



**HAL**  
open science

# FISTA is an automatic geometrically optimized algorithm for strongly convex functions

J-F Aujol, Charles Dossal, Aude Rondepierre

► **To cite this version:**

J-F Aujol, Charles Dossal, Aude Rondepierre. FISTA is an automatic geometrically optimized algorithm for strongly convex functions. *Mathematical Programming, Series A*, 2023, 204, pp.449-491. 10.1007/s10107-023-01960-6 . hal-03491527v2

**HAL Id: hal-03491527**

**<https://hal.science/hal-03491527v2>**

Submitted on 13 Nov 2024

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

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

Public Domain

# FISTA is an automatic geometrically optimized algorithm for strongly convex functions

J-F. Aujol\*      Ch. Dossal†      A. Rondepierre‡

## Abstract

In this work, we are interested in the famous FISTA algorithm. We show that FISTA is an automatic geometrically optimized algorithm for functions satisfying a quadratic growth assumption. This explains why FISTA works better than the standard Forward-Backward algorithm (FB) in such a case, although FISTA is known to have a polynomial asymptotic convergence rate while FB is exponential. We provide a simple rule to tune the  $\alpha$  parameter within the FISTA algorithm to reach an  $\varepsilon$ -solution with an optimal number of iterations. These new results highlight the efficiency of FISTA algorithms, and they rely on new non asymptotic bounds for FISTA.

**Keywords:** Nesterov acceleration, ODE, first order scheme, optimization.

## 1 Introduction

Let  $F = f + h$  be a composite convex function defined from  $\mathbb{R}^N$  to  $\mathbb{R}$ , such that  $\nabla f$  is  $L$ -Lipschitz, and the proximal operator of  $h$  can be easily computed. Throughout the paper, the function  $F$  is assumed to satisfy a quadratic growth property  $\mathcal{G}_\mu^2$  for some parameter  $\mu > 0$  i.e.:

$$\exists \mu > 0, \forall x \in \mathbb{R}^N, F(x) - F^* \geq \frac{\mu}{2} d(x, X^*)^2 \quad (1)$$

where:  $F^* = \inf F$  and where  $X^*$  is the set of minimizers of  $F$ . Moreover in the new results presented in this article it is assumed that the function  $F$  has a unique minimizer and this growth condition becomes a strong minimizer property :

$$\exists \mu > 0, \forall x \in \mathbb{R}^N, F(x) - F^* \geq \frac{\mu}{2} \|x - x^*\|^2. \quad (2)$$

These two conditions are actually in force when  $F$  is  $\mu$ -strongly convex, but they are more general as will be shown later. Under these hypotheses, the Forward-Backward algorithm (FB) defines a sequence  $(x_n)_{n \geq 0}$  satisfying  $F(x_n) - F^* = \mathcal{O}(e^{-\frac{\mu}{L}n})$ , i.e the number  $n_\varepsilon^{FB}$  of iterations necessary to get an  $\varepsilon$ -solution satisfies  $n_\varepsilon^{FB} = \mathcal{O}(\frac{L}{\mu} \log(\frac{1}{\varepsilon}))$ .

Several accelerations of FB have been proposed when  $F$  is convex or strongly convex. The FISTA algorithm proposed by Beck and Teboulle [11] using an inertial scheme from Yurii Nesterov [21], was built to improve the convergence rate of FB under a convexity assumption. Nesterov and many others [21, 27, 9, 8, 26, 29, 24, 23, 15] proposed various schemes dedicated to the strongly convex case. All these schemes necessitate an estimation  $\tilde{\mu}$  of  $\mu$  such that  $\tilde{\mu} \leq \mu$ , and they reach an  $\varepsilon$ -solution with at most  $n_\varepsilon = \mathcal{O}\left(\sqrt{\frac{L}{\tilde{\mu}}} \log\left(\frac{1}{\varepsilon}\right)\right)$  iterations. The exact decay rate depends on the considered algorithm and on the precise geometrical assumptions made on the objective  $F$ , see Table 1 for details.

---

\*Univ. Bordeaux, CNRS, Bordeaux INP, IMB, UMR 5251, F-33400 Talence, France. Jean-Francois.Aujol@math.u-bordeaux.fr

†IMT, Univ. Toulouse, INSA Toulouse, France. dossal@insa-toulouse.fr

‡IMT, Univ. Toulouse, INSA Toulouse, France. aude.ronddepierre@insa-toulouse.fr

	Geometry of $F$	References	Convergence rate for $F(x_n) - F^*$	Number of iterations to reach an $\varepsilon$ solution
FB	Convex	[23, 11]	$\frac{2L\ x_0 - x^*\ ^2}{n}$	$\frac{4L^2}{\varepsilon^2} \ x_0 - x^*\ ^2$
FISTA with $\alpha = 3$	Convex	[23, 11, 14]	$\frac{2L\ x_0 - x^*\ ^2}{(n+1)^2}$	$\frac{2L}{\varepsilon} \ x_0 - x^*\ $
FB	Convex and $\mathcal{G}_\mu^2$	[17]	$(1+\kappa)^{-n} (F(x_0) - F^*)$	$\mathcal{O}\left(\frac{1}{\kappa} \log\left(\frac{1}{\varepsilon}\right)\right)$
NSC	Strongly convex Requires an estimation of $\mu$	[23]	$2(1 - \sqrt{\kappa})^n (F(x_0) - F^*)$	$\mathcal{O}\left(\frac{1}{\sqrt{\kappa}} \log\left(\frac{1}{\varepsilon}\right)\right)$
$N\mu$	Strongly convex	[22]	$\mathcal{O}\left(e^{-\frac{2}{ \log \kappa } \sqrt{\kappa} n}\right)$	$\mathcal{O}\left(\frac{ \log(\kappa) }{\sqrt{\kappa}} \log\left(\frac{1}{\varepsilon}\right)\right)$
FISTA $\alpha \geq 3$	Convex and $\mathcal{G}_\mu^2$ Uniqueness of the minimizer	[5, 6]	$\mathcal{O}\left(n^{-\frac{2\alpha}{3}}\right)$	Unknown
FISTA $\alpha = 3 \log\left(\frac{3\sqrt{LM_0}}{e\sqrt{2\varepsilon}}\right)$	Convex and $\mathcal{G}_\mu^2$ Uniqueness of the minimizer	ADR	$\mathcal{O}\left(e^{-Cn\sqrt{\kappa}}\right)$	$\mathcal{O}\left(\frac{1}{\sqrt{\kappa}} \log\left(\frac{1}{\varepsilon}\right)\right)$
Optimal FISTA restart	Strongly convex Requires an estimation of $\mu$	[20]	$\mathcal{O}\left(e^{-\frac{1}{\varepsilon} \sqrt{\kappa} n}\right)$	$\mathcal{O}\left(\frac{1}{\sqrt{\kappa}} \log\left(\frac{1}{\varepsilon}\right)\right)$
FISTA restart	Convex and $\mathcal{G}_\mu^2$	[10]	$\mathcal{O}\left(e^{-\frac{1}{12} \sqrt{\kappa} n}\right)$	$\mathcal{O}\left(\frac{1}{\sqrt{\kappa}} \log\left(\frac{1}{\varepsilon}\right)\right)$

Table 1: Convergence rates for Forward-Backward, FISTA, the Nesterov variant for strongly convex functions and FISTA restart depending on the geometry of the objective function  $F$ , and their interpretation in terms of  $\varepsilon$ -solution

In a high dimensional setting  $\kappa := \frac{\mu}{L} \ll 1$  and the square root gain may be crucial. It explains why inertial algorithms are widely used, and why they behave numerically better than FB, especially when  $\tilde{\mu}$  is close to  $\mu$ .

The main contribution of this paper is to show that the version of FISTA proposed by Chambolle and Dossal [14] and Su, Boyd and Candès [28] can reach an  $\varepsilon$ -solution with at most  $n_\varepsilon^{FISTA} = \mathcal{O}\left(\sqrt{\frac{\mu}{L}} \log\left(\frac{1}{\varepsilon}\right)\right)$  iterations under some quadratic growth condition, without any estimation of  $\mu$ . This result applies to the LASSO problem which may not be strongly convex and for which the estimation of the growth parameter  $\mu$  can not be tackled easily. This bound on  $n_\varepsilon^{FISTA}$  especially explains the better performance of FISTA comparing to FB on problems such like the LASSO problem on which FB is known to reach an exponential rate and where the previous bounds on FISTA indicate that its rate was only polynomial.

Even though FISTA was built to produce acceleration in a convex setting, it has two main advantages comparing to algorithms built for strongly convex functions:

- There is no need to estimate the growth parameter  $\mu$  and the convergence rate does not suffer from any underestimation of the growth parameter  $\mu$ . Indeed, the bound on  $n_\varepsilon^{FISTA}$  depends on the true  $\mu$  and not on an estimation  $\tilde{\mu}$  of the true  $\mu$ .
- The bound on  $n_\varepsilon^{FISTA}$  applies under a quadratic growth condition which is a weaker assumption than strong convexity and extends the field of applications of FISTA.

In Section 2, notations and definitions are given, the various algorithms and the notion of  $\varepsilon$ -solutions are detailed. The main contribution, Theorem 3, is stated and the comparison with the state of the art is done. Section 3 is devoted to the non asymptotic study of the FISTA algorithm which provides finite time bounds allowing to prove Theorem 3.

## 2 Number of iterations to reach an $\varepsilon$ -solution

### 2.1 Framework and notations

In this paper we focus on the class of composite functions:  $F = f + h$  where  $f$  is a convex, differentiable function having a  $L$ -Lipschitz gradient and  $h$  is a proper lower semicontinuous (l.s.c.) convex function whose proximal operator is known. The proximal operator of  $h$  is denoted by  $\text{prox}_h$  and defined by:

$$\text{prox}_h(x) = \operatorname{argmin}_{y \in \mathbb{R}^N} \left( h(y) + \frac{1}{2} \|y - x\|^2 \right). \quad (3)$$

For this class of functions a classical minimization algorithm is the Forward-Backward algorithm (FB) whose iterations are described by:

$$x_{n+1} = \text{prox}_{sh}(x_n - s\nabla f(x_n)), \quad s \in \left(0, \frac{2}{L}\right). \quad (4)$$

Beck and Teboulle, based on the ideas of the Nesterov acceleration, propose an accelerated version FISTA (Fast Iterative Shrinkage-Thresholding Algorithm) [11]:

$$y_n = x_n + \frac{t_n - 1}{t_{n+1}}(x_n - x_{n-1}), \quad x_{n+1} = \text{prox}_{sh}(y_n - s\nabla f(y_n)). \quad (5)$$

where the sequence  $(t_n)_{n \in \mathbb{N}}$  is recursively defined by:  $t_1 = 1$  and  $t_{n+1} = (1 + \sqrt{1 + 4t_n^2})/2$ . In this paper we consider the variant of FISTA proposed by Chambolle and Dossal in [14] and also denoted by FISTA:

$$y_n = x_n + \frac{n}{n + \alpha}(x_n - x_{n-1}), \quad x_{n+1} = \text{prox}_{sh}(y_n - s\nabla f(y_n)) \quad (6)$$

that ensures in addition the weak convergence of the iterates (when  $\alpha > 3$ ). We assume moreover that the function  $F$  satisfies some global quadratic growth property  $\mathcal{G}_\mu^2$ . Classically the quadratic

growth condition  $\mathcal{G}_\mu^2$  can be seen as a relaxation of the strong convexity, and it is equivalent in the convex setting to a global Łojasiewicz property with an exponent  $\frac{1}{2}$  [18, 19, 12].

For this class of functions, the Forward-Backward algorithm ensures an exponential decay whereas FISTA classically only ensures a polynomial asymptotic convergence rate. The main contribution of this paper is to provide a non-asymptotic analysis of the FISTA algorithm and to compare the convergence rate in finite time to state-of-the-art algorithms like Forward-Backward and the Nesterov accelerated algorithm for strongly convex functions [23]. Analyzing these algorithms in finite time provides a different insight on these convergence rates.

More precisely, let  $\varepsilon > 0$  be the expected accuracy. The minimizers of a composite function  $F$  can be characterized by the optimality condition  $0 \in \partial F(x)$ , or equivalently  $g(x) = 0$  where:

$$g(x) = L(x - x^+) := L\left(x - \operatorname{prox}_{\frac{1}{L}h}\left(x - \frac{1}{L}\nabla f(x)\right)\right), \quad x \in \mathbb{R}^N, \quad (7)$$

denotes the composite gradient mapping and:  $x^+ := \operatorname{prox}_{\frac{1}{L}h}\left(x - \frac{1}{L}\nabla f(x)\right)$ . This last formulation is convenient for defining an approximate solution to the composite problem, and thus to deduce a tractable stopping criterion for a dedicated optimization algorithm:

**Definition 1 ( $\varepsilon$ -solution)** *Let  $\varepsilon$  be the expected accuracy. The iterate  $x_n$  is said to be an  $\varepsilon$ -solution of the problem  $\min_{x \in \mathbb{R}^N} F(x)$  if:*

$$\|g(x_n)\| \leq \varepsilon. \quad (8)$$

Observe that in the differentiable case (i.e. when  $h = 0$ ), we have:  $g(x) = \nabla f(x)$  so that an  $\varepsilon$ -solution is nothing more than an iterate  $x_n$  satisfying:

$$\|g(x_n)\| = \|\nabla F(x_n)\| \leq \varepsilon. \quad (9)$$

The notion of  $\varepsilon$ -solution can be seen as a good stopping criterion for an algorithm solving the composite optimization problem for mainly three reasons: first it is numerically quantifiable. Secondly controlling the norm of the composite gradient mapping is roughly equivalent to having a control on the values of the objective function. Indeed using the following property of the composite gradient mapping proven in [22, Theorem 1] and [10]:

$$\forall x \in \mathbb{R}^N, \quad \frac{1}{2L}\|g(x)\|^2 \leq F(x) - F(x^+), \quad (10)$$

we can prove that the composite gradient mapping is controlled by the values of the objective function:

$$\forall x \in \mathbb{R}^N, \quad \frac{1}{2L}\|g(x)\|^2 \leq F(x) - F^*. \quad (11)$$

Conversely, as shown in [10, Lemma 3.1], we also have:

$$\forall x \in \mathbb{R}^N, \quad F(x^+) - F^* \leq \frac{2}{\mu}\|g(x)\|^2. \quad (12)$$

Thirdly, using (8) as a stopping criterion will enable us to analyze and compare algorithms in terms of the number of iterations needed to reach a given accuracy  $\varepsilon$ .

## 2.2 Analysing state-of-the-art optimization algorithms in terms of $\varepsilon$ -solution

### 2.2.1 FB and FISTA without the growth condition $\mathcal{G}_\mu^2$

Let us first recall the Forward-Backward algorithm (FB) described by Algorithm 1:

---

**Algorithm 1** FB: Forward-Backward algorithm to minimize  $F = f + h$ .

---

- *Initialization:*  $x_0 \in \mathbb{R}^N$ ,  $\varepsilon > 0$ .
- *Iterations* ( $n \geq 0$ ): update  $x_n$  as follows:

$$x_{n+1} = \text{prox}_{\frac{1}{L}h} \left( x_n - \frac{1}{L} \nabla f(x_n) \right) \quad (13)$$

until  $\|g(x_n)\| \leq \varepsilon$ .

---

The FB algorithm provides the following bound when  $F$  is convex [23, 11]:

$$F(x_n) - F^* \leq \frac{2L\|x_0 - x^*\|^2}{n}. \quad (14)$$

Using (11), this implies that a number of iterations of the order  $\mathcal{O}\left(\frac{L^2}{\varepsilon^2}\right)$  is required to get an  $\varepsilon$ -solution.

A. Beck and M. Teboulle propose in [11] an accelerated version of FB, known as FISTA (Fast Iterative Shrinkage-Thresholding Algorithm). In this paper we focus on the version proposed by Chambolle and Dossal [14] and Su, Boyd and Candes [28] and simply called FISTA from now on.

---

**Algorithm 2** FISTA: Nesterov accelerated algorithm for convex functions  $F = f + h$

---

- *Initialization:*  $x_0 \in \mathbb{R}^N$ ,  $x_{-1} = x_0$ ,  $\varepsilon > 0$ ,  $\alpha \geq 3$ .
- *Iterations* ( $n \geq 0$ ): update  $x_n$  and  $y_n$  as follows:

$$\begin{cases} y_n = x_n + \frac{n}{n+\alpha}(x_n - x_{n-1}) \\ x_{n+1} = \text{prox}_{\frac{1}{L}h}(y_n - \frac{1}{L}\nabla f(y_n)) \end{cases} \quad (15)$$

until  $\|g(x_n)\| \leq \varepsilon$ .

---

The FISTA algorithm provides the following bound when  $F$  is convex [23, 11, 14] and  $\alpha \geq 3$ :

$$F(x_n) - F^* \leq \frac{L(\alpha-1)^2\|x_0 - x^*\|^2}{2(n+\alpha-2)^2}. \quad (16)$$

Using (11), this implies that a number of iterations of the order  $\mathcal{O}\left(\frac{L}{\varepsilon}\right)$  is required to get an  $\varepsilon$ -solution.

