

# ECE 8201: Low-dimensional Signal Models for High-dimensional Data Analysis

Lecture 5: FISTA

Yuejie Chi  
The Ohio State University



THE OHIO STATE UNIVERSITY

---

## Reference

---

- Beck, A., & Teboulle, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1), 183-202.

See also:

- Nesterov, Y. (2007). Gradient methods for minimizing composite objective function.
- Lecture notes by L. Vandenberghe.  
<http://www.seas.ucla.edu/~vandenbe/236C/lectures/fgrad.pdf>.

# How to solve composite optimization problems?

---

General composite optimization problem:

$$\text{(COP)} : \quad \hat{\boldsymbol{x}} = \underset{\boldsymbol{x}}{\operatorname{argmin}} \{F(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x})\}$$

- $f(\boldsymbol{x})$  is convex and differentiable,
- $g(\boldsymbol{x})$  is convex, possibly non-differentiable

Examples:

- LASSO:  $f(\boldsymbol{x}) = \frac{1}{2}\|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2$ , and  $g(\boldsymbol{x}) = \lambda\|\boldsymbol{x}\|_1$ . (focus of this lecture)
- Nuclear norm minimization (later for matrix completion):

$$f(\boldsymbol{X}) = \|\mathcal{P}_\Omega(\boldsymbol{Y} - \boldsymbol{X})\|_F^2, \quad g(\boldsymbol{X}) = \lambda\|\boldsymbol{X}\|_*$$

where  $\|\boldsymbol{X}\|_* = \sum_{i=1}^{\min(m,n)} \sigma_i(\boldsymbol{X})$ , the sum of the singular values of  $\boldsymbol{X} \in \mathbb{R}^{m \times n}$ .

# Motivation

---

Standard methods (e.g. subgradient methods) for solving COP has very slow convergence rate (need  $O(1/\epsilon^2)$  iterations to reach  $\epsilon$  accuracy).

We would discuss an algorithm called FISTA that

- is iterative, and has low computational cost (first-order algorithm, which requires computation of a single gradient per iteration);
- has quadratic convergence rate;
- performs well in practice and works for a large class of problems.

FISTA stands for **F**ast **I**terative **S**hrinkage-**T**hresholding **A**lgorithm.

# Gradient descent

---

Consider the unconstrained minimization of a continuously differentiable function  $f(\mathbf{x})$  as

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x})$$

using gradient descent: start with an initialization  $\mathbf{x}_0 \in \mathbb{R}^n$ , and iterate

$$\mathbf{x}_k = \mathbf{x}_{k-1} - t_k \nabla f(\mathbf{x}_{k-1})$$

where  $t_k$  is a suitable step-size at step  $k$ .

**Key observation:** we can view the gradient descent step as solving a *proximal regularization* of the linearized function  $f$  at  $\mathbf{x}_{k-1}$ ,

$$\mathbf{x}_k = \underset{\mathbf{x}}{\operatorname{argmin}} \left\{ f(\mathbf{x}_{k-1}) + \langle \mathbf{x} - \mathbf{x}_{k-1}, \nabla f(\mathbf{x}_{k-1}) \rangle + \frac{1}{2t_k} \|\mathbf{x} - \mathbf{x}_{k-1}\|_2^2 \right\}.$$

## Generalized gradient descent

---

In the COP,

$$\hat{\boldsymbol{x}} = \operatorname{argmin}_{\boldsymbol{x}} f(\boldsymbol{x}) + g(\boldsymbol{x})$$

we would like to generalize the proximal regularization idea, by extending the update rule as

$$\boldsymbol{x}_k = \operatorname{argmin}_{\boldsymbol{x}} \left\{ f(\boldsymbol{x}_{k-1}) + \langle \boldsymbol{x} - \boldsymbol{x}_{k-1}, \nabla f(\boldsymbol{x}_{k-1}) \rangle + \frac{1}{2t_k} \|\boldsymbol{x} - \boldsymbol{x}_{k-1}\|_2^2 + g(\boldsymbol{x}) \right\}.$$

