^2 + \|y^* - y^0\|_{N/L+B^*B}^2}{(k+1)^2}.
 \end{aligned} \tag{35}$$

5 The Strongly Convex Case

The case where $F, G > 0$ is a bit trickier, if one wants to exploit it to gain a better (linear) convergence. The main observation is that in (13), the (unknown) points x and y on the left-hand side of the inequality are evaluated, respectively, in the M/K and $N/L + B^*B$ norms, while on the right-hand side, they are evaluated in the $M/K + tF$ and $N/L + B^*B + tG$ norms, respectively. (In fact, one should also, as in Sect. 3, use the term involving $Ax + By$ to transfer some control from x to y or conversely, leading to more tedious calculations—this would be necessary for instance if only one of the metrics F, G were positive, which to simplify we do not assume here.) It follows that if one can choose t such that, for some $\omega < 1$,

$$\begin{aligned}
 \frac{1}{K}M + tF &\geq \omega^{-1} \frac{1}{K}M, & \frac{1}{L}N + B^*B \\
 + tG &\geq \omega^{-1} (\frac{1}{L}N + B^*B), & t^2 \geq \omega^{-1}(t(t-1)),
 \end{aligned}$$

then the energy can be reduced by a constant ratio at each iteration. The last inequality suggests that ω should be simply equal to $1 - 1/t$, and the optimal $t \geq 1$ is the smallest value such that

$$t(t - 1)F \geq \frac{1}{K}M, \quad t(t - 1)G \geq \frac{1}{L}N + B^*B. \tag{36}$$

In practice, M and N are often chosen of the form $I/\tau - A^*A$ and $I/\sigma - B^*B$, respectively, so that the descent steps in x, y can be computed explicitly. One needs $\tau \leq \|A^*A\|^{-1}$ and $\sigma \leq \|B^*B\|^{-1}$ in order for M, N to be nonnegative. The condition on G then boils down to $t(t - 1)G \geq I/(\sigma L) + (1 - 1/L)B^*B$, which is ensured as soon as $t - 1 \geq \sqrt{\|G^{-1}\|/\sigma}$, while the condition on F is ensured if $t - 1 \geq \sqrt{\|F^{-1}\|/(K\tau)}$ (and is thus in general easier to ensure). In any case, the geometric ratio involves the square root of the condition number of the problems in x and y , which indicates that the accelerated algorithm we can derive should have good performances.

Let us now derive an implementable algorithm. We assume now that (36) holds and denote $\omega = 1 - 1/t$. A more general derivation in the spirit of [12, 19] with variable t could be derived as in Sect. 3, but the calculations would be much more tedious. Estimate (13) can be rewritten:

$$\begin{aligned} \omega t^2 \left(\mathcal{E}(x^k, y^k) - \mathcal{E}(x, y) \right) &+ \frac{1}{2K} \|x + (t - 1)x^k - t\bar{x}^k\|_M^2 \\ &+ \frac{1}{2L} \|y + (t - 1)y^k - t\bar{y}^k\|_N^2 + \frac{1}{2} \|y + (t - 1)y^k - t\bar{y}^k\|_{B^*B}^2 \\ \geq t^2 \left(\mathcal{E}(x^{k+1}, y^{k+1}) - \mathcal{E}(x, y) \right) &+ \frac{1}{2K} \|x + (t - 1)x^k - t\hat{x}^{k+1}\|_M^2 \\ &+ \frac{1}{2L} \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_N^2 + \frac{1}{2} \|x + (t - 1)x^k - t\hat{x}^{k+1}\|_F^2 \\ &+ \frac{1}{2} \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_{B^*B+G}^2 + \frac{t-1}{2} \left(\|x - x^k\|_F^2 + \|y - y^k\|_G^2 \right) \end{aligned} \tag{37}$$

First, one observes that using (36),

$$\begin{aligned} &\frac{1}{2K} \|x + (t - 1)x^k - t\hat{x}^{k+1}\|_M^2 \\ &+ \frac{1}{2} \|x + (t - 1)x^k - t\hat{x}^{k+1}\|_F^2 + \frac{t-1}{2} \|x - x^k\|_F^2 \\ \geq \frac{1}{2K} \left(\|x + (t - 1)x^k - t\hat{x}^{k+1}\|_M^2 \right. \\ &\quad \left. + \frac{1}{t(t-1)} \|x + (t - 1)x^k - t\hat{x}^{k+1}\|_M^2 + \frac{1}{t} \|x - x^k\|_M^2 \right) \\ &= \frac{t}{2K(t-1)} \left(\frac{t-1}{t} \|x + (t - 1)x^k - t\hat{x}^{k+1}\|_M^2 \right. \\ &\quad \left. + \frac{1}{t^2} \|x + (t - 1)x^k - t\hat{x}^{k+1}\|_M^2 + \frac{t-1}{t^2} \|x - x^k\|_M^2 \right) \\ &\geq \frac{t}{2K(t-1)} \left\| x + \frac{(t-1)^2}{t} x^k - (t - 1)\hat{x}^{k+1} - \frac{1}{t}\hat{x}^{k+1} \right\|_M^2. \end{aligned} \tag{38}$$

Hence a good choice for \bar{x}^k is a choice which ensures that

$$x + (t - 1)x^k - t\bar{x}^k = x + \frac{(t-1)^2}{t} x^{k-1} - (t - 1)\hat{x}^k - \frac{1}{t}\hat{x}^k,$$

yielding

$$\begin{aligned} \bar{x}^k &= x^k + \left(\frac{t - 1}{t} \right)^2 (x^k - x^{k-1}) + \frac{t - 1}{t} (\hat{x}^k - x^k) \\ &\quad + \frac{1}{t^2} (\hat{x}^k - x^k). \end{aligned} \tag{39}$$

The situation is slightly more complicated for the y variable, but the computations are very similar. One has, using (36),

$$\begin{aligned} &\frac{1}{2L} \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_N^2 \\ &\quad + \frac{1}{2} \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_{B^*B+G}^2 + \frac{t-1}{2} \|y - y^k\|_G^2 \\ \geq \frac{1}{2L} \left(\|y + (t - 1)y^k - t\hat{y}^{k+1}\|_N^2 \right. \\ &\quad \left. + \frac{1}{t(t-1)} \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_N^2 + \frac{1}{t} \|y - y^k\|_N^2 \right) \\ &\quad + \frac{1}{2} \left(\left(1 + \frac{1}{t(t-1)} \right) \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_{B^*B}^2 \right. \\ &\quad \left. + \frac{1}{t} \|y - y^k\|_{B^*B}^2 \right). \end{aligned}$$

