

MATH 590 2023 FALL REPORT

HONGDA LI

November 19, 2023

Abstract

In this paper, we review the paper written by Walkington [1] on the topic of proximal gradient with Nesterov accelerations. We compare the performance of the FISTA method and some of its variants with numerical experiments on the total variation minimization problem; in addition, we propose a heuristic estimation of the strong convexity parameter and demonstrate that it converges faster when applied. We give a literature review on the frontier for both the theories and applications around FISTA. We correct one misconception that occurred in Walkington [1] regarding Nesterov's proof of lower bound on the optimality of first-order algorithms. We present a better proof of linear convergence of FISTA under strong convexity assumption from Beck [2, theorem 10.7.7] by eliminating an identity used in their proof. Finally, we use the Forward Backward Envelope to adapt smooth non-smooth additive objective to Nesterov's generic accelerated gradient algorithm [3, 2.2.7].

1 Preliminaries

In this section, We initiate the discussion by denoising a one-dimensional signal. The dual objective function of the problem derived in this section motivates the use of Accelerated Proximal Gradient with smooth, non-smooth composite objective function. We summarized it from Walkington [1], sections 1 and 4.

1.1 Signal Denoising in One Dimension

Let a one dimensional signal be $u : [0, 1] \mapsto \mathbb{R}$ and u . Regularizing the derivative of the signal using the L1 norm helps recover the digital signal if it's a piecewise constant function. This is called a Total Variation (TV) minimization. Let \hat{u} denote an observation of u corrupted by noise. The denoised signal is the minimizer of $f(u)$,

$$f(u) = \int_0^1 \frac{1}{2}(u - \hat{u})^2 + \alpha |u'| dt.$$

Practical implementations for modern computing devices would necessitate discretization of the integral.

We use the trapezoidal rule and second-order forward difference for the derivative. Let $\hat{u} \in \mathbb{R}^{N+1}$, a vector in the form of $\hat{u} = [\hat{u}_0 \ \cdots \ \hat{u}_N]$, let $t_0 < \cdots < t_N$ be a sequence of time corresponded to each observation of \hat{u}_i . The time intervals are $h_i = t_i - t_{i-1}$ for $i = 1, \dots, N$, not necessarily equally spaced, making this formulation below is slightly more general than Walkington[1]. We derive the

approximation of the integral. Denote $s_i = u_i - \hat{u}_i$.

$$\begin{aligned}
\frac{1}{2} \int_0^1 (u - \hat{u})^2 dt + \alpha \int_0^1 |u'| dt &\approx \frac{1}{2} \sum_{i=0}^N \left(\frac{s_i^2 + s_{i+1}^2}{2} \right) h_{i+1} + \alpha \sum_{i=1}^N \left| \frac{u_i - u_{i-1}}{h_{i+1}} \right| \\
&\triangleright \text{let } C \in \mathbb{R}^{N \times (N+1)} \text{ be upper bi-diagonal with } (1, -1) \\
&= \frac{1}{2} \left(\frac{s_0^2 h_1}{2} + \frac{s_N^2 h_N}{2} + \sum_{i=1}^{N-1} s_i^2 h_i \right) + \alpha \|Cu\|_1 \\
&\triangleright \text{using } D \in \mathbb{R}^{N \times (N+1)}, \\
&\triangleright D := \text{diag}(h_1/2, h_1, h_2, \dots, h_N, h_N/2) \\
&= \frac{1}{2} \langle u - \hat{u}, D(u - \hat{u}) \rangle + \alpha \|Cu\|_1.
\end{aligned}$$

28 The above formulation suggests a smooth, non-smooth additive composite objective for $f(u)$. The
 29 Proximal Gradient method and its variants can solve this optimization problem. Unfortunately,
 30 the non-smooth part $\alpha \|Cu\|_1$ presents computational difficulty if matrix C is unfriendly for the
 31 prox operator. One way to bypass the difficulty involves reformulating with $p = Cu$ and solving
 32 the dual problem.

33 Dual Reformulation

34 Let $p = Cu$, $C \in \mathbb{R}^{(N+1) \times N}$ with $D \in \mathbb{R}^{(N+1) \times (N+1)}$, we reformulate it into

$$\min_{u \in \mathbb{R}^{N+1}} \left\{ \underbrace{\frac{1}{2} \langle (u - \hat{u}), D(u - \hat{u}) \rangle}_{f(u)} + \underbrace{\alpha \|p\|_1}_{h(p)} \mid p = Cu \right\},$$

35 producing Lagrangian of the form

$$\mathcal{L}((u, p), \lambda) = f(u) + h(p) + \langle \lambda, p - Cu \rangle.$$

