



Lecture 4. Gradient Descent Method II

Advanced Optimization (Fall 2024)

Peng Zhao

zhaop@lamda.nju.edu.cn

Nanjing University

Outline

- GD for Smooth Optimization
 - Smooth and Convex Functions
 - Smooth and Strongly Convex Functions
- Momentum and Acceleration
 - Polyak's Momentum
 - Nesterov's Accelerated GD
- Extension to Composite Optimization
 - Proximal Gradient and Accelerated One

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>
	σ -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>
	σ -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:

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

□

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

only exploited convexity, but haven't used smoothness

Refined Result for Smooth Optimization

- Recall the first-order characterization of smooth functions

Smoothness

Theorem 2 (*First-order* Characterizations of L -smoothness). Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a convex function, differentiable over \mathcal{X} . Then the following claims are equivalent:

- (i) f is L -smooth.
- (ii) $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2$ for all $\mathbf{x}, \mathbf{y} \in \mathcal{X}$.
- (iii) $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_*^2$ for all $\mathbf{x}, \mathbf{y} \in \mathcal{X}$.
- (iv) $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{1}{L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_*^2$ for all $\mathbf{x}, \mathbf{y} \in \mathcal{X}$.
- (v) $f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \geq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) - \frac{L}{2} \lambda(1 - \lambda) \|\mathbf{x} - \mathbf{y}\|^2$ for any $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ and $\lambda \in [0, 1]$.

co-coercivity

Proofs can be found below Theorem 5.8 of Amir Beck's book.

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 β -co-coercive (or β -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*.

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

only exploited convexity, but haven't used smoothness

Refined Result for Smooth Optimization

Proof:

$$\begin{aligned}\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 &= \|\Pi_{\mathcal{X}}[\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)] - \mathbf{x}^*\|^2 \quad (\text{GD}) \\ &\leq \|\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) - \mathbf{x}^*\|^2 \quad (\text{Pythagoras Theorem}) \\ &= \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\eta_t \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle + \eta_t^2 \|\nabla f(\mathbf{x}_t)\|^2 \\ &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \left(\eta_t^2 - \frac{2\eta_t}{L} \right) \|\nabla f(\mathbf{x}_t)\|^2 \\ &\quad \text{exploiting coercivity of smoothness and unconstrained first-order optimality} \\ \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle &= \langle \nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}^*), \mathbf{x}_t - \mathbf{x}^* \rangle \geq \frac{1}{L} \|\nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}^*)\|^2 = \frac{1}{L} \|\nabla f(\mathbf{x}_t)\|^2 \\ \Rightarrow \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \left(\eta_t^2 - \frac{2\eta_t}{L} \right) \|\nabla f(\mathbf{x}_t)\|^2 \\ &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 \quad (\text{by picking } \eta_t = \eta = \frac{1}{L} \text{ to minimize the r.h.s}) \\ &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 \leq \dots \leq \|\mathbf{x}_1 - \mathbf{x}^*\|^2\end{aligned}$$

Smooth and Convex

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

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

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

Cautious: This derivation even doesn't require convexity!!

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

Smooth and Convex

Proof:

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

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

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

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

(by optimizer's decreasing property, 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 $\beta \triangleq \frac{1}{2L\|\mathbf{x}_1 - \mathbf{x}^*\|^2}$.

$$\Rightarrow \delta_{t+1} \leq \delta_t - \beta \delta_t^2$$

$$\Rightarrow \frac{1}{\delta_t} \leq \frac{1}{\delta_{t+1}} - \beta \quad (\text{noting that } \delta_t^2 \geq \delta_t \cdot \delta_{t+1} \text{ and then dividing } \delta_t \delta_{t+1} \text{ from both sides})$$

$$\Rightarrow \sum_{t=1}^{T-1} \beta \leq \frac{1}{\delta_T} - \frac{1}{\delta_1} \leq \frac{1}{\delta_T}$$

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

□

Key Lemma for Smooth GD

- During the proof, we have obtained an important lemma for ***smooth*** optimization, that is, ***one-step improvement***

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq \left(-\eta_t + \frac{L}{2}\eta_t^2 \right) \|\nabla f(\mathbf{x}_t)\|^2 \quad \Rightarrow \quad f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \mathcal{O}\left(\frac{1}{T}\right).$$

last-iterated convergence

- 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

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

This lemma usually implies convergence like $f(\bar{\mathbf{x}}_T) - f^* \leq \dots$ with $\bar{\mathbf{x}}_T \triangleq \sum_{t=1}^T \frac{\eta_t \mathbf{x}_t}{\sum_{t=1}^T \eta_t}$ (or other average).

average-iterated convergence

One-Step Improvement Lemma for Smooth GD

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

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

Function progress is proportional to the square of gradient magnitude (*consider due reasons*).

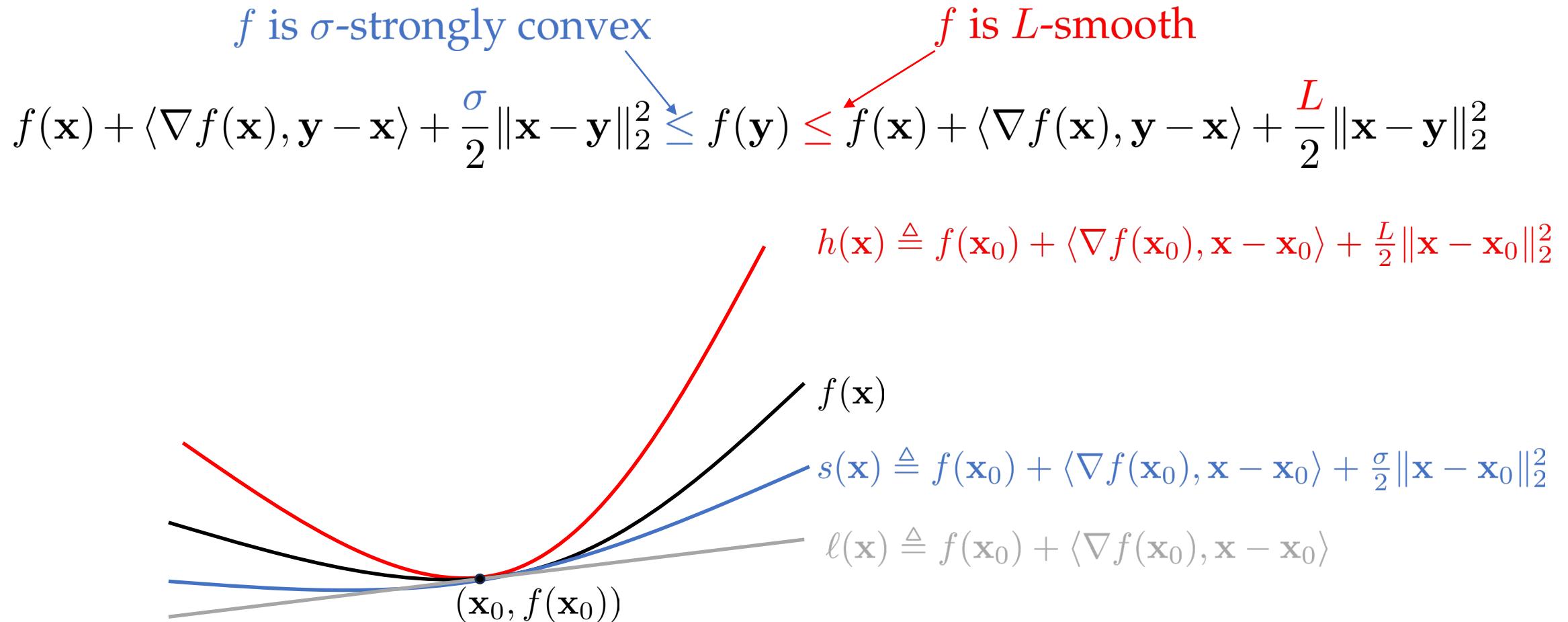
Smooth and Strongly Convex

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

Definition 5 (Strong Convexity). A function f is σ -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



Smooth and Strongly Convex

Theorem 2. Suppose the function $f : \mathbb{R}^d \mapsto \mathbb{R}$ is σ -strongly-convex and differentiable, and also L -smooth. 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:

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

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 σ -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 σ -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 σ -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:

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

exploiting co-coercivity of smooth and strongly convex function

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

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

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.