As in (38) (replacing x with y and M with N),

$$\begin{aligned} &\frac{1}{2L} \left(\|y + (t - 1)y^k - t\hat{y}^{k+1}\|_N^2 \right. \\ &\quad \left. + \frac{1}{t(t-1)} \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_N^2 + \frac{1}{t} \|y - y^k\|_N^2 \right) \\ \geq \frac{t}{2L(t-1)} \|y + \frac{(t-1)^2}{t} y^k - (t - 1)\hat{y}^{k+1} - \frac{1}{t}\hat{y}^{k+1}\|_N^2, \end{aligned}$$

while the second expression is a bit simpler:

$$\begin{aligned} &\frac{1}{2} \left(\left(1 + \frac{1}{t(t-1)} \right) \|y + (t - 1)y^k - t\hat{y}^{k+1}\|_{B^*B}^2 \right. \\ &\quad \left. + \frac{1}{t} \|y - y^k\|_{B^*B}^2 \right) \\ \geq \frac{t}{2(t-1)} \|y + \frac{(t-1)^2}{t} y^k - (t - 1)\hat{y}^{k+1} - \frac{1}{t}\hat{y}^{k+1}\|_{B^*B}^2 \\ &= \frac{t}{2(t-1)} \|y + \frac{(t-1)^2}{t} y^k - \left(1 + \frac{(t-1)^2}{t} \right) \hat{y}^{k+1}\|_{B^*B}^2. \end{aligned}$$

We will therefore choose \bar{y}^k, \bar{y}'^k to satisfy

$$\begin{aligned} y + (t - 1)y^k - t\bar{y}^k &= y + \frac{(t-1)^2}{t} y^{k-1} - (t - 1)\hat{y}^k - \frac{1}{t}\hat{y}^k, \\ y + (t - 1)y^k - t\bar{y}'^k &= y + \frac{(t-1)^2}{t} (y^{k-1} - \hat{y}^k) - \hat{y}^k, \end{aligned}$$

which is ensured provided

$$\begin{aligned} \bar{y}^k &= y^k + \left(\frac{t - 1}{t} \right)^2 (y^k - y^{k-1}) \\ &\quad + \frac{t - 1}{t} (\hat{y}^k - y^k) + \frac{1}{t^2} (\hat{y}^k - y^k). \end{aligned} \tag{40}$$

$$\bar{y}'^k = y^k + \left(\frac{t - 1}{t} \right)^2 (\hat{y}^k - y^{k-1}) + \frac{1}{t} (\hat{y}^k - y^k). \tag{41}$$

With choices (39), (40), (41), inequality (37) becomes

$$\begin{aligned} &\omega t^2 \left(\mathcal{E}(x^k, y^k) - \mathcal{E}(x, y) \right) \\ &+ \frac{1}{2K} \|x + \frac{(t-1)^2}{t} x^{k-1} - (t-1)\hat{x}^k - \frac{1}{t}\hat{x}^{k+1}\|_M^2 \\ &+ \frac{1}{2L} \|y + \frac{(t-1)^2}{t} y^{k-1} - (t-1)\hat{y}^k - \frac{1}{t}\hat{y}^{k+1}\|_N^2 \\ &+ \frac{1}{2} \|y + \frac{(t-1)^2}{t} (y^{k-1} - \hat{y}^{k-1}) - \hat{y}^{k+1}\|_{B^*B}^2 \\ &\geq t^2 \left(\mathcal{E}(x^{k+1}, y^{k+1}) - \mathcal{E}(x, y) \right) \\ &+ \frac{1}{2K\omega} \|x + \frac{(t-1)^2}{t} x^k - (t-1)\hat{x}^{k+1} - \frac{1}{t}\hat{x}^{k+2}\|_M^2 \\ &+ \frac{1}{2L\omega} \|y + \frac{(t-1)^2}{t} y^k - (t-1)\hat{y}^{k+1} - \frac{1}{t}\hat{y}^{k+2}\|_N^2 \\ &+ \frac{1}{2\omega} \|y + \frac{(t-1)^2}{t} y^k - (1 + \frac{(t-1)^2}{t})\hat{y}^{k+1}\|_{B^*B}^2, \end{aligned}$$

so that one has (assuming $\bar{x}^0 = x^0, \bar{y}^0 = \bar{y}^0 = y^0$)

$$\begin{aligned} \mathcal{E}(x^k, y^k) - \mathcal{E}(x, y) &\leq \omega^k \left(\mathcal{E}(x^0, y^0) - \mathcal{E}(x, y) \right) \\ &+ \frac{1}{t^2\omega} \left(\frac{1}{2K} \|x - x^0\|_M^2 + \frac{1}{2} \|y - y^0\|_{N/L+B^*B}^2 \right). \end{aligned}$$

Hence one has in this case a linear convergence rate.

Algorithm 3 Accelerated alternating descent method (strongly convex case)

Input: Metrics M, N and F, G , number of inner loops $K, L \geq 1$.
 Choose $(x^0, y^0) \in \mathcal{X} \times \mathcal{Y}$, t such that (36) holds. Let $\bar{x}^0 = x^0, \bar{y}^0 = \bar{y}^0 = y^0$.
for all $k \geq 1$ **do**
 Find $(\hat{x}^k, \hat{y}^k, \hat{x}^{k+1}, \hat{y}^{k+1}) = T_{K,L}(\bar{x}^{k-1}, \bar{y}^{k-1}, \bar{y}^{k-1})$ (cf Eq. (2-5)).
 Choose a point (x^k, y^k) such that $\mathcal{E}(x^k, y^k) \leq \mathcal{E}(\hat{x}^k, \hat{y}^k)$, for instance $(x^k, y^k) = (\hat{x}^k, \hat{y}^k)$,
 then compute $(\bar{x}^k, \bar{y}^k, \bar{y}^k)$ according to (39), (40), (41).
end for

Theorem 3 Let (x^k, y^k) be computed using Algorithm 3, starting from initial points (x^0, y^0) , and let (x^*, y^*) be a minimizer of \mathcal{E} . Then the energy decays with the linear rate:

$$\begin{aligned} &\mathcal{E}(x^k, y^k) - \mathcal{E}(x^*, y^*) \\ &\leq \omega^k \left(\mathcal{E}(x^0, y^0) - \mathcal{E}(x^*, y^*) \right) \\ &+ \frac{1}{t^2\omega} \left(\frac{1}{2K} \|x^* - x^0\|_M^2 + \frac{1}{2} \|y^* - y^0\|_{N/L+B^*B}^2 \right) \end{aligned} \quad (42)$$

where $\omega = 1 - 1/t$.

6 Experiments: A Toy Model

Before implementing the methods discussed in this paper—or variants—on relatively large-scale problems, we consider as a toy model the minimization of the elementary signal denoising energy:

$$\min_{x,y} \sum_{i=1}^{N-1} \left(\lambda |x_i| + \frac{\mu}{2} |x_i - (Dy)_i|^2 \right) + \frac{1}{2} \sum_{i=1}^N (y_i - g_i)^2. \quad (43)$$

here, $g = (g_i)_{i=1}^N$ is a noisy signal, $y = (y_i)_{i=1}^N$ its reconstruction and $x = (x_i)_{i=1}^{N-1}$ is an approximation of the discrete derivative Dy . D is an operator from \mathbb{R}^N to \mathbb{R}^{N-1} , defined by $(Dy)_i = (y_{i+1} - y_i)$, $i = 1, \dots, N - 1$; it has norm $\|D\| \leq 2$. If we eliminate x (by direct minimization) in (43), we get the so-called Huber-TV regularization problem for the signal y .