The dual is

$$\begin{aligned}
-g(\lambda) &:= \inf_{(u, p) \in \mathbb{R}^{N+1} \times \mathbb{R}^N} \{ \mathcal{L}((u, p), \lambda) \} \\
&= \inf_{(u, p) \in \mathbb{R}^{N+1} \times \mathbb{R}^N} \{ f(u) + h(p) + \langle \lambda, p - Cu \rangle \} \\
&= \inf_{u \in \mathbb{R}^{N+1}} \left\{ f(u) - \langle \lambda, Cu \rangle + \inf_{p \in \mathbb{R}^N} \{ h(p) + \langle \lambda, p \rangle \} \right\} \\
&\leq -f^*(-C^T \lambda) - h^*(p).
\end{aligned}$$

36 The theorem of strong duality applies hence equality. With the assumption that D is positive
 37 definite, we have

$$-g(\lambda) = -\frac{1}{2} \|C^T \lambda\|_{D^{-1}}^2 - \langle \hat{u}, C^T \lambda \rangle - \delta_{[-\alpha, \alpha]^N}(p).$$

Observe that the above admits a hyperbox indicator function that makes the prox operator friendlier because the proximal operator of the indicator is projection; in the case of projecting onto the box, the operator is simple. Given dual variable λ , primal is obtained by

$$\begin{aligned}
u &= \operatorname{argmin}_u \mathcal{L}((u, p), \lambda) \\
\partial_u \mathcal{L}((u, p), \lambda) &= D(u - \hat{u}) - C^T \lambda = \mathbf{0} \\
\implies u &= \hat{u} + D^{-1} C^T \lambda.
\end{aligned}$$

38 $-g(\lambda)$ is easier to optimize, and obtaining the primal solution is also simple since D^{-1} is a diagonal
 39 matrix.

40 1.2 FISTA has Worse Convergence Guarantee for Strongly Convex Ob- 41 jectives

42 The dual objective for a total variation minimization problem is a strongly convex and Lipschitz
 43 smooth function because of the norm induced by the positive definite matrix D^{-1} . It's in a form
 44 where FISTA proposed by [4] can solve with a convergence rate of $\mathcal{O}(1/k^2)$ on the objective value
 45 of the function. However, highlighted in Walkington[1], the proximal gradient method without
 46 acceleration achieves $\mathcal{O}((1 - 1/\kappa)^k)$ convergence rate. Which is faster. The parameter κ is the
 47 condition number; in this case, it would be L/α , where L is the Lipschitz smooth constant of $g(u)$
 48 and α is the strong convexity constant for $g(u)$.

49 We emphasize that the Proximal Gradient without acceleration has better theoretical conver-
 50 gence results than the accelerated version for the class of strongly convex objectives. It sparks
 51 the discussion in this paper on the variants of FISTA, hoping to provide some insights on why
 52 Nesterov's momentum-based method's inability to adapt the convergence rate with objective has
 53 strong convexity. For the terminologies, we use FISTA to refer to the proximal gradient method
 54 presented by Beck and Teboulle[4]. We use the Accelerated Proximal Gradient method (APG) to
 55 refer to the first-order acceleration algorithms developed/inspired by FISTA.

56 Finally, whether the original FISTA[5] or Nesterov Accelerated gradient from 1983 has linear
 57 convergence with the presence of strong convexity (or potentially other weaker conditions) is not
 58 known during our research and literature review.

59 1.3 Outline of the Paper

60 Section 2 consists of 3 parts. The first part reviews the literature on the problem of Total Variation
 61 (TV) minimization for image/signal denoising and deblurring. Presenting FISTA and its variants
 62 is the second part. The third part reviews the algorithmic tricks and improvements applied to
 63 the APG. Section 3 addresses a mistake made in Walkington's writing [1, theorem 2.4]. We will
 64 discuss a first-order method and how a different function achieves the lower complexity bound on
 65 the objective value and iterates for a fixed iteration. We discuss how omitting the details of this
 66 theorem creates potential misconceptions of other frontier research ideas. Section 4 presents a proof
 67 that I adapted from Amir Beck's writing [2, theorem 10.7.7]. The proof is slightly more general,
 68 removing one equality to strengthen interpretability and generality. Section 6 presents plots of
 69 convergence and results of applying variants of APG to the TV problem.

70 2 Literatures Review

71 2.1 Total Variation Minimizations

72 Rudin-Osher and Fatemi introduced the Total Variation (TV) minimization method in [6]. They
 73 pioneer the theories of TV minimization by solving PDE. They discussed the empirical observation
 74 that the L1 regularization term produces sharper images. Walkington [1] gives a basic formulation
 75 of one-dimensional signal denoising. However, it's essential to keep in mind that this is a problem
 76 that motivates a variety of modern computational methods and theories. We will list some of them
 77 for context.

