

Introductory Lectures on Optimization

Acceleration Methods (2)

Hui Qian
qianhui@zju.edu.cn

College of Computer Science, Zhejiang University

November 6, 2024

Outline

- 1 Accelerated GD: Basic Scheme
 - Difference between Lower Bounds and Real Efficiency
 - Estimate Sequences
 - Optimal Scheme
- 2 Accelerated GD: Theoretical Analysis and Variants
 - Analysis of Optimal Scheme
 - Variant of Optimal Scheme
- 3 Reference

Part I

Accelerated Gradient Descent

Unconstrained Minimization Problem

Consider the following unconstrained minimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^N} f(\mathbf{x}),$$

where f is strongly convex: $f \in \mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$, $\mu \geq 0$.

This family of classes also contains the class of convex functions with Lipschitz continuous gradient, since

$$\mathcal{S}_{0,L}^{1,1}(\mathbb{R}^n) \equiv \mathcal{F}_L^{1,1}(\mathbb{R}^n).$$

Remark. The gradient is Lipschitz continuous.

Efficiency Estimation

For the **Gradient Descent** method, we proved the following convergence rates:

$$\begin{aligned}\mathcal{F}_L^{1,1}(\mathbb{R}^n) : f(\mathbf{x}_k) - f^* &\leq \frac{2L \|\mathbf{x}_0 - \mathbf{x}^*\|}{k + 4}, \\ \mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n) : f(\mathbf{x}_k) - f^* &\leq \frac{L}{2} \left(\frac{L - \mu}{L + \mu} \right)^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.\end{aligned}$$

These estimates differ from our lower complexity bounds by an order of magnitude.

- Generally speaking, this does not mean that the Gradient Method is **not optimal** (it may be that the lower bounds are **too optimistic**).
- However, we will see that in our case the lower bounds are **exact** up to a constant factor.

We prove this by constructing a method with rate of convergence proportional to these bounds.

Efficiency Estimation

The upper bounds vs. the lower bounds

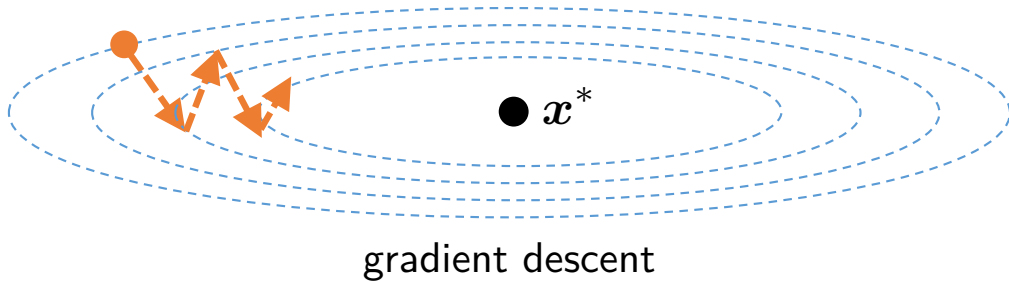
$$\mathcal{F}_L^{1,1}(\mathbb{R}^n) : f(x_k) - f^* \leq \frac{2L \|x_0 - x^*\|}{k+4}, \quad (\text{upper bound of GD})$$

$$\mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n) : f(x_k) - f^* \leq \frac{L}{2} \left(\frac{L-\mu}{L+\mu} \right)^{2k} \|x_0 - x^*\|^2, \quad (\text{upper bound of GD})$$

$$\mathcal{F}_L^{1,1}(\mathbb{R}^n) : f(x_k) - f^* \geq \frac{3L \|x_0 - x^*\|^2}{32(k+1)^2}, \quad (\text{lower bound of class})$$

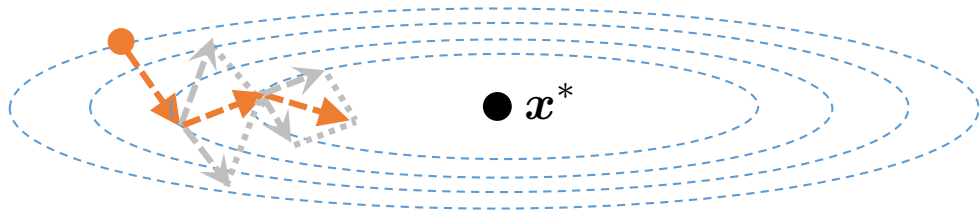
$$\mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n) : f(x_k) - f^* \geq \frac{\mu}{2} \left(\frac{\sqrt{Q_f} - 1}{\sqrt{Q_f} + 1} \right)^{2k} \|x_0 - x^*\|^2. \quad (\text{lower bound of class})$$

Efficiency Estimation



This picture comes from ELE 522: Large-Scale Optimization for Data Science of Princeton University.

Efficiency Estimation



heavy-ball method

This picture comes from ELE 522: Large-Scale Optimization for Data Science of Princeton University.

Historical Notes on Accelerated GD

Nesterov first discovered an accelerated gradient descent method for smooth functions in 1983[Nesterov, 1983], and proposed a more generalized method in 1988 [Nesterov and Nemirovsky, 1988]. Nesterov's technique boils down to a mechanism called [Estimate Sequences](#).

- In 2004 and 2007, Nesterov generalized the acceleration to non-smooth objective functions [Nesterov, 2013] and combinatorial functions [Nesterov, 2007], respectively.
- In 2008, [FISTA](#)(Fast Iterative Shrinkage-Thresholding Algorithm) [Beck and Teboulle, 2009].
- In 2009, Unified Analysis of Smooth and Nonsmooth Function Classes[Tseng and Yun, 2009].

Since then, various first-order algorithms have had their own accelerated versions.

Historical Notes on Accelerated GD

Since then, new approaches have been proposed, with the interpretation and study