



# Lecture 5. Gradient Descent Method II

Advanced Optimization (Fall 2022)

**Peng Zhao**

[zhaop@lamda.nju.edu.cn](mailto:zhaop@lamda.nju.edu.cn)

Nanjing University

# Outline

- GD for Smooth Optimization
  - Smooth and Convex Functions
  - Smooth and Strongly Convex Functions
- Nesterov's Accelerated GD
- Extension to Composite Optimization

# Part 1. GD for Smooth Optimization

- Smooth and Convex
- Smooth and Strongly Convex
- Extension to Constrained Case

# Overview

Table 1: A summary of convergence rates of GD for different function families, where we use  $\kappa \triangleq L/\sigma$  to denote the condition number.

Function Family		Step Size	Output Sequence	Convergence Rate	
G-Lipschitz	convex	$\eta = \frac{D}{G\sqrt{T}}$	$\bar{\mathbf{x}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$	$\mathcal{O}(1/\sqrt{T})$	<i>last lecture</i>
	$\sigma$ -strongly convex	$\eta_t = \frac{2}{\sigma(t+1)}$	$\bar{\mathbf{x}}_T = \sum_{t=1}^T \frac{2t}{T(T+1)} \mathbf{x}_t$	$\mathcal{O}(1/T)$	
$L$ -smooth	convex	$\eta = \frac{1}{L}$	$\bar{\mathbf{x}}_T = \mathbf{x}_T$	$\mathcal{O}(1/T)$	<i>this lecture</i>
	$\sigma$ -strongly convex	$\eta = \frac{2}{\sigma+L}$	$\bar{\mathbf{x}}_T = \mathbf{x}_T$	$\mathcal{O}\left(\exp\left(-\frac{T}{\kappa}\right)\right)$	

For simplicity, we mostly focus on *unconstrained* domain, i.e.,  $\mathcal{X} = \mathbb{R}^d$ .

# Convex and Smooth

**Theorem 1.** Suppose the function  $f : \mathbb{R}^d \mapsto \mathbb{R}$  is convex and differentiable, and also  $L$ -smooth. GD updates by  $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)$  with step size  $\eta_t = \frac{1}{L}$ , and then GD enjoys the following convergence guarantee:

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{2L \|\mathbf{x}_1 - \mathbf{x}^*\|^2}{T - 1} = \mathcal{O}\left(\frac{1}{T}\right).$$

**Note:** we are working on *unconstrained* setting and using a *fixed* step size tuning.

# The First Gradient Descent Lemma

**Lemma 1.** Suppose that  $f$  is proper, closed and convex; the feasible domain  $\mathcal{X}$  is nonempty, closed and convex. Let  $\{\mathbf{x}_t\}_{t=1}^T$  be the sequence generated by the gradient descent method,  $\mathcal{X}^*$  be the optimal set of the optimization problem and  $f^*$  be the optimal value. Then for any  $\mathbf{x}^* \in \mathcal{X}^*$  and  $t \geq 0$ ,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\eta_t(f(\mathbf{x}_t) - f^*) + \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2.$$

**Proof:**

$$\begin{aligned} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 &= \|\Pi_{\mathcal{X}}[\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)] - \mathbf{x}^*\|^2 \quad (\text{GD}) \\ &\leq \|\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) - \mathbf{x}^*\|^2 \quad (\text{Pythagoras Theorem}) \\ &= \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\eta_t \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle + \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2 \\ &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\eta_t(f(\mathbf{x}_t) - f^*) + \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2 \\ &\quad (\text{convexity: } f(\mathbf{x}_t) - f^* = f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle) \quad \square \end{aligned}$$

# Refined Result for Smooth Optimization

*Proof:*  $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 = \|\Pi_{\mathcal{X}}[\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)] - \mathbf{x}^*\|^2$  (GD)

$$\leq \|\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) - \mathbf{x}^*\|^2 \text{ (Pythagoras Theorem)}$$
$$= \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\eta_t \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle + \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2$$
$$\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\eta_t (f(\mathbf{x}_t) - f^*) + \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2$$

(convexity:  $f(\mathbf{x}_t) - f^* = f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle$ )

*haven't used smoothness*

**Lemma 2** (co-coercivity). Let  $f$  be convex and  $L$ -smooth over  $\mathbb{R}^d$ . Then for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , one has

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{1}{L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2$$

# Co-coercive Operator

**Lemma 2** (co-coercivity). Let  $f$  be convex and  *$L$ -smooth* over  $\mathbb{R}^d$ . Then for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , one has

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{1}{L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2$$

**Definition 1** (co-coercive operator). An operator  $C$  is called  $\beta$ -co-coercive (or  $\beta$ -inverse-strongly monotone, for  $\beta > 0$ , if for any  $x, y \in \mathcal{H}$ ,

$$\langle Cx - Cy, x - y \rangle \geq \beta \|Cx - Cy\|^2.$$

The co-coercive condition is relatively standard in *operator splitting* literature and *variational inequalities*.

# Smooth and Convex

*Proof:*  $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 = \|\Pi_{\mathcal{X}}[\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)] - \mathbf{x}^*\|^2$  (GD)

$$\leq \|\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) - \mathbf{x}^*\|^2 \text{ (Pythagoras Theorem)}$$
$$= \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\eta_t \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle + \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2$$
$$\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \left( \eta_t^2 - \frac{2\eta_t}{L} \right) \|\nabla f(\mathbf{x}_t)\|^2$$

*exploiting coercivity of smoothness and unconstrained first-order optimality*

$$\langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle = \langle \nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}^*), \mathbf{x}_t - \mathbf{x}^* \rangle \geq \frac{1}{L} \|\nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}^*)\|^2 = \frac{1}{L} \|\nabla f(\mathbf{x}_t)\|^2$$

$$\begin{aligned} \Rightarrow \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \left( \eta_t^2 - \frac{2\eta_t}{L} \right) \|\nabla f(\mathbf{x}_t)\|^2 \\ &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 && \text{by picking } \eta_t = \eta = \frac{1}{L} \text{ to minimize the r.h.s} \\ &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 \leq \dots \leq \|\mathbf{x}_1 - \mathbf{x}^*\|^2 && \text{which already implies the convergence} \end{aligned}$$

# Smooth and Convex

*Proof:* Now, we consider the function-value level,

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) = f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) + f(\mathbf{x}_t) - f(\mathbf{x}^*)$$

$$\begin{aligned} & f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \\ &= f(\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)) - f(\mathbf{x}_t) \\ &\leq \langle \nabla f(\mathbf{x}_t), -\eta_t \nabla f(\mathbf{x}_t) \rangle + \frac{L}{2} \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2 \quad (\text{smoothness}) \\ &= \left( -\eta_t + \frac{L}{2} \eta_t^2 \right) \|\nabla f(\mathbf{x}_t)\|^2 \\ &= -\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \quad (\text{recall that we have picked } \eta_t = \eta = \frac{1}{L}) \end{aligned}$$

*one-step  
improvement*

$$\Rightarrow f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 + f(\mathbf{x}_t) - f(\mathbf{x}^*)$$

# Smooth and Convex

*Proof:*

$$\Rightarrow f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 + f(\mathbf{x}_t) - f(\mathbf{x}^*)$$

*Next step:* relating  $\|\nabla f(\mathbf{x}_t)\|$  to function-value gap to form a telescoping structure.

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle \leq \|\nabla f(\mathbf{x}_t)\| \|\mathbf{x}_t - \mathbf{x}^*\| \Rightarrow \|\nabla f(\mathbf{x}_t)\|^2 \geq \frac{(f(\mathbf{x}_t) - f(\mathbf{x}^*))^2}{\|\mathbf{x}_t - \mathbf{x}^*\|^2}$$

$$\begin{aligned} \Rightarrow f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) &\leq -\frac{(f(\mathbf{x}_t) - f(\mathbf{x}^*))^2}{2L \|\mathbf{x}_t - \mathbf{x}^*\|^2} + f(\mathbf{x}_t) - f(\mathbf{x}^*) \\ &\leq -\frac{(f(\mathbf{x}_t) - f(\mathbf{x}^*))^2}{2L \|\mathbf{x}_1 - \mathbf{x}^*\|^2} + f(\mathbf{x}_t) - f(\mathbf{x}^*) \end{aligned}$$

(by optimizer's convergence, i.e.,  $\|\mathbf{x}_t - \mathbf{x}^*\| \leq \|\mathbf{x}_1 - \mathbf{x}^*\|$ )

# Smooth and Convex

*Proof:*  $f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq -\frac{1}{2L\|\mathbf{x}_1 - \mathbf{x}^*\|^2} (f(\mathbf{x}_t) - f(\mathbf{x}^*))^2 + f(\mathbf{x}_t) - f(\mathbf{x}^*)$

Define  $\Delta_t \triangleq f(\mathbf{x}_t) - f(\mathbf{x}^*)$  and  $\beta_t \triangleq \frac{1}{2L\|\mathbf{x}_1 - \mathbf{x}^*\|^2}$ .

$$\Rightarrow \Delta_{t+1} \leq \Delta_t - \beta_t \Delta_t^2 \quad \Rightarrow \frac{\beta_t \Delta_t}{\Delta_{t+1}} \leq \frac{1}{\Delta_{t+1}} - \frac{1}{\Delta_t}$$

$$\Rightarrow \beta_t \leq \frac{1}{\Delta_{t+1}} - \frac{1}{\Delta_t} \quad \Rightarrow \sum_{t=1}^{T-1} \beta_t \leq \frac{1}{\Delta_T} - \frac{1}{\Delta_1} \leq \frac{1}{\Delta_T}$$

$$\Rightarrow \Delta_T \triangleq f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{1}{\sum_{t=1}^{T-1} \beta_t} = \frac{2L \|\mathbf{x}_1 - \mathbf{x}^*\|^2}{T-1}.$$

□

# Key Lemma for Smooth GD

- During the proof, we have obtained an important lemma for smooth optimization.

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq \left( -\eta_t + \frac{L}{2} \eta_t^2 \right) \|\nabla f(\mathbf{x}_t)\|^2 \quad \textit{one-step improvement}$$

Compare a similar result that holds for convex and Lipschitz functions.

**Lemma 2.** *Under the same assumptions as Theorem 1. Let  $\{\mathbf{x}_t\}_{t=1}^T$  be the sequence generated by GD. Then we have*

$$\sum_{t=1}^T \eta_t (f(\mathbf{x}_t) - f^*) \leq \frac{1}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 + \frac{1}{2} \sum_{t=1}^T \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2.$$

*average-iterated* convergence vs *last-iterated* convergence

# Key Lemma for Smooth GD

- One-step improvement for *smooth* GD under *unconstrained* setting.

**Lemma 3** (one-step improvement). *Suppose the function  $f : \mathbb{R}^d \mapsto \mathbb{R}$  is convex and differentiable, and also  $L$ -smooth. Consider the following unconstrained GD update:  $\mathbf{x}' = \mathbf{x} - \eta \nabla f(\mathbf{x})$ . Then,*

$$f(\mathbf{x}') - f(\mathbf{x}) \leq \left( -\eta + \frac{L}{2} \eta^2 \right) \|\nabla f(\mathbf{x})\|^2.$$

*In particular, when choosing  $\eta = \frac{1}{L}$ , we have*

$$f\left(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})\right) - f(\mathbf{x}) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x})\|^2.$$

# Smooth and Strongly Convex

- Recall the definition of strongly convex functions (*first-order* version).

