## AN ABSTRACT CONVERGENCE FRAMEWORK WITH APPLICATION TO INERTIAL INEXACT FORWARD–BACKWARD METHODS

SILVIA BONETTINI<sup>\*</sup>, PETER OCHS<sup>†</sup>, MARCO PRATO<sup>†</sup>, AND SIMONE REBEGOLDI<sup>‡</sup>

Abstract. In this paper we introduce a novel abstract descent scheme suited for the minimization of proper and lower semicontinuous functions. The proposed abstract scheme generalizes a set of properties that are crucial for the convergence of several first-order methods designed for nonsmooth nonconvex optimization problems. Such properties guarantee the convergence of the full sequence of iterates to a stationary point, if the objective function satisfies the Kurdyka–Lojasiewicz property. The abstract framework allows for the design of new algorithms. We propose two inertial-type algorithms with (implementable) inexactness criteria for the main iteration update step. The first algorithm, i<sup>2</sup>Piano, exploits large steps by adjusting a local Lipschitz constant. The second algorithm, iPila, overcomes the main drawback of line-search based methods by enforcing a descent only on a merit function instead of the objective function, which even allows for the escape of local minimizers. Both algorithms are proved to enjoy the full convergence guarantees of the abstract descent scheme. The efficiency of the proposed algorithms is demonstrated on an exemplary image deblurring problem in presence of data corrupted by impulse noise, where we can appreciate the benefits of performing a linesearch along the descent direction inside an inertial scheme.

Key words. forward-backward methods, inertial methods, linesearch, nonconvex optimization

AMS subject classifications. 65K05, 90C30

**1.** Introduction. The design of efficient first-order descent methods is vital for tackling composite optimization problems of the form

(1.1) 
$$\min_{x \in \mathbb{R}^n} f(x), \quad f(x) = f_0(x) + f_1(x),$$

where  $f_1$  is convex and  $f_0$  is continuously differentiable on an open set containing the domain of  $f_1$ . Such problems are frequently encountered in image processing and machine learning applications [4, 18, 19], where one of the two terms is usually a data fidelity term and the other one encodes some apriori information on the ground truth [4]. Popular and effective first-order methods aimed at solving (1.1) include forward-backward (FB) methods [25, 24, 31, 12], whose structure consists in the alternation of a gradient step on  $f_0$  followed by a proximal minimization step on  $f_1$ , block coordinate methods [9, 14, 23, 26], Douglas-Rachford methods [25, 28] and several others.

In recent years, the convergence of first-order descent methods in nonconvex settings has been carefully addressed by relying on the so-called Kurdyka–Lojasiewicz (KL) inequality [2, 8, 27]. This analytical property is satisfied by a large number of objective functions arising in signal processing and machine learning, such as real analytic or semialgebraic functions (see e.g. [8, 6]), thus making quite natural to consider the KL inequality as a standard blanket assumption whenever the objective function is nonconvex. Combining the KL inequality with some crucial properties of descent methods allows to prove the convergence of the iterates to a stationary point of the objective function, provided that the sequence is bounded. The convergence of descent methods under the KL assumption was first considered in [1], where the authors prove the convergence of linesearch and trust-region methods for real analytic objective functions. The key idea in [1] is to combine the KL inequality with some strong descent conditions holding for the iterates of classical gradient methods. In [2, 9, 26, 10], the authors extend this seminal idea by providing the first abstract descent schemes in the KL framework, namely a set of abstract properties ensuring the convergence of a generic iterative scheme to a stationary point if combined with the KL inequality. Such properties include a sufficient decrease condition on the function values, a relative error condition on the norm of a subgradient at the current iterate, and a continuity condition of the iterates with respect to the objective function. In [31, 30], the authors modify the abstract scheme proposed in [2] in order to include an inertial term inside a classical FB splitting scheme and devise the so-called iPiano (inertial Proximal algorithm for nonconvex optimization). Indeed, as

<sup>\*</sup>Dipartimento di Scienze Fisiche, Informatiche e Matematiche, Università di Modena e Reggio Emilia, Via Campi 213/b, 41125 Modena, Italy (silvia.bonettini@unimore.it, marco.prato@unimore.it).

<sup>&</sup>lt;sup>†</sup>Department of Mathematics, University of Tübingen, Auf der Morgenstelle 10, 72076 Tübingen, Germany (ochs@math.uni-tuebingen.de).

<sup>&</sup>lt;sup>‡</sup>Dipartimento di Ingegneria Industriale, Università di Firenze, Via di S. Marta 3, 50139 Firenze, Italy (simone.rebegoldi@unifi.it).

a generalization of the Heavy-Ball method [33, 32], iPiano is not a monotone algorithm, and the sufficient decrease condition cannot be directly imposed on the objective function; however, one can define an appropriate surrogate function  $\mathcal{F}$  (often called Lyapunov function) in such a way that the decrease condition holds for  $\mathcal{F}$  in place of the objective function. Similarly, the authors in [15, 16] define a surrogate function in order to adapt the abstract descent scheme to their proposed method VMILA (Variable Metric Inexact Linesearch Algorithm); yet, unlike in [31], the surrogate function is introduced in order to include an implementable inexactness criterion for the computation of the proximal point, and thus avoid the actual implementation of the relative error condition, which seems rather difficult to impose in practice [15, 29]. Inexactness in FB methods may arise when the proximal operator of the function  $f_1$  cannot be computed in closed form or the exact proximal point is too costly to compute [3, 12, 17, 37].

In this paper, we propose a novel abstract descent scheme for proving the convergence of iterative methods under the KL assumption. The proposed scheme defines two separate surrogate functions, in order to treat separately the possible non-monotonicity of the algorithm and the inexact computation of the iterate, and then imposes a set of abstract properties holding for the surrogate functions evaluated at the iterates. Our approach can be considered as an extension of the abstract scheme in [30], where the inexactness of the iterates is treated by employing the same implementation mechanism used in [16] for VMILA. We show that the iterates generated by our abstract scheme converge to a stationary point of one of the surrogate functions, provided that this function satisfies the KL inequality on its domain. Furthermore, we devise two novel inertial FB algorithms, which are encompassed by the proposed framework. The first one is denominated i<sup>2</sup>Piano (inertial inexact Proximal algorithm for nonconvex optimization) and can be considered as an inexact version of the iPiano algorithm equipped with a backtracking procedure based on a local version of the Descent Lemma. The second one is denoted iPila (inertial Proximal inexact line-search algorithm), which features an inertial-like step followed by a linesearch procedure along the descent direction of a suitable merit function. The main advantage of iPila resides precisely in its linesearch strategy, which allows to compute the inexact proximal point only once per iteration, unlike the backtracking procedure of i<sup>2</sup>Piano.

The paper is organized as follows. In Section 2 some basic notions on variational analysis and the definition of the KL property are reported. In Section 3 the proposed abstract scheme is presented and its convergence properties analysed. Section 4 is devoted to the design of the two algorithms i<sup>2</sup>Piano and iPila and their inclusion in the abstract framework presented in Section 3. Finally, a numerical illustration on an image deblurring problem in presence of impulse noise is reported in Section 5.

**2. Preliminaries.** In the remainder of the paper, we denote with  $\overline{\mathbb{R}} = \mathbb{R} \cup \{-\infty, +\infty\}$  the extended real numbers set and  $\mathbb{R}^{n \times n}$  the set of  $n \times n$  real-valued matrices, while  $\|\cdot\|$  denotes the Euclidean norm. Given a function  $\mathcal{F} : \mathbb{R}^n \to \overline{\mathbb{R}}$  and denoting with dom $(\mathcal{F}) = \{x \in \mathbb{R}^n : \mathcal{F}(x) < +\infty\}$  the domain of  $\mathcal{F}$ , we say that  $\mathcal{F}$  is proper if dom $(\mathcal{F}) \neq \emptyset$  and  $\mathcal{F}$  is finite on dom $(\mathcal{F})$ . The distance operator of a point  $x \in \mathbb{R}^n$  to a set  $\Omega \subset \mathbb{R}^n$  is defined as

$$\operatorname{dist}(x,\Omega) = \inf_{y \in \Omega} \|x - y\|.$$

Observe that, if  $\Omega = \Omega_1 \times \Omega_2$ , where  $\Omega_1 \subset \mathbb{R}^{n_1}$ ,  $\Omega_2 \subset \mathbb{R}^{n_2}$ , with  $n_1 + n_2 = n$ , then for all  $x = (x_1, x_2) \in \mathbb{R}^n$ , with  $x_1 \in \mathbb{R}^{n_1}$ ,  $x_2 \in \mathbb{R}^{n_2}$ , we have

(2.1) 
$$\operatorname{dist}(x,\Omega) = \sqrt{\operatorname{dist}(x_1,\Omega_1)^2 + \operatorname{dist}(x_2,\Omega_2)^2}.$$

This follows by observing that  $||x - z||^2 = ||x_1 - z_1||^2 + ||x_2 - z_2||^2$  for all  $z = (z_1, z_2) \in \Omega$ .

DEFINITION 2.1. Given a proper, lower semicontinuous function  $\mathcal{F} : \mathbb{R}^n \to \overline{\mathbb{R}}$ , the Fréchet subdifferential of  $\mathcal{F}$  at  $\overline{z} \in \operatorname{dom}(\mathcal{F})$  is defined as the set [35, Definition 8.3(a)]

$$\hat{\partial}\mathcal{F}(\overline{z}) = \left\{ w \in \mathbb{R}^n : \liminf_{u \to \overline{z}, u \neq \overline{z}} \frac{\mathcal{F}(u) - \mathcal{F}(\overline{z}) - (u - \overline{z})^T w}{\|u - \overline{z}\|} \ge 0 \right\}$$

Furthermore, the limiting subdifferential of  $\mathcal{F}$  at  $\overline{z}$  is given by [35, Definition 8.3(b)]

$$\partial \mathcal{F}(\overline{z}) = \{ w \in \mathbb{R}^n : \exists \ z^{(k)} \to \overline{z}, \ \mathcal{F}(z^{(k)}) \to \mathcal{F}(\overline{z}), \ w^{(k)} \in \hat{\partial} \mathcal{F}(z^{(k)}) \to w \ as \ k \to \infty \}$$

DEFINITION 2.2. Let  $\mathcal{F} : \mathbb{R}^n \to \overline{\mathbb{R}}$  be a proper, lower semicontinuous function. A point  $\overline{z} \in \mathbb{R}^n$  is stationary for  $\mathcal{F}$  if  $0 \in \partial \mathcal{F}(\overline{z})$ .

Let us introduce the lazy slope of  $\mathcal{F}$  at  $\overline{z}$ , which is given by [26, p. 877]

$$\|\partial \mathcal{F}(\overline{z})\|_{-} = \inf_{v \in \partial \mathcal{F}(\overline{z})} \|v\|.$$

It is then easy to prove the following sufficient criterion for establishing if a point  $\overline{z} \in \mathbb{R}^n$  is stationary for the function  $\mathcal{F}$ .

LEMMA 2.3. [26, Lemma 2.1] Let  $\mathcal{F} : \mathbb{R}^n \to \overline{\mathbb{R}}$  be a proper, lower semicontinuous function and  $\overline{z} \in \mathbb{R}^n$ . If there exists  $\{z^{(k)}\}_{k\in\mathbb{N}} \subset \mathbb{R}^n$  such that  $z^{(k)} \to \overline{z}$ ,  $\mathcal{F}(z^{(k)}) \to \mathcal{F}(\overline{z})$  and  $\liminf_{k\to\infty} \|\partial \mathcal{F}(z^{(k)})\|_{-} = 0$ , then  $0 \in \partial \mathcal{F}(\overline{z})$ .

DEFINITION 2.4. Let  $\mathcal{F} : \mathbb{R}^n \longrightarrow \overline{\mathbb{R}}$  be a proper, lower semicontinuous function. The function  $\mathcal{F}$  is said to have the KL property at  $\overline{z} \in \operatorname{dom}(\partial \mathcal{F})$  if there exist  $v \in (0, +\infty]$ , a neighborhood U of  $\overline{z}$ , a continuous concave function  $\phi : [0, v) \longrightarrow [0, +\infty)$  with  $\phi(0) = 0$ ,  $\phi \in C^1(0, v)$ ,  $\phi'(s) > 0$  for all  $s \in (0, v)$ , such that the following inequality is satisfied

$$\phi'(\mathcal{F}(z) - \mathcal{F}(\overline{z})) \|\partial \mathcal{F}(z)\|_{-} \ge 1$$

for all  $z \in U \cap \{z \in \mathbb{R}^n : \mathcal{F}(\overline{z}) < \mathcal{F}(z) < \mathcal{F}(\overline{z}) + v\}$ . If  $\mathcal{F}$  satisfies the KL property at each point of dom $(\partial \mathcal{F})$ , then  $\mathcal{F}$  is called a KL function.

The function  $\phi$  in the previous definition is called *desingularization function* and it depends on the point  $\bar{z}$ . In [9, Lemma 6], the following *uniformized* version of the KL property is introduced, where the KL inequality holds with the same desingularization function for all points in a suitable neighborhood of a compact set where the function is constant.

LEMMA 2.5. Let  $\mathcal{F} : \mathbb{R}^n \to \overline{\mathbb{R}}$  be a proper, lower semicontinuous function and  $X \subset \mathbb{R}^n$  a compact set. Suppose that  $\mathcal{F}$  satisfies the KL property at each point belonging to X and that  $\mathcal{F}$  is constant over X, i.e.,  $\mathcal{F}(\bar{x}) = \overline{\mathcal{F}} \in \mathbb{R}$  for all  $\bar{x} \in X$ . Then, there exists  $\mu, v > 0$  and a function  $\phi$  as in Definition 2.4 such that

(2.2) 
$$\phi'(\mathcal{F}(z) - \bar{\mathcal{F}}) \|\partial \mathcal{F}(z)\|_{-} \ge 1, \quad \forall z \in \bar{B}$$

where the set  $\overline{B}$  is defined as

(2.3) 
$$\overline{B} = \{ z \in \mathbb{R}^n : \operatorname{dist}(z, X) < \mu \text{ and } \overline{\mathcal{F}} < \mathcal{F}(z) < \overline{\mathcal{F}} + v \}.$$

3. Abstract algorithm scheme. In the following, we are interested in proving the convergence of an abstract descent algorithm to a stationary point of a proper, lower semicontinuous function  $\mathcal{F}$ . Such abstract algorithm is defined through a specific set of properties that are shared by several first-order methods designed for nonsmooth nonconvex optimization, including gradient descent methods [1], forward-backward methods [12, 22, 31] and block coordinate methods [2, 14, 26]. Similarly to other abstract descent algorithms in the KL framework, the two main ingredients guaranteeing the convergence of our scheme are the *sufficient decrease condition* and the *relative error condition*, the latter being related to the minimization subproblem that one has to (inexactly) solve at each iteration of a first-order method. However, unlike in previous works in the literature, we require that the relative error condition is satisfied at a point that might be different from the actual iterate generated by the method. As we will see in Section 4, this simple modification allows to circumvent the issue of the actual implementation of the relative error condition, allowing to include inexact forward-backward methods equipped with an implementable inexactness criterion for the solution of the minimization subproblem.