### 2.2.2 FB and FISTA with a quadratic growth condition $\mathcal{G}_\mu^2$

Assume now that  $F$  additionally satisfies some quadratic growth condition  $\mathcal{G}_\mu^2$ :

$$\exists \mu > 0, \forall x \in \mathbb{R}^N, F(x) - F^* \geq \frac{\mu}{2}d(x, X^*)^2 \quad (17)$$

where:  $F^* = \inf F$  and where  $X^*$  is the set of minimizers of  $F$ . Classically the FB method provides then an  $\varepsilon$ -solution in  $\mathcal{O}\left(\frac{1}{\kappa} \log\left(\frac{1}{\varepsilon}\right)\right)$  iterations [17], which is of course much better than the previous case (without the  $\mathcal{G}_\mu^2$  assumption). More precisely, from [17, Theorem 4.1], we can derive the following result:

**Theorem 1** *Let  $F = f + g$  where  $f$  is a convex differentiable function having a  $L$ -Lipschitz gradient for some  $L > 0$ , and  $g$  a proper convex l.s.c. function. Assume additionally that  $F$  satisfies a quadratic growth condition  $\mathcal{G}_\mu^2$  for some real parameter  $\mu > 0$ .*

Let  $\varepsilon > 0$  and:

$$n_\varepsilon^{FB} = \frac{1}{|\log(1 - \kappa)|} \log\left(\frac{2LM_0}{\varepsilon^2}\right) \quad (18)$$

where  $\kappa = \frac{\mu}{L}$  and  $M_0 = F(x_0) - F^*$  denotes the potential energy at initial time. Let  $(x_n)_{n \in \mathbb{R}^N}$  be a sequence of iterates generated by the FB algorithm. If  $n \geq n_\varepsilon^{FB}$  then the iterate  $x_n$  is a  $\varepsilon$ -solution.

To provide bounds on the number of iterations to get an  $\varepsilon$  solution, asymptotic bounds on  $n$  are not sufficient. The dependencies of these bounds on  $\alpha$  and  $\kappa = \frac{\mu}{L}$  are also crucial. The main results of the paper (Theorems 3 and 4) are based on new explicit and non asymptotic bounds developed in Section 3.

### 2.2.3 Algorithms devoted to strongly convex functions

Consider now the Nesterov scheme designed for strongly convex functions [23, Algorithm 2.2.11] which is known to ensure obtaining an  $\varepsilon$ -solution at most in  $\mathcal{O}\left(\frac{1}{\sqrt{\kappa}} \log\left(\frac{1}{\varepsilon}\right)\right)$  iterations [23, Theorem 2.2.3]. More precisely:

---

**Algorithm 3** NSC: Nesterov accelerated algorithm for strongly convex functions

---

- Initialization:  $x_0 \in \mathbb{R}^N$ ,  $x_{-1} = x_0$ ,  $s \in (0, \frac{1}{L})$ .

- Iterations ( $n \geq 0$ ): update  $x_n$  and  $y_n$  as follows:

$$\begin{cases} y_n = x_n + \frac{1 - \sqrt{\kappa}}{1 + \sqrt{\kappa}}(x_n - x_{n-1}) \\ x_{n+1} = y_n - \frac{1}{L}\nabla F(y_n) \end{cases} \quad (19)$$

until  $\|g(x_n)\| \leq \varepsilon$ .

---

**Theorem 2** Let  $F$  be a convex differentiable function having a  $L$ -Lipschitz gradient and admitting a unique minimizer  $x^*$ . Assume additionally that  $F$  is  $\mu$ -strongly convex for some real parameter  $\mu > 0$ . Let  $\varepsilon > 0$ . Then for  $\kappa = \frac{\mu}{L}$  small enough,

$$\forall n \in \mathbb{N}, F(x_n) - F^* \leq 2(1 - \sqrt{\kappa})^n (F(x_0) - F^*),$$

which means that an  $\varepsilon$ -solution can be obtained in at most:

$$n_\varepsilon^{NSC} = \frac{1}{|\log(1 - \sqrt{\kappa})|} \log\left(\frac{4LM_0}{\varepsilon^2}\right). \quad (20)$$

where  $M_0 = F(x_0) - F^*$  denotes the potential energy at initial time.

A crucial point to keep in mind with this result, is that iterations (19) depends on  $\mu$ , and thus,  $\mu$  must be estimated a priori. In practice,  $\mu$  may be unknown and only an estimation  $\tilde{\mu}$  of  $\mu$  can be used to define the sequence  $(x_n)_{n \geq 0}$ . To apply the previous theorem we must have  $\tilde{\mu} \leq \mu$  and the previous bound becomes :

$$n_\varepsilon^{NSC} = \frac{1}{\left|\log\left(1 - \sqrt{\frac{\tilde{\mu}}{L}}\right)\right|} \log\left(\frac{4LM_0}{\varepsilon^2}\right), \quad (21)$$

where  $\tilde{\mu} \leq \mu$  is the one used to define  $\kappa := \frac{\tilde{\mu}}{L}$  in (19). If  $\tilde{\mu} \ll \mu$ , this bound may be much higher than the one given in Theorem 3.

The field of accelerated methods for strongly convex functions is a very active one. The fastest one is the one proposed in [29] (which improves NSC with a 2 factor within the exponential decay).

The references [26] (with the additional hypothesis of the differentiability of the function  $F$ ), and [9] (without additional assumption) propose a  $\sqrt{2}$  factor improvement with respect to NSC. Note that the case of  $F = f + g$  with  $F$  convex satisfying a growth condition  $\mathcal{G}_\mu^2$  for some real parameter  $\mu > 0$  has recently been addressed in [27, 9, 8].

Notice that an accelerated scheme for strongly convex function, where the strong convexity parameter  $\mu$  is automatically estimated, has been proposed by Nesterov in [22]. Such an approach does not suffer from the estimation issue for  $\mu$ . This algorithm is not based on the NSC scheme, and we refer the reader to Section 5.3 of [22] where the detailed algorithm is given. Y. Nesterov shows that the number of iterations to reach an  $\varepsilon$ -solution is at most:

$$n_\varepsilon^{N\mu} = \mathcal{O}\left(\frac{-\log(\kappa)}{\sqrt{\kappa}}\right) + \mathcal{O}\left(\frac{\log(\kappa)}{\sqrt{\kappa}} \log(\kappa\varepsilon)\right), \quad (22)$$

where  $\kappa := \frac{\mu}{L}$ . When  $\varepsilon = o(\kappa)$ , then  $n_\varepsilon^{N\mu} = \mathcal{O}\left(\frac{\log(\kappa)}{\sqrt{\kappa}} \log(\varepsilon)\right)$ .

#### 2.2.4 FISTA restart

Restarting FISTA is another way to get a linear convergence in the case when  $F$  is convex and satisfies the growth condition  $\mathcal{G}_\mu^2$  (see e.g. [20] or [25]). Indeed the inertia introduced by FISTA allows fast convergences, but it also generates oscillations. Restarting FISTA is roughly equivalent to set the inertia to zero at each restart which helps to reduce oscillations. A classical strategy consists in restarting FISTA at regular time [21]. In the differentiable convex case, restarting every  $\lfloor 2e\sqrt{\frac{L}{\mu}} \rfloor$  iterations ensures that  $F(x_k) - F^* = \mathcal{O}\left(e^{-\frac{1}{e}\sqrt{\frac{\mu}{L}k}}\right)$  for the class of  $\mu$ -strongly convex functions [20], and this result can be extended to composite functions satisfying a quadratic growth condition [25, 20]. Note that in these schemes an estimation of the growth parameter is required which can be restrictive in applications since this parameter is rarely known.

In particular, there have been recent works where the growth parameter  $\mu$  is estimated on the fly by the algorithms, see [2, 10, 16]. Adaptive restart schemes take advantage of each iteration to estimate the geometry of  $F$ , and enables to fit the parameters of the algorithm progressively. In [16] Fercoq and Qu provide an adaptive restart scheme for FISTA for the set of functions satisfying  $\mathcal{G}_\mu^2$  which builds a sequence of estimates of the growth parameter  $\mu$ . It requires a prior estimate  $\mu_0$  and the convergence rate of the method is given by

$$\mathcal{O}\left(e^{-\frac{\sqrt{2}-1}{2\sqrt{e}(2-\sqrt{\frac{\mu}{\mu_0}})}\sqrt{\frac{\mu}{L}k}}\right)$$

if  $\mu_0 \geq \mu$ . Note that this rate significantly depends on  $\frac{\mu}{\mu_0}$  and that this method might be less effective if  $\mu_0$  is highly overestimated.

In [2] and [1], Alamo et al. introduce adaptive restart schemes for FISTA under a quadratic growth assumption on  $F$  and without any prior estimate of  $\mu$ . In [2] they provide a fast exponential decay  $\mathcal{O}\left(e^{-\frac{1}{16}\sqrt{\frac{\mu}{L}k}}\right)$ , but the computational cost is quite heavy since their scheme relies on computations of  $F(x_k)$  at least at half of the iterations. On the other hand the scheme introduced in [1] is based on gradient information and therefore it does not require to compute  $F$ . It provides the convergence rate

$$\mathcal{O}\left(e^{-\frac{1}{4e(1+\sqrt{\mu+1})}\frac{\mu}{L}k}\right)$$

which can be seen as a low exponential decay as  $e^{-\frac{1}{4e(1+\sqrt{\mu+1})}\frac{\mu}{L}k} > e^{-\frac{1}{8e}\frac{\mu}{L}k}$  for all  $k > 0$ .

Let us finally focus on a recent convergence result obtained in [10]: a new restart scheme that does not require a priori knowledge of  $\mu$ , is proposed and it provides a fast exponential decay on the values  $F(x_n) - F^*$  in  $\mathcal{O}\left(e^{-\frac{1}{12}\sqrt{\frac{\mu}{L}n}}\right)$ . The authors prove that this method has a fast convergence rate in terms of iterations and computational time as well, and the number of iterations to get an  $\varepsilon$ -solution is bounded by:

$$n_\varepsilon^{Restart} = \frac{7.2}{\sqrt{\kappa}} \left(4.5 + \log\left(1 + 1.3\frac{LM_0}{\varepsilon^2}\right)\right) \quad (23)$$



so that  $n_\varepsilon^{Restart}$  is of order:

$$\mathcal{O}\left(\frac{1}{\sqrt{\kappa}} \log\left(\frac{1}{\varepsilon}\right)\right). \quad (24)$$

### 2.3 FISTA is an automatic geometrically optimized algorithm

FISTA with a varying parameter  $\alpha$  applied to strongly convex functions, or satisfying a quadratic growth condition, has already been studied and is known to have a polynomial decay. In [28, Theorem 9], the authors proved that if  $F = f + h$  and if  $f$  is strongly convex then for  $\alpha \geq \frac{9}{2}$

$$F(x_n) - F^* \leq C(\alpha)L \sqrt{\frac{L}{\mu} \frac{\|x_0 - x^*\|^2}{n^3}} \leq C(\alpha) \left(\frac{\sqrt{L}}{\sqrt{\mu n}}\right)^3 (F(x_0) - F^*). \quad (25)$$

In [5] and in [6], Attouch et al and Aujol et al. proved the following asymptotic decay :

$$F(x_n) - F^* = \mathcal{O}\left(\frac{1}{n^{\frac{2\alpha\gamma}{\gamma+2}}}\right) \quad (26)$$

which coincides with the previous bound (25) if  $\alpha = \frac{9}{2}$  and  $\gamma = 1$ , when  $F$  satisfies some flatness hypothesis  $\mathcal{H}_\gamma$ :

$$\forall x \in \mathbb{R}^N, F(x) - F^* \leq \frac{1}{\gamma} \langle \nabla F(x), x - x^* \rangle. \quad (27)$$

To provide bounds on the number of iterations to get an  $\varepsilon$  solution, asymptotic bounds on  $n$  are not sufficient. The dependencies of these bounds on  $\alpha$  and  $\kappa = \frac{\mu}{L}$  are also crucial. The following Theorem is based on new explicit and non asymptotic bounds developed in Part 3, and given by Theorem 6: for  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$  and  $\kappa$  small enough:

$$\forall n \geq \frac{3\alpha}{\sqrt{\kappa}}, F(x_n) - F^* \leq \frac{9}{4} e^{-2} M_0 \left(\frac{8e}{3n\sqrt{\kappa}} \alpha\right)^{\frac{2\alpha}{3}}. \quad (28)$$

Our main result is to prove that under some quadratic growth condition, the number of iterations of FISTA to reach an  $\varepsilon$ -solution is actually in  $\mathcal{O}\left(\frac{1}{\sqrt{\kappa}} \log\left(\frac{1}{\varepsilon}\right)\right)$ :

**Theorem 3** *Let  $F = f + h$  where  $f$  is a convex differentiable function having a  $L$ -Lipschitz gradient for some  $L > 0$ , and  $h$  a proper convex l.s.c. function. Assume that  $F$  admits a unique minimizer  $x^*$  and satisfies a quadratic growth condition  $\mathcal{G}_\mu^2$  for some real parameter  $\mu > 0$ .*

*Let  $\varepsilon > 0$  and*

$$\alpha_{1,\varepsilon} := 3 \log\left(\frac{3\sqrt{LM_0}}{e\sqrt{2\varepsilon}}\right) \quad \text{and} \quad n_{1,\varepsilon}^{FISTA} := \frac{8e^2}{3\sqrt{\kappa}} \alpha_{1,\varepsilon} = \frac{8e^2}{\sqrt{\kappa}} \log\left(\frac{3\sqrt{LM_0}}{e\sqrt{2\varepsilon}}\right), \quad (29)$$

*where  $\kappa = \frac{\mu}{L}$  and  $M_0 = F(x_0) - F^*$  denotes the potential energy at initial time. Let  $(x_n)_{n \in \mathbb{R}^N}$  be a sequence of iterates generated by the FISTA algorithm with parameter  $\alpha_{1,\varepsilon}$ . There exists  $\kappa_0 \in (0, 1)$  (independent from  $\varepsilon$ ) such that for any  $\kappa \leq \kappa_0$ , if  $n \geq n_{1,\varepsilon}^{FISTA}$  then the iterate  $x_n$  is an  $\varepsilon$ -solution.*

Unlike Algorithm 3, the parameter of FISTA does not depend on  $\mu$  and it follows that  $n_{1,\varepsilon}^{FISTA}$  depends on the *real* value of  $\mu$ , and not on any estimation. To set  $\alpha$ , one only needs to define an accuracy  $\varepsilon$ , the value of  $L$  (which is supposed to be known, also to define the step), and the value of  $M_0$ . The value of  $M_0$  should be chosen to bound the potential energy of the system. In several situations, simple bounds can be found for  $M_0$  for instance when  $F^*$  is known (least square problems) or can be estimated. Moreover since  $M_0$  appears in the logarithm,  $n_{1,\varepsilon}^{FISTA}$  is not very sensitive to  $M_0$ . But we must keep in mind that a bound on  $M_0$  must be given to set  $\alpha$ . It is not surprising that  $\alpha$  depends on the ratio  $\frac{M_0}{\sqrt{\varepsilon}}$  because this ratio measures the decay we want to reach on the value of  $F(x_n) - F^*$ .

Theorem 3 does not provide a simple explicit rule to fix  $\alpha$  according to  $\varepsilon$  since  $M_0$  is never known precisely and since the bound on  $n_{1,\varepsilon}^{FISTA}$  provided by Theorem 3 is probably not very tight. The condition  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$  is also technical and other choices may have been done. Our goal was mainly to provide explicit bound on the number of iterations needed to get an  $\varepsilon$  solution and to prove that for a suitable choice of  $\alpha$ , this number is quite optimal. The computations are technical and we believe that the bounds could be improved but we prefer to avoid more technical computations.

Consider now the subclass of convex differentiable functions ( $h = 0$ ) additionally satisfying a quadratic flatness property  $\mathcal{H}_2$ : for any minimizer  $x^*$ , we assume that:

$$\forall x \in \mathbb{R}^N, F(x) - F^* \leq \frac{1}{2} \langle \nabla F(x), x - x^* \rangle. \quad (30)$$

To have a better intuition of the geometry of functions satisfying  $\mathcal{H}_2$ , observe that the flatness property (30) implies that for any minimizer  $x^*$ , there exists a real constant  $M > 0$  such that:

$$\forall x \in \mathbb{R}^N, F(x) - F^* \leq M \|x - x^*\|^2, \quad (31)$$

see [6, Lemma 2.2]. Thus any convex differentiable function satisfying both  $\mathcal{G}_\mu^2$  and  $\mathcal{H}_2$  for some  $\mu > 0$  can be seen as almost quadratic. It naturally includes quadratic functions and least square problems. For this subclass of functions, Theorem 3 can be improved:

**Theorem 4** *Let  $F$  be a convex differentiable function having a  $L$ -Lipschitz gradient and admitting a unique minimizer  $x^*$ . Assume that  $F$  satisfies a quadratic growth condition  $\mathcal{G}_\mu^2$  for some real parameter  $\mu > 0$  and a flatness assumption  $\mathcal{H}_2$ .*

*Let  $\varepsilon > 0$  and*

$$\alpha_{2,\varepsilon} = 2 \log \left( \frac{4\sqrt{2LM_0}}{3e\varepsilon} \right) \quad \text{and} \quad n_{2,\varepsilon}^{FISTA} = \frac{19e^2}{8\sqrt{\kappa}} \alpha_{2,\varepsilon} = \frac{19e^2}{4\sqrt{\kappa}} \log \left( \frac{4\sqrt{2LM_0}}{3e\varepsilon} \right) \quad (32)$$

where  $\kappa = \frac{\mu}{L}$  and  $M_0 = F(x_0) - F^*$  denotes the potential energy at initial time. Let  $(x_n)_{n \in \mathbb{R}^N}$  be a sequence of iterates generated by the FISTA algorithm with parameter  $\alpha_{2,\varepsilon}$ . There exists  $\kappa_0 \in (0, 1)$  (independent from  $\varepsilon$ ) such that for any  $\kappa \leq \kappa_0$ , if  $n \geq n_{2,\varepsilon}^{FISTA}$  then the iterate  $x_n$  is a  $\varepsilon$ -solution.