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

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

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

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

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

□

Constrained Optimization

- For unconstrained optimization, the key technical lemma is

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

where $\nabla f(\mathbf{x})$ is used to measure the function progress.

- For constrained optimization, a *generalized* one-step improvement:

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

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

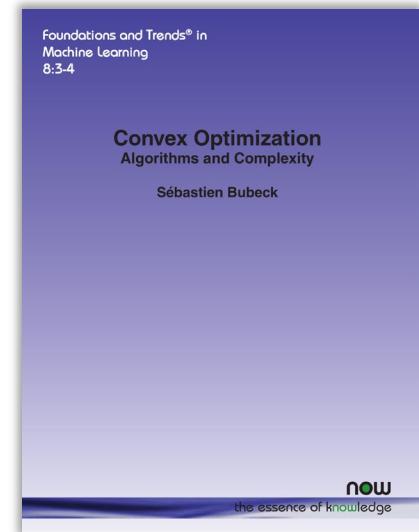
- $g(\mathbf{x}_t)$ is used to qualify the progress; and in the unconstrained case, $g(\mathbf{x}_t) = \nabla f(\mathbf{x}_t)$.
- comparator \mathbf{u} is introduced because (projected) GD is not necessarily “descent”.

Constrained Optimization

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

Detailed proofs for the constrained optimization will not be presented. The proof follows the same vein yet requires some additional 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
Sébastien Bubeck
Foundations and Trends in ML, 2015

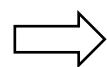
Lower Bound

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	$\mathcal{O}(1/\sqrt{T})$	$\Omega(1/\sqrt{T})$
	σ -strongly convex	$\mathcal{O}(1/T)$	$\Omega(1/T)$
L -smooth	convex	$\mathcal{O}(1/T)$	$\Omega(1/T^2)$
	σ -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!*

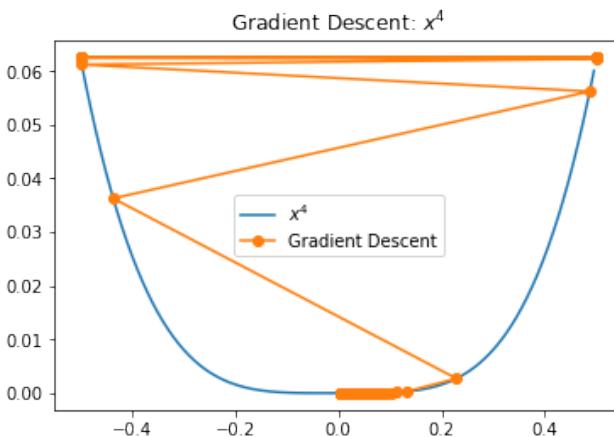
Part 2. Momentum and Acceleration

- Polyak's Momentum
- Nesterov's Accelerated GD
- Smooth and Convex
- Smooth and Strongly Convex

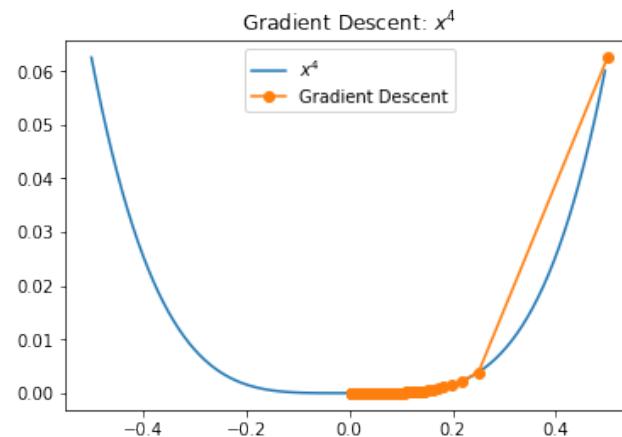
Polyak's Momentum

- GD method (with a fixed step size): $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$, e.g., $\eta = \frac{1}{L}$
- The problem: *pathological curvature*

Consider deploying GD on a quartic function $f(x) = x^4$.



(a) large step size



(b) large step size

Motivation:

- Ensure smaller steps in regions of high curvature to dampen oscillations.
- Ensure larger steps and accelerate in regions of low curvature.

Source: <https://boostedml.com/2020/07/gradient-descent-and-momentum-the-heavy-ball-method.html>

Polyak's Momentum

- GD with momentum:

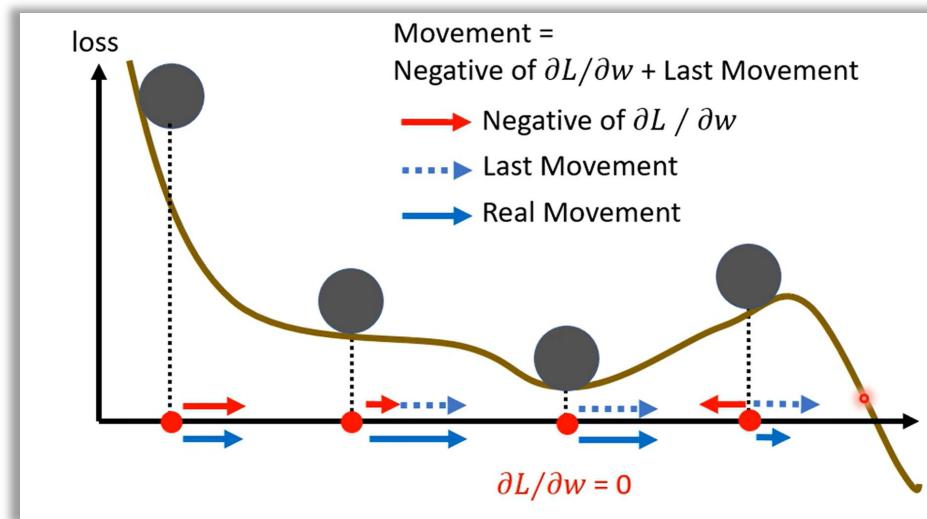
$$\mathbf{x}_{t+1} = \underbrace{\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)}_{\text{GD Update}} + \underbrace{\beta(\mathbf{x}_t - \mathbf{x}_{t-1})}_{\text{momentum}}$$

where β is a hyperparameter (usually $\beta \in [0, 1]$ though not limited to it), which scales down the previous step adaptively.

- If the current gradient step is in the same direction as the previous step (e.g., in the region of low curvature), then move a little further in that direction;
- If it's in the opposite direction (e.g., in the region of high curvature), move less far.
- Also known as the “**heavy ball method**” (think of the physical intuition).