This setting corresponds to our general setting with $A = \sqrt{\mu}I$, $B = \sqrt{\mu}D$ (of squared norm 4μ), $F = 0, G = I$. Observe that here still, (36) will be satisfied if $M = 0$, that is, if we minimize the problem exactly with respect to x at each iteration, hence the analysis is valid even though the l^1 norm is not strongly convex.

Figure 1, left, shows the typical result of this minimization for a given signal (here $N = 10,000, \mu = 150, \lambda = 300$).

Energy (43) can be minimized by many techniques. The most natural approach consists in eliminating the variable x and implementing a (strongly convex) ‘FISTA’ [12], which is basically the same as alternatively minimizing exactly the problem with respect to x and then performing one descent step ($L = 1$) with respect to y , with an appropriate acceleration. This gives the solid curve in Fig. 1, right.

While this method is (as expected) very efficient, we find that taking $L = 3$ (i.e., performing 3 successive steps of descent in the y variable before updating it again) reduces the total number of iterations (we observe however that the outer iterations are significantly slower so that the overall gain is negligible for this problem).

The dotted line tagged ‘ $L = 3, no\ averaging$ ’ in Fig. 1 is obtained by removing the averaging process in the algorithm. We observe that (despite we have no theoretical explanation), the behavior of the method is exactly the same as with averaging (while the global execution time is decreased and less memory is used). In any case, the performance of the accelerated methods is significantly better than the performance of the nonaccelerated version (dashed line).

Eventually, the reason for which we do not plot the rate of a basic ‘FISTA’ method applied to the variable (u, v) is that we do not know of any result which deals with the case where none of the smooth $(\mu \|x - Dy\|^2/2)$ and nonsmooth $(\lambda \|x\|_1 + \|y - g\|^2/2)$ term is strongly convex (here only the sum is), and shows linear convergence. A standard ($O(1/k^2)$) FISTA method turns out to be much slower than the four methods in Fig. 1, and even by tuning by hand the parameters to take into account the strong convexity we could not obtain competitive rates.

This experiment (and a few other we tried) shows that in practice the averaging required by our algorithms does not

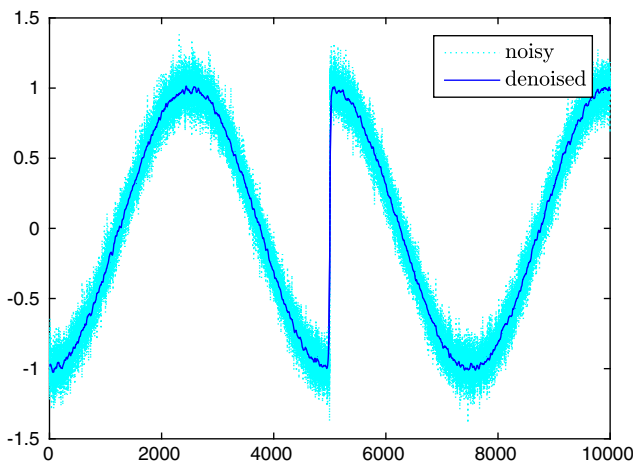


Fig. 1 A signal denoising test

seem necessary to ensure fast convergence. Actually, this is something we had observed before starting this theoretical study, which was aimed at proving it, however without success. In practice, for the large-scale problems we will consider in the next section, we will replace the averaged point by the current (last) iterate, hence saving a bit of time and a lot of memory: experimentally, we will still observe a fast convergence.

7 Application: Even/Odd Splitting of the Total Variation

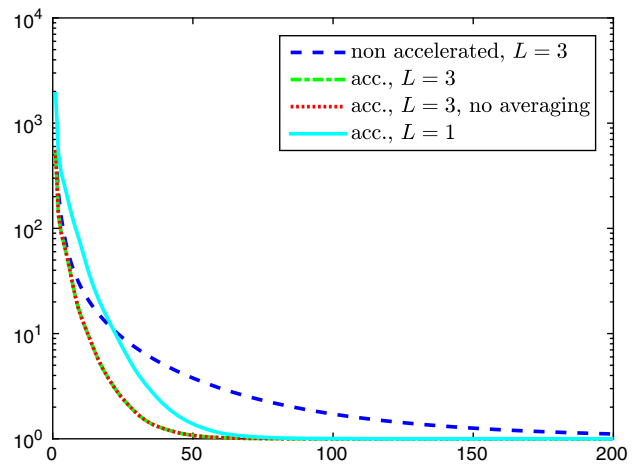
7.1 Description

We now consider the computation of the proximity operator of the total variation, using a splitting proposed in [11]. The idea (which we explain in dimension 2, but could be extended to any dimension) is to consider separately the pixels $(i, j) + \{0, 1\}^2$ for (i, j) even and for (i, j) odd. Given $\mathbf{u} = (u_{i,j})_{1 \leq i \leq n, 1 \leq j \leq m}$ an image, We let for each (i, j)

$$TV_{i,j}^4(\mathbf{u}) = \sqrt{2} \sqrt{(u_{i+1,j} - u_{i,j})^2 + (u_{i+1,j+1} - u_{i,j+1})^2 + (u_{i+1,j+1} - u_{i+1,j})^2 + (u_{i,j+1} - u_{i,j})^2}.$$

Then (here $[\cdot]$ is the integer part)

$$J(\mathbf{u}) = \sum_{i=1}^{[(n-1)/2]} \sum_{j=1}^{[(m-1)/2]} TV_{2i,2j}^4(\mathbf{u}) + \sum_{i=1}^{[n/2]-1} \sum_{j=1}^{[m/2]-1} TV_{2i+1,2j+1}^4(\mathbf{u}). \tag{44}$$



We will denote by $J^e(\mathbf{u})$ the first sum above, and by $J^o(\mathbf{u})$ the second one. It is possible to show that this is an approximation of the isotropic total variation in a variational sense, see ‘Appendix’ for a sketch of proof. Given $\varepsilon > 0$ a smoothing parameter, we will also consider the ‘Huber’ variant $J_\varepsilon(\mathbf{u}) = J_\varepsilon^e(\mathbf{u}) + J_\varepsilon^o(\mathbf{u})$ defined similarly, but replacing $TV_{i,j}^4$ with

$$TV_{i,j}^{4,\varepsilon}(\mathbf{u}) = \begin{cases} TV_{i,j}^4(\mathbf{u}) - \varepsilon & \text{if } TV_{i,j}^4(\mathbf{u}) \geq 2\varepsilon \\ \frac{(TV_{i,j}^4(\mathbf{u}))^2}{4\varepsilon} & \text{else.} \end{cases}$$

We will show how to compute, using the approach described so far, the proximity operator of these functions $J = J_0$ and J_ε , which is defined as the solution of the following problem:

$$\min_{\mathbf{u}} J_\varepsilon(\mathbf{u}) + \frac{1}{2\lambda} \|\mathbf{u} - \mathbf{u}^\dagger\|^2. \tag{45}$$