**Proof of Theorems 3 and 4** The proof of Theorems 3 and 4 is based on new non-asymptotic bounds for FISTA provided by Theorems 6 and 7 and of the form:

$$F(x_n) - F^* \leq \left( \frac{C_\gamma(\alpha)}{n\sqrt{\kappa}} \right)^{\frac{2\alpha\gamma}{\gamma+2}} M_0$$

where  $M_0 = F(x_0) - F^*$ , which means that an  $\varepsilon$ -solution can be obtained at most in

$$\frac{1}{\sqrt{\kappa}} C_\gamma(\alpha) \left( \frac{2LM_0}{\varepsilon^2} \right)^{\frac{2+\gamma}{2\alpha\gamma}}$$

iterations where:

$$C_1(\alpha) = \frac{8e}{3} \alpha \left( \frac{3}{2e} \right)^{\frac{3}{\alpha}}, \quad C_2(\alpha) = \frac{5e}{2} \alpha \left( \frac{4}{3e} \right)^{\frac{2}{\alpha}}.$$

Thus with these bounds, the optimized numbers  $n_{\gamma,\varepsilon}$ ,  $\gamma \in \{1, 2\}$  of iterations to reach an  $\varepsilon$ -solution with FISTA, are respectively obtained for:

$$\alpha_{1,\varepsilon} = 3 \log \left( \frac{3\sqrt{LM_0}}{e\sqrt{2}\varepsilon} \right), \quad \alpha_{2,\varepsilon} = 2 \log \left( \frac{4\sqrt{2LM_0}}{3e\varepsilon} \right). \quad (33)$$

For these choices of  $\alpha$ , the number of iterations to reach an  $\varepsilon$ -solution is respectively given by:

$$n_{1,\varepsilon}^{FISTA} = \frac{8e^2}{3\sqrt{\kappa}}\alpha_{1,\varepsilon} = \frac{8e^2}{\sqrt{\kappa}} \log\left(\frac{3\sqrt{LM_0}}{e\sqrt{2\varepsilon}}\right) \quad (34)$$

$$n_{2,\varepsilon}^{FISTA} = \frac{5e^2}{2\sqrt{\kappa}}\alpha_{2,\varepsilon} = \frac{5e^2}{\sqrt{\kappa}} \log\left(\frac{4\sqrt{2LM_0}}{3e\varepsilon}\right). \quad (35)$$

■

## 2.4 Comparisons

### 2.4.1 Comparison with FB

Let us now compare the Forward-Backward algorithm to the FISTA scheme for a given accuracy  $\varepsilon > 0$ . For a condition number  $\kappa$  small enough and choosing:

$$\alpha = \alpha_{1,\varepsilon} = 3 \log\left(\frac{3\sqrt{LM_0}}{e\sqrt{2\varepsilon}}\right),$$

we easily check that FISTA requires fewer iterations than the FB algorithm to reach an  $\varepsilon$ -solution of  $\min_{x \in \mathbb{R}^N} F(x)$ . Indeed,

$$n_{1,\varepsilon}^{FISTA} = \frac{4e^2}{\sqrt{\kappa}} \log\left(\frac{9LM_0}{2e^2\varepsilon^2}\right) \leq n_{\varepsilon}^{FB} = \frac{1}{|\log(1-\kappa)|} \log\left(\frac{2LM_0}{\varepsilon^2}\right). \quad (36)$$

Observe that the smaller  $\kappa$  is, the better FISTA is compared to FB. Note also that since  $n_{2,\varepsilon}^{FISTA} \leq n_{1,\varepsilon}^{FISTA}$ , the comparison between FISTA and FB remains in favor of FISTA for nearly quadratic functions (namely for differentiable functions satisfying both a quadratic growth condition and some flatness assumption).

Even though FB is known to have an exponential convergence [17], and FISTA a polynomial one [6], choosing the  $\alpha$  parameter for FISTA enables to have a much better convergence rate in term of  $\varepsilon$ -solution than FB.

### 2.4.2 Comparison with NSC

On this subsection only the comparison between FISTA and the inertial algorithm dedicated to strongly convex functions proposed by Nesterov [23, Algorithm 2.2.11] is detailed but comparisons with any other algorithms built for strongly convex functions described in [23, 27, 9, 8, 29] will lead to the same conclusions.

The Nesterov algorithm for strongly convex functions (Algorithm 3) necessitates an estimation  $\tilde{\mu}$  of the strong convexity parameter  $\mu$  that is usually not known. If  $\mu$  is underestimated (i.e.  $\tilde{\mu} \leq \mu$ ), the NSC algorithm will be miscalibrated and therefore slowed down. One of the strengths of FISTA, revealed by our non-asymptotic analysis, is that FISTA is self-adaptative and will have better performances than its NSC variant for strongly convex functions when the strong convexity parameter is not known. Indeed, by choosing the  $\alpha$  friction parameter only according to the desired accuracy  $\varepsilon$ , FISTA will generate a sequence of iterates until reaching an  $\varepsilon$ -solution without the need of any estimation of the strong convexity parameter  $\mu$ . Note also that the number of iterations to reach an  $\varepsilon$ -solution for FISTA actually depends on the true value of  $\mu$  whether it is known or not.

If the strong convexity parameter  $\mu$  of  $F$  is known, it is clear that for small enough  $\kappa$ , Algorithm 3 is faster than FISTA i.e.

$$n_{\varepsilon}^{NSC} = \frac{1}{|\log(1-\sqrt{\kappa})|} \log\left(\frac{4LM_0}{\varepsilon^2}\right) \leq n_{1,\varepsilon}^{FISTA} = \frac{4e^2}{\sqrt{\kappa}} \log\left(\frac{9LM_0}{2e^2\varepsilon^2}\right). \quad (37)$$

Hence the number of iterations needed to reach an  $\varepsilon$ -solution is smaller for the scheme built for strongly convex functions and using explicitly this parameter  $\mu$  at each step (19) (since  $\kappa = \frac{\mu}{L}$ ) of

the algorithm than the one necessary for FISTA to get the same accuracy. This better behavior was indeed expected.

On the other hand if  $\mu$  is not perfectly known, which is often the case in large dimensional problems,  $\mu$  should be estimated by  $\tilde{\mu}$  and to ensure that the exponential decay of these inertial algorithms are in force,  $\tilde{\mu}$  must be chosen such that  $\tilde{\mu} \leq \mu$  and

$$n_\varepsilon^{NSC} = \frac{1}{\left| \log(1 - \sqrt{\frac{\tilde{\mu}}{L}}) \right|} \log\left(\frac{4LM_0}{\varepsilon^2}\right) \geq \frac{1}{\left| \log(1 - \sqrt{\kappa}) \right|} \log\left(\frac{4LM_0}{\varepsilon^2}\right) \quad (38)$$

If  $\tilde{\mu}$  is close to  $\mu$ , one can expect that  $n_\varepsilon^{NSC} \leq n_{1,\varepsilon}^{FISTA}$ . But if  $\tilde{\mu} \leq \frac{1}{16\varepsilon^4}\mu$  then

$$n_\varepsilon^{NSC} \geq n_{1,\varepsilon}^{FISTA} \quad (39)$$

and thus if  $\mu$  is not known with a good accuracy, it may be better to use FISTA.

In practice, FISTA may outperform Algorithm 3 even for much smaller underestimation of  $\mu$ . The bound given in Theorem 3 may be pessimistic. We illustrate this lower performance of this inertial algorithm with fixed friction term in the numerical experiments comparing to FISTA in the subsection dedicated to numerical experiments. These experiments confirm that even for  $\tilde{\mu} = \frac{\mu}{10}$  the loss may be huge and FISTA is actually better for a large set of accuracies  $\varepsilon$ .

Indeed iterations of FISTA use the value of a parameter  $\alpha$  which is defined from  $\varepsilon$  and not from an estimation of  $\mu$  which implies that  $n_{1,\varepsilon}^{FISTA}$  does not depend on any estimation of  $\mu$ . The fact that  $n_\varepsilon^{NSC}$  may be larger than  $n_{1,\varepsilon}^{FISTA}$  when  $\mu$  is not well estimated is a small surprise since FISTA was not built for strongly convex functions. Moreover bounds on  $n_{1,\varepsilon}^{FISTA}$  apply under a quadratic growth property and then extend the potential application of this result to a larger set of functions  $F$  such as the LASSO :

$$F(x) = \frac{1}{2} \|Ax - y\|^2 + \lambda \|x\|_1 \quad (40)$$

widely used in Statistics and Image and Signal processing. This function  $F$  is not strongly convex when  $\lambda > 0$  but it satisfies some growth properties, see for example [13].

To be fair in the comparison between both algorithm, we must emphasize that for FISTA, the parameter  $\alpha$  depends on the targeted accuracy  $\varepsilon$ , and FISTA may be better than the other inertial algorithms for this specific accuracy. If  $n$  goes to infinity, the algorithms built with a fixed inertia (here  $\frac{1-\sqrt{\kappa}}{1+\sqrt{\kappa}}$ ) may have a better behavior than FISTA, since for any fixed  $\alpha$ , FISTA as a polynomial decay rate, see for example [3, 4].

Let us finish this section by comparing with the scheme introduced by Nesterov in [22] for strongly convex functions and that automatically estimates the strong convexity parameter. Referring to (22) we observe that:

$$n_\varepsilon^{N\mu} = \mathcal{O}\left(\frac{|\log(\kappa)|}{\sqrt{\kappa}} \log\left(\frac{1}{\varepsilon}\right)\right). \quad (41)$$

However, we have for FISTA:

$$n_{1,\varepsilon}^{FISTA} = \frac{4e^2}{\sqrt{\kappa}} \log\left(\frac{9LM_0}{2e^2\varepsilon^2}\right) = \mathcal{O}\left(\frac{1}{\sqrt{\kappa}} \log\left(\frac{1}{\varepsilon}\right)\right). \quad (42)$$

Hence for small  $\kappa$ , we lose a factor  $|\log(\kappa)|$  which makes the scheme of [22] less efficient.

### 2.4.3 Comparison with FISTA restart

Under the assumption of quadratic growth it appears that the FISTA restart scheme has a convergence rate similar to FISTA with optimal parameter. In both cases the bounds on the number of iterations needed to get an  $\varepsilon$ -solution is proportional to  $\frac{1}{\sqrt{\kappa}} \log\left(\frac{\sqrt{ML}}{\varepsilon}\right)$ . Preliminary numerical

experiments seem to indicate that the best solution may depend on the function  $F$  to minimize and both approaches deserve to be tested. In Theorem 1,  $F$  is supposed to have a unique minimizer, while one can notice that such an assumption is not needed for FISTA restart, even if in practice, it does not seem to have any impact on the convergence rate. One can also observe that the parameter of FISTA must be chosen according to the wished accuracy  $\varepsilon$ ; for this accuracy FISTA will be efficient and almost optimal all along the trajectory. The main inconvenient of FISTA restart scheme is its relative complexity where FISTA is really simple to implement.

When  $\varepsilon$  tends to 0, according to (23):

$$n_{\varepsilon}^{Restart} \sim \frac{14.4}{\sqrt{\kappa}} \log \left( \frac{\sqrt{LM_0}}{\varepsilon} \right) \quad (43)$$

while

$$n_{1,\varepsilon}^{FISTA} \sim \frac{8e^2}{\sqrt{\kappa}} \log \left( \frac{\sqrt{LM_0}}{\varepsilon} \right), \quad n_{2,\varepsilon}^{FISTA} \sim \frac{19e^2}{4\sqrt{\kappa}} \log \left( \frac{\sqrt{LM_0}}{\varepsilon} \right). \quad (44)$$

Hence

$$n_{1,\varepsilon}^{FISTA} \approx \frac{59.2}{\sqrt{\kappa}} \log \left( \frac{\sqrt{LM_0}}{\varepsilon} \right), \quad n_{2,\varepsilon}^{FISTA} \approx \frac{35.1}{\sqrt{\kappa}} \log \left( \frac{\sqrt{LM_0}}{\varepsilon} \right), \quad (45)$$

i.e. the bounds for FISTA restart [10] are always slightly better than the one proposed here in the paper. However, we will see in the next subsection that this slight theoretical edge for the worst case analysis may not always prevail in practice. Moreover, it can be argued that FISTA restart needs in general more calls to the function to minimize than classical FISTA (which has a negative impact on its speed), and that it is slightly more complicated to implement.

#### 2.4.4 Numerical comparisons

Let us first consider a function  $F(x) = \|Ax - b\|^2$  where  $A$  is  $100 \times 100$  gaussian matrix with independent and identically distributed components. On that example  $\kappa \approx 4.7 \times 10^{-7}$ . On Figure 1, are given the values of  $\log(\|\nabla F(x_n)\|)$  along the trajectory for various algorithms.

The blue curve corresponds to the Gradient descend, the red curve to FISTA with  $\alpha = 8$ , the yellow curve to FISTA with  $\alpha = 30$ , the green one to the Algorithm 3 of Nesterov for strongly convex functions using the exact value of  $\mu$  computed from the realisation of the matrix  $A$ . The grey curve corresponds to Algorithm 3 with  $\tilde{\mu} = \frac{1}{10}\mu$ . The graph is displayed in Figure 1. From this graph, several comments can be done :

- For small  $\varepsilon$ , Algorithm 3 with the precise value of  $\mu$  (the green curve), seems to be the best one, which was expected.
- If  $\mu$  is underestimated (grey curve), even of a 10 factor, the decay may be much smaller with the FISTA schemes.
- For  $\varepsilon \leq e^{-15}$ , the parameter  $\alpha = 8$  (red curve) seems to be better than  $\alpha = 30$  for FISTA.
- For  $\varepsilon \geq e^{-15}$ , the parameter  $\alpha = 30$  (yellow curve) seems to be better than  $\alpha = 8$  for FISTA.

More complete numerical experiments show in practice that the optimal value of  $\alpha$  for a given accuracy, that is the value of  $\alpha$  ensuring the minimum number iteration to reach an  $\varepsilon$ -solution, is actually a non increasing function of the accuracy  $\varepsilon$ , which illustrates Theorem 3.

We then illustrate the paper with a second example, where an inpainting problem is solved using a LASSO formulation.

$$F(x) = \frac{1}{2} \|Mx - Mx^o\|^2 + \lambda \|Tx\|_1 \quad (46)$$

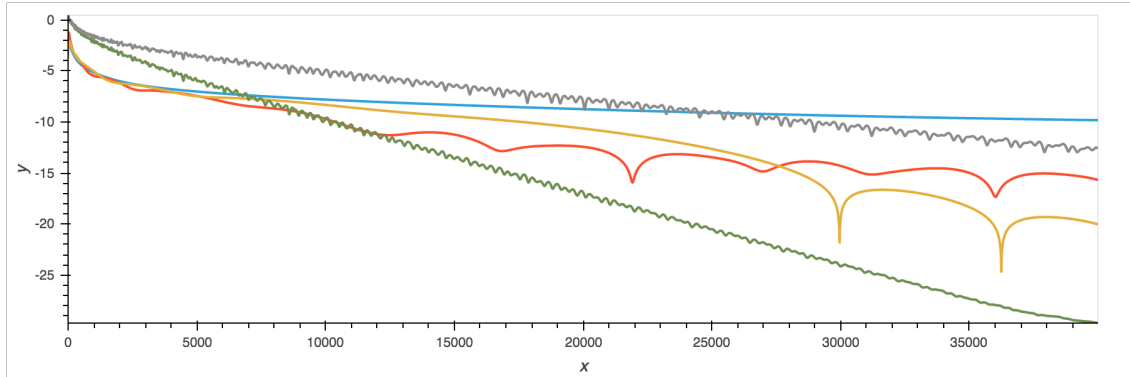


Figure 1: Example on the least square problem. Blue is the gradient descent (FB), grey is NSC with a 10 factor underestimate of  $\mu$ , red is FISTA with  $\alpha = 8$ , yellow is FISTA with  $\alpha = 30$ , and green is NSC with the exact value of  $\mu$ .

