

# Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problem

Hongda Li

UBC Okanagan

October 17, 2023

## Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problem

A. Beck and M. Teboulle,

SIAM J. IMAGING SCIENCES, Vol.2, 2009.

SIAM J. IMAGING SCIENCES  
Vol. 2, No. 1, pp. 183–202

© 2009 Society for Industrial and Applied Mathematics

### A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems\*

Amir Beck<sup>†</sup> and Marc Teboulle<sup>‡</sup>

**Abstract.** We consider the class of iterative shrinkage-thresholding algorithms (ISTA) for solving linear inverse problems arising in signal/image processing. This class of methods, which can be viewed as an extension of the classical gradient algorithm, is attractive due to its simplicity and thus is adequate for solving large-scale problems even with dense matrix data. However, such methods are also known to converge quite slowly. In this paper we present a new fast iterative shrinkage-thresholding algorithm (FISTA) which preserves the computational simplicity of ISTA but with a global rate of convergence which is proven to be significantly better, both theoretically and practically. Initial promising numerical results for wavelet-based image deblurring demonstrate the capabilities of FISTA which is shown to be faster than ISTA by several orders of magnitude.

**Key words.** iterative shrinkage-thresholding algorithm, deconvolution, linear inverse problem, least squares and  $\ell_1$  regularization problems, optimal gradient method, global rate of convergence, two-step iterative algorithms, image deblurring

**AMS subject classifications.** 90C25, 90C06, 65F22

**DOI.** 10.1137/080716542

A. Beck's Book: First Order Methods in Optimizations, MOS-SIAM Series on Optimization



Amir Beck and Marc Teboulle.

A fast iterative shrinkage-thresholding algorithm for linear inverse problems.

*SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.



Amir Beck.

*First-Order Methods in Optimization*.

Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017.

- 1 References
- 2 Introduction
  - Context
  - The proximal operator
- 3 Proximal gradient and accelerated proximal gradient
  - Proximal gradient
  - The accelerated proximal gradient
- 4 The momentum term
  - Questions to answer
  - Bounded sequence
  - Some quantities
  - A Sketch of the Proof
- 5 Numerical experiments
  - LASSO
  - Image deconvolution with noise

# Sum of two functions

Consider a function  $f$  that we can write into the sum of two functions.  $\mathbb{E}$  denotes the Euclidean space  $\mathbb{R}^n$  where  $n \in \mathbb{N}$  and  $n$  is finite.

$$\min_{x \in \mathbb{E}} g(x) + h(x), \quad (1)$$

whenever  $g(x)$  is Lipschitz smooth and  $h(x)$ , closed convex and proper, we can use Beck and Teboulle's FISTA algorithm.

1.  $h$  can be nonsmooth.
2. Taking the proximal operator of  $h$  has to be possible and easy to implement (more on that later).
3. Under the right conditions, the FISTA algorithm converges with  $\mathcal{O}(1/k^2)$ .

1. In Beck's and Teboulle's paper [1], they popularized the use of Nesterov Momentum for nonsmooth functions.
2. They proved the convergence with the convexity assumption on  $g, h$ .
3. The FISTA algorithm is provably faster than the alternative algorithms: ISTA, TWIST.

# A major assumption

## Assumption (Assumption A.1)

- 1  $g : \mathbb{E} \mapsto \mathbb{R}$  is **strongly smooth** with constant  $L_g$
- 2  $h : \mathbb{E} \mapsto \bar{\mathbb{R}}$  is **closed convex and proper**.
- 3 Let  $f := g + h$
- 4  $\text{ri} \circ \text{dom}(g) \cap \text{ri} \circ \text{dom}(h) \neq \emptyset$
- 5 A set of minimizers exists for the function  $f$  and the set is bounded.  
Denote the minimizer using  $\bar{x}$ .