CONDITIONS 3.1 (Abstract algorithm scheme). Let  $\mathcal{F} : \mathbb{R}^n \times \mathbb{R}^m \to \overline{\mathbb{R}}$  be a proper, lower semicontinuous function and  $\Phi : \mathbb{R}^n \times \mathbb{R}^q \to \overline{\mathbb{R}}$  a proper, lower semicontinuous, bounded from below function. Consider two sequences  $\{x^{(k)}\}_{k\in\mathbb{N}}$ ,  $\{u^{(k)}\}_{k\in\mathbb{N}}$  in  $\mathbb{R}^n$ , a sequence  $\{\rho^{(k)}\}_{k\in\mathbb{N}}$  in  $\mathbb{R}^m$ , a sequence  $\{s^{(k)}\}_{k\in\mathbb{N}}$  in  $\mathbb{R}^q$  and a sequence of nonnegative real numbers  $\{d_k\}_{k\in\mathbb{N}}$  such that the following relations are satisfied. [H1] There exists a sequence of positive real numbers  $\{a_k\}_{k\in\mathbb{N}}$  such that

$$\Phi(x^{(k+1)}, s^{(k+1)}) + a_k d_k^2 \le \Phi(x^{(k)}, s^{(k)}), \quad \forall \ k \ge 0.$$

[H2] There exists a sequence of nonnegative real numbers  $\{r_k\}_{k\in\mathbb{N}}$  with  $\lim_{k\to\infty} r_k = 0$  such that

$$\Phi(x^{(k+1)}, s^{(k+1)}) \le \mathcal{F}(u^{(k)}, \rho^{(k)}) \le \Phi(x^{(k)}, s^{(k)}) + r_k, \quad \forall \ k \ge 0.$$

[H3] There exist b > 0, a sequence of positive real numbers  $\{b_k\}_{k \in \mathbb{N}}$ , a summable sequence of nonnegative real numbers  $\{\zeta_k\}_{k \in \mathbb{N}}$ , a non-empty finite index set  $I \subset \mathbb{Z}$  and  $\theta_i \ge 0$ ,  $i \in I$  with  $\sum_{i \in I} \theta_i = 1$  such that, setting  $d_j = 0$  for  $j \le 0$ , we have

$$b_{k+1} \| \partial \mathcal{F}(u^{(k)}, \rho^{(k)}) \|_{-} \le b \sum_{i \in I} \theta_i d_{k+1-i} + \zeta_{k+1}, \quad \forall \ k \ge 0.$$

[H4] If  $\{(x^{(k_j)}, \rho^{(k_j)})\}_{j \in \mathbb{N}}$  is a subsequence of  $\{(x^{(k)}, \rho^{(k)})\}_{k \in \mathbb{N}}$  converging to some  $(x^*, \rho^*) \in \mathbb{R}^n \times \mathbb{R}^m$ , then we have for  $\{u^{(k_j)}\}_{j \in \mathbb{N}}$ :

$$\lim_{j \to \infty} \|u^{(k_j)} - x^{(k_j)}\| = 0, \quad \lim_{j \to \infty} \mathcal{F}(u^{(k_j)}, \rho^{(k_j)}) = \mathcal{F}(x^*, \rho^*).$$

[H5] There exists a positive real number p > 0 and  $k' \in \mathbb{Z}$  such that

$$\|x^{(k+1)} - x^{(k)}\| \le pd_{k+k'}, \quad \forall \ k \ge 0.$$

[H6] The sequences  $\{a_k\}_{k\in\mathbb{N}}, \{b_k\}_{k\in\mathbb{N}}$  satisfy the following conditions

$$\sum_{k=0}^{+\infty} b_k = +\infty, \quad \sup_{k\in\mathbb{N}} \frac{1}{b_k a_k} < +\infty, \quad \inf_{k\in\mathbb{N}} a_k > 0.$$

The abstract scheme given in Conditions 3.1 can be seen as a further extension of the one proposed in [30], which is indeed recovered by setting  $\Phi = \mathcal{F}$ ,  $u^{(k)} = x^{(k+1)}$  and  $\rho^{(k)} \equiv s^{(k)}$ . In the following, we discuss in detail conditions [H1]-[H6] and their relation with the abstract scheme in [30].

- Condition [H1] requires the sufficient decrease of a proper, lower semicontinuous function  $\Phi$  between two successive iterates. The quantity  $a_k d_k^2$  measures the amount of the decrease, where  $d_k$  is thought as a generalization of the Euclidean norm, whereas the parameter  $s^{(k)}$  allows for some flexibility in the asymptotic behaviour of the function  $\Phi$ . Note that, in earlier works based on the KL property [2, 9, 12, 26], condition [H1] is usually presented by setting  $d_k = ||x^{(k+1)} x^{(k)}||_2$ ,  $\rho^{(k)} \equiv s^{(k)} \equiv 0$  and  $\Phi(x, s) = f(x)$ , being f the function to minimize. The generalized condition reported here is almost identical to the one introduced in the more recent work [30], with the only difference that here the sufficient decrease is required on a function  $\Phi$  that may be different from the function  $\mathcal{F}$  appearing in [H3].
- Condition [H3] is the so-called relative error condition, which is related to the (possibly) inexact solution of the minimization subproblem performed at each iteration of a first-order method. In the previous literature [2, 9, 12, 26], such condition is usually employed by setting  $u^{(k)} = x^{(k+1)}$ ,  $d_k = ||x^{(k+1)} x^{(k)}||_2$ ,  $\rho^{(k)} \equiv 0$ ,  $I = \{1\}$ ,  $\theta_1 = 1$  and  $\mathcal{F}(u, \rho) = f(u)$ , being f the function to minimize. In [30], a general positive term  $d_k$ , a finite index set I, a variable parameter  $s^{(k)}$  and a generic surrogate function  $\mathcal{F}$  are employed, while keeping  $u^{(k)} = x^{(k+1)}$  and  $\rho^{(k)} \equiv 0$ . Here we also allow the sequence  $\{u^{(k)}\}_{k\in\mathbb{N}}$  to be distinct from  $\{x^{(k)}\}_{k\in\mathbb{N}}$  and the parameters  $\{\rho^{(k)}\}_{k\in\mathbb{N}}$  to vary at each iteration. The reason to do so comes from the fact that condition [H3] is hard to enforce algorithmically on  $x^{(k+1)}$  when the minimization subproblem is solved inexactly, as noted in [15, 16, 29]. However, if a specific, implementable inexactness criterion is adopted for the solution of the subproblem, then the same condition holds for a surrogate function  $\mathcal{F}$  evaluated at a different iterate  $(u^{(k)}, \rho^{(k)})$ . For instance, this is observed in the convergence analysis of the so-called VMILA algorithm, a variable metric linesearch based forward-backward method studied in [12, 15, 16]. In [12], VMILA is included in the KL framework by setting

$$\begin{split} \Phi(x,s) &= \mathcal{F}(x,\rho) = f(x), \ u^{(k)} = x^{(k+1)}, \ d_k = \|x^{(k+1)} - x^{(k)}\|_2, \ \rho^{(k)} \equiv s^{(k)} \equiv 0 \ \text{and noting} \\ \text{that, in so doing, the relative error condition holds only by exactly computing the proximal \\ \text{operator. In [15], the authors include VMILA in the abstract scheme in a different way, \\ \text{by using } \Phi(x,s) = f(x), \ \text{a surrogate function } \mathcal{F} \ \text{defined upon the concept of forward-} \\ \text{backward envelope of } f \ [36], \ \text{the iterate } u^{(k)} \ \text{as the inexact proximal-gradient point } \tilde{y}^{(k)}, \\ d_k = \|x^{(k+1)} - x^{(k)}\|_2 \ \text{and } \rho^{(k)} \ \text{as the error parameter due to the computation of } \tilde{y}^{(k)}. \\ \text{Finally, in [16], VMILA is framed by setting } \Phi(x,s) = f(x), \ \mathcal{F}(u,\rho) = f(u) + \rho^2/2, \ \text{the iterate } u^{(k)} \ \text{as the exact proximal-gradient point } y^{(k)} \ \text{and } d_k = -h^{(k)}(\tilde{y}^{(k)}), \ \text{where } h^{(k)} \\ \text{is the function to minimize when computing the approximation } \tilde{y}^{(k)} \ \text{of the exact point } y^{(k)} \ \text{no ther words, from the analysis in [16], it turns out that we are able to enforce the relative error condition at the exact point <math>y^{(k)}, \ \text{which we do not need to compute explicitly, provided that the approximation } \tilde{y}^{(k)} \ \text{is computed using a specific criterion.} \end{split}$$

- Condition [H2] is crucial in the convergence proof of Theorem 3.3. Indeed it ensures that the sequence  $\{\mathcal{F}(u^{(k)}, \rho^{(k)})\}_{k \in \mathbb{N}}$  converges to a limit value  $\mathcal{F}^* \in \mathbb{R}$ , thus allowing to apply the uniformized KL property at the point  $(u^{(k)}, \rho^{(k)})$  for all sufficiently large k and, furthermore, it enables the combination of [H3] and [H1] with the KL inequality. Imposing condition [H2] is required only when  $\Phi \neq \mathcal{F}$ , this is why it does not appear in [30].
- Condition [H4] is the analogue of the so-called *continuity condition* in [30]. Here we impose the property for all converging subsequences  $(x^{(k_j)}, \rho^{(k_j)})$ , whereas in [30] it is only required the existence of one such subsequence. This is because, unlike in [30], we need to ensure that the distance between  $\{(u^{(k)}, \rho^{(k)})\}_{k \in \mathbb{N}}$  and the limit set of  $\{(x^{(k)}, \rho^{(k)})\}_{k \in \mathbb{N}}$  converges to 0 (see Lemma 3.2(*iii*)).
- Condition [H5] is also called the *distance condition*. It states the connection between the general term  $d_k$  and the Euclidean norm, which is fundamental in order to prove the finite length of the sequence  $\{x^{(k)}\}_{k\in\mathbb{N}}$ . Note that the distance condition given in [30] is slightly more general than [H5]; however, that condition alone allows to prove only the finite length of the sequence  $\{d_k\}_{k\in\mathbb{N}}$ , which in general does not imply the convergence of the sequence  $\{x^{(k)}\}_{k\in\mathbb{N}}$ . In order to obtain the strongest result, condition [H5] is then imposed in [30, Theorem 10].
- Condition [H6] is the same as the parameter condition in [30]. These requirements on the sequences  $\{a_k\}_{k\in\mathbb{N}}, \{b_k\}_{k\in\mathbb{N}}$  were first introduced in [26] in order to generalize the abstract descent scheme in [2].

In the remainder of this section, we will denote with  $\{x^{(k)}\}_{k\in\mathbb{N}}, \{u^{(k)}\}_{k\in\mathbb{N}}, \{\rho^{(k)}\}_{k\in\mathbb{N}}, \{s^{(k)}\}_{k\in\mathbb{N}}$ the sequences complying with Conditions 3.1. Furthermore, let us define the set of all limit points of the sequence  $\{(x^{(k)}, \rho^{(k)})\}_{k\in\mathbb{N}}$ :

$$\Omega^*(x^{(0)}, \rho^{(0)}) = \{ (x^*, \rho^*) \in \mathbb{R}^n \times \mathbb{R}^m : \exists \{k_j\}_{j \in \mathbb{N}} \subset \mathbb{N} \text{ such that } (x^{(k_j)}, \rho^{(k_j)}) \to (x^*, \rho^*) \}.$$

Note that the set  $\Omega^*(x^{(0)}, \rho^{(0)})$  can be written as

$$\Omega^*(x^{(0)}, \rho^{(0)}) = X^*(x^{(0)}) \times R^*(\rho^{(0)})$$

where  $X^*(x^{(0)}) = \{x^* \in \mathbb{R}^n : \exists \{k_j\}_{j \in \mathbb{N}} \subset \mathbb{N} \text{ such that } x^{(k_j)} \to x^*\} \subset \mathbb{R}^n, R^*(\rho^{(0)}) = \{\rho^* \in \mathbb{R}^m : \exists \{k_j\}_{j \in \mathbb{N}} \subset \mathbb{N} \text{ such that } \rho^{(k_j)} \to \rho^*\} \subset \mathbb{R}^m.$ 

LEMMA 3.2. Let Conditions 3.1 be satisfied. Suppose that  $\{(x^{(k)}, \rho^{(k)})\}_{k \in \mathbb{N}}$  is a bounded sequence. Then the following facts hold true.

(i)  $\Omega^*(x^{(0)}, \rho^{(0)})$  is nonempty and compact.

(ii) There exists  $\mathcal{F}^* \in \mathbb{R}$  such that  $\lim_{k \to \infty} \Phi(x^{(k)}, s^{(k)}) = \lim_{k \to \infty} \mathcal{F}(u^{(k)}, \rho^{(k)}) = \mathcal{F}^*.$ 

(iii) We have

$$\lim_{k \to \infty} \operatorname{dist}((x^{(k)}, \rho^{(k)}), \Omega^*(x^{(0)}, \rho^{(0)})) = \lim_{k \to \infty} \operatorname{dist}((u^{(k)}, \rho^{(k)}), \Omega^*(x^{(0)}, \rho^{(0)})) = 0.$$

(iv) We have  $\mathcal{F}(x^*, \rho^*) = \mathcal{F}^*, \, \forall \, (x^*, \rho^*) \in \Omega^*(x^{(0)}, \rho^{(0)}).$ 

*Proof.* (i) Since the sequence  $\{(x^{(k)}, \rho^{(k)})\}_{k \in \mathbb{N}}$  is bounded, it admits at least a limit point and, hence,  $\Omega^*(x^{(0)}, \rho^{(0)})$  is nonempty. Compactness can be proved by observing that  $\Omega^*(x^{(0)}, \rho^{(0)})$  is

a countable intersection of compact sets (see [9, Lemma 5]). (ii) From [H1] we have that the sequence  $\{\Phi(x^{(k)}, s^{(k)})\}_{k \in \mathbb{N}}$  is nonincreasing and, since  $\Phi$  is bounded from below, there exists  $\mathcal{F}^* \in \mathbb{R}$  such that

$$\lim_{k \to \infty} \Phi(x^{(k)}, s^{(k)}) = \mathcal{F}^*$$

The previous relation combined with [H2] proves Part (ii). (iii) Since  $\{(x^{(k)}, \rho^{(k)})\}_{k \in \mathbb{N}}$  is bounded and by definition of  $\Omega^*(x^{(0)}, \rho^{(0)})$ , we have

$$\lim_{k \to \infty} \operatorname{dist}((x^{(k)}, \rho^{(k)}), \Omega^*(x^{(0)}, \rho^{(0)})) = 0.$$

Observing that  $\Omega^*(x^{(0)}, \rho^{(0)}) = X^*(x^{(0)}) \times R^*(\rho^{(0)})$  and recalling (2.1), the previous limit implies

$$\lim_{k \to \infty} \operatorname{dist}(x^{(k)}, X^*(x^{(0)})) = 0, \quad \lim_{k \to \infty} \operatorname{dist}(\rho^{(k)}, R^*(\rho^{(0)})) = 0.$$

Combining the boundedness of  $\{x^{(k)}\}_{k\in\mathbb{N}}$  with the definition of  $X^*(x^{(0)})$  and property [H4], we obtain  $\lim_{k\to\infty} \operatorname{dist}(u^{(k)}, X^*(x^{(0)})) = 0$ , which together with the second limit above yields

$$\lim_{k \to \infty} \operatorname{dist}((u^{(k)}, \rho^{(k)}), \Omega^*(x^{(0)}, \rho^{(0)})) = 0.$$

(iv) This point follows directly from part (ii) and [H4].

THEOREM 3.3. Let Conditions 3.1 be satisfied, and suppose that  $\{(x^{(k)}, \rho^{(k)})\}_{k \in \mathbb{N}}$  is a bounded sequence and that  $\mathcal{F}$  is a KL function. Then the following statements are true.

(i) The sequence  $\{d_k\}_{k\in\mathbb{N}}$  is summable, i.e., it satisfies

$$\sum_{k=0}^{+\infty} d_k < +\infty$$

(ii) The sequence  $\{x^{(k)}\}_{k\in\mathbb{N}}$  has finite length, i.e., it satisfies

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

and thus  $\{x^{(k)}\}_{k\in\mathbb{N}}$  is a convergent sequence.

(iii) If also  $\{\rho^{(k)}\}_{k\in\mathbb{N}}$  converges, then the sequence  $\{(x^{(k)},\rho^{(k)})\}_{k\in\mathbb{N}}$  converges to a stationary point for  $\mathcal{F}$ .

*Proof.* (i) By Lemma 3.2(*i*)-(*iv*), the function  $\mathcal{F}$  is constant over the compact set  $\Omega^*(x^{(0)}, \rho^{(0)})$ , therefore we can apply Lemma 2.5. Let  $v, \mu, \phi, \bar{B}$  be as in Lemma 2.5. Thanks to Lemma 3.2(*ii*)-(*iii*) and [H2], there exists a positive integer  $k_0$  such that

(3.1) 
$$\Phi(x^{(k)}, s^{(k)}) + r_k < \mathcal{F}^* + \upsilon, \quad \operatorname{dist}((u^{(k)}, \rho^{(k)}), \Omega^*(x^{(0)}, \rho^{(0)})) < \mu$$

for all  $k \ge k_0$ . Without loss of generality, up to a translation of the iteration index, we can assume  $k_0 = 0$ .

Let us now set  $c = \sup_k 1/(a_k b_k)$ , where  $c < +\infty$  due to [H6],  $\zeta'_k = \zeta_k/b$  and

$$\phi_k = \frac{b}{c} (\phi(\Phi(x^{(k)}, s^{(k)}) - \mathcal{F}^*) - \phi(\Phi(x^{(k+1)}, s^{(k+1)}) - \mathcal{F}^*))$$

and prove that

(3.2) 
$$2d_k \le \phi_k + \sum_{i \in I} \theta_i d_{k-i} + \zeta'_k, \quad \forall \ k \ge 1.$$

We first observe that the definition of  $\phi_k$  is well posed, since  $\mathcal{F}^*$  is the limit of the nonincreasing sequence  $\{\Phi(x^{(k)}, s^{(k)})\}_{k \in \mathbb{N}}$  (see Lemma 3.2) and, hence, we have  $\Phi(x^{(k)}, s^{(k)}) \geq \mathcal{F}^*, \forall k \in \mathbb{N}$ . Moreover, the monotonicity of both function  $\phi$  and sequence  $\{\Phi(x^{(k)}, s^{(k)})\}_{k \in \mathbb{N}}$  implies that  $\phi_k \geq 0$ ,  $\forall k \in \mathbb{N}$ .

Let us now consider the two cases  $d_k = 0$  and  $d_k > 0$  separately.