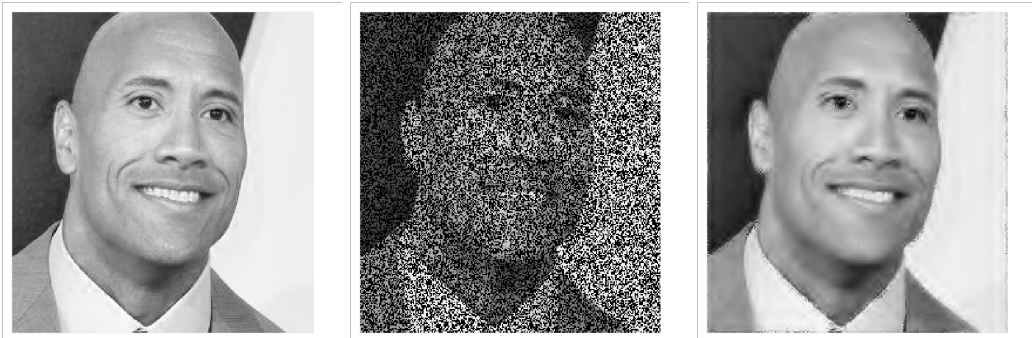


Figure 2: An inpainting example (left: original image, middle: degraded image, right: inpainted result)

where  $x^o$  is a target image,  $M$  a random masking operator and  $T$  an orthogonal wavelet transform. Figure 2 shows an example, with from the left to right, the original image  $x^o$ , the masked image  $Mx^o$  and the solution  $x^*$  minimizer of  $F$ .

Figure 3 displays the curves of the trajectory of  $\log(L\|x_n - x_n^+\|)$  for several algorithms. The blue curve corresponds to the Forward Backward algorithm. The red curve corresponds to FISTA-restart scheme described in [10], the yellow curve to FISTA with  $\alpha = 3$  (which is the classical FISTA algorithm), the green curve to FISTA with  $\alpha = 12$  and the grey curve to FISTA with  $\alpha = 30$ .

Several comments can be done: first, all inertial schemes are better than FB (blue curve) for any precision  $\varepsilon$ . If we compare the three FISTA algorithms we observe that for large  $\varepsilon$ ,  $\alpha = 3$  seems to be better. For small  $\varepsilon$ ,  $\alpha = 30$  seems to be the better choice and in between  $\alpha = 12$  is better than the two others. That is what was expected from Theorem 3: the optimal value of  $\alpha$  is a non increasing function of  $\varepsilon$ . Note that the restart FISTA [10] behaves quite well for any precision and its efficiency seems to be close to FISTA with the best parameter for all accuracy. The optimal value of  $\alpha$  given in Theorem 3, depends on the accuracy  $\varepsilon$ , and on  $M_0$  which depends on the distance of the initial point  $x^0$  to the minimizer. In practical situations only an upper bound on  $M_0$  can be known. The numerical experiments illustrate the fact that the optimal  $\alpha$  is a non increasing function of the desired accuracy. However, there is no explicit formula for  $\alpha$  that does not depend on the distance of  $x_0$  to the minimizer. Theorem 3 explains the behaviour of  $\alpha$  with respect to the accuracy  $\varepsilon$ , even though it does not provide an explicit formula.

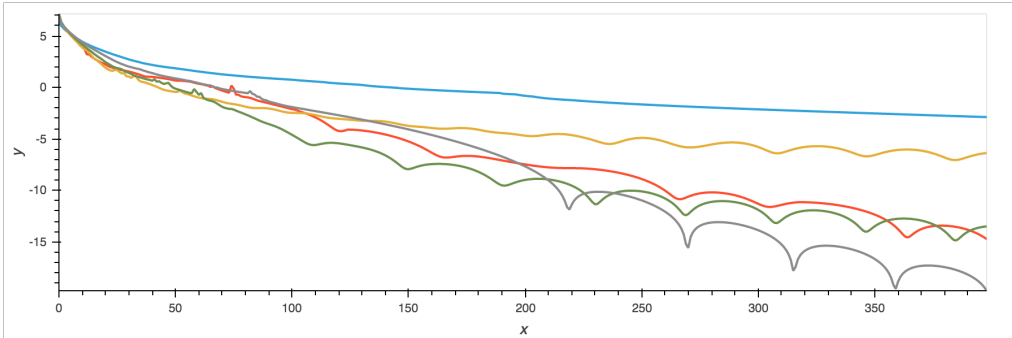


Figure 3: Example on a LASSO problem: FB is in blue, FISTA-restart in red, FISTA with  $\alpha = 3$  in yellow, FISTA with  $\alpha = 12$  in green, FISTA with  $\alpha = 30$  in grey.

### 3 Discrete non asymptotic analysis of FISTA

Let  $F = f + h$  be a convex composite function where  $f$  is a convex, differentiable function having a  $L$ -Lipschitz gradient and  $h$  is a l.s.c. convex function whose proximal operator is known. Furthermore, we assume that  $F$  has a unique minimizer  $x^*$ . Let  $F^* = F(x^*) = \inf F$ .

In this section we provide a complete non-asymptotic analysis of FISTA [14]:

$$y_n = x_n + \frac{n}{n+\alpha}(x_n - x_{n-1}), \quad x_{n+1} = \text{prox}_{sh}(y_n - s\nabla f(y_n)) \quad (47)$$

for this class of convex composite functions  $F$  satisfying additionally some global quadratic growth property  $\mathcal{G}_\mu^2$  (that becomes a strong minimizer property since we assume the uniqueness of the minimizer)

$$\exists \mu > 0, \forall x \in \mathbb{R}^N, F(x) - F^* \geq \frac{\mu}{2} \|x - x^*\|^2. \quad (48)$$

Our main contribution in this section is to provide non-asymptotic bounds on the values  $F(x_n) - F^*$  along the iterates generated by FISTA, see Theorem 6. We also prove that this bound could be slightly improved when  $F$  is nearly quadratic (i.e. satisfying both a quadratic growth condition and a flatness condition  $\mathcal{H}_2$ ). Our analysis is based on Lyapunov energies:

$$E_n = 2sn^2(F(x_n) - F^*) + \|\lambda(x_{n-1} - x^*) + n(x_n - x_{n-1})\|^2$$

that can be seen as discretizations of the Lyapunov energy introduced in the continuous setting:

$$\mathcal{E}(t) = t^2(F(x(t)) - F^*) + \frac{1}{2} \|\lambda(x(t) - x^*) + t\dot{x}(t)\|^2,$$

where  $x^*$  denotes some minimizer the objective function  $F$ . In Section 3.1, assuming that  $F$  is a convex differentiable function, we provide the convergence rate analysis of the trajectories of the ODE:

$$\ddot{x}(t) + \frac{\alpha}{t}\dot{x}(t) + \nabla F(x(t)) = 0 \quad (49)$$

associated to the Nesterov scheme. For the class of strongly convex functions, Su Boyd and Candès [28, Theorem 8] prove that the ODE (49) achieves a convergence rate in  $\mathcal{O}\left(t^{-\frac{2\alpha}{3}}\right)$ . We prove that this result can be generalized to the class of convex functions satisfying some quadratic growth condition. More importantly, the analysis of the continuous case will give us a guideline for the convergence rate analysis in the discrete setting: the main idea is to identify the key inequalities whose discrete equivalents will be the milestones for our analysis of FISTA in Section 3.2.

### 3.1 The continuous case, a guideline to analyse the discrete

Let  $F : \mathbb{R}^N \rightarrow \mathbb{R}$  be a convex differentiable function admitting a unique minimizer  $x^*$ . In this section, we study the convergence rates in finite time for the values  $F(x(t)) - F^*$  along the trajectories of the well-known ordinary differential equation:

$$\ddot{x}(t) + \frac{\alpha}{t}\dot{x}(t) + \nabla F(x(t)) = 0 \quad (50)$$

for any  $t \geq t_0$  with  $t_0 > 0$ , associated to the Nesterov scheme. We assume that, for any initial conditions  $(x_0, v_0) \in \mathbb{R}^N \times \mathbb{R}^N$ , the Cauchy problem associated with the ODE (49) admits a unique global solution satisfying  $(x(t_0), \dot{x}(t_0)) = (x_0, v_0)$ .

In this section  $F$  is assumed to satisfy some general flatness assumption  $\mathcal{H}_\gamma$  ensuring that  $F$  is not too sharp: let  $\gamma \geq 1$ . For any minimizer  $x^*$ , we assume that:

$$\forall x \in \mathbb{R}^N, F(x) - F^* \leq \frac{1}{\gamma} \langle \nabla F(x), x - x^* \rangle. \quad (51)$$

Note that any convex differentiable function automatically satisfies  $\mathcal{H}_1$ . To have a better intuition of the geometry of functions satisfying  $\mathcal{H}_\gamma$  for some  $\gamma > 0$ , observe that the flatness property (51) implies that for any minimizer  $x^* \in X^*$ , there exists a real constant  $M > 0$  such that:

$$\forall x \in \mathbb{R}^N, F(x) - F^* \leq M \|x - x^*\|^\gamma, \quad (52)$$

see [6, Lemma 2.2]. For this class of functions, we have:

**Theorem 5** *Let  $F : \mathbb{R}^N \rightarrow \mathbb{R}$  be a convex differentiable function admitting a unique minimizer  $x^*$ . Assume that  $F$  satisfies both a quadratic growth condition  $\mathcal{G}_\mu^2$  and a flatness condition  $\mathcal{H}_\gamma$  for some  $\mu > 0$  and  $\gamma \geq 1$ . If  $\alpha > 1 + \frac{2}{\gamma}$  and for  $\mu$  small enough, we have:*

$$\forall t \geq \frac{\alpha r^*}{(\gamma + 2)\sqrt{\mu}} \geq t_0, F(x(t)) - F^* \leq C_1 e^{\frac{2\gamma}{\gamma+2} C_2 (\alpha - 1 - \frac{2}{\gamma})} E_m(t_0) \left( \frac{\alpha r^*}{t(\gamma + 2)\sqrt{\mu}} \right)^{\frac{2\alpha\gamma}{\gamma+2}}$$

where  $E_m(t) = F(x(t)) - F^* + \frac{1}{2} \|\dot{x}(t)\|^2$  denotes the mechanical energy of the system,  $r^* \simeq 3$  is the unique positive real root of the polynomial:  $r \mapsto r^3 - r^2 - 2(1 + \sqrt{2})r - 4$  and

$$C_1 = 1 + \frac{2}{r^*} + \frac{4}{r^{*2}}, \quad C_2 = \frac{1}{r^*} + \frac{1 + \sqrt{2}}{r^{*2}} + \frac{4}{3r^{*3}}.$$

Theorem 5 is a generalization of [28, Theorem 8] to the class of convex functions satisfying some quadratic growth condition, and provides in particular an explicit bound on  $F(x(t)) - F^*$  decaying like  $t^{-\frac{2\alpha}{3}}$  when  $\gamma = 1$  (respectively like  $\frac{1}{t^\alpha}$  when  $\gamma = 2$ ). This bound depends on the growth parameter  $\mu$  and the friction coefficient  $\alpha$  and is valid for sufficiently large enough  $t$  namely for

$$t \geq t_{\gamma, \alpha, \mu} := \frac{\alpha r^*}{(\gamma + 2)\sqrt{\mu}}. \quad (53)$$

Actually this restriction is not really a problem. First because it is possible to reduce  $t_{\gamma, \alpha, \mu}$  in the proof, which would lead to not as good asymptotic bounds and secondly because inequality such that  $F(x(t)) - F^* \leq Kt^{-\frac{2\alpha}{3}}$  may not provide interesting bounds for small  $t$ . It turns out that for most inertial algorithm, the decay of the  $F(x(t)) - F^*$  is not significant if  $t \ll \frac{1}{\sqrt{\mu}}$ . Indeed, even if finite time bounds are valid from  $t = t_0 = 0$ , they only provide accurate bounds for  $t \geq \frac{1}{\sqrt{\mu}}$ .

We can also observe that the bound given by Theorem 5 for a given  $t$ , is not a decaying function of  $\alpha$  which explains why it is not relevant to choose  $\alpha$  as large as possible if we consider the ODE on an interval  $[t_0, T]$ .

The proof of Theorem 5 is detailed in Appendix A. We only detail hereafter the sketch of the proof of Theorem 5 in the convex case (i.e. for  $\gamma = 1$ ) which will give us the main steps to follow



for the study of the discrete case: considering the following Lyapunov energy which is a variation of the one used in [28]:

$$\mathcal{E}(t) = t^2(F(x(t)) - F^*) + \frac{1}{2} \|\lambda(x(t) - x^*) + t\dot{x}(t)\|^2, \quad \lambda = \frac{2\alpha}{3}, \quad (54)$$

we first establish some differential inequality of the form:

$$\mathcal{E}'(t) + \frac{\lambda - 2}{2} \mathcal{E}(t) \leq \varphi(t) \mathcal{E}(t), \quad (55)$$

where the function  $\varphi$  is polynomial of order greater than 2 in  $\frac{1}{t}$ . Let  $\Phi(t) = \int_t^{+\infty} \varphi(s) ds$ . The function  $t \mapsto \mathcal{E}(t)t^{\lambda-2}e^{\Phi}(t)$  is so non-increasing. Integrating (55) between  $t_1$  and  $t$  for any  $t_1 \geq t_0$ , we get:

$$\forall t \geq t_0, \quad \mathcal{E}(t) \leq \mathcal{E}(t_1) \left(\frac{t_1}{t}\right)^{\lambda-2} e^{\Phi(t_1) - \Phi(t)}. \quad (56)$$

Choosing  $t_1$  ensuring a control as tight as possible on the energy  $\mathcal{E}$  and noticing that  $\mathcal{E}(t) \geq t^2(F(x(t)) - F^*)$ , we eventually obtain the explicit control given in Theorem 5 on the values of  $F$  along the trajectories of the ODE (49).

Finally observe that the choice of the parameter  $\alpha$  can be optimized for a given  $t$  by choosing:

$$\alpha_{opt} = t \frac{(\gamma + 2)\sqrt{\mu}}{r^*} e^{-1-C_2} \quad (57)$$

which implies a fast exponential decay rate on the values:

$$\begin{aligned} F(x(t)) - F^* &\leq C_1 E_m(t_0) \left( e^{-1-C_2} \right)^{\frac{2\gamma\alpha_{opt}t}{\gamma+2}} e^{\frac{2\gamma}{\gamma+2} C_2 (\alpha_{opt}t - 1 - \frac{2}{\gamma})} \\ &\leq C_1 E_m(t_0) e^{-(1+\frac{2}{\gamma})C_2} e^{-\frac{2\gamma}{r^*} e^{-1-C_2} \sqrt{\mu}t}. \end{aligned}$$

This optimal choice depends also on  $\mu$  that can be unknown in practice but we can remark that if  $\alpha \approx t\sqrt{\mu}$ , FISTA ensures a fast exponential decay which ensures the best possible decay for the bound given in Theorem 5.

### 3.2 Non-asymptotic bounds for convex composite functions

Let

$$E_n = 2sn^2(F(x_n) - F^*) + \|\lambda(x_{n-1} - x^*) + n(x_n - x_{n-1})\|^2 \quad (58)$$

be a well-chosen discretization of the Lyapunov energy (78) introduced in the continuous setting. Our main result provides non-asymptotic bounds on the value  $F(x_n) - F^*$  along the iterates generated by FISTA:

**Theorem 6** *Let  $F = f + g$  where  $f$  is a convex differentiable function having a  $L$ -Lipschitz gradient for some  $L > 0$ , and  $g$  a proper convex l.s.c. function. Assume that  $F$  admits a unique minimizer  $x^*$  and satisfies a quadratic growth condition  $\mathcal{G}_\mu^2$  for some real parameter  $\mu > 0$ .*

*Let  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$  and  $\kappa = \frac{\mu}{L}$ . Then there exist  $\kappa_0 > 0$  such that for any  $0 < \kappa \leq \kappa_0$ , the sequence  $(x_n)_{n \in \mathbb{N}}$  generated by FISTA satisfies:*

$$\forall n \geq \frac{3\alpha}{\sqrt{\kappa}}, \quad F(x_n) - F^* \leq \frac{9}{4} e^{-2} M_0 \left( \frac{8e}{3\sqrt{\kappa}} \alpha \right)^{\frac{2\alpha}{3}} n^{-\frac{2\alpha}{3}}. \quad (59)$$

where  $M_0 = F(x_0) - F^*$  denotes the potential energy of the system at initial time.

Observe that in finite time (i.e. for a given number of iterations  $n$ ), a fast exponential decay may be obtained from Theorem 6 by choosing  $\alpha$  that minimizes the function:

$$\alpha \mapsto \frac{2\alpha}{3} \log \left( e \frac{8\alpha}{3n\sqrt{\kappa}} \right).$$

A straightforward computation shows that the minimum value is reached for:

$$\alpha^* := \frac{3n\sqrt{\kappa}}{8e^2} \quad (60)$$

and with this choice of parameters we finally deduce:

$$F(x_n) - F^* \leq \frac{9}{4}e^{-2} \exp\left(-\frac{n\sqrt{\kappa}}{4e^2}\right) M_0. \quad (61)$$

The sketch of the proof of Theorem 6 is given in Subsection 3.3, while the proofs of the technical lemmas are detailed in Appendix B.

