
A Multi-step Inertial Forward–Backward Splitting Method for Non-convex Optimization

Jingwei Liang and Jalal M. Fadili

Normandie Univ, ENSICAEN, CNRS, GREYC
{Jingwei.Liang, Jalal.Fadili}@greyc.ensicaen.fr

Gabriel Peyré

CNRS, DMA, ENS Paris
Gabriel.Peyre@ens.fr

Abstract

We propose a multi-step inertial Forward–Backward splitting algorithm for minimizing the sum of two non-necessarily convex functions, one of which is proper lower semi-continuous while the other is differentiable with a Lipschitz continuous gradient. We first prove global convergence of the algorithm with the help of the Kurdyka–Łojasiewicz property. Then, when the non-smooth part is also partly smooth relative to a smooth submanifold, we establish finite identification of the latter and provide sharp local linear convergence analysis. The proposed method is illustrated on several problems arising from statistics and machine learning.

1 Introduction

1.1 Non-convex non-smooth optimization

Non-smooth optimization has proved extremely useful to all quantitative disciplines of science including statistics and machine learning. A common trend in modern science is the increase in size of datasets, which drives the need for more efficient optimization schemes. For large-scale problems with non-smooth and possibly non-convex terms, it is possible to generalize gradient descent with the Forward–Backward (FB) splitting scheme [3] (a.k.a proximal gradient descent), which includes projected gradient descent as a sub-case.

Formally, we equip \mathbb{R}^n the n -dimensional Euclidean space with the standard inner product $\langle \cdot, \cdot \rangle$ and associated norm $\| \cdot \|$ respectively. Our goal is the generic minimization of composite objectives of the form

$$\min_{x \in \mathbb{R}^n} \{ \Phi(x) \stackrel{\text{def}}{=} R(x) + F(x) \}, \quad (\mathcal{P})$$

where we have

- (A.1) $R : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is the *penalty function* which is proper lower semi-continuous (lsc), and bounded from below;
- (A.2) $F : \mathbb{R}^n \rightarrow \mathbb{R}$ is the *loss function* which is finite-valued, differentiable and its gradient ∇F is L -Lipschitz continuous.

Throughout, no convexity is imposed neither on R nor on F .

The class of problems we consider is that of non-smooth and non-convex optimization problems. Here are some examples that are of particular relevance to problems in regression, machine learning and classification.

Example 1.1 (Sparse regression). Let $A \in \mathbb{R}^{m \times n}$, $y \in \mathbb{R}^m$, $\mu > 0$, and $\|x\|_0$ is the ℓ_0 pseudo-norm (see Example 4.1). Consider (see e.g. [11])

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - Ax\|^2 + \mu \|x\|_0. \quad (1.1)$$

Example 1.2 (Principal component pursuit (PCP)). The PCP problem [9] aims at decomposing a given matrix into *sparse* and *low-rank* components

$$\min_{(x_s, x_l) \in (\mathbb{R}^{n_1 \times n_2})^2} \frac{1}{2} \|y - x_s - x_l\|_F^2 + \mu_1 \|x_s\|_0 + \mu_2 \text{rank}(x_l), \quad (1.2)$$

where $\|\cdot\|_F$ is the Frobenius norm and μ_1 and $\mu_2 > 0$.

Example 1.3 (Sparse Support Vector Machines). One would like to find a linear decision function which minimizes the objective

$$\min_{(b, x) \in \mathbb{R} \times \mathbb{R}^n} \frac{1}{m} \sum_{i=1}^m G(\langle x, z_i \rangle + b, y_i) + \mu \|x\|_0, \quad (1.3)$$

where for $i = 1, \dots, m$, $(z_i, y_i) \in \mathbb{R}^n \times \{\pm 1\}$ is the training set, and G is a smooth loss function with Lipschitz-continuous gradient such as the squared hinge loss $G(\hat{y}_i, y_i) = \max(0, 1 - \hat{y}_i y_i)^2$ or the logistic loss $G(\hat{y}_i, y_i) = \log(1 + e^{-\hat{y}_i y_i})$.

(Inertial) Forward–Backward The Forward–Backward splitting method for solving (\mathcal{P}) reads

$$x_{k+1} \in \text{prox}_{\gamma_k R}(x_k - \gamma_k \nabla F(x_k)), \quad (1.4)$$

where $\gamma_k > 0$ is a descent step-size, and

$$\text{prox}_{\gamma R}(\cdot) \stackrel{\text{def}}{=} \text{Argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|x - \cdot\|^2 + \gamma R(x), \quad (1.5)$$

denotes the proximity operator of R . $\text{prox}_{\gamma R}(x)$ is non-empty under **(A.1)** and is set-valued in general. Lower-boundedness of R can be relaxed by requiring e.g. coercivity of the objective in (1.5).

Since the pioneering work of Polyak [24] on the *heavy-ball method* approach to gradient descent, several works have adapted this methodology to various optimization schemes. For instance, the inertial proximal point algorithm [1, 2], or the inertial FB methods [22, 21, 4, 20]. The FISTA scheme [5, 10] also belongs to this class. See [20] for a detailed account.