This can be simplified (by ignoring constant terms) as

$$\boldsymbol{x}_k = \operatorname{argmin}_{\boldsymbol{x}} \left\{ \frac{1}{2t_k} \|\boldsymbol{x} - (\boldsymbol{x}_{k-1} - t_k \nabla f(\boldsymbol{x}_{k-1}))\|_2^2 + g(\boldsymbol{x}) \right\} \quad (*)$$

# Proximal mapping

---

**Definition 1.** *The proximal mapping (operator) of a convex function  $g(\mathbf{x})$  is written as*

$$\text{prox}_g(\mathbf{x}) = \underset{\mathbf{u}}{\text{argmin}} \left\{ \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|_2^2 + g(\mathbf{u}) \right\}.$$

- $g(\mathbf{x}) = 0$ :  $\text{prox}_g(\mathbf{x}) = \mathbf{x}$ .
- $g(\mathbf{x}) = I_C(\mathbf{x})$  is an indicator function of a convex set  $C$ , then

$$\text{prox}_g(\mathbf{x}) = \underset{\mathbf{u} \in C}{\text{argmin}} \|\mathbf{u} - \mathbf{x}\|_2^2$$

- $g(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$ :  $\text{prox}_g(\mathbf{x})$  is the shrinkage (soft-thresholding) operator and can be decomposed entry-wise:

$$\text{prox}_g(x_i) := \mathcal{T}_\lambda(x_i) = \begin{cases} x_i - \lambda, & x_i \geq \lambda \\ 0, & |x_i| < \lambda \\ x_i + \lambda, & x_i \leq -\lambda \end{cases}$$

# Generalized gradient descent and ISTA

---

- The generalized gradient descent (\*) can be regarded as a proximal mapping:

$$\begin{aligned}\mathbf{x}_k &= \operatorname{argmin}_{\mathbf{x}} \left\{ \frac{1}{2t_k} \|\mathbf{x} - (\mathbf{x}_{k-1} - t_k \nabla f(\mathbf{x}_{k-1}))\|_2^2 + g(\mathbf{x}) \right\} \\ &= \operatorname{prox}_{t_k g}(\mathbf{x}_{k-1} - t_k \nabla f(\mathbf{x}_{k-1}))\end{aligned}$$

- When  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ , and  $g(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$ , this gives the update rule for ISTA (Iterative Shrinkage-Thresholding Algorithm), or proximal gradient descent:

$$\begin{aligned}\mathbf{x}_k &= \operatorname{prox}_{\lambda t_k \|\mathbf{x}\|_1}(\mathbf{x}_{k-1} - t_k \nabla f(\mathbf{x}_{k-1})) \\ &= \operatorname{prox}_{\lambda t_k \|\mathbf{x}\|_1}(\mathbf{x}_{k-1} - t_k \nabla f(\mathbf{x}_{k-1})) \\ &= \mathcal{T}_{\lambda t_k}(\mathbf{x}_{k-1} - t_k \nabla f(\mathbf{x}_{k-1}))\end{aligned}$$

where  $\nabla f(\mathbf{x}_{k-1}) = \mathbf{A}^\top(\mathbf{A}\mathbf{x}_{k-1} - \mathbf{y})$ . This can be efficiently computed.