Polyak's Momentum

- Provable benefit: can achieve *accelerated rate* for optimizing the *quadratic functions* (but fail for more general cases like smooth and convex/strongly convex functions). Details are omitted [[more details](#)].
- Other benefit: help jump out of the local region (can be useful for non-convex opt)

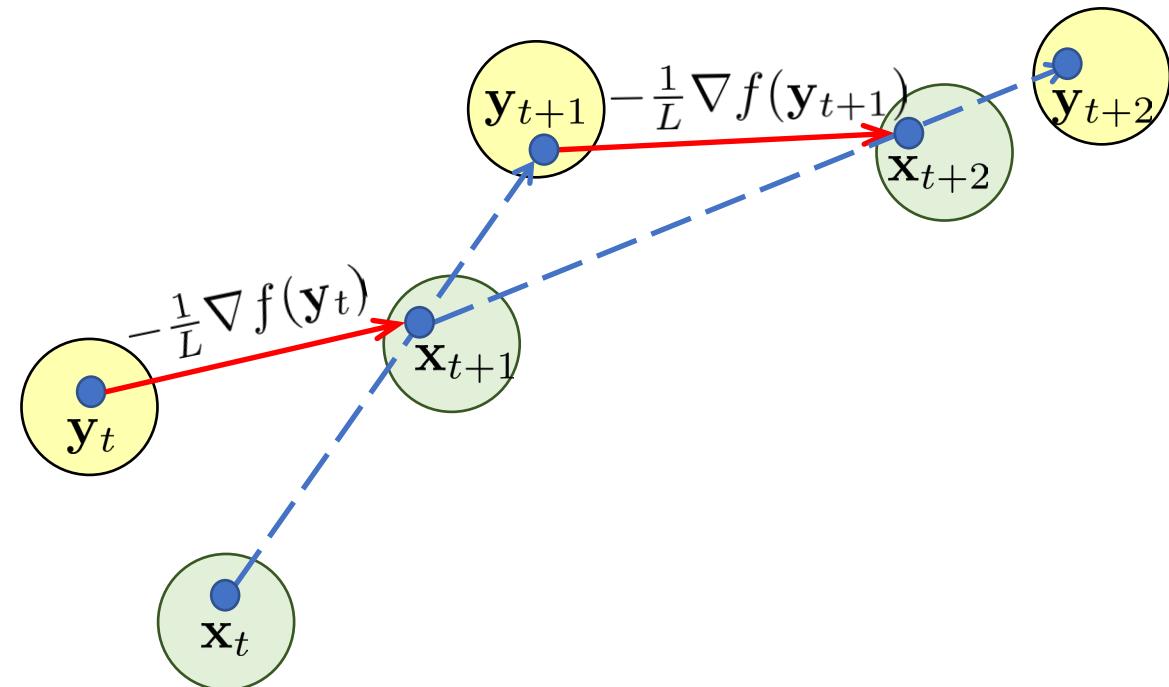


Source: Hung-yi Lee ML 2021 Spring course Lecture on batch and momentum

Nesterov's Accelerated GD

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

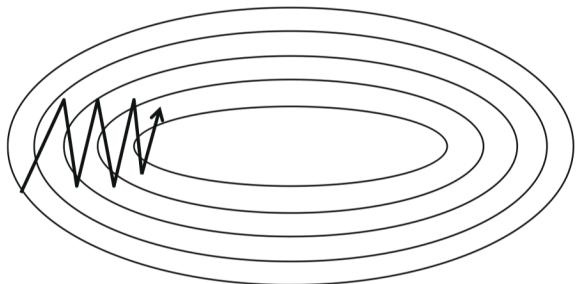
$$\mathbf{y}_{t+1} = \mathbf{x}_{t+1} + \beta_t (\mathbf{x}_{t+1} - \mathbf{x}_t)$$



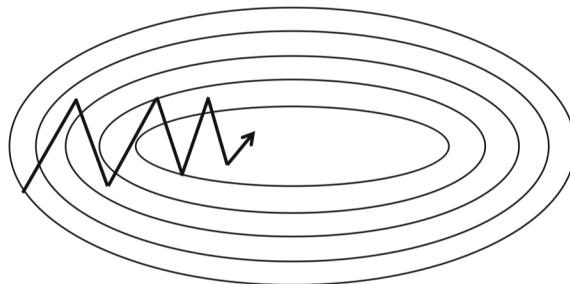
- Define $\mathbf{x}_1 = \mathbf{y}_1$.
- $\beta_t > 0$ is a *time-varying* mixing rate of \mathbf{x}_t and \mathbf{x}_{t+1} ; $\beta_t = 0$ recovers vanilla GD.
- AGD can be also thought a version of GD with *momentum*.

Nesterov's Accelerated GD

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



GD

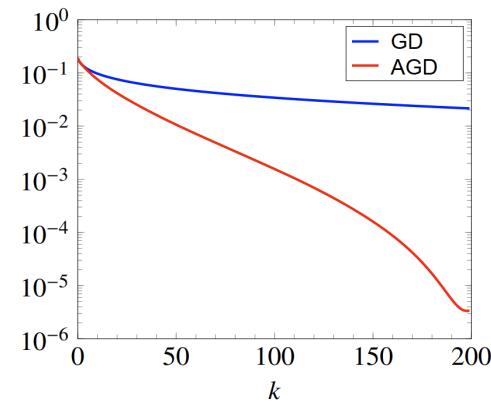


Accelerated GD

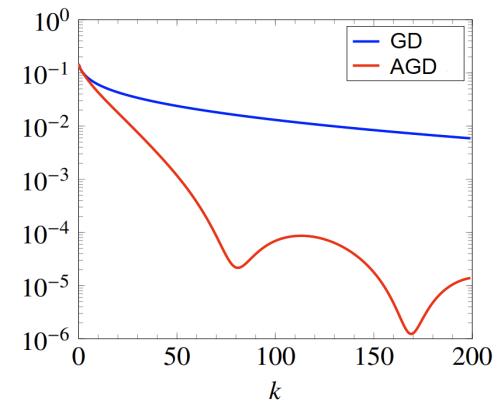
Example

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

- 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>

Polyak's Momentum v.s. Nesterov's AGD

- Polyak's Momentum:

$$\mathbf{x}_{t+1} = \underbrace{\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)}_{\text{GD Update}} + \underbrace{\beta(\mathbf{x}_t - \mathbf{x}_{t-1})}_{\text{momentum}}$$

- Nesterov's AGD:

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

$$\mathbf{x}_{t+1} = \underbrace{\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t + \beta_t(\mathbf{x}_t - \mathbf{x}_{t-1}))}_{\text{GD Update}} + \underbrace{\beta(\mathbf{x}_t - \mathbf{x}_{t-1})}_{\text{momentum}}$$

Main difference: separate *the gradient calculation state* and *the momentum state*.

Convergence of Nesterov's Accelerated GD

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

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

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

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

Note: for simplicity, we are working on *unconstrained*.

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:

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

a comparator variable \mathbf{u} is introduced here,
because now AGD is not necessary “descent” due to the momentum

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

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \quad \text{GD for smooth and convex functions}$$

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

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

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 $\mathbf{u} = \mathbf{x}_t$:

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

(ii) Plugging in $\mathbf{u} = \mathbf{x}^*$:

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

Proof of AGD Convergence

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

Proof: (continued proving Theorem 3)

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

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

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

$$\begin{aligned}& (\lambda_t - 1)(f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t)) + f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \\&= \lambda_t(f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*)) - (\lambda_t - 1)(f(\mathbf{x}_t) - f(\mathbf{x}^*))\end{aligned}$$

Define $\delta_t \triangleq f(\mathbf{x}_t) - f(\mathbf{x}^*)$, then we have

$$\text{LHS} = \lambda_t \delta_{t+1} - (\lambda_t - 1) \delta_t$$

Goal: design a telescoping series

Proof of AGD Convergence

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

Proof: (continued proving Theorem 3)

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

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

RHS of $(\lambda_t - 1)(\text{i}) + (\text{ii})$ equals:

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

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

That is

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

Proof of AGD Convergence

Proof: (continued proving Theorem 3)

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

Cautious: many terms of interest have already appeared in the following inequality.

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

gradient inner product *optimal point*
_____ _____
optimality gap *linear combination*
telescoping structure *related to momentum* *gradient norm*

Proof of AGD Convergence

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

Proof: (continued proving Theorem 3)

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

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

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

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

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

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

Proof of AGD Convergence

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

Proof: (continued proving Theorem 3)

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

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

Goal: design a telescoping series

Proof of AGD Convergence

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

Proof: (continued proving Theorem 3)

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

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

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

telescope

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

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

Proof of AGD Convergence

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

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): $\lambda_0 = 0$

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

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

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

Proof

Proof: (continued proving Theorem 3)

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

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

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

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

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

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

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

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

Requirement (3): $\lambda_0 = 0$

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

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

Smooth and Strongly Convex

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

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

satisfies

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

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

core technique: estimate sequence (developed by Yurii Nesterov)

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

It can be proved that the above construction satisfies the two properties.