The non-convex case In the context of non-convex optimization, [3] was the first to establish convergence of the FB iterates when the objective Φ satisfies the Kurdyka-Łojasiewicz property¹. Following their footprints, [8, 23] established convergence of the special inertial schemes in [22] in the non-convex setting.

1.2 Contributions

In this paper, we introduce a novel inertial scheme (Algorithm 1) and study its global and local properties to solve the non-smooth and non-convex optimization problem (\mathcal{P}) . More precisely, our main contributions can be summarized as follows.

A globally convergent general inertial scheme We propose a general multi-step inertial FB (MiFB) algorithm to solve (\mathcal{P}) . This algorithm is very flexible as it allows higher memory and even *negative* inertial parameters (unlike previous work [20]). Global convergence of any bounded sequence of iterates to a critical point is proved when the objective Φ is lower-bounded and satisfies the Kurdyka-Łojasiewicz property.

Local convergence properties under partial smoothness Under the additional assumptions that the smooth part is locally C^2 around a critical point x^* (where $x_k \rightarrow x^*$), and that the non-smooth component R is partly smooth (see Definition 3.1) relative to an active submanifold \mathcal{M}_{x^*} , we show that \mathcal{M}_{x^*} can be identified in finite time, i.e. $x_k \in \mathcal{M}_{x^*}$ for all k large enough. Building on this finite identification result, we provide a sharp local linear convergence analysis and we characterize precisely the corresponding convergence rate which, in particular, reveals the role of \mathcal{M}_{x^*} . Moreover, this local convergence analysis naturally opens the door to higher-order acceleration, since under some circumstances, the original problem (\mathcal{P}) is eventually equivalent to locally minimizing Φ on \mathcal{M}_{x^*} , and partial smoothness implies that Φ is actually C^2 on \mathcal{M}_{x^*} .

¹We are aware of the works existing on convergence of the objective sequence $\Phi(x_k)$ of FB, including rates, in the non-smooth and non-convex setting. But given that, in general, this does not say anything about convergence of the sequence of iterates x_k , they are irrelevant to our discussion.

Algorithm 1: A Multi-step Inertial Forward–Backward (MiFB)

Initial: $s \geq 1$ is an integer, $I = \{0, 1, \dots, s-1\}$, $x_0 \in \mathbb{R}^n$ and $x_{-s} = \dots = x_{-1} = x_0$.

repeat

Let $0 < \underline{\gamma} \leq \gamma_k \leq \bar{\gamma} < \frac{1}{L}$, $\{a_{0,k}, a_{1,k}, \dots\} \in]-1, 2]^s$, $\{b_{0,k}, b_{1,k}, \dots\} \in]-1, 2]^s$:

$$y_{a,k} = x_k + \sum_{i \in I} a_{i,k} (x_{k-i} - x_{k-i-1}), \quad (1.6)$$

$$y_{b,k} = x_k + \sum_{i \in I} b_{i,k} (x_{k-i} - x_{k-i-1}),$$

$$x_{k+1} \in \text{prox}_{\gamma_k R}(y_{a,k} - \gamma_k \nabla F(y_{b,k})). \quad (1.7)$$

$k = k + 1$;

until convergence;

1.3 Notations

Throughout the paper, \mathbb{N} is the set of non-negative integers. For a nonempty closed convex set $\Omega \subset \mathbb{R}^n$, $\text{ri}(\Omega)$ is its relative interior, and $\text{par}(\Omega) = \mathbb{R}(\Omega - \Omega)$ is the subspace parallel to it.

Let $R : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be a lsc function, its domain is defined as $\text{dom}(R) \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n : R(x) < +\infty\}$, and it is said to be proper if $\text{dom}(R) \neq \emptyset$. We need the following notions from variational analysis, see e.g. [25] for details. Given $x \in \text{dom}(R)$, the Fréchet subdifferential $\partial^F R(x)$ of R at x , is the set of vectors $v \in \mathbb{R}^n$ that satisfies $\liminf_{z \rightarrow x, z \neq x} \frac{1}{\|x-z\|} (R(z) - R(x) - \langle v, z-x \rangle) \geq 0$. If $x \notin \text{dom}(R)$, then $\partial^F R(x) = \emptyset$. The limiting-subdifferential (or simply subdifferential) of R at x , written as $\partial R(x)$, is defined as $\partial R(x) \stackrel{\text{def}}{=} \{v \in \mathbb{R}^n : \exists x_k \rightarrow x, R(x_k) \rightarrow R(x), v_k \in \partial^F R(x_k) \rightarrow v\}$. Denote $\text{dom}(\partial R) \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n : \partial R(x) \neq \emptyset\}$. Both $\partial^F R(x)$ and $\partial R(x)$ are closed, with $\partial^F R(x)$ convex and $\partial^F R(x) \subset \partial R(x)$ [25, Proposition 8.5]. Since R is lsc, it is (subdifferentially) regular at x if and only if $\partial^F R(x) = \partial R(x)$ [25, Corollary 8.11].

An lsc function R is r -prox-regular at $\bar{x} \in \text{dom}(R)$ for $\bar{v} \in \partial R(\bar{x})$ if $\exists r > 0$ such that $R(x') > R(\bar{x}) + \langle \bar{v}, x' - \bar{x} \rangle - \frac{1}{2r} \|x - x'\|^2 \forall x, x'$ near \bar{x} , $R(x)$ near $R(\bar{x})$ and $v \in \partial R(x)$ near \bar{v} .