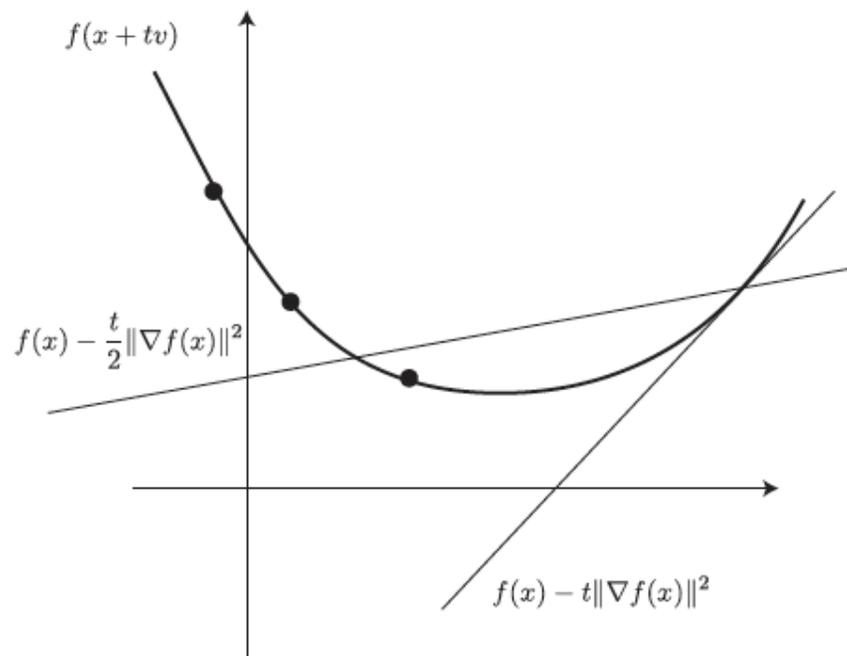
## Choice of step size

---

- Constant step-size:  $t_k = t$
- Backtracking line search: start with  $t_0$  and do  $t = \beta t$  until

$$f(\mathbf{x} - t\nabla f(\mathbf{x})) \leq f(\mathbf{x}) - \alpha t \|\nabla f(\mathbf{x})\|_2^2$$

with  $0 < \alpha, \beta < 1$ , e.g.  $\alpha = 1/2$ .



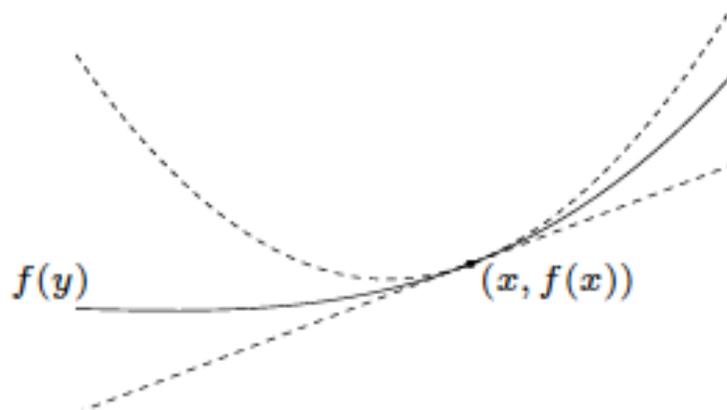
# Assumptions

---

- $g : \mathbb{R}^n \mapsto \mathbb{R}$  is a continuous convex function, possibly nonsmooth;
- $f : \mathbb{R}^n \mapsto \mathbb{R}$  is a smooth convex function that is continuously differentiable with *Lipschitz constant*:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L_f \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

Example: For LASSO problems, we have  $L_f = \sigma_{\max}(\mathbf{A}^\top \mathbf{A})$ .



- The optimal value of  $F = f + g$  is  $F^*$  with optimal solution  $\mathbf{x}^*$ .

# Convergence of ISTA

---

**Theorem 1. [Convergence for generalized gradient descent]** *Fix step size*  
 $t_k = t \leq 1/L,$

$$F(\mathbf{x}_k) - F^* \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2tk}$$

*Similar results hold with backtracking for step size.*

- Similar to the convergence of gradient descent
- The best possible is  $O(1/k^2)$  for first-order methods – can we achieve it?

The answer is yes, with minimal additional computational cost.

## Accelerated Gradient Descent

---

ISTA reaches an accuracy within  $O(1/k)$  after  $k$  steps; this is not optimal (which is  $O(1/k^2)$ ). The methods of Nesterov meet the optimal bound with the same computational cost (one gradient computation per iteration).

- We will first examine Nesterov's acceleration method (1983) for smooth convex functions;
- We then extend it to optimizing composite functions, using FISTA (Beck and Teboulle, 2009), which extends Nesterov's method.