If  $d_k = 0$ , inequality (3.2) holds trivially. Otherwise, for any iteration index  $k \ge 1$  such that  $d_k > 0$ , taking into account [H1] and [H2], we can write

$$\mathcal{F}^* \le \Phi(x^{(k+1)}, s^{(k+1)}) < \Phi(x^{(k)}, s^{(k)}) \le \mathcal{F}(u^{(k-1)}, \rho^{(k-1)}).$$

On the other hand, the rightmost inequality in [H2] gives

$$\mathcal{F}(u^{(k-1)}, \rho^{(k-1)}) \le \Phi(x^{(k-1)}, s^{(k-1)}) + r_{k-1}$$

which, in view of (3.1), implies that  $(u^{(k-1)}, \rho^{(k-1)}) \in \overline{B}$ . Then, we can write the KL inequality related to the function  $\mathcal{F}$  at  $(u^{(k-1)}, \rho^{(k-1)})$ :

$$\phi'(\mathcal{F}(u^{(k-1)},\rho^{(k-1)}) - \mathcal{F}^*) \ge \frac{1}{\|\partial \mathcal{F}(u^{(k-1)},\rho^{(k-1)})\|_{-}}.$$

Furthermore, combining the previous inequality with [H3] yields

$$\phi'(\mathcal{F}(u^{(k-1)}, \rho^{(k-1)}) - \mathcal{F}^*) \ge \frac{1}{\frac{b}{b_k} \sum_{i \in I} \theta_i d_{k-i} + \frac{1}{b_k} \zeta_k}$$

Since  $\phi$  is concave,  $\phi'$  is nonincreasing. Therefore, [H2] implies

$$\phi'(\Phi(x^{(k)}, s^{(k)}) - \mathcal{F}^*) \ge \phi'(\mathcal{F}(u^{(k-1)}, \rho^{(k-1)}) - \mathcal{F}^*).$$

Exploiting again the concavity of  $\phi$ , we have

 $\phi(\Phi(x^{(k)}, s^{(k)}) - \mathcal{F}^*) - \phi(\Phi(x^{(k+1)}, s^{(k+1)}) - \mathcal{F}^*) \ge \phi'(\Phi(x^{(k)}, s^{(k)}) - \mathcal{F}^*)(\Phi(x^{(k)}, s^{(k)}) - \Phi(x^{(k+1)}, s^{(k+1)})).$ Combining the last three relations with [H1] leads to

$$\begin{split} \phi(\Phi(x^{(k)}, s^{(k)}) - \mathcal{F}^*) - \phi(\Phi(x^{(k+1)}, s^{(k+1)}) - \mathcal{F}^*) &\geq \frac{\Phi(x^{(k)}, s^{(k)}) - \Phi(x^{(k+1)}, s^{(k+1)})}{\frac{b}{b_k} \sum_{i \in I} \theta_i d_{k-i} + \frac{1}{b_k} \zeta_k} \\ &\geq \frac{a_k d_k^2}{\frac{b}{b_k} \sum_{i \in I} \theta_i d_{k-i} + \frac{1}{b_k} \zeta_k}. \end{split}$$

Recalling the definition of  $\phi_k$  and  $\zeta'_k$ , the above inequality implies the following one

$$d_k^2 \le \phi_k \left( \sum_{i \in I} \theta_i d_{k-i} + \zeta'_k \right).$$

Taking the square root of both sides and using the inequality  $2\sqrt{uv} \le u+v$  on the right-hand-side, we obtain (3.2).

Summing (3.2) from 1 to k leads to

(3.3) 
$$2\sum_{j=1}^{k} d_j \le \sum_{j=1}^{k} \phi_j + \sum_{j=1}^{k} \sum_{i \in I} \theta_i d_{j-i} + \sum_{j=1}^{k} \zeta'_j.$$

We now observe that

$$\sum_{j=1}^{k} \phi_j = \frac{b}{c} (\phi(\Phi(x^{(1)}, s^{(1)}) - \mathcal{F}^*) - \phi(\Phi(x^{(k+1)}, s^{(k+1)}) - \mathcal{F}^*))$$
$$\leq \frac{b}{c} \phi(\Phi(x^{(1)}, s^{(1)}) - \mathcal{F}^*),$$

where the rightmost inequality follows from the positive sign of  $\phi$ . Furthermore, the second sum in the right-hand side of (3.3) can be rewritten as below

$$\sum_{j=1}^{k} \sum_{i \in I} \theta_i d_{j-i} = \sum_{i \in I} \sum_{j=1}^{k} \theta_i d_{j-i} = \sum_{i \in I} \sum_{r=1-i}^{k-i} \theta_i d_r$$
$$\leq \sum_{i \in I} \sum_{r=1-i}^{0} \theta_i d_r + \left(\sum_{i \in I} \theta_i\right) \sum_{r=1}^{k} d_r + \sum_{i \in I} \sum_{r=k+1}^{k-i} \theta_i d_r$$
$$= \sum_{i \in I} \sum_{r=1-i}^{0} \theta_i d_r + \sum_{r=1}^{k} d_r + \sum_{i \in I} \sum_{r=k+1}^{k-i} \theta_i d_r$$

where we have used the change of variable r = j - 1 and the property  $\sum_{i \in I} \theta_i = 1$ . Note that the sums appearing in the previous relation are assumed to be zero whenever the start index of the summation is larger than the termination index. Therefore we can write

$$2\sum_{j=1}^{k} d_j \le \sum_{i \in I} \sum_{r=1-i}^{0} \theta_i d_r + \sum_{r=1}^{k} d_r + \sum_{i \in I} \sum_{r=k+1}^{k-i} \theta_i d_r + \frac{b}{c} \phi(\Phi(x^{(1)}, s^{(1)}) - \mathcal{F}^*) + \sum_{j=1}^{k} \zeta_j'$$

which clearly implies

(3.4) 
$$\sum_{j=1}^{k} d_j \leq \sum_{i \in I} \sum_{r=1-i}^{0} \theta_i d_r + \sum_{i \in I} \sum_{r=k+1}^{k-i} \theta_i d_r + \frac{b}{c} \phi(\Phi(x^{(1)}, s^{(1)}) - \mathcal{F}^*) + \sum_{j=1}^{k} \zeta_j'.$$

At this point, observe that the first two sums in the right-hand side of (3.4) are finite linear combinations of the terms  $\{d_r\}_{r\in\mathbb{N}}$ . Conditions [H1] and [H6] ensure that  $d_k \to 0$ , hence those sums are converging to 0 for  $k \to \infty$ . Noting also that  $\{\zeta'_k\}_{k\in\mathbb{N}}$  is summable and taking the limit of (3.4) for  $k \to \infty$ , we obtain

(3.5) 
$$\sum_{k=0}^{\infty} d_k < +\infty.$$

(ii) Combining (3.5) with /H5/, we also obtain

$$\sum_{k=0}^{\infty} \|x^{(k+1)} - x^{(k)}\| \le p \sum_{k=0}^{\infty} d_{k+k'} < +\infty.$$

which implies that the sequence  $\{x^{(k)}\}_{k\in\mathbb{N}}$  converges to a point  $x^* \in \mathbb{R}^n$ . (iii) Let  $(x^*, \rho^*) \in \mathbb{R}^n \times \mathbb{R}^m$  be the unique limit point of the sequence  $\{(x^{(k)}, \rho^{(k)})\}_{k\in\mathbb{N}}$ , namely  $(x^{(k)}, \rho^{(k)}) \to (x^*, \rho^*)$ . By using [H4], it follows that  $(u^{(k)}, \rho^{(k)}) \to (x^*, \rho^*)$  and  $\mathcal{F}(u^{(k)}, \rho^{(k)}) \to \mathcal{F}(x^*, \rho^*)$ . Furthermore, summing [H3] for  $k = 0, \ldots, K$  yields

$$\sum_{k=0}^{K} b_{k+1} \|\partial \mathcal{F}(u^{(k)}, \rho^{(k)})\|_{-} \le b \sum_{k=0}^{K} \sum_{i \in I} \theta_{i} d_{k+1-i} + \sum_{k=0}^{K} \zeta_{k}.$$

Taking the limit for  $K \to \infty$ , using (3.5) and recalling that  $\{\zeta_k\}_{k \in \mathbb{N}}$  is summable, we obtain

$$\sum_{k=0}^{\infty} b_{k+1} \|\partial \mathcal{F}(u^{(k)}, \rho^{(k)})\|_{-} < +\infty.$$

Since [H6] requires  $\sum_k b_k = +\infty$ , the previous relation implies that

$$\liminf_{k \to \infty} \|\partial \mathcal{F}(u^{(k)}, \rho^{(k)})\|_{-} = 0.$$

In conclusion, the sequence  $\{(u^{(k)}, \rho^{(k)})\}_{k \in \mathbb{N}}$  satisfies all the hypotheses of Lemma 2.3, which means that  $0 \in \partial \mathcal{F}(x^*, \rho^*)$ .

4. Applications of the abstract scheme. In this section, we show how we can devise some brand new forward–backward–type algorithms satisfying Conditions 3.1 and, hence, guarantee their convergence to a stationary point in virtue of Theorem 3.3. In particular, from now on, we address the problem

(4.1) 
$$\min_{x \in \mathbb{R}^n} f(x), \quad f(x) = f_0(x) + f_1(x),$$

where we assume that  $f_0, f_1$  are as follows:

[A1]  $f_1 : \mathbb{R}^n \to \mathbb{R} \cup \{+\infty\}$  is a proper, lower semicontinuous, convex function;

[A2]  $f_0: \mathbb{R}^n \to \mathbb{R}$  is continuously differentiable on an open set  $\Omega_0 \supset \overline{\operatorname{dom}(f_1)}$ .

[A3]  $f_0$  has L-Lipschitz continuous gradient on dom $(f_1)$ , i.e.,

$$\|\nabla f_0(x) - \nabla f_0(y)\| \le L \|x - y\|, \quad \forall x, y \in \operatorname{dom}(f_1),$$

for some L > 0.

[A4] f is bounded from below. Under the above assumptions, for any  $z \in \text{dom}(f_1)$ , the following subdifferential calculus rules

hold [35, Proposition 8.12, Exercise 8.8(c)]  $\Im f(x) = \{x \in \mathbb{D}^n : f(x) > f(x) + \langle x \rangle, x \in \mathbb{D}^n \}$ 

(4.2) 
$$\partial f_1(z) = \{ w \in \mathbb{R}^n : f_1(y) \ge f_1(z) + \langle w, y - z \rangle, \ \forall y \in \mathbb{R}^n \}$$
$$\partial f(z) = \{ \nabla f_0(z) \} + \partial f_1(z).$$

One of the most popular algorithms for solving problem (4.1) is the inertial, or heavy ball, proximalgradient method iPiano [31, 30], which is defined by the iteration

(4.3) 
$$x^{(k+1)} = \operatorname{prox}_{\alpha_k f_1}(x^{(k)} - \alpha_k \nabla f_0(x^{(k)}) + \beta_k(x^{(k)} - x^{(k-1)})),$$

where  $\alpha_k, \beta_k$  are suitably chosen parameters. By definition of the proximity operator, the above updating rule consists in defining the new point as the unique solution of the minimization problem

(4.4) 
$$\min_{y \in \mathbb{R}^n} f_1(y) - f_1(x^{(k)}) + \langle \nabla f_0(x^{(k)}) - \frac{\beta_k}{\alpha_k} (x^{(k)} - x^{(k-1)}), y - x^{(k)} \rangle + \frac{1}{2\alpha_k} \|y - x^{(k)}\|^2.$$

We will refer to the minimizer of this problem as the *inertial proximal gradient point*. If  $\beta_k = 0$ , we recover the standard proximal gradient point, otherwise the inertial step  $x^{(k)} - x^{(k-1)}$  is included in the argument of the proximal gradient operator, with the aim of improving the convergence behaviour of the overall method.

In the following, we address the key challenge of designing algorithms that inexactly compute the inertial proximal gradient point with implementable conditions that still preserve the convergence guarantees of Theorem 3.3. More precisely, we propose two new inexact inertial-type algorithms, where the second one also features a linesearch procedure along a descent direction of a suitable merit function. The convergence analysis of both algorithms can be performed in the abstract framework provided by Conditions 3.1. We stress that, due to the inexactness in the computation of the inertial proximal gradient point, our algorithms cannot be cast in the abstract frameworks proposed in previous works.

We start our presentation by defining the inexactness criterion for the inertial proximal gradient point, which is a generalization of the one proposed in [11, 12, 15].

**4.1. Inexact inertial proximal gradient point.** Given two positive parameters  $\alpha, \beta$ , consider the function  $h : \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R} \cup \{+\infty\}$  defined as follows:

(4.5) 
$$h(y;x,s) = f_1(y) - f_1(x) + \langle \nabla f_0(x) - \frac{\beta}{\alpha}(x-s), y-x \rangle + \frac{1}{2\alpha} \|y-x\|^2.$$

Clearly, the inertial proximal gradient point (4.3) is the minimizer of the above function with respect to the first argument, with  $\alpha = \alpha_k$ ,  $\beta = \beta_k$ ,  $x = x^{(k)}$ ,  $s = x^{(k-1)}$ . Given (x, s), we denote by  $\hat{y}$  the (exact) minimizer of the function in (4.5)

(4.6) 
$$\hat{y} = \operatorname*{argmin}_{y \in \mathbb{R}^n} h(y; x, s) = \operatorname{prox}_{\alpha f_1} (x - \alpha \nabla f_0(x) + \beta (x - s)).$$

The point  $\hat{y}$  is the unique point satisfying the optimality condition

(4.7) 
$$0 \in \partial h(\hat{y}; x, s) \Leftrightarrow -\frac{1}{\alpha}(\hat{y} - x + \alpha \nabla f_0(x) - \beta(x - s)) \in \partial f_1(\hat{y}).$$

Borrowing the ideas in [11, 16], we define an approximation of  $\hat{y}$  as any point  $\tilde{y} \in \text{dom}(f_1)$  such that

(4.8) 
$$h(\tilde{y};x,s) - h(\hat{y};x,s) \le -\frac{\tau}{2}h(\tilde{y};x,s).$$

The above condition is equivalent to the following one:

(4.9) 
$$h(\tilde{y}; x, s) \le \left(\frac{2}{2+\tau}\right) h(\hat{y}; x, s) \le 0,$$

where the rightmost inequality is a consequence of the fact that  $\hat{y}$  is a minimizer of  $h(\cdot; x, s)$  and h(x; x, s) = 0. Therefore, we have  $h(\tilde{y}; x, s) \leq 0$  and condition (4.8) can be rewritten in equivalent way as

(4.10) 
$$0 \in \partial_{\epsilon} h(\tilde{y}; x, s), \text{ with } \epsilon = -\frac{\tau}{2} h(\tilde{y}; x, s),$$

where  $\partial_{\epsilon}h(\tilde{y}; x, s) = \{w \in \mathbb{R}^n : h(y; x, s) \ge h(\tilde{y}; x, s) + \langle w, y - \tilde{y} \rangle - \epsilon, \forall y \in \mathbb{R}^n\}$  is the  $\epsilon$ -subdifferential in the first argument of the convex function  $h(\cdot; x, s)$  at point  $\tilde{y}$  [38, p. 82]. Therefore, the point  $\tilde{y}$ is defined upon a relaxation of the optimality condition (4.7), where the subdifferential of  $h(\cdot; x, s)$ is replaced by the  $\epsilon$ -subdifferential and the accuracy parameter  $\epsilon$  is chosen in a specific way, which is crucial for preserving the theoretical convergence properties.

Even if the inclusion (4.10) is implicit, a point  $\tilde{y} \in \text{dom}(f_1)$  satisfying (4.10) can be actually computed in practice with a well defined, explicit primal-dual procedure, as explained in [11, 17]: an iterative optimization method is applied to the dual of the function  $h(\cdot; x, s)$ , until the difference between the primal and the dual function is below the tolerance  $\epsilon$ . A more detailed discussion about the inexactness criterion and on its practical realization can be found also in [16, Section 3.1].

**4.1.1. Preliminary results.** We collect new several basic results to incorporate inexactness into the design and analysis of our proposed algorithms. The following lemma is a consequence of the strong convexity of the function  $h(\cdot; x, s)$ , and will be often employed in the following.

LEMMA 4.1. Suppose that Assumptions [A1]-[A2] hold true. For a given pair  $(x, s) \in \text{dom}(f_1) \times \mathbb{R}^n$ , let  $\hat{y}, \tilde{y}$  be defined as in (4.6),(4.8). Then, the following inequalities hold.