**Definition 5** (Strong Convexity). A function  $f$  is  $\sigma$ -strongly convex if, for any  $\mathbf{x} \in \text{dom}(\partial f)$ ,  $\mathbf{y} \in \text{dom}(f)$  and  $\mathbf{g} \in \partial f(\mathbf{x})$ ,

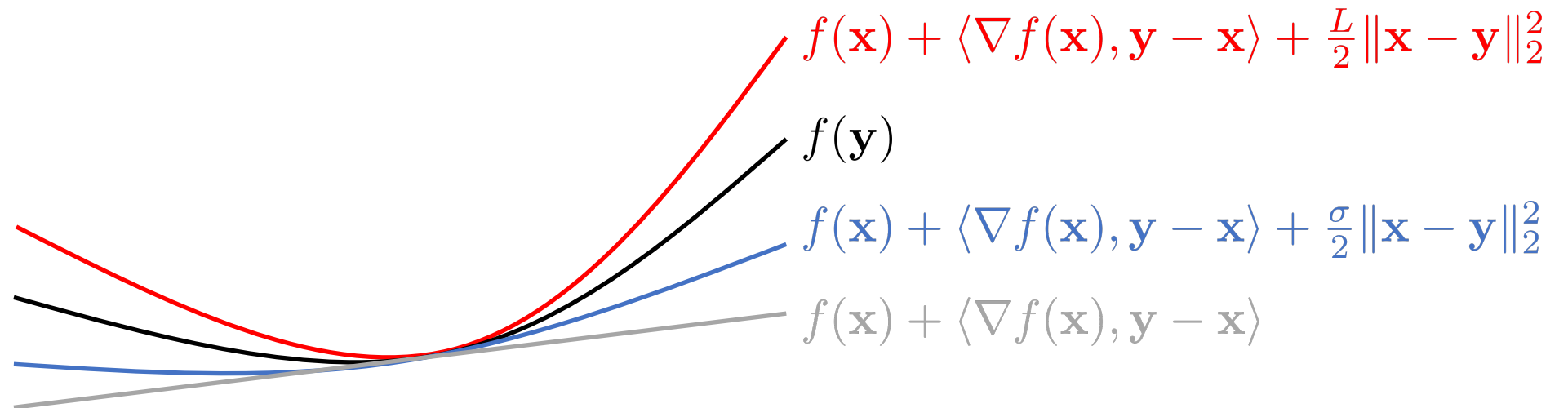
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle + \frac{\sigma}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

# Smooth and Strongly Convex

$f$  is  $\sigma$ -strongly convex

$f$  is  $L$ -smooth

$$f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\sigma}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \leq f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$



# Smooth and Strongly Convex

**Theorem 2.** Suppose the function  $f : \mathbb{R}^d \mapsto \mathbb{R}$  is  $\sigma$ -strongly-convex and differentiable, and also  $L$ -smooth; and the feasible domain  $\mathcal{X} \subseteq \mathbb{R}^d$  is compact and convex with a diameter  $D > 0$ . Then, setting  $\eta_t = \frac{2}{\sigma+L}$ , GD satisfies

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2} \exp\left(-\frac{4(T-1)}{\kappa+1}\right) \|\mathbf{x}_1 - \mathbf{x}^*\|^2 = \mathcal{O}\left(\exp\left(-\frac{T}{\kappa}\right)\right),$$

where  $\kappa \triangleq L/\sigma$  denotes the condition number of  $f$ .

**Note:** we are working on *unconstrained* setting and using a *fixed* step size tuning.

# Smooth and Strongly Convex

*Proof:*  $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 = \|\Pi_{\mathcal{X}}[\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)] - \mathbf{x}^*\|^2$  (GD)

$$\leq \|\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) - \mathbf{x}^*\|^2 \text{ (Pythagoras Theorem)}$$
$$= \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\eta_t \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle + \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2$$

*how to exploiting the **strong convexity** and **smoothness** simultaneously*

**Lemma 4** (co-coercivity of smooth and strongly convex function). *Let  $f$  be  $L$ -smooth and  $\sigma$ -strongly convex on  $\mathbb{R}^d$ . Then for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , one has*

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{\sigma L}{\sigma + L} \|\mathbf{x} - \mathbf{y}\|^2 + \frac{1}{\sigma + L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2.$$

# Coercivity of Smooth and Strongly Convex Function

**Lemma 4** (co-coercivity of smooth and strongly convex function). *Let  $f$  be  $L$ -smooth and  $\sigma$ -strongly convex on  $\mathbb{R}^d$ . Then for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , one has*

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{\sigma L}{\sigma + L} \|\mathbf{x} - \mathbf{y}\|^2 + \frac{1}{\sigma + L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2.$$

**Proof:** Define  $h(\mathbf{x}) \triangleq f(\mathbf{x}) - \frac{\sigma}{2} \|\mathbf{x}\|^2$ . Then,  $h$  enjoys the following properties:

- $h$  is convex: by  $\sigma$ -strong convexity (see previous lecture).
- $h$  is  $(L - \sigma)$ -smooth.  $\nabla^2 h(\mathbf{x}) = \nabla^2 f(\mathbf{x}) - \sigma I \preceq (L - \sigma)I$ .

$$\implies \langle \nabla h(\mathbf{x}) - \nabla h(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{1}{L - \sigma} \|\nabla h(\mathbf{x}) - \nabla h(\mathbf{y})\|^2 \quad \text{by co-coercivity of smooth and convex functions}$$

Then, rearranging the terms finishes the proof. □

# Smooth and Strongly Convex

*Proof:*  $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 = \|\Pi_{\mathcal{X}}[\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)] - \mathbf{x}^*\|^2$  (GD)

$$\leq \|\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) - \mathbf{x}^*\|^2 \text{ (Pythagoras Theorem)}$$
$$= \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\eta_t \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle + \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2$$
$$\leq \left(1 - \frac{2\eta_t \sigma L}{L + \sigma}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \left(\eta_t^2 - \frac{2\eta_t}{L + \sigma}\right) \|\nabla f(\mathbf{x}_t)\|^2$$

*exploiting co-coercivity of smoothness and strong convexity*

$$\langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle = \langle \nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}^*), \mathbf{x}_t - \mathbf{x}^* \rangle \geq \frac{1}{L + \sigma} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L\sigma}{L + \sigma} \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

$$\Rightarrow \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \left(1 - \frac{2\eta_t \sigma L}{L + \sigma}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \left(\eta_t^2 - \frac{2\eta_t}{L + \sigma}\right) \|\nabla f(\mathbf{x}_t)\|^2$$

*serving as the “one-step improvement” in the analysis*

# Smooth and Strongly Convex

*Proof:* 
$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \left(1 - \frac{2\eta_t \sigma L}{L + \sigma}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \left(\eta_t^2 - \frac{2\eta_t}{L + \sigma}\right) \|\nabla f(\mathbf{x}_t)\|^2$$

The step size configuration:

(i) first, we need  $1 - \frac{2\eta_t \sigma L}{L + \sigma} < 1$  to ensure the contraction property;

(ii) second, we hope  $(\eta_t^2 - \frac{2\eta_t}{L + \sigma}) \leq 0$ , or it becomes 0 is enough.

$\Rightarrow$  a feasible (simple) setting:  $\eta_t = \eta = \frac{2}{L + \sigma}$

$$\Rightarrow \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \left(1 - \frac{4\sigma L}{(L + \sigma)^2}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2 = \left(\frac{L - \sigma}{L + \sigma}\right)^2 \|\mathbf{x}_t - \mathbf{x}^*\|^2 = \left(\frac{\kappa - 1}{\kappa + 1}\right)^2 \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

$$\Rightarrow \|\mathbf{x}_T - \mathbf{x}^*\|^2 \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^{2(T-1)} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 \leq \exp\left(-\frac{4(T-1)}{\kappa + 1}\right) \|\mathbf{x}_1 - \mathbf{x}^*\|^2$$

# Smooth and Strongly Convex

*Proof:*  $\|\mathbf{x}_T - \mathbf{x}^*\|^2 \leq \left(\frac{\kappa-1}{\kappa+1}\right)^{2(T-1)} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 \leq \exp\left(-\frac{4(T-1)}{\kappa+1}\right) \|\mathbf{x}_1 - \mathbf{x}^*\|^2$

*Next step:* relating  $\|\mathbf{x}_T - \mathbf{x}^*\|^2$  to  $f(\mathbf{x}_T) - f(\mathbf{x}^*)$ .

$$f(\mathbf{x}_t) \leq f(\mathbf{x}^*) + \langle \nabla f(\mathbf{x}^*), \mathbf{x}_t - \mathbf{x}^* \rangle + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2 = f(\mathbf{x}^*) + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

(in unconstrained case,  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ )

$$\Rightarrow f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2} \exp\left(-\frac{4(T-1)}{\kappa+1}\right) \|\mathbf{x}_1 - \mathbf{x}^*\|^2 = \mathcal{O}\left(\exp\left(-\frac{T}{\kappa}\right)\right).$$

□

# Constrained Optimization

- A *generalized* one-step improvement lemma for smooth optimization.

**Lemma 5.** Suppose  $f$  is  $L$ -smooth. Let  $\mathbf{x}, \mathbf{u} \in \mathcal{X}$ ,  $\mathbf{x}_{t+1} = \Pi_{\mathcal{X}}[\mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)]$ , and  $g(\mathbf{x}) = L(\mathbf{x} - \mathbf{x}_{t+1})$ . Then the following holds true:

$$f(\mathbf{x}_{t+1}) - f(\mathbf{u}) \leq \langle g(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|g(\mathbf{x}_t)\|^2.$$

comparator  $\mathbf{u}$  is introduced because now GD is not necessary “descent” due to the projection

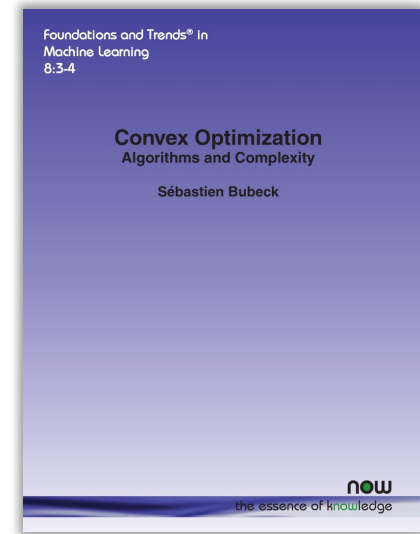
- In unconstrained case,  $g(\mathbf{x}_t) = \nabla f(\mathbf{x}_t)$ .
- In unconstrained case, setting  $\mathbf{u} = \mathbf{x}_t$  recovers the one-step improvement:  
$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$$

# Constrained Optimization

**Same** convergence rates as unconstrained case can be obtained in the constrained setting for smooth optimization.

Detailed proofs for constrained case are not presented in our course. The proof follows the same vein via some twists, we refer anyone interested to the following parts in **Bubeck's book**:

- *Constrained* + smooth + convex: **Section 3.2**
- *Constrained* + smooth + strongly convex: **Section 3.4.2**



**Convex Optimization:  
Algorithms and Complexity**  
Sebastien Bubeck

Foundations and Trends in ML, 2015

# Lower Bound

Lower bounds reflect the **difficulty** of the problem, regardless of algorithms.

*notice: this lower bound only holds for first-order methods*

Table 1: A summary of convergence rates of GD for different function families.

Function Family		Convergence Rate	Lower Bound	Optimal?
$G$ -Lipschitz	convex	$\mathcal{O}(1/\sqrt{T})$	$\Omega(1/\sqrt{T})$	✓
	$\sigma$ -strongly convex	$\mathcal{O}(1/T)$	$\Omega(1/T)$	✓
$L$ -smooth	convex	$\mathcal{O}(1/T)$	$\Omega(1/T^2)$	✗
	$\sigma$ -strongly convex	$\mathcal{O}\left(\exp\left(-\frac{T}{\kappa}\right)\right)$	$\Omega\left(\exp\left(-\frac{T}{\sqrt{\kappa}}\right)\right)$	✗

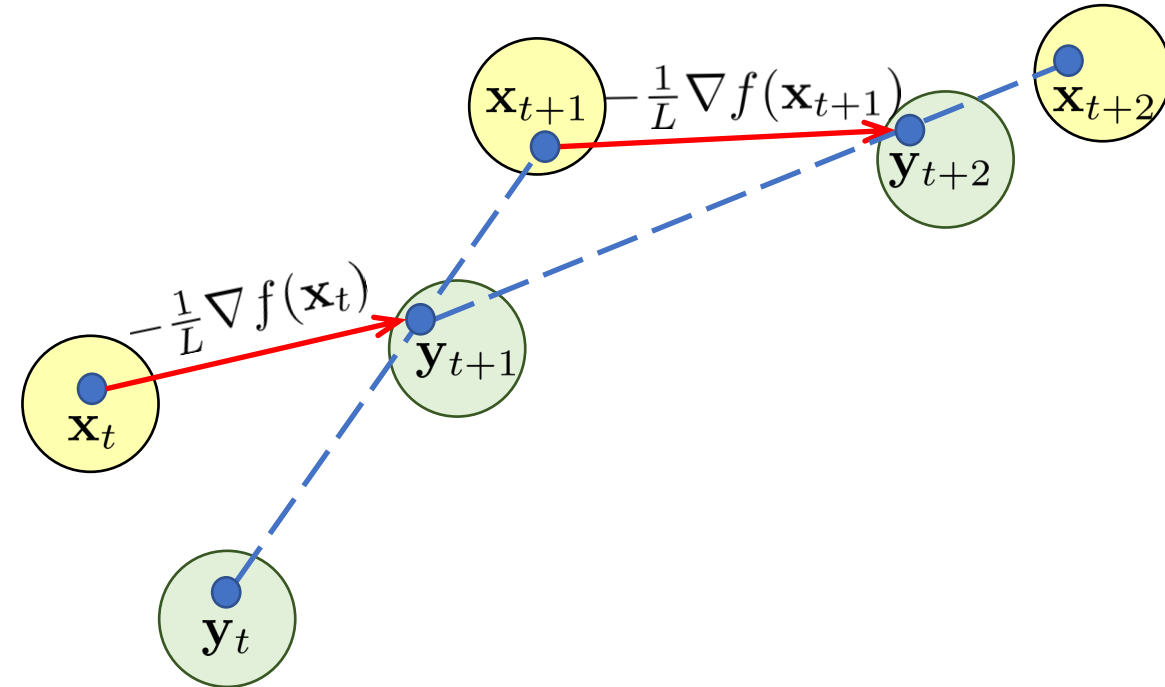
⇒ GD is **suboptimal** in *smooth* convex optimization!