78 Goldstein et al. in [7, 3.2.1] showcased the dual reformulation of a 2D signal recovery with $\|\nabla u\|$
 79 as the regularizations term. We note that this norm is without the squared. A more hardcore,

80 detailed coverage of reformulating the dual with L1 penalty terms for 2D signal recovery is in [5].
81 For a complete survey of the state of arts computational methods applied to TV minimizations, see
82 Chambolle [8]. For a detailed exposition of mathematical theories regarding variational analysis on
83 different types of TV problems and statistical inferences-based interpretations of the TV regular-
84 ization term, consult the work by Chambolle et al.[9]. For frontier work of applying non-convex
85 penalty term and its theoretical guarantee consult [10], [11].

86 Variants of FISTA

87 Different variants of FISTA differ by the sequence involved for their momentum method. Choosing
88 different parameters in [algorithm 1](#) produces variants of FISTA.

Algorithm 1 Generic FISTA

```

1: Input:  $(g, h, x^{(0)})$ 
2:  $y^{(0)} = x^{(0)}, \kappa = L/\sigma$ 
3: for  $k = 0, 1, \dots$  do
4:    $x^{(k+1)} = T_L y^{(k)}$ 
5:    $y^{(k+1)} = x^{(k+1)} + \theta_{k+1}(x^{(k+1)} - x^{(k)})$ 
6:   Execute subroutine  $\mathcal{S}$ .
7: end for

```

89 The scope of [algorithm 1](#) considers the additive composition of convex smooth and nonsmooth
90 $f = g + h$ function with g being a L -smooth function. Changing T_L, θ_{k+1} and \mathcal{S} , produce different
91 variants of FISTA.

- 92 1. Original FISTA proposed by Beck [4] considers $\theta_{k+1} = (t_k - 1)/t_{k+1}$, $t_{k+1}(t_{k+1} - 1) = t_k^2$, with
93 $T_L x = \text{prox}_{L^{-1}h}(x - L^{-1}\nabla g(x))$ and $t_0 = 1$. It achieves $\mathcal{O}(1/k^2)$ on the objective value. Our
94 literature review didn't discover proofs for the convergence of the iterates. We also didn't find
95 proofs for a convergence rate faster than $\mathcal{O}((1 - 1/\kappa)^k)$ under strong convexity.
- 96 2. This is a variant where $\theta_{k+1} = (n + a - 1)/a$, for $a > 2$. T_L is the same as (1.). Proved in
97 Chambolle, Dossal [12], its iterates of this version of FISTA exhibit weak convergence. T_L is
98 the same as (1.).
- 99 3. Using $\theta_{k+1} = (\sqrt{\kappa} - 1)/(\sqrt{\kappa} + 1)$ where $\kappa = L/\sigma$, with σ being the strong convexity constant
100 produces V-FISTA in Beck [2, 10.7.7][1, 3.3]. T_L is the same as (1.).
- 101 4. A modification we proposed is based on (3.), but it estimates σ , the strong convexity constant
102 based $x^{(k)}, x^{(k+1)}, \nabla f(x^{(k+1)}), f(x^{(k)})$ using

$$\sigma \approx \langle \nabla f(x^{(k+1)}) - \nabla f(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle / \|x^{(k+1)} - x^{(k)}\|^2.$$

103 It yields excellent results for the numerical experiments.

- 104 5. MFISTA in Beck[5] is produced by adding \mathcal{S} to be the procedure

$$(y^{(k+1)}, t_{k+1}) = \begin{cases} (x^{(k+1)}, 1) & f(y^{(k+1)}) > f(x^{(k+1)}), \\ (y^{(k+1)}, t_{k+1}) & \text{else.} \end{cases}$$

105 This condition asserts a monotone decrease in the objective value. If the objective with
106 momentum increases, it resets the momentum for the next iteration. It has a convergence rate
107 of $\mathcal{O}(1/k^2)$ back then. Our review of the literature didn't confirm the existence of a proof
108 where it has a faster convergence than $\mathcal{O}((1 - 1/\kappa)^k)$ under the presence of strong convexity.

For our problem posed in section 1, the V-FISTA algorithm variant (3.), when applied to the dual objective, achieves a convergence rate of $\mathcal{O}((1 - 1/\sqrt{\kappa})^k)$ for the function objective. For larger κ , this produces a significantly better convergence rate than gradient descent, which is $\mathcal{O}((1 - 1/\kappa)^k)$, and FISTA, which is $\mathcal{O}(1/k^2)$.

Variants (3.) is simple; unfortunately, obtaining σ in itself could be prohibitively expensive and require knowledge about the inverse of hessian (in the case of our TV Minimization problem). Underestimation of σ slows down the convergence rate. This observation sparks research interest in a method that achieves approximately $\mathcal{O}((1 - 1/\sqrt{k})^k)$ linear convergence rate for strongly convex objectives and still retains $\mathcal{O}(1/k^2)$ for Lipschitz-smooth functions in general. One of the other interests is designing a unified theoretical framework to describe all variants of APG.