Assume now that  $F$  is a nearly quadratic convex differentiable function ( $h = 0$ ) i.e. satisfying both a quadratic growth condition and a flatness condition  $\mathcal{H}_2$  i.e. for any minimizer  $x^*$ :

$$\forall x \in \mathbb{R}^N, F(x) - F^* \leq \frac{1}{2} \langle \nabla F(x), x - x^* \rangle. \quad (62)$$

Attouch et al. [5] and Su Boyd and Candes [28] proved that for strongly convex functions  $F(x_n) - F(x^*) = \mathcal{O}(n^{-2\alpha/3})$ . Under the flatness assumption (62), Theorem 7 improves this decay to  $\mathcal{O}(n^{-\alpha})$  for nearly quadratic convex and only differentiable functions:

**Theorem 7** *Let  $F : \mathbb{R}^N \rightarrow \mathbb{R}$  be a convex differentiable function having a  $L$ -Lipschitz gradient and a unique minimizer  $x^*$ . Assume additionally that  $F$  satisfies a quadratic growth condition  $\mathcal{G}_\mu^2$  for some real parameter  $\mu > 0$  and  $L > 0$  and a flatness condition  $\mathcal{H}_2$ .*

*Let  $\alpha \geq 4 + 2\sqrt{2}$ . Then there exist  $\kappa_0 > 0$  and a real constant  $C_3 > 0$  such that for any  $0 < \kappa \leq \kappa_0$ , the sequence  $(x_n)_{n \in \mathbb{N}}$  generated by FISTA satisfies*

$$\forall n \geq \frac{5\alpha}{2\sqrt{\kappa}}, F(x_n) - F^* \leq \frac{16}{9} \left( e \frac{5\alpha}{2n\sqrt{\kappa}} \right)^\alpha e^{-2} M_0 \quad (63)$$

where  $M_0 = F(x_0) - F^*$  denotes the potential energy of the system at initial time.

The proof of Theorem 7 is based on a variant of the Lyapunov energy (58):

$$E_n = 2s n^2 (F(x_n) - F^*) + \left\| \frac{\alpha}{2} (x_{n-1} - x^*) + \left(n - \frac{\alpha}{4}\right) (x_n - x_{n-1}) \right\|^2 \quad (64)$$

and follows the exact same steps than those of the proof of Theorem 6. The sketch of the proof of Theorem 7 is given in Appendix C, while the proofs of the associated technical lemmas are omitted here, but can be found in [7].

### 3.3 Sketch of the proof of Theorem 6

The proof of Theorem 6 is based on the Lyapunov energy (58) which can be seen as a discretization of the Lyapunov energy (54) introduced in the continuous setting. Stating:

$$w_n = 2s(F(x_n) - F^*), \quad h_n = \|x_n - x^*\|^2, \quad \delta_n = \|x_n - x_{n-1}\|^2, \quad \alpha_n = \frac{n}{n + \alpha}, \quad \lambda = \frac{2\alpha}{3}, \quad (65)$$

the energy  $E_n$  can be rewritten as:

$$E_n = n^2 w_n + (\lambda^2 - \lambda n) h_{n-1} + (n^2 - \lambda n) \delta_n + \lambda n h_n \quad (66)$$

As in the continuous setting, the first step of the proof consists in establishing some discrete version of the differential inequality (55):

**Lemma 1** Let  $\kappa = \frac{\mu}{L}$ . There exists  $\kappa_0 > 0$  such that for any  $\kappa \leq \kappa_0$  and for any  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$ , there exists some real constants  $\tilde{c}_1$  and  $\tilde{c}_2$  such that:

$$\forall n \geq \frac{4\alpha}{3\sqrt{\kappa}}, E_{n+1} - \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) E_n \leq C_1(\alpha, \kappa) \frac{E_n}{n^2} + C_2(\alpha, \kappa) \frac{E_{n+1}}{(n+1)^2} \quad (67)$$

with:

$$C_1(\alpha, \kappa) = \frac{5}{4} \sqrt{\frac{2}{\kappa}} \left[ \frac{2}{9} (\alpha - 3) \left( \frac{8}{5} \alpha - 3 \right) - 1 \right] (1 + \sqrt{\kappa})^2 (1 + \tilde{c}_1 \sqrt{\kappa} + \tilde{c}_2 \kappa) \quad (68)$$

$$C_2(\alpha, \kappa) = \frac{5}{4} \sqrt{\frac{2}{\kappa}} \left( \frac{2\alpha}{3} - 1 \right) (1 + \sqrt{\kappa})^2 (1 + 2\sqrt{2\kappa}). \quad (69)$$

The proof of Lemma 1 is detailed in Appendix B.1. The next step consists in integrating the inequality (67):

**Lemma 2** Let  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$  and  $n_0 \geq \frac{4\alpha}{3\sqrt{\kappa}}$ . If the energy  $E_n$  satisfies (67) then we have:

$$\forall n \geq n_0, E_n \leq E_{n_0} \left( \frac{n}{n_0} \right)^{-\left(\frac{2\alpha}{3} - 2\right)} e^{\Phi(n_0)} \quad (70)$$

with

$$\Phi(n_0) = \frac{5}{6n_0} \sqrt{\frac{2}{\kappa}} (\alpha - 3) \left( \frac{16}{15} \alpha - 1 \right) (1 + C_3 \kappa^{1/4}). \quad (71)$$

The proof of Lemma 2 is detailed in Appendix B.2. Remembering that  $F(x_n) - F^* \leq \frac{1}{2sn^2} E_n$  for any  $n$ , we thus have:

$$\forall n \geq \frac{4\alpha}{3\sqrt{\kappa}}, F(x_n) - F^* \leq \frac{E_{n_0}}{2s} \left( n_0^{\frac{2\alpha}{3} - 2} e^{\Phi(n_0)} \right) n^{-\frac{2\alpha}{3}}. \quad (72)$$

A good choice of  $n_0$  is one ensuring a control as tight as possible on the values  $F(x_n) - F^*$ . For that,  $n_0$  is chosen such that it minimizes the function  $f : x \mapsto x^{\frac{2\alpha}{3} - 2} e^{\Phi(x)}$ . A straightforward computation gives:

$$n_0 = \frac{5}{4} \sqrt{\frac{2}{\kappa}} \left( \frac{16}{15} \alpha - 1 \right) (1 + C_3 \kappa^{1/4}). \quad (73)$$

Note that the optimized value of  $n_0$  satisfies:  $n_0 \geq \frac{4\alpha}{3\sqrt{\kappa}}$  for any  $\kappa > 0$ . Observe that  $f(n_0) = (e n_0)^{\frac{2\alpha}{3} - 2}$  and that for any  $\kappa$  small enough:

$$\frac{5}{4} \sqrt{\frac{2}{\kappa}} \left( \frac{16}{15} \alpha - 1 \right) (1 + C_3 \kappa^{1/4}) \leq \frac{3\alpha}{\sqrt{\kappa}}. \quad (74)$$

We deduce:

$$\forall n \geq \frac{3\alpha}{\sqrt{\kappa}}, F(x_n) - F^* \leq \frac{E_{n_0}}{2s} n^{-\frac{2\alpha}{3}} (e n_0)^{\left(\frac{2\alpha}{3} - 2\right)} = \frac{E_{n_0}}{2sn_0^2} e^{-2} \left( \frac{e n_0}{n} \right)^{\frac{2\alpha}{3}}. \quad (75)$$

Observe now that the energy of the system is controlled by the mechanical energy of the system at initial time  $t_0$ :

**Lemma 3** Let  $M_n$  the mechanical energy:  $M_n = F(x_n) - F^* + \frac{1}{2s} \|x_n - x_{n-1}\|^2$ . Then we have

$$\frac{E_n}{2sn^2} \leq \left( 1 + \frac{4\alpha^2}{9\kappa n^2} + \frac{4\alpha}{3\sqrt{\kappa}n} \right) M_n = \left( 1 + \frac{2\alpha}{3\sqrt{\kappa}n} \right)^2 M_n \quad (76)$$

Applying Lemma 3 whose proof is detailed in Appendix 3, to uniformly bound the energy  $E_{n_0}$  and the fact that  $n_0 \geq \frac{4\alpha}{3\sqrt{\kappa}}$ , we have:

$$\frac{E_{n_0}}{2sn_0^2} \leq \left(1 + \frac{2\alpha}{3\sqrt{\kappa}n_0}\right)^2 M_{n_0} \leq \frac{9}{4}M_{n_0}. \quad (77)$$

Since the mechanical energy associated to the Nesterov scheme is non-increasing (see [14, Corollary 2]) and  $x_{-1} = x_0$ , we then get:

$$\begin{aligned} \forall n \geq \frac{3\alpha}{\sqrt{\kappa}}, \quad F(x_n) - F^* &\leq \frac{9}{4}e^{-2}M_0 \left(\frac{5e}{2\sqrt{\kappa}}\left(\frac{16}{15}\alpha - 1\right)\right)^{\frac{2\alpha}{3}} n^{-\frac{2\alpha}{3}} \\ &\leq \frac{9}{4}e^{-2}M_0 \left(\frac{8e}{3\sqrt{\kappa}}\alpha\right)^{\frac{2\alpha}{3}} n^{-\frac{2\alpha}{3}}. \end{aligned}$$

where  $M_0 = F(x_0) - F^*$ .

## Acknowledgments

J.-F. Aujol acknowledges the support of the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No777826. The authors acknowledge the support of the French Agence Nationale de la Recherche (ANR) under reference ANR-PRC-CE23 MaSDOL and the support of FMJH Program PGMO 2019-0024 and from the support to this program from EDF-Thales-Orange.

## A The continuous case - Proof of Theorem 5

Our analysis is based on the following Lyapunov energy:

$$\mathcal{E}(t) = t^2(F(x(t)) - F^*) + \frac{1}{2}\|\lambda(x(t) - x^*) + t\dot{x}(t)\|^2, \quad \lambda = \frac{2\alpha}{\gamma+2} \quad (78)$$

where the parameter  $\lambda$  is chosen accordingly to [6]. Remember that the expected asymptotic convergence rate is polynomial in  $\mathcal{O}\left(t^{-\frac{2\alpha\gamma}{\gamma+2}}\right)$  [6] with an exponent equal to  $\lambda\gamma$ . Differentiating the Lyapunov energy  $\mathcal{E}$ , we easily prove that:

$$\begin{aligned} \mathcal{E}'(t) + \frac{\gamma\lambda - 2}{t}\mathcal{E}(t) &= \lambda\gamma t \left( F(x(t)) - F^* - \frac{1}{\gamma}\langle \nabla F(x(t)), x(t) - x^* \rangle \right) \\ &\quad + \frac{\lambda^2(\gamma\lambda - 2)}{2t}\|x(t) - x^*\|^2 + (\lambda^2(\gamma + 1) - \lambda - \alpha\lambda)\langle x(t) - x^*, \dot{x}(t) \rangle \\ &\quad + t(\lambda + 1 - \alpha + \frac{\gamma\lambda - 2}{2})\|\dot{x}(t)\|^2. \end{aligned}$$

Using the flatness assumption and replacing  $\lambda = \frac{2\alpha}{\gamma+2}$ , we finally get:

$$\mathcal{E}'(t) + \frac{\gamma\lambda - 2}{t}\mathcal{E}(t) \leq K(\alpha) \left( \frac{2\alpha}{(\gamma+2)t}\|x(t) - x^*\|^2 + \langle x(t) - x^*, \dot{x}(t) \rangle \right) \quad (79)$$

where:  $K(\alpha) = \frac{2\alpha\gamma}{(\gamma+2)^2}(\alpha - 1 - \frac{2}{\gamma})$ . We now need to control the scalar product whose sign is unknown. Combining the following two inequalities:

$$|\langle x(t) - x^*, \dot{x}(t) \rangle| \leq \frac{\sqrt{\mu}}{2}\|x(t) - x^*\|^2 + \frac{1}{2\sqrt{\mu}}\|\dot{x}(t)\|^2 \quad (80)$$

where the coefficients to bound the scalar product  $\sqrt{\mu}$  are chosen to get the tightest control on the energy, and

$$t^2 \|\dot{x}(t)\|^2 \leq \left(1 + \theta \frac{\alpha}{t\sqrt{\mu}}\right) \|\lambda(x(t) - x^*) + t\dot{x}(t)\|^2 + \lambda^2 \left(1 + \frac{t\sqrt{\mu}}{\theta\alpha}\right) \|x(t) - x^*\|^2 \quad (81)$$

for any  $\theta > 0$ , we get:

$$\begin{aligned} \mathcal{E}'(t) + \frac{\gamma\lambda - 2}{t} \mathcal{E}(t) &\leq K(\alpha) \left[ \frac{\sqrt{\mu}}{2} + \frac{2\alpha}{(\gamma+2)t} \left(1 + \frac{1}{(\gamma+2)\theta}\right) + \frac{2\alpha^2}{(\gamma+2)^2\sqrt{\mu}t^2} \right] \|x(t) - x^*\|^2 \\ &\quad + \frac{K(\alpha)}{2\sqrt{\mu}t^2} \left(1 + \theta \frac{\alpha}{t\sqrt{\mu}}\right) \|\lambda(x(t) - x^*) + t\dot{x}(t)\|^2 \end{aligned} \quad (82)$$

$$\begin{aligned} &\leq \frac{2}{\mu} K(\alpha) \left[ \frac{\sqrt{\mu}}{2} + \frac{2\alpha}{(\gamma+2)t} \left(1 + \frac{1}{(\gamma+2)\theta}\right) + \frac{2\alpha^2}{(\gamma+2)^2\sqrt{\mu}t^2} \right] (F(x(t)) - F^*) \\ &\quad + \frac{K(\alpha)}{2\sqrt{\mu}t^2} \left(1 + \theta \frac{\alpha}{t\sqrt{\mu}}\right) \|\lambda(x(t) - x^*) + t\dot{x}(t)\|^2 \end{aligned} \quad (83)$$

using the growth condition  $\mathcal{G}_\mu^2$ . We then choose the parameter  $\theta$  to make equal the coefficients before  $\frac{1}{t^3}$  in  $t^2(F(x(t)) - F^*)$  and  $\frac{1}{2} \|\lambda(x(t) - x^*) + t\dot{x}(t)\|^2$ , i.e. such that:

$$\frac{2}{\mu} \frac{2\alpha}{(\gamma+2)} \left(1 + \frac{1}{(\gamma+2)\theta}\right) = \frac{\theta\alpha}{\mu} \quad (84)$$

or equivalently:

$$(\gamma+2)^2\theta^2 - 4(\gamma+2)\theta - 4 = 0. \quad (85)$$

A straightforward computation shows that this last equation has exactly one positive root:

$$\theta = \frac{2}{\gamma+2}(1 + \sqrt{2}). \quad (86)$$

For these choice of parameters, we have:

$$\mathcal{E}'(t) + \frac{\gamma\lambda - 2}{t} \mathcal{E}(t) \leq \frac{K(\alpha)}{\mu t^2} \left( \sqrt{\mu} + \frac{2\alpha}{(\gamma+2)t} (1 + \sqrt{2}) + \frac{4\alpha^2}{(\gamma+2)^2\sqrt{\mu}t^2} \right) \mathcal{E}(t). \quad (87)$$

Let us now define:

$$\varphi(t) := \frac{K(\alpha)}{\mu t^2} \left( \sqrt{\mu} + \frac{2\alpha}{(\gamma+2)t} (1 + \sqrt{2}) + \frac{4\alpha^2}{(\gamma+2)^2\sqrt{\mu}t^2} \right) \quad (88)$$

and:  $\Phi(t) = \int_t^{+\infty} \varphi(x) dx$ . We so have:

$$\forall t \geq t_0, \quad \mathcal{E}'(t) + \frac{\gamma\lambda - 2}{t} \mathcal{E}(t) \leq \varphi(t) \mathcal{E}(t).$$

Consequently the function  $t \mapsto \mathcal{E}(t)t^{\lambda\gamma-2}e^{\Phi(t)}$  is non-increasing, and for any  $t_1 \in \mathbb{R}$ , we get:

$$\forall t \geq t_1, \quad \mathcal{E}(t) \leq \mathcal{E}(t_1) \left(\frac{t_1}{t}\right)^{\lambda\gamma-2} e^{\Phi(t_1) - \Phi(t)}. \quad (89)$$

A good choice of  $t_1$  is one ensuring a control as tight as possible on the energy  $\mathcal{E}$ . For that,  $t_1$  is chosen such that  $t_1$  minimizes the function  $u \mapsto u^{\lambda\gamma-2}e^{\Phi(u)}$  i.e. such that  $t_1$  satisfies the equation:

$$\frac{\lambda\gamma - 2}{u} - \varphi(u) = 0 \quad (90)$$

Noticing that:  $\lambda\gamma - 2 = \frac{\gamma+2}{\alpha}K(\alpha)$  and simplifying the equation by  $K(\alpha)$ , the equation can be rewritten as:

$$\frac{\gamma+2}{\alpha u} = \frac{1}{\mu u^2} \left( \sqrt{\mu} + \frac{2\alpha}{(\gamma+2)u}(1+\sqrt{2}) + \frac{4\alpha^2}{(\gamma+2)^2\sqrt{\mu}u^2} \right). \quad (91)$$

Introducing  $r = (\gamma+2)\frac{\sqrt{\mu}}{\alpha}u$ , we finally have to solve:

$$r^3 - r^2 - 2(1+\sqrt{2})r - 4 = 0. \quad (92)$$

A straightforward computation shows that the polynomial  $r \mapsto r^3 - r^2 - 2(1+\sqrt{2})r - 4$  has only one real root:  $r^* \simeq 3$  (for which Python gives us an analytical value).