## Nesterov's ACG for convex smooth function

---

Consider minimizing a convex smooth function  $f(\mathbf{x})$  with Lipschitz constant  $L$ :

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x})$$

Nesterov's Accelerated Gradient Descent performs attains a rate of  $O(1/k^2)$ . It proceeds as below:

- Start with an initialization  $\mathbf{x}_0 = \mathbf{x}_{-1}$ ,  $\theta_0 = 0$ ;
- for  $k = 1, 2, \dots$ ,

$$\theta_k = \frac{1 + \sqrt{1 + 4\theta_{k-1}^2}}{2},$$

$$\mathbf{y}_k = \mathbf{x}_{k-1} + \left( \frac{\theta_{k-1} - 1}{\theta_k} \right) (\mathbf{x}_{k-1} - \mathbf{x}_{k-2})$$

$$\mathbf{x}_k = \mathbf{y}_k - t_k \nabla f(\mathbf{y}_k)$$

Remark: other choice of the momentum term with  $\theta_k = \frac{k+1}{2}$ :

$$\mathbf{y}_k = \mathbf{x}_{k-1} + \frac{k-2}{k+1}(\mathbf{x}_{k-1} - \mathbf{x}_{k-2})$$

**Theorem 2. [Nesterov 1983]** *The Nesterov's AGD satisfies*

$$f(\mathbf{y}_k) - f(\mathbf{x}^*) \leq \frac{2\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{Lk^2}$$

Achieves the optimal rate!

# FISTA

---

The FISTA algorithm with step size  $t_k$  (e.g.  $t_k = \frac{1}{L}$ , where  $L_f$  is the Lipschitz constant of  $f$ ):

- Initialization:  $\mathbf{x}_0 = \mathbf{x}_{-1} \in \mathbb{R}^n$ ,  $\theta_0 = 1$ ,
- For  $k = 1, 2, \dots$ ,

$$\theta_k = \frac{1 + \sqrt{1 + 4\theta_{k-1}^2}}{2}$$

$$\mathbf{y}_k = \mathbf{x}_{k-1} + \left( \frac{\theta_{k-1} - 1}{\theta_k} \right) (\mathbf{x}_{k-1} - \mathbf{x}_{k-2})$$

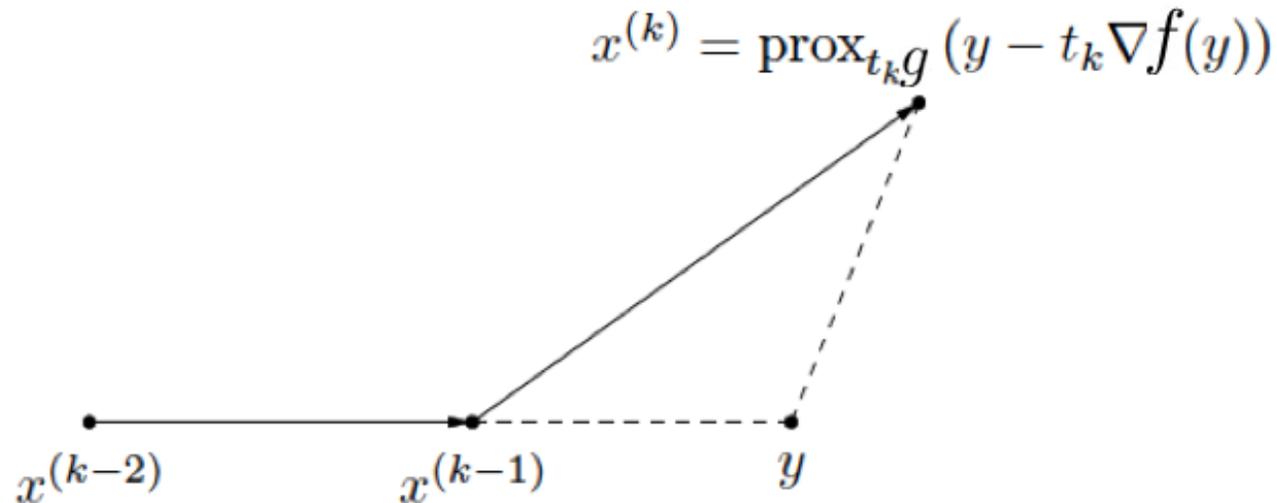
$$\mathbf{x}_k = \text{prox}_{t_k g} (\mathbf{y}_k - t_k \nabla f(\mathbf{y}_k))$$

FISTA is computationally efficient when the proximal operator can be computed efficiently (e.g. LASSO).

# Interpretation

---

- first iteration is a proximal gradient step at  $y_1 = x_0$
- next iterations are proximal gradient steps at extrapolated points  $y_k$ ,  $k \geq 2$ , with the linear combinations carefully chosen.