Given i, j , we denote by $D_{i+1/2,j}\mathbf{u} = u_{i+1,j} - u_{i,j}$ if $1 \leq i \leq n - 1, 1 \leq j \leq m$, and $D_{i,j+1/2}\mathbf{u} = u_{i,j+1} - u_{i,j}$ if $1 \leq i \leq n, 1 \leq j \leq m - 1$. Then, we call $D^o\mathbf{u}$ the ‘odd’ part of $D\mathbf{u}$ and $D^e\mathbf{u}$ the even part, that is

$$D^o\mathbf{u} = ((D_{i+1/2,j}\mathbf{u}, D_{i,j+1/2}\mathbf{u}, D_{i+1/2,j+1}\mathbf{u}, D_{i+1,j+1/2}\mathbf{u}))_{i,j \text{ odd}}$$

and $D^e\mathbf{u}$ is define in the same way but for even indices i, j . It follows that

$$J_\varepsilon^o(\mathbf{u}) = \sup \left\{ \langle \xi, D^o\mathbf{u} \rangle - \frac{\varepsilon}{2} \|\xi\|^2 : \|\langle \xi_{i+1/2,j}, \xi_{i,j+1/2}, \xi_{i+1/2,j+1}, \xi_{i+1,j+1/2} \rangle\|^2 \leq 2 \quad \forall (i, j) \text{ odd} \right\}$$

and the same holds for J^e , replacing D^o with D^e and ‘odd’ with ‘even.’ We will denote

$$\begin{aligned} \xi^o &= ((\xi_{i+1/2,j}, \xi_{i,j+1/2}, \xi_{i+1/2,j+1}, \xi_{i+1,j+1/2}))_{i,j \text{ odd}}, \\ \xi^e &= ((\xi_{i+1/2,j}, \xi_{i,j+1/2}, \xi_{i+1/2,j+1}, \xi_{i+1,j+1/2}))_{i,j \text{ even}}. \end{aligned}$$

The dual of problem (45) reads

$$\min_{(\xi^e, \xi^o)} \|D^{o,*}\xi^o + D^{e,*}\xi^e - \mathbf{u}^\dagger\|^2 + f(\xi^e) + g(\xi^o), \quad (46)$$

where $D^{\bullet,*}$ is the adjoint of D^\bullet ,

$$f(\xi^e) = \begin{cases} \frac{\varepsilon}{2\lambda} \|\xi^e\|^2 & \text{if for all } i, j \text{ even, } \|(\xi_{i+1/2,j}, \xi_{i,j+1/2}, \xi_{i+1/2,j+1}, \xi_{i+1,j+1/2})\|_2^2 \leq 2\lambda^2, \\ +\infty & \text{else} \end{cases}$$

and $g(\xi^o)$ is defined similarly.

We find that (46) is a particular case of (1) (the extra term \mathbf{u}^\dagger in (46) does not change anything to the analysis, and could in fact be transferred to the functions f, g). In that case, A and B have the same norm (which is exactly 2, as these operators can be thought as independent cyclic one-dimensional finite differences over 4 points). Moreover, the functions f, g are (ε/λ) -strongly convex.

7.2 Alternating Minimizations

For this problem, one may to implement an alternating minimization scheme. An approach to do it is detailed in [11] and consists in solving, for each odd or even square, a reduced total variation minimization problem over a cycle of 4 points. This can be done at the expense of a few Newton iterations to find the Lagrange multiplier associated to the constraint on ξ . It follows that one can use Algorithm 1, yielding a $O(1/k^2)$ (for $\varepsilon = 0$) or a linear (for $\varepsilon > 0$) convergence rate. As one has $A^*A \leq 4I$ and $B^*B \leq 4I$, the parameters are $\gamma = \delta = \varepsilon/(4\lambda)$. In particular,

$$q = \frac{\gamma}{1 + \gamma} = \frac{\varepsilon}{4\lambda + \varepsilon}, \quad \gamma' = \frac{\varepsilon}{4\lambda} + \frac{\varepsilon}{4\lambda + \varepsilon},$$

which allow to implement the rules (19), (20), (21).

7.3 Alternating Descent

For alternating descent, one considers metrics $M = I/\tau - A^*A$ and $N = I/\sigma - B^*B$ which are nonnegative as soon as $\tau \leq 1/4, \sigma \leq 1/4$. In the nonstrongly convex case, one could then use Algorithm 2.

On the other hand, if $\varepsilon > 0$, in order to ensure (36) it is enough to have (for $\sigma = \tau = 1/4$)

$$\max \left\{ \frac{4}{K}, \frac{4}{L} + \left(1 - \frac{1}{L}\right) B^*B \right\} \leq t(t-1) \frac{\varepsilon}{\lambda}$$

which is ensured as soon as $t(t-1) \geq 4\lambda/\varepsilon$, hence one should take $t = (1 + \sqrt{1 + 16\lambda/\varepsilon})/2$. The linear convergence should then follow with the rate

$$\omega = 1 - \frac{1}{t} = 1 - \frac{\varepsilon}{8\lambda} \sqrt{1 + 16\frac{\lambda}{\varepsilon}} + \frac{\varepsilon}{8\lambda} \approx 1 - \frac{1}{2} \sqrt{\frac{\varepsilon}{\lambda}}$$

when $\varepsilon \ll \lambda$. In practice, we implemented both the over-relaxation rule with constant steps and the one in Sect. 7.2 (however for both ‘odd’ and ‘even’ variables) and found a very slight advantage for the latter one.

7.4 Experiments

7.4.1 Comparison Between the Algorithms

A first round of experiments simply compares 4 different methods for solving (45) with as input the image in Fig. 2, left:

- The accelerated alternating minimization method (AAMM) of Algorithm 1 where the subproblems are solved almost exactly using an exact inversion with a Lagrange multiplier computed by 4 iterations of a Newton method [11] (which we found was yielding the same result as with more iterations);
- The alternating descent method (AADM) of Algorithms 2–3 where, to simplify, we have used only the points (\hat{x}, \hat{y}) (and not the averages), and we have used the overrelaxation for both updates x, y , as in Algorithms 2 and 3;
- An inexact implementation of Algorithm 1 (AAMM-inexact) where the (almost) exact minimizations of (AAMM) are replaced with a fixed number of descent steps, as in AADM. (The main difference with (AADM) being that the overrelaxation is only implemented on the second variable);
- The ‘FISTA’ method [3] (FISTA) (with parameter updates which take into account the strong convexity of the objective when $\varepsilon > 0$, as explained in [12, 19]). This corresponds to a proximal gradient descent on the

Fig. 2 A
 $360^2 = 129,600$ -pixel image
 (with values in $[0, 255]$) and the
 solution of (45) for $\lambda = 30$,
 $\varepsilon = 1$



(partially smooth, and strongly convex for $\varepsilon > 0$) objective (46), jointly in the variables (ξ^e, ξ^o) , as classically implemented to solve such problems.