To address the first issue, people use the idea of restarting FISTA. Earlier attempts proved an asymptotic fast linear convergence rate under quadratic growth conditions by triggering the restart of FISTA based on the gradient mapping norm. See [13][14]. Aujol et al. [15] proposed an automatic restart algorithm that achieves fast linear convergence without knowing the prior strong convexity (or weaker quadratic growth condition) parameter σ . Their bound is not asymptotic. Later, they developed the idea into a parameter-free algorithm in their work [16]. The interests gather around restarting FISTA because spending too much computational effort would be competing against the Proximal Quasi-Newton method, questioning the use of momentum in the first place.

On the theoretical side, Su et al. [17] identifies a second-order differential equation with the exact limit of (1.). A dynamical system understanding of FISTA and APG, in general, enables a wider variety of mathematical tools. For example, in Attouch and Peypouquet [18], they showed a $o(1/k^2)$ convergence rate of variant (2.) based on the ODE understanding. We emphasize that it's the little-o and not the big-O. In Nesterov [3], he proposed a generic algorithm that can derive variants (1.), (3.), and more using the idea of Estimating Sequences and Functions. He constructed a proof of convergence on his generic algorithm, demonstrating both linear and sub-linear convergence rates (depending on the parameter) on the functions' objective value without assuming the minimizers' existence. For Nesterov's involvement in proving convergence of APG in the non-convex settings, consult [19]. For a theoretical underpinning of Nesterov's generic method in his book, consult Ahn and Sra [20]. They derived a lot of variants of APGs using the Proximal Point method of Rockafellar and discussed the unified theme of a "similar triangle" behind the Nesterov APG method.

Finally, the most recent hardcore idea in theory and practice is from Jang et al. [21]. They squeezed out a constant from the convergence rate of FISTA by formulating the search for a faster algorithm as a QCQP. They call their algorithm OptISTA. Their approach is based on the performance estimation problem (PEP).

3 Nesterov's Lower Bound Clarified

Nesterov discussed his claim of the lower convergence rate for the first-order method on differentiable function in his book [22]. Walkington [1] rephrased his work with one crucial mistake in understanding Nesterov's claim.

A precise understanding is required to prevent confusion and lack of forethought in further research. We detail Nesterov's claim and provide context for understanding the mistakes in Walkington. We comment on how the claim relates to other works at the end of the section. We use the following notations

1. Let $\mathcal{A}_f^k x^{(0)}$ denotes the solution of the k-th iterate $x^{(k)}$ generated by an algorithm $\mathcal{A} \in \text{GA}^{1\text{st}}$, with initial guess $x^{(0)}$, objective function f . With this notation, $x^{(1)} = \mathcal{A}_f x^{(0)}$ and $(\mathcal{A}_f^k x^{(0)})_{k \in \mathbb{N}}$ denotes the sequence generated by $\mathcal{A} \in \text{GA}^{1\text{st}}$, with f .

2. Let $\mathcal{F}_L^{1,1}$ denote the set of convex functions with L -Lipschitz gradient mapping from \mathbb{R}^n to \mathbb{R} .

3.1 First-order Method

The following is rephrased from Assumption [3, 2.1.4]. We came up with the two examples to illustrate the definition for better understanding.

Definition 1 (First Order Method). We are in \mathbb{R}^n for now. Given $x^{(0)} \in \mathbb{R}^n$, an iterative algorithm generates sequence of $(x^{(n)})_{n \in \mathbb{N}}$ in the space. All $\mathcal{A} \in \text{GA}^{1\text{st}}$ satisfy that

$$x^{(j+1)} := \mathcal{A}_f^{j+1} x^{(0)} \in \left\{ x^{(0)} \right\} + \text{span} \left\{ \nabla f \left(x^{(i)} \right) \right\}_{i=1}^{j-1} \quad \forall f, \forall 1 \leq j \leq k-1.$$

Example 3.1 (Fixed Step Descent). The method of fixed-step gradient descent, $x^{(k+1)} = x^{(k)} - L^{-1} \nabla f(x^{(k)})$ is $\bar{\mathcal{A}} \in \text{GA}^{1\text{st}}$ achieves a maximal decrease in objective value for all $f \in \mathcal{F}_L^{1,1}$ given $x^{(k)}$, it can be understood as

$$\bar{\mathcal{A}} \in \operatorname{argmin}_{\mathcal{A} \in \text{GA}^{1\text{st}}} \max_{f \in \mathcal{F}_L^{1,1}} \left\{ f \left(\mathcal{A}_f x^{(k)} \right) \right\}.$$

This method is memoryless because it only matters what $x^{(k)}$, prior iterate $x^{(i)}, 1 \leq i \leq k-1$ plays no role.