A necessary condition for x to be a minimizer of R is $0 \in \partial R(x)$. The set of critical points of R is $\text{crit}(R) = \{x \in \mathbb{R}^n : 0 \in \partial R(x)\}$.

2 Global convergence of MiFB

2.1 Kurdyka–Łojasiewicz property

Let $R : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be a proper lsc function. For η_1, η_2 such that $-\infty < \eta_1 < \eta_2 < +\infty$, define the set $[\eta_1 < R < \eta_2] \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n : \eta_1 < R(x) < \eta_2\}$.

Definition 2.1. R is said to have the Kurdyka–Łojasiewicz property at $\bar{x} \in \text{dom}(R)$ if there exists $\eta \in]0, +\infty]$, a neighbourhood U of \bar{x} and a continuous concave function $\varphi : [0, \eta[\rightarrow \mathbb{R}_+$ such that

- (i) $\varphi(0) = 0$, φ is C^1 on $]0, \eta[$, and for all $s \in]0, \eta[$, $\varphi'(s) > 0$;
- (ii) for all $x \in U \cap [R(\bar{x}) < R < R(\bar{x}) + \eta]$, the Kurdyka–Łojasiewicz inequality holds

$$\varphi'(R(x) - R(\bar{x})) \text{dist}(0, \partial R(x)) \geq 1. \quad (2.1)$$

Proper lsc functions which satisfy the Kurdyka–Łojasiewicz property at each point of $\text{dom}(\partial R)$ are called KL functions.

Roughly speaking, KL functions become sharp up to reparameterization via φ , called a desingularizing function for R . Typical KL functions are the class of semi-algebraic functions, see [6, 7]. For instance, the ℓ_0 pseudo-norm and the rank function (see Example 1.1-1.3 and Section 4.1) are indeed KL.

2.2 Global convergence

Let $\mu, \nu > 0$ be two constants. For $i \in I$ and $k \in \mathbb{N}$, define the following quantities,

$$\beta_k \stackrel{\text{def}}{=} \frac{1 - \gamma_k L - \mu - \nu \gamma_k}{2\gamma_k}, \quad \underline{\beta} \stackrel{\text{def}}{=} \liminf_{k \in \mathbb{N}} \beta_k \quad \text{and} \quad \alpha_{i,k} \stackrel{\text{def}}{=} \frac{sa_{i,k}^2}{2\gamma_k \mu} + \frac{sb_{i,k}^2 L^2}{2\nu}, \quad \bar{\alpha}_i \stackrel{\text{def}}{=} \limsup_{k \in \mathbb{N}} \alpha_{i,k}. \quad (2.2)$$

Theorem 2.2 (Convergence of MiFB (Algorithm 1)). For problem (\mathcal{P}) , suppose that (A.1)-(A.2) hold, and moreover Φ is a proper lsc KL function. For Algorithm 1, choose $\mu, \nu, \gamma_k, a_{i,k}, b_{i,k}$ such that

$$\delta \stackrel{\text{def}}{=} \underline{\beta} - \sum_{i \in I} \bar{\alpha}_i > 0. \quad (2.3)$$

Then each bounded sequence $\{x_k\}_{k \in \mathbb{N}}$ generated by MiFB satisfies

- (i) $\{x_k\}_{k \in \mathbb{N}}$ has finite length, i.e. $\sum_{k \in \mathbb{N}} \|x_k - x_{k-1}\| < +\infty$;
- (ii) There exists a critical point $x^* \in \text{crit}(\Phi)$ such that $\lim_{k \rightarrow \infty} x_k = x^*$.
- (iii) If Φ has the KL property at a global minimizer x^* , then starting sufficiently close from x^* , any sequence $\{x_k\}_{k \in \mathbb{N}}$ converges to a global minimum of Φ and satisfies (i).

The proof is detailed in the supplementary material.

Remark 2.3.

- (i) The convergence result holds true for any real Hilbert space. The boundedness of $\{x_k\}_{k \in \mathbb{N}}$ is automatically ensured under standard assumptions such as coercivity of Φ .
- (ii) It is known from [13] that if the desingularizing function $\varphi = \frac{C}{\theta} t^\theta$, $C > 0$ and $\theta \in [\frac{1}{2}, 1[$, then global linear convergence of the objective and the iterates can be derived. However, we will not pursue this further since our main interest is local linear convergence.
- (iii) Unlike existing work, *negative* inertial parameters are allowed by Theorem 2.2.
- (iv) When $a_{i,k} \equiv 0$ and $b_{i,k} \equiv 0$, i.e. the case of FB splitting, condition (2.3) holds naturally as long as $\bar{\gamma} < \frac{1}{L}$ which recovers the case of [3];
- (v) From (2.2) and (2.3), we conclude the following:
 - (a) $s = 1$: if $b_{0,k} \equiv b$, $a_{0,k} \equiv a$ (i.e. constant inertial parameters), then (2.3) implies that a, b belong to an ellipsoid: $\frac{a^2}{2\gamma\mu} + \frac{b^2}{2\nu/L^2} < \underline{\beta} = \frac{1-\bar{\gamma}L-\mu-\nu\bar{\gamma}}{2\bar{\gamma}}$.
 - (b) When $s \geq 2$, for each $i \in I$, let $b_{i,k} = a_{i,k} \equiv a_i$ (i.e. constant symmetric inertial parameters), then (2.3) means that the a_i 's live in a ball: $(\frac{1}{2\gamma\mu} + \frac{1}{2\nu/L^2}) \sum_{i \in I} a_i^2 < \underline{\beta}$.