Defining  $t_1 = \frac{\alpha}{(\gamma+2)\sqrt{\mu}}r^*$ , the control on the energy is given by:

$$\forall t \geq t_1, \mathcal{E}(t) \leq \mathcal{E}\left(\frac{\alpha}{(\gamma+2)\sqrt{\mu}}r^*\right) \left(\frac{\alpha r^*}{t(\gamma+2)\sqrt{\mu}}\right)^{\gamma\lambda-2} e^{\Phi(t_1)-\Phi(t)}. \quad (93)$$

Observe now that the term  $\mathcal{E}\left(\frac{\alpha}{(\gamma+2)\sqrt{\mu}}r^*\right)$  can be bounded by the mechanical energy of the system:

$$E_m(t) = F(x(t)) - F^* + \frac{1}{2}\|\dot{x}(t)\|^2 \quad (94)$$

Note that this energy is non-increasing since:  $E'_m(t) = \langle \nabla F(x(t)) + \ddot{x}(t), \dot{x}(t) \rangle = -\frac{\alpha}{t}\|\dot{x}(t)\|^2 \leq 0$ , hence  $E_m$  is uniformly bounded on  $[t_0, +\infty[$ . We then have:

$$\begin{aligned} \mathcal{E}(t_1) &= t_1^2(F(x(t_1)) - F^*) + \frac{1}{2} \left\| \frac{2\alpha}{\gamma+2}(x(t_1) - x^*) + t_1\dot{x}(t_1) \right\|^2 \\ &= t_1^2(F(x(t_1)) - F^* + \frac{1}{2}\|\dot{x}(t_1)\|^2) + \frac{2\alpha^2}{(\gamma+2)^2}\|x(t_1) - x^*\|^2 + \frac{2\alpha}{\gamma+2}t_1\langle x(t_1) - x^*, \dot{x}(t_1) \rangle \\ &= t_1^2 E_m(t_1) + \frac{2\alpha^2}{(\gamma+2)^2}\|x(t_1) - x^*\|^2 + \frac{2\alpha}{\gamma+2}t_1\langle x(t_1) - x^*, \dot{x}(t_1) \rangle \end{aligned}$$

Using again (80) to control the scalar product combined with the quadratic growth condition  $\mathcal{G}_\mu^2$ , we can prove that:

$$\begin{aligned} 2\langle x(t_1) - x^*, \dot{x}(t_1) \rangle &\leq \sqrt{\mu}\|x(t_1) - x^*\|^2 + \frac{1}{\sqrt{\mu}}\|\dot{x}(t_1)\|^2 \\ &\leq \frac{2}{\sqrt{\mu}}(F(x(t_1)) - F^*) + \frac{1}{\sqrt{\mu}}\|\dot{x}(t_1)\|^2 = \frac{2}{\sqrt{\mu}}E_m(t_1) \end{aligned}$$

Noticing that the quadratic growth condition also implies:

$$\|x(t_1) - x^*\|^2 \leq \frac{2}{\mu}(F(x(t_1)) - F^*) \leq \frac{2}{\mu}E_m(t_1)$$

and remembering that  $t_1 = \frac{\alpha}{(\gamma+2)\sqrt{\mu}}r^*$ , we finally get:

$$\mathcal{E}(t_1) \leq t_1^2 E_m(t_1) + \frac{2\alpha^2}{(\gamma+2)^2}\|x(t_1) - x^*\|^2 + \frac{2\alpha}{(\gamma+2)\sqrt{\mu}}t_1 E_m(t_1) \quad (95)$$

$$\leq \left[ t_1^2 + \frac{4\alpha^2}{(\gamma+2)^2\mu} + \frac{2\alpha}{(\gamma+2)\sqrt{\mu}}t_1 \right] E_m(t_1) = \left( 1 + \frac{2}{r^*} + \frac{4}{r^{*2}} \right) t_1^2 E_m(t_1) \quad (96)$$

$$\leq \left( 1 + \frac{2}{r^*} + \frac{4}{r^{*2}} \right) t_1^2 E_m(t_0) \quad (97)$$

Observe that the primitive  $\Phi(t) = \int_t^{+\infty} \varphi(x)dx$  of  $\varphi$  has a simple analytic expression showing that  $\Phi$  is non-positive and:

$$\Phi(t_1) = (\gamma + 2) \frac{K(\alpha)}{\alpha} \left( \frac{1}{r^*} + \frac{1 + \sqrt{2}}{r^{*2}} + \frac{4}{3r^{*3}} \right) \quad (98)$$

We finally obtain the following control on the values:

$$F(x(t)) - F^* \leq C_1 E_m(t_0) \left( \frac{\alpha r^*}{t(\gamma + 2)\sqrt{\mu}} \right)^{\frac{2\alpha\gamma}{\gamma+2}} e^{\frac{2\gamma}{\gamma+2} C_2(\alpha-1-\frac{2}{\gamma})} \quad (99)$$

where:  $C_1 = 1 + \frac{2}{r^*} + \frac{4}{r^{*2}}$ ,  $C_2 = \frac{1}{r^*} + \frac{1+\sqrt{2}}{r^{*2}} + \frac{4}{3r^{*3}}$ .

## B Technical Lemmas for Theorem 6

The proof of Theorem 6 is based on the following Lyapunov energy:

$$E_n = 2sn^2(F(x_n) - F^*) + \|\lambda(x_{n-1} - x^*) + n(x_n - x_{n-1})\|^2 \quad (100)$$

which can be rewritten as:

$$E_n = n^2 w_n + (\lambda^2 - \lambda n) h_{n-1} + (n^2 - \lambda n) \delta_n + \lambda n h_n \quad (101)$$

using the reduced notations (65).

### B.1 Proof of Lemma 1.

**First step:** using the reduced notations (65), we prove that:

$$\begin{aligned} E_{n+1} - \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) E_n &\leq \frac{4\alpha K(\alpha) h_n}{3} + A_1(n, \alpha) \delta_n + B_1(n, \alpha) (h_{n-1} - h_n) \\ &\quad + B_3(n, \alpha) (h_{n+1} - h_n - \delta_{n+1}) \end{aligned} \quad (102)$$

with:

$$\begin{aligned} A_1(n, \alpha) &= \frac{17\alpha^2}{9} - \frac{8\alpha}{3} + 2 - \alpha \frac{(10\alpha^2 - 18\alpha + 9)n + 7\alpha^3 - 12\alpha^2 + 6\alpha}{3(n + \alpha)^2}, \\ B_1(n, \alpha) &= -\frac{2}{9}\alpha^2 + \frac{4}{3}\alpha - 1 + \frac{1}{3} \frac{3\alpha - 2\alpha^3}{n + \alpha} + \frac{1}{27} \frac{8\alpha^3 - 24\alpha^2}{n}, \quad B_3(n, \alpha) = \frac{2}{3}\alpha - 1. \end{aligned}$$

Indeed:

$$\begin{aligned} E_{n+1} - \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) E_n &= (n+1)^2 w_{n+1} - \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) n^2 w_n \\ &\quad + ((n+1)^2 - \lambda(n+1)) \delta_{n+1} - \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) (n^2 - \lambda n) \delta_n \\ &\quad + \left(\lambda^2 - \lambda(n+1) - \lambda n \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right)\right) h_n + \lambda(n+1) h_{n+1} \\ &\quad - (\lambda^2 - \lambda n) \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) h_{n-1} \end{aligned} \quad (103)$$

Observe now that:

$$\begin{aligned} (n+1)^2 w_{n+1} - \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) n^2 w_n &= \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) n^2 (w_{n+1} - w_n) + \left((n+1)^2 - n^2 \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right)\right) w_{n+1} \\ &= n \left(n - \left(\frac{2\alpha}{3} - 2\right)\right) (w_{n+1} - w_n) + \left(\frac{2\alpha}{3} n + 1\right) w_{n+1} \end{aligned}$$

Combining the two following inequalities

$$w_{n+1} - w_n \leq \alpha_n^2 \delta_n - \delta_{n+1} \quad (104)$$

from [14] and:

$$w_{n+1} \leq \|x_n + \alpha_n(x_n - x_{n-1}) - x^*\|^2 - \|x_{n+1} - x^*\|^2$$

from [3], or equivalently with our notations:

$$w_{n+1} \leq (1 + \alpha_n)h_n - \alpha_n h_{n-1} - h_{n+1} + (\alpha_n + \alpha_n^2)\delta_n \quad (105)$$

we then deduce:

$$\begin{aligned} & (n+1)^2 w_{n+1} - \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) n^2 w_n \\ & \leq n \left(n - \frac{2\alpha}{3} + 2\right) (\alpha_n^2 \delta_n - \delta_{n+1}) + \left(\frac{2\alpha}{3}n + 1\right) ((1 + \alpha_n)h_n - \alpha_n h_{n-1} - h_{n+1} + (\alpha_n + \alpha_n^2)\delta_n) \end{aligned}$$

It follows:

$$E_{n+1} - \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) E_n \leq A_1(n, \alpha)\delta_n + A_2(n, \alpha)\delta_{n+1} + B_1(n, \alpha)h_{n-1} + B_2(n, \alpha)h_n + B_3(n, \alpha)h_{n+1} \quad (106)$$

where:

$$\begin{aligned} A_1(n, \alpha) &= \frac{17\alpha^2}{9} - \frac{8\alpha}{3} + 2 - \alpha \frac{(10\alpha^2 - 18\alpha + 9)n + 7\alpha^3 - 12\alpha^2 + 6\alpha}{3(n + \alpha)^2}, & A_2(n, \alpha) &= 1 - \frac{2\alpha}{3} \\ B_1(n, \alpha) &= -\frac{2}{9}\alpha^2 + \frac{4}{3}\alpha - 1 + \frac{1}{3} \frac{3\alpha - 2\alpha^3}{n + \alpha} + \frac{1}{27} \frac{8\alpha^3 - 24\alpha^2}{n}, \end{aligned}$$

and

$$B_2(n, \alpha) = \frac{2}{9}\alpha^2 - 2\alpha + 2 - \frac{1}{3} \frac{3\alpha - 2\alpha^3}{n + \alpha}, \quad B_3(n, \alpha) = \frac{2}{3}\alpha - 1. \quad (107)$$

Observe now that:  $A_2(n, \alpha) = -B_3(n, \alpha)$  and:

$$B_1(n, \alpha) + B_2(n, \alpha) + B_3(n, \alpha) = \frac{8\alpha^2}{27} \frac{\alpha - 3}{n} = \frac{4\alpha K(\alpha)}{3n}.$$

so that (106) becomes:

$$\begin{aligned} E_{n+1} - \left(1 - \frac{\frac{2\alpha}{3} - 2}{n}\right) E_n & \leq \frac{4\alpha K(\alpha)}{3} \frac{h_n}{n} + A_1(n, \alpha)\delta_n + B_1(n, \alpha)(h_{n-1} - h_n) \\ & \quad + B_3(n, \alpha)(h_{n+1} - h_n - \delta_{n+1}) \end{aligned} \quad (108)$$

**Step 2:** First observe that combining the growth condition  $\mathcal{G}_\mu^2$  with the control of the values by the energy (namely:  $E_n \geq n^2 w_n$  for all  $n$ ), we have:

$$\forall n \in \mathbb{N}^*, \quad \frac{h_n}{n} \leq \frac{w_n}{\kappa n} \leq \frac{E_n}{\kappa n^3} \leq \frac{E_n}{\kappa n (n - \frac{2\alpha}{3})^2},$$

so that applying the following Lemma whose proof is detailed in Appendix B.4:

**Lemma 4** For all  $n \geq 1$  and any  $(A, B) \in \mathbb{R}^2$

$$A\delta_n + B(h_{n-1} - h_n) \leq \left(2|A + B| + \frac{\sqrt{2}|B|}{\sqrt{s\mu}}\right) \left(1 + \frac{4\alpha^2}{9s\mu n^2}\right) \frac{E_n}{(n - \frac{2\alpha}{3})^2}.$$



we can prove that:

$$\frac{4\alpha K(\alpha)}{3} \frac{h_n}{n} + A_1(n, \alpha)\delta_n + B_1(n, \alpha)(h_{n-1} - h_n) \leq \frac{\tilde{C}_1(n, \alpha, \kappa)E_n}{(n - \frac{2\alpha}{3})^2} \quad (109)$$

and:

$$B_3(n, \alpha)(h_{n+1} - h_n - \delta_{n+1}) \leq \frac{\tilde{C}_2(n, \alpha, \kappa)E_{n+1}}{(n + 1 - \frac{2\alpha}{3})^2} \quad (110)$$

where:

$$\tilde{C}_1(n, \alpha, \kappa) = 2 \left| \frac{5}{3}\alpha^2 - \frac{4\alpha}{3} + 1 + R(n, \alpha) \right| + \sqrt{2} \left( \frac{\left| -\frac{2\alpha^2}{9} + \frac{4\alpha}{3} - 1 + Q(n, \alpha) \right|}{\sqrt{\kappa}} \right) \left( 1 + \frac{4\alpha^2}{9\kappa n^2} \right) + \frac{4\alpha K(\alpha)}{3\kappa n} \quad (111)$$

with:

$$\begin{aligned} |R(\alpha, n)| &= \left| A_1(n, \alpha) + B_1(n, \alpha) - \left( \frac{5}{3}\alpha^2 - \frac{4\alpha}{3} + 1 \right) \right| \leq \frac{8\alpha^3}{n} \\ |Q(\alpha, n)| &= \frac{\alpha^3}{3n} \left| n \frac{3 - 2\alpha^2}{\alpha^2(n + \alpha)} + 8 \frac{\alpha - 3}{9\alpha} \right| \leq \frac{\alpha^3}{n}, \end{aligned}$$

and:

$$\tilde{C}_2(n, \alpha, \kappa) = \left( \frac{2\alpha}{3} - 1 \right) \left( 4 + \frac{\sqrt{2}}{\sqrt{\kappa}} \right) \left( 1 + \frac{4\alpha^2}{9\kappa(n+1)^2} \right) \quad (112)$$

Finally observe that since  $\kappa \in [0, 1]$ , for all  $n \geq \frac{4\alpha}{3\sqrt{\kappa}}$ , we have:

$$\frac{1}{n - \frac{2\alpha}{3}} \leq \frac{1}{n} (1 + \sqrt{\kappa}) \quad \text{and} \quad \frac{1}{n + 1 - \frac{2\alpha}{3}} \leq \frac{1}{n + 1} (1 + \sqrt{\kappa}) \quad (113)$$

hence:

$$\forall n \geq \frac{4\alpha}{3\sqrt{\kappa}}, \quad E_{n+1} - \left( 1 - \frac{2\alpha - 2}{n} \right) E_n \leq (1 + \sqrt{\kappa})^2 \left( \tilde{C}_1(n, \alpha, \kappa) \frac{E_n}{n^2} + \tilde{C}_2(n, \alpha, \kappa) \frac{E_{n+1}}{(n+1)^2} \right). \quad (114)$$

**Step 3:** The last step is to uniformly bound the coefficients  $\tilde{C}_1(n, \alpha, \kappa)$  and  $\tilde{C}_2(n, \alpha, \kappa)$  with respect to  $n$ . For any  $n \geq \frac{4\alpha}{3\sqrt{\kappa}}$  and  $\alpha \geq 3$ , we have:

$$\begin{aligned} \tilde{C}_2(n, \alpha, \kappa) &= \left( \frac{2\alpha}{3} - 1 \right) \left( 4 + \frac{\sqrt{2}}{\sqrt{\kappa}} \right) \left( 1 + \frac{4\alpha^2}{9\kappa(n+1)^2} \right) \\ &\leq \frac{5}{4} \sqrt{\frac{2}{\kappa}} \left( \frac{2\alpha}{3} - 1 \right) (1 + 2\sqrt{2\kappa}) \end{aligned}$$

The calculations to bound the coefficient  $\tilde{C}_1(n, \alpha, \kappa)$  are similar but a little more painful. For all  $n \geq \frac{4\alpha}{3\sqrt{\kappa}}$ , we have:

$$4\alpha \frac{K(\alpha)}{3\kappa n} \leq \frac{2\alpha(\alpha - 3)}{9\sqrt{\kappa}} \quad \text{and} \quad \frac{4\alpha^2}{9\kappa n^2} \leq \frac{1}{4}$$

so that for all  $n \geq \frac{4\alpha}{3\sqrt{\kappa}}$ :

$$\begin{aligned} \tilde{C}_1(n, \alpha, \kappa) &= 2 \left| \frac{5}{3}\alpha^2 - \frac{4\alpha}{3} + 1 + R(n, \alpha) \right| + \sqrt{\frac{2}{\kappa}} \left| -\frac{2\alpha^2}{9} + \frac{4\alpha}{3} - 1 + Q(n, \alpha) \right| \left( 1 + \frac{4\alpha^2}{9\kappa n^2} \right) + 4\alpha \frac{K(\alpha)}{3\kappa n} \\ &\leq \frac{5}{4} \sqrt{\frac{2}{\kappa}} \left[ \left| -\frac{2\alpha^2}{9} + \frac{4\alpha}{3} - 1 + Q(n, \alpha) \right| + \frac{4\sqrt{2}\alpha(\alpha - 3)}{45} + \frac{4}{5} \left| \frac{5}{3}\alpha^2 - \frac{4\alpha}{3} + 1 + R(n, \alpha) \right| \sqrt{2\kappa} \right] \end{aligned}$$