Example 3.2 (Steepest Descent). Fix some $f, x^{(k)}$, the method of steepest descent would be $\bar{\mathcal{A}} \in \text{GA}^{1\text{st}}$ and it's

$$\bar{\mathcal{A}} \in \operatorname{argmin}_{\mathcal{A} \in \text{GA}^{1\text{st}}} \left\{ f \left(\mathcal{A}_f x^{(k)} \right) \right\}.$$

This method is also memoryless.

Other methods of $\text{GA}^{1\text{st}}$ include Conjugate Gradient, Quasi-Newton, and Gradient Descent with Momentum.

3.2 Lower Complexity Bounds for L -Lipschitz Smooth Function

The following is Nesterov [3, Thm 2.1.7].

Theorem 3.2.1 (Nesterov's Claim of Lower Bound). For any $1 \leq k \leq 1/2(n-1)$, for all $x^{(0)} \in \mathbb{R}^n$, there exists a Lipschitz smooth convex function in \mathbb{R}^n such that for all algorithm from $\text{GA}^{1\text{st}}$, we have the lower bound for the optimality gap for the function values and its iterates:

$$f \left(x^{(k)} \right) - f^* \geq \frac{3L \|x - x^*\|^2}{32(k+1)^2}, \quad \|x^{(k)} - x^*\|^2 \geq \frac{1}{8} \|x^{(0)} - x^*\|^2.$$

Where x^* is the minimizer of f , so that $f(x^*) = \inf_x f(x)$.

Remark 3.2.1. We emphasize that in [theorem 3.2.1](#) fixes each k and finds a function f such that the lower bound applies at the k -th iterations. Mathematically, it would mean

$$\begin{aligned} \forall 1 \leq k \leq \frac{n+1}{2}, x^{(0)} \in \mathbb{R}^n \exists f \in \mathcal{F}_L^{1,1} \text{ s.t.: } \min_{\mathcal{A} \in \text{GA}^{1\text{st}}} \left\{ f \left(\mathcal{A}_f^k x^{(0)} \right) \right\} - f^* &\geq \frac{3L \|x^{(0)} - x^*\|^2}{32(k+1)^2} \\ \forall 1 \leq k \leq \frac{n+1}{2}, x^{(0)} \in \mathbb{R}^n \max_{f \in \mathcal{F}_L^{1,1}} \min_{\mathcal{A} \in \text{GA}^{1\text{st}}} \left\{ f \left(\mathcal{A}_f^k x^{(0)} \right) - f^* \right\} &\geq \frac{3L \|x^{(0)} - x^*\|^2}{32(k+1)^2} \\ \forall x^{(0)} \in \mathbb{R}^n \min_{1 \leq k \leq 1/2(n+1)} \max_{f \in \mathcal{F}_L^{1,1}} \min_{\mathcal{A} \in \text{GA}^{1\text{st}}} \left\{ f \left(\mathcal{A}_f^k x^{(0)} \right) - f^* \right\} &\geq \frac{3L \|x^{(0)} - x^*\|^2}{32(1/2(n+1))^2}, \end{aligned}$$

A function f_k provides the lower bound for fixed $1 \leq k \leq (n+1)/2$. k parameterized f_k , which Nesterov did in his proof. We emphasize that f_k is different depending on what k is. In addition, observe that minimizer x^* is assumed to exist. We believe that the theorem is generalizable to infinite dimensional Hilbert spaces.

We now quote Walkington [1, theorem 2.4]

Theorem 3.2.2 (Walkington’s Claim of Lower Bound). Let X be an infinite-dimensional Hilbert Space and set $x^{(0)} = \mathbf{0}$. There exists a convex function $f : X \mapsto \mathbb{R}$ with Lipschitz gradient and minimum $f(x_*) > -\infty$ such that for any sequence satisfying

$$x_{i+1} \in \text{Span} \left\{ \nabla f(x^{(0)}), \nabla f(x^{(1)}), \dots, \nabla f(x^{(i)}) \right\}, \quad i = 0, 1, 2, \dots,$$

there holds

$$\min_{1 \leq i \leq n} f(x_i) - f(x_*) \geq \frac{3L\|x_1 - x_*\|^2}{32(n+1)^2},$$

where L is the Lipschitz constant of the gradient.

Remark 3.2.2. Theorem 3.2.2 and theorem 3.2.1 is completely different. The former claims there exists a single function from $\mathcal{F}_L^{1,1}(\mathcal{H})$ introduces the lower bound for all values of k , and all algorithms from GA^{1st}, but the latter didn’t claim that. The difference would remain in infinite dimension Hilbert space if we were to generalize theorem 3.2.1. There is no proof after Walkington’s claim; we can’t know if he had his way of proving the latter claim. It makes us think it is likely a missed detail in his writing.

3.3 Discussion

Walkington cited Bubeck [23, thm 3.14], and Nesterov’s old 2004 book[24] for the lower bound claim. Bubeck has the correct claim, and it’s the same as Nesterov. Attouch’s claim in [18, thm 1] doesn’t contradict theorem 3.2.1 because he fixed function Φ function and the existence of minimizer x^* is not assumed.