An empirical approach for inertial parameters Besides Theorem 2.2, we also provide an empirical bound for the choice of the inertial parameters. Consider the setting: $\gamma_k \equiv \gamma \in]0, 1/L[$ and $b_{i,k} = a_{i,k} \equiv a_i \in]-1, 2[$, $i \in I$. We have the following empirical bound for the summand $\sum_{i \in I} a_i$:

$$\sum_i a_i \in]0, \min\{1, \frac{1/L-\gamma}{|2\gamma-1/L|}\} [. \quad (2.4)$$

To ensure the convergence $\{x_k\}_{k \in \mathbb{N}}$, an online updating rule should be applied together with the empirical bound. More precisely, choose a_i according to (2.4). Then for each $k \in \mathbb{N}$, let $b_{i,k} = a_{i,k}$ and choose $a_{i,k}$ such that $\sum_i a_{i,k} = \min\{\sum_i a_i, c_k\}$ where $c_k > 0$ is such that $\{c_k \sum_{i \in I} \|x_{k-i} - x_{k-i-1}\|\}_{k \in \mathbb{N}}$ is summable. For instance, $c_k = \frac{c}{k^{1+q} \sum_{i \in I} \|x_{k-i} - x_{k-i-1}\|}$, $c > 0, q > 0$.

Note that the allowed choices of the summand $\sum_i a_i$ by (2.4) is larger than those of Theorem 2.2. For instance, (2.4) allows $\sum_i a_i = 1$ for $\gamma \in]0, \frac{2}{3L}[$. While for Theorem 2.2, $\sum_i a_i = 1$ can be reached only when $\gamma \rightarrow 0$.

3 Local convergence properties of MiFB

3.1 Partial smoothness

Let $\mathcal{M} \subset \mathbb{R}^n$ be a C^2 -smooth submanifold, let $\mathcal{T}_{\mathcal{M}}(x)$ the tangent space of \mathcal{M} at any point $x \in \mathcal{M}$.

Definition 3.1. The function $R : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is C^2 -partly smooth at $\bar{x} \in \mathcal{M}$ relative to \mathcal{M} for $\bar{v} \in \partial R(\bar{x}) \neq \emptyset$ if \mathcal{M} is a C^2 -submanifold around \bar{x} , and

- (i) (Smoothness): R restricted to \mathcal{M} is C^2 around \bar{x} ;
- (ii) (Regularity): R is regular at all $x \in \mathcal{M}$ near \bar{x} and R is r -prox-regular at \bar{x} for \bar{v} ;
- (iii) (Sharpness): $\mathcal{T}_{\mathcal{M}}(\bar{x}) = \text{par}(\partial R(\bar{x}))^\perp$;
- (iv) (Continuity): The set-valued mapping ∂R is continuous at \bar{x} relative to \mathcal{M} .

We denote the class of partly smooth functions at x relative to \mathcal{M} for v as $\text{PSF}_{x,v}(\mathcal{M})$. Partial smoothness was first introduced in [15] and its directional version stated here is due to [18, 12]. Prox-regularity is sufficient to ensure that the partly smooth submanifolds are locally unique [18, Corollary 4.12], [12, Lemma 2.3 and Proposition 10.12].

3.2 Finite activity identification

One of the key consequences of partial smoothness is finite identification of the partial smoothness submanifold associated to R for problem (P) . This is formalized in the following statement.

Theorem 3.2 (Finite activity identification). *Suppose that Algorithm 1 is run under the conditions of Theorem 2.2, such that the generated sequence $\{x_k\}_{k \in \mathbb{N}}$ converges to a critical point $x^* \in \text{crit}(\Phi)$. Assume that $R \in \text{PSF}_{x^*, -\nabla F(x^*)}(\mathcal{M}_{x^*})$ and the non-degeneracy condition*

$$-\nabla F(x^*) \in \text{ri}(\partial R(x^*)), \quad (\text{ND})$$

holds. Then, $x_k \in \mathcal{M}_{x^}$ for all k large enough.*

See the supplementary material for the proof. This result generalizes that of [20] to the non-convex case and multiple inertial steps.

3.3 Local linear convergence

Given $\gamma \in]0, \frac{1}{L}[$ and a critical point $x^* \in \text{crit}(\Phi)$, let \mathcal{M}_{x^*} be a C^2 -smooth submanifold and $R \in \text{PSF}_{x^*, -\nabla F(x^*)}(\mathcal{M}_{x^*})$. Denote $T_{x^*} \stackrel{\text{def}}{=} \mathcal{T}_{\mathcal{M}_{x^*}}(x^*)$ and the following matrices which are all symmetric,

$$H \stackrel{\text{def}}{=} \gamma P_{T_{x^*}} \nabla^2 F(x^*) P_{T_{x^*}}, \quad G \stackrel{\text{def}}{=} \text{Id} - H, \quad Q \stackrel{\text{def}}{=} \gamma \nabla_{\mathcal{M}_{x^*}}^2 \Phi(x^*) P_{T_{x^*}} - H, \quad (3.1)$$

where $\nabla_{\mathcal{M}_{x^*}}^2 \Phi$ is the Riemannian Hessian of Φ along the submanifold \mathcal{M}_{x^*} (readers may refer to the supplementary material from more details on differential calculus on Riemannian manifolds).