(4.11) 
$$\frac{1}{2\alpha} \|\hat{y} - x\|^2 \le \left(1 + \frac{\tau}{2}\right) \left(-h(\tilde{y}; x, s)\right)$$

(4.12) 
$$\frac{1}{2\alpha} \|\tilde{y} - \hat{y}\|^2 \le \frac{\gamma}{2} (-h(\tilde{y}; x, s))$$

(4.13) 
$$\frac{\theta}{2\alpha} \|\tilde{y} - x\|^2 \le (-h(\tilde{y}; x, s)), \quad \text{with } \theta = 1/\left(\sqrt{1 + \frac{\tau}{2}} + \sqrt{\frac{\tau}{2}}\right)^2 \le 1.$$

*Proof.* Inequalities (4.11)–(4.12) follow by combining the strong convexity of the function  $h(\cdot; x, s)$  and condition (4.8) as in [16, Lemma 2]. As for (4.13), we have

$$\begin{split} \frac{1}{2\alpha} \|\tilde{y} - x\|^2 &= \frac{1}{2\alpha} \|\tilde{y} - \hat{y} + \hat{y} - x\|^2 = \frac{1}{2\alpha} \|\tilde{y} - \hat{y}\|^2 + \frac{1}{2\alpha} \|\hat{y} - x\|^2 + \frac{1}{\alpha} \langle \tilde{y} - \hat{y}, \hat{y} - x \rangle \\ &\leq \frac{1}{2\alpha} \|\tilde{y} - \hat{y}\|^2 + \frac{1}{2\alpha} \|\hat{y} - x\|^2 + \frac{1}{\alpha} \|\tilde{y} - \hat{y}\| \cdot \|\hat{y} - x\| \\ &\leq \left(1 + \frac{\tau}{2}\right) \left(-h(\tilde{y}; x, s)\right) + \frac{\tau}{2} \left(-h(\tilde{y}; x, s)\right) + 2\sqrt{1 + \frac{\tau}{2}} \sqrt{\frac{\tau}{2}} \left(-h(\tilde{y}; x, s)\right) \\ &= \left(\sqrt{1 + \frac{\tau}{2}} + \sqrt{\frac{\tau}{2}}\right)^2 \left(-h(\tilde{y}; x, s)\right), \end{split}$$

where the last inequality follows from the application of (4.11)-(4.12).

The next lemma provides a subgradient  $\hat{v} \in \partial f(\hat{y})$  whose norm is bounded from above by a quantity containing  $\sqrt{-h(\tilde{y}; x, s)}$ . Its proof is omitted since it is almost identical to the one of Lemma 3 in [16].

LEMMA 4.2. Suppose Assumptions [A1]–[A3] hold true. Let x be a point in dom(f<sub>1</sub>) and let  $\hat{y}, \tilde{y}$  be defined as in (4.6)–(4.8). Moreover, assume that  $\alpha \in [\alpha_{min}, \alpha_{max}]$ , with  $0 < \alpha_{min} \leq \alpha_{max}$  and  $\beta \in [0, \beta_{max}]$ , with  $\beta_{max} \geq 0$ . Then, there exists a subgradient  $\hat{v} \in \partial f(\hat{y})$  such that

(4.14) 
$$\|\hat{v}\| \le p(\|\hat{y} - x\| + \|x - s\|)$$

(4.15) 
$$\leq q(\sqrt{-h(\tilde{y};x,s) + \|x-s\|}),$$

where the two constants p, q depend only on  $\alpha_{min}, \alpha_{max}, \beta_{max}$  and on the Lipschitz constant L.

The following lemma is the equivalent of Lemma 4-5 in [16].

LEMMA 4.3. Suppose Assumptions [A1]–[A3] hold true and assume  $0 < \alpha_{min} \leq \alpha \leq \alpha_{max}$ ,  $\beta \in [0, \beta_{max}]$ . Let (x, s) be a point in dom $(f_1) \times \mathbb{R}^n$  and let  $\hat{y}, \tilde{y}$  be defined as in (4.6)–(4.8), for some  $\tau \geq 0$ . Then, there exists  $c, d, \bar{c}, \bar{d} \in \mathbb{R}$  depending only on  $\alpha_{min}, \alpha_{max}, \beta_{max}, \tau$  such that

(4.16) 
$$f(\hat{y}) \ge f(\tilde{y}) + ch(\tilde{y}; x, s) - d||x - s||^2$$

(4.17)  $f(\hat{y}) \le f(x) - \bar{c}h(\tilde{y}; x, s) + \bar{d} ||x - s||^2.$ 

*Proof.* From the Descent Lemma [5, Proposition A.24] we have

(4.18) 
$$f_0(\hat{y}) \ge f_0(\tilde{y}) - \langle \nabla f_0(\hat{y}), \tilde{y} - \hat{y} \rangle - \frac{L}{2} \| \tilde{y} - \hat{y} \|^2$$

The inclusion  $0 \in \partial_{\epsilon} h(\tilde{y}; x, s)$  in (4.10) implies that there exists a vector  $e \in \mathbb{R}^n$  with

(4.19) 
$$\frac{1}{2\alpha} \|e\|^2 \le \epsilon$$

such that

$$-\frac{1}{\alpha}(\tilde{y} - x + \alpha \nabla f_0(x) - \beta(x - s) + e) \in \partial_{\epsilon} f_1(\tilde{y})$$

(see [12] and references therein). The definition of  $\epsilon$ -subdifferential implies

$$(4.20) \quad f_1(\hat{y}) \ge f_1(\tilde{y}) - \frac{1}{\alpha} \langle \tilde{y} - x, \hat{y} - \tilde{y} \rangle + \frac{\beta}{\alpha} \langle \hat{y} - \tilde{y}, x - s \rangle - \langle \nabla f_0(x), \hat{y} - \tilde{y} \rangle - \frac{1}{\alpha} \langle e, \hat{y} - \tilde{y} \rangle - \epsilon.$$

Summing inequalities (4.18) and (4.20) yields

(4.21) 
$$f(\hat{y}) \ge f(\tilde{y}) - \langle \nabla f_0(x) - \nabla f_0(\hat{y}), \hat{y} - \tilde{y} \rangle + \frac{\beta}{\alpha} \langle \hat{y} - \tilde{y}, x - s \rangle$$
$$-\frac{1}{\alpha} \langle \tilde{y} - x, \hat{y} - \tilde{y} \rangle - \frac{1}{\alpha} \langle e, \hat{y} - \tilde{y} \rangle - \frac{L}{2} \| \tilde{y} - \hat{y} \|^2 - \epsilon.$$

Now we consider each term at the right-hand-side in the above inequality so as to obtain a lower bound. Using the Cauchy-Schwarz inequality, Assumption [A3], (4.11) and (4.12) we obtain

(4.22) 
$$\begin{aligned} \langle \nabla f_0(x) - \nabla f_0(\hat{y}), \hat{y} - \tilde{y} \rangle &\leq \| \nabla f_0(x) - \nabla f_0(\hat{y}) \| \| \hat{y} - \tilde{y} \| \\ &\leq L \alpha_{max} \sqrt{2\tau \left( 1 + \frac{\tau}{2} \right)} (-h(\tilde{y}; x, s)) \end{aligned}$$

Similarly, using again the Cauchy-Schwarz inequality, (4.12) and (4.13), we can write

(4.23) 
$$\frac{1}{\alpha} \langle \tilde{y} - x, \hat{y} - \tilde{y} \rangle \leq \frac{1}{\alpha_{min}} \|\tilde{y} - x\| \|\hat{y} - \tilde{y}\| \leq \frac{\alpha_{max}}{\alpha_{min}} \sqrt{\frac{2\tau}{\theta}} (-h(\tilde{y}; x, s)).$$

Moreover, from (4.19) and (4.10) we obtain  $||e|| \leq \sqrt{2\alpha_{max}\epsilon} \leq \sqrt{\alpha_{max}\tau(-h(\tilde{y};x,s))}$  which, using also (4.12), yields

(4.24) 
$$\frac{1}{\alpha} \langle e, \hat{y} - \tilde{y} \rangle \leq \frac{1}{\alpha} ||e|| ||\hat{y} - \tilde{y}|| \leq \frac{\alpha_{max}\tau}{\alpha_{min}} (-h(\tilde{y}; x, s)).$$

Finally, using (4.12), we can also write

$$(4.25) \quad \frac{\beta}{\alpha} \langle \hat{y} - \tilde{y}, x - s \rangle \ge -\frac{\beta}{2\alpha} (\|\hat{y} - \tilde{y}\|^2 + \|x - s\|^2) \ge -\frac{\beta_{max}}{2\alpha_{min}} \left( -\alpha_{max} \tau h(\tilde{y}; x, s) + \|x - s\|^2 \right).$$

Combining (4.21) with (4.13), (4.22), (4.23), (4.24), (4.25) and (4.10), gives (4.16) with

$$c = L\alpha_{max}\sqrt{2\tau\left(1+\frac{\tau}{2}\right)} + \frac{\alpha_{max}}{\alpha_{min}}\sqrt{\frac{2\tau}{\theta}} + \frac{\alpha_{max}\tau}{\alpha_{min}} + \frac{L\alpha_{max}\tau}{2} + \frac{\tau}{2} + \frac{\beta_{max}\alpha_{max}\tau}{2\alpha_{min}}, \quad d = \frac{\beta_{max}}{2\alpha_{min}}.$$

As for (4.17), using the Descent Lemma we obtain

$$f_0(y) \le f_0(x) + \langle \nabla f_0(x), y - x \rangle + \frac{L}{2} ||y - x||^2,$$

for all  $x, y \in \text{dom}(f_1)$ . Summing  $f_1(y)$  on both sides yields

$$\begin{split} f(y) &\leq f(x) + f_1(y) - f_1(x) + \langle \nabla f_0(x), y - x \rangle + \frac{L}{2} \|y - x\|^2 \\ &\leq f(x) + h(y; x, s) + \frac{L}{2} \|y - x\|^2 + \frac{\beta}{\alpha} \langle x - s, y - x \rangle \\ &\leq f(x) + h(y; x, s) + \frac{L}{2} \|y - x\|^2 + \frac{\beta}{2\alpha} (\|x - s\|^2 + \|y - x\|^2). \end{split}$$

From the previous inequality with  $y = \hat{y}$ , recalling that  $h(\hat{y}; x, s) \leq 0$  and combining with (4.11) yields (4.17), where the constants are set as  $\bar{c} = (L\alpha_{max} + \beta_{max})(1 + \tau/2), \ \bar{d} = \beta_{max}/(2\alpha_{min}).$ 

4.2. i<sup>2</sup>Piano: inertial inexact Proximal algorithm for nonconvex optimization. In this section we propose a generalization of the inertial method in [30], introducing the possibility of an inexact computation of the inertial proximal gradient point. We will refer to the new algorithm as i<sup>2</sup>Piano (inertial inexact Proximal algorithm for nonconvex optimization). Let us describe the i<sup>2</sup>Piano iteration as follows. STEP 1–4 determine the stepsize  $\alpha_k$  and inertial parameter  $\beta_k$  at iteration k. Given the parameters  $\alpha_k, \beta_k$ , STEP 5 seeks to find a possibly inexact inertial proximal point, i.e., an inexact minimizer of the function

(4.26) 
$$h^{(k)}(y;x,s) = f_1(y) - f_1(x) + \langle \nabla f_0(x) - \frac{\beta_k}{\alpha_k}(x-s), y-x \rangle + \frac{1}{2\alpha_k} \|y-x\|^2.$$

According to (4.8)–(4.10), the i<sup>2</sup>Piano iterate is any point  $x^{(k+1)} = \tilde{y}^{(k)}$  such that

(4.27) 
$$0 \in \partial_{\epsilon_k} h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)}), \quad \text{with } \epsilon_k = -\frac{\tau}{2} h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)})$$

for some fixed nonnegative constant  $\tau$ . When  $\tau = 0$ , we recover the exact inertial proximal gradient point provided by the iPiano method [30, 31]. As explained at the beginning of Section 4, STEP 5 of i<sup>2</sup>Piano can be practically implemented with an inner loop consisting of an iterative optimization method applied to the dual of problem  $\min_{y \in \mathbb{R}^n} h^{(k)}(y; x^{(k)}, x^{(k-1)})$ , until the duality gap is smaller than  $\epsilon_k$ . In the implementation of i<sup>2</sup>Piano, besides the stepsize  $\alpha_k$  and the inertial parameter  $\beta_k$ , a further parameter,  $L_k$ , is introduced (cf. STEP 6) with the aim of estimating a local Lipschitz constant of  $\nabla f_0$  that allows us to take larger steps. In particular,  $L_k$  is successively increased by a factor  $\eta > 1$  until the following descent condition holds

(4.28) 
$$f_0(x^{(k+1)}) \le f_0(x^{(k)}) + \langle \nabla f_0(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle + \frac{L_k}{2} \|x^{(k+1)} - x^{(k)}\|^2.$$

i<sup>2</sup>Piano inertial inexact Proximal algorithm for nonconvex optimization.

$$\begin{split} \overline{\text{Choose } x^{(-1)}, x^{(0)} \in \text{dom}(f_1), \delta \geq \gamma > 0, \eta > 1, 0 < L_{min} \leq L_{max}, \tau \geq 0. \text{ Set } \theta = 2/(\sqrt{2+\tau} + \sqrt{\tau})^2 \text{ and choose } \omega \in [0, 1) \text{ if } \tau > 0, \omega \in [0, 1] \text{ if } \tau = 0. \\ \text{FOR } k = 0, 1, \dots \\ \text{STEP 1. Choose } L_k \in [L_{min}, L_{max}]. \\ \text{STEP 2. Set } b_k = \frac{L_k + 2\delta}{L_k + 2\gamma}. \\ \text{STEP 3. Set } \beta_k = \frac{1 + \theta\omega}{2} \cdot \frac{b_k - 1}{b_k - \frac{1}{2}}. \\ \text{STEP 4. Set } \alpha_k = \frac{1 + \theta\omega - 2\beta_k}{L_k + 2\gamma}. \\ \text{STEP 5. Compute } \tilde{y}^{(k)} \text{ such that} \\ 0 \in \partial_{\epsilon_k} h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, x^{(k-1)}), \quad \text{with } \epsilon_k = -\frac{\tau}{2} h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, x^{(k-1)}) \\ \text{STEP 6. Check the local descent:} \\ \text{IF } f_0(\tilde{y}^{(k)}) \leq f_0(x^{(k)}) + \langle \nabla f_0(x^{(k)}), \tilde{y}^{(k)} - x^{(k)} \rangle + \frac{L_k}{2} \|\tilde{y}^{(k)} - x^{(k)}\|^2 \\ - \text{ Set } x^{(k+1)} = \tilde{y}^{(k)}. \\ \text{ELSE} \\ - \text{ Set } L_k = \eta L_k. \\ - \text{ Go to STEP 2.} \end{split}$$

**4.2.1. Convergence analysis.** Our aim now is to frame i<sup>2</sup>Piano in the abstract scheme defined by Conditions 3.1 to enjoy the favourable convergence guarantees that are provided by Theorem 3.3. We start the convergence analysis by showing that i<sup>2</sup>Piano is well-posed and that its parameters satisfy some useful relations.

LEMMA 4.4. The loop between STEP 2 and STEP 6 terminates in a finite number of steps. In particular, there exists  $\overline{L} > 0$  such that  $L_k \leq \overline{L}, \forall k \geq 0$ . Moreover, we have

(4.29) 
$$0 \le \beta_k \le \frac{1+\theta\omega}{2}, \quad \forall \ k \ge 0$$

and there exist two positive constants  $\alpha_{min}, \alpha_{max}$  with  $0 < \alpha_{min} \leq \alpha_{max}$  such that  $\alpha_k \in [\alpha_{min}, \alpha_{max}], \forall k \geq 0$ . We also have

(4.30) 
$$\frac{1+\theta\omega}{2\alpha_k} - \frac{L_k}{2} - \frac{\beta_k}{2\alpha_k} = \delta$$

(4.31) 
$$\delta - \frac{\beta_k}{2\alpha_k} = \gamma.$$

Proof. Since  $\eta > 1$ , after a finite number of steps the tentative value of  $L_k$  satisfies  $L_k \ge L$ , where L is the Lipschitz constant of  $\nabla f_0$ . Then, from the Descent Lemma, the inequality at STEP 6 is satisfied. From  $\delta \ge \gamma$  we have  $b_k \ge 1$ , which implies (4.29). A simple inspection shows that the following equalities hold:

$$b_k = \frac{1 + \theta\omega - \beta_k}{1 + \theta\omega - 2\beta_k} \Rightarrow \frac{L_k + 2\delta}{L_k + 2\gamma} = \frac{1 + \theta\omega - \beta_k}{1 + \theta\omega - 2\beta_k}$$

which leads to rewriting the parameter  $\alpha_k$  as

(4.32) 
$$\alpha_k = \frac{1 + \theta\omega - \beta_k}{L_k + 2\delta}$$

Then there holds  $\alpha_k \geq \alpha_{min}$  with  $\alpha_{min} = (1 + \theta\omega)/(2(\overline{L} + 2\delta))$  and, since  $\theta\omega \leq 1$ , we also have  $\alpha_k \leq \alpha_{max}$  with  $\alpha_{max} = 2/L_{min}$ . Moreover, we have

$$\frac{1+\theta\omega}{\alpha_k} - \frac{L_k}{2} - \frac{\beta_k}{2\alpha_k} = \frac{1+\theta\omega - \beta_k}{2\alpha_k} - \frac{L_k}{2} = \frac{L_k + 2\delta}{2} - \frac{L_k}{2} = \delta$$

and

14

$$\delta - \frac{\beta_k}{2\alpha_k} = \frac{1+\theta\omega}{2\alpha_k} - \frac{L_k}{2} - \frac{\beta_k}{\alpha_k} = \frac{1+\theta\omega-2\beta_k}{2\alpha_k} - \frac{L_k}{2} = \frac{L_k+2\gamma}{2} - \frac{L_k}{2} = \gamma.$$

Notice that the case  $\tau = 0$ , which corresponds to the exact computation of the inertial proximal gradient point at STEP 5, implies  $\theta = 1$  and, choosing  $\omega = 1$ , the parameters settings in i<sup>2</sup>Piano are exactly the same as in [31]. The need of introducing the parameter  $\omega$  is mainly technical and will be explained in the following.