Smooth and Strongly Convex

- Proof sketch

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

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

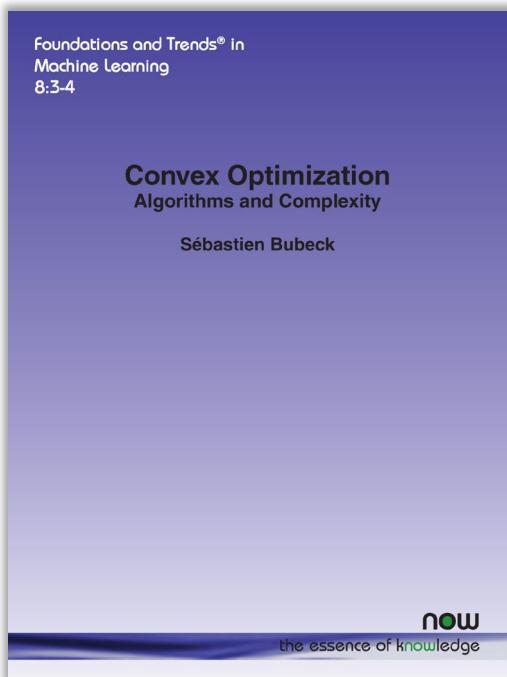
$$\Phi_{t+1}(\mathbf{x}) \triangleq (1 - \theta)\Phi_t(\mathbf{x}) + \theta \left(f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \frac{\sigma}{2} \|\mathbf{x} - \mathbf{x}_t\|^2 \right)$$

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\stackrel{(ii)}{\leq} \min_{\mathbf{x} \in \mathbb{R}^d} \Phi_t(\mathbf{x}) - f(\mathbf{x}^*) \leq \Phi_t(\mathbf{x}^*) - f(\mathbf{x}^*) && \text{(by property (ii))} \\ &\stackrel{(i)}{\leq} (1 - \theta)^t (\Phi_1(\mathbf{x}^*) - f(\mathbf{x}^*)) && \text{(by property (i))} \\ &= (1 - \theta)^t \left(f(\mathbf{x}_1) + \frac{\sigma}{2} \|\mathbf{x}^* - \mathbf{x}_1\|^2 - f(\mathbf{x}^*) \right) && \text{(definition of } \Phi_1) \\ &\lesssim (\sigma + L) \|\mathbf{x}^* - \mathbf{x}_1\|^2 \exp(-\theta t) && \text{(smoothness)} \end{aligned}$$

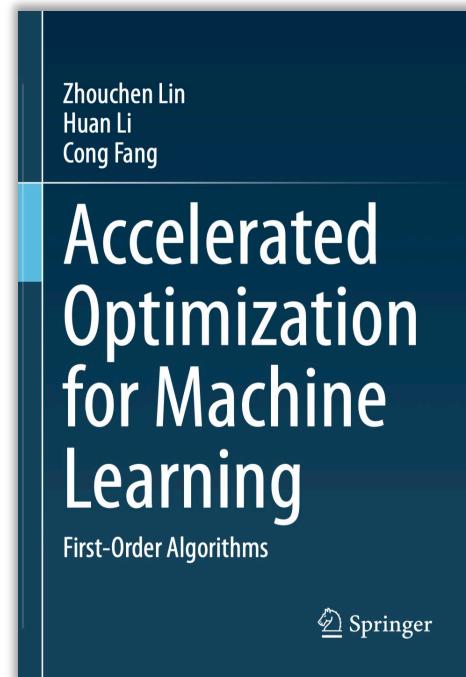
Estimate Sequence

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

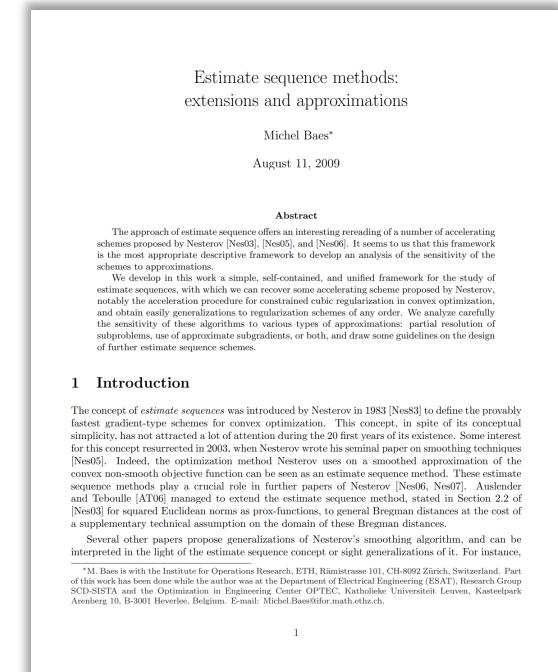
References:



Chapter 3.7



Chapter 2.1



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

More Explanations for Nesterov's AGD

- Ordinary Differentiable Equations
 - Su, W., Boyd, S., & Candes, E. A differential equation for modeling Nesterov's accelerated gradient method: theory and insights. *NIPS* 2014.
 - Even, M., Berthier, R., Bach, F., Flammarion, N., Gaillard, P., Hendrikx, H., Taylor, A. A continuized view on nesterov acceleration for stochastic gradient descent and randomized gossip. *NeurIPS* 2021.
- Variational Analysis
 - Wibisono, A., Wilson, A. C., & Jordan, M. I. A variational perspective on accelerated methods in optimization. *PNAS* 2016, 113(47), E7351-E7358.

More Explanations for Acceleration

- Linear Coupling of GD and MD
 - Allen-Zhu, Z., & Orecchia, L. Linear coupling: An ultimate unification of gradient and mirror descent. *ITCS* 2017.
 - Cutkosky A. Chapter 14 Momentum & Chapter 15 Acceleration. *Lecture Notes for EC525: Optimization for Machine Learning*, 2022.
- Online Learning with Suitable Optimism
 - Kavis, A., Levy, K. Y., Bach, F., & Cevher, V. UnixGrad: A universal, adaptive algorithm with optimal guarantees for constrained optimization. *NeurIPS* 2019.
 - Kreisler, I., Ivgi, M., Hinder, O., & Carmon, Y. Accelerated Parameter-Free Stochastic Optimization. *COLT* 2024.

History Bits

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

UDC 51

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

YU. E. NESTEROV

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 $\{x_k\}_0^\infty$. This property allows us to reduce the amount of computation. At the same time, it is possible to obtain an estimate of convergence rate which is improved for the class of problems under consideration (see [1]).

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 interior, we use the following method.

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

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

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

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

b) Put

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

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

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

The way in which the one-dimensional search (4) is halted is different from that in [2]. The difference is only that in (4) the subdivision in the interval $[x_{k-1}, x_k]$ (and not with 1 as in [2]). In view of this (see the proof of Theorem 1) the sequence $\{x_k\}_0^\infty$ is constructed by method (3)–(5), no more subdivisions will be made. The recalculations of the points y_k in (5) are not needed.

1980 Mathematics Subject Classification. Primary 90C25.

372

Let us also remark that method (3)–(5) does not guarantee the sequences $\{x_k\}_0^\infty$ and $\{y_k\}_0^\infty$.

THEOREM 1. Let $f(x)$ be a convex function in $C^{1,1}(E)$. Then a sequence $\{x_k\}_0^\infty$ 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 ϵ with respect to the functional f , it is necessary and sufficient to have

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

b) to evaluate the objective function no more than NF times.

Here and in what follows, $\lfloor \cdot \rfloor$ is the integer part of \cdot .