To state our local linear convergence result, the following assumptions will play a key role.

Restricted injectivity Besides the local C^2 -smoothness assumption on F , following the idea of [19, 20], we assume the restricted injectivity condition,

$$\ker(\nabla^2 F(x^*)) \cap T_{x^*} = \{0\}. \quad (\text{RI})$$

Positive semi-definiteness of Q Assume that Q is *positive semi-definite*, i.e. $\forall h \in T_{x^*}$,

$$\langle h, Qh \rangle \geq 0. \quad (3.2)$$

Under (3.2), $\text{Id} + Q$ is symmetric positive definite, hence invertible, we denote $P \stackrel{\text{def}}{=} (\text{Id} + Q)^{-1}$.

Convergent parameters The parameters of MiFB (Algorithm 1), are convergent, i.e.

$$a_{i,k} \rightarrow a_i, \quad b_{i,k} \rightarrow b_i, \quad \forall i \in I \quad \text{and} \quad \gamma_k \rightarrow \gamma \in [\underline{\gamma}, \min\{\bar{\gamma}, \bar{r}\}], \quad (3.3)$$

where $\bar{r} < r$, and r is the prox-regularity modulus of R (see Definition 3.1).

Remark 3.3.

- (i) Condition (3.2) can be met by various non-convex functions, such as polyhedral functions, including the ℓ_0 pseudo-norm. For the rank function, it is also observed that this condition holds in our numerical experiments of Section 4.
- (ii) Condition (3.3) asserts that both the inertial parameters $(a_{i,k}, b_{i,k})$ and the step-size γ_k should converge to some limit points, and this condition cannot be relaxed in general.
- (iii) It can be shown that conditions (3.2) and (RI) together imply that x^* is a local minimizer of Φ in (P) , and Φ grows at least quadratically near x^* . The arguments to prove this are essentially adapted from those used to show [20, Proposition 4.1(ii)].

We need the following notations:

$$\begin{aligned} M_0 &\stackrel{\text{def}}{=} (a_0 - b_0)P + (1 + b_0)PG, \quad M_s \stackrel{\text{def}}{=} -(a_{s-1} - b_{s-1})P - b_{s-1}PG, \\ M_i &\stackrel{\text{def}}{=} -((a_{i-1} - a_i) - (b_{i-1} - b_i))P - (b_{i-1} - b_i)PG, \quad i = 1, \dots, s-1, \\ M &\stackrel{\text{def}}{=} \begin{bmatrix} M_0 & \cdots & M_{s-1} & M_s \\ \text{Id} & \cdots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \cdots & \text{Id} & 0 \end{bmatrix}, \quad d_k \stackrel{\text{def}}{=} \begin{pmatrix} x_k - x^* \\ \vdots \\ x_{k-s} - x^* \end{pmatrix}. \end{aligned} \quad (3.4)$$

Theorem 3.4 (Local linear convergence). *Suppose that Algorithm 1 is run under the setting of Theorem 3.2. Moreover, assume that (RI), (3.2) and (3.3) hold. Then for all k large enough,*

$$d_{k+1} = Md_k + o(\|d_k\|). \quad (3.5)$$

If $\rho(M) < 1$, then given any $\rho \in]\rho(M), 1[$, there exists $K \in \mathbb{N}$ such that $\forall k \geq K$,

$$\|x_k - x^*\| = O(\rho^{k-K}). \quad (3.6)$$

In particular, if $s = 1$, then $\rho(M) < 1$ if R is locally polyhedral around x^ or if $a_0 = b_0$.*

See the supplementary material for the proof.

Remark 3.5.

- (i) When $s = 1$, $\rho(M)$ can be given explicitly in terms of the parameters of the algorithm (i.e. a_0 , b_0 and γ), see [20, Section 4.2] for details. However, the spectral analysis of M becomes much more complicated to get for $s \geq 2$, where the analysis of at least cubic equations are involved. Therefore, for the sake of brevity, we shall skip the detailed discussion here.
- (ii) When $s = 1$, it was shown in [20] that the optimal convergence rate that can be obtained by 1-step inertial scheme with fixed γ is $\rho_{s=1}^* = 1 - \sqrt{1 - \tau\gamma}$, where from condition (RI), continuity of $\nabla^2 F$ at x^* implies that there exists $\tau > 0$ and a neighbourhood of x^* such that $\langle h, \nabla^2 F(x^*)h \rangle \geq \tau\|h\|^2$, for all $h \in T_{x^*}$. As we will see in the numerical experiments of Section 4, such a rate can be improved by our multi-step inertial scheme. Taking $s = 2$ for example, we will show that for a certain class of functions, the optimal local linear rate is close to or even is $\rho_{s=2}^* = 1 - \sqrt[3]{1 - \tau\gamma}$, which is obviously faster than $\rho_{s=1}^*$.
- (iii) Though it can be satisfied for many problems in practice, the restricted injectivity (RI) can be removed for some penalties R , for instance, when R is locally polyhedral near x^* .