The reasons for which we did not use the complete set of variables $(\hat{x}, \hat{y}, \hat{x}', \hat{y}')$ in our implementation of Algorithms 2 and 3 are explained in Sect. 6, where we did not see a different behavior between the original algorithms and the inexact ones where we approximate the points (\hat{x}', \hat{y}') with the nonaveraged corresponding points (\hat{x}, \hat{y}) . Computing the averages would be much more expensive in memory and computational time per iteration—while with this approximation we only need a number of variable of the same order as for the FISTA method, and slightly larger than (AAMM) [(and (AAMM-inexact)] which overrelax only one of the two variables (observe that for a single descent step, our implementation is the correct implementation of the algorithms). In these alternating descent algorithms, we also found the constant step update rule (36) (for $\varepsilon > 0$) slightly worse than the variable rule (19) of Algorithm 1, so in the end we used the latter rule for both methods.

These first experiments were conducted on a Dell Laptop under Ubuntu Linux, with an Intel Core i7-3740QM CPU (6Mb cache) with 4 cores and 8 threads, at 3.70GHz. The programs were implemented in C with `omp` parallelization over 8 threads, which roughly divides the running time by 8 as the operations which are run in parallel are truly independent and take about the same time (they consist in a fixed number of similar operations). For each experiment, our programs were calling the optimization 10 times in a row and we then divided the total elapsed time by 10. There is still some variability which depends on many factors (some which we cannot really control, such as the temperature of the CPU, other easier to understand and deal with such as the total load of the system), we tried to run all the experiments in the same conditions. The results are shown in Table 1. The

number of iterations and time shown are to reach a gap G such that $\sqrt{G/N} \leq 0.1$, where N is the size of the problem (here $N = 129,600$). This implies in particular that the RMSE between the computed solution and the exact one is less than 0.1 (we use a standard a posteriori estimator for this RMSE, see for instance [12, Example 3.1]).

The results are almost as expected. The exact minimization works best, except, strangely, when $\varepsilon = 0$. Computing one step of descent (with a complete overrelaxation in both variables x and y) is quite efficient for this particular problem: even if one needs to perform much more many iterations, these are very fast (in these examples, about 1 versus 1.4 ms for (FISTA) and 1.6 ms for the Newton iterations) which makes the strategy competitive. We recall however that this approach requires more memory. If, as expected, the method (AAMM-inexact) gives terrible results when the number of inner loops is too small (it is improperly overrelaxed in only one of the two variables), for more than 5–6 iteration it starts to compare with the exact (AAMM) method, which means it probably also almost achieves the exact minimization in each variable (consider that the dimension of each subproblem is 4). Surprisingly, for $\varepsilon = 0$ (the nonsmoothed total variation), it converges even faster than the exact minimization approach, and we do not have a reasonable explanation for this.¹

¹ Consider however that as our implementations are here the same c program where depending on our choice either the descent step or the exact minimization is called, this should not be a bug. This is confirmed both by the fact that for $\varepsilon > 0$, when the subproblems are easier and hence it is even more likely that the descent steps will converge to the exact solution in few iterations, the exact and inexact method need nearly the same number of iterations, and the fact that increasing the number of descent steps yield eventually a number of outer iterations equal to the (AAMM) algorithm.

Table 1 Comparison of different strategies

ε	Method	(FISTA)	(AAMM)	(AADM)			(AAMM-inexact)		
				# descent steps			# descent steps		
				1	3	5	1	3	5
0	#iter.	495	146	273	169	153	1654	271	130
	t (ms)	681	232	251	228	281	1115	320	204
0.1	#iter.	142	57	89	62	59	513	100	58
	t (ms)	203	91	89	91	108	333	130	104
1	#iter.	69	27	57	31	28	174	44	29
	t (ms)	101	40	62	50	57	127	54	43

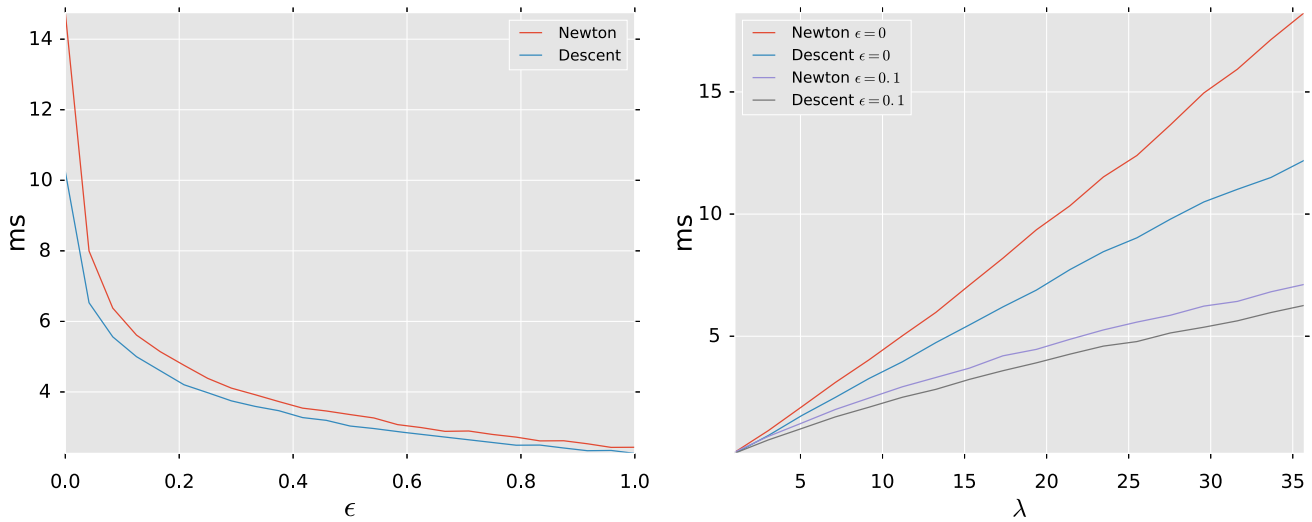


Fig. 3 Left influence of ϵ . Right influence of λ

7.4.2 GPU Implementation

Thanks to the good parallelization properties of the odd/even splitting, it is easy to implement such a scheme on a GPU architecture. The practitioner should download the source code available at <http://github.com/svaiter/ftvp> to test against its image database. This repository contains a C/CUDA library together with a Python 3 binding. All the computation are performed on an Amazon EC2 g2.2xlarge instance on Linux Ubuntu Server 14.04 LTS with CUDA 6.5.

If not specified otherwise, the parameter of all simulations is as follows. We used a standard image of size 512×512 which a dynamic inside the range $[0, 255]$. Our stopping criterion is as before by checking that the square root of the dual over the size of the image is less than 0.1 which is an upper bound of the root mean-square error (RMSE). The dual gap is computed at each iteration. If such a bound is not obtained after 10,000 iterations, we stop the alternating minimization. In term of distributed computing, we choose to use thread blocks of size 16×16 .

