

# Local Linear Convergence of Inertial Forward–Backward Splitting for Low Complexity Regularization

Jingwei Liang and Jalal M. Fadili

CNRS, GREYC, ENSICAEN, Université de Caen  
Email: {Jingwei.Liang, Jalal.Fadili}@ensicaen.fr

Gabriel Peyré

CNRS, Ceremade, Université Paris-Dauphine  
Email: Gabriel.Peyre@ceremade.dauphine.fr

**Abstract**—In this abstract, we consider the inertial Forward-Backward (iFB) splitting method and its special cases (Forward-Backward/ISTA and FISTA). Under the assumption that the non-smooth part of the objective is *partly smooth* relative to an active smooth manifold, we show that iFB-type methods (i) identify the active manifold in finite time, then (ii) enter a local linear convergence regime that we characterize precisely. This gives a grounded and unified explanation to the typical behaviour that has been observed numerically for many low-complexity regularizers, including  $\ell_1$ ,  $\ell_{1,2}$ -norms, total variation (TV) and nuclear norm to name a few. The obtained results are illustrated by concrete examples.

## I. INTRODUCTION

Consider the following structured optimization problem

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

where  $J \in \Gamma_0(\mathbb{R}^n)$ , the set of proper, lower semi-continuous and convex functions,  $F$  is convex,  $C^{1,1}(\mathbb{R}^n)$  with  $\nabla F$  being  $\beta$ -Lipschitz continuous. We assume that  $\text{Argmin } \Phi \neq \emptyset$ .

In this paper, we consider a generic form of inertial Forward-Backward for solving  $(\mathcal{P})$  which reads,

$$\begin{aligned} y_a^k &= x^k + a_k(x^k - x^{k-1}), & y_b^k &= x^k + b_k(x^k - x^{k-1}), \\ x^{k+1} &= \text{Prox}_{\gamma_k J}(y_a^k - \gamma_k \nabla F(y_b^k)), \end{aligned} \quad (\text{I.1})$$

where  $a_k \in [0, \bar{a}]$  and  $b_k \in [0, \bar{b}]$ ,  $(\bar{a}, \bar{b}) \in [0, 1]^2$ , and the step-size  $0 < \underline{\gamma} \leq \gamma_k \leq \bar{\gamma} < 2/\beta$ . For  $\gamma > 0$ , the proximity operator is defined as  $\text{Prox}_{\gamma J}(x) = \arg\min_{z \in \mathbb{R}^n} \frac{1}{2\gamma} \|z - x\|^2 + J(z)$ .

iFB (I.1) covers various special cases in the literature, including the (unrelaxed) Forward-Backward (FB) [1] and FISTA [2]. In the original FISTA, only convergence of the objective function is guaranteed. Recently in [5], the iterates are proved to be convergent under  $a_k = b_k = (t_{k-1} - 1)/t_k$  where  $t_k = (k + p - 1)/p$ ,  $p \geq 2$ .

## II. PARTLY SMOOTH FUNCTIONS AND FINITE IDENTIFICATION

The class of partly smooth functions [3], is specialized here to functions in  $\Gamma_0(\mathbb{R}^n)$ . Denote  $\text{par}(C)$  the linear subspace parallel to the non-empty convex set  $C \subset \mathbb{R}^n$ , and  $\text{ri}(C)$  its relative interior.

**Definition II.1.** Let  $J \in \Gamma_0(\mathbb{R}^n)$  and  $x \in \mathbb{R}^n$  such that  $\partial J(x) \neq \emptyset$ .  $J$  is *partly smooth* at  $x$  relative to a set  $\mathcal{M}$  containing  $x$  if

- (Smoothness)  $\mathcal{M}$  is a  $C^2$ -manifold,  $J|_{\mathcal{M}}$  is  $C^2$  around  $x$ ;
- (Sharpness) The tangent space  $\mathcal{T}_{\mathcal{M}}(x) = T_x \stackrel{\text{def}}{=} \text{par}(\partial J(x))^\perp$ ;
- (Continuity) The  $\partial J$  is continuous at  $x$  relative to  $\mathcal{M}$ .

Examples of such functions are given in Section IV, see also [4].

**Theorem II.2 (Finite activity identification).** Suppose  $x^k$  converges to a minimizer  $x^*$  of  $(\mathcal{P})$  such that  $J$  is partly smooth at  $x^*$  relative to  $\mathcal{M}_{x^*}$ , and

$$-\nabla F(x^*) \in \text{ri}(\partial J(x^*)), \quad (\text{II.1})$$

then there exists a  $K > 0$  such that for all  $k \geq K$ ,  $x^k \in \mathcal{M}_{x^*}$ . If moreover  $\mathcal{M}_{x^*}$  is affine/linear, then  $y_a^k, y_b^k \in \mathcal{M}_{x^*}$  for  $k > K$ .