4 Numerical experiments

In this section, we illustrate our results with some numerical experiments carried out on the problems in Example 1.1, 1.2 and 1.3.

4.1 Examples of KL and partly smooth functions

All the objectives Φ in the above mentioned examples are continuous KL functions. Indeed, in Example 1.1 and 1.2, Φ is the sum of semi-algebraic functions which is also semi-algebraic. In Example 1.3, Φ is also algebraic when G is the squared hinge loss, and definable in an o-minimal structure for the logistic loss (see e.g. [26] for material on o-minimal structures).

Moreover, R is partly smooth in all these examples as we show now.

Example 4.1 (ℓ_0 pseudo-norm). The ℓ_0 pseudo-norm is locally constant. Moreover, it is regular on \mathbb{R}^n ([14, Remark 2]) and its subdifferential is given by (see [14, Theorem 1])

$$\partial\|x\|_0 = \text{span}((e_i)_{i \in \text{supp}(x)^c}),$$

where $(e_i)_{i=1, \dots, n}$ is the standard basis, and $\text{supp}(x) = \{i : x_i \neq 0\}$. The proximity operator of ℓ_0 -norm is given by hard-thresholding,

$$\text{prox}_{\gamma\|x\|_0}(z) = \begin{cases} z & \text{if } |z| > \sqrt{2\gamma}, \\ \text{sign}(z)[0, z] & \text{if } |z| = \sqrt{2\gamma}, \\ 0 & \text{if } |z| < \sqrt{2\gamma}. \end{cases}$$

It can then be easily verified that the ℓ_0 pseudo-norm is partly smooth at any x relative to the subspace

$$\mathcal{M}_x = T_x = \{z \in \mathbb{R}^n : \text{supp}(z) \subset \text{supp}(x)\}.$$

It is also prox-regular at x for any bounded $v \in \partial\|x\|_0$. Note also condition (ND) is automatically verified and that the Riemannian gradient and Hessian along T_x of $\|\cdot\|_0$ vanish.

Example 4.2 (Rank). The rank function is the spectral extension of ℓ_0 pseudo-norm to matrix-valued data $x \in \mathbb{R}^{n_1 \times n_2}$ [17]. Consider a singular value decomposition (SVD) of x , i.e. $x = U \text{diag}(\sigma(x)) V^*$, where $U = \{u_1, \dots, u_n\}$, $V = \{v_1, \dots, v_n\}$ are orthonormal matrices, and

$\sigma(x) = (\sigma_i(x))_{i=1,\dots,n}$ is the vector of singular values. By definition, $\text{rank}(x) \stackrel{\text{def}}{=} \|\sigma(x)\|_0$. Thus the rank function is partly smooth relative at x to the set of fixed rank matrices

$$\mathcal{M}_x = \{z \in \mathbb{R}^{n_1 \times n_2} : \text{rank}(z) = \text{rank}(x)\},$$

which is a C^2 -smooth submanifold [16]. The tangent space of \mathcal{M}_x at x is

$$\mathcal{T}_{\mathcal{M}_x}(x) = T_x = \{z \in \mathbb{R}^{n_1 \times n_2} : u_i^* z v_j = 0, \text{ for all } r < i \leq n_1, r < j \leq n_2\},$$

The rank function is also regular its subdifferential reads

$$\partial \text{rank}(x) = U \partial(\|\sigma(x)\|_0) V^* = U \text{span}((e_i)_{i \in \text{supp}(\sigma(x))^c}) V^*,$$

which is a vector space (see [14, Theorem 4 and Proposition 1]). The proximity operator of rank function amounts to applying hard-thresholding to the singular values. Observe that by definition of \mathcal{M}_x , the Riemannian gradient and Hessian of the rank function along \mathcal{M}_x also vanish.

For Example 1.2, it is worth noting from the above examples and separability of the regularizer that the latter is also partly smooth relative to the cartesian product of the partial smoothness submanifolds of ℓ_0 and the rank function.

4.2 Experimental results

For the problem in Example 1.1, we generated $y = Ax_{\text{ob}} + \omega$ with $m = 48$, $n = 128$, the entries of A are i.i.d. zero-mean and unit variance Gaussian, x_{ob} is 8-sparse, and $\omega \in \mathbb{R}^m$ is an additive noise with small variance.

For the problem in Example 1.2, we generated $y = x_s + x_l + \omega$, with $n_1 = n_2 = 50$, x_s is 250-sparse, and the rank of x_l is 5, and ω is an additive noise with small variance.

For Example 1.3, we generated $m = 64$ training samples with $n = 96$ -dimensional feature space.

For all presented numerical results, 3 different settings were tested:

- the FB method, with $\gamma_k \equiv 0.3/L$, noted as “FB”;
- MiFB with $s = 1$, $b_k = a_k \equiv a$ and $\gamma_k \equiv 0.3/L$, noted as “1-iFB”;
- MiFB with $s = 2$, $b_{i,k} = a_{i,k} \equiv a_i$, $i = 0, 1$ and $\gamma_k \equiv 0.3/L$, noted as “2-iFB”.