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

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

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

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

$$\begin{aligned} \|p_{k+1} - x_{k+1} + x^*\|^2 &= \|p_k - x_k + x^*\|^2 + 2(a_{k+1} - 1) \\ &\quad + 2a_{k+1}\alpha_{k+1}\langle f'(y_{k+1}), x^* \rangle. \end{aligned}$$

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

$$\begin{aligned} \langle f'(y_{k+1}), y_{k+1} - x^* \rangle &\geq f(x_{k+1}) - f^* \\ 0.5\alpha_{k+1}\|f'(y_{k+1})\|^2 &\leq f(y_{k+1}) - f(x_{k+1}) \\ &\quad - a_{k+1}^2\langle f'(y_{k+1}), x^* \rangle. \end{aligned}$$

We substitute these two inequalities into the preceding

$$\begin{aligned} \|p_{k+1} - x_{k+1} + x^*\|^2 - \|p_k - x_k + x^*\|^2 &\leq 2(a_{k+1} - 1) \\ &\quad - 2a_{k+1}\alpha_{k+1}\langle f(x_{k+1}) - f^*, p_{k+1} - a_{k+1}^2\rangle \\ &\leq -2a_{k+1}\alpha_{k+1}\langle f(x_{k+1}) - f^*, p_{k+1} - a_{k+1}^2\rangle \\ &= 2a_{k+1}a_{k+1}^2\langle f(x_k) - f^*, p_{k+1} - a_{k+1}^2\rangle \\ &\leq 2a_k a_k^2\langle f(x_k) - f^*, p_{k+1} - a_{k+1}^2\rangle. \end{aligned}$$

Thus

$$\begin{aligned} 2\alpha_{k+1}a_{k+1}^2\langle f(x_{k+1}) - f^*, p_{k+1} - a_{k+1}^2\rangle &\leq 2a_k a_k^2\langle f(x_k) - f^*, p_{k+1} - a_{k+1}^2\rangle \\ &\leq 2a_k a_k\langle f(x_k) - f^*, p_k - x_k + x^*\rangle \\ &\leq 2a_0 a_0^2\langle f(x_0) - f^*, p_0 - x_0 + x^*\rangle \leq \|y_0 - z\|. \end{aligned}$$

It remains to observe that $a_{k+1} \geq a_k + 0.5 \geq 1 + 0.50$.

It follows from the estimate of the convergence rate that method (3)–(5) needs to achieve accuracy ϵ will be necessary at least once each iteration, one gradient and at least two values of

be calculated. Let us remark, however, that to each additional function corresponds a halving of α_k . Therefore the total number of iterations 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 f , then we can take $\alpha_k \equiv L^{-1}$ in the method (3)–(5) for any $k \geq 0$. In this case the halving of α_k is not necessary.

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

Assume that $f(x) - f^* \geq 0.5m\|x - x^*\|^2$ for all $x \in E$. Then the 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. Then, by (3)–(5), one has $N \leq 4\sqrt{L/m} - 1$. 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 $4\sqrt{L/m} - 1$ iteration the function decreases by a factor of 2. Thus the method (3)–(5) 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 $\tilde{f}(x) = (f_i(x))$, $i = 1, \dots, m$, are continuous differentiable functions on E . The set X is assumed to be nonempty. In addition to this, we will also assume that the functions $\{f_i(\cdot), \tilde{f}(\cdot)\}$ 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 .

As is well known, the identity $F(u) \equiv \max\{\langle \lambda, u \rangle | \lambda \in \partial F(0)\}$ holds for functions that are positive homogeneous of degree one. Then 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 implies the existence of a saddle point (λ^*, x^*) for problem (9). The solution of problem (9) can be written as $\Omega^* = \Lambda^* \times X^*$, where

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

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 $\bar{L} = (L^{(1)}, \dots, L^{(m)})$.

Consider the function

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

where

$$\begin{aligned} \tilde{f}(y, z) &= (f^{(1)}(y, z), \dots, f^{(m)}(y, z)), \\ f^{(k)}(y, z) &= f_k(y) + \langle f'(y), z - y \rangle, \end{aligned}$$

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 n-dimensional case, the function $\Phi^*(y, A)$ is 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

$$(11) \quad k = 0, \quad a_0 = 1, \quad x_{-1} = y_0, \quad A_{-1} = F(y_0),$$

where $\bar{L}_0 = (L_0^{(1)}, \dots, L_0^{(m)})$, $L_0^{(k)} = \|f_k(y_0) - f_k(x_0)\|/\|y_0 - x_0\|$, $z \in E$, $z \neq y_0$.

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

$$(13) \quad \begin{aligned} a_{k+1} &= (1 + \sqrt{4a_k^2 + 1})/2, \\ y_{k+1} &= x_k + (a_k - 1)(x_k - x_{k-1})/a_{k+1}. \end{aligned}$$

It is not hard to see that the method (3)–(5) is simply the method (11)–(13) for the unconstrained minimization problem $\min\{\Phi^*(y, A) | y \in E\}$ in (8).

The author expresses his sincere thanks to V. G. Kostyuk for stimulating his interest in the questions considered here.

Central Economico-Mathematical Institute
Academy of Sciences of the USSR

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.

Received 19/JULY/82

Translated by A. ROSA

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.

УДК 51

Ю.Е. НЕСТЕРОВ

МЕТОД РЕШЕНИЯ ЗАДАЧИ ВЫПУКЛОГО ПРОИЗВОДСТВА СО СКОРОСТЬЮ СХОДИМОСТИ $O(1/k^2)$

(Представлено академиком Л.В. Канторовичем)

1. В статье предлагается метод решения задачи в гильбертовом пространстве E . В отличие от более раннего программирования, предлагавшегося ранее, этот метод не использует последовательность точек $\{x_k\}_{k=0}^\infty$, которая не является особенностю позволяющей свести к минимуму вычислительные затраты. В то же время для такого метода удается получить оценку скорости сходимости.

2. Рассмотрим начальную задачу безусловной минимизации $f(x)$. Мы будем предполагать, что функция $f(x)$ принадлежит классу конечнозначимых и непрерывных функций, что существует константа $L > 0$, для которой при всех $x, y \in E$

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

Из неравенства (1) следует, что при всех $x, y \in E$

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

Для решения задачи $\min\{f(x) | x \in E\}$ с непустым множеством X^* предлагается следующий метод.

0) Выбираем точку $y_0 \in E$. Полагаем

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

где z – любая точка из E , $z \neq y_0$, $f'(z) \neq f'(y_0)$.

1) k -я Итерация.

а) Вычисляем наименьший номер $i \geq 0$, для которого

$$(4) \quad f(y_k) - f(y_{k-1}) - 2^{-i} \alpha_{k-1} f'(y_k) \geq 2^{-i-1} \alpha_{k-1} \|f'(y_k)\|^2.$$

б) Полагаем

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

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

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

Способ прерывания одномерного поиска (4) аналогичен в [2]. Разница лишь в том, что в (4) дробление производится, начиная с a_{k-1} (а не с единицей, как в [2]). Доказательство теоремы 1) при построении методом (3)–(5) будет сделано не более $O(\log_2 L)$. Пересечение с помощью "овражного" шага. Отметим также, что получившееся монотонное убывание функции $f(x)$ на последовательности $\{y_k\}_{k=0}^\infty$.

Теорема 1. Пусть выпуклая функция $f(x) \in E$ последовательность $\{x_k\}_{k=0}^\infty$ построена методом (3)–(5),

$$(6) \quad f(x_k) - f^* \leq C/(k+2)^2,$$

$$\text{где } C = 4L \|y_0 - x^*\|^2, \quad f^* = f(x^*), \quad x^* \in X^*;$$

2) для достижения точности ϵ по функционалу необходимо

а) вычислить градиент целевой функции не более NG

б) вычислить значение целевой функции не

$$+ \lceil \log_2(2L\alpha_{-1}) \rceil + 1 \text{ раз.}$$

Здесь и далее $\lceil \cdot \rceil$ – целая часть числа (\cdot) .