## Definition (Strong smoothness)

A differentiable function  $g$  is called strongly smooth with a constant  $\alpha$  if it satisfies:

$$|g(y) - g(x) - \langle \nabla g(x), y - x \rangle| \leq \frac{\alpha}{2} \|x - y\|^2 \quad \forall x, y \in \mathbb{E}. \quad (2)$$

## Remark

When  $g$  is convex, then the absolute value can be removed; the above condition is equivalent to:

$$\|\nabla g(x) - \nabla g(y)\| \leq \alpha \|y - x\| \quad \forall x, y \in \mathbb{E},$$

we assume  $\|\cdot\|$  is the euclidean norm for simplicity. In Beck's book, thm 5.8 [2].



# Proximal operator definition

## Definition (The Proximal Operator)

For a function  $f$  with  $\alpha > 0$ , the proximal operator is defined as:

$$\text{prox}_{f,\alpha}(x) := \arg \min_y \left\{ f(y) + \frac{1}{2\alpha} \|y - x\|^2 \right\}.$$

## Remark

*The proximal operator is a singled-valued mapping when  $f$  is convex, closed, and proper.*

# Set projections

Observe that when  $f$  is an indicator function  $\delta_Q$  defined as:

$$\delta_Q(x) := \begin{cases} 0 & x \in Q, \\ \infty & x \notin Q. \end{cases}$$

the proximal operator of  $\delta_Q$  is

$$P(x) = \underset{\delta_Q, \alpha}{\operatorname{prox}}(x) = \underset{y}{\operatorname{argmin}} \left\{ \delta_Q(y) + \frac{1}{2\alpha} \|x - y\|^2 \right\} = \underset{y \in Q}{\operatorname{argmin}} \|x - y\|^2,$$

it searches for the closest point to the set  $Q$  for all  $\alpha > 0$ , and it is called a projection. The point is unique when  $Q \neq \emptyset$  is convex and closed.

# Example of prox operator

## Definition (Soft thresholding)

For some  $x \in \mathbb{R}$ , the proximal operator of the absolute value is:

$$\text{prox}_{\lambda \|\cdot\|_1, t}(x) = \text{sign}(x) \max(|x| - t\lambda, 0).$$

One could interpret the sign operator as projecting  $x$  onto the interval  $[-1, 1]$  and the  $\max(|x| - t\lambda, 0)$  as the distance of the point  $x$  to the interval  $[-t\lambda, t\lambda]$ .

# The proximal gradient algorithm

## The proximal gradient method

---

### Algorithm Proximal gradient with fixed step-sizes

---

```
1: Input:  $g, h$ , smooth and nonsmooth,  $L$  stepsize,  $x^{(0)}$  an initial guess of solution.  
2: for  $k = 1, 2, \dots, N$  do  
3:    $x^{(k+1)} = \arg \min_y \{h(y) + \langle \nabla g(x^{(k)}), y - x^{(k)} \rangle + \frac{L}{2} \|y - x^{(k)}\|^2\}$   
4:   if  $x^{(k+1)}, x^{(k)}$  close enough then  
5:     Break  
6:   end if  
7: end for
```

---

1. It takes the lowest point on the upper bounding function to go next.
2. It is a fixed-point iteration.

# The upper bounding function

Observe that when  $g$  is Lipschitz smooth with constant  $L_g$  then fix any point  $x$  and for all  $y \in \mathbb{E}$ :

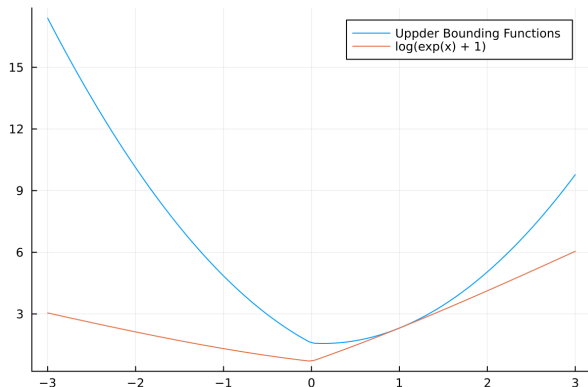
$$g(x) + h(x) \leq g(x) + \nabla g(x)^T (y - x) + \frac{\beta}{2} \|y - x\|^2 + h(y) =: m_x(y|\beta).$$

Moreover, we have:

$$\underbrace{\text{prox}_{h,\beta^{-1}}(x - \beta^{-1} \nabla g(x))}_{=:\mathcal{P}_{\beta^{-1}}^{g,h}(x)} = \arg \min_y \{m_x(y|\beta)\}.$$

The proximal gradient algorithm performs fixed point iterations on the proximal gradient operator.