The use of Huber-TV induces better performances, in term of execution time or raw number of iterations. We first study the influence of ϵ in Fig. 3. We compare both the case where

the inner iterations are done with a Newton step and with a simple descent, both with 5 steps. For every experience in the following, we consider 20 repetitions of the experiment, and average the time obtained. Moreover, all time benchmarked are reported minus the memory initialization time. We fix the value of $\lambda = 30.0$. Note that choosing ϵ too big is however problematic in term of quality of approximation of the true total variation regularization.

A similar study can be performed for the influence of λ , see Fig. 3. Again, we compare both the case where the inner iterations are done with a Newton step and with a descent, both with 5 steps. We let vary λ over $[1, 36]$ and fix the value of $\epsilon = 0$ (exact-TV) and also $\epsilon = 0.1$. Note that the execution time scales nicely with the dimension of the image. For instance, running our algorithm for $\epsilon = 0.1$ and $\lambda = 20.0$ took 800 ms for a 2048×2048 image and 4s for a 4096×4096 image.

7.4.3 Color TV

For color images, we can implement the same method. The difficulty now is that the ‘exact’ minimization approach of [11] becomes heavier to program and solve, as the sub-

Fig. 4 A
 $3264 \times 2448 = 7,990,272$ -pixel
 image and a
 $360 \times 360 = 129,600$ -pixel
 crop (with RGB values in
 $[0, 255]$), and (below) the
 solutions of (45) for $\lambda = 10$,
 $\varepsilon = 0.1$



problems are now in dimension 12, involving a ‘Laplacian’ matrix of rank 7. Table 1 suggests that performing a sufficient number of descent steps [method (AAMM-inexact)] yields essentially the same results as an exact minimization, in roughly the same time. We thus present the result of such an implementation. We have just extended the program implementing (AAMM-inexact) to work with RGB images, and tested it first on a 360×360 crop and then on the $3264 \times 2448 = 7,990,272$ pixels image of Fig. 4.² The results, shown in Table 2, show that the method is also efficient with this inexact implementation (with 5 descent steps). The left part of the table shows the execution time for the small image, on the same computer as in Table 1. The time spent in each iteration is about 3–4 times longer than for gray-level images.

On the right, we display typical execution times for the large image (of almost 8×10^6 pixels), on a slightly faster computer (with an Intel Xeon E5-2643 CPU (20Mb cache) at 3.40 GHz, which has 12 threads).

8 Conclusion

In this paper, we have studied the acceleration of alternating minimization or descent schemes for problems with two

Table 2 Color results

ε	λ	Small image			Large image		
		1	10	50	1	10	50
0	#iter.	21	50	280	17	41	258
	t (s)	0.165	0.243	1.050	3.5	6.3	32.0
0.1	#iter.	15	31	134	13	26	116
	t (s)	0.154	0.210	0.526	3.0	4.5	14.7
1	#iter.	8	14	43	8	14	44
	t (s)	0.124	0.150	0.206	2.4	3.1	6.1

variables with a quadratic coupling, as already considered in [11]. We have extended some of these results to strongly convex problems and have investigated the case of partial descent steps, showing that (theoretically) acceleration is also possible in this setting. A natural development would be to analyze better the behavior of the inexact variant, which we use in practice and which seems to be quite efficient in our application. The correct framework for this analysis should probably be the framework of inexact accelerated schemes, as studied in [1, 21]; however, for this we would need to better estimate the errors which are introduced by the method (AAMM-inexact) and which seem much smaller than one could naturally expect.

² Image belongs to the authors.

Acknowledgements This work is supported by the ANR via the international project ‘EANOI’ (Efficient Algorithms for Nonsmooth Optimization in Imaging), FWF No. I1148 / ANR-12-IS01-0003. A. Chambolle also benefits from support of the ‘Programme Gaspard Monge pour l’Optimisation et la Recherche Opérationnelle’ (PGMO), through the ‘MAORI’ group, as well as the ‘GdR MIA’ of the CNRS. He also warmly thanks Churchill College and DAMTP, Centre for Mathematical Sciences, University of Cambridge, for their kind hospitality during the completion of this work, thanks to a support of the French Embassy in the UK and the Cantab Capital Institute for Mathematics of Information.

Appendix: An Approximation Result

In this appendix, we show that although this is not totally obvious at first glance, the discrete energy $J_\varepsilon(\mathbf{u})$ is an approximation of the isotropic total variation. The result is more precisely as follows. To simplify we work in the domain $\Omega = (0, 1)^2$ (extension to more general regular domains is not difficult) and we define, for $N \geq 1$ an integer, the functional, defined for $u \in L^1(\Omega)$,

$$J_{\varepsilon,N}(u) = \begin{cases} \frac{1}{N} J_{\varepsilon/N}^{N,N}(\mathbf{u}) & \text{if } \mathbf{u} = (u_{i,j})_{1 \leq i,j \leq N}, u(x) = \sum_{i=1}^N \sum_{j=1}^N u_{i,j} \chi_{(\frac{i-1}{N}, \frac{j}{N}) \times (\frac{j-1}{N}, \frac{j}{N})}(x) \text{ a.e.,} \\ +\infty & \text{else.} \end{cases}$$

here, $J_{\varepsilon/N}^{N,N}$ is a notation for the energy (44) in case $m = n = N$ (and with the smoothing parameter ε/N). We also denote $\Phi_\varepsilon(p) := |p|^2/(2\varepsilon)$ if $|p| \leq \varepsilon$, $|p| - \varepsilon/2$ else and recall that for $u \in BV(\Omega)$ a function with bounded variation $|Du|(\Omega) < +\infty$ [17,22], $\int_\Omega \Phi_\varepsilon(Du) = \int_\Omega \Phi_\varepsilon(\nabla u) dx + |D^s u|$ where $Du = \nabla u dx + D^s u$ is the Radon-Nikodym decomposition of Du as an absolutely continuous and singular part, see [15]. We introduce the functional

$$J_\varepsilon(u) = \begin{cases} \Phi_\varepsilon(Du)(\Omega) & \text{if } u \in BV(\Omega), \\ +\infty & \text{if } u \in L^1(\Omega) \setminus BV(\Omega). \end{cases}$$

Then, one can show that J_ε can also be defined by duality, as follows:

$$J_\varepsilon(u) = \sup \left\{ \int_\Omega u(x) \operatorname{div} \varphi(x) dx - \frac{\varepsilon}{2} \int_\Omega |\varphi(x)|^2 dx : \varphi \in C_c^\infty(\Omega; \mathbb{R}^2), |\varphi(x)| \leq 1 \quad \forall x \in \Omega \right\} \quad (47)$$

One has the following result:

Theorem 4 *As $N \rightarrow \infty$, $J_{\varepsilon,N}$ Γ -converges to J_ε . Moreover, if for some sequence $(u^N) \in L^1(\Omega)^N$, $J_{\varepsilon,N}(u^N) \leq C < +\infty$, then there exists $u \in BV(\Omega)$, a subsequence $(u^{N_k})_k$ and a sequence of constants $(a_k)_k$ such that $u^{N_k} - a_k \rightarrow u$ in $L^1(\Omega)$.*

For the proper definition and main properties of Γ -convergence, see for instance [6, 14]. The theorem establishes that images minimizing $J_{\varepsilon/N}^{N,N}$ (+ other terms such as a quadratic penalization) should be close if N is large to minimizers of the isotropic ‘Huber-total variation’ J_ε , in the continuum. The proof is easy, however not really found in this form in the literature, as far as we know. The closest results are maybe the Γ -convergence theorems of Cai et al. [7] in the context of wavelet-based approximations of the total variation.