Доказательство. Пусть $y_k(a) = y_k - af'(y_k)$

получаем $f(y_k) - f(y_k(a)) \geq 0.5a(2 - aL) \|f'(y_k)\|^2$.

$$C - 2^{-i} \alpha_{k-1} \geq 0.5m \|x_k - x^*\|^2, \quad \text{где } m > 0,$$

и пусть константа n .

Введем в метод (3)–(5) следующее правило прерывания:

$$(7) \quad k \geq 2\sqrt{2/(m\alpha_{-1})} + 1 \text{ раз.}$$

Пусть прерывание произошло на N -м шаге. Так как

$$> 0.5L^{-1}, \quad \text{то } N \leq \lceil 4\sqrt{L/m} \rceil - 1.$$

В то же время

В заключение этого раздела покажем, как можно решить задачу минимизации сильно выпуклой функции $f(x)$ при всех $x \in E$. Предположим, что для функции $f(x)$ в пределах E выполнено условие $f(x) \geq 0.5m \|x - x^*\|^2$, где $m > 0$, и пусть константа n . Введем в метод (3)–(5) следующее правило прерывания: Останавливаемся, если

$$(7) \quad k \geq 2\sqrt{2/(m\alpha_{-1})} + 1 \text{ раз.}$$

Пусть прерывание произошло на N -м шаге. Так как $\alpha_{-1} > 0.5L^{-1}$, то $N \leq \lceil 4\sqrt{L/m} \rceil - 1$. В то же время

$$f(x_N) - f^* \leq \frac{\|y_0 - x^*\|^2}{\alpha_N(N+2)^2} \leq 0.25m \|y_0 - x^*\|^2 \leq 0$$

После того как получена точка x_N , необходимо оценить с помощью методом (3)–(5), (7) из точек x_N как из начальной точки x_0 .

В результате получаем, что за каждые $\lceil 4\sqrt{L/m} \rceil$ итерации убывает вдвое. Таким образом, метод (3)–(5) является неупущающим (с точностью до безразмерной константы) на классе сильно выпуклых функций из $C^{1,1}$.

3. Рассмотрим следующую экстремальную задачу:

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

где Q – выпуклое замкнутое множество из E , $F(u)$, $u \in Q$ – положительно-однородная степень единицы функция

$\dots, f_m(x)$ – вектор выпуклых непрерывно дифференцируемых функций X^* решений задачи (8) всегда предполагает, что система функций

имеет следующим свойством:

(*) Если существует вектор $\lambda \in \partial F(0)$ такой, что

ненайда функция.

Через $\partial F(0)$ в (*) обозначен субдифференциал функции

Как известно, для выпуклых положительно-однородных функций справедливо тождество $F(u) = \max\{\langle \lambda, u \rangle | \lambda \in \partial F(0)\}$.

Предположим (*) следует выпуклость функции $F(\bar{x})$.

Задачу (8) можно записать в минимаксной форме:

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

Можно показать, что из непустоты множества X^* и предположения о задаче (9) седловой точки $(\bar{\lambda}^*, \bar{x}^*)$. Поэтому задача (9) представимо в виде $\Omega^* = \bar{\lambda}^* \times X^*$, где

$\bar{\lambda} \in \partial F(0)$, $\Psi(\bar{\lambda}) = \min\{\langle \lambda, f(x) \rangle | x \in Q\}$. Задачу

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

мы будем называть задачей, двойственной к (8).

Пусть в задаче (8) функции $f_k(x)$, $k = 1, 2, \dots, C^{1,1}(E)$ с константами $L^{(k)} \geq 0$. Обозначим $\bar{L} = (L^{(1)}, L^{(2)}, \dots, L^{(m)})$.

Рассмотрим функцию $\Phi(y, A, z) = F(\bar{f}(y, z)) + C$

$$= (f^{(1)}(y, z), f^{(2)}(y, z), \dots, f^{(m)}(y, z)), f^{(k)}(y, z) = f_k(z),$$

\dots, m, A – положительная константа. Обозначим

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

3. 174

Отметим, что отображение $y \rightarrow T(y, A)$ является естественным для решения задачи "градиентного" отображения, введенного в [1] методом минимизации функций вида $\max_{1 \leq k \leq m} f_k(x)$. Для (8) предположим, что для функции $f(x)$ при всех $x \in E$ выполнено условие $f(x) \geq 0.5m \|x - x^*\|^2$, где $m > 0$, и пусть константа n . Введем в метод (3)–(5) следующее правило прерывания:

$$(10) \quad \Phi^*(y, A) + A(y - T(y, A), x - y) + 0.5A \|y - T(y, A)\|^2 \leq 0$$

причем если $A \geq F(L)$, то

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

Для решения задачи (8) предлагается следующий метод:

0) Выбираем точку $y_0 \in E$. Полагаем

$$(11) \quad k = 0, \quad a_0 = 1, \quad x_{-1} = y_0, \quad A_0 = F(\bar{L}_0),$$

где $\bar{L}_0 = (L_0^{(1)}, L_0^{(2)}, \dots, L_0^{(m)})$, $L_0^{(k)} = \|f'_k(y_0) - f'_k(x^*)\|$, точка из E , $z \neq y_0$.

1) k -я Итерация.

а) Вычисляем наименьший номер $i \geq 0$, для которого равенство

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

б) Полагаем $A_k = 2^i A_{k-1}$, $x_k = T(y_k, A_k)$,

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

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

Нетрудно заметить, что метод (3)–(5) является записью метода (11)–(13) для задачи безусловной минимизации $m = 1$, $F(y) = y$, $Q = E$.

Теорема 2. Если последовательность $\{x_k\}_{k=0}^\infty$

$$(13), \text{ то:}$$

1) для любого $k \geq 0$ $F(\bar{f}(x_k)) - F(\bar{f}(x^*))$

$$= 4F(\bar{L}) \|y_0 - x^*\|^2, \quad x^* \in X^*$$

2) для достижения точности ϵ по функционалу необ

ходимо:

а) решить вспомогательную задачу $\min\{\Phi(y, z) | \sqrt{C_1/\epsilon} + 1 \leq \max\{\log_2(F(\bar{L})/A_{k-1}), 0\} | z\}$,

б) вычислить набор градиентов $f'_1(y), f'_2(y), \dots, f'_m(y)$,

$$\sqrt{C_1/\epsilon} | z |$$

в) вычислить вектор-функцию $\bar{f}(x)$ не более $2\sqrt{C_1/\epsilon}$.

Теорема 2 доказывается практически так же, как только вместо неравенства (2) использовать неравенство вектора $\alpha_k f'(y_k)$ будет вектор $y_0 - T(y_k, A_k)$, а ана

логично так же, как и в методе (3)–(5), в методе информацией о константе $F(\bar{L})$ и параметре сильной выпуклости α_k . Поэтому, чтобы $y_0 \in Q$.

В заключение отметим два важных частных случаев вспомогательной задачи $\min\{\Phi(y, A, x) | z \in Q\}$.

а) Минимизация гладкой выпуклой функции на простом множестве мы понимаем такое множество, для которого

дифференцирования записывается в явном виде. В этом случае в

и в методе (11)–(13)

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

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

б) Безусловная минимизация (в задаче (8) $Q \equiv E$). В этом случае вспомогательная задача $\min\{\Phi(y, A, x) | x \in E\}$ эквивалентна следующей двойственной задаче:

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

При этом $T(y, A) = y - A^{-1} \sum_{k=1}^m \lambda^{(k)} f'_k(y)$, $\lambda^{(k)} = \min\{\lambda | \Phi(y, A, x) \leq 0 | x \in E\}$.

шения задачи (14) при физически заданных простыми случаях задача (14)

Автор выражает

признательность за интерес к

разработке и выражает

аппринципии

Центральный экономико-математический институт

Академия наук СССР, Москва

ЛИТЕРАТУРА

1. Немировский А.С., Юдин Д.Б. Сложность задач и эффективность методов оптимизации. М.: Наука, 1979. 2. Шпеничный Б.Н., Данилин Ю.М. Численные методы в экстремальных задачах. М.: Наука, 1975.