# Upper bounding function example



**Figure:** The upper bounding function that the proximal gradient algorithm is minimizing for each iteration. In this case  $g(x) = \log(1 + \exp(x))$ ,  $h(x) = |x|$

# Facts about proximal gradient descent

With our Assumption A.1:

1. It converges monotonically with a  $\mathcal{O}(1/k)$  rate as shown in Beck's Book [2] with a step size  $L^{-1}$  such that  $L > L_g$ .
2. When  $h$  is the indicator function, it is just the projected subgradient method. When  $h = 0$ , this is the smooth gradient descent method where the norm of the fixed point error converges with  $\mathcal{O}(1/k)$  [2].
3. The fixed point of the proximal gradient operator is the minimizer of  $f$ .

# The accelerated proximal gradient method

## Momentum template method

**Algorithm** Template proximal gradient method with momentum

- 1: **Input:**  $x^{(0)}, x^{(-1)}, L, h, g$ ; 2 initial guesses and stepsize  $L$
- 2:  $y^{(0)} = x^{(0)} + \theta_k(x^{(0)} - x^{(-1)})$
- 3: **for**  $k = 1, \dots, N$  **do**
- 4:    $x^{(k)} = \text{prox}_{h, L^{-1}}(y^{(k)} + L^{-1}\nabla g(y^{(k)})) = \mathcal{P}_{L^{-1}}^{g, h}(y^{(k)})$
- 5:    $y^{(k+1)} = x^{(k)} + \theta_k(x^{(k)} - x^{(k-1)})$
- 6: **end for**

In the case of FISTA, we use:

$$t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}, \theta_k = \frac{t_k - 1}{t_{k+1}}, t_0 = 1, x^{(-1)} = x^{(0)}$$



# Facts about the accelerated proximal gradient method

1. When  $h = 0$ , the algorithm is Nesterov's famous accelerated gradient method proposed back in 1983.
2. It is no longer a monotone method.
3. It has a convergence rate of  $\mathcal{O}(1/k^2)$  under Assumption A.1, proved by Beck, Teboulle in the FISTA paper [1].

# Some important questions to address for the second half

1. Why does the sequence of  $t_k, \theta_k$  makes sense?
2. What ideas are involved in proving that the convergence rate is  $\mathcal{O}(1/k^2)$ ?
3. If the above is true, what secret sauce cooks up the sequence  $t_k, \theta_k$ ?
  - unfortunately, the secret sauce is not the Nesterov momentum term  $t_k$  but rather an inequality involving two sequences that give us a bound.

# Two bounded sequence

The proof for the convergence rate for deriving the momentum sequence hinges on the following lemma about two sequences of numbers:

## Lemma (Two Bounded Sequences)

*Consider the sequences  $a_k, b_k \geq 0$  for  $k \in \mathbb{N}$  with  $a_1 + b_1 \leq c$ . Inductively the two sequences satisfy  $a_k - a_{k+1} \leq b_{k+1} - b_k$ , which describes a sequence with oscillations bounded by the difference of another sequence. Consider the telescoping sum:*

$$\begin{aligned} a_k - a_{k+1} &\geq b_{k+1} - b_k \quad \forall k \in \mathbb{N} \\ \implies -\sum_{k=1}^N a_{k+1} - a_k &\geq \sum_{k=1}^N b_{k+1} - b_k \\ -(a_{N+1} - a_1) &\geq b_{N+1} - b_1 \\ c \geq a_1 + b_1 &\geq b_{N+1} + a_{N+1} \\ \implies c &\geq a_{N+1}. \end{aligned}$$

# The Nesterov momentum sequence

In the algorithm we listed, the Nesterov momentum sequence consists of  $\theta_k, t_k$  with  $t_0 = 1$  given by

$$t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}, \theta_k = \frac{t_k - 1}{t_{k+1}}, t_0 = 1 \quad (3)$$

It is shown that the sequence of  $t_k$  grows linearly and  $t_k > (1 + k)/2$ .