4 FISTA Under Strong Convexity

In this section, we show several claims on the convergence proof of the algorithm of V-FISTA. Beck [2, 10.7.7] inspires the works. We removed one identity from their proof to reveal more transparency and interpretability. The original author intends to convince the reader with as few words as possible. We intend to educate and share thoughts. We present the essential claims here to expedite understanding of the big picture. The proofs with details are in the appendix.

4.1 Setting up the Stage

Starting with algorithm 1, $\theta_k = (t_k - 1)/(t_k + 1)$. We consider $f = g + h$, g is Lipschitz smooth with constant L , and strongly convex with constant σ . h is convex. The parameters, θ_k , and t_k will be determined as we review the proof. $t_0 = 1$ is the base case for t_k sequence; it represents the fact that there are no accelerations on the first step of the algorithm; its value depends on what we want it to be.

Here is a list of quantities we constructed for a better exposition.

1. $s^{(k)} = x^{(k)} - x^{(k-1)}$, the velocity vector, for all $k \geq 1$.
2. $e^{(k)} = x^{(k)} - \bar{x}$, the error vector at the k th iteration, for $k \geq 0$.
3. $\theta_k = (t_k - 1)/(t_k + 1)$, which is the momentum step size.
4. $u^{(k)} = \bar{x} + t_k(x^{(k-1)} - x^{(k)}) - x^{(k-1)}$, the error term extrapolated with the velocity. We take $u^{(0)} = \bar{x} - x^{(0)}$.
5. $\delta_k = f(x^{(k)}) - f_{\text{opt}}$, with $f_{\text{opt}} = f(\bar{x}) = \inf_x f(x)$.
6. The quantity R_k plays a crucial role in the proof; it's

$$R_k = \frac{\sigma(t_{k+1} - 1)}{2} \|x^{(k)} - \bar{x}\|^2 - \frac{L - \sigma}{2} \|\bar{x} + t_{k+1}(x^{(k)} - y^{(k)}) - x^{(k)}\|^2$$

7. $\kappa = L/\sigma$, the condition number, it would be that $\kappa \geq 0$.
8. $q = \sigma/L$, the reciprocal of the condition number. It would be that $q \in (0, 1)$.

4.2 Convergence Claim

Lemma 4.2.1. If there exists a sequence $(t_k)_{k \in \mathbb{N}}, C_k$ such that

$$\begin{cases} \frac{C_k}{t_{k+1}^2} = \frac{L(1-t_{k+1}^{-1})}{2t_k^2} \\ R_k + C_k \|u^{(k)}\|^2 \geq 0, \end{cases} \quad (4.2.1)$$

then

$$\delta_{k+1} \leq \left(\prod_{i=0}^k (1 - t_i^{-1}) \right) \left(\delta_0 + \frac{L}{2t_0} \|u^{(0)}\|^2 \right).$$

Proposition 4.1. The choice of $t_k = t_{k+1} = \sqrt{L/\sigma} \forall k \geq 0$ makes the proposed condition in [proposition 4.2.1](#) true. Hence, the V-FISTA variant (3.) has a convergence rate bound

$$\delta_{k+1} \leq \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \left(\delta_0 + \frac{L}{2t_0} \|u^{(0)}\|^2 \right).$$

Proposition 4.2. If the sequence $(t_k)_{k \geq \mathbb{N}}$ by t_k satisfies $t_{k+1} = 1 + \frac{t_k^2(1-q)}{t_k+1}$, and $t_{k+1} \geq t_k + 1 - t_k^2 q$, for all $k \geq 0$, then [proposition 4.1](#) stays true.

5 Adding Non-smoothness in Nesterov's Generic Method

5.1 Setting up the Stage

6 Numerical Experiments

6.1 Subsection

A Appendix

A.1 Proof for Lemma 4.2.1

For better notation we use T to denote the proximal gradient operator T_L appeared in [algorithm 1](#). We need the following lemma to start the proof.

Lemma A.1.1 (Proximal Gradient Lemma). With $f = g + h$, where h is convex, closed and proper, with g be L -Lipschitz smooth with a constant of L , let $y \in \mathbb{E}$, defining $y^+ = T(y)$ being the proximal gradient step, then for any $x \in \mathbb{E}$, we have:

$$f(x) - f(y^+) \geq \frac{L}{2} \|x - y^+\|^2 - \frac{L}{2} \|x - y\|^2 + D_g(x, y),$$

Where $D_g(x, y) := g(x) - g(y) - \langle \nabla g(y), x - y \rangle$ is the Bregman Divergence for the smooth part of the sum: g .

See [2, remark 10.17] for this theorem.