УДК 515.1

Е.И. НОЧКА

К ТЕОРИИ МЕРОМОРФНЫХ КРИВЫХ

(Представлено академиком В.С. Владимировым 18 V 1982)

1. Пусть задана мероморфная кривая, т.е. мероморфное отображение

$$f: C \rightarrow \mathbb{C}^n, \quad f = (f_1, f_2, \dots, f_n)$$

является редуцированным представлением кривой \tilde{f} . Характеристическую функцию \tilde{f} определим, следуя А. Картану [1]:

$$T(\tilde{f}, r) = \frac{1}{2\pi} \int_0^{2\pi} |\log|f(re^{i\gamma})|^2 d\gamma - \log|f(0)|^2.$$

Пусть A – гиперплоскость в \mathbb{C}^n и a – единичный вектор такой, что равенство $\langle a, A \rangle = 0$ (скобки обозначают эрмитово скалярное произведение) есть уравнение гиперплоскости A в однородных координатах; обозначим $f_A = (f, a)$.

МАТЕМАТИКА

Received 19/JULY/82

Поступило 19 VII 1982

History Bits

- Polyak's Momentum, credit goes to Polyak, date back to 1960s

B. T. Polyak. Some methods of speeding up the convergence of iteration methods. USSR Computational Mathematics and Mathematical Physics, 4(5):1–17, 1964.



Boris T. Polyak
1935-2023

Math. Program., Ser. B 91: 401–416 (2002)

Digital Object Identifier (DOI) 10.1007/s101070100258

B.T. Polyak

History of mathematical programming in the USSR: analyzing the phenomenon*

Received: January 29, 2001 / Accepted: May 17, 2001
Published online October 2, 2001 – © Springer-Verlag 2001

Abstract. I am not a historian; these are just reminiscences of a person involved in the development of optimization theory and methods in the former USSR. I realize that my point of view may be very personal; however, I am trying to present as broad and unbiased picture as I can.

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

How to effectively leverage the (partial) smoothness to improve convergence?

Recall Non-composite Optimization

Recall how we ***invent*** GD for unconstrained non-composite optimization.

- **Idea: surrogate optimization**

We aim to find a sequence of *local upper bounds* U_1, \dots, U_T , where the surrogate function $U_t : \mathbb{R}^d \mapsto \mathbb{R}$ may depend on \mathbf{x}_t such that

(i) $f(\mathbf{x}_t) = U_t(\mathbf{x}_t);$

(ii) $f(\mathbf{x}) \leq U_t(\mathbf{x})$ holds for all $\mathbf{x} \in \mathbb{R}^d;$

(iii) $U_t(\mathbf{x})$ should be simple enough to minimize.



Then, our proposed algorithm would be $\mathbf{x}_{t+1} = \arg \min_{\mathbf{x}} U_t(\mathbf{x})$

Recall Non-composite Optimization

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

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

→ to minimize $f(\mathbf{x})$, it suffices to minimize the *surrogate* sequence $\{U_t(\mathbf{x})\}_{t=1}^T$.

Claim. GD for smooth functions can be equivalently represented by

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

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

Recall Non-composite Optimization

Claim. GD for smooth functions can be equivalently represented by

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

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

Proof:

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

□

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*.

A natural idea for surrogate objective:

Following previous argument (for non-composite optimization), to minimize $F \triangleq f + h$, it's natural to optimize surrogate sequence $\{U_t(\mathbf{x})\}_{t=1}^T$ defined as

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

Composite Optimization

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

⇒ to minimize $F(\mathbf{x}) = f(\mathbf{x}) + h(\mathbf{x})$, it suffices to minimize $U_t(\mathbf{x}) \triangleq u_t(\mathbf{x}) + h(\mathbf{x})$.

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

Composite Optimization

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

surrogate objective

⇒ to minimize $F(\mathbf{x}) = f(\mathbf{x}) + h(\mathbf{x})$, it suffices to minimize $U_t(\mathbf{x}) \triangleq u_t(\mathbf{x}) + h(\mathbf{x})$.

$$\arg \min_{\mathbf{x}} U_t(\mathbf{x}) = \arg \min_{\mathbf{x}} \left\{ \frac{L}{2} \left(-2 \left\langle \mathbf{x}_t - \frac{\nabla f(\mathbf{x}_t)}{L}, \mathbf{x} \right\rangle + \|\mathbf{x}\|^2 \right) + h(\mathbf{x}) \right\}$$

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

this will be abstracted as an operator, a subproblem to optimize

Composite Optimization

- Iteratively solve the surrogate optimization problem.

Deploying the following update rule:

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

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

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

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 on \mathbf{x} is the operator given by

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

Proximal Gradient Method

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

An equivalent notation: $\mathbf{x}_{t+1} = \mathcal{P}_L^h(\mathbf{x}_t) \triangleq \text{prox}_{\frac{1}{L}h} \left(\mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \right).$

Proximal Gradient

Proximal Gradient Method

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

- 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

$$\text{PG: } \mathbf{x}_{t+1} = \text{prox}_{\frac{1}{L}h}\left(\mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)\right)$$

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

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 **σ -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 $\mathbf{u} \in \mathcal{X}$,

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

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

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

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

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

Proof of PG Convergence

Proof:

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

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

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

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:

- $F(\mathbf{x}_T) - \frac{1}{T-1} \sum_{t=1}^{T-1} F(\mathbf{x}_t) \leq 0$
 - $\frac{1}{T-1} \sum_{t=1}^{T-1} F(\mathbf{x}_t) - F(\mathbf{x}^*) \leq \frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2(T-1)}$
- $$\iff F(\mathbf{x}_T) - F(\mathbf{x}^*) \leq \frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2(T-1)}$$

□

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:* $F(\mathbf{x}) \leq U_t(\mathbf{x})$ for any $\mathbf{x} \in \mathcal{X} \Rightarrow F(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq U_t(\mathbf{x}_{t+1}) - F(\mathbf{u})$

analyzing this quantity

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

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

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

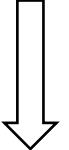
Proof:

What we have: $F(\mathbf{x}) \leq U_t(\mathbf{x})$ for any $\mathbf{x} \in \mathcal{X} \Rightarrow F(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq U_t(\mathbf{x}_{t+1}) - F(\mathbf{u})$

analyzing this quantity

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

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

 *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: $F(\mathbf{x}) \leq U_t(\mathbf{x})$ for any $\mathbf{x} \in \mathcal{X} \Rightarrow F(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq U_t(\mathbf{x}_{t+1}) - F(\mathbf{u})$

analyzing this quantity

$$\left\{ \begin{array}{l} U_t(\mathbf{x}_{t+1}) - F(\mathbf{u}) \leq \langle \nabla f(\mathbf{x}_t) + \nabla h(\mathbf{x}_{t+1}), \mathbf{x}_{t+1} - \mathbf{u} \rangle + \frac{1}{2L} \|g(\mathbf{x}_t)\|^2 \\ \text{and the fact that } \nabla f(\mathbf{x}_t) + \nabla h(\mathbf{x}_{t+1}) = -L(\mathbf{x}_{t+1} - \mathbf{x}_t) = -g(\mathbf{x}_t) \end{array} \right.$$

$$\begin{aligned} \implies U_t(\mathbf{x}_{t+1}) - F(\mathbf{u}) &\leq \langle g(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{u} \rangle + \frac{1}{2L} \|g(\mathbf{x}_t)\|^2 \\ &= \langle g(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle - \frac{1}{2L} \|g(\mathbf{x}_t)\|^2 \end{aligned}$$

□

One-Step Improvement Lemma

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

unconstrained, GD

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

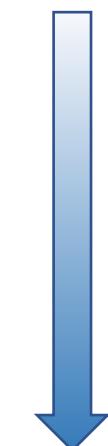
specialized

unconstrained, AGD

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

constrained, GD

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



general

composite, GD/AGD

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

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 achieve AGD in composite optimization?
→ This induces the Accelerated Proximal Gradient (**APG**) method.