# Part 2. Nesterov's Accelerated GD

- Algorithm
- Smooth and Convex
- Smooth and Strongly Convex

# Nesterov's Accelerated GD

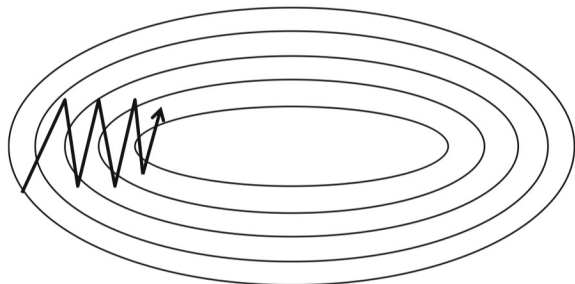
$$\begin{aligned}\mathbf{y}_{t+1} &= \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \\ \mathbf{x}_{t+1} &= (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t\end{aligned}$$



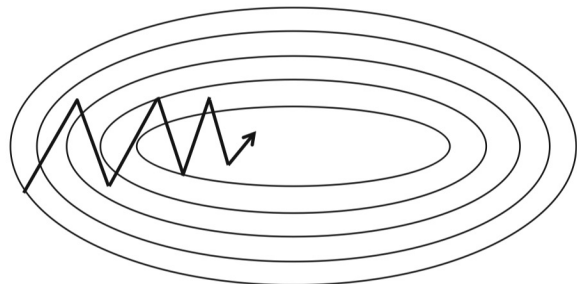
- Define  $\mathbf{x}_1 = \mathbf{y}_1$ .
- $\alpha_t < 0$  is a *time-varying* mixing rate of  $\mathbf{y}_t$  and  $\mathbf{y}_{t+1}$ .
- $\mathbf{x}_{t+1} = \mathbf{y}_{t+1} + \alpha_t(\mathbf{y}_t - \mathbf{y}_{t+1})$  is an *extrapolated* point, i.e., with *momentum*.

# Nesterov's Accelerated GD

- a momentum term is added to improve convergence
- the descent property is relaxed and not ensured now



GD

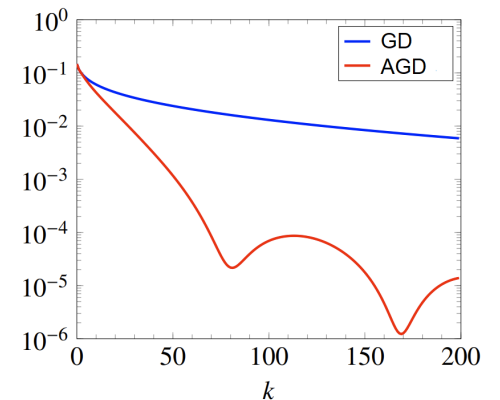
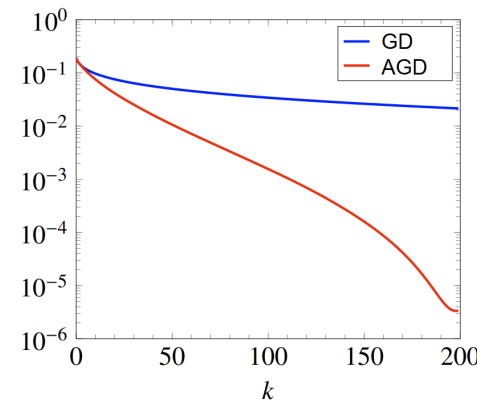


Accelerated GD

## Example

$$\text{minimize} \quad \log \sum_{i=1}^p \exp(a_i^T x + b_i)$$

- two randomly generated problems with  $p = 2000, n = 1000$
- same fixed step size used for gradient method and FISTA
- figures show  $(f(x^{(k)}) - f^*)/f^*$



Accelerated proximal gradient methods

7.9

<https://www.seas.ucla.edu/~vandenbe/236C/lectures/fgrad.pdf>

# Convergence of Nesterov's Accelerated GD

**Theorem 3.** *Let  $f$  be convex and  $L$ -smooth. Nesterov's accelerated GD is configured as*

$$\mathbf{y}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t), \quad \mathbf{x}_{t+1} = (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t,$$

where  $\lambda_0 = 0$ ,  $\lambda_t = \frac{1 + \sqrt{1 + 4\lambda_{t-1}^2}}{2}$ , and  $\alpha_t = \frac{1 - \lambda_t}{\lambda_{t+1}}$ . Then, we have

$$f(\mathbf{y}_T) - f(\mathbf{x}^*) \leq \frac{2L \|\mathbf{x}_1 - \mathbf{x}^*\|^2}{T^2} = \mathcal{O}\left(\frac{1}{T^2}\right).$$

# Proof of AGD Convergence

*Proof:* first we prove the following *generalized one-step improvement lemma*.

**Lemma 6.** For any  $\mathbf{u} \in \mathcal{X}$ , if  $\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$ , then the following holds true:

$$f(\mathbf{x}_{t+1}) - f(\mathbf{u}) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$$

comparator  $\mathbf{u}$  is introduced because now GD is not necessary “descent” due to the momentum

Setting  $\mathbf{u} = \mathbf{x}_t$  recovers the one-step improvement used in earlier analysis.

# Generalized One-Step Improvement

**Lemma 6.** For any  $\mathbf{u} \in \mathcal{X}$ , if  $\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$ , then the following holds true:

$$f(\mathbf{x}_{t+1}) - f(\mathbf{u}) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$$

Setting  $\mathbf{u} = \mathbf{x}_t$  recovers the one-step improvement used in earlier analysis.

**Proof:**

$$\begin{aligned} f(\mathbf{x}_{t+1}) - f(\mathbf{u}) &= f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) + f(\mathbf{x}_t) - f(\mathbf{u}) \\ &\leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle \quad (\text{smoothness and convexity}) \\ &= \langle \nabla f(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{u} \rangle + \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \quad (\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)) \\ &= \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \end{aligned}$$

□

# Proof of AGD Convergence

$$\begin{aligned}\mathbf{y}_{t+1} &= \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \\ \mathbf{x}_{t+1} &= (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t\end{aligned}$$

*Proof:* (continued proving Theorem 3)

**Lemma 6.** For any  $\mathbf{u} \in \mathcal{X}$ , if  $\mathbf{x}' = \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})$ , then the following holds true:

$$f(\mathbf{x}') - f(\mathbf{u}) \leq \langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{u} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x})\|^2.$$

(i) Plugging in  $\mathbf{u} = \mathbf{y}_t$ :  $f(\mathbf{y}_{t+1}) - f(\mathbf{y}_t) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{y}_t \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$

(ii) Plugging in  $\mathbf{u} = \mathbf{x}^*$ :  $f(\mathbf{y}_{t+1}) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$

LHS of  $(\lambda_t - 1)(i) + (ii)$  equals:

$$(\lambda_t - 1)(f(\mathbf{y}_{t+1}) - f(\mathbf{y}_t)) + f(\mathbf{y}_{t+1}) - f(\mathbf{x}^*) = \lambda_t(f(\mathbf{y}_{t+1}) - f(\mathbf{x}^*)) - (\lambda_t - 1)(f(\mathbf{y}_t) - f(\mathbf{x}^*))$$

Define  $\delta_t \triangleq f(\mathbf{y}_t) - f(\mathbf{x}^*)$ , LHS =  $\lambda_t \delta_{t+1} - (\lambda_t - 1) \delta_t$  *Goal: design a telescoping series*

# Proof of AGD Convergence

$$\begin{aligned} \mathbf{y}_{t+1} &= \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \\ \mathbf{x}_{t+1} &= (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t \end{aligned}$$

*Proof:* (continued proving Theorem 3)

(i) Plugging in  $\mathbf{u} = \mathbf{y}_t$ :  $f(\mathbf{y}_{t+1}) - f(\mathbf{y}_t) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{y}_t \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$

(ii) Plugging in  $\mathbf{u} = \mathbf{x}^*$ :  $f(\mathbf{y}_{t+1}) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$

RHS of  $(\lambda_t - 1)(i) + (i)$  equals:

$$\begin{aligned} &(\lambda_t - 1) \left( \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{y}_t \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \right) + \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \\ &= \langle \nabla f(\mathbf{x}_t), \lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^* \rangle - \frac{\lambda_t}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \end{aligned}$$

That is

$$\lambda_t \delta_{t+1} - (\lambda_t - 1) \delta_t \leq \langle \nabla f(\mathbf{x}_t), \lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^* \rangle - \frac{\lambda_t}{2L} \|\nabla f(\mathbf{x}_t)\|^2$$

# Proof of AGD Convergence

$$\begin{aligned} \mathbf{y}_{t+1} &= \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \\ \mathbf{x}_{t+1} &= (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t \end{aligned}$$

*Proof:* (continued proving Theorem 3)

$$\lambda_t \delta_{t+1} - (\lambda_t - 1) \delta_t \leq \langle \nabla f(\mathbf{x}_t), \lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^* \rangle - \frac{\lambda_t}{2L} \|\nabla f(\mathbf{x}_t)\|^2$$

$$\Rightarrow \lambda_t^2 \delta_{t+1} - \lambda_t(\lambda_t - 1) \delta_t \leq \frac{1}{2L} (2 \langle \lambda_t \nabla f(\mathbf{x}_t), L(\lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*) \rangle - \|\lambda_t \nabla f(\mathbf{x}_t)\|^2)$$

**Requirement (1):**  $\lambda_t(\lambda_t - 1) = \lambda_{t-1}^2$

$$\Rightarrow \lambda_t^2 \delta_{t+1} - \lambda_{t-1}^2 \delta_t \leq \frac{1}{2L} (2 \langle \lambda_t \nabla f(\mathbf{x}_t), L(\lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*) \rangle - \|\lambda_t \nabla f(\mathbf{x}_t)\|^2)$$

Denote by  $\mathbf{a} \triangleq \lambda_t \nabla f(\mathbf{x}_t)$ ,  $\mathbf{b} \triangleq L(\lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*)$ .

$$\Rightarrow \lambda_t^2 \delta_{t+1} - \lambda_{t-1}^2 \delta_t \leq \frac{1}{2L} (2 \langle \mathbf{a}, \mathbf{b} \rangle - \|\mathbf{b}\|^2) \leq \frac{1}{2L} (\|\mathbf{b}\|^2 - \|\mathbf{b} - \mathbf{a}\|^2)$$

# Proof of AGD Convergence

$$\begin{aligned} \mathbf{y}_{t+1} &= \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \\ \mathbf{x}_{t+1} &= (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t \end{aligned}$$

*Proof:* (continued proving Theorem 3)

Denote by  $\mathbf{a} \triangleq \lambda_t \nabla f(\mathbf{x}_t)$ ,  $\mathbf{b} \triangleq L(\lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*)$ .