Proof. from lemma,

$$F(x) - F \circ T(y) \geq \frac{L}{2} \|x - T(y)\|^2 - \frac{L}{2} \|x - y\|^2 + D_g(x, y) \quad (\text{A.1.1})$$

$$\triangleright \text{strong convexity of } g \text{ makes } D_g(x, y) \geq \frac{\sigma}{2} \|y - x\|^2 \quad (\text{A.1.2})$$

$$\geq \frac{L}{2} \|x - Ty\|^2 - \frac{L - \sigma}{2} \|x - y\|^2. \quad (\text{A.1.3})$$

With $k \geq 0$, make a sequence $(t_k)_{k \in \mathbb{N}}$, we consider:

$$1. \ x = t_{k+1}^{-1} \bar{x} + (1 - t_{k+1}^{-1}) x^{(k)}, y = y^{(k)}.$$

$$2. \ \bar{x} \in \underset{x}{\operatorname{argmin}} f(x) \text{ and } f(\bar{x}) = f_{\text{opt}}.$$

substituting the above into A.1 then the RHS yields

$$\begin{aligned} x - T(y) &= t_{k+1}^{-1} \bar{x} + (1 - t_{k+1}^{-1}) x^{(k)} - Ty^{(k)} \\ &= t_{k+1}^{-1} \bar{x} + (1 - t_{k+1}^{-1}) x^{(k)} - x^{(k+1)} \\ &= t_{k+1}^{-1} \left(\bar{x} + (t_{k+1} - 1) x^{(k)} - t_{k+1} x^{(k+1)} \right) \\ &= t_{k+1}^{-1} \left(\bar{x} + t_{k+1} \left(x^{(k)} - x^{(k+1)} \right) - x^{(k)} \right) \\ x - y &= t_{k+1}^{-1} \bar{x} + (1 - t_{k+1}^{-1}) x^{(k)} - y^{(k)} \\ &= t_{k+1}^{-1} \left(\bar{x} + (t_{k+1} - 1) x^{(k)} - t_{k+1} y^{(k)} \right) \\ &= t_{k+1}^{-1} \left(\bar{x} + t_{k+1} \left(x^{(k)} - y^{(k)} \right) - x^{(k)} \right), \end{aligned}$$

□

A.2 Cute Subsection

References

- [1] W. Noel, “Nesterov’s Method for Convex Optimization,” *SIAM Review*, vol. 65, no. 2, pp. 539–562. [Online]. Available: <https://epubs-siam-org.eu1.proxy.openathens.net/doi/epdf/10.1137/21M1390037>
- [2] A. Beck, *First-Order Methods in Optimization / SIAM Publications Library*, ser. MOS-SIAM Series in Optimization. SIAM. [Online]. Available: <https://epubs.siam.org/doi/book/10.1137/1.9781611974997>

- [3] Y. Nesterov, “Lecture on Convex Optimizations Chapter 2, Smooth Convex Optimization,” in *Lectures on Convex Optimization*, ser. Springer Optimization and Its Applications, Y. Nesterov, Ed. Cham: Springer International Publishing, 2018, pp. 59–137. [Online]. Available: https://doi.org/10.1007/978-3-319-91578-4_2
- [4] A. Beck and M. Teboulle, “A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems,” *SIAM Journal on Imaging Sciences*, vol. 2, no. 1, pp. 183–202, Jan. 2009. [Online]. Available: <http://epubs.siam.org/doi/10.1137/080716542>
- [5] —, “Fast Gradient-Based Algorithms for Constrained Total Variation Image Denoising and Deblurring Problems,” *IEEE Transactions on Image Processing*, vol. 18, no. 11, pp. 2419–2434, Nov. 2009, conference Name: IEEE Transactions on Image Processing. [Online]. Available: <https://ieeexplore.ieee.org/document/5173518>
- [6] L. I. Rudin, S. Osher, and E. Fatemi, “Nonlinear total variation based noise removal algorithms,” *Physica D: Nonlinear Phenomena*, vol. 60, no. 1, pp. 259–268, Nov. 1992. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/016727899290242F>
- [7] T. Goldstein, C. Studer, and R. Baraniuk, “A Field Guide to Forward-Backward Splitting with a FASTA Implementation,” Dec. 2016, arXiv:1411.3406 [cs]. [Online]. Available: <http://arxiv.org/abs/1411.3406>
- [8] A. Chambolle and T. Pock, “An introduction to continuous optimization for imaging,” *Acta Numerica*, vol. 25, pp. 161–319, 2016, publisher: Cambridge University Press (CUP). [Online]. Available: <https://hal.science/hal-01346507>
- [9] A. Chambolle, V. Caselles, D. Cremers, M. Novaga, and T. Pock, “An Introduction to Total Variation for Image Analysis,” in *Theoretical Foundations and Numerical Methods for Sparse Recovery*, M. Fornasier, Ed. DE GRUYTER, Jul. 2010, pp. 263–340. [Online]. Available: <https://www.degruyter.com/document/doi/10.1515/9783110226157.263/html>
- [10] C. An, H.-N. Wu, and X. Yuan, “Enhanced total variation minimization for stable image reconstruction,” *Inverse Problems*, vol. 39, no. 7, p. 075005, Jul. 2023, arXiv:2201.02979 [cs, eess, math]. [Online]. Available: <http://arxiv.org/abs/2201.02979>
- [11] —, “The springback penalty for robust signal recovery,” *Applied and Computational Harmonic Analysis*, vol. 61, pp. 319–346, Nov. 2022, arXiv:2110.06754 [cs, math]. [Online]. Available: <http://arxiv.org/abs/2110.06754>
- [12] A. Chambolle and C. Dossal, “On the Convergence of the Iterates of the “Fast Iterative Shrinkage/Thresholding Algorithm”,” *Journal of Optimization Theory and Applications*, vol. 166, no. 3, pp. 968–982, Sep. 2015. [Online]. Available: <https://doi.org/10.1007/s10957-015-0746-4>
- [13] T. Alamo, P. Krupa, and D. Limon, “Gradient Based Restart FISTA,” in *2019 IEEE 58th Conference on Decision and Control (CDC)*, Dec. 2019, pp. 3936–3941, iSSN: 2576-2370. [Online]. Available: <https://ieeexplore.ieee.org/document/9029983>
- [14] O. Fercoq and Z. Qu, “Adaptive restart of accelerated gradient methods under local quadratic growth condition,” *IMA Journal of Numerical Analysis*, vol. 39, no. 4, pp. 2069–2095, Oct. 2019, arXiv:1709.02300 [math]. [Online]. Available: <http://arxiv.org/abs/1709.02300>