We now prove that condition [H1] holds for i<sup>2</sup>Piano when the corresponding surrogate function  $\Phi$  is defined as follows:

(4.33) 
$$\Phi: \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}, \quad \Phi(x,s) = f(x) + \delta ||x - s||^2.$$

**PROPOSITION 4.5.** Let  $\{x^{(k)}\}_{k\in\mathbb{N}}$  be the sequence generated by  $i^2 Piano$ . Then there holds

$$(4.34) \quad \Phi(x^{(k+1)}, x^{(k)}) \le \Phi(x^{(k)}, x^{(k-1)}) - \gamma \|x^{(k)} - x^{(k-1)}\|^2 + (1-\omega)h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)}).$$

Consequently, there exist  $\{s^{(k)}\}_{k\in\mathbb{N}}, \{d_k\}_{k\in\mathbb{N}}, \{a_k\}_{k\in\mathbb{N}}$  such that [H1] holds with  $\Phi$  as in (4.33).

*Proof.* By summing the quantity  $f_1(x^{(k+1)})$  to both sides of inequality (4.28) we obtain

$$\begin{split} f(x^{(k+1)}) &\leq f(x^{(k)}) + f_1(x^{(k+1)}) - f_1(x^{(k)}) + \langle \nabla f_0(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle + \frac{L_k}{2} \|x^{(k+1)} - x^{(k)}\|^2 \\ &= f(x^{(k)}) + h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)}) - \left(\frac{1}{2\alpha_k} - \frac{L_k}{2}\right) \|x^{(k+1)} - x^{(k)}\|^2 \\ &+ \frac{\beta_k}{\alpha_k} \langle x^{(k+1)} - x^{(k)}, x^{(k)} - x^{(k-1)} \rangle \\ &\leq f(x^{(k)}) + h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)}) - \left(\frac{1}{2\alpha_k} - \frac{L_k}{2}\right) \|x^{(k+1)} - x^{(k)}\|^2 \\ &+ \frac{\beta_k}{2\alpha_k} (\|x^{(k+1)} - x^{(k)}\|^2 + \|x^{(k)} - x^{(k-1)}\|^2) \\ &\leq f(x^{(k)}) + (1 - \omega)h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)}) - \frac{\theta\omega}{2\alpha_k} \|x^{(k+1)} - x^{(k)}\|^2 \\ &- \left(\frac{1}{2\alpha_k} - \frac{L_k}{2}\right) \|x^{(k+1)} - x^{(k)}\|^2 + \frac{\beta_k}{2\alpha_k} (\|x^{(k+1)} - x^{(k)}\|^2 + \|x^{(k)} - x^{(k-1)}\|^2) \\ &= f(x^{(k)}) - \left(\frac{1 + \theta\omega}{2\alpha_k} - \frac{L_k}{2} - \frac{\beta_k}{2\alpha_k}\right) \|x^{(k+1)} - x^{(k)}\|^2 \\ &+ (1 - \omega)h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)}) + \frac{\beta_k}{2\alpha_k} \|x^{(k)} - x^{(k-1)}\|^2 \end{split}$$

where the first equality is obtained by adding and subtracting to the right-hand-side the quantity  $||x^{(k+1)} - x^{(k)}||^2/(2\alpha_k) + \beta_k/\alpha_k \langle x^{(k+1)} - x^{(k)}, x^{(k)} - x^{(k-1)} \rangle$ , the subsequent inequality follows from the basic relation  $2\langle a, b \rangle \leq ||a||^2 + ||b||^2$  and the next one from (4.13). Recalling (4.30), the above inequality can be conveniently rewritten as

$$\begin{aligned} f(x^{(k+1)}) + \delta \|x^{(k+1)} - x^{(k)}\|^2 &\leq f(x^{(k)}) + \delta \|x^{(k)} - x^{(k-1)}\|^2 \\ &+ \left(\frac{\beta_k}{2\alpha_k} - \delta\right) \|x^{(k)} - x^{(k-1)}\|^2 + (1-\omega)h(x^{(k+1)}; x^{(k)}, x^{(k-1)}). \end{aligned}$$

Finally, exploiting (4.31), we obtain condition [H1] with  $\Phi$  given in (4.33),  $a_k = 1$  and

$$(4.35) s^{(k)} = x^{(k-1)}$$

(4.36) 
$$d_k^2 = \gamma \|x^{(k)} - x^{(k-1)}\|^2 - (1-\omega)h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)}).$$

Under Assumption  $[A_4]$ ,  $\Phi$  is bounded from below, hence, condition [H1] implies

(4.37) 
$$\lim_{k \to \infty} \|x^{(k)} - x^{(k-1)}\| = 0$$

and, if  $\omega < 1$ , also

(4.38) 
$$\lim_{k \to \infty} h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)}) = 0.$$

The choice  $\omega < 1$  is enforced when  $\tau > 0$  in order to obtain (4.38). If we take  $\omega = 1$ , with the same arguments as above we still obtain (4.37). It is also worth noticing that, for large values of  $\tau$ , that is when a coarser accuracy is allowed in the computation of  $\tilde{y}^{(k)}$ , the parameter  $\theta$  can be very small and this also influences the choice of  $\alpha_k$  and  $\beta_k$ .

In order to prove condition [H2], let us now introduce the second surrogate function as follows

(4.39) 
$$\mathcal{F}: \mathbb{R}^n \times \mathbb{R} \to \mathbb{R}, \quad \mathcal{F}(u,\rho) = f(u) + \frac{1}{2}\rho^2$$

The following result is proved using similar arguments as the ones used in Lemma 4-5 in [16].

PROPOSITION 4.6. Let  $\{x^{(k)}\}_{k\in\mathbb{N}}$  be the sequence generated by  $i^2$  Piano and, for each k, let the point  $\hat{x}^{(k+1)}$  be defined as

(4.40) 
$$\hat{x}^{(k+1)} = \operatorname*{argmin}_{y \in \mathbb{R}^n} h^{(k)}(y; x^{(k)}, x^{(k-1)}).$$

Then, there exist  $\{\rho^{(k)}\}_{k\in\mathbb{N}}$ ,  $\{r_k\}_{k\in\mathbb{N}}$  such that condition [H2] holds with  $\Phi$  defined in (4.33),  $\mathcal{F}$  defined in (4.39),  $s^{(k)}$  defined in (4.35) and  $u^{(k)} = \hat{x}^{(k+1)}$ .

*Proof.* When  $\omega = 1$  we necessarily have  $\tau = 0$ , which means  $x^{(k+1)} = \hat{x}^{(k+1)}$ . Therefore, [H2] directly follows from [H1] with  $u^{(k)} = x^{(k+1)}$ ,  $\rho^{(k)} = \sqrt{2\delta} ||x^{(k)} - x^{(k-1)}||$ ,  $r_k = 0$ . Consider now the case  $\omega < 1$ . From (4.16) and (4.17), we directly obtain

$$\begin{aligned} f(\hat{x}^{(k+1)}) &\geq f(x^{(k+1)}) + ch^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)}) - d \|x^{(k)} - x^{(k-1)}\|^2 \\ f(\hat{x}^{(k+1)}) &\leq f(x^{(k)}) - \bar{c}h(x^{(k+1)}; x^{(k)}, x^{(k-1)}) + \bar{d} \|x^{(k)} - x^{(k-1)}\|^2, \end{aligned}$$

where  $c, d, \bar{c}, \bar{d}$  are defined as in Lemma 4.3 and do not depend on k. Combining the two inequalities above we obtain

$$\begin{aligned} f(x^{(k+1)}) + \delta \|x^{(k+1)} - x^{(k)}\|^2 &\leq \\ &\leq f(\hat{x}^{(k+1)}) + \delta \|x^{(k+1)} - x^{(k)}\|^2 + c(-h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)})) + d\|x^{(k)} - x^{(k-1)}\|^2 \\ &\leq f(x^{(k)}) - (c + \bar{c})h(x^{(k+1)}; x^{(k)}, x^{(k-1)}) + (d + \bar{d})\|x^{(k)} - x^{(k-1)}\|^2 + \delta \|x^{(k+1)} - x^{(k)}\|^2. \end{aligned}$$

Recalling the definition of  $\mathcal{F}$  in (4.39), the above inequalities can be rewritten as

(4.41) 
$$f(x^{(k+1)}) + \delta \|x^{(k+1)} - x^{(k)}\|^2 \le \mathcal{F}(\hat{x}^{(k+1)}, \rho^{(k)}) \le f(x^{(k)}) + \delta \|x^{(k)} - x^{(k-1)}\|^2 + r_k$$

where  $\rho^{(k)}, r_k$  are given by

$$(4.42) \quad \rho^{(k)} = \sqrt{2}(\delta \| x^{(k+1)} - x^{(k)} \|^2 + c(-h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)})) + d \| x^{(k)} - x^{(k-1)} \|^2)^{\frac{1}{2}}$$
$$r_k = -(c + \bar{c})h(x^{(k+1)}; x^{(k)}, x^{(k-1)}) + (d + \bar{d}) \| x^{(k)} - x^{(k-1)} \|^2 + \delta \| x^{(k+1)} - x^{(k)} \|^2.$$

From (4.37)–(4.38) we obtain  $\lim_{k\to\infty} r_k = 0$ , hence, assumption [H2] is satisfied with the above settings and with  $u^{(k)} = \hat{x}^{(k+1)}$ .

Next we show that condition [H3] holds for i<sup>2</sup>Piano. The following result combines elements of Lemma 3 in [16] and Lemma 17 in [30].

PROPOSITION 4.7. There exist b > 0,  $I \subset \mathbb{Z}$ ,  $\{\theta_i\}_{i \in I}$  such that condition [H3] holds with  $\{u^{(k)}\}_{k \in \mathbb{N}}, \{\rho^{(k)}\}_{k \in \mathbb{N}}, \{d_k\}_{k \in \mathbb{N}}$  defined as in Proposition 4.6 and in (4.36),  $\zeta_{k+1} = 0$  and  $b_{k+1} = 1$ .

*Proof.* From the separable structure of  $\mathcal{F}$ , if  $v \in \partial f(u)$  and  $\rho \in \mathbb{R}$ , we have that  $(v, \rho) \in \mathcal{F}$  $\partial \mathcal{F}(u,\rho)$ . In particular,

(4.43) 
$$\|\partial \mathcal{F}(u,\rho)\|_{-} \le \|v\| + |\rho|, \text{ for all } v \in \partial f(u), \rho \in \mathbb{R}.$$

If  $\omega = 1$ , we are in the case  $\tau = 0$ . This means that  $x^{(k+1)} = \hat{x}^{(k+1)}$ ,  $\rho^{(k)} = \sqrt{2\delta} \|x^{(k)} - x^{(k-1)}\|$ . Moreover, (4.36) reduces to  $d_k = \sqrt{\gamma} \|x^{(k)} - x^{(k-1)}\| = \sqrt{\gamma/(2\delta)}\rho^{(k)}$ . From (4.14), we have that there exists a subgradient  $\hat{v}^{(k+1)} \in \partial f(\hat{x}^{(k+1)})$  and a positive constant p such that

$$\|\hat{v}^{(k+1)}\| \le p(\|x^{(k+1)} - x^{(k)}\| + \|x^{(k)} - x^{(k-1)}\|) = \frac{p}{\sqrt{\gamma}}(d_{k+1} + d_k).$$

Hence,  $\|\hat{v}^{(k+1)}\| + \rho^{(k)} \leq pd_{k+1}/\sqrt{\gamma} + d_k(p + \sqrt{2\delta})/\sqrt{\gamma}$  and, recalling (4.43), [H3] follows with  $b = (p + \sqrt{2\delta})/\sqrt{\gamma}, I = \{0, 1\}, \theta_0 = \theta_1 = \frac{1}{2}.$ 

Consider now the case  $\omega < 1$ . From (4.15), there exists a subgradient  $\hat{v}^{(k+1)} \in \partial f(\hat{x}^{(k+1)})$  and a positive constant q such that

(4

$$+q\sqrt{\frac{2\alpha_{max}}{\theta(1-\omega)}}\sqrt{-(1-\omega)h^{(k-1)}(x^{(k)};x^{(k-1)},x^{(k-2)})} + \gamma \|x^{(k-1)} - x^{(k-1)}\| + \gamma \|x^{(k-1)}\| + \gamma$$

(4.45) 
$$\leq \frac{q}{\sqrt{1-\omega}}d_k + q\sqrt{\frac{2\alpha_{max}}{\theta(1-\omega)}}d_{k-1}.$$

From (4.42) and (4.13) we have

$$(4.46) \ \rho^{(k)} \le \sqrt{2} \left( -(2\delta\alpha_{max}/\theta + c)h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)}) + d\|x^{(k)} - x^{(k-1)}\|^2 \right)^{\frac{1}{2}} \le r \cdot d_k$$

where  $r = \sqrt{2} \max \left\{ \frac{2\delta \alpha_{max}}{\theta} + c \right\} / (1-\omega), \frac{d}{\gamma} \right\}^{\frac{1}{2}}$ . Then, combining (4.45) with (4.46), in view of (4.43) we obtain

$$\|\partial \mathcal{F}(\hat{x}^{(k+1)}, \rho^{(k)})\|_{-} \le \frac{b}{2}(d_k + d_{k-1}), \quad \text{with } b = 2\max\left\{\frac{q}{\sqrt{1-\omega}} + r, q\sqrt{\frac{2\alpha_{max}}{\theta(1-\omega)}}\right\}$$

and the thesis follows with  $I = \{1, 2\}, \theta_1 = \theta_2 = \frac{1}{2}$ .

The following lemma holds for all methods whose iterates satisfy [H1], [H2], [H3] with  $\mathcal{F}$  defined as in (4.39), and it is crucial to ensure condition  $[H_4]$  for i<sup>2</sup>Piano.

PROPOSITION 4.8. Let Assumptions [A1]-[A4] be satisfied and assume that  $\{x^{(k)}\}_{k\in\mathbb{N}}$  satisfies [H1], [H2], [H3] with  $\mathcal{F}$  defined as in (4.39). If  $\{(x^{(k_j)}, \rho^{(k_j)})\}_{j\in\mathbb{N}}$  is a subsequence of  $\{(x^{(k)}, \rho^{(k)})\}_{k\in\mathbb{N}}$  and  $[x^{(k)}, \rho^{(k)}] \in \mathbb{R}^n \times \mathbb{R}^m$  such that  $\lim_{j\to\infty} ||u^{(k_j)} - x^{(k_j)}|| = 0$ , then  $\rho^* = 0$  and  $[x^{(k)}, \rho^{(k)}] \in \mathbb{R}^n \times \mathbb{R}^n$  such that  $\lim_{j\to\infty} ||u^{(k_j)} - x^{(k_j)}|| = 0$ , then  $\rho^* = 0$  and  $[x^{(k)}, \rho^{(k)}] \in \mathbb{R}^n$ .  $\lim_{j \to \infty} \mathcal{F}(u^{(k_j)}, \rho^{(k_j)}) = \mathcal{F}(x^*, \rho^*).$ 

*Proof.* Recalling that [H1] implies  $\lim_{k\to\infty} d_k = 0$ , from [H3] and thanks to the separable structure of  $\mathcal{F}$ , we obtain that there exists  $\hat{v}^{(k)} \in \partial f(u^{(k)})$  such that  $\lim_{k \to \infty} \|\hat{v}^{(k)}\| = \lim_{k \to \infty} \rho^{(k)} = 0$ . In particular, in view of (4.2), we can write  $\hat{v}^{(k)} = \nabla f_0(u^{(k)}) + w^{(k)}$ , where  $w^{(k)} \in \partial f_1(u^{(k)})$ . Therefore, by continuity of  $\nabla f_0$ , the following implication holds

$$\lim_{k \to \infty} \hat{v}^{(k)} = 0 \Rightarrow \lim_{j \to \infty} w^{(k_j)} = -\nabla f_0(x^*).$$

Adding the quantity  $\frac{1}{2}(\rho^{(k_j)})^2$  to both sides of the subgradient inequality yields

$$f_1(x^*) + \frac{1}{2}(\rho^{(k_j)})^2 \ge f_1(u^{(k_j)}) + \langle w^{(k_j)}, x^* - u^{(k_j)} \rangle + \frac{1}{2}(\rho^{(k_j)})^2$$
$$= \mathcal{F}(u^{(k_j)}, \rho^{(k_j)}) - f_0(u^{(k_j)}) + \langle w^{(k_j)}, x^* - u^{(k_j)} \rangle$$

where the last equality is obtained by adding and subtracting  $f_0(u^{(k_j)})$  to the right-hand-side. Taking limits on both sides we obtain

$$f_1(x^*) + \frac{1}{2}(\rho^*)^2 \ge \lim_{j \to \infty} \mathcal{F}(u^{(k_j)}, \rho^{(k_j)}) - f_0(x^*),$$

which, rearranging terms, gives  $\lim_{j\to\infty} \mathcal{F}(u^{(k_j)}, \rho^{(k_j)}) \leq \mathcal{F}(x^*, \rho^*)$ . On the other side, by assumption,  $\mathcal{F}$  is lower semicontinuous, therefore  $\lim_{j\to\infty} \mathcal{F}(u^{(k_j)}, \rho^{(k_j)}) \geq \mathcal{F}(x^*, \rho^*)$ , which completes the proof.