$$\begin{aligned} & \lambda_t^2 \delta_{t+1} - \lambda_{t-1}^2 \delta_t \\ & \leq \frac{1}{2L} (L^2 \|\lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*\|^2 - \|L(\lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*) - \lambda_t \nabla f(\mathbf{x}_t)\|^2) \\ & = \frac{L}{2} \left( \|\lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*\|^2 - \left\| \lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^* - \lambda_t \frac{\nabla f(\mathbf{x}_t)}{L} \right\|^2 \right) \\ & = \frac{L}{2} (\|\lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*\|^2 - \|\lambda_t \mathbf{y}_{t+1} - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*\|^2) \end{aligned}$$

*Goal: design a telescoping series*

# Proof of AGD Convergence

$$\begin{aligned}\mathbf{y}_{t+1} &= \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \\ \mathbf{x}_{t+1} &= (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t\end{aligned}$$

*Proof:* (continued proving Theorem 3)

$$\lambda_t^2 \delta_{t+1} - \lambda_{t-1}^2 \delta_t \leq \frac{L}{2} (\|\lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*\|^2 - \|\lambda_t \mathbf{y}_{t+1} - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*\|^2)$$

*Requirement (2):*  $\lambda_t \mathbf{y}_{t+1} - (\lambda_t - 1) \mathbf{y}_t = \lambda_{t+1} \mathbf{x}_{t+1} - (\lambda_{t+1} - 1) \mathbf{y}_{t+1}$

$$\lambda_t^2 \delta_{t+1} - \lambda_{t-1}^2 \delta_t \leq \frac{L}{2} (\|\lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*\|^2 - \|\lambda_{t+1} \mathbf{x}_{t+1} - (\lambda_{t+1} - 1) \mathbf{y}_{t+1} - \mathbf{x}^*\|^2)$$

*telescope*

Define  $\mathbf{z}_t \triangleq \lambda_t \mathbf{x}_t - (\lambda_t - 1) \mathbf{y}_t - \mathbf{x}^*$

$$\lambda_t^2 \delta_{t+1} - \lambda_{t-1}^2 \delta_t \leq \frac{L}{2} (\|\mathbf{z}_t\|^2 - \|\mathbf{z}_{t+1}\|^2) \Rightarrow \lambda_{T-1}^2 \delta_T - \lambda_0^2 \delta_1 = \frac{L}{2} (\|\mathbf{z}_1\|^2 - \|\mathbf{z}_T\|^2)$$

# Proof of AGD Convergence

$$\begin{aligned} \mathbf{y}_{t+1} &= \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \\ \mathbf{x}_{t+1} &= (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t \end{aligned}$$

*Proof:* (continued proving Theorem 3)

$$\lambda_{T-1}^2 \delta_T - \lambda_0^2 \delta_1 = \frac{L}{2} (\|\mathbf{z}_1\|^2 - \|\mathbf{z}_T\|^2)$$

*Requirement (3):*  $\lambda_0 = 0$

$$\lambda_{T-1}^2 \delta_T \leq \frac{L}{2} \|\mathbf{z}_1\|^2 \Rightarrow \delta_T \leq \frac{L \|\mathbf{z}_1\|^2}{2\lambda_{T-1}^2} = \frac{L \|\lambda_1 \mathbf{x}_1 - (\lambda_1 - 1) \mathbf{y}_1 - \mathbf{x}^*\|^2}{2\lambda_{T-1}^2}$$

*Requirement (4):*  $\mathbf{x}_1 = \mathbf{y}_1$

$$\lambda_{T-1}^2 \delta_T \leq \frac{L}{2} \|\mathbf{z}_1\|^2 \Rightarrow \delta_T \leq \frac{L \|\mathbf{z}_1\|^2}{2\lambda_{T-1}^2} = \frac{L \|\mathbf{x}_1 - \mathbf{x}^*\|^2}{2\lambda_{T-1}^2}$$

# Proof

*Proof:* (continued proving Theorem 3)

**Theorem 3.** Let  $f$  be convex and  $L$ -smooth. Nesterov's accelerated GD is configured as

$$\mathbf{y}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t), \quad \mathbf{x}_{t+1} = (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t,$$

where  $\lambda_0 = 0$ ,  $\lambda_t = \frac{1 + \sqrt{1 + 4\lambda_{t-1}^2}}{2}$ , and  $\alpha_t = \frac{1 - \lambda_t}{\lambda_{t+1}}$ . Then, we have

$$f(\mathbf{y}_T) - f(\mathbf{x}^*) \leq \frac{2L \|\mathbf{x}_1 - \mathbf{x}^*\|^2}{T^2} = \mathcal{O}\left(\frac{1}{T^2}\right).$$

*Requirement (1):*  $\lambda_t(\lambda_t - 1) = \lambda_{t-1}^2$

$$\lambda_t = \frac{1 + \sqrt{1 + 4\lambda_{t-1}^2}}{2}$$

*Requirement (2):*  $\lambda_t \mathbf{y}_{t+1} - (\lambda_t - 1) \mathbf{y}_t = \lambda_{t+1} \mathbf{x}_{t+1} - (\lambda_{t+1} - 1) \mathbf{y}_{t+1}$

$$\mathbf{x}_{t+1} = \mathbf{y}_{t+1} - \frac{1 - \lambda_t}{\lambda_{t+1}} (\mathbf{y}_t - \mathbf{y}_{t+1}) \Rightarrow \alpha_t = \frac{1 - \lambda_t}{\lambda_{t+1}}$$

*Requirement (3):*  $\lambda_0 = 0$

*Requirement (4):*  $\mathbf{x}_1 = \mathbf{y}_1$

$$\lambda_t = \frac{1 + \sqrt{1 + 4\lambda_{t-1}^2}}{2} \Rightarrow \lambda_t \geq \frac{t + 1}{2} \Rightarrow \delta_T \leq \frac{L \|\mathbf{x}_1 - \mathbf{x}^*\|^2}{2\lambda_{T-1}^2} \leq \frac{2L \|\mathbf{x}_1 - \mathbf{x}^*\|^2}{T^2} = \mathcal{O}\left(\frac{1}{T^2}\right) \quad \square$$

# Smooth and Strongly Convex

**Theorem 4.** *Let  $f$  be  $\sigma$ -strongly convex and  $L$ -smooth, then Nesterov's accelerated gradient descent:*

$$\mathbf{y}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t), \quad \mathbf{x}_{t+1} = \mathbf{y}_{t+1} + \frac{\sqrt{\gamma} - 1}{\sqrt{\gamma} + 1} (\mathbf{y}_{t+1} - \mathbf{y}_t)$$

*satisfies*

$$f(\mathbf{y}_T) - f(\mathbf{x}^*) \leq \frac{\sigma + L}{2} \|\mathbf{x}^* - \mathbf{x}_1\|^2 \exp\left(-\frac{T}{\sqrt{\gamma}}\right),$$

*where  $\gamma \triangleq L/\sigma$  denotes the condition number.*

*core technique: estimate sequence*

# Smooth and Strongly Convex

- Proof sketch

*Core technique:* construct an estimate sequence (*developed by Yurii Nesterov*)

$$\Phi_1(\mathbf{x}) \triangleq f(\mathbf{x}_1) + \frac{\sigma}{2} \|\mathbf{x} - \mathbf{x}_1\|^2$$

$$\Phi_{t+1}(\mathbf{x}) \triangleq (1 - \theta)\Phi_t(\mathbf{x}) + \theta \left( f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \frac{\sigma}{2} \|\mathbf{x} - \mathbf{x}_t\|^2 \right)$$

The estimate sequence  $\{\Phi_t\}_{t=1}^T$  is required to satisfy some nice properties:

(i)  $\Phi_{t+1}(\mathbf{x}) - f(\mathbf{x}) \leq (1 - \theta)^t (\Phi_1(\mathbf{x}) - f(\mathbf{x})) \Rightarrow$  approximate  $f$  well.

(ii)  $f(\mathbf{y}_t) \leq \min_{\mathbf{x} \in \mathbb{R}^d} \Phi_t(\mathbf{x}) \Rightarrow$  useful when giving the convergence rate.

*It can be proven that the above construction satisfies the two properties.*

# Smooth and Strongly Convex

- Proof sketch

*Core technique:* construct an estimate sequence (*developed by Yurii Nesterov*)

$$\Phi_1(\mathbf{x}) \triangleq f(\mathbf{x}_1) + \frac{\sigma}{2} \|\mathbf{x} - \mathbf{x}_1\|^2$$

$$\Phi_{t+1}(\mathbf{x}) \triangleq (1 - \theta)\Phi_t(\mathbf{x}) + \theta \left( f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \frac{\sigma}{2} \|\mathbf{x} - \mathbf{x}_t\|^2 \right)$$

$$f(\mathbf{y}_t) - f(\mathbf{x}^*) \stackrel{(ii)}{\leq} \min_{\mathbf{x} \in \mathbb{R}^d} \Phi_t(\mathbf{x}) - f(\mathbf{x}^*) \leq \Phi_t(\mathbf{x}^*) - f(\mathbf{x}^*) \quad (\text{by property (ii)})$$

$$\stackrel{(i)}{\leq} (1 - \theta)^t (\Phi_1(\mathbf{x}^*) - f(\mathbf{x}^*)) \quad (\text{by property (i)})$$

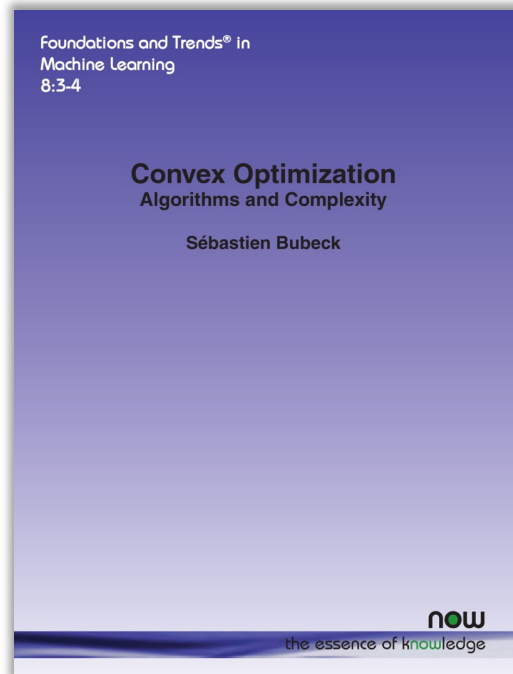
$$= (1 - \theta)^t \left( f(\mathbf{x}_1) + \frac{\sigma}{2} \|\mathbf{x}^* - \mathbf{x}_1\|^2 - f(\mathbf{x}^*) \right) \quad (\text{definition of } \Phi_1)$$

$$\lesssim (\sigma + L) \|\mathbf{x}^* - \mathbf{x}_1\|^2 \exp(-\theta t) \quad (\text{smoothness})$$

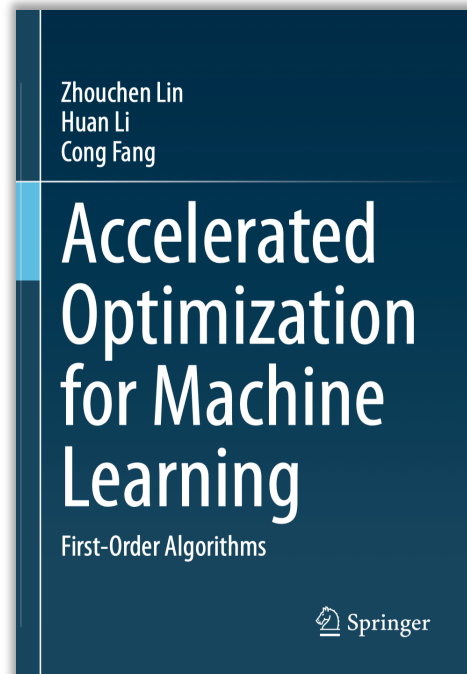
# Estimate Sequence

- Admittedly, how to construct estimate sequence is highly *tricky*

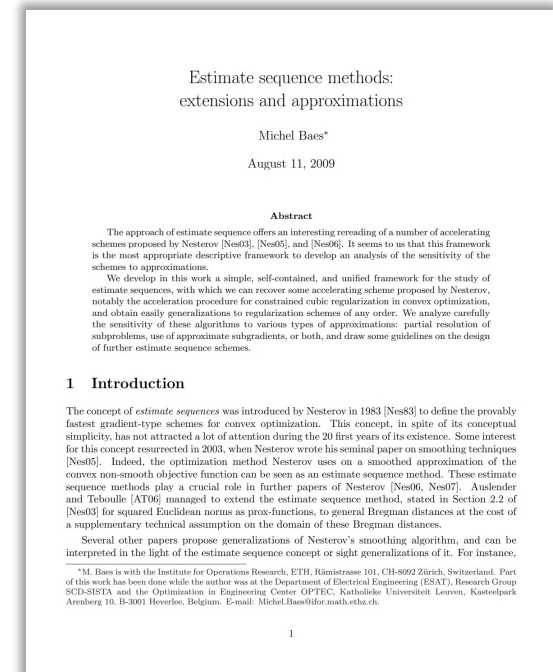
## References:



Chapter 3.7



Chapter 2.1



M. Baes, Estimate sequence methods: extensions and approximations. Technical report, ETH, Zürich (2009)

# References for Nesterov's Accelerated GD

Nesterov's four ideas (three acceleration methods):

- Y. Nesterov (1983), A method for solving a convex programming problem with convergence rate  $O(1/k^2)$
- Y. Nesterov (1988), On an approach to the construction of optimal methods of minimization of smooth convex functions
- Y. Nesterov (2005), Smooth minimization of non-smooth functions
- Y. Nesterov (2007), Gradient methods for minimizing composite objective function



**Yurii Nesterov**  
1956 –  
UCLouvain, Belgium

# Nesterov, Y. (1983), A method of solving a convex programming problem with convergence rate $O(1/k^2)$ , Soviet Mathematics Doklady 27(2), 372–376.