# Some quantities and one lemma

1.  $v^{(k)} = x^{(k)} - x^{(k-1)}$  is the velocity term.
2.  $\bar{v}^{(k)} = \theta_k v^{(k)}$  is the weighed velocity term.
3.  $e^{(k)} := x^{(k)} - \bar{x}$ , where  $\bar{x} \in \arg \min_x (f(x))$ , where  $\bar{x}$  is fixed.
4.  $\Delta_k := f(x^{(k)}) - f(\bar{x})$  which represent the optimality gap at step  $k$ .

## Lemma

*With assumption A.1, and stepsize  $\beta^{-1} > L_g$  for algorithm 1, let  $y \in \mathbb{E}$  and define  $y^+ = \mathcal{P}_{\beta^{-1}}^{g,h}(y)$  we have for any  $x \in \mathbb{E}$ :*

$$f(x) - f(y^+) \geq \frac{\beta}{2} \|y^+ - y\|^2 + \beta \langle y - x, y^+ - y \rangle.$$

# Some quantities and one lemma

1.  $v^{(k)} = x^{(k)} - x^{(k-1)}$  is the velocity term.
2.  $\bar{v}^{(k)} = \theta_k v^{(k)}$  is the weighed velocity term.
3.  $e^{(k)} := x^{(k)} - \bar{x}$ , where  $\bar{x} \in \arg \min_x (f(x))$ , where  $\bar{x}$  is fixed.
4.  $\Delta_k := f(x^{(k)}) - f(\bar{x})$  which represent the optimality gap at step  $k$ .

## Lemma

*With assumption A.1, and stepsize  $\beta^{-1} > L_g$  for algorithm 1, let  $y \in \mathbb{E}$  and define  $y^+ = \mathcal{P}_{\beta^{-1}}^{g,h}(y)$  we have for any  $x \in \mathbb{E}$ :*

$$f(x) - f(y^+) \geq \frac{\beta}{2} \|y^+ - y\|^2 + \beta \langle y - x, y^+ - y \rangle.$$

# Some quantities and one lemma

1.  $v^{(k)} = x^{(k)} - x^{(k-1)}$  is the velocity term.
2.  $\bar{v}^{(k)} = \theta_k v^{(k)}$  is the weighed velocity term.
3.  $e^{(k)} := x^{(k)} - \bar{x}$ , where  $\bar{x} \in \arg \min_x (f(x))$ , where  $\bar{x}$  is fixed.
4.  $\Delta_k := f(x^{(k)}) - f(\bar{x})$  which represent the optimality gap at step  $k$ .

## Lemma

*With assumption A.1, and stepsize  $\beta^{-1} > L_g$  for algorithm 1, let  $y \in \mathbb{E}$  and define  $y^+ = \mathcal{P}_{\beta^{-1}}^{g,h}(y)$  we have for any  $x \in \mathbb{E}$ :*

$$f(x) - f(y^+) \geq \frac{\beta}{2} \|y^+ - y\|^2 + \beta \langle y - x, y^+ - y \rangle.$$

# Some quantities and one lemma

1.  $v^{(k)} = x^{(k)} - x^{(k-1)}$  is the velocity term.
2.  $\bar{v}^{(k)} = \theta_k v^{(k)}$  is the weighed velocity term.
3.  $e^{(k)} := x^{(k)} - \bar{x}$ , where  $\bar{x} \in \arg \min_x (f(x))$ , where  $\bar{x}$  is fixed.
4.  $\Delta_k := f(x^{(k)}) - f(\bar{x})$  which represent the optimality gap at step  $k$ .