We are now ready to prove that i<sup>2</sup>Piano complies with Conditions 3.1 and, hence, its iterates converge to a stationary point.

THEOREM 4.9. Suppose that  $\mathcal{F}$  is a KL function and assume that the sequence  $\{x^{(k)}\}_{k\in\mathbb{N}}$  generated by  $i^2$  Piano is bounded. Then,  $\{x^{(k)}\}_{k\in\mathbb{N}}$  converges to a stationary point of f.

*Proof.* By Propositions 4.5-4.6-4.7, we know that conditions [H1]-[H2]-[H3] hold for i<sup>2</sup>Piano. Furthermore, if  $\omega = 1$ , then  $u^{(k)} = \hat{x}^{(k+1)} = x^{(k+1)}$ , and (4.37) directly implies that  $\lim_{k\to\infty} ||u^{(k)} - x^{(k)}|| = 0$ . If  $\omega < 1$ , using (4.11) and (4.38) we have

$$\lim_{k \to \infty} \|u^{(k)} - x^{(k)}\| = \lim_{k \to \infty} \|\hat{x}^{(k+1)} - x^{(k)}\| \le \lim_{k \to \infty} \sqrt{-2\alpha_{max} \left(1 + \frac{\tau}{2}\right) h^{(k)}(x^{(k+1)}; x^{(k)}, x^{(k-1)})} = 0.$$

Then, in both cases, the assumptions of Proposition 4.8 are satisfied and, therefore, condition [H4] holds. Finally, condition [H5] follows from (4.36), while condition [H6] is trivially satisfied, since both sequences  $\{a_k\}_{k\in\mathbb{N}}$  and  $\{b_k\}_{k\in\mathbb{N}}$  are constant. Then, Theorem 3.3 applies and guarantees that the sequence  $\{(x^{(k)}, \rho^{(k)})\}_{k\in\mathbb{N}}$  converges to a stationary point  $(x^*, \rho^*)$  of  $\mathcal{F}$ . Note that, since  $\mathcal{F}$  is the sum of separable functions, its subdifferential can be written as  $\partial \mathcal{F}(x, \rho) = \partial f(x) \times \{\rho\}$ . Then,  $(x^*, \rho^*)$  is stationary for  $\mathcal{F}$  if and only if  $\rho^* = 0$  and  $0 \in \partial f(x^*)$ . Hence,  $x^*$  is a stationary point for f and  $\{x^{(k)}\}_{k\in\mathbb{N}}$  converges to it.

Remark 4.10. We underline that  $\mathcal{F}$  is a KL function if, for instance, f and  $\frac{1}{2} \|\cdot\|^2$  are definable in the same o-minimal structure [7, Definition 7]. Indeed functions definable in an o-minimal structure satisfy the KL property on their domain [7, Theorem 11] and o-minimal structures are closed with respect to the sum, see [7, Remark 5] and references therein. Examples of functions definable in an o-minimal structure are semialgebraic, subanalytic and real analytic functions.

Remark 4.11. Theorem 4.9 requires the boundedness of the iterates as hypothesis. A standard way to assert such a condition is when the Lyapunov function  $\Phi$  defined in (4.33) is coercive, since this assumption combined with the descent property (4.34) guarantees that the sequence  $\{(x^{(k)}, x^{(k-1)})\}_{k\in\mathbb{N}}$  is included in a (bounded) level set of the coercive function  $\Phi$ .

4.3. iPila: inertial Proximal inexact line-search algorithm. In the following we introduce a novel algorithm combining a line-search along the descent direction and an inertial proximal-gradient step as a special case of our abstract scheme. A line-search procedure for the objective function f requires a descent direction  $d \in \mathbb{R}^n$ , i.e., a vector such that the directional derivative  $f'(x; d) = \lim_{\lambda \downarrow 0} (f(x+\lambda) - f(x))/\lambda$  is negative. As explained extensively in [11, 12], the inexact proximal-gradient point provides a descent direction for the objective function f. Indeed, if  $\tilde{y}$  satisfies (4.8), where the inertial parameter  $\beta$  in (4.5) is equal to zero, then the vector  $\tilde{y} - x$  is a descent direction for f at x. Unfortunately, this is not true, in general, when  $\beta > 0$ . In this section we show that in the general case  $\beta \geq 0$ , the point  $\tilde{y}$  can be still used to define a descent direction for a suitable merit function. Then, we propose a line-search procedure along this direction that enable us to define a descent algorithm such that the merit function monotonically

decreases along the iterates. Differently from the backtracking procedure in i<sup>2</sup>Piano, the proposed line–search requires to solve the minimization subproblem (4.4) only once per iteration. Finally, we show that the new algorithm can be analyzed in the framework of Section 3, in order to prove the convergence of the iterates to a stationary point of f. In this case, unlike i<sup>2</sup>Piano, one of the two merit functions involved in the abstract scheme defined by Conditions 3.1 will play an active role in the algorithm, as it will be explicitly computed at each iteration to determine the new point. In particular, we define the merit function  $\Phi$  appearing in |H1|-|H2| as follows:

$$\Phi: \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R} \cup \{+\infty\}, \quad \Phi(x,s) = f(x) + \frac{1}{2} \|x - s\|^2,$$

where the variable s will be considered as an actual optimization variable, independent on x. The composite structure of  $\Phi$  is  $\Phi(x,s) = \Phi_0(x,s) + \Phi_1(x,s)$ , with

$$\Phi_0(x,s) = f_0(x) + \frac{1}{2} ||x - s||^2, \quad \Phi_1(x,s) = f_1(x).$$

The function  $\Phi_0$  is differentiable with gradient

$$\nabla \Phi_0(x,s) = \begin{pmatrix} \nabla_x \Phi_0(x,s) \\ \nabla_s \Phi_0(x,s) \end{pmatrix} = \begin{pmatrix} \nabla f_0(x) + x - s \\ s - x \end{pmatrix}.$$

It is easy to see that  $\nabla \Phi_0(x,s)$  is Lipschitz continuous and in particular it holds that

(4.47) 
$$\|\nabla\Phi_0(x,s) - \nabla\Phi_0(\bar{x},\bar{s})\| \le M \left\| \begin{pmatrix} x - \bar{x} \\ s - \bar{s} \end{pmatrix} \right\|, \quad \forall x, \bar{x}, s, \bar{s},$$

with M = L + 2, where L is the Lipschitz constant of  $\nabla f_0$ . Given a vector  $d \in \mathbb{R}^{2n}$ 

(4.48) 
$$d = \begin{pmatrix} d_x \\ d_s \end{pmatrix},$$

the directional derivative of  $\Phi$  at the point (x, s) with respect to the direction d can be written as

$$\Phi'(x,s;d_x,d_s) = \Phi'_0(x,s;d_x,d_s) + \Phi'_1(x,s;d_x,d_s) = f'_1(x;d_x) + \langle \nabla \Phi_0(x,s),d \rangle$$
  
=  $f'_1(x;d_x) + \langle \nabla f_0(x) + x - s, d_x \rangle + \langle s - x, d_s \rangle,$ 

which always exists thanks to the convexity of  $f_1$ . A vector  $d \in \mathbb{R}^{2n}$  is called a descent direction for  $\Phi$  at (x, s) when

$$\Phi'(x,s;d_x,d_s) < 0.$$

Assume now that the vector  $d_x$  has the form  $d_x = y - x$ , where y is a point belonging to the domain of  $f_1$ ; then, from [34, Theorem 23.1] we have

(4.49) 
$$\Phi'(x,s;y-x,d_s) \le f_1(y) - f_1(x) + \langle \nabla f_0(x) + x - s, d_x \rangle + \langle s - x, d_s \rangle.$$

The above inequality holds independently on the form of  $d_s$ .

**4.3.1. Descent direction for the merit function.** Assume that  $(x^{(k)}, s^{(k)})$  is a given point in dom $(f_1) \times \mathbb{R}^n$ , while  $\alpha_k$ ,  $\beta_k$  are two given parameters. Let the function  $h^{(k)}(y; x, s)$  be defined as in (4.26). Given a tolerance parameter  $\tau > 0$ , according to (4.8)–(4.10), we denote by  $\tilde{y}^{(k)}$  any point in dom $(f_1)$  satisfying

(4.50) 
$$0 \in \partial_{\epsilon_k} h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}), \quad \text{with } \epsilon_k = -\frac{\tau}{2} h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}).$$

For a given  $\gamma_k \geq 0$ , consider

(4.51) 
$$d^{(k)} = \begin{pmatrix} d_x^{(k)} \\ d_s^{(k)} \end{pmatrix},$$

where

(4.52) 
$$d_x^{(k)} = \tilde{y}^{(k)} - x^{(k)}$$

(4.53) 
$$d_s^{(k)} = \left(1 + \frac{\beta_k}{\alpha_k}\right) (\tilde{y}^{(k)} - x^{(k)}) + \gamma_k (x^{(k)} - s^{(k)}).$$

In the following lemma we show that  $d^{(k)}$  is a descent direction for  $\Phi$  at  $(x^{(k)}, s^{(k)})$ .

LEMMA 4.12. Let  $d^{(k)} \in \mathbb{R}^{2n}$  be defined according to (4.52)-(4.53), where  $0 < \alpha_k \leq \alpha_{max}$ ,  $\beta_k \geq 0, \ \gamma_k \geq \gamma_{min}$ , for some given positive real constants  $\alpha_{max}, \gamma_{min} > 0$ . Define  $\Delta_k \in \mathbb{R}_{\leq 0}$  as

(4.54) 
$$\Delta_k = h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}) - \gamma_k \|x^{(k)} - s^{(k)}\|^2$$

Then, we have

(4.55) 
$$\Phi'(x^{(k)}, s^{(k)}; d_x^{(k)}, d_s^{(k)}) \le \Delta_k$$
  
(4.56) 
$$\le -a \|\tilde{y}^{(k)} - x^{(k)}\|^2 - \gamma_{min} \|x^{(k)} - s^{(k)}\|^2$$

where a > 0 is a positive constant. Therefore,  $\Phi'(x^{(k)}, s^{(k)}; d_x^{(k)}, d_s^{(k)}) < 0$  whenever  $\tilde{y}^{(k)} \neq x^{(k)}$  or  $x^{(k)} \neq s^{(k)}$ .

*Proof.* We first observe that (4.49) with  $x = x^{(k)}$ ,  $s = s^{(k)}$ ,  $y = \tilde{y}^{(k)}$ ,  $d_x = d_x^{(k)}$ ,  $d_s = d_s^{(k)}$ , gives:

$$\Phi'(x^{(k)}, s^{(k)}; d_x^{(k)}, d_s^{(k)}) 
(4.57) \leq f_1(\tilde{y}^{(k)}) - f_1(x^{(k)}) + \langle \nabla f_0(x^{(k)}) + x^{(k)} - s^{(k)}, d_x^{(k)} \rangle + \langle s^{(k)} - x^{(k)}, d_s^{(k)} \rangle 
= f_1(\tilde{y}^{(k)}) - f_1(x^{(k)}) + \langle \nabla f_0(x^{(k)}), \tilde{y}^{(k)} - x^{(k)} \rangle + \langle x^{(k)} - s^{(k)}, y^{(k)} - x^{(k)} \rangle + 
+ \left(1 + \frac{\beta_k}{\alpha_k}\right) \langle s^{(k)} - x^{(k)}, \tilde{y}^{(k)} - x^{(k)} \rangle - \gamma_k \|x^{(k)} - s^{(k)}\|^2 
= f_1(\tilde{y}^{(k)}) - f_1(x^{(k)}) + \langle \nabla f_0(x^{(k)}) - \frac{\beta_k}{\alpha_k}(x^{(k)} - s^{(k)}), \tilde{y}^{(k)} - x^{(k)} \rangle - \gamma_k \|x^{(k)} - s^{(k)}\|^2 
(4.58) \leq h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}) - \gamma_k \|x^{(k)} - s^{(k)}\|^2 
\leq -\frac{\theta}{2\alpha_k} \|\tilde{y}^{(k)} - x^{(k)}\|^2 - \gamma_k \|x^{(k)} - s^{(k)}\|^2$$

where the last inequality follows from (4.13). Then, the thesis follows from  $\alpha_k \leq \alpha_{max}$  with  $a = \theta/2\alpha_{max}$ .

A useful consequence of the previous lemma is stated in the following corollary.

COROLLARY 4.13. Assume that  $\alpha_k \in [\alpha_{min}, \alpha_{max}], \ \beta_k \in [0, \beta_{max}], \ \gamma_k \in [0, \gamma_{max}], \ with \ 0 < \alpha_{min} \le \alpha_{max}, \ \beta_{max} \ge 0, \ \gamma_{max} > 0.$  Then, there exists a positive constant C such that (4.59)  $\Delta_k \le -C \|d^{(k)}\|^2.$ 

*Proof.* Setting  $\delta_k = 1 + \frac{\beta_k}{\alpha_k}$ , the bounds on the parameters imply that  $1 \leq \delta_k \leq \overline{\delta}$ , with  $\overline{\delta} = 1 + \frac{\beta_{max}}{\alpha_{min}}$ . By definition of  $d^{(k)}$  in (4.51) we have

$$\begin{split} \|d^{(k)}\|^2 &= \|d^{(k)}_x\|^2 + \|d^{(k)}_s\|^2 \\ &= \|\tilde{y}^{(k)} - x^{(k)}\|^2 + \|\delta_k(\tilde{y}^{(k)} - x^{(k)}) + \gamma_k(x^{(k)} - s^{(k)})\|^2 \\ &= (1 + \delta^2_k)\|\tilde{y}^{(k)} - x^{(k)}\|^2 + \gamma^2_k\|x^{(k)} - s^{(k)}\|^2 + 2\delta_k\gamma_k\langle\tilde{y}^{(k)} - x^{(k)}, x^{(k)} - s^{(k)}\rangle \\ &\leq (1 + \delta^2_k + \delta_k\gamma_k)\|\tilde{y}^{(k)} - x^{(k)}\|^2 + (\gamma^2_k + \delta_k\gamma_k)\|x^{(k)} - s^{(k)}\|^2 \\ &\leq (1 + \bar{\delta}^2 + \bar{\delta}\gamma_{max})\|\tilde{y}^{(k)} - x^{(k)}\|^2 + \gamma_k(\gamma_{max} + \bar{\delta})\|x^{(k)} - s^{(k)}\|^2 \\ &\leq \frac{1}{C}(a\|\tilde{y}^{(k)} - x^{(k)}\|^2 + \gamma_k\|x^{(k)} - s^{(k)}\|^2) \end{split}$$

where  $C = 1/\max\{(1 + \bar{\delta}^2 + \bar{\delta}\gamma_{max})/a, \gamma_{max} + \bar{\delta}\}$ . Multiplying both sides of the last inequality above by C and combining with (4.56) gives (4.59).

Line–Search Armijo line–search

 $\overline{\text{INPUT: } (x^{(k)}, s^{(k)}) \in \mathbb{R}^n \times \mathbb{R}^n, d^{(k)}, \Delta_k \text{ as in } (4.54), \sigma, \delta \in (0, 1)}$   $\operatorname{Set} \lambda^+ = 1.$   $\operatorname{WHILE} \Phi(x^{(k)} + \lambda^+ d^{(k)}_x, s^{(k)} + \lambda^+ d^{(k)}_s) > \Phi(x^{(k)}, s^{(k)}) + \sigma \lambda^+ \Delta_k$   $\operatorname{Set} \lambda^+ = \delta \lambda^+$   $\operatorname{Set} \lambda^+_k = \lambda^+$   $\overline{\text{END}}$   $\operatorname{OUTPUT:} \lambda^+_k.$ 

**4.3.2.** Outline of the algorithm. In this section we present a new forward-backwardtype method, which combines the inertial approach with the line-search procedure at the basis of VMILA [12, 16]. The convergence guarantees of the new method are based on a line-search along the descent direction defined in the previous section, ensuring the sufficient decrease of the merit function  $\Phi(x, s)$ . The line-search consists in a backtracking procedure with a generalized Armijo inequality as stopping rule. In particular, if the vector  $d^{(k)}$  in (4.51) is a descent direction for  $\Phi$  at  $(x^{(k)}, s^{(k)})$ , the line-search algorithm computes a positive parameter  $\lambda_k^+$  satisfying

(4.60) 
$$\Phi(x^{(k)} + \lambda_k^+ d_x^{(k)}, s^{(k)} + \lambda_k^+ d_s^{(k)}) \le \Phi(x^{(k)}, s^{(k)}) + \sigma \lambda_k^+ \Delta_k, \quad \sigma \in (0, 1).$$

The well posedness of the line–search procedure and its main properties are summarized in the following lemma.