Докл. Акад. Наук СССР  
Том 269 (1983), № 3

## A METHOD OF SOLVING A CONVEX PROGRAMMING PROBLEM WITH CONVERGENCE RATE $O(1/k^2)$

UDC 51

YU. E. NESTEROV

1. In this note we propose a method of solving a convex programming problem in a Hilbert space  $E$ . Unlike the majority of convex programming methods this method constructs a minimizing sequence of points  $\{x_k\}$ . This property allows us to reduce the amount of computation. At the same time, it is possible to obtain an estimate of convergence rate improved for the class of problems under consideration (see [2]).

2. Consider first the problem of unconstrained minimization. We will assume that  $f(x)$  belongs to the class  $C^{1,1}(E)$ , i.e.  $L > 0$  such that for all  $x, y \in E$

$$(1) \quad \|f'(x) - f'(y)\| \leq L\|x - y\|.$$

From (1) it follows that for all  $x, y \in E$

$$(2) \quad f(y) \leq f(x) + \langle f'(x), y - x \rangle + 0.5L\|y - x\|^2.$$

To solve the problem  $\min\{f(x) | x \in E\}$  with a nonempty set  $E$  the following method.

0) Select a point  $y_0 \in E$ . Put

$$(3) \quad k = 0, \quad a_0 = 1, \quad x_{-1} = y_0, \quad \alpha_{-1} = \|y_0 - z\|/L,$$

where  $z$  is an arbitrary point in  $E$ ,  $z \neq y_0$  and  $f'(z) \neq f'(y_0)$ .

1)  $k$ th iteration. a) Calculate the smallest index  $i \geq 0$  for which

$$(4) \quad f(y_k) - f(y_k - 2^{-i}\alpha_{k-1}f'(y_k)) \geq 2^{i-1}\alpha_k,$$

b) Put

$$\alpha_k = 2^{-i}\alpha_{k-1}, \quad x_k = y_k - \alpha_k f'(y_k).$$

$$(5) \quad a_{k+1} = (1 + \sqrt{4a_k^2 + 1})/2,$$

$$y_{k+1} = x_k + (a_k - 1)(x_k - x_{k-1})/a_{k+1}.$$

The way in which the one-dimensional search (4) is halted [2]. The difference is only that in (4) the subdivision in the  $k$ th iteration is made with  $\alpha_{k-1}$  (and not with 1 as in [2]). In view of this (see the proof of Theorem 1) the sequence  $\{x_k\}$  is constructed by method (3)–(5), no more iterations will be made. The recalculation of the points  $y_k$  in (5) is

1980 Mathematics Subject Classification. Primary 90C25.

372

Let us also remark that method (3)–(5) does not guarantee the convergence of the sequences  $\{x_k\}$  and  $\{y_k\}$ .

THEOREM 1. Let  $f(x)$  be a convex function in  $C^{1,1}(E)$ . The sequence  $\{x_k\}$  is constructed by method (3)–(5), then

1) For any  $k \geq 0$ :

$$(6) \quad f(x_k) - f^* \leq C/(k+1),$$

where  $C = 4L\|y_0 - x^*\|^2$  and  $f^* = f(x^*)$ ,  $x^* \in X^*$ .

2) In order to achieve accuracy  $\epsilon$  with respect to the function

a) to compute the gradient of the objective function no more than

b) to evaluate the objective function no more than  $N$

Here and in what follows,  $[\cdot]$  is the integer part of  $\cdot$ .

PROOF. Let  $y_k(\alpha) = y_k - \alpha f'(y_k)$ . From (2) we obtain

$$f(y_k) - f(y_k(\alpha)) \geq 0.5\alpha(2 - L\alpha).$$

Consequently, as soon as  $2^{-i}\alpha_{k-1}$  becomes less than  $1/L$ ,  $\alpha_k$  will not be further decreased. Thus  $\alpha_k \geq 0.5L^{-1}$ .

Let  $p_k = (a_k - 1)(x_{k-1} - x_k)$ . Then  $p_{k+1} - x_k$  Consequently,

$$\|p_{k+1} - x_{k+1} + x^*\|^2 = \|p_k - x_k + x^*\|^2 + 2\langle p_k - x_k + x^*, p_{k+1} - x_{k+1} + x^* \rangle + \|p_{k+1} - x_{k+1} + x^*\|^2.$$

Using inequality (4) and the convexity of  $f(x)$ , we obtain

$$\langle f'(y_{k+1}), y_{k+1} - x^* \rangle \geq f(y_{k+1}) - f^*.$$

$$0.5\alpha_{k+1}\|f'(y_{k+1})\|^2 \leq f(y_{k+1}) - f(x_{k+1})$$

$$- \alpha_{k+1}\langle f'(y_{k+1}), y_{k+1} - x^* \rangle.$$

We substitute these two inequalities into the preceding

$$\|p_{k+1} - x_{k+1} + x^*\|^2 - \|p_k - x_k + x^*\|^2 \leq 2\langle p_k - x_k + x^*, p_{k+1} - x_{k+1} + x^* \rangle +$$

$$- 2\alpha_{k+1}\alpha_{k+1}\langle f'(y_{k+1}), y_{k+1} - x^* \rangle + (a_{k+1}^2 - a_k^2)\|p_k - x_k + x^*\|^2 \leq$$

$$\leq -2\alpha_{k+1}\alpha_{k+1}\langle f'(y_{k+1}), y_{k+1} - x^* \rangle + 2(a_{k+1}^2 - a_k^2)\|p_k - x_k + x^*\|^2 \leq$$

$$= 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

$$\leq 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), y_{k+1} - x^* \rangle - 2\alpha_{k+1}\alpha_{k+1}^2\langle f'(y_{k+1}), p_k - x_k + x^* \rangle \leq$$

be calculated. Let us remark, however, that to each addition of a point corresponds a halving of  $\alpha_k$ . Therefore the total number of additions does not exceed  $\log_2(2L\alpha_{-1}) + 1$ . This completes the proof of Theorem 1.

If the Lipschitz constant  $L$  is known for the gradient of the function  $f(x)$ , then in the method (3)–(5) for any  $k \geq 0$ . In fact, to hold, and therefore Theorem 1 remains valid. Let  $y_0 = x^* + \sqrt{2L/\epsilon}$ ,  $N = 1/\epsilon$ .

To conclude this section we will show how one may make the problem of minimizing a strictly convex function.

Assume that  $f(x) - f^* \geq 0.5m\|x - x^*\|^2$  for all  $x \in E$  and constant  $m$  is known.

We introduce the following halting rule in the method (3)–(5):

c) We stop when

$$(7) \quad k \geq 2\sqrt{2/(m\alpha_k)} - 2.$$

Suppose that the halting has occurred in the  $N$ th step.

(3)–(5), one has  $N \leq 4\sqrt{2/(m\alpha_k)} - 1$ . At the same time,

$$f(x_N) - f^* \leq \frac{2\|y_0 - x^*\|^2}{\alpha_N(N+2)^2} \leq 0.25m\|y_0 - x^*\|^2.$$

After the point  $x_N$  has been obtained, it is necessary to begin calculating, by the method (3)–(5), (7), from the point  $x_N$ .

As a result we obtain that after each  $4\sqrt{2/(m\alpha_k)} - 1$  iterations the value of the function decreases by a factor of 2. Thus the number of iterations cannot be improved (up to a dimensionless constant) among the class of strictly convex functions in  $C^{1,1}(E)$  (see [1]).

3. Consider the following extremal problem:

$$(8) \quad \min\{F(\tilde{f}(x)) | x \in Q\}$$

where  $Q$  is a convex closed set in  $E$ ,  $F(u)$ , with  $u \in R^m$ , is a positive homogeneous of degree one, and  $\tilde{f}(x) = (f_1(x), \dots, f_m(x))$  is a vector of continuously differentiable functions on  $E$ . The set  $X$  is assumed to be nonempty. In addition to this, we will assume that the functions  $\{F(\cdot), \tilde{f}(\cdot)\}$  has the following property:

(\*) If there exists a vector  $\lambda \in \partial F(0)$  such that  $\lambda^{(k)} < 0$ .

The notation  $\partial F(0)$  means the subdifferential of the function  $F$  at the point 0.

As is well known, the identity  $F(u) \equiv \max\{\langle \lambda, u \rangle | \lambda \in \partial F(0)\}$  holds for all  $u \in R^m$ . The convexity of the function  $F(\tilde{f}(x))$  on all of  $E$ .

Problem (8) can be written in minimax form:

$$(9) \quad \min\{\max\{\langle \lambda, \tilde{f}(x) \rangle | \lambda \in \partial F(0)\} | x \in Q\}$$

One can show that the fact that the set  $X^*$  is nonempty is equivalent to the existence of a saddle point  $(\lambda^*, x^*)$  for problem (9).

Of problem (9) can be written as  $\Omega^* = \Lambda^* \times X^*$ , where

$$\Lambda^* = \text{Arg max}\{\Psi(\lambda) | \lambda \in \partial F(0)\}, \quad \Psi(\lambda) =$$

The problem

$$\max\{\Psi(\lambda) | \lambda \in \partial F(0) \cap \text{dom } \Psi\}$$

will be called the problem dual to (8).

Suppose the functions  $f_k(x)$ ,  $k = 1, \dots, m$ , in problem (8) with constants  $L^{(k)} \geq 0$ . Let  $L = (L^{(1)}, \dots, L^{(m)})$ .

Consider the function

$$\Phi(y, A, z) = F(\tilde{f}(y, z)) + 0.5A\|y - z\|^2$$

where

$$\tilde{f}(y, z) = (f^{(1)}(y, z), \dots, f^{(m)}(y, z)),$$

$$f^{(k)}(y, z) = f_k(y) + \langle f'_k(y), z - y \rangle, \quad k = 1, \dots, m,$$

and  $A$  is a positive constant. Let

$$\Phi^*(y, A) = \min\{\Phi(y, A, z) | z \in Q\}, \quad T(y, A) = \arg \min\{\Phi(y, A, z) | z \in Q\}.$$

Observe that the mapping  $y \rightarrow T(y, A)$  is a natural generalization of the "gradient" mapping introduced in [1] in connection with the minimizing functions of the form  $\max_{1 \leq k \leq m} f_k(x)$ . For the case of the "gradient" mapping of [1] we have

$$(10) \quad \Phi^*(y, A) + A\langle y - T(y, A), x - y \rangle + 0.5A\|y - T(y, A)\|^2$$

for all  $x \in Q$ ,  $y \in E$  and  $A \geq 0$ , and if  $A \geq F(L)$ , then

$$\Phi^*(y, A) \geq F(\tilde{f}(T(y, A))).$$

To solve problem (8) we propose the following method.