## Lemma

*With assumption A.1, and stepsize  $\beta^{-1} > L_g$  for algorithm 1, let  $y \in \mathbb{E}$  and define  $y^+ = \mathcal{P}_{\beta^{-1}}^{g,h}(y)$  we have for any  $x \in \mathbb{E}$ :*

$$f(x) - f(y^+) \geq \frac{\beta}{2} \|y^+ - y\|^2 + \beta \langle y - x, y^+ - y \rangle.$$



# Form matching to the sequences

From lemma 2.3 in Beck's FISTA paper, one can obtain the following two expressions using the template method and the quantities introduced:

Substitute  $x = x^{(k)}, y = y^{(k+1)}$  lemma 2.3

$$2L^{-1}(\Delta_k - \Delta_{k+1}) \geq \|v^{(k+1)} - \bar{v}^{(k)}\|^2 + 2\langle v^{(k+1)} - \bar{v}^{(k)}, \bar{v}^{(k)} \rangle \quad (*)$$

Substitute  $x = \bar{x}, y = y^{(k+1)}$  lemma 2.3

$$-2L^{-1}\Delta_{k+1} \geq \|v^{(k+1)} - \bar{v}^{(k)}\|^2 + 2\langle v^{(k+1)} - \bar{v}^{(k)}, e^{(k)} + \bar{v}^{(k)} \rangle. \quad (*)$$

# The sequences $t_k$

For now, we don't know what the sequence  $t_k$  is, but we can assume that  $t_k > 1$  for all  $k$  and using  $(*)$ ,  $(\star)$  and consider  $(t_{k+1}^2 - t_{k+1})(*) + t_{k+1}(\star)$  which gives us:

$$\begin{aligned} & 2L^{-1}(\underbrace{(t_{k+1}^2 - t_{k+1})\Delta_k}_{a_k} - \underbrace{t_{k+1}^2\Delta_{k+1}}_{a_{k+1}}) \\ & \geq \dots \text{Non-Trivial Amount of Math is Skipped} \dots \\ & \geq \underbrace{\|t_{k+1}v^{(k+1)} + e^{(k)}\|^2}_{b_{k+1}} - \underbrace{\|e^{(k-1)} + (t_{k+1}\theta_k + 1)v^{(k)}\|^2}_{b_k}, \end{aligned} \tag{\star\star}$$

If the form were to match the Two Bounded Sequences, it has to be the case that  $t_{k+1}\theta_k + 1 = t_k$  and  $t_{k+1}^2 - t_{k+1} = t_k^2$ . In this case, the Nesterov Momentum terms satisfy the conditions perfectly.

# Using the two bounded sequences

It is not hard to show that the base case  $a_1 + b_1 < c$  is bounded using Assumption A1, The Two Bounded Sequences give

$$a_N \leq c \quad (4)$$

$$t_N^2 \Delta_N \leq c \quad (5)$$

$$\Delta_N \leq \frac{c}{t_N^2}, \quad (6)$$

And recall that  $t_k \geq (k+1)/2$ , we can conclude that  $\Delta_N$  convergences with rante  $\mathcal{O}(1/N^2)$ .

## The LASSO problem

LASSO minimizes the 2-norm objective with one norm penalty.

$$\min_x \left\{ \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 \right\}$$

And the prox for  $\|\cdot\|_1$  is given by:

$$(\text{prox}_{\lambda\|\cdot\|_1, t}(x))_i = \text{sign}(x_i) \max(|x_i| - t\lambda, 0),$$

For our experiment:

1.  $A$  has diagonal elements that are numbers equally spaced on the interval  $[0, 2]$ .
2. Vector  $b$  is the diagonal of  $A$  and every odd index is changed into  $\epsilon \sim N(0, 10^{-3})$ .