Proof It is enough to prove: (i) that if $u^N \in L^1(\Omega)$ is such that $\ell = \liminf_N J_{\varepsilon,N}(u^N) < \infty$, then not only one can extract u^{N_k} which converges to some u , but in addition $J_\varepsilon(u) \leq \ell$; (ii) that given u with finite total variation, one can build a sequence u^N with $\limsup_N J_{\varepsilon,N}(u^N) \leq J_\varepsilon(u)$.

For point (i), we first consider a subsequence (u^{N_k}) such that $\ell = \lim_k J_{\varepsilon,N_k}(u^{N_k})$. Then, we see that since for all k (large enough) $J_{\varepsilon,N_k}(u^{N_k}) < +\infty$, by definition u^{N_k} is piecewise constant and can be written

$$u^{N_k}(x) = \sum_{i=1}^{N_k} \sum_{j=1}^{N_k} u_{i,j}^k \chi_{(\frac{i-1}{N_k}, \frac{j}{N_k}) \times (\frac{j-1}{N_k}, \frac{j}{N_k})}(x)$$

for some matrix $\mathbf{u}^k = (u_{i,j}^k)_{1 \leq i,j \leq N_k}$. Then we observe that for some constant $\sigma > 0$,

$$\begin{aligned} J_{\varepsilon,N_k}(u^{N_k}) + \frac{\varepsilon}{2} &\geq J_{0,N_k}(u^{N_k}) \\ &\geq \sigma \frac{1}{N_k} \sum_{i,j} (|u_{i+1,j}^k - u_{i,j}^k| + |u_{i,j+1}^k - u_{i,j}^k|) \\ &= \sigma |Du^{N_k}|(\Omega). \end{aligned}$$

Hence $|Du^{N_k}|(\Omega)$ is bounded, showing that $(u^{N_k} - a_k)_k$ is precompact in $L^1(\Omega)$, where a_k is the average of the function u^{N_k} in Ω . Without loss of generality, we assume $a_k = 0$ and we denote by u the limit of a subsequence (which for convenience we do not relabel). We must now show that $J_\varepsilon(u) \leq \ell$.

Let $\delta > 0$, and let $\varphi = (\varphi^1, \varphi^2) \in C_c^\infty(\Omega; \mathbb{R}^2)$ be a smooth vector field with $|\varphi(x)|^2 = \varphi^1(x)^2 + \varphi^2(x)^2 \leq 1 - \delta$ for all $x \in \Omega$. Observe that

$$\begin{aligned} \int_\Omega u^{N_k}(x) \operatorname{div} \varphi(x) dx &= \sum_{i,j} u_{i,j}^k \int_{(\frac{i-1}{N_k}, \frac{j}{N_k}) \times (\frac{j-1}{N_k}, \frac{j}{N_k})} \operatorname{div} \varphi(x) dx \\ &= \sum_{i,j} (u_{i+1,j}^k - u_{i,j}^k) \varphi_{i+\frac{1}{2},j}^1 \end{aligned}$$

$$+ (u_{i,j+1}^k - u_{i,j}^k) \varphi_{i,j+\frac{1}{2}}^2$$

where $\varphi_{i+\frac{1}{2},j}^1$ is the flux of φ through the vertical segment $\{\frac{i}{N_k}\} \times (\frac{j-1}{N_k}, \frac{j}{N_k})$ and $\varphi_{i,j+\frac{1}{2}}^2$ is the flux through the horizontal segment $(\frac{i-1}{N_k}, \frac{i}{N_k}) \times \{\frac{j}{N_k}\}$.

Assume (i, j) are both odd or even. Denote by $\bar{x} = (i/N_k, j/N_k)$: as φ is smooth, one clearly has that $N_k \varphi_{i+\frac{1}{2},j}^1 = \varphi^1(\bar{x}) + O(1/N_k)$, etc., and, in fact,

$$\begin{aligned} N_k^2 \mathcal{N}_{i,j}^2 &:= (N_k \varphi_{i+\frac{1}{2},j}^1)^2 + (N_k \varphi_{i,j+\frac{1}{2}}^2)^2 \\ &\quad + (N_k \varphi_{i+\frac{1}{2},j+1}^1)^2 + (N_k \varphi_{i+1,j+\frac{1}{2}}^2)^2 \\ &\leq 2(1 - \delta) + O\left(\frac{1}{N_k^2}\right) \leq 2 \end{aligned}$$

if N_k is large enough. As a consequence

$$\begin{aligned} &(u_{i+1,j}^k - u_{i,j}^k) \varphi_{i+\frac{1}{2},j}^1 \\ &\quad + (u_{i,j+1}^k - u_{i,j}^k) \varphi_{i,j+\frac{1}{2}}^2 \\ &\quad + (u_{i+1,j+1}^k - u_{i,j+1}^k) \varphi_{i+\frac{1}{2},j+1}^1 \\ &\quad + (u_{i+1,j+1}^k - u_{i+1,j}^k) \varphi_{i+1,j+\frac{1}{2}}^2 - \frac{\varepsilon}{2} \mathcal{N}_{i,j}^2 \\ &\leq \frac{1}{N_k} T V_{i,j}^{4,\varepsilon/N_k}(\mathbf{u}^k). \end{aligned} \tag{48}$$

Thanks to the smoothness of φ , one can check easily that

$$\sum_{(i,j) \text{ even}} \mathcal{N}_{i,j}^2 + \sum_{(i,j) \text{ odd}} \mathcal{N}_{i,j}^2 \rightarrow \int_{\Omega} |\varphi(x)|^2 dx$$

as $k \rightarrow \infty$, hence, summing (48) over all (i, j) both odd or both even, we find (using also the fact that φ has compact support) that

$$\begin{aligned} &\int_{\Omega} u^{N_k}(x) \operatorname{div} \varphi(x) dx - \frac{\varepsilon}{2} \int_{\Omega} |\varphi(x)|^2 dx \\ &\quad + o(1) \leq \frac{1}{N_k} J_{\varepsilon/N_k}^{N_k, N_k}(\mathbf{u}^k) = \mathcal{F}_{\varepsilon, N_k}(u^{N_k}). \end{aligned}$$

In the limit, we find that

$$\int_{\Omega} u(x) \operatorname{div} \varphi(x) dx - \frac{\varepsilon}{2} \int_{\Omega} |\varphi(x)|^2 dx \leq \ell.$$

Thanks to (47), we deduce that $\mathcal{F}_{\varepsilon}(u) \leq \ell$.