0) Select a point  $y_0 \in E$ . Put

$$(11) \quad k = 0, \quad a_0 = 1, \quad x_{-1} = y_0, \quad \alpha_{-1} = \|y_0 - z\|/L,$$

where  $z = (L^{(1)}y_0, \dots, L^{(m)}y_0)$ ,  $L^{(k)} = \|f'_k(y_0) - f'_k(z)\|/\|y_0 - z\|$  in  $E$ ,  $z \neq y_0$ .

1)  $k$ th iteration. a) Calculate the smallest index  $i \geq 0$  for which

$$(12) \quad \Phi^*(y_k, 2^i A_{k-1}) \geq F(\tilde{f}(T(y_k, 2^i A_{k-1})))$$

b) Put  $A_k = 2A_{k-1}$ ,  $x_k = T(y_k, A_k)$  and

$$(13) \quad a_{k+1} = (1 + \sqrt{4a_k^2 + 1})/2,$$

$$y_{k+1} = x_k + (a_k - 1)(x_k - x_{k-1})/a_{k+1}.$$

It is not hard to see that the method (3)–(5) is simply the method (11)–(13) for the unconstrained minimization problem and  $Q = E$  in (8).

THEOREM 2. If the sequence  $\{x_k\}$  is constructed by method (11)–(13), then

1) For any  $k \geq 0$

$$F(\tilde{f}(x_k)) - F(\tilde{f}(x^*)) \leq C_1/(k+1),$$

where  $C_1 = 4F(L)\|y_0 - x^*\|^2$ ,  $x^* \in X^*$ .

2) To obtain accuracy  $\epsilon$  with respect to the functional, one needs

a) to solve an auxiliary problem  $\min\{\Phi(y, A, x) | x \in Q\}$  no more than

$$[\sqrt{C_1/\epsilon}] + [\max\{\log_2(F(L)/A_{-1}), 0\}]$$

$$times,$$

b) to evaluate the collection of gradients  $f'_1(y), \dots, f'_m(y)$  no more than  $[\sqrt{C_1/\epsilon}]$  times, and

c) to evaluate the vector-valued function  $\tilde{f}(x)$  at most

$$2[\sqrt{C_1/\epsilon}] + [\max\{\log_2(F(L)/A_{-1}), 0\}]$$

$$times.$$

Theorem 2 is proved in essentially the same way as Theorem 1. It is only necessary to use (10) instead of (2), while the analogue of  $\alpha_k f'_k(y_k)$  will be the vector  $y_k - T(y_k, A_k)$ , and the analogue of  $\alpha_k$  the values of  $A_k^{-1}$ .

Just as in the method (3)–(5), in the method (11)–(13) one can take into account information about the constant  $F(L)$  and the parameter of strict convexity of the function  $F(\tilde{f}(x)) - m$  (for this, of course, we must have  $y_0 \in Q$ ).

In conclusion let us mention two important special cases of problem (8) in which the auxiliary problem  $\min\{\Phi(y, A, x) | x \in Q\}$  turns out to be rather simple.

a) Minimization of a smooth function on a simple set. By a simple set we understand a set for which the projection operator can be written in explicit form. In this case  $m = 1$  and  $F(y) = y$  in problem (8), and

$$\Phi^*(y, A) = f(y) - 0.5A^{-1}\|f'(y)\|^2 + 0.5A\|T(y, A) - y + A^{-1}f'(y)\|^2,$$

in the method (11)–(13), where

$$T(y, A) = \arg \min\{\|y - A^{-1}f'(y) - z\| | z \in Q\}.$$

b) Unconstrained minimization (in problem (8),  $Q = E$ ). In this case the auxiliary problem  $\min\{\Phi(y, A, x) | x \in E\}$  is equivalent to the following dual problem:

$$(14) \quad \max\left\{-0.5A^{-1}\left\|\sum_{k=1}^m \lambda^{(k)} f'_k(y)\right\|^2 + \sum_{k=1}^m \lambda^{(k)} f_k(y) \mid (\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(m)}) \in \partial F(0)\right\}.$$

Here

$$T(y, A) = y - A^{-1} \sum_{k=1}^m \lambda^{(k)}(y) f'_k(y),$$

where the  $\lambda^{(k)}(y)$ ,  $k = 1, \dots, m$ , are determined by the condition that the set  $\partial F(0)$  is used.

Such cases problem (14) is the stationary problem.

The author expresses his sincere thanks to V. A. Izrael for his interest in the questions considered here.

Central Economic-Mathematical Institute  
Academy of Sciences of the USSR

Received 19/JULY/82

BIBLIOGRAPHY

1. A. S. Nemirovskii and D. B. Yudin, Complexity of problems and efficiency of optimization methods, "Nauka", Moscow, 1979, (Russian).

2. B. N. Pshenichnyi and Yu. M. Danilin, Numerical methods in extremal problems, "Nauka", Moscow, 1975; French transl., "Mir", Moscow, 1977.

Translated by A. ROSA

978

978

978

978

978

978

978

978

978

978

978

978

978

978

978

978

Nesterov, Y. (1983), A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ , Soviet Mathematics Doklady 27(2), 372–376.

546

# More Explanations for Nesterov's AGD

- Ordinary Differentiable Equations

- Su, W., Boyd, S., & Candes, E. (2014). A differential equation for modeling Nesterov's accelerated gradient method: theory and insights. *Advances in neural information processing systems* 27 (NIPS).
- Berthier, R., Bach, F., Flammarion, N., Gaillard, P., & Taylor, A. (2021). A continuized view on Nesterov acceleration. *arXiv preprint arXiv:2102.06035*.

- Variational Analysis

- Wibisono, A., Wilson, A. C., & Jordan, M. I. (2016). A variational perspective on accelerated methods in optimization. *Proceedings of the National Academy of Sciences (PNAS)*, 113(47), E7351-E7358.

- Linear Coupling of GD and MD

- Allen-Zhu, Z., & Orecchia, L. (2017). Linear coupling: An ultimate unification of gradient and mirror descent. *The 8th Innovations in Theoretical Computer Science Conference (ITCS)*.

# Part 3. Extension to Composite Optimization

- Composite Optimization
- Proximal Gradient Method (PG)
- Accelerated Proximal Gradient Method (APG)
- Application to LASSO

# Composite Optimization

- Problem setup

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) \triangleq f(\mathbf{x}) + h(\mathbf{x})$$

where  $f$  is *smooth* (namely, gradient Lipschitz) while  $h$  is *not smooth*.

- The composite optimization problem is common in practice.

**Example 1.** The objective of *LASSO*:  $F(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}^\top X - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1$ , where  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ ,  $\mathbf{y} = [y_1, \dots, y_n]^\top$ .

# Recall Non-composite Optimization

- Consider  $\min_{\mathbf{x}} f(\mathbf{x})$ , and assume  $f$  is  $L$ -smooth.

$$\text{By smoothness: } f(\mathbf{x}) \leq \underbrace{f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2}_{\triangleq Q(\mathbf{x}; \mathbf{y}) \text{ surrogate objective}}$$

$\Rightarrow$  to minimize  $f(\mathbf{x})$ , it suffices to minimize the *surrogate* objective  $Q(\mathbf{x}; \mathbf{y})$ .

**Claim.** GD for smooth functions can be equivalently represented by

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} Q(\mathbf{x}; \mathbf{x}_t) = \Pi_{\mathcal{X}} \left[ \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right],$$

where  $Q(\mathbf{x}; \mathbf{x}_t) \triangleq f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_t\|^2$  is a quadratic upper bound of  $f$  at  $\mathbf{x}_t$ .

# Another View of GD Method

**Claim.** GD for smooth functions can be equivalently represented by

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} Q(\mathbf{x}; \mathbf{x}_t) = \Pi_{\mathcal{X}} \left[ \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right],$$

where  $Q(\mathbf{x}; \mathbf{x}_t) \triangleq f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_t\|^2$  is a quadratic upper bound of  $f$  at  $\mathbf{x}_t$ .

**Proof:**

$$\begin{aligned} \mathbf{x}_{t+1} &= \arg \min_{\mathbf{x} \in \mathcal{X}} Q(\mathbf{x}; \mathbf{x}_t) = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \langle \nabla f(\mathbf{x}_t), \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}\|^2 - L \langle \mathbf{x}, \mathbf{x}_t \rangle \right\} \text{(remove irrelative terms)} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \frac{L}{2} \left( -2 \left\langle \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t), \mathbf{x} \right\rangle + \|\mathbf{x}\|^2 \right) \right\} \text{(rearrange)} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \frac{L}{2} \left\| \mathbf{x} - \left( \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right) \right\|^2 = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\| \mathbf{x} - \left( \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right) \right\| = \Pi_{\mathcal{X}} \left[ \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right] \end{aligned}$$

□

# Another View of GD Method

**Claim.** GD for smooth functions can be equivalently represented by

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} Q(\mathbf{x}; \mathbf{x}_t) = \Pi_{\mathcal{X}} \left[ \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right],$$

where  $Q(\mathbf{x}; \mathbf{x}_t) \triangleq f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_t\|^2$  is a quadratic upper bound of  $f$  at  $\mathbf{x}_t$ .

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} Q(\mathbf{x}; \mathbf{x}_t) = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \boxed{f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle} + \boxed{\frac{L}{2} \|\mathbf{x} - \mathbf{x}_t\|^2} \right\}$$

*linear approximation of  $f$  at  $\mathbf{x}_t$*  *prevent  $\mathbf{x}_t$  from getting too far*

# Composite Optimization

- Problem setup

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) \triangleq f(\mathbf{x}) + h(\mathbf{x})$$

where  $f$  is *smooth* (namely, gradient Lipschitz) while  $h$  is *not smooth*.

*An idea:*

Following previous argument (for non-composite optimization), to minimize  $F \triangleq f + h$ , it suffices to minimize

$$Q(\mathbf{x}; \mathbf{x}_t) \triangleq f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_t\|^2 + h(\mathbf{x})$$

# Formulation

By smoothness:  $f(\mathbf{x}) \leq \underbrace{f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2}_{\triangleq q(\mathbf{x}; \mathbf{y})}$

*surrogate objective*

$\Rightarrow$  to minimize  $F(\mathbf{x}) = f(\mathbf{x}) + h(\mathbf{x})$ , it suffices to minimize  $Q(\mathbf{x}; \mathbf{y}) \triangleq q(\mathbf{x}; \mathbf{y}) + h(\mathbf{x})$ .

$$\begin{aligned} \arg \min_{\mathbf{x}} Q(\mathbf{x}; \mathbf{y}) &= \arg \min_{\mathbf{x}} \left\{ f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 + h(\mathbf{x}) \right\} \\ &= \arg \min_{\mathbf{x}} \left\{ \langle \nabla f(\mathbf{y}), \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}\|^2 - L \langle \mathbf{x}, \mathbf{y} \rangle + h(\mathbf{x}) \right\} \\ &= \arg \min_{\mathbf{x}} \left\{ \frac{L}{2} \left( -2 \left\langle \mathbf{y} - \frac{\nabla f(\mathbf{y})}{L}, \mathbf{x} \right\rangle + \|\mathbf{x}\|^2 \right) + h(\mathbf{x}) \right\} \end{aligned}$$

# Composite Optimization

By smoothness:  $f(\mathbf{x}) \leq f(\mathbf{y}) + \underbrace{\langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2}_{\triangleq q(\mathbf{x}; \mathbf{y})}$

*surrogate objective*

$\Rightarrow$  to minimize  $F(\mathbf{x}) = f(\mathbf{x}) + h(\mathbf{x})$ , it suffices to minimize  $Q(\mathbf{x}; \mathbf{y}) \triangleq q(\mathbf{x}; \mathbf{y}) + h(\mathbf{x})$ .

$$\begin{aligned} \arg \min_{\mathbf{x}} Q(\mathbf{x}; \mathbf{y}) &= \arg \min_{\mathbf{x}} \left\{ \frac{L}{2} \left( -2 \left\langle \mathbf{y} - \frac{\nabla f(\mathbf{y})}{L}, \mathbf{x} \right\rangle + \|\mathbf{x}\|^2 \right) + h(\mathbf{x}) \right\} \\ &= \arg \min_{\mathbf{x}} \left\{ \frac{L}{2} \left\| \mathbf{x} - \left( \mathbf{y} - \frac{\nabla f(\mathbf{y})}{L} \right) \right\|^2 + h(\mathbf{x}) \right\} \end{aligned}$$

*an operator is defined for this (sub-)optimization problem*

# Composite Optimization

- Iteratively solve the surrogate optimization problem.