LEMMA 4.14. Let the assumptions of Lemma 4.12 be satisfied. Then, the Armijo backtracking line-search algorithm terminates in a finite number of steps, and there exists  $\lambda_{min} > 0$  such that the parameter  $\lambda_k^+$  computed with the line-search algorithm satisfies

(4.61) 
$$\lambda_k^+ \ge \lambda_{min}$$

*Proof.* Since  $\Phi_0$  has *M*-Lipschitz continuous gradient, with M = L + 2 (see (4.47)), we can apply the Descent Lemma obtaining

$$\Phi_{0}(x^{(k)} + \lambda d_{x}^{(k)}, s^{(k)} + \lambda d_{s}^{(k)}) \leq \Phi_{0}(x^{(k)}, s^{(k)}) + \lambda \langle \nabla \Phi_{0}(x^{(k)}, s^{(k)}), d^{(k)} \rangle + \frac{M}{2} \lambda^{2} \|d^{(k)}\|^{2} \\
\leq \Phi_{0}(x^{(k)}, s^{(k)}) + \lambda \langle \nabla \Phi_{0}(x^{(k)}, s^{(k)}), d^{(k)} \rangle - \frac{M}{2C} \lambda^{2} \Delta_{k},$$
(4.62)

where the second inequality follows from (4.59). From the Jensen's inequality applied to the convex function  $f_1$  we also obtain

(4.63) 
$$\Phi_1(x^{(k)} + \lambda d_x^{(k)}, s^{(k)} + \lambda d_s^{(k)}) = f_1(x^{(k)} + \lambda d_x^{(k)}) = f_1(\lambda \tilde{y}^{(k)} + (1 - \lambda)x^{(k)})$$
$$\leq (1 - \lambda)f_1(x^{(k)}) + \lambda f_1(\tilde{y}^{(k)}).$$

Summing (4.62) with (4.63) gives

$$\begin{split} \Phi(x^{(k)} + \lambda d_x^{(k)}, s^{(k)} + \lambda d_s^{(k)}) &\leq \\ &\leq \Phi(x^{(k)}, s^{(k)}) + \lambda \left( f_1(\tilde{y}^{(k)}) - f_1(x^{(k)}) + \langle \nabla \Phi_0(x^{(k)}, s^{(k)}), d^{(k)} \rangle \right) - \frac{M}{2C} \lambda^2 \Delta_k \\ &\leq \Phi(x^{(k)}, s^{(k)}) + \lambda \Delta_k - \frac{M}{2C} \lambda^2 \Delta_k, \end{split}$$

where the last inequality follows from (4.57)-(4.58). The above relation implies

$$\Phi(x^{(k)} + \lambda d_x^{(k)}, s^{(k)} + \lambda d_s^{(k)}) \le \Phi(x^{(k)}, s^{(k)}) + \lambda(1 - \rho\lambda)\Delta_k, \quad \forall \lambda \in [0, 1]$$

with  $\rho = \frac{M}{2C}$ . Moreover, comparing the above inequality with the Armijo condition (4.60) shows that the last one is surely fulfilled when  $\lambda_k^+$  satisfies  $1 - \rho \lambda_k^+ \ge \sigma$ , that is when  $\lambda_k^+ \le (1 - \sigma)/\rho$ . Since  $\lambda_k^+$  in the backtracking procedure is obtained starting from 1 and by successive reductions

iPila : inertial Proximal inexact line-search algorithm

 $\overline{\text{INPUT: } (x^{(0)}, s^{(0)}) \in \text{dom}(f_1) \times \mathbb{R}^n, \sigma \in (0, 1), 0 < \alpha_{\min} \le \alpha_{\max}, \beta_{\max} > 0, 0 < \gamma_{\min} \le \gamma_{\max}, \beta_{\max} < 0, 0 < \gamma_{\max} < 0, 0 < \gamma_{\max} < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 < 0, 0 <$  $\tau \ge 0.$ FOR k = 0, 1, ...

STEP 1. Choose  $\alpha_k \in [\alpha_{min}, \alpha_{max}], \beta_k \in [0, \beta_{max}]$ STEP 2. Compute  $\tilde{y}^{(k)}$  such that

 $0 \in \partial_{\epsilon_k} h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}), \quad \text{with } \epsilon_k = -\frac{\tau}{2} h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}).$ 

STEP 3. Choose  $\gamma_k \in [\gamma_{min}, \gamma_{max}]$ . STEP 4. Compute  $\Delta_k = h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}) - \gamma_k ||x^{(k)} - s^{(k)}||^2$ .

STEP 5. Compute the search direction

$$d_x^{(k)} = \tilde{y}^{(k)} - x^{(k)}$$
  
$$d_s^{(k)} = \left(1 + \frac{\beta_k}{\alpha_k}\right) (\tilde{y}^{(k)} - x^{(k)}) + \gamma_k (x^{(k)} - s^{(k)})$$

STEP 6. Compute  $\lambda_k \in (0, 1]$  such that

$$\Phi(x^{(k)} + \lambda_k d_x^{(k)}, s^{(k)} + \lambda_k d_s^{(k)}) \le \Phi(x^{(k)}, s^{(k)}) + \sigma \lambda_k \Delta_k$$

with the line–search backtracking algorithm. STEP 7. Define the new point as

$$(x^{(k+1)}, s^{(k+1)}) = \begin{cases} (\tilde{y}^{(k)}, x^{(k)}) & \text{if } \Phi(\tilde{y}^{(k)}, x^{(k)}) \le \Phi(x^{(k)}, s^{(k)}) + \sigma \lambda_k \Delta_k \\ (x^{(k)} + \lambda_k d_x^{(k)}, s^{(k)} + \lambda_k d_s^{(k)}) & \text{otherwise} \end{cases}$$

END

of a factor  $\delta < 1$ , we have  $\lambda_k^+ \geq \delta^M$ , where M is the smallest nonnegative integer such that  $\delta^M \leq (1-\sigma)/\rho$ . Therefore, (4.61) is satisfied with  $\lambda_{min} = \delta^M$ .

The descent direction and the backtracking procedure described above are at the basis of the new algorithm, named iPila (inertial Proximal inexact line-search algorithm), which formally consists in a descent method for the merit function  $\Phi(x, s)$ . In particular, it generates a sequence of iterates  $\{(x^{(k)}, s^{(k)})\}_{k \in \mathbb{N}}$  and a sequence of steplength parameters  $\{\lambda_k\}_{k \in \mathbb{N}}$  fulfilling the following decrease condition

$$\Phi(x^{(k+1)}, s^{(k+1)}) \le \Phi(x^{(k)}, s^{(k)}) + \sigma \lambda_k \Delta_k \quad \text{and} \quad \Phi(x^{(k+1)}, s^{(k+1)}) \le \Phi(\tilde{y}^{(k)}, x^{(k)}).$$

The connection with the inertial methods is in fact that, when  $(x^{(k+1)}, s^{(k+1)}) = (\tilde{y}^{(k)}, x^{(k)})$  is selected at STEP 7, the following iteration will consist of an actual inertial step. In practice, the condition at STEP 7 can be considered as an alternative acceptance rule for the inexact inertial proximal gradient point, having a similar role than the condition at STEP 6 of  $i^2$ Piano (see also (4.28)). The main difference is that here the acceptance condition is based on the Armijo inequality, while the one in i<sup>2</sup>Piano is based on the Descent Lemma.

Notice also that the inexact evaluation of the proximity operator in iPila is required only once per iteration, unlike in i<sup>2</sup>Piano, where it is needed at each step of the loop for selecting the parameter  $L_k$ , until inequality (4.28) is satisfied.

The Armijo condition results also in a larger freedom of choosing the parameters  $\alpha_k, \beta_k$ , which here satisfy very minimal conditions. A possible strategy to choose these parameters preserving both the theoretical prescriptions and the benefits deriving from the presence of an inertial step is described in Section 5.

**4.3.3. Convergence analysis.** In this section we first show that conditions [H1]–[H3] are satisfied for iPila.

PROPOSITION 4.15. Let  $\{(x^{(k)}, s^{(k)})\}_{k \in \mathbb{N}}$  be the sequence generated by *iPila*. Then, condition [H1] holds with  $d_k = \sqrt{-\Delta_k}$ ,  $a_k = \sigma \lambda_{min}$ . Moreover, under Assumption [A4], we also have

(4.64) 
$$0 = \lim_{k \to \infty} \|x^{(k)} - s^{(k)}\| = \lim_{k \to \infty} h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}) = \lim_{k \to \infty} \|\tilde{y}^{(k)} - x^{(k)}\|.$$

Proof. From the updating rule at STEP 7 and from Lemma 4.14, we have

$$\Phi(x^{(k+1)}, s^{(k+1)}) \le \Phi(x^{(k)}, s^{(k)}) + \sigma \lambda_k \Delta_k \le \Phi(x^{(k)}, s^{(k)}) + \sigma \lambda_{\min} \Delta_k.$$

Then, condition [H1] is satisfied with  $d_k^2 = -\Delta_k$ ,  $a_k = \sigma \lambda_{min}$ . Since from Assumption [A4] f is bounded from below,  $\Phi$  is bounded from below as well. Therefore, [H1] implies  $-\sum_{k=0}^{\infty} \Delta_k < \infty$  which, in turn, yields  $\lim_{k\to\infty} \Delta_k = 0$ . Recalling (4.56), this implies (4.64).

In the following we describe the setup for proving [H2], with the second auxiliary function defined as in (4.39).

PROPOSITION 4.16. Let  $\{(x^{(k)}, s^{(k)})\}_{k \in \mathbb{N}}$  be the sequence generated by Algorithm iPila with  $\gamma_{min} > 0$  and let  $\mathcal{F}$  be defined as in (4.39). Then, there exist  $\{\rho^{(k)}\}_{k \in \mathbb{N}}, \{r_k\}_{k \in \mathbb{N}} \subset \mathbb{R}$ , with  $\lim_{k \to \infty} r_k = 0$ , such that [H2] holds with  $u^{(k)} = \hat{y}^{(k)}$ , where  $\hat{y}^{(k)}$  is the exact minimizer of  $h^{(k)}(y; x^{(k)}, s^{(k)})$ , i.e.,

$$\hat{y}^{(k)} = \operatorname*{argmin}_{y \in \mathbb{R}^n} h^{(k)}(y; x^{(k)}, s^{(k)}).$$

Proof. From STEP 7, we have

$$\begin{split} \Phi(x^{(k+1)}, s^{(k+1)}) &\leq \Phi(\tilde{y}^{(k)}, x^{(k)}) = f(\tilde{y}^{(k)}) + \frac{1}{2} \|\tilde{y}^{(k)} - x^{(k)}\|^2 \\ &\leq f(\hat{y}^{(k)}) - \left(c + \frac{\alpha_{max}}{\theta}\right) h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}) + d\|x^{(k)} - s^{(k)}\|^2, \end{split}$$

where the last inequality follows from (4.16) and (4.13). Setting

(4.65) 
$$\rho^{(k)} = \sqrt{2} \left( -\left(c + \frac{\alpha_{max}}{\theta}\right) h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}) + d\|x^{(k)} - s^{(k)}\|^2 \right)^{\frac{1}{2}},$$

we obtain  $\Phi(x^{(k+1)}, s^{(k+1)}) \leq \mathcal{F}(\hat{y}^{(k)}, \rho^{(k)})$ , which represents the left-most inequality in [H2], with  $u^{(k)} = \hat{y}^{(k)}$ . On the other hand, from inequality (4.17) we obtain

$$\mathcal{F}(\hat{y}^{(k)},\rho^{(k)}) = f(\hat{y}^{(k)}) + \frac{1}{2}(\rho^{(k)})^2 \le f(x^{(k)}) - \bar{c}h^{(k)}(\tilde{y}^{(k)};x^{(k)},s^{(k)}) + \bar{d}\|x^{(k)} - s^{(k)}\|^2 + \frac{1}{2}(\rho^{(k)})^2.$$

Setting  $r_k = (\rho^{(k)})^2 / 2 - \bar{c}h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}) + (\bar{d} - \frac{1}{2}) \|x^{(k)} - s^{(k)}\|^2$ , we can write

$$\mathcal{F}(\hat{y}^{(k)}, \rho^{(k)}) \le f(x^{(k)}) + \frac{1}{2} \|x^{(k)} - s^{(k)}\|^2 + r_k = \Phi(x^{(k)}, s^{(k)}) + r_k.$$

From (4.64) we have that  $\lim_{k\to\infty} r_k = 0$  and this proves [H2].

PROPOSITION 4.17. Let  $\{(x^{(k)}, s^{(k)})\}_{k \in \mathbb{N}}$  be the sequence generated by *iPila* with  $\gamma_{min} > 0$ . Then, there exists a positive constant b such that [H3] is satisfied with  $I = \{1\}, \theta_1 = 1, \zeta_k = 0$ .

*Proof.* From (4.15) we know that there exists a subgradient  $\hat{v}^{(k)} \in \partial f(\hat{y}^{(k)})$  such that

(4.66) 
$$\|\hat{v}^{(k)}\| \le q\sqrt{-h^{(k)}(\tilde{y}^{(k)};x^{(k)},s^{(k)})} + q\|x^{(k)} - s^{(k)}\|$$

and, reasoning as in the proof of Proposition 4.8, it follows that

(4.67) 
$$\|\partial \mathcal{F}(\hat{y}^{(k)}, \rho^{(k)})\|_{-} \le \left\| \begin{pmatrix} \hat{v}^{(k)} \\ \rho^{(k)} \end{pmatrix} \right\| \le \|\hat{v}^{(k)}\| + |\rho^{(k)}|.$$

Let us analyze the two terms at the right-hand side of the inequality above, showing that both can be bounded from above with a multiple of  $\sqrt{-\Delta_k}$ . From (4.66) we obtain

$$\begin{aligned} \|\hat{v}^{(k)}\| &\leq q\sqrt{-h^{(k)}(\tilde{y}^{(k)};x^{(k)},s^{(k)}) + \gamma_{min}\|x^{(k)} - s^{(k)}\|^2} \\ &+ \frac{q}{\sqrt{\gamma_{min}}}\sqrt{\gamma_{min}\|x^{(k)} - s^{(k)}\|^2 - h^{(k)}(\tilde{y}^{(k)};x^{(k)},s^{(k)})}. \end{aligned}$$

which, setting  $A = q \left( 1 + 1/\sqrt{\gamma_{min}} \right)$  and using (4.56), yields

(4.68) 
$$\|\hat{v}^{(k)}\| \le A\sqrt{-h^{(k)}(\tilde{y}^{(k)};x^{(k)},s^{(k)})} + \gamma_{min}\|x^{(k)} - s^{(k)}\|^2 \le A\sqrt{-\Delta_k}$$

On the other hand, from definition (4.65)

$$\rho^{(k)} \le B\sqrt{-h^{(k)}(\tilde{y}^{(k)}; x^{(k)}, s^{(k)}) + \gamma_k \|x^{(k)} - s^{(k)}\|^2} = B\sqrt{-\Delta_k}$$

where  $B = \sqrt{2} \max \{c + \alpha_{max}/\theta, d/\gamma_{min}\}^{\frac{1}{2}}$ . Therefore, combining the last inequality above with (4.68) and (4.67), yields  $\|\partial \mathcal{F}(\hat{y}^{(k)}, \rho^{(k)})\|_{-} \leq (A+B)\sqrt{-\Delta_k}$  which proves that [H3] is satisfied with  $I = \{1\}, \ \theta_1 = 1, \ \zeta_k = 0, \ b = A + B, \ b_k = 1$ .

We are now ready for presenting the main convergence result for Algorithm iPila.

THEOREM 4.18. Suppose that  $\mathcal{F}$  is a KL function. Moreover, assume that the sequence  $\{x^{(k)}\}_{k\in\mathbb{N}}$  generated by iPila is bounded. Then,  $\{x^{(k)}\}_{k\in\mathbb{N}}$  converges to a stationary point of f.