Nesterov's Accelerated GD

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

Accelerated Proximal Gradient

$$\mathbf{x}_{t+1} = \text{prox}_{\frac{1}{L} h} \left(\mathbf{y}_t - \frac{1}{L} \nabla f(\mathbf{y}_t) \right), \quad \mathbf{y}_{t+1} = \mathbf{x}_{t+1} + \beta_t (\mathbf{x}_{t+1} - \mathbf{x}_t)$$

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 σ -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:** ℓ_1 -regularized least squares

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

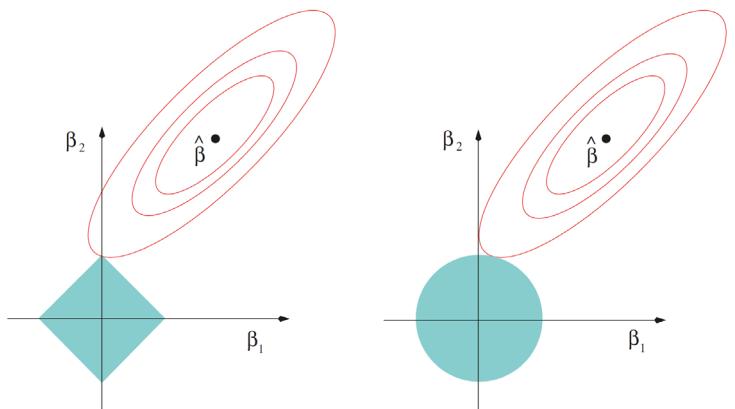
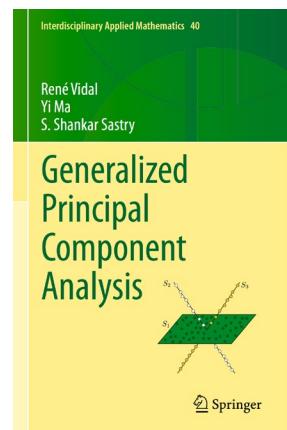
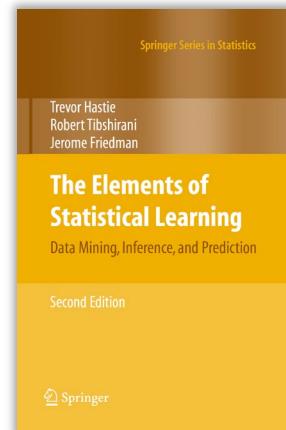
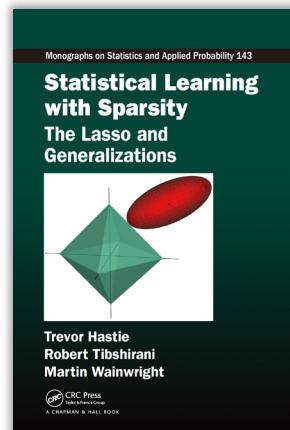


FIGURE 3.11. Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions $|\beta_1| + |\beta_2| \leq t$ and $\beta_1^2 + \beta_2^2 \leq t^2$, respectively, while the red ellipses are the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

61368 1996

Application to LASSO

- **LASSO:** ℓ_1 -regularized least squares

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

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**

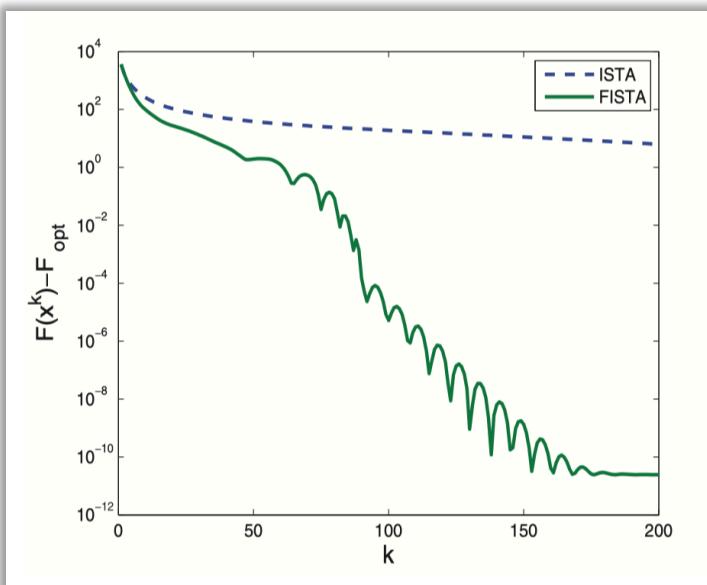
Closed-form solution:

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

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

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[†] and Marc Teboulle[‡]

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 over denoising 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) $\mathbf{Ax} = \mathbf{b} + \mathbf{w},$

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/sims/2-1/1654.html>

[‡]Department of Industrial Engineering and Management, Technion—Israel Institute of Technology, Haifa 32000, Israel (becka@ie.technion.ac.il).

[‡]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

13994

2009

A Beck, M Teboulle

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

Summary

Table 1: A summary of convergence rates of GD method for smooth optimization.

Algorithm	Function Family	Step Size	Output Sequence	Convergence Rate	Remark
GD	L -smooth and convex	$\eta = \frac{1}{L}$	$\bar{\mathbf{x}}_T \triangleq \mathbf{x}_T$	$\mathcal{O}(1/T)$	suboptimal
	L -smooth and σ -strongly convex	$\eta = \frac{2}{\sigma+L}$	$\bar{\mathbf{x}}_T \triangleq \mathbf{x}_T$	$\mathcal{O}\left(\exp\left(-\frac{T}{\kappa}\right)\right)$	suboptimal
AGD	L -smooth and convex	$\mathbf{x}_{t+1} = \mathbf{y}_t - \frac{1}{L} \nabla f(\mathbf{y}_t), \quad \mathbf{y}_{t+1} = \mathbf{x}_{t+1} + \beta_t (\mathbf{x}_{t+1} - \mathbf{x}_t)$	$\bar{\mathbf{x}}_T \triangleq \mathbf{x}_T$	$\mathcal{O}(1/T^2)$	optimal
	L -smooth and σ -strongly convex	$\mathbf{x}_{t+1} = \mathbf{y}_t - \frac{1}{L} \nabla f(\mathbf{y}_t), \quad \mathbf{y}_{t+1} = \mathbf{x}_{t+1} + \frac{\sqrt{\gamma}-1}{\sqrt{\gamma}+1} (\mathbf{x}_{t+1} - \mathbf{x}_t)$	$\bar{\mathbf{x}}_T \triangleq \mathbf{x}_T$	$\mathcal{O}\left(\exp\left(-\frac{T}{\sqrt{\kappa}}\right)\right)$	optimal
PG	$F(\mathbf{x}) \triangleq f(\mathbf{x}) + h(\mathbf{x})$ f and h are convex	$\mathbf{x}_{t+1} = \mathcal{P}_{\mathcal{L}}^h(\mathbf{x}_t) \triangleq \text{prox}_{\frac{1}{L}h}(\mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t))$	$\bar{\mathbf{x}}_T \triangleq \mathbf{x}_T$	$\mathcal{O}(1/T)$	suboptimal
APG	f is L -smooth but h is not smooth	$\mathbf{x}_{t+1} = \mathcal{P}_{\mathcal{L}}^h(\mathbf{y}_t), \quad \mathbf{y}_{t+1} = \mathbf{x}_{t+1} + \beta_t (\mathbf{x}_{t+1} - \mathbf{x}_t)$	$\bar{\mathbf{x}}_T \triangleq \mathbf{x}_T$	$\mathcal{O}(1/T^2)$	optimal

Summary