Assuming now that  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$ , we have:  $\left| -\frac{2\alpha^2}{9} + \frac{4\alpha}{3} - 1 \right| = \frac{2}{9}(\alpha - 3)^2 - 1$ , and:

$$\tilde{C}_1(n, \alpha, \kappa) \leq \frac{5}{4} \sqrt{\frac{2}{\kappa}} \left[ \frac{2}{9}(\alpha - 3)^2 - 1 + \frac{6\alpha(\alpha - 3)}{45} + |Q(n, \alpha)| + \frac{4}{5} \left| \frac{5}{3}\alpha^2 - \frac{4\alpha}{3} + 1 + R(n, \alpha) \right| \sqrt{2\kappa} \right]$$

Let The coefficient  $\tilde{C}_1(n, \alpha, \kappa)$  can be rewritten as:

$$\forall n \geq \frac{4\alpha}{3\sqrt{\kappa}}, \tilde{C}_1(n, \alpha, \kappa) \leq \frac{5}{4} \sqrt{\frac{2}{\kappa}} P(\alpha) \left[ 1 + \left| \frac{Q(n, \alpha)}{P(\alpha)} \right| + \left( \frac{5\alpha^2 - 4\alpha + 3}{3P(\alpha)} + \left| \frac{R(n, \alpha)}{P(\alpha)} \right| \right) \sqrt{2\kappa} \right].$$

Studying the variations of the functions  $\alpha \mapsto \frac{\alpha^2}{P(\alpha)}$  and  $\alpha \mapsto \frac{5\alpha^2 - 4\alpha + 3}{P(\alpha)}$ , we easily prove that they are uniformly bounded for any real  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$ , so that there exists a real constant  $B > 0$  such that:

$$\forall n \geq \frac{4\alpha}{3\sqrt{\kappa}}, \left| \frac{Q(n, \alpha)}{P(\alpha)} \right| \leq \frac{\alpha^3}{nP(\alpha)} \leq B\sqrt{\kappa}.$$

Likewise:

$$\forall n \geq \frac{4\alpha}{3\sqrt{\kappa}}, \left| \frac{R(n, \alpha)}{P(\alpha)} \right| \leq 8 \frac{\alpha^3}{nP(\alpha)} \leq B\sqrt{\kappa}.$$

It finally exists some real constants  $\tilde{c}_1$  and  $\tilde{c}_2$  such that for any  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$  and any  $n \geq \frac{4\alpha}{3\sqrt{\kappa}}$ ,

$$\tilde{C}_1(n, \alpha, \kappa) \leq \frac{5}{4} \sqrt{\frac{2}{\kappa}} P(\alpha) (1 + \tilde{c}_1 \sqrt{\kappa} + \tilde{c}_2 \kappa). \quad (115)$$

Combining (113) and (115), the inequality (67) holds as expected for any  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$  and without any condition on  $\kappa$ :

$$\forall n \geq \frac{4\alpha}{3\sqrt{\kappa}}, E_{n+1} - \left( 1 - \frac{\frac{2\alpha}{3} - 2}{n} \right) E_n \leq \frac{C_1(\alpha, \kappa) E_n}{n^2} + \frac{C_2(\alpha, \kappa) E_{n+1}}{(n+1)^2}$$

with:

$$C_1(\alpha, \kappa) = \frac{5}{4} \sqrt{\frac{2}{\kappa}} \left[ \frac{2}{9}(\alpha - 3) \left( \frac{8}{5}\alpha - 3 \right) - 1 \right] (1 + \sqrt{\kappa})^2 (1 + \tilde{c}_1 \sqrt{\kappa} + \tilde{c}_2 \kappa) \quad (116)$$

$$C_2(\alpha, \kappa) = \frac{5}{4} \sqrt{\frac{2}{\kappa}} \left( \frac{2\alpha}{3} - 1 \right) (1 + \sqrt{\kappa})^2 (1 + 2\sqrt{2\kappa}). \quad (117)$$

■

## B.2 Proof of Lemma 2

Assume that the energy  $E_n$  satisfies:

$$E_{n+1} - \left( 1 - \frac{\frac{2\alpha}{3} - 2}{n} \right) E_n \leq \frac{C_1(\alpha, \kappa) E_n}{n^2} + \frac{C_2(\alpha, \kappa) E_{n+1}}{(n+1)^2}$$

i.e.:

$$\left( 1 - \frac{C_2(\alpha, \kappa)}{(n+1)^2} \right) E_{n+1} - \left( 1 - \frac{\frac{2\alpha}{3} - 2}{n} + \frac{C_1(\alpha, \kappa)}{n^2} \right) E_n \leq 0. \quad (118)$$

Let  $n_0 \geq \frac{4\alpha}{3\sqrt{\kappa}}$ . We then deduce:

$$\forall n \geq n_0, \log(E_{n+1}) - \log(E_{n_0}) \leq \sum_{k=n_0}^n \log \left( \frac{1 - \frac{\frac{2\alpha}{3} - 2}{k} + \frac{C_1(\alpha, \kappa)}{k^2}}{1 - \frac{C_2(\alpha, \kappa)}{(k+1)^2}} \right). \quad (119)$$

Using now the following classical inequalities:

$$\forall x > -1, \quad \frac{x}{x+1} \leq \log(1+x) \leq x, \quad (120)$$

we get:

$$\log \left( 1 - \frac{\frac{2\alpha}{3} - 2}{k} + \frac{C_1(\alpha, \kappa)}{k^2} \right) \leq -\frac{\frac{2\alpha}{3} - 2}{k} + \frac{C_1(\alpha, \kappa)}{k^2} \quad (121)$$

and

$$-\log \left( 1 - \frac{C_2(\alpha, \kappa)}{(k+1)^2} \right) \leq \frac{C_2(\alpha, \kappa)}{(k+1)^2 - C_2(\alpha, \kappa)} \quad (122)$$

We therefore get:

$$\log \left( \frac{1 - \frac{\frac{2\alpha}{3} - 2}{k} + \frac{C_1(\alpha, \kappa)}{k^2}}{1 - \frac{C_2(\alpha, \kappa)}{(k+1)^2}} \right) \leq -\frac{\frac{2\alpha}{3} - 2}{k} + \frac{C_1(\alpha, \kappa)}{k^2} + \frac{C_2(\alpha, \kappa)}{(k+1)^2 - C_2(\alpha, \kappa)} \quad (123)$$

Hence:

$$\log(E_{n+1}) - \log(E_{n_0}) \leq \sum_{k=n_0}^n \left( -\frac{\frac{2\alpha}{3} - 2}{k} + \frac{C_1(\alpha, \kappa)}{k^2} + \frac{C_2(\alpha, \kappa)}{(k+1)^2 - C_2(\alpha, \kappa)} \right) \quad (124)$$

We are now going to make use of the fact that the functions  $x \mapsto \frac{1}{x}$ ,  $x \mapsto \frac{1}{x^2}$  and  $x \mapsto \frac{C_2(\alpha, \kappa)}{x^2 - C_2(\alpha, \kappa)}$  are decreasing functions on  $(C_2, +\infty)$ . Observe that all coefficients in the very last inequality are actually non negative since  $\alpha \geq \alpha_0 > 3$ . We then have:

$$\int_k^{k+1} \frac{dx}{x} \leq \frac{1}{k}, \quad \frac{1}{k^2} \leq \int_{k-1}^k \frac{dx}{x^2} \quad (125)$$

and:

$$\frac{C_2(\alpha, \kappa)}{(k+1)^2 - C_2(\alpha, \kappa)} \leq \int_k^{k+1} \frac{C_2(\alpha, \kappa)}{x^2 - C_2(\alpha, \kappa)} dx \quad (126)$$

so that:

$$\log(E_{n+1}) - \log(E_{n_0}) \leq -\left(\frac{2\alpha}{3} - 2\right) \int_{n_0}^{n+1} \frac{dx}{x} + C_1(\alpha, \kappa) \int_{n_0-1}^n \frac{dx}{x^2} + C_2(\alpha, \kappa) \int_{n_0}^{n+1} \frac{dx}{x^2 - C_2(\alpha, \kappa)}$$

Noticing that:

$$\frac{1}{x^2 - C_2(\alpha, \kappa)} = \frac{1}{2\sqrt{C_2(\alpha, \kappa)}} \left( \frac{1}{x - \sqrt{C_2(\alpha, \kappa)}} - \frac{1}{x + \sqrt{C_2(\alpha, \kappa)}} \right),$$

we eventually get:

$$\begin{aligned} \log(E_{n+1}) - \log(E_{n_0}) &\leq -\left(\frac{2\alpha}{3} - 2\right) \log \left( \frac{n+1}{n_0} \right) + C_1(\alpha, \kappa) \left( \frac{1}{n_0-1} - \frac{1}{n} \right) \\ &\quad + \frac{\sqrt{C_2(\alpha, \kappa)}}{2} \log \left( \frac{(n+1 - \sqrt{C_2(\alpha, \kappa)})(n_0 + \sqrt{C_2(\alpha, \kappa)})}{(n+1 + \sqrt{C_2(\alpha, \kappa)})(n_0 - \sqrt{C_2(\alpha, \kappa)})} \right) \end{aligned} \quad (127)$$

i.e.:

$$\begin{aligned} \log(E_{n+1}) - \log(E_{n_0}) &\leq -\left(\frac{2\alpha}{3} - 2\right) \log \left( \frac{n+1}{n_0} \right) + C_1(\alpha, \kappa) \left( \frac{1}{n_0-1} - \frac{1}{n} \right) \\ &\quad + \frac{\sqrt{C_2(\alpha, \kappa)}}{2} \left( \log \left( \frac{n+1 - \sqrt{C_2(\alpha, \kappa)}}{n+1 + \sqrt{C_2(\alpha, \kappa)}} \right) + \log \left( \frac{n_0 + \sqrt{C_2(\alpha, \kappa)}}{n_0 - \sqrt{C_2(\alpha, \kappa)}} \right) \right) \end{aligned} \quad (128)$$

Taking the exponential, we get:

$$E_{n+1} \leq E_{n_0} \left( \frac{n+1}{n_0} \right)^{-\left(\frac{2\alpha}{3}-2\right)} \exp(\tilde{\Phi}(n_0) - \tilde{\Phi}(n+1)) \quad (129)$$

with:

$$\tilde{\Phi}(n) = \frac{C_1(\alpha, \kappa)}{n-1} + \frac{\sqrt{C_2(\alpha, \kappa)}}{2} \log \left( \frac{n + \sqrt{C_2(\alpha, \kappa)}}{n - \sqrt{C_2(\alpha, \kappa)}} \right).$$

Let us finally compute a more tractable bound on the function  $\tilde{\Phi}(n)$ : using the inequality  $\log(1+x) \leq x$  for  $x \leq 1$ , we have:

$$0 \leq \log \left( \frac{n + \sqrt{C_2(\alpha, \kappa)}}{n - \sqrt{C_2(\alpha, \kappa)}} \right) = \log \left( 1 + \frac{2\sqrt{C_2(\alpha, \kappa)}}{n - \sqrt{C_2(\alpha, \kappa)}} \right) \leq \frac{2\sqrt{C_2(\alpha, \kappa)}}{n - \sqrt{C_2(\alpha, \kappa)}} \quad (130)$$

Hence we deduce that:

$$0 \leq \frac{\sqrt{C_2(\alpha, \kappa)}}{2} \log \left( \frac{n + \sqrt{C_2(\alpha, \kappa)}}{n - \sqrt{C_2(\alpha, \kappa)}} \right) \leq \frac{C_2(\alpha, \kappa)}{n - \sqrt{C_2(\alpha, \kappa)}} \quad (131)$$

Now, using the definition of the coefficients  $C_1(\alpha, \kappa)$  and  $C_2(\alpha, \kappa)$  given in Lemma 1, we get:

$$\begin{aligned} 0 \leq \tilde{\Phi}(n) &\leq \frac{C_1(\alpha, \kappa)}{n-1} + \frac{C_2(\alpha, \kappa)}{n - \sqrt{C_2(\alpha, \kappa)}} \leq \frac{2C_1(\alpha, \kappa)}{n} + \frac{C_2(\alpha, \kappa)}{n - \sqrt{C_2(\alpha, \kappa)}} \quad (132) \\ &\leq \frac{5}{4n} \sqrt{\frac{2}{\kappa}} (1 + \sqrt{\kappa})^2 \left[ 2P(\alpha) (1 + \tilde{c}_1 \sqrt{\kappa} + \tilde{c}_2 \kappa) + \left( \frac{2\alpha}{3} - 1 \right) \frac{1 + 2\sqrt{2\kappa}}{1 - \frac{\sqrt{C_2(\alpha, \kappa)}}{n}} \right] \end{aligned}$$

where:  $P(\alpha) = \frac{2}{9}(\alpha-3) \left( (1 + \frac{2\sqrt{2}}{5})\alpha - 3 \right) - 1$ . Observe then that for  $\kappa$  small enough and  $n \geq \frac{4\alpha}{3\sqrt{\kappa}}$ ,

$$\frac{1}{1 - \frac{\sqrt{C_2(\alpha, \kappa)}}{n}} \leq \frac{1}{1 - \frac{3\sqrt{C_2(\alpha, \kappa)}}{4\alpha} \sqrt{\kappa}} \leq 1 + 2 \frac{\sqrt{C_2(\alpha, \kappa)}}{\alpha} \sqrt{\kappa}$$

so that there exists a real constant  $\tilde{c}_3$  such that for  $\kappa$  small enough and  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$  we have:

$$\frac{1}{1 - \frac{\sqrt{C_2(\alpha, \kappa)}}{n}} \leq 1 + \tilde{c}_3 \kappa^{1/4}.$$

Therefore we finally get for any  $n \geq \frac{4\alpha}{3\sqrt{\kappa}}$ :

$$\tilde{\Phi}(n) \leq \frac{5}{4n} \sqrt{\frac{2}{\kappa}} (1 + \sqrt{\kappa})^2 \left( 2P(\alpha) (1 + \tilde{c}_1 \sqrt{\kappa} + \tilde{c}_2 \kappa) + \left( \frac{2\alpha}{3} - 1 \right) (1 + 2\sqrt{2\kappa}) (1 + \tilde{c}_3 \kappa^{1/4}) \right). \quad (133)$$

We then deduce that there exists  $C_3 > 0$  (independent to  $\alpha$ ) such that

$$\begin{aligned} \forall n \geq \frac{4\alpha}{3\sqrt{\kappa}}, \tilde{\Phi}(n) &\leq \frac{5}{4n} \sqrt{\frac{2}{\kappa}} \left( 2P(\alpha) + \frac{2\alpha}{3} - 1 \right) (1 + C_3 \kappa^{1/4}) \\ &\leq \frac{5}{4n} \sqrt{\frac{2}{\kappa}} \left( 2P(\alpha) + \frac{2\alpha}{3} \right) (1 + C_3 \kappa^{1/4}) \end{aligned}$$

where:  $2P(\alpha) + \frac{2\alpha}{3} = \frac{2}{3}(\alpha-3) \left( \frac{16}{15}\alpha - 1 \right)$ . Let us introduce:

$$\Phi(n) = \frac{5}{6n} \sqrt{\frac{2}{\kappa}} (\alpha-3) \left( \frac{16}{15}\alpha - 1 \right) (1 + C_3 \kappa^{1/4}).$$

Let  $\alpha \geq 3 + \frac{3}{\sqrt{2}}$  and  $n_0 \geq \frac{4\alpha}{3\sqrt{\kappa}}$ . As expected we finally get:

$$\forall n \geq n_0, E_{n+1} \leq E_{n_0} \left( \frac{n+1}{n_0} \right)^{-\left(\frac{2\alpha}{3}-2\right)} e^{\Phi(n_0)} \quad (134)$$

■

### B.3 Proof of Lemma 3

Let  $M_n$  the mechanical energy:

$$M_n = F(x_n) - F^* + \frac{1}{2s} \|x_n - x_{n-1}\|^2.$$

Let us prove that for any  $n \in \mathbb{N}$ , we have:

$$\frac{E_n}{2sn^2} \leq \left( 1 + \frac{4\alpha^2}{9\kappa n^2} + \frac{4\alpha}{3\sqrt{\kappa}n} \right) M_n = \left( 1 + \frac{2\alpha}{3\sqrt{\kappa}n} \right)^2 M_n \quad (135)$$