We now must prove (ii). We only sketch the proof, which is very simple: one first observes that as any $u \in BV(\Omega)$ can be approximated by a sequence (u_n) with $u_n \in C^{\infty}(\bar{\Omega})$, $u_n \rightarrow u$ in $L^1(\Omega)$ and $\int_{\Omega} \Phi_{\varepsilon}(\nabla u_n(x)) dx = \mathcal{F}_{\varepsilon}(u)$, it is

enough to show the result for a smooth function and use then a diagonal argument.

But if u is smooth, letting simply for each N , $u_{i,j}^N = u((i-1/2)/N, (j-1/2)/N)$, one first observes that

$$u^N(x) := \sum_{i,j} u_{i,j}^N \chi_{(\frac{i-1}{N}, \frac{j}{N}) \times (\frac{j-1}{N}, \frac{j}{N})}(x) \rightarrow u(x)$$

uniformly in Ω , and then that $\mathcal{F}_{\varepsilon, N}(u^N)$ is a finite-difference approximation of $\int_{\Omega} \Phi_{\varepsilon}(u(x)) dx$, which converges to this limit as $N \rightarrow \infty$. \square

References

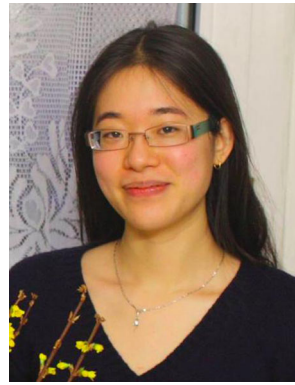
1. Aujol, J.-F., Dossal, C.: Stability of over-relaxations for the forward-backward algorithm, application to FISTA. *SIAM J. Optim.* **25**(4), 2408–2433 (2015)
2. Beck, A.: On the convergence of alternating minimization for convex programming with applications to iteratively reweighted least squares and decomposition schemes. *SIAM J. Optim.* **25**(1), 185–209 (2015)
3. Beck, A., Teboulle, M.: A fast iterative shrinkage–thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.* **2**(1), 183–202 (2009)
4. Beck, A., Tetrushvili, L.: On the convergence of block coordinate descent type methods. *SIAM J. Optim.* **23**(4), 2037–2060 (2013)
5. Boyle, J.P., Dykstra, R.L.: A method for finding projections onto the intersection of convex sets in Hilbert spaces. In: Dykstra, R., Robertson, T., Wright, F.T. (eds) *Advances in Order Restricted Statistical Inference* (Iowa City, Iowa, 1985), vol. 37 of *Lecture Notes in Statistics*, pp. 28–47. Springer, Berlin (1986)
6. Braides, A.: *Gamma-Convergence for Beginners*. Number 22 in *Oxford Lecture Series in Mathematics and Its Applications*. Oxford University Press, Oxford (2002)
7. Cai, J.-F., Dong, B., Osher, S., Shen, Z.: Image restoration: total variation, wavelet frames, and beyond. *J. Am. Math. Soc.* **25**(4), 1033–1089 (2012)
8. Chambolle, A., Dossal, C.: On the convergence of the iterates of the “fast iterative shrinkage/thresholding algorithm”. *J. Optim. Theory Appl.* **166**(3), 968–982 (2015)
9. Chambolle, A., Pock, T.: A first-order primal-dual algorithm for convex problems with applications to imaging. *J. Math. Imaging Vis.* **40**(1), 120–145 (2011)
10. Chambolle, A., Pock, T.: On the ergodic convergence rates of a first-order primal-dual algorithm. *Math. Program.* **159**, 253–287 (2016)
11. Chambolle, A., Pock, T.: A remark on accelerated block coordinate descent for computing the proximity operators of a sum of convex functions. *SMAI J. Comput. Math.* **1**, 29–54 (2015)
12. Chambolle, A., Pock, T.: An introduction to continuous optimization for imaging. *Acta Numer.* **25**, 161–319, 5 (2016)
13. Combettes, P.L., Pesquet, J.-C.: Proximal splitting methods in signal processing. In: Bauschke, H.H., Burachik, R., Combettes, P.L., Elser, V., Luke, D.R., Wolkowicz, H. (eds.) *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, vol. 49 of *Springer Optimization and Applications*, pp. 185–212. Springer, New York (2011)
14. Dal Maso, G.: *An Introduction to Γ -Convergence*. Birkhäuser, Boston (1993)
15. Demengel, F., Temam, R.: Convex functions of a measure and applications. *Indiana Univ. Math. J.* **33**(5), 673–709 (1984)

16. Deutsch, F., Hundal, H.: The rate of convergence of Dykstra's cyclic projections algorithm: the polyhedral case. *Numer. Funct. Anal. Optim.* **15**(5–6), 537–565 (1994)
17. Evans, L.C., Gariepy, R.F.: *Measure Theory and Fine Properties of Functions*. CRC Press, Boca Raton (1992)
18. Nemirovski, A.S., Yudin, D.: Informational complexity of mathematical programming. *Izv. Akad. Nauk SSSR Tekhn. Kibernet.* **1**, 88–117 (1983)
19. Nesterov, Y.: *Introductory Lectures on Convex Optimization: A Basic Course*, vol. 87 of Applied Optimization. Kluwer, Boston (2004)
20. Shefi, R., Teboulle, M.: On the rate of convergence of the proximal alternating linearized minimization algorithm for convex problems. *EURO J. Comput. Optim.* **4**(1), 27–46 (2016)
21. Villa, S., Salzo, S., Baldassarre, L., Verri, A.: Accelerated and inexact forward–backward algorithms. *SIAM J. Optim.* **23**(3), 1607–1633 (2013)
22. Ziemer, W.P.: *Weakly Differentiable Functions. Sobolev Spaces and Functions of Bounded Variation*. Springer, New York (1989)



Antonin Chambolle has studied at Ecole Normale Supérieure in Paris and obtained a Ph.D. in 1993 in Applied Mathematics with Jean-Michel Morel at Université Paris-Dauphine. He has then worked as a CNRS researcher, a postdoc (in SISSA, Trieste, Italy), and is a currently a CNRS research director in Applied Mathematics at Ecole Polytechnique, Palaiseau. His research interests focus on calculus of variations and optimization for free boundary problems which arise in mathematics, mechanics or image processing.

blems which arise in mathematics, mechanics or image processing.



Pauline Tan has been a student of the Ecole Normale Supérieure de Cachan, France. She has obtained a Ph.D. in Applied Mathematics in 2016, from Ecole Polytechnique in Paris, where she was working under the supervision of Antonin Chambolle and Pascal Monasse. She is currently a postdoc at ONERA (The French Aerospace Lab). Her interests are in Applied Mathematics for imaging and image analysis.



Samuel Vaïter has studied Applied Mathematics and Theoretical Computer Science in Lyon and Paris. He has obtained in 2014 a Ph.D. in Applied Mathematics from Université Paris-Dauphine, where he was a student of Gabriel Peyré. He has then worked as a postdoc in CMAP, Ecole Polytechnique, Palaiseau, and has now a CNRS research position at the Institut de Mathématiques de Bourgogne in Dijon, France. His current research interests focus on varia-

tional regularization in signal and image processing, convex analysis, sparsity and risk estimation.