- 281 [15] J.-F. Aujol, C. H. Dossal, H. Labarrière, and A. Rondepierre, “FISTA restart using
282 an automatic estimation of the growth parameter,” May 2022. [Online]. Available:
283 <https://hal.science/hal-03153525>
- 284 [16] J.-F. Aujol, L. Calatroni, C. Dossal, H. Labarrière, and A. Rondepierre, “Parameter-Free
285 FISTA by Adaptive Restart and Backtracking,” *arXiv.org*, Jul. 2023. [Online]. Available:
286 <https://arxiv.org/abs/2307.14323v1>
- 287 [17] W. Su, S. Boyd, and E. J. Candes, “A Differential Equation for Modeling Nesterov’s
288 Accelerated Gradient Method: Theory and Insights,” *arXiv.org*, Mar. 2015. [Online].
289 Available: <https://arxiv.org/abs/1503.01243v2>
- 290 [18] H. Attouch and J. Peypouquet, “The Rate of Convergence of Nesterov’s Accelerated Forward-
291 Backward Method is Actually Faster Than $1/k^2$,” *SIAM Journal on Optimization*, vol. 26,
292 no. 3, pp. 1824–1834, Jan. 2016, publisher: Society for Industrial and Applied Mathematics.
293 [Online]. Available: <https://epubs.siam.org/doi/10.1137/15M1046095>
- 294 [19] I. Necoara, Y. Nesterov, and F. Glineur, “Linear convergence of first order methods for
295 non-strongly convex optimization,” *Mathematical Programming*, vol. 175, no. 1, pp. 69–107,
296 May 2019. [Online]. Available: <https://doi.org/10.1007/s10107-018-1232-1>
- 297 [20] K. Ahn and S. Sra, “Understanding Nesterov’s Acceleration via Proximal Point Method,”
298 Jun. 2022, arXiv:2005.08304 [cs, math]. [Online]. Available: <http://arxiv.org/abs/2005.08304>
- 299 [21] U. Jang, S. D. Gupta, and E. K. Ryu, “Computer-Assisted Design of Accelerated Composite
300 Optimization Methods: OptISTA,” May 2023, arXiv:2305.15704 [math]. [Online]. Available:
301 <http://arxiv.org/abs/2305.15704>
- 302 [22] Y. Nesterov, *Lectures on Convex Optimization*, ser. Springer Optimization and Its
303 Applications. Cham: Springer International Publishing, 2018, vol. 137. [Online]. Available:
304 <http://link.springer.com/10.1007/978-3-319-91578-4>
- 305 [23] S. Bubeck, “Convex Optimization: Algorithms and Complexity,” Nov. 2015, arXiv:1405.4980
306 [cs, math, stat]. [Online]. Available: <http://arxiv.org/abs/1405.4980>
- 307 [24] Y. Nesterov, *Introductory Lectures on Convex Optimization*, ser. Applied Optimization, P. M.
308 Pardalos and D. W. Hearn, Eds. Boston, MA: Springer US, 2004, vol. 87. [Online]. Available:
309 <http://link.springer.com/10.1007/978-1-4419-8853-9>