# Results

The plot of  $\Delta_k$ :

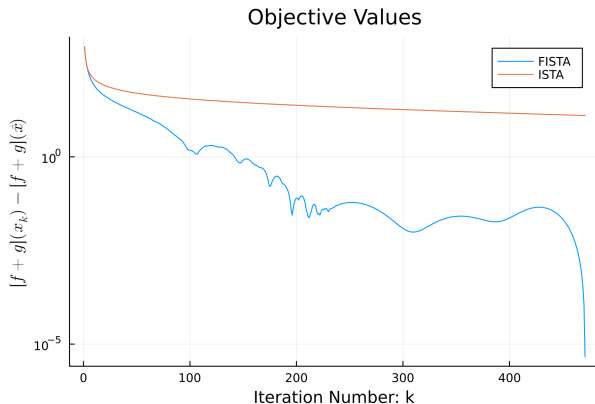
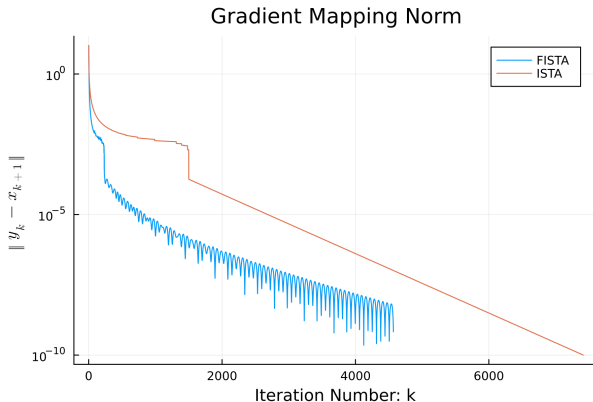


Figure: The left is the objective value of the function during all iterations.

# Results

The plot of  $\|y^{(k)} - x^{(k)}\|_\infty$ :



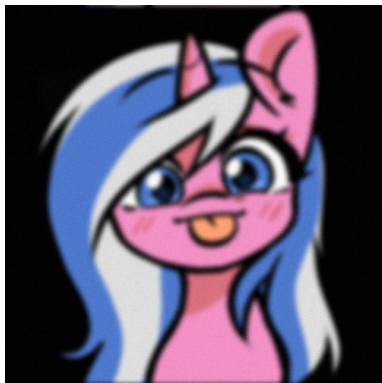
# Experiment Setup

Given an image that is convoluted by a Gaussian kernel with some Gaussian noise, we want to recover the image, given the parameters for convolutions.

1. Gaussian blur with a discrete 15 by 15 kernel is a linear transform represented by a sparse matrix  $A$  in the computer.
2. When an image is 500 by 500 with three color channels,  $A$  is  $750000 \times 750000$ .
3. Let the noise be on all normalized colors values with  $N(0, 10^{-2})$
4. We let  $\lambda = \alpha \times (3 \times 500^2)^{-1}$ .
5. Implemented in Julia, the code is too long to be shown here.

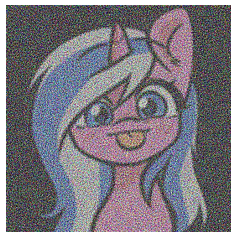
# The blurred image

We consider blurring the image of a pink unicorn that I own.

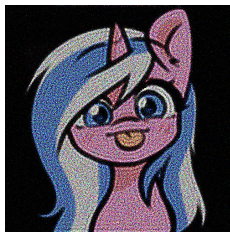


**Figure:** The image is blurred by the Gaussian blurred matrix  $A$  with a tiny amount of noise on the level of  $2 \times 10^{-2}$  that is barely observable. Zoom in to observe the tiny amount of Gaussian noise on top of the blur.

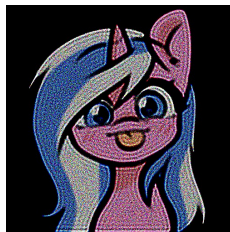




(a)



(b)



(c)

**Figure:** (a)  $\alpha = 0$ , without any one norm penalty, is not robust to the additional noise. (b)  $\alpha = 0.01$ , there is a tiny amount of  $\lambda$ . (c)  $\alpha = 0.1$ , it is more penalty compared to (a).