First remark that:

$$\begin{aligned} b_n &= \left\| \frac{2\alpha}{3}(x_{n-1} - x^*) + n(x_n - x_{n-1}) \right\|^2 = \left\| \frac{2\alpha}{3}(x_n - x^*) + \left( n - \frac{2\alpha}{3} \right) (x_n - x_{n-1}) \right\|^2 \\ &= \frac{4\alpha^2}{9} \|x_n - x^*\|^2 + \left( n - \frac{2\alpha}{3} \right)^2 \|x_n - x_{n-1}\|^2 + \frac{4\alpha}{3} \left( n - \frac{2\alpha}{3} \right) \langle x_n - x^*, x_n - x_{n-1} \rangle \\ &\leq \frac{4\alpha^2}{9} \|x_n - x^*\|^2 + n^2 \|x_n - x_{n-1}\|^2 + \frac{4\alpha}{3} \left( n - \frac{2\alpha}{3} \right) \langle x_n - x^*, x_n - x_{n-1} \rangle \end{aligned}$$

Using a discrete version of the inequality (80), we have:

$$|\langle x_n - x^*, x_n - x_{n-1} \rangle| \leq \frac{\sqrt{\kappa}}{2} \|x_n - x^*\|^2 + \frac{1}{2\sqrt{\kappa}} \|x_n - x_{n-1}\|^2 \quad (136)$$

so that:

$$b_n \leq \frac{4\alpha^2}{9} \|x_n - x^*\|^2 + n^2 \|x_n - x_{n-1}\|^2 + \frac{2\alpha n}{3} \left( \sqrt{\kappa} \|x_n - x^*\|^2 + \frac{1}{\sqrt{\kappa}} \|x_n - x_{n-1}\|^2 \right) \quad (137)$$

Hence:

$$\begin{aligned} \frac{E_n}{2sn^2} &= F(x_n) - F^* + \frac{1}{2sn^2} b_n \\ &= M_n + \frac{2\alpha^2}{9sn^2} \|x_n - x^*\|^2 + \frac{\alpha}{3sn} \left( \sqrt{\kappa} \|x_n - x^*\|^2 + \frac{1}{\sqrt{\kappa}} \|x_n - x_{n-1}\|^2 \right) \end{aligned}$$

Using now the quadratic growth condition  $\mathcal{G}_\mu^2$  and remembering that:  $s\mu = \kappa$ , we get:

$$\frac{E_n}{2sn^2} \leq \left( 1 + \frac{4\alpha^2}{9\kappa n^2} + \frac{4\alpha}{3\sqrt{\kappa}n} \right) M_n = \left( 1 + \frac{2\alpha}{3\sqrt{\kappa}n} \right)^2 M_n$$

### B.4 Proof of Lemma 4

Let us prove that for all  $n \geq 1$  and any  $(A, B) \in \mathbb{R}^2$

$$A\delta_n + B(h_{n-1} - h_n) \leq \left( 2|A + B| + \frac{\sqrt{2}|B|}{\sqrt{s\mu}} \right) \left( 1 + \frac{4\alpha^2}{9s\mu n^2} \right) \frac{E_n}{\left( n - \frac{2\alpha}{3} \right)^2}. \quad (138)$$

Firstly notice that

$$A\delta_n + B(h_{n-1} - h_n) = (A + B)\delta_n + B(h_{n-1} - h_n - \delta_n) \quad (139)$$

and for any  $\theta > 0$

$$|h_{n-1} - h_n - \delta_n| = 2|\langle x_n - x_{n-1}, x_n - x^* \rangle| \leq \frac{h_n}{\theta} + \theta\delta_n. \quad (140)$$

Combining the last two inequalities, it follows that for any  $\theta > 0$ :

$$A\delta_n + B(h_{n-1} - h_n) \leq (A + B + \theta|B|)\delta_n + \frac{|B|}{\theta}h_n \quad (141)$$

To bound the coefficient of  $\delta_n$  we use a specific expression of  $b_n$ :

$$b_n = \left\| \frac{2\alpha}{3}(x_n - x^*) + (n - \frac{2\alpha}{3})(x_n - x_{n-1}) \right\|^2 \quad (142)$$

Applying the inequality  $\|u\|^2 \leq 2\|u + v\|^2 + 2\|v\|^2$  to  $u = (n - \frac{2\alpha}{3})(x_n - x_{n-1})$  and  $v = \frac{2\alpha}{3}(x_n - x^*)$ , we get:

$$(n - \frac{2\alpha}{3})^2\delta_n \leq 2b_n + \frac{8\alpha^2}{9}h_n. \quad (143)$$

It follows that

$$\delta_n \leq \frac{2}{(n - \frac{2\alpha}{3})^2}b_n + \frac{8\alpha^2}{9(n - \frac{2\alpha}{3})^2}h_n. \quad (144)$$

and thus

$$A\delta_n + B(h_{n-1} - h_n) \leq (|A + B| + \theta|B|)\frac{2}{(n - \frac{2\alpha}{3})^2}b_n + \left( \frac{|B|}{\theta} + \frac{8\alpha^2}{9(n - \frac{3\alpha}{4})^2} \right)h_n \quad (145)$$

Using now the growth condition  $h_n \leq \frac{1}{s\mu}w_n$  for all  $n \in \mathbb{N}$ , we get:

$$A\delta_n + B(h_{n-1} - h_n) \leq (|A + B| + \theta|B|)\frac{2}{(n - \frac{2\alpha}{3})^2}b_n + \left( \frac{|B|}{s\mu\theta} + \frac{8\alpha^2}{9s\mu(n - \frac{2\alpha}{3})^2} \right)w_n \quad (146)$$

Choosing  $\theta = \frac{1}{\sqrt{2s\mu}}$  we finally deduce:

$$A\delta_n + B(h_{n-1} - h_n) \leq (2|A + B| + \frac{\sqrt{2}|B|}{\sqrt{s\mu}})\frac{b_n}{(n - \frac{2\alpha}{3})^2} + \left( \frac{\sqrt{2}|B|}{\sqrt{s\mu}} + (2|A + B| + \frac{\sqrt{2}|B|}{\sqrt{s\mu}})\frac{4\alpha^2}{9s\mu(n - \frac{2\alpha}{3})^2} \right)w_n \quad (147)$$

and

$$A\delta_n + B(h_{n-1} - h_n) \leq \left( 2|A + B| + \frac{\sqrt{2}|B|}{\sqrt{s\mu}} \right) \left( 1 + \frac{4\alpha^2}{9s\mu n^2} \right) \frac{E_n}{(n - \frac{2\alpha}{3})^2}, \quad (148)$$

which concludes the proof of the lemma.

## C Sketch of the proof of Theorem 7

The proof of Theorem 7 follows the same line than the proof of Theorem 6, and is based on the following Lyapunov energy:

$$E_n = 2sn^2(F(x_n) - F^* + \left\| \frac{\alpha}{2}(x_{n-1} - x^*) + (n - \frac{\alpha}{4})(x_n - x_{n-1}) \right\|^2). \quad (149)$$

As in the proof of Theorem 6, the first step of this proof consists in establishing some discrete version of the differential inequality (87):

**Lemma 5** Let  $\alpha > 4 + 2\sqrt{2}$  and  $\kappa = \frac{\mu}{L}$ . There exists  $\kappa_0 > 0$  such that for any  $\kappa \leq \kappa_0$ , there exists some real constants  $\tilde{c}_1$  and  $\tilde{c}_2$  such that:

$$\forall n \geq \frac{3\alpha}{2\sqrt{\kappa}}, E_{n+1} - \left(1 - \frac{\alpha-2}{n}\right) E_n \leq C_1(\alpha, \kappa) \frac{E_n}{n^2} + C_2(\alpha, \kappa) \frac{E_{n+1}}{(n+1)^2} \quad (150)$$

where:

$$C_1(\alpha, \kappa) = \frac{19}{36\sqrt{2\kappa}} (\alpha-2)(2\alpha-1) [1 + \tilde{c}_1\sqrt{\kappa} + \tilde{c}_2\kappa] (1 + \sqrt{\kappa})^2 \quad (151)$$

$$C_2(\alpha, \kappa) = \frac{19(\alpha-2)^2}{72\sqrt{2\kappa}} (1 + 11\sqrt{\kappa}) (1 + \sqrt{\kappa})^2. \quad (152)$$

As in the proof of Theorem 7, the next step consists in integrating the inequality (150):

**Lemma 6** Let  $\alpha \geq 4 + 2\sqrt{2}$  and  $n_0 \geq \frac{3\alpha}{2\sqrt{\kappa}}$ . If  $E_n$  satisfies (150) then there exists a real constant  $C_3 > 0$  such that:

$$\forall n \geq n_0, E_n \leq E_{n_0} \left(\frac{n_0}{n}\right)^{\alpha-2} e^{\Phi(n_0)} \quad (153)$$

with

$$\Phi(n) = \frac{19(\alpha-2)(3\alpha-2)}{24n\sqrt{2\kappa}} \left(1 + C_3\kappa^{1/4}\right). \quad (154)$$

The proofs of Lemma 5 and Lemma 6 are very similar to those of Lemma 1 and Lemma 2 and are omitted here. They can be found in [7].

A good choice for  $n_0$  is one ensuring a control as tight as possible on the values  $F(x_n) - F^*$ . For that  $n_0$  is chosen such that it minimizes the function  $f : x \mapsto x^{\alpha-2} e^{\Phi(x)}$ . A straightforward computation gives:

$$n_0 := \frac{19(3\alpha-2)}{24\sqrt{2\kappa}} \left(1 + C_3\kappa^{1/4}\right). \quad (155)$$

Observe that  $f(n_0) = (en_0)^{\alpha-2}$  and that for any  $\alpha \geq 4 + 2\sqrt{2}$ , the optimized value of  $n_0$  satisfies:  $n_0 > \frac{3\alpha}{2\sqrt{\kappa}}$  without any condition on  $\kappa$  and that reducing  $\kappa_0$  if needed, we get:

$$\forall \alpha \geq 4 + 2\sqrt{2}, \frac{3\alpha}{2\sqrt{\kappa}} \leq \frac{19(3\alpha-2)}{24\sqrt{2\kappa}} \left(1 + C_3\kappa^{1/4}\right) \leq \frac{5\alpha}{\sqrt{2\kappa}}. \quad (156)$$

Hence:

$$\forall n \geq \frac{5\alpha}{\sqrt{2\kappa}}, F(x_n) - F(x^*) \leq \frac{E_n}{2sn^2} \leq \frac{E_{n_0}}{2sn_0^2} \left(\frac{n_0}{n}\right)^\alpha e^{\alpha-2} \quad (157)$$

i.e.:

$$\forall n \geq \frac{5\alpha}{\sqrt{2\kappa}}, F(x_n) - F(x^*) \leq \frac{E_{n_0}}{2se^2n_0^2} \left(e \frac{5\alpha}{2n\sqrt{\kappa}}\right)^\alpha \quad (158)$$

Uniformly bounding the energy  $E_{n_0}$  and noticing that:  $\frac{\alpha}{2n_0\sqrt{\kappa}} \leq \frac{1}{3}$ , we have:

$$\frac{E_{n_0}}{2sn_0^2} \leq \left(1 + \frac{\alpha}{2n_0\sqrt{\kappa}}\right)^2 M_{n_0} \leq \frac{16}{9} M_{n_0}$$

where  $M_n$  denotes the potential energy:  $M_n = F(x_n) - F^* + \frac{1}{2}\|x_n - x_{n-1}\|^2$ . Since the mechanical energy associated to the Nesterov scheme is non-increasing (see [14, Corollary 2]) and  $x_{-1} = x_0$ , we then get:

$$\forall n \geq \frac{5\alpha}{\sqrt{2\kappa}}, F(x_n) - F(x^*) \leq \frac{16}{9} \left(e \frac{5\alpha}{2n\sqrt{\kappa}}\right)^\alpha e^{-2} M_0. \quad (159)$$

## References

- [1] Alamo, T., Krupa, P., Limon, D.: Gradient based restart FISTA. In: 2019 IEEE 58th Conference on Decision and Control (CDC), pp. 3936–3941. IEEE (2019)
- [2] Alamo, T., Limon, D., Krupa, P.: Restart fista with global linear convergence. In: 2019 18th European Control Conference (ECC), pp. 1969–1974. IEEE (2019)
- [3] Apidopoulos, V., Aujol, J.F., Dossal, C., Rondepierre, A.: Convergence rates of an inertial gradient descent algorithm under growth and flatness conditions. *Mathematical Programming* (2020)
- [4] Attouch, H., Chbani, Z.: Fast inertial dynamics and FISTA algorithms in convex optimization. Perturbation aspects. *arXiv preprint arXiv:1507.01367* (2015)
- [5] Attouch, H., Chbani, Z., Peypouquet, J., Redont, P.: Fast convergence of inertial dynamics and algorithms with asymptotic vanishing viscosity. *Mathematical Programming* **168**(1-2), 123–175 (2018)
- [6] Aujol, J.F., Dossal, C., Rondepierre, A.: Optimal convergence rates for Nesterov acceleration. *SIAM Journal on Optimization* **29**(4), 3131–3153 (2019)
- [7] Aujol, J.F., Dossal, C., Rondepierre, A.: FISTA is an automatic geometrically optimized algorithm for strongly convex functions. *HAL preprint* (2021)
- [8] Aujol, J.F., Dossal, C., Rondepierre, A.: Convergence rates of the heavy-ball method for quasi-strongly convex optimization. *SIAM Journal on Optimization* **32**, 1817–1842 (2022)
- [9] Aujol, J.F., Dossal, C., Rondepierre, A.: Convergence rates of the Heavy-Ball method under the Łojasiewicz property. *Mathematical Programming* (2022)
- [10] Aujol, J.F., Dossal, C.H., Labarrière, H., Rondepierre, A.: FISTA restart using an automatic estimation of the growth parameter (2021). URL <https://hal.archives-ouvertes.fr/hal-03153525>. Preprint
- [11] Beck, A., Teboulle, M.: A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences* **2**(1), 183–202 (2009)
- [12] Bolte, J., Daniilidis, A., Lewis, A., Shiotani, M.: Clarke subgradients of stratifiable functions. *SIAM Journal on Optimization* **18**(2), 556–572 (2007)
- [13] Bolte, J., Nguyen, T., Peypouquet, J., Suter, B.: From error bounds to the complexity of first-order descent methods for convex functions. *Math. Program.* **165**(2), 471–507 (2017)
- [14] Chambolle, A., Dossal, C.: On the convergence of the iterates of the “fast iterative shrinkage/thresholding algorithm”. *Journal of Optimization Theory and Applications* **166**(3), 968–982 (2015)
- [15] Chambolle, A., Pock, T.: An introduction to continuous optimization for imaging. *Acta Numerica* **25**, 161–319 (2016)
- [16] Fercoq, O., Qu, Z.: Adaptive restart of accelerated gradient methods under local quadratic growth condition. *IMA Journal of Numerical Analysis* **39**(4), 2069–2095 (2019)
- [17] Garrigos, G., Rosasco, L., Villa, S.: Convergence of the Forward-Backward algorithm: Beyond the worst case with the help of geometry. *Mathematical Programming* (2022)
- [18] Łojasiewicz, S.: Une propriété topologique des sous-ensembles analytiques réels. In: *Les Équations aux Dérivées Partielles* (Paris, 1962), pp. 87–89. Éditions du Centre National de la Recherche Scientifique, Paris (1963)



- [19] Łojasiewicz, S.: Sur la géométrie semi- et sous-analytique. *Annales de l'Institut Fourier. Université de Grenoble* **43**(5), 1575–1595 (1993)
- [20] Necoara, I., Nesterov, Y., Glineur, F.: Linear convergence of first order methods for non-strongly convex optimization. *Mathematical Programming* **175**(1-2), 69–107 (2019)
- [21] Nesterov, Y.: A method of solving a convex programming problem with convergence rate  $o(\frac{1}{k^2})$ . In: *Soviet Mathematics Doklady*, vol. 27(2), pp. 372–376 (1983)
- [22] Nesterov, Y.: Gradient methods for minimizing composite functions. *Mathematical Programming* **140**(1), 125–161 (2013)
- [23] Nesterov, Y.: *Introductory lectures on convex optimization: A basic course*, vol. 87. Springer Science & Business Media (2013)
- [24] Ochs, P., Brox, T., Pock, T.: iPiasco: inertial proximal algorithm for strongly convex optimization. *Journal of Mathematical Imaging and Vision* **53**(2), 171–181 (2015)
- [25] O’Donoghue, B., Candes, E.: Adaptive restart for accelerated gradient schemes. *Foundations of computational mathematics* **15**(3), 715–732 (2015)
- [26] Park, C., Park, J., Ryu, E.K.: Factor- $\sqrt{2}$  acceleration of accelerated gradient methods. *arXiv preprint arXiv:2102.07366* (2021)
- [27] Siegel, J.: Accelerated first-order methods: Differential equations and Lyapunov functions. *arXiv preprint arXiv:1903.05671* (2019)
- [28] Su, W., Boyd, S., Candes, E.J.: A differential equation for modeling Nesterov’s accelerated gradient method: theory and insights. *Journal of Machine Learning Research* **17**(153), 1–43 (2016)
- [29] Taylor, A., Drori, Y.: An optimal gradient method for smooth strongly convex minimization. *Mathematical Programming* (2022)