Condition (II.1) can be viewed as a geometric generalization of the strict complementarity of non-linear programming, and is almost necessary for the finite identification of  $\mathcal{M}_{x^*}$  [3].

## III. LOCAL LINEAR CONVERGENCE

We now turn to the local linear convergence of the iFB-type methods with partly smooth functions. For space limitations, we mainly focus on the case where  $a_k = b_k$ , and denote  $d^{k+1} = \begin{pmatrix} x^{k+1} - x^* \\ x^k - x^* \end{pmatrix}$ .

**Theorem III.1.** We assume the conditions of Theorem II.2 hold. If moreover  $F$  is  $C^2$  near  $x^*$  and there exists  $\alpha \geq 0$  such that  $\text{P}_{T_{x^*}} \nabla^2 F(x^*) \text{P}_{T_{x^*}} \succ \alpha \text{Id}$ . Then for all  $k$  large enough, we have

- 1)  $Q$ -linear rate: if  $0 < \underline{\gamma} \leq \gamma_k \leq \bar{\gamma} < \min(2\alpha\beta^{-2}, 2\beta^{-1})$ , then given any  $\rho$  such that  $1 > \rho \geq \tilde{\rho}_k$ , the iterates satisfy

$$\|x^{k+1} - x^*\|^2 \leq \|d^{k+1}\|^2 \leq \rho \|d^k\|^2,$$

where  $\eta = \max\{q(\underline{\gamma}), q(\bar{\gamma})\} \in [0, 1]$ ,  $q(\gamma) = 1 - 2\alpha\gamma + \beta^2\gamma^2$ ,

$$\tilde{\rho}_k = \begin{cases} \frac{(1+a_k)\eta + \sqrt{(1+a_k)^2\eta^2 - 4a_k\eta}}{2}, & \eta \in [\frac{4a_k}{(1+a_k)^2}, 1], \\ \sqrt{a_k\eta}, & \eta \in [0, \frac{4a_k}{(1+a_k)^2}]. \end{cases}$$

- 2)  $R$ -linear rate: if  $\mathcal{M}_{x^*}$  is affine/linear, then

$$\|x^{k+1} - x^*\|^2 \leq \|d^{k+1}\|^2 \leq \rho_k \|d^k\|^2,$$

where  $\rho_k \in [0, 1]$

$$\rho_k = \begin{cases} \frac{|(1+a_k)\eta_k| + \sqrt{(1+a_k)^2\eta_k^2 - 4a_k\eta_k}}{2}, & \eta_k \in ]-1, 0] \cup [\frac{4a_k}{(1+a_k)^2}, 1], \\ \sqrt{a_k\eta_k}, & \eta_k \in [0, \frac{4a_k}{(1+a_k)^2}], \end{cases}$$

and  $\eta_k \in ]-1, 1[$  is an eigenvalue of  $\text{Id} - \gamma_k \text{P}_T \int_0^1 \nabla^2 F(x^* + t(y_a^k - x^*)) dt \text{P}_T$ .

## IV. NUMERICAL EXPERIMENTS

**Example IV.1 ( $\ell_1$ -norm).** The  $\ell_1$ -norm is partly smooth relative to  $\mathcal{M} = T_x = \{u \in \mathbb{R}^n : \text{supp}(u) \subseteq \text{supp}(x)\}$ .

**Example IV.2 ( $\ell_{1,2}$ -norm).**  $\ell_{1,2}$ -norm is partly smooth relative to  $\mathcal{M} = T_x = \{u \in \mathbb{R}^n : \text{supp}_B(u) \subseteq \text{supp}_B(x)\}$ , where  $\text{supp}_B(x) = \bigcup\{b : x_b \neq 0\}$ , and  $\bigcup_{b \in B} b = \{1, \dots, n\}$ .

**Example IV.3 (TV semi-norm).** The TV semi-norm  $\|x\|_{\text{TV}} = \|\nabla x\|_1$  is partly smooth relative to the subspace  $\mathcal{M} = T_x = \{u \in \mathbb{R}^n : \text{supp}(\nabla u) \subseteq I\}$ ,  $I = \text{supp}(\nabla x)$ .

**Example IV.4 (Nuclear norm).** The nuclear norm is partly smooth relative to the manifold of fixed rank matrices,  $\mathcal{M} = \{z \in \mathbb{R}^{n_1 \times n_2} : \text{rank}(z) = r\}$ .

We now consider the problem  $\min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - Ax\|^2 + \lambda J(x)$ , where  $y \in \mathbb{R}^m$  is the observation,  $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is drawn from the standard Gaussian ensemble, and  $\lambda > 0$  is the regularization parameter. The convergence profiles are depicted in Figure 1.

