

# Convergence of ISTA (Beck & Teboulle 2009)

---

## Andersen Ang

ECS, Uni. Southampton, UK  
andersen.ang@soton.ac.uk

Homepage [angms.science](http://angms.science)

Version: July 17, 2023

First draft: Dec 20, 2019

### Content

Problem setup:  $\min f(x) + g(x)$   
ISTA / Proximal gradient method  
Prerequisite  
A key lemma  
Convergence rate of ISTA  $\mathcal{O}(\frac{1}{k})$

Reference: Amir Beck and Marc Teboulle. "A fast iterative shrinkage-thresholding algorithm for linear inverse problems." SIAM journal on imaging sciences, 2009.

# Problem setup: smooth convex optimisation with convex possibly-nonsmooth regularizer

$$(\mathcal{P}) : \operatorname{argmin}_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}. \quad (\text{a convex optimization problem})$$

▶ We consider Euclidean space

▶  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

▶  $f$  is  $L$ -smooth

▶  $f$  is continuously differentiable

▶  $\nabla f$  is  $L$ -Lipschitz

$$\left\{ (\forall \mathbf{a} \in \operatorname{dom} f)(\forall \mathbf{b} \in \operatorname{dom} f) \left\{ f(\mathbf{a}) - f(\mathbf{b}) \leq \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle + \frac{L}{2} \|\mathbf{a} - \mathbf{b}\|_2^2 \right\} \right\}.$$

$f \in \mathcal{C}^1$ , i.e.,  $\nabla f(\mathbf{x})$  exists for all  $\mathbf{x} \in \operatorname{dom} f$

$L > 0$  is the least upper bound in  $\frac{\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|}{\|\mathbf{x} - \mathbf{y}\|} \leq L$

▶  $f$  is convex

all local minima of  $\mathcal{P}$  are global minima

$$\left\{ (\forall \mathbf{x} \in \operatorname{dom} f)(\forall \mathbf{y} \in \operatorname{dom} f) \left\{ f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \right\} \right\}$$

▶  $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$

$$\overline{\mathbb{R}} := [-\infty, +\infty]$$

▶  $g$  is convex, closed, proper and possibly nonsmooth

▶ closed & proper is for proximal operator to work

▶ possibly nonsmooth means we need to replace gradient with subgradient

▶ Details of convexity,  $L$ -smoothness, closeness, proper, see [here](#)

▶  $(\mathcal{P})$  is assumed solvable: the solution set  $\mathcal{X}^* = \operatorname{argmin}_{\mathbf{x}} F \neq \emptyset$

▶ Notation :  $\mathbf{x}^* \in \mathcal{X}^*$  and  $F^* = F(\mathbf{x}^*) \notin \{\pm\infty\}$

## Idea on solving $(\mathcal{P})$

- **Key observation**  $f$  is  $L$ -smooth :  $(\forall \mathbf{x} \in \text{dom} f)(\forall \mathbf{y} \in \text{dom} f)$  we have

$$f(\mathbf{x}) \leq q(\mathbf{x}; \mathbf{y}) := f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2.$$

where  $q(\mathbf{x})$  is a local quadratic over-approximation of  $f$  at a point  $\mathbf{y}$ . Here  $\mathbf{y}$  is treated as a “parameter”.

- **Quadratic over-estimator  $Q$**

By  $f(\mathbf{x}) \leq q(\mathbf{x})$ , for  $F(\mathbf{x})$  we have

$$\begin{aligned} f(\mathbf{x}) &\leq q(\mathbf{x}; \mathbf{y}) && \text{by } \blacksquare \\ \iff f(\mathbf{x}) + g(\mathbf{x}) &\leq q(\mathbf{x}; \mathbf{y}) + g(\mathbf{x}) && \text{add } g \\ \iff F(\mathbf{x}) &\leq q(\mathbf{x}; \mathbf{y}) + g(\mathbf{x}) && \text{by definition } F := f + g \\ \iff F(\mathbf{x}) &\leq Q(\mathbf{x}; \mathbf{y}) && \text{let } Q := q + g \\ \iff F(\mathbf{x}) &\leq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + g(\mathbf{x}) && \text{expand } q \end{aligned}$$

- **Majorization-minimization:**

minimizing  $Q(\mathbf{x}; \mathbf{y}) \xrightarrow{Q \text{ is global upper bound of } F} \text{indirectly minimize } F(\mathbf{x})$

## Discussion on minimizing $Q$

$$Q(\mathbf{x}; \mathbf{y}) := f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + g(\mathbf{x}) \quad \text{is a strictly convex function}$$

- ▶  $f(\mathbf{y})$  is a constant it is a convex function
- ▶  $\langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$  is a linear function on  $\mathbf{x}$  it is a convex function
- ▶  $\frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$  is a quadratic function on  $\mathbf{x}$  it is a strictly convex function by  $L > 0$
- ▶  $g(\mathbf{x})$  it is a convex function by assumption in problem setup

$$\operatorname{argmin}_{\mathbf{x}} \left\{ Q(\mathbf{x}; \mathbf{y}) := f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + g(\mathbf{x}) \right\}$$

- ▶ This is a convex optimization problem: all local minimizers are global
- ▶ This problem is a strictly convex: the global minimizer is unique
- ▶ This problem has nonempty solution set: because we assumed  $g$  is closed and proper

Conclusion :  $\operatorname{argmin}_{\mathbf{x}} Q(\mathbf{x})$  is a  $\mathbb{R}^n \rightarrow \mathbb{R}^n$  map.

## Minimizing $Q \equiv$ computing a proximal operator

$$\begin{aligned} & \operatorname{argmin}_{\mathbf{x}} Q(\mathbf{x}; \mathbf{y}) \\ \equiv & \operatorname{argmin}_{\mathbf{x}} f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + g(\mathbf{x}) && \text{by definition of } Q \\ \iff & \operatorname{argmin}_{\mathbf{x}} f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} \rangle - \langle \nabla f(\mathbf{y}), \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x}\|_2^2 - L \langle \mathbf{x}, \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{y}\|_2^2 + g(\mathbf{x}) && \text{expand} \\ \iff & \operatorname{argmin}_{\mathbf{x}} \langle \nabla f(\mathbf{y}), \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}\|_2^2 - L \langle \mathbf{x}, \mathbf{y} \rangle + g(\mathbf{x}) && \text{argmin ignore constant} \\ \iff & \operatorname{argmin}_{\mathbf{x}} \langle \nabla f(\mathbf{y}) - L\mathbf{y}, \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}\|_2^2 + g(\mathbf{x}) && \text{combine inner product} \\ \iff & \operatorname{argmin}_{\mathbf{x}} \frac{L}{2} \left( -2 \langle \mathbf{y} - \frac{1}{L} \nabla f(\mathbf{y}), \mathbf{x} \rangle + \|\mathbf{x}\|_2^2 \right) + g(\mathbf{x}) && \text{factor out } \frac{L}{2} \\ \iff & \operatorname{argmin}_{\mathbf{x}} \frac{L}{2} \left( \left\| \mathbf{y} - \frac{1}{L} \nabla f(\mathbf{y}) \right\|_2^2 - 2 \langle \mathbf{y} - \frac{1}{L} \nabla f(\mathbf{y}), \mathbf{x} \rangle + \|\mathbf{x}\|_2^2 \right) + g(\mathbf{x}) && \text{argmin ignore constant} \\ \iff & \operatorname{argmin}_{\mathbf{x}} \frac{L}{2} \left\| \mathbf{x} - \left( \mathbf{y} - \frac{1}{L} \nabla f(\mathbf{y}) \right) \right\|_2^2 + g(\mathbf{x}) && \text{completing squares} \\ \iff & \operatorname{prox}_g^L \left( \mathbf{y} - \frac{1}{L} \nabla f(\mathbf{y}) \right) && \text{proximal operator of } g \end{aligned}$$

## The update operator

$$\begin{aligned} p_L(\mathbf{y}) &= \operatorname{prox}_g^L \left( \mathbf{y} - \frac{\nabla f(\mathbf{y})}{L} \right) := \operatorname{argmin}_{\mathbf{x}} \frac{L}{2} \left\| \mathbf{x} - \left( \mathbf{y} - \frac{\nabla f(\mathbf{y})}{L} \right) \right\|_2^2 + g(\mathbf{x}) \\ &= \operatorname{argmin}_{\mathbf{x}} \frac{L}{2} \left\| \mathbf{x} - \left( \mathbf{x}_{k-1} - \frac{\nabla f(\mathbf{x}_{k-1})}{L} \right) \right\|_2^2 + g(\mathbf{x}) \quad \text{put } \mathbf{y} = \mathbf{x}_{k-1} \end{aligned}$$

- ▶  $p_L(\mathbf{y})$  can be seen as an update operator of  $\mathcal{P}$  at a point  $\mathbf{y}$
- ▶  $p_L(\mathbf{y})$  forms the basis of the iterative algorithm of proximal gradient method.
- ▶  $p_L(\mathbf{y})$  itself is an optimization problem that is generally not proximable:
  - ▶  $p_L(\mathbf{y})$  may have no closed-form solution.
  - ▶  $p_L(\mathbf{y})$  may be expensive to compute
- ▶ When  $g(\mathbf{x}) = \|\mathbf{x}\|_1$  and  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$ ,  $p_L$  is the update step of ISTA (Iterative Soft Thresholding Algorithm). In this case  $p_L$  has a closed-form solution. [Details here.](#)

On solving  $p_L(\mathbf{y}) = \operatorname{prox}_g^L\left(\mathbf{y} - \frac{\nabla f(\mathbf{y})}{L}\right) := \operatorname{argmin}_{\mathbf{x}} \frac{L}{2} \left\| \mathbf{x} - \left(\mathbf{y} - \frac{\nabla f(\mathbf{y})}{L}\right) \right\|_2^2 + g(\mathbf{x})$ .

- ▶ How to solve  $p_L$ : by 1st-order optimality condition
  - ▶  $p_L$  contains non-differentiable  $g$ , so we have to use subgradient instead of gradient

▶ **Subgradient 1st-order optimality condition**

A point  $\mathbf{x}^*$  is the minimizer of  $p_L(\mathbf{y})$  if and only if (necessary & sufficient)

	$\mathbf{x}^* \in \partial p_L(\mathbf{y})$	subgradient 1st-order optimality condition
$\iff$	$\mathbf{x}^* = \partial p_L(\mathbf{y})$	strictly convex cost so solution is unique
$\iff$	$\mathbf{0} \in \partial \left\{ \frac{L}{2} \left\  \mathbf{x}^* - \left(\mathbf{y} - \frac{\nabla f(\mathbf{y})}{L}\right) \right\ _2^2 + g(\mathbf{x}^*) \right\}$	write down the subdifferential
$\iff$	$\mathbf{0} \in \partial \left\{ \frac{L}{2} \left\  \mathbf{x}^* - \left(\mathbf{y} - \frac{\nabla f(\mathbf{y})}{L}\right) \right\ _2^2 \right\} + \partial g(\mathbf{x}^*)$	$\operatorname{dom} g \subset \operatorname{int} \operatorname{dom} f$ so $\partial$ is a linear operator
$\iff$	$\mathbf{0} \in \nabla \left\{ \frac{L}{2} \left\  \mathbf{x}^* - \left(\mathbf{y} - \frac{\nabla f(\mathbf{y})}{L}\right) \right\ _2^2 \right\} + \partial g(\mathbf{x}^*)$	$\partial$ on a differentiable function reduces to $\nabla$
$\iff$	$\mathbf{0} \in L \left\{ \mathbf{x}^* - \left(\mathbf{y} - \frac{\nabla f(\mathbf{y})}{L}\right) \right\} + \partial g(\mathbf{x}^*)$	compute $\nabla$
$\iff$	$\mathbf{0} \in L(\mathbf{x}^* - \mathbf{y}) + \nabla f(\mathbf{y}) + \partial g(\mathbf{x}^*)$	

Let  $\gamma(\mathbf{y}) \in \partial g(\mathbf{y})$ , then one has  $\mathbf{x} = p_L(\mathbf{y})$  if and only if

- ▶ there exists  $\gamma(\mathbf{y}) \in \partial g(\mathbf{y})$  is the sub-differential of  $g$  s.t.  $\underbrace{\nabla f(\mathbf{y}) + L(\mathbf{x} - \mathbf{y}) + \gamma(\mathbf{y})}_{= \mathbf{0}} = \mathbf{0}$

## A key lemma

For all  $\mathbf{y}$ ,  $L \geq L_f > 0$  s.t.  $F(p_L(\mathbf{y})) \leq Q(p_L(\mathbf{y}); \mathbf{y})$ , then for all  $\mathbf{x}$  we have

$$F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2 + L \langle \mathbf{y} - \mathbf{x}, p_L(\mathbf{y}) - \mathbf{y} \rangle.$$

**Proof.** Consider for all  $\mathbf{y}$ ,  $L > L_f > 0$  s.t.  $F(p_L(\mathbf{y})) \leq Q(p_L(\mathbf{y}); \mathbf{y})$ ,

$$\begin{aligned} F(p_L(\mathbf{y})) \leq Q(p_L(\mathbf{y}); \mathbf{y}) &\iff -F(p_L(\mathbf{y})) \geq -Q(p_L(\mathbf{y}); \mathbf{y}) \\ &\iff F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq F(\mathbf{x}) - Q(p_L(\mathbf{y}); \mathbf{y}). \end{aligned} \quad (1)$$

Now we bound  $F$  using convexity, first

$$\begin{aligned} f(\mathbf{x}) &\geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle & (*) \quad & f \text{ is convex and differentiable} \\ g(\mathbf{x}) &\geq g(p_L(\mathbf{y})) + \langle \gamma(\mathbf{y}), \mathbf{x} - p_L(\mathbf{y}) \rangle & (**) \quad & g \text{ is convex and not differentiable} \\ \underbrace{f(\mathbf{x}) + g(\mathbf{x})}_{F(\mathbf{x})} &\geq f(\mathbf{y}) + g(p_L(\mathbf{y})) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \langle \gamma(\mathbf{y}), \mathbf{x} - p_L(\mathbf{y}) \rangle \end{aligned}$$

Hence (1) becomes

$$F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq f(\mathbf{y}) + g(p_L(\mathbf{y})) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \langle \gamma(\mathbf{y}), \mathbf{x} - p_L(\mathbf{y}) \rangle - Q(p_L(\mathbf{y}); \mathbf{y}) \quad (1')$$

Next we write down  $Q$ : by the definition of  $Q$ , put  $\mathbf{x} = p_L(\mathbf{y})$ , we have

$$Q(p_L(\mathbf{y}); \mathbf{y}) = f(\mathbf{y}) + \langle p_L(\mathbf{y}) - \mathbf{y}, \nabla f(\mathbf{y}) \rangle + \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2 + g(p_L(\mathbf{y}))$$

Now (1') becomes

$$F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq \underbrace{f(\mathbf{y}) + g(p_L(\mathbf{y})) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \langle \gamma(\mathbf{y}), \mathbf{x} - p_L(\mathbf{y}) \rangle}_{-f(\mathbf{y}) - \langle p_L(\mathbf{y}) - \mathbf{y}, \nabla f(\mathbf{y}) \rangle - \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2 - g(p_L(\mathbf{y}))} \quad (1'')$$

$$F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq \underbrace{f(\mathbf{y})}_{\text{red}} + \underbrace{g(p_L(\mathbf{y}))}_{\text{blue}} + \underbrace{\langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle}_{\text{green}} + \langle \gamma(\mathbf{y}), \mathbf{x} - p_L(\mathbf{y}) \rangle - \underbrace{f(\mathbf{y})}_{\text{red}} - \underbrace{\langle p_L(\mathbf{y}) - \mathbf{y}, \nabla f(\mathbf{y}) \rangle}_{\text{green}} - \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2 - \underbrace{g(p_L(\mathbf{y}))}_{\text{blue}} \quad (1'')$$

$$\iff F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq \langle \nabla f(\mathbf{y}), \mathbf{x} - p_L(\mathbf{y}) \rangle + \langle \gamma(\mathbf{y}), \mathbf{x} - p_L(\mathbf{y}) \rangle - \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2$$

$$\iff F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq \langle \underbrace{\nabla f(\mathbf{y}) + \gamma(\mathbf{y})}_{\text{wavy}} , \mathbf{x} - p_L(\mathbf{y}) \rangle - \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2$$

$$\iff F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq \langle \underbrace{L(\mathbf{y} - \mathbf{x})}_{\text{wavy}} , \mathbf{x} - p_L(\mathbf{y}) \rangle - \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2 \quad \text{1st-order optimality}$$

$$\iff F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq L \langle \mathbf{y} - p_L(\mathbf{y}), \mathbf{x} - p_L(\mathbf{y}) \rangle - \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2 \quad \text{put } \mathbf{x} = p_L(\mathbf{y})$$

Tricky step:  $(a - c)(b - c) = (b - a)(c - b) + (c - b)^2$ , then

$$L \langle \mathbf{y} - p_L(\mathbf{y}), \mathbf{x} - p_L(\mathbf{y}) \rangle = L \langle \mathbf{x} - \mathbf{y}, p_L(\mathbf{y}) - \mathbf{y} \rangle + L \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2$$

Therefore

$$F(\mathbf{x}) - F(p_L(\mathbf{y})) \leq L \langle \mathbf{x} - \mathbf{y}, p_L(\mathbf{y}) - \mathbf{y} \rangle + \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2. \quad \square$$

# ISTA and FISTA

## ISTA / Proximal gradient descent

- ▶ Iterate the following step

$$\mathbf{x}_k = p_L(\mathbf{x}_{k-1}).$$

with  $L = L_f$  known.

- ▶ If  $L_f$  is unknown, backtracking can be used.
- ▶ ISTA  $\in$  proximal gradient algorithm. [Details](#).

## FISTA / accelerated proximal gradient method

- ▶ Iterate the following step

- ▶  $\mathbf{x}_k = p_L(\mathbf{y}_{k-1})$

- ▶  $t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$

- ▶  $\mathbf{y}_{k+1} = \mathbf{x}_k + \frac{t_k - 1}{t_{k+1}}(\mathbf{x}_k - \mathbf{x}_{k-1})$

- ▶ FISTA = Nesterov's acceleration applied on ISTA
- ▶ FISTA  $\in$  accelerated proximal gradient algorithm
- ▶ See [here](#) for the convergence of FISTA

## Convergence of proximal gradient descent / ISTA

- ▶ Proximal gradient descent is *monotone*:

1.  $Q(\mathbf{x}; \mathbf{y}) = f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$   
so if  $\mathbf{x} = \mathbf{y}$  then  $Q_L(\mathbf{x}; \mathbf{y}) = F(\mathbf{x})$
2.  $p_L(\mathbf{y}) := \operatorname{argmin}_{\mathbf{x}} Q(\mathbf{x}; \mathbf{y})$  by definition
3.  $Q$  is local quadratic over-estimator of  $F$ , so  $F(\mathbf{x}) \leq Q(\mathbf{x})$

Hence we have

$$F(\mathbf{x}_k) \stackrel{3}{\leq} Q(\mathbf{x}_k, \mathbf{x}_{k-1}) \stackrel{2}{\leq} Q(\mathbf{x}_{k-1}, \mathbf{x}_{k-1}) \stackrel{1}{\leq} F(\mathbf{x}_{k-1})$$

- ▶ If  $F(\mathbf{x}) \geq 0$  for all  $\mathbf{x}$ , together with the monotonicity, the sequence of the objective function value converge.
- ▶ As  $(\mathcal{P})$  is a convex problem, it converges to global minima with optimal objective value  $F^*$ .

### Theorem (Convergence rate of ISTA using constant step size)

$$F(\mathbf{x}_k) - F^* \leq \frac{L_f R_0^2}{2k} \quad \forall k, R_0 = \|\mathbf{x}_0 - \mathbf{x}^*\|_2$$

- ▶ Here the gradient stepsize is  $\frac{1}{L}$
- ▶ As ISTA  $\in$  proximal gradient algorithm, the convergence properties of proximal gradient algorithm applies to ISTA. See [here](#) for another approach to show the  $1/k$  convergence rate.
- ▶ The proof starting in the next slide follows [Beck & Teboulle 2009](#).

$$F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2 + L \langle \mathbf{y} - \mathbf{x}, p_L(\mathbf{y}) - \mathbf{y} \rangle \text{ (key lemma)}$$

**Proof** In the key lemma, put  $\mathbf{x} = \mathbf{x}^*$  and  $\mathbf{y} = \mathbf{x}_k$ , then based on ISTA,  $p_L(\mathbf{y}) = \mathbf{x}_{k+1}$ :

$$F(\mathbf{x}^*) - F(\mathbf{x}_{k+1}) \geq \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 + L \langle \mathbf{x}_k - \mathbf{x}^*, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle,$$

Based on [observation](#), consider completing squares

$$\begin{aligned} \frac{2}{L} (F^* - F(\mathbf{x}_{k+1})) &\geq \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 + 2 \langle \mathbf{x}_k - \mathbf{x}^*, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle, \\ &= \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 + 2 \langle \mathbf{x}_k - \mathbf{x}^*, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \\ &= \|\mathbf{x}_{k+1} - \mathbf{x}_k + \mathbf{x}_k - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \\ &= \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \end{aligned}$$

Take summation from  $k = 0$  to  $k - 1$

$$\frac{2}{L} \left( kF^* - \sum_{i=0}^{k-1} F(\mathbf{x}_{i+1}) \right) \geq \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2 - \underbrace{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}_{R_0^2} \quad (2)$$

# Proof ... 2/3

$$F(\mathbf{x}) - F(p_L(\mathbf{y})) \geq \frac{L}{2} \|p_L(\mathbf{y}) - \mathbf{y}\|_2^2 + L(\mathbf{y} - \mathbf{x}, p_L(\mathbf{y}) - \mathbf{y}) \text{ (key lemma)}$$

In the key lemma, put  $\mathbf{x} = \mathbf{y} = \mathbf{x}_k$ ,  $p_L(\mathbf{y}) = \mathbf{x}_{k+1}$

$$\begin{aligned}
 & F(\mathbf{x}_k) - F(\mathbf{x}_{k+1}) \geq \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \\
 \Leftrightarrow & \frac{2}{L} (F(\mathbf{x}_k) - F(\mathbf{x}_{k+1})) \geq \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \\
 \Leftrightarrow & \frac{2}{L} k (F(\mathbf{x}_k) - F(\mathbf{x}_{k+1})) \geq k \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \quad \text{Tricky: multiply } k \\
 \Leftrightarrow & \frac{2}{L} \sum_{i=0}^{k-1} i (F(\mathbf{x}_i) - F(\mathbf{x}_{i+1})) \geq \sum_{i=0}^{k-1} i \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \quad \text{summation from } k=0 \text{ to } k-1 \\
 \Leftrightarrow & \frac{2}{L} \sum_{i=0}^{k-1} \left\{ i (F(\mathbf{x}_i) - F(\mathbf{x}_{i+1})) + F(\mathbf{x}_{i+1}) - F(\mathbf{x}_{i-1}) \right\} \geq \sum_{i=0}^{k-1} i \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \quad \text{tricky step} \\
 \Leftrightarrow & \frac{2}{L} \sum_{i=0}^{k-1} \left\{ i F(\mathbf{x}_i) - (i+1) F(\mathbf{x}_{i+1}) + F(\mathbf{x}_{i+1}) \right\} \geq \sum_{i=0}^{k-1} i \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \\
 \Leftrightarrow & \frac{2}{L} \left( \underbrace{\sum_{i=0}^{k-1} (i F(\mathbf{x}_i) - (i+1) F(\mathbf{x}_{i+1}))}_{=-k F(\mathbf{x}_k) \text{ by telescoping sum}} + \sum_{i=0}^{k-1} F(\mathbf{x}_{i+1}) \right) \geq \sum_{i=0}^{k-1} i \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2
 \end{aligned}$$

Therefore

$$\frac{2}{L} \left( -k F(\mathbf{x}_k) + \sum_{i=0}^{k-1} F(\mathbf{x}_{i+1}) \right) \geq \sum_{i=0}^{k-1} i \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \tag{3}$$

Proof ... 3/3

$$\frac{2}{L} \left( kF^* - \sum_{i=0}^{k-1} F(\mathbf{x}_{i+1}) \right) \geq \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2 - \underbrace{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}_{R_0^2} \quad (2)$$

$$\frac{2}{L} \left( -kF(\mathbf{x}_k) + \sum_{i=0}^{k-1} F(\mathbf{x}_{i+1}) \right) \geq \sum_{i=0}^{k-1} i \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \quad (3)$$

$$\frac{2}{L} \left( kF^* - \sum_{i=0}^{k-1} F(\mathbf{x}_{i+1}) \right) \geq \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2 - R_0^2 \quad (2)$$

$$\frac{2}{L} \left( -kF(\mathbf{x}_k) + \sum_{i=0}^{k-1} F(\mathbf{x}_{i+1}) \right) \geq \sum_{i=0}^{k-1} i \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \quad (3)$$

$$\frac{2}{L} (kF^* - kF(\mathbf{x}_k)) \geq \sum_{i=0}^{k-1} i \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 + \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2 - R_0^2 \quad (2) + (3)$$

$$\frac{2k}{L} (F(\mathbf{x}_k) - F^*) \leq \underbrace{R_0^2 - \sum_{i=0}^{k-1} i \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2}_{\geq 0} - \underbrace{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2}_{\geq 0} \quad \text{multiply } -1$$

$$\leq R_0^2$$

$$F(\mathbf{x}_k) - F^* \leq \frac{LR_0^2}{2k} \quad \square$$

End of document