Deploying the following update rule:

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} Q(\mathbf{x}; \mathbf{x}_t) = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \frac{L}{2} \left\| \mathbf{x} - \left( \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right) \right\|^2 + h(\mathbf{x}) \right\}$$

**Definition 2** (proximal mapping). Given a function  $h : \mathbb{R}^d \mapsto \mathbb{R}$ , the *proximal mapping* (or called *proximal operator*) of  $h$  is the operator given by

$$\mathbf{prox}_h(\mathbf{x}) \triangleq \arg \min_{\mathbf{u} \in \mathbb{R}^d} \left\{ h(\mathbf{u}) + \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|^2 \right\} \text{ for any } \mathbf{x} \in \mathbb{R}^d.$$

# Proximal Gradient

**Definition 2** (proximal mapping). Given a function  $h : \mathbb{R}^d \mapsto \mathbb{R}$ , the *proximal mapping* (or called *proximal operator*) of  $h$  is the operator given by

$$\mathbf{prox}_h(\mathbf{x}) \triangleq \arg \min_{\mathbf{u} \in \mathbb{R}^d} \left\{ h(\mathbf{u}) + \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|^2 \right\} \text{ for any } \mathbf{x} \in \mathbb{R}^d.$$

## Proximal Gradient Method.

$$\begin{aligned} \mathbf{x}_{t+1} &= \mathcal{P}_L^h(\mathbf{x}_t) \triangleq \mathbf{prox}_{\frac{1}{L}h} \left( \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right) \\ &= \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \frac{L}{2} \left\| \mathbf{x} - \left( \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right) \right\|^2 + h(\mathbf{x}) \right\}. \end{aligned}$$

# Proximal Gradient

## Proximal Gradient Method.

$$\begin{aligned}\mathbf{x}_{t+1} &= \mathcal{P}_L^h(\mathbf{x}_t) \triangleq \mathbf{prox}_{\frac{1}{L}h} \left( \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right) \\ &= \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \frac{L}{2} \left\| \mathbf{x} - \left( \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right) \right\|^2 + h(\mathbf{x}) \right\}.\end{aligned}$$

- In LASSO, where  $h(\mathbf{x}) = \|\mathbf{x}\|_1$ ,  $\mathcal{P}_L^h$  is easy to compute and has closed form solution.
- Algorithmically, PG induces famous algorithms for solving LASSO problem, which are called **ISTA** (GD-type) and **FISTA** (Nesterov's AGD-type).

# Convergence of Proximal Gradient

## *Smooth Optimization*

problem:  $\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$

assumption:  $f$  is  $L$ -smooth

$$\text{GD: } \mathbf{x}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$$

$$\text{Convergence: } f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \mathcal{O}\left(\frac{1}{T}\right)$$

## *Smooth **Composite** Optimization*

problem:  $\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) \triangleq f(\mathbf{x}) + h(\mathbf{x})$

assumption:  $f$  is  $L$ -smooth,  $h$  not

$$\text{PG: } \mathbf{x}_{t+1} = \mathcal{P}_L^h(\mathbf{x}_t)$$

$$\text{Convergence: } F(\mathbf{x}_T) - F(\mathbf{x}^*) \leq ?$$

# Convergence of Proximal Gradient

**Theorem 5.** Suppose that  $f$  and  $h$  are convex and  $f$  is  $L$ -smooth. Setting the parameters properly, Proximal Gradient (PG) enjoys

$$F(\mathbf{x}_T) - F(\mathbf{x}^*) \leq \frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2(T-1)} = \mathcal{O}\left(\frac{1}{T}\right)$$

Proximal gradient can also achieve an  $\mathcal{O}(1/T)$  convergence rate, which is the **same** as the non-composite optimization counterpart.

The result can be further boosted to  $\mathcal{O}(\exp(-T/\kappa))$  when the function  $f$  is  **$\sigma$ -strongly convex** (where  $\kappa = L/\sigma$  is the condition number).

# Convergence of Proximal Gradient

- Generalized one-step improvement lemma on  $F \triangleq f + h$

**Lemma 7.** Suppose that  $f$  and  $h$  are convex and  $f$  is  $L$ -smooth. Let  $\mathbf{x}_{t+1} = \mathcal{P}_L^h(\mathbf{x}_t)$  and  $g(\mathbf{x}) \triangleq L(\mathbf{x} - \mathbf{x}_{t+1})$ . Then for any  $\mathbf{u} \in \mathcal{X}$ ,

$$F(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq \langle g(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|g(\mathbf{x}_t)\|^2.$$

Suppose the above lemma holds for a moment, we now prove the  $\mathcal{O}(1/T)$  convergence rate of **PG**.

# Proof of PG Convergence

**Proof:**

Setting  $\mathbf{u} = \mathbf{x}^*$  in Lemma 7:

**Lemma 7.** Suppose that  $f$  and  $h$  are convex and  $f$  is  $L$ -smooth. Let  $\mathbf{x}_{t+1} = \mathcal{P}_L^h(\mathbf{x}_t)$  and  $g(\mathbf{x}) \triangleq L(\mathbf{x} - \mathbf{x}_{t+1})$ . Then for any  $\mathbf{u} \in \mathcal{X}$ ,

$$F(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq \langle g(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|g(\mathbf{x}_t)\|^2.$$

$$F(\mathbf{x}_{t+1}) - F(\mathbf{x}^*) \leq \langle g(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle - \frac{1}{2L} \|g(\mathbf{x}_t)\|^2$$

$$\begin{aligned} \Rightarrow F(\mathbf{x}_{t+1}) - F(\mathbf{x}^*) &\leq L \langle \mathbf{x}_t - \mathbf{x}_{t+1}, \mathbf{x}_t - \mathbf{x}^* \rangle - \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \quad (g(\mathbf{x}_t) \triangleq L(\mathbf{x}_t - \mathbf{x}_{t+1})) \\ &= \frac{L}{2} (2 \langle \mathbf{x}_t - \mathbf{x}_{t+1}, \mathbf{x}_t - \mathbf{x}^* \rangle - \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2) \\ &= \frac{L}{2} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \quad (2 \langle \mathbf{a}, \mathbf{b} \rangle - \|\mathbf{a}\|^2 = \|\mathbf{b}\|^2 - \|\mathbf{b} - \mathbf{a}\|^2) \end{aligned}$$

$$\Rightarrow \sum_{t=1}^{T-1} F(\mathbf{x}_t) - (T-1)F(\mathbf{x}^*) \leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

# Proof of PG Convergence

*Proof:*

$$\Rightarrow \frac{1}{T-1} \sum_{t=1}^{T-1} F(\mathbf{x}_t) - F(\mathbf{x}^*) \leq \frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2(T-1)}$$

which already gives an  $\mathcal{O}(1/T)$  convergence rate of  $\bar{\mathbf{x}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$ .

*What we want:*  $F(\mathbf{x}_T) - F(\mathbf{x}^*)$

*Next step:* analyzing  $F(\mathbf{x}_T) - \frac{1}{T-1} \sum_{t=1}^{T-1} F(\mathbf{x}_t)$ .

Setting  $\mathbf{u} = \mathbf{x}_t$  in Lemma 7:  $F(\mathbf{x}_{t+1}) - F(\mathbf{x}_t) \leq -\frac{1}{2L} \|g(\mathbf{x}_t)\|^2 \leq 0$ .

$$\Rightarrow \sum_{t=1}^T t(F(\mathbf{x}_{t+1}) - F(\mathbf{x}_t)) \leq 0$$

# Proof of PG Convergence

*Proof:*

*What we want:*  $F(\mathbf{x}_T) - F(\mathbf{x}^*) \Rightarrow$  *Next step:* analyzing  $F(\mathbf{x}_T) - \frac{1}{T-1} \sum_{t=1}^{T-1} F(\mathbf{x}_t)$ .

$$\begin{aligned} \sum_{t=1}^{T-1} t(F(\mathbf{x}_{t+1}) - F(\mathbf{x}_t)) &= \sum_{t=1}^{T-1} t(F(\mathbf{x}_{t+1}) - F(\mathbf{x}_t)) + F(\mathbf{x}_t) - F(\mathbf{x}_t) \\ &= \sum_{t=1}^{T-1} \left( tF(\mathbf{x}_{t+1}) - (t-1)F(\mathbf{x}_t) \right) - \sum_{t=1}^{T-1} F(\mathbf{x}_t) = (T-1)F(\mathbf{x}_T) - \sum_{t=1}^{T-1} F(\mathbf{x}_t) \leq 0 \end{aligned}$$

*What we have:*

$$\begin{aligned} - F(\mathbf{x}_T) - \frac{1}{T-1} \sum_{t=1}^{T-1} F(\mathbf{x}_t) &\leq 0 \\ - \frac{1}{T-1} \sum_{t=1}^{T-1} F(\mathbf{x}_t) - F(\mathbf{x}^*) &\leq \frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2(T-1)} \end{aligned} \quad \Rightarrow \quad F(\mathbf{x}_T) - F(\mathbf{x}^*) \leq \frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2(T-1)} \quad \square$$

# Proof of One-Step Improvement Lemma

**Lemma 7.** Suppose that  $f$  and  $h$  are convex and  $f$  is  $L$ -smooth. Let  $\mathbf{x}_{t+1} = \mathcal{P}_L^h(\mathbf{x}_t)$  and  $g(\mathbf{x}_t) \triangleq L(\mathbf{x}_t - \mathbf{x}_{t+1})$ . Then for any  $\mathbf{u} \in \mathcal{X}$ ,

$$F(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq \langle g(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|g(\mathbf{x}_t)\|^2.$$

**Proof:** *What we have:*  $F(\mathbf{x}) \leq Q(\mathbf{x}; \mathbf{y})$  for any  $\mathbf{y} \in \mathcal{X} \Rightarrow F(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq Q(\mathbf{x}_{t+1}; \mathbf{x}_t) - F(\mathbf{u})$   
*analyzing this quantity*

$$\begin{cases} F(\mathbf{u}) = f(\mathbf{u}) + h(\mathbf{u}) \geq \cancel{f(\mathbf{x}_t)} + \langle \nabla f(\mathbf{x}_t), \mathbf{u} - \mathbf{x}_t \rangle + \cancel{h(\mathbf{x}_{t+1})} + \langle \nabla h(\mathbf{x}_{t+1}), \mathbf{u} - \mathbf{x}_{t+1} \rangle & (\text{convexity}) \\ Q(\mathbf{x}_{t+1}; \mathbf{x}_t) = \cancel{f(\mathbf{x}_t)} + \langle \nabla f(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|_2^2 + \cancel{h(\mathbf{x}_{t+1})} \end{cases}$$

$$\Rightarrow Q(\mathbf{x}_{t+1}; \mathbf{x}_t) - F(\mathbf{u}) \leq \underbrace{\langle \nabla f(\mathbf{x}_t) + \nabla h(\mathbf{x}_{t+1}), \mathbf{x}_{t+1} - \mathbf{u} \rangle + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|_2^2}_{= \frac{1}{2L} \|g(\mathbf{x}_t)\|^2 \text{ } (g(\mathbf{x}_t) \triangleq L(\mathbf{x}_t - \mathbf{x}_{t+1}))}$$

*Next step:* relate  $\nabla f(\mathbf{x}_t) + \nabla h(\mathbf{x}_{t+1})$  to  $g(\mathbf{x}_t)$ .