Tightness of theoretical prediction The convergence profiles of $\|x_k - x^*\|$ are shown in Figure 1. As it can be seen from all the plots, finite identification and local linear convergence indeed occur. The positions of the *green dots* indicate the iteration from which x_k numerically identifies the submanifold \mathcal{M}_{x^*} . The solid lines (“P”) represents practical observations, while the dashed lines (“T”) denotes theoretical predictions.

As the Riemannian Hessians of ℓ_0 and the rank both vanish in all examples, our predicted rates coincide exactly with the observed ones (same slopes for the dashed and solid lines).

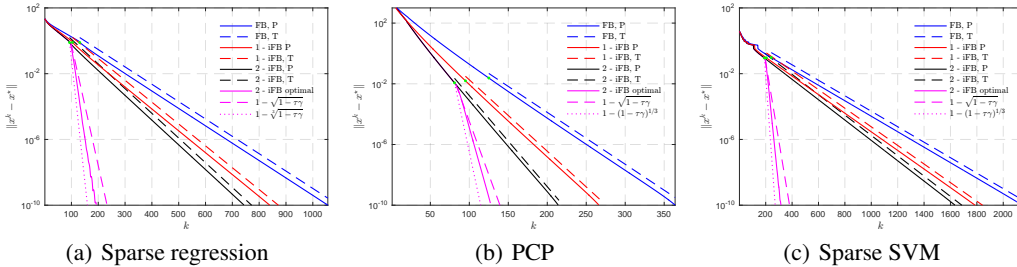


Figure 1: Finite identification and local linear convergence of MiFB under different inertial settings in terms of $\|x_k - x^*\|$. “P” stands for practical observation and “T” indicates the theoretical estimate. We fix $\gamma_k \equiv 0.3/L$ for all tests. For the 2 inertial schemes, inertial parameters are first chosen such that (2.3) holds. The position of the green dot in each plot indicates the iteration beyond which identification of \mathcal{M}_{x^*} occurs.

Comparison of the methods Under the tested settings, we draw the following remarks on the comparison of the inertial schemes:

- The MiFB scheme is much faster than FB both globally and locally. Finite activity identification also occurs earlier for MiFB than for FB;
- Comparing the two MIFB inertial schemes, “2-iFB” outperforms “1-iFB”, showing the advantages of a 2-step inertial scheme over the 1-step one.

Optimal first-order method To highlight the potential of multiple steps in MiFB, for the “2-iFB” scheme, we also added an example where we locally optimized the rate for the inertial parameters. See the *magenta* lines all the examples, where the solid line corresponds to the observed profile for the optimal inertial parameters, the *dashed* line stands for the rate $1 - \sqrt{1 - \tau\gamma}$, and the *dotted* line is that of $1 - \sqrt[3]{1 - \tau\gamma}$, which shows indeed that a faster linear rate can be obtained owing to multiple inertial parameters.

We refer to [20, Section 4.5] for the optimal choice of inertial parameters for the case $s = 1$.

The empirical bound (2.4) and inertial steps s We now present a short comparison of the empirical bound (2.4) of inertial parameters and different choices of s under bigger choice of $\gamma = 0.8/L$. MiFB with 3 inertial steps, *i.e.* $s = 3$, is added which is noted as “3-iFB”, see the *magenta* line in Figure 2.

Similar to the above experiments, we choose $b_{i,k} = a_{i,k} \equiv a_i, i \in I$, and “Thm 2.2” means that a_i ’s are chosen according to Theorem 2.2, while “Bnd (2.4)” means that a_i ’s are chosen based on the empirical bound (2.4). We can infer from Figure 2 the following. Compared to the results in Figure 1, a bigger choice of γ leads to faster convergence. Yet still, under the same choice of γ , MiFB is faster than FB both locally and globally; For either “Thm 2.2” or “Bnd (2.4)”, the performance of the three MiFB schemes are close, this is mainly due to the fact that values of the sum $\sum_{i \in I} a_i$ for each scheme are close. Then between “Thm 2.2” and “Bnd (2.4)”, “Bnd (2.4)” shows faster convergence result, since the allowed value of $\sum_{i \in I} a_i$ of (2.4) is bigger than that of Theorem 2.2. It should be noted that, when $\gamma \in]0, \frac{2}{3L}]$, the largest value of $\sum_{i \in I} a_i$ allowed by (2.4) is 1. If we choose $\sum_{i \in I} a_i$ equal or very close to 1, then it can be observed in practice that MiFB locally oscillates, which is a well-known property of the FISTA scheme [5, 10]. We refer to [20, Section 4.4] for discussions of the properties of such oscillation behaviour.