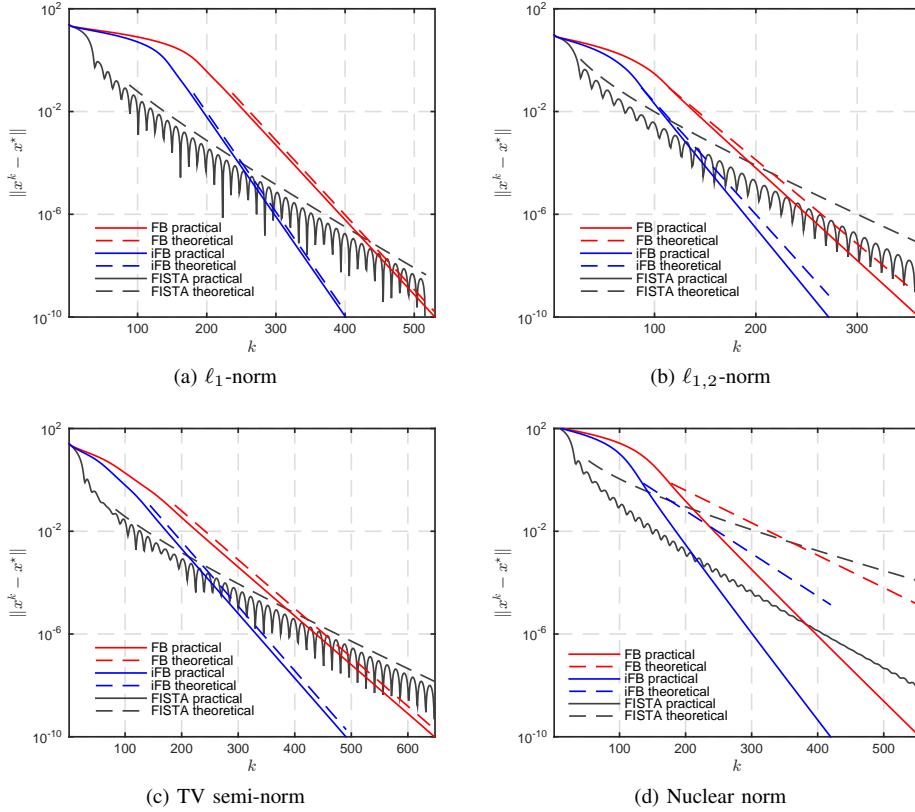


Fig. 1: Local linear convergence of iFB-type methods in terms of  $\|x^k - x^*\|$ . The forward model of the problem of interests reads  $y = Ax_0 + \varepsilon$ ,  $\varepsilon \sim \mathcal{N}(0, \delta^2)$ . (a)  $\ell_1$ -norm,  $(m, n) = (48, 128)$ ,  $x_0$  is 8-sparse; (b)  $\ell_{1,2}$ -norm,  $(m, n) = (60, 128)$ ,  $x_0$  has 3 non-zero blocks with block-size 4; (c) 1D TV semi-norm,  $(m, n) = (48, 128)$ ,  $\nabla x_0$  is 8-sparse; (d) Nuclear norm,  $(m, n) = (1425, 2500)$ ,  $x_0 \in \mathbb{R}^{50 \times 50}$  and  $\text{rank}(x_0) = 5$ . The red, black and blue lines are respectively the results of FB, FISTA [5] and iFB (with  $a_k = b_k \equiv \sqrt{5} - 2.01$ ). All algorithms were tested with  $\gamma_k \equiv 1/\|A\|^2$ . The solid lines are the practical observed profiles and the dashed ones the theoretical predictions. The beginning of the dashed lines are the points when  $x^k$  identifies the manifold  $\mathcal{M}_{x^*}$ . As one can observe, FISTA has the fastest manifold identification, however, locally it is the slowest for all tested examples. Indeed, when the manifold is affine, it can be shown from Theorem III.1 that  $\rho_k \in ]\eta_k, \sqrt{\eta_k}]$  for  $a_k > \eta_k$ , i.e. FISTA is locally slower than FB.

#### ACKNOWLEDGMENT

This work has been partly supported by the European Research Council (ERC project SIGMA-Vision). JF was partly supported by Institut Universitaire de France.

#### REFERENCES

- [1] P. L. Lions and B. Mercier, *Splitting algorithms for the sum of two nonlinear operators*, SIAM Journal on Numerical Analysis, 16(6):964–979, 1979.
- [2] 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.
- [3] A. S. Lewis, *Active sets, nonsmoothness, and sensitivity*, SIAM Journal on Optimization, 13(3):702–725, 2003.
- [4] S. Vaiter, G. Peyré, and M. J. Fadili, *Partly smooth regularization of inverse problems*, arXiv:1405.1004, 2014.
- [5] A. Chambolle and D. Dossal, *How to make sure the iterates of FISTA converge*, 2014.