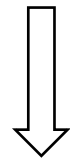
# Proof of One-Step Improvement Lemma

*Proof:*

*What we have:*  $F(\mathbf{x}) \leq Q(\mathbf{x}; \mathbf{y})$  for any  $\mathbf{y} \in \mathcal{X} \Rightarrow F(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq Q(\mathbf{x}_{t+1}; \mathbf{x}_t) - F(\mathbf{u})$   
*analyzing this quantity*

$$\Rightarrow Q(\mathbf{x}_{t+1}; \mathbf{x}_t) - F(\mathbf{u}) \leq \langle \nabla f(\mathbf{x}_t) + \nabla h(\mathbf{x}_{t+1}), \mathbf{x}_{t+1} - \mathbf{u} \rangle + \frac{1}{2L} \|g(\mathbf{x}_t)\|^2$$

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x}} \left\{ \underbrace{h(\mathbf{x}) + \frac{L}{2} \left\| \mathbf{x} - \left( \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right) \right\|^2}_{\triangleq H(\mathbf{x})} \right\}$$



by *Fermat's optimality condition*

**Theorem 8** (Fermat's Optimality Condition). Let  $f : \mathbb{R}^d \rightarrow (-\infty, \infty]$  be a proper convex function. Then

$$\mathbf{x}^* \in \operatorname{argmin}\{f(\mathbf{x}) \mid \mathbf{x} \in \mathbb{R}^d\}$$

if and only if  $\mathbf{0} \in \partial f(\mathbf{x}^*)$ .

$$\mathbf{0} = \nabla H(\mathbf{x}_{t+1}) = \nabla h(\mathbf{x}_{t+1}) + L(\mathbf{x}_{t+1} - \mathbf{x}_t) + \nabla f(\mathbf{x}_t)$$

# Proof of One-Step Improvement Lemma

*Proof:*

*What we have:*  $F(\mathbf{x}) \leq Q(\mathbf{x}; \mathbf{y})$  for any  $\mathbf{y} \in \mathcal{X} \Rightarrow F(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq Q(\mathbf{x}_{t+1}; \mathbf{x}_t) - F(\mathbf{u})$   
*analyzing this quantity*

$$\Rightarrow Q(\mathbf{x}_{t+1}; \mathbf{x}_t) - F(\mathbf{u}) \leq \langle \nabla f(\mathbf{x}_t) + \nabla h(\mathbf{x}_{t+1}), \mathbf{x}_{t+1} - \mathbf{u} \rangle + \frac{1}{2L} \|g(\mathbf{x}_t)\|^2$$

$$\Rightarrow \mathbf{0} = \nabla H(\mathbf{x}_{t+1}) = \nabla h(\mathbf{x}_{t+1}) + L(\mathbf{x}_{t+1} - \mathbf{x}_t) + \nabla f(\mathbf{x}_t)$$

$$\Rightarrow g(\mathbf{x}_t) = L(\mathbf{x}_t - \mathbf{x}_{t+1}) = \nabla h(\mathbf{x}_{t+1}) + \nabla f(\mathbf{x}_t)$$

$$\begin{aligned} \text{Finally we have } Q(\mathbf{x}_{t+1}; \mathbf{x}_t) - F(\mathbf{u}) &\leq \langle g(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{u} \rangle + \frac{1}{2L} \|g(\mathbf{x}_t)\|^2 \\ &= \langle g(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|g(\mathbf{x}_t)\|^2 \end{aligned} \quad \square$$

# One-Step Improvement Lemma

- A *fundamental* result for GD of smoothed optimization.

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2$$

$$f(\mathbf{x}_{t+1}) - f(\mathbf{u}) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2$$

$$f(\mathbf{x}_{t+1}) - f(\mathbf{u}) \leq \langle g(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|g(\mathbf{x}_t)\|^2$$

$$F(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq \langle g(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|g(\mathbf{x}_t)\|^2$$

*specialized*

*general*

*Corollary:* the proof of **PG** can also be used to prove the  $\mathcal{O}(1/T)$  convergence rate of GD.

# Accelerated Proximal Gradient Method

- A natural idea

Can we extend the Nesterov's AGD to the composite optimization?

⇒ This induces the Accelerated Proximal Gradient (**APG**) method.

## Nesterov's Accelerated GD

$$\mathbf{y}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t), \quad \mathbf{x}_{t+1} = (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t$$

## Accelerated Proximal Gradient

$$\mathbf{y}_{t+1} = \mathcal{P}_L^h(\mathbf{x}_t), \quad \mathbf{x}_{t+1} = (1 - \alpha_t) \mathbf{y}_{t+1} + \alpha_t \mathbf{y}_t$$

The convergence rates can be similarly obtained. *Proofs are omitted.*

# Accelerated Proximal Gradient Method

**Theorem 6.** Suppose that  $f$  and  $h$  are convex and  $f$  is  $L$ -smooth. Setting the parameters properly, APG enjoys

$$F(\mathbf{x}_T) - F(\mathbf{x}^*) \leq \frac{2L}{(T+1)^2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Suppose that  $h$  is convex and  $f$  is  $\sigma$ -strongly convex and  $L$ -smooth. Setting the parameters properly, APG enjoys

$$F(\mathbf{x}_T) - F(\mathbf{x}^*) \leq \exp\left(-\frac{T}{\sqrt{\kappa}}\right) \left(F(\mathbf{x}_0) - F(\mathbf{x}^*) + \frac{\sigma}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2\right),$$

where  $\kappa \triangleq L/\sigma$  denotes the condition number.

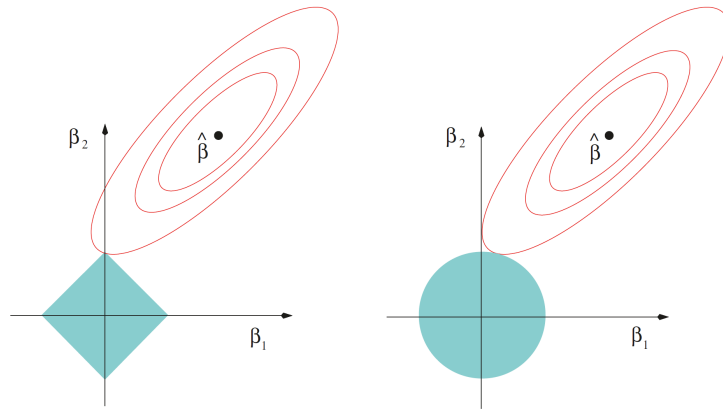
*The convergence rates can be obtained same as those in non-composite optimization.*

# Application to LASSO

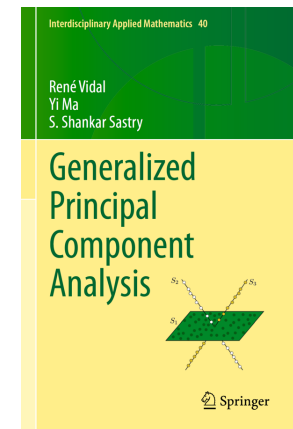
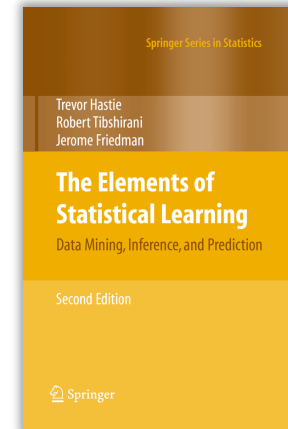
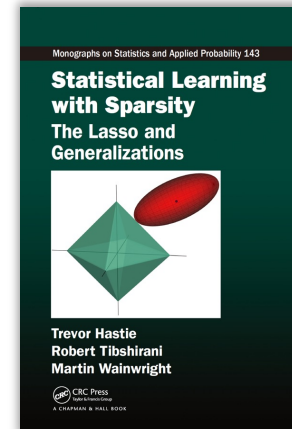
- **LASSO:**  $\ell_1$ -regularized least squares

$$F(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}^\top X - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1$$

commonly encountered in  
*signal/image processing.*



**FIGURE 3.11.** Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions  $|\beta_1| + |\beta_2| \leq t$  and  $\beta_1^2 + \beta_2^2 \leq t^2$ , respectively, while the red ellipses are the contours of the least squares error function.



Regression shrinkage and selection via the lasso

R Tibshirani

Journal of the Royal Statistical Society. Series B (Methodological), 267-288

47812

1996

# Application to LASSO

- **LASSO:**  $\ell_1$ -regularized least squares

$$F(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}^\top X - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1$$

commonly encountered in *signal/image processing*.

⇒ **composite optimization:** first part is *smooth*, the other one is *non-smooth*

- **ISTA** (Iterative Shrinkage-Thresholding Algorithm): **PG** for LASSO
- **FISTA** (Fast ISTA): **APG** for LASSO

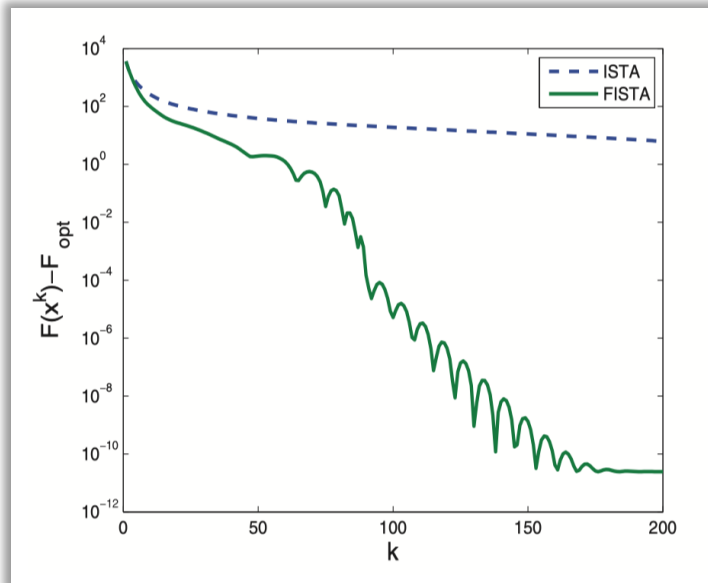
Closed-form solution:

$$(x_+ \triangleq \max\{x, 0\})$$

$$[\mathcal{P}_L^h(\mathbf{w}_t)]_i = \text{sign} \left( \left[ \mathbf{w}_t - \frac{1}{L} \nabla f(\mathbf{w}_t) \right]_i \right) \left( \left| \left[ \mathbf{w}_t - \frac{1}{L} \nabla f(\mathbf{w}_t) \right]_i \right| - \frac{\lambda}{L} \right)_+$$

# Application to LASSO

- Comparison of ISTA and FISTA



Comparison of ISTA and FISTA.

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  $l_1$  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

**1. Introduction.** Linear inverse problems arise in a wide range of applications such as astrophysics, signal and image processing, statistical inference, and optics, to name just a few. The interdisciplinary nature of inverse problems is evident through a vast literature which includes a large body of mathematical and algorithmic developments; see, for instance, the monograph [13] and the references therein.

A basic linear inverse problem leads us to study a discrete linear system of the form

$$(1.1) \quad \mathbf{Ax} = \mathbf{b} + \mathbf{w},$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$  are known,  $\mathbf{w}$  is an unknown noise (or perturbation) vector, and  $\mathbf{x}$  is the “true” and unknown signal/image to be estimated. In image blurring problems, for example,  $\mathbf{b} \in \mathbb{R}^m$  represents the blurred image, and  $\mathbf{x} \in \mathbb{R}^n$  is the unknown true image, whose size is assumed to be the same as that of  $\mathbf{b}$  (that is,  $m = n$ ). Both  $\mathbf{b}$  and  $\mathbf{x}$  are formed by stacking the columns of their corresponding two-dimensional images. In these applications, the matrix  $\mathbf{A}$  describes the blur operator, which in the case of spatially invariant blurs represents a two-dimensional convolution operator. The problem of estimating  $\mathbf{x}$  from the observed blurred and noisy image  $\mathbf{b}$  is called an *image deblurring* problem.

\*Received by the editors February 25, 2008; accepted for publication (in revised form) October 23, 2008; published electronically March 4, 2009. This research was partially supported by the Israel Science Foundation, ISF grant 489-06.

<http://www.siam.org/journals/siims/2-1/71654.html>

<sup>†</sup>Department of Industrial Engineering and Management, Technion-Israel Institute of Technology, Haifa 32000, Israel (becka@ie.technion.ac.il).

<sup>‡</sup>School of Mathematical Sciences, Tel Aviv University, Tel Aviv 69978, Israel (teboulle@post.tau.ac.il).

183

A fast iterative shrinkage-thresholding algorithm for linear inverse problems

A Beck, M Teboulle

SIAM journal on imaging sciences 2 (1), 183-202

11526

2009

# Summary