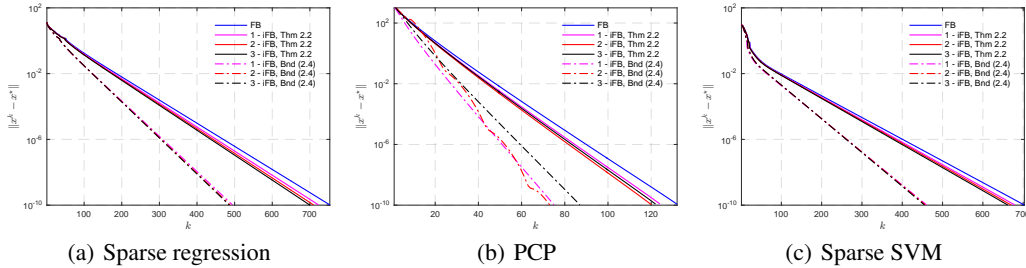


Figure 2: Comparison of MiFB under different inertial settings. We fix $\gamma_k \equiv 0.8/L$ for all tests. For the three inertial schemes, the inertial parameters were chosen such that (2.3) holds.

Acknowledgments

This work was partly supported by the European Research Council (ERC project SIGMA-Vision).

References

- [1] F. Alvarez. On the minimizing property of a second order dissipative system in Hilbert spaces. *SIAM Journal on Control and Optimization*, 38(4):1102–1119, 2000.
- [2] F. Alvarez and H. Attouch. An inertial proximal method for maximal monotone operators via discretization of a nonlinear oscillator with damping. *Set-Valued Analysis*, 9(1-2):3–11, 2001.
- [3] 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. *Mathematical Programming*, 137(1-2):91–129, 2013.
- [4] H. Attouch, J. Peypouquet, and P. Redont. A dynamical approach to an inertial Forward–Backward algorithm for convex minimization. *SIAM J. Optim.*, 24(1):232–256, 2014.

- [5] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [6] J. Bolte, A. Daniilidis, and A. Lewis. The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. *SIAM Journal on Optimization*, 17(4):1205–1223, 2007.
- [7] J. Bolte, A. Daniilidis, O. Ley, and L. Mazet. Characterizations of Łojasiewicz inequalities: subgradient flows, talweg, convexity. *Transactions of the American Mathematical Society*, 362(6):3319–3363, 2010.
- [8] R. I. Boş, E. R. Csetnek, and S. C. László. An inertial Forward–Backward algorithm for the minimization of the sum of two nonconvex functions. *EURO Journal on Computational Optimization*, pages 1–23, 2014.
- [9] E. J. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis? *Journal of the ACM (JACM)*, 58(3):11, 2011.
- [10] A. Chambolle and C. Dossal. On the convergence of the iterates of the “Fast Iterative Shrinkage/Thresholding Algorithm”. *Journal of Optimization Theory and Applications*, pages 1–15, 2015.
- [11] D. L. Donoho, M. Elad, and V. N. Temlyakov. Stable recovery of sparse overcomplete representations in the presence of noise. *IEEE Trans. Inform. Theory*, 52(1):6–18, 2006.
- [12] D. Drusvyatskiy and A. S. Lewis. Optimality, identifiability, and sensitivity. *Mathematical Programming*, pages 1–32, 2013.
- [13] P. Frankel, G. Garrigos, and J. Peypouquet. Splitting methods with variable metric for kurdyka–łojasiewicz functions and general convergence rates. *Journal of Optimization Theory and Applications*, 165(3):874–900, 2015.
- [14] H. Y. Le. Generalized subdifferentials of the rank function. *Optimization Letters*, 7(4):731–743, 2013.
- [15] A. S. Lewis. Active sets, nonsmoothness, and sensitivity. *SIAM J. on Optimization*, 13(3):702–725, 2003.
- [16] A. S. Lewis and J. Malick. Alternating projections on manifolds. *Mathematics of Operations Research*, 33(1):216–234, 2008.
- [17] A. S. Lewis and H. S. Sendov. Twice differentiable spectral functions. *SIAM Journal on Matrix Analysis and Applications*, 23(2):368–386, 2001.
- [18] A. S. Lewis and S. Zhang. Partial smoothness, tilt stability, and generalized Hessians. *SIAM Journal on Optimization*, 23(1):74–94, 2013.
- [19] J. Liang, J. Fadili, and G. Peyré. Local linear convergence of Forward–Backward under partial smoothness. In *Advances in Neural Information Processing Systems*, pages 1970–1978, 2014.
- [20] J. Liang, J. Fadili, and G. Peyré. Activity identification and local linear convergence of Forward–Backward-type methods. arXiv:1503.03703, 2015.
- [21] D. A. Lorenz and T. Pock. An inertial Forward–Backward algorithm for monotone inclusions. *Journal of Mathematical Imaging and Vision*, 51(2):311–325, 2014.
- [22] A. Moudafi and M. Oliny. Convergence of a splitting inertial proximal method for monotone operators. *Journal of Computational and Applied Mathematics*, 155(2):447–454, 2003.
- [23] P. Ochs, Y. Chen, T. Brox, and T. Pock. iPiano: inertial proximal algorithm for nonconvex optimization. *SIAM Journal on Imaging Sciences*, 7(2):1388–1419, 2014.
- [24] B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- [25] R. T. Rockafellar and R. Wets. *Variational analysis*, volume 317. Springer Verlag, 1998.
- [26] L. van den Dries. *Tame topology and o-minimal structures*, volume 248 of *Mathematical Society Lecture Notes*. Cambridge University Press, New York, 1998.