Proof. By Proposition 4.15, 4.16, and 4.17, we know that conditions [H1]-[H2]-[H3] hold for iPila. From (4.11) we have  $\lim_{k\to\infty} ||u^{(k)} - x^{(k)}|| = \lim_{k\to\infty} ||\hat{y}^{(k)} - x^{(k)}|| = 0$ . Hence we can apply Proposition 4.8 and conclude that [H4] holds. Moreover, condition [H5] holds as a consequence of Lemma 4.12, since  $d_k = \sqrt{-\Delta_k}$  and  $||\tilde{y}^{(k)} - x^{(k)}|| \ge ||x^{(k+1)} - x^{(k)}||/\lambda_k \ge ||x^{(k+1)} - x^{(k)}||$  (see STEP 7). Finally, condition [H6] is trivially satisfied, since both sequences  $\{a_k\}_{k\in\mathbb{N}}$  and  $\{b_k\}_{k\in\mathbb{N}}$  are constant. Then Theorem 3.3 applies and guarantees that the sequence  $\{(x^{(k)}, \rho^{(k)})\}_{k\in\mathbb{N}}$  converges to a stationary point  $(x^*, \rho^*)$  of  $\mathcal{F}$ . Note that, since  $\mathcal{F}$  is the sum of separable functions, its subdifferential can be written as  $\partial \mathcal{F}(x, \rho) = \partial f(x) \times \{\rho\}$ . Then,  $(x^*, \rho^*)$  is stationary for  $\mathcal{F}$  if and only if  $\rho^* = 0$  and  $0 \in \partial f(x^*)$ . Hence  $x^*$  is a stationary point for f and  $\{x^{(k)}\}_{k\in\mathbb{N}}$  converges to it.  $\Box$ 

We refer the reader to Remark 4.10-4.11 for conditions on f guaranteeing that  $\mathcal{F}$  is a KL function and that the sequence  $\{x^{(k)}\}_{k\in\mathbb{N}}$  is bounded.

5. Numerical illustration: image denoising and deblurring in presence of impulse noise. In image restoration problems the goal is to recover a good quality image from a noisy blurred one. Following the variational approach, the clean image is obtained by solving an optimization problem with the structure (4.1), where the objective function includes a measure of the data fidelity and a regularization/penalization term, incorporating all a priori information on the desired solution. The data discrepancy is usually selected according to the noise statistics: in particular, when the data are affected by impulse noise, the variational model is defined as

$$\min_{x \in \mathbb{R}^n_{\geq 0}} \|Hx - g\|_1 + \mathcal{R}(x),$$

where  $g \in \mathbb{R}^n$  is the noisy blurred data,  $H \in \mathbb{R}^{n \times n}$  represents the blurring operator,  $\mathcal{R} : \mathbb{R}^n \to \mathbb{R}$  is the regularization term, and  $\mathbb{R}^n_{\geq 0} = \{x \in \mathbb{R}^n : x_i \geq 0, i = 1, ..., n\}$  is the nonnegative orthant. In the recent literature, several nonconvex regularization functionals have been proposed to overcome the known drawbacks of classical approaches, for example those based on the Total Variation function. In this paper we consider as regularization function the one proposed in [20, 21]:

$$\mathcal{R}(x) = \rho \sum_{\ell=1}^{q} \sum_{i=1}^{n} \log(1 + (K_{\ell}x)_{i}^{2}),$$

where the matrices  $K_{\ell} \in \mathbb{R}^{n \times n}$  correspond to the convolution with a given filter  $k_{\ell}$ , while  $\rho, \theta_{\ell}$  are positive parameters. In particular, the set of 48 filters  $k_{\ell}$  of size 7 × 7 and corresponding coefficients  $\theta_{\ell}$  are computed with a supervised machine learning technique. It is possible to prove that  $\mathcal{R}(x)$ has Lipschitz-continuous gradient, using the same arguments as in [16]. Then, setting  $f_0(x) = \mathcal{R}(x)$ complies with assumptions [A2]-[A3]. The nonnegativity constraint, expressed by means of the indicator function of the nonnegative orthant  $\iota_{\geq 0}(x)$ , can be included in the convex, nonsmooth term of the objective function, i.e.  $f_1(x) = ||Hx - g||_1 + \iota_{\geq 0}(x)$ . The regularization parameter  $\rho$ has been manually tuned in order to have a good quality restoration. Its value has been set equal to 0.08 for all the runs. Note that  $f_1$  does not have a closed-form proximal operator, therefore it is necessary to employ implementable inexactness criteria as the one proposed in Section 4.1. The proposed algorithms have been implemented in Matlab on a laptop equipped with a 2.5 GHzIntel Core i7-6500 processor and 16GB of RAM; the Matlab code is available online at [13]. The parameters of Algorithm i<sup>2</sup>Piano have been set as  $\delta = 0.5$ ,  $\gamma = 0.2$ ,  $\eta = 1.5$ ,  $\omega = 0.95$ . The estimate of the Lipschitz constant  $L_k$  is updated in a nondecreasing way. In particular, the initial value  $L_0$  is set as an input parameter: then, at STEP 1 of each iteration, the first tentative value is set as  $L_k = L_{k-1}$ . This value is possibly increased until inequality (4.28) is met. Actually, more sophisticated updating rules for this parameter could be adopted; however the objective function of the considered image restoration problem is very costly to evaluate, therefore a more conservative parameters selection rule has shown to be more convenient. As for Algorithm iPila, the parameters settings aim to mimic that of the inertial method i<sup>2</sup>Piano. Indeed, introducing the additional parameters  $\delta, \gamma, L_k, b_k > 0$ , with  $b_k = \frac{L_k + 2\delta}{L_k + 2\gamma}$ , we set  $\alpha_k, \beta_k, \gamma_k$  as follows

$$\beta_k = \frac{b_k - 1}{b_k - \frac{1}{2}}, \qquad \qquad \alpha_k = 2 \frac{1 - \beta_k}{L_k + 2\gamma}, \qquad \qquad \gamma_k = \gamma.$$

This choice is motivated by the following arguments. If  $L_k$  is a good local approximation of the Lipschitz constant satisfying condition

(5.1) 
$$f_0(\tilde{y}^{(k)}) \le f_0(x^{(k)}) + \langle \nabla f_0(x^{(k)}), \tilde{y}^{(k)} - x^{(k)} \rangle + \frac{L_k}{2} \| \tilde{y}^{(k)} - x^{(k)} \|^2$$

then reasoning as in the proof of Proposition 4.5, and choosing  $\delta = 0.5$  and  $\gamma_k = \gamma$ , we obtain

$$\Phi(\tilde{y}^{(k)}, x^{(k)}) \le \Phi(x^{(k)}, s^{(k)}) + \Delta_k + \frac{1}{2\alpha_k} \|\tilde{y}^{(k)} - x^{(k)}\|^2.$$

Hence, if  $L_k$  satisfies (5.1), the point  $(\tilde{y}^{(k)}, x^{(k)})$  will be likely accepted at STEP 7, as also confirmed by the numerical experience. This reasoning suggests to implement algorithm iPila as follows. We check the condition  $\Phi(\tilde{y}^{(k)}, x^{(k)}) \leq \Phi(x^{(k)}, s^{(k)}) + \sigma \Delta_k$  right after STEP 4: if the condition holds, then the steps from 4 to 7 are skipped in order to avoid unnecessary computations, and the next point is directly defined as  $(x^{(k+1)}, s^{(k+1)}) = (\tilde{y}^{(k)}, x^{(k)})$ ; otherwise, the value  $L_k$  is increased by a factor  $\eta = 1.5$ , and the line–search in steps 5–7 is performed in order to compute the next point. By possibly increasing  $L_k$ , we aim at improving the chances that  $(\tilde{y}^{(k+1)}, x^{(k+1)})$  is accepted at the next iteration, thus reducing the computational time due to the line–search reductions steps. The parameter in the Armijo condition, is set to  $\sigma = 10^{-4}$ .

For both i<sup>2</sup>Piano and iPila, the inexact proximal point  $\tilde{y}^{(k)}$  is computed by approximately solving the dual of problem  $\min_{y \in \mathbb{R}^n} h^{(k)}(y; x^{(k)}, s^{(k)})$ , which is a quadratic problem with simple constraints, with FISTA (more details can be found in [16] and references therein). The accuracy of the approximation is controlled by the parameter  $\tau$ : in our experiments we set  $\tau = 10^6$ , which corresponds to a good balancing of the computational complexity among inner and outer iterations. An extensive performance assessment of the algorithms with respect to this and other parameters is out of the scope of this paper, and it will be subject of future research.

The deblurring test problem has been obtained by first artificially blurring a good quality image, then simulating impulse noise on the 15% of the pixels with imnoise. The clean and the noisy image are in Figure 1 (a) and (b). Assuming reflective boundary conditions, matrix-vector multiplications involving H and  $H^T$  can be implemented efficiently with the DCT transform. In particular, each inner (dual) iteration requires the computation of two matrix vector product of this kind. Indeed, in our experiments, only one or two inner iterations per outer iteration are,



FIG. 1. Image deblurring test problem: (a) Original image  $(512 \times 512 \text{ pixels})$ ; (b) Noisy image, PSNR = 13.19; (c) Restored image, PSNR = 23.72



FIG. 2. Image deblurring test problem: decrease of the objective function (a) with respect to the computational time and (b) with different choices of  $L_0$ .

in general, needed to satisfy the inner stopping criterion. We compare our proposed algorithms to the variable metric line-search based method denominated VMILAn [12, 16] with its standard parameters settings. The objective function decrease for all three algorithms i<sup>2</sup>Piano, iPila and VMILAn is reported in Figure 2. Panel (a) reports the values  $f(x^{(k)})$  with respect to the computational time, where the initial estimate of the Lipschitz constant has been set equal to 0.001 for both i<sup>2</sup>Piano and iPila. The decrease of the objective function during the iterates of i<sup>2</sup>Piano and iPila with different choices of the parameter  $L_0$  is also presented in panel (b), where the horizontal axis still refers to the computational time. The picture shows that iPila is quite insensitive to the choice of  $L_0$ . In general, the numerical results show that i<sup>2</sup>Piano and iPila are able to solve challenging problems and are competitive with state-of-the-art methods. More effective rules for selecting their parameters will be subject of future work.

Acknowledgments. Silvia Bonettini, Marco Prato and Simone Rebegoldi are members of the INdAM research group GNCS. Peter Ochs acknowledges funding by the German Research Foundation (DFG Grant OC 150/1-1).

## REFERENCES

- P. A. ABSIL, R. MAHONY, AND B. ANDREWS, Convergence of the iterates of descent methods for analytic cost functions, SIAM J. Optim., 16 (2005), pp. 531–547.
- [2] H. ATTOUCH, J. BOLTE, AND B. F. SVAITER, Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized Gauss-Seidel methods, Math. Program., 137 (2013), pp. 91–129.
- [3] F. BACH, R. JENATTON, J. MAIRAL, AND G. OBOZINSKI, Structured sparsity through convex optimization, Stat. Sci., 27 (2012), pp. 450–468.
- M. BERTERO, P. BOCCACCI, AND V. RUGGIERO, Inverse Imaging with Poisson Data, IOP Publishing, Bristol, 2018.

- [5] D. BERTSEKAS, Nonlinear programming, Athena Scientific, Belmont, 1999.
- J. BOLTE, A. DANIILIDIS, O. LEY, AND L. MAZET, Characterizations of Lojasiewicz inequalities: Subgradient flows, talweg, convexity, Trans. Am. Math. Soc., 362 (2010), pp. 3319–3363.
- J. BOLTE, A. DANIILIDIS, AND S. M., Clarke subgradients of stratifiable functions, SIAM J. Optim., 10 (2007), pp. 556–572.
- [8] J. BOLTE, A. DANILIDIS, AND A. LEWIS, The lojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems, SIAM J. Optim., 17 (2007), pp. 1205–1223.
- [9] J. BOLTE, S. SABACH, AND M. TEBOULLE, Proximal alternating linearized minimization for nonconvex and nonsmooth problems, Math. Program., 146 (2014), pp. 459–494.
- [10] J. BOLTE, S. SABACH, M. TEBOULLE, AND Y. VAISBOURD, First order methods beyond convexity and Lipschitz gradient continuity with applications to quadratic inverse problems, SIAM J. Optim., 28 (2018), pp. 2131– 2151.
- [11] S. BONETTINI, I. LORIS, F. PORTA, AND M. PRATO, Variable metric inexact line-search based methods for nonsmooth optimization, SIAM J. Optim., 26 (2016), pp. 891–921.
- [12] S. BONETTINI, I. LORIS, F. PORTA, M. PRATO, AND S. REBEGOLDI, On the convergence of a linesearch based proximal-gradient method for nonconvex optimization, Inverse Probl., 33 (2017), p. 055005.
- [13] S. BONETTINI, P. OCHS, M. PRATO, AND S. REBEGOLDI, inertial inexact Proximal algorithm for nonconvex optimization (i2Piano) and inertial Proximal inexact linesearch algorithm (iPila) software. http://www.oasis.unimore.it/site/home/software.html, 2021.
- [14] S. BONETTINI, M. PRATO, AND S. REBEGOLDI, A block coordinate variable metric linesearch based proximal gradient method, Comput. Optim. Appl., 71 (2018), pp. 5–52.
- [15] S. BONETTINI, M. PRATO, AND S. REBEGOLDI, Convergence of inexact forward-backward algorithms using the forward-backward envelope, SIAM J. Optim., 30 (2020), pp. 3069–3097.
- [16] S. BONETTINI, M. PRATO, AND S. REBEGOLDI, New convergence results for the inexact variable metric forwardbackward method, Appl. Math. Comput., 392 (2021), p. 125719.
- [17] S. BONETTINI, S. REBEGOLDI, AND V. RUGGIERO, Inertial variable metric techniques for the inexact forwardbackward algorithm, SIAM J. Sci. Comput., 40 (2018), pp. A3180–A3210.
- [18] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, Optimization methods for large-scale machine learning, SIAM Rev., 60 (2018), pp. 223–311.
- [19] A. CHAMBOLLE AND T. POCK, An introduction to continuous optimization for imaging, Acta Numer., 25 (2016), pp. 161–319.
- [20] Y. CHEN, T. POCK, R. RANFTL, AND H. BISCHOF, Revisiting loss-specific training of filter-based MRFs for image restoration, in Pattern Recognition, J. Weickert, M. Hein, and B. Schiele, eds., Berlin, Heidelberg, 2013, Springer Berlin Heidelberg, pp. 271–281.
- [21] Y. CHEN, R. RANFTL, AND T. POCK, Insights into analysis operator learning: From patch-based sparse models to higher order mrfs, IEEE Trans. Image Process., 23 (2014), pp. 1060–1072.
- [22] E. CHOUZENOUX, J.-C. PESQUET, AND A. REPETTI, Variable metric forward-backward algorithm for minimizing the sum of a differentiable function and a convex function, J. Optim. Theory Appl., 162 (2014), pp. 107–132.
- [23] E. CHOUZENOUX, J.-C. PESQUET, AND A. REPETTI, A block coordinate variable metric forward-backward algorithm, Journal of Global Optimization, 66 (2016), pp. 457–485.
- [24] P. COMBETTES AND V. R. WAJS, Signal recovery by proximal forward-backward splitting, Multiscale Model. Simul., 4 (2005), pp. 1168–1200.
- [25] P. L. COMBETTES AND J.-C. PESQUET, Proximal splitting methods in signal processing, in Fixed-point algorithms for inverse problems in science and engineering, H. H. Bauschke, R. S. Burachik, P. L. Combettes, V. Elser, D. R. Luke, and H. Wolkowicz, eds., Springer Optimization and Its Applications, Springer, New York, NY, 2011, pp. 185–212.
- [26] P. FRANKEL, G. GARRIGOS, AND J. PEYPOUQUET, Splitting methods with variable metric for Kurdyka-Lojasiewicz functions and general convergence rates, J. Opt. Theory Appl., 165 (2015), pp. 874–900.
- [27] K. KURDYKA, On gradients of functions definable in o-minimal structures, Ann. Inst. Fourier, 48 (1998), pp. 769–783.
- [28] G. LI AND T. K. PONG, Douglas-Rachford splitting for nonconvex optimization with application to nonconvex feasibility problems, Math. Program., 159 (2016), pp. 371–401.
- [29] D. NOLL, Convergence of non-smooth descent methods using the kurdyka-Lojasiewicz inequality, J. Opt. Theory Appl., 160 (2014), pp. 553–572.
- [30] P. OCHS, Unifying abstract inexact convergence theorems and block coordinate variable metric iPiano, SIAM J. Optim., 29 (2019), pp. 541–570.
- [31] P. OCHS, Y. CHEN, T. BROX, AND T. POCK, *iPiano: Inertial proximal algorithm for non-convex optimization*, SIAM J. Imaging Sci., 7 (2014), pp. 1388–1419.
- [32] B. POLYAK, Introduction to optimization, Optimization Software Inc., Publication Division, New York, 1987.
- [33] B. T. POLYAK, Some methods of speeding up the convergence of iteration methods, USSR Comput. Math. Math. Phys., 4 (1964), pp. 1–17.
- [34] R. T. ROCKAFELLAR, Convex Analysis, Princeton University Press, Princeton, NJ, 1970.
- [35] R. T. ROCKAFELLAR, R. J.-B. WETS, AND M. WETS, Variational Analysis, vol. 317 of Grundlehren der Mathematischen Wissenschaften, Springer, Berlin, 1998.
- [36] L. STELLA, A. THEMELIS, AND P. PATRINOS, Forward-backward quasi-Newton methods for nonsmooth optimization problems., Comput. Optim. Appl., 67 (2017), pp. 443–487.
- [37] S. VILLA, S. SALZO, L. BALDASSARRE, AND A. VERRI, Accelerated and inexact forward-backward algorithms, SIAM J. Optim., 23 (2013), pp. 1607–1633.
- [38] A. ZALINESCU, Convex analysis in general vector spaces, World Scientific Publishing Co. Inc., 2002.