## Case Study: LASSO

---

For LASSO: set  $\mathbf{y}_1 = \mathbf{x}_0 \in \mathbb{R}^n$ ,  $\theta_1 = 0$ , and  $t_k = 1/\sigma_{\max}(\mathbf{A}^\top \mathbf{A})$  (constant step-size), iterate

$$\theta_k = \frac{1 + \sqrt{1 + 4\theta_{k-1}^2}}{2}$$
$$\mathbf{y}_k = \mathbf{x}_{k-1} + \left( \frac{\theta_{k-1} - 1}{\theta_k} \right) (\mathbf{x}_{k-1} - \mathbf{x}_{k-2})$$
$$\mathbf{x}_k = \mathcal{T}_{\lambda t_k} (\mathbf{y}_k - t_k \mathbf{A}^\top (\mathbf{A} \mathbf{y}_k - \mathbf{y}))$$

The main computation cost to apply  $\mathbf{A}$  and  $\mathbf{A}^\top$ ; no matrix inversion is needed.

# Convergence of FISTA

## Theorem 3.

$$F(\mathbf{x}_k) - F(\mathbf{x}^*) \leq \frac{2L\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{(k+1)^2} \sim O\left(\frac{1}{k^2}\right)$$

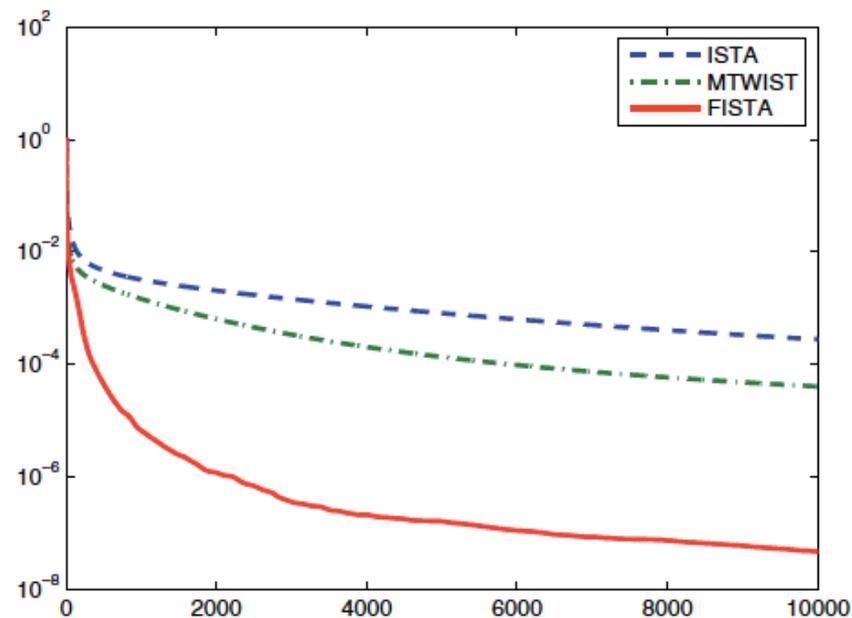


Figure 5. Comparison of function value errors  $F(\mathbf{x}_k) - F(\mathbf{x}^*)$  of ISTA, MTWIST, and FISTA.

## Proof of Theorem 3

---

- Introduce another sequence  $\mathbf{v}_k$ , which satisfies

$$\mathbf{v}_k := \mathbf{x}_{k-1} + \theta_k(\mathbf{x}_k - \mathbf{x}_{k-1})$$

$$\mathbf{y}_k = \frac{1}{\theta_k}\mathbf{v}_{k-1} + \left(1 - \frac{1}{\theta_k}\right)\mathbf{x}_{k-1}$$

- Two useful facts:

1.  $\mathbf{v}_k = \mathbf{v}_{k-1} + \theta_k(\mathbf{x}_k - \mathbf{y}_k)$

2.  $\left(1 - \frac{1}{\theta_k}\right)\theta_k^2 = \theta_{k-1}^2$

## Important inequalities

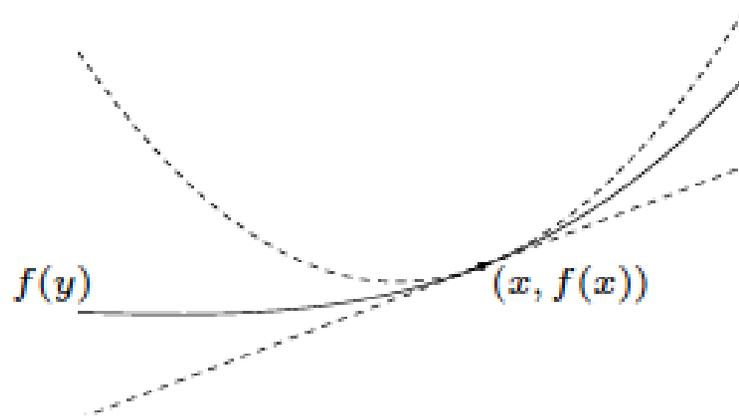
---

Upper bound of  $f$  from Lipschitz property:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L_f \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

we have

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L_f}{2} \|\mathbf{y} - \mathbf{x}\|_2^2, \quad \forall \mathbf{x}, \mathbf{y}$$



## Important inequalities

---

Upper bound of  $g$  from definition of proximal operator:

$$g(\mathbf{y}) \leq g(\mathbf{z}) + \frac{1}{t}(\mathbf{w} - \mathbf{y})^\top(\mathbf{y} - \mathbf{z}), \quad \forall \mathbf{w}, \mathbf{z}, \mathbf{y} = \text{prox}_{tg}(\mathbf{w})$$

Proof: since  $\mathbf{y} = \text{prox}_{tg}(\mathbf{w})$  minimizes  $tg(\mathbf{u}) + \frac{1}{2}\|\mathbf{w} - \mathbf{u}\|_2^2$  by definition, we have

$$0 \in t\partial g(\mathbf{y}) + (\mathbf{y} - \mathbf{w})$$

i.e.

$$\frac{1}{t}(\mathbf{w} - \mathbf{y}) \in \partial g(\mathbf{y}),$$

By the definition of subgradient we have  $\forall \mathbf{z}$ ,

$$g(\mathbf{z}) \geq g(\mathbf{y}) + \frac{1}{t}(\mathbf{w} - \mathbf{y})^\top(\mathbf{z} - \mathbf{y})$$

## Progress in one iteration

---

Define  $\mathbf{x}^+ = \mathbf{x}_k$ ,  $\mathbf{x} = \mathbf{x}_{k-1}$ ,  $\mathbf{y} = \mathbf{y}_k$ ,  $\theta = \theta_k$ ,  $\mathbf{v} = \mathbf{v}_{k-1}$ ,  $\mathbf{v}^+ = \mathbf{v}_k$ ,

- upper bound from Lipschitz property: if  $0 < t \leq 1/L$ ,

$$f(\mathbf{x}^+) \leq f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x}^+ - \mathbf{y}) + \frac{1}{2t} \|\mathbf{y} - \mathbf{x}^+\|_2^2$$

- upper bound from the definition of prox-operator ( $\mathbf{x}^+ = \text{prox}_{tg}(\mathbf{y} - t\nabla f(\mathbf{y}))$ ):

$$g(\mathbf{x}^+) \leq g(\mathbf{z}) + \nabla f(\mathbf{y})^\top (\mathbf{z} - \mathbf{x}^+) + \frac{1}{t} (\mathbf{x}^+ - \mathbf{y})^\top (\mathbf{z} - \mathbf{x}^+), \quad \forall \mathbf{z}$$

- add the upper bounds and use convexity of  $f$ :

$$F(\mathbf{x}^+) \leq F(\mathbf{z}) + \frac{1}{t} (\mathbf{x}^+ - \mathbf{y})^\top (\mathbf{z} - \mathbf{x}^+) + \frac{1}{2t} \|\mathbf{y} - \mathbf{x}^+\|_2^2, \quad \forall \mathbf{z}$$

- make convex combination of upper bounds for  $z = \mathbf{x}$  and  $z = \mathbf{x}^*$ :

$$\begin{aligned}
F(\mathbf{x}^+) - F^* - \left(1 - \frac{1}{\theta}\right) (F(\mathbf{x}) - F^*) &= F(\mathbf{x}^+) - \frac{1}{\theta} F^* - \left(1 - \frac{1}{\theta}\right) F(\mathbf{x}) \\
&\leq \frac{1}{t} (\mathbf{x}^+ - \mathbf{y})^\top \left( \frac{1}{\theta} \mathbf{x}^* + \left(1 - \frac{1}{\theta}\right) \mathbf{x} - \mathbf{x}^+ \right) + \frac{1}{2t} \|\mathbf{y} - \mathbf{x}^+\|_2^2 \\
&= \frac{1}{2t} \left( \|\mathbf{y} - \frac{1}{\theta} \mathbf{x}^* + \left(1 - \frac{1}{\theta}\right) \mathbf{x}\|_2^2 - \|\mathbf{x}^+ - \frac{1}{\theta} \mathbf{x}^* + \left(1 - \frac{1}{\theta}\right) \mathbf{x}\|_2^2 \right) \\
&= \frac{1}{2\theta^2 t} \left( \|\mathbf{v} - \mathbf{x}^*\|_2^2 - \|\mathbf{v}^+ - \mathbf{x}^*\|_2^2 \right)
\end{aligned}$$

We now have, at the  $k$ th iteration:

$$\begin{aligned}
\theta_k^2 t (F(\mathbf{x}_k) - F^*) + \frac{1}{2} \|\mathbf{v}_k - \mathbf{x}^*\|_2^2 &\leq (\theta_k^2 - \theta_{k-1}^2) t (F(\mathbf{x}_{k-1}) - F^*) + \frac{1}{2} \|\mathbf{v}_{k-1} - \mathbf{x}^*\|_2^2 \\
&= \theta_{k-1}^2 t (F(\mathbf{x}_{k-1}) - F^*) + \frac{1}{2} \|\mathbf{v}_{k-1} - \mathbf{x}^*\|_2^2
\end{aligned}$$

Applying the above relationship recursively, we obtain

$$\begin{aligned}\theta_k^2 t (F(\mathbf{x}_k) - F^*) + \frac{1}{2} \|\mathbf{v}_k - \mathbf{x}^*\|_2^2 &\leq \theta_0^2 t (F(\mathbf{x}_0) - F^*) + \frac{1}{2} \|\mathbf{v}_0 - \mathbf{x}^*\|_2^2 \\ &= \frac{1}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2\end{aligned}$$

therefore, plug in  $t = \frac{1}{L}$ ,

$$F(\mathbf{x}_k) - F^* \leq \frac{1}{2\theta_k^2 t} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \leq \frac{2L}{(k+1)^2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$

## Alternative formulation

---

Alternative formulation:

- Initialization:  $\mathbf{y}_1 = \mathbf{x}_0 \in \mathbb{R}^n$ , and  $L_f$  is the Lipschitz constant;
- Fix step size  $t_k = \frac{1}{L}$ .
- For  $k = 1, 2, \dots$ ,

$$\begin{aligned}\mathbf{x}_k &= \text{prox}_{t_k g}(\mathbf{y}_k - t_k \nabla f(\mathbf{y}_k)) \\ \mathbf{y}_{k+1} &= \mathbf{x}_k + \left(\frac{k-2}{k+1}\right) (\mathbf{x}_k - \mathbf{x}_{k-1})\end{aligned}$$

Convergence speed  $O(1/k^2)$  in  $k$  steps.

## Computational-Statistical Trade-off

---

If there is indeed a ground truth  $\boldsymbol{x}^*$  and we wish  $\hat{\boldsymbol{x}}$  is close to  $\boldsymbol{x}^*$ ; we have a sequence of  $\{\boldsymbol{x}_k\}$  and hope  $\boldsymbol{x}_k$  converges to  $\hat{\boldsymbol{x}}$ . At a fixed  $k$ , we may bound

$$\|\boldsymbol{x}_k - \boldsymbol{x}^*\|_2 \leq \underbrace{\|\boldsymbol{x}_k - \hat{\boldsymbol{x}}\|_2}_{\text{computational error}} + \underbrace{\|\hat{\boldsymbol{x}} - \boldsymbol{x}^*\|_2}_{\text{statistical error}}$$

Active research in studying the computational-statistical trade-offs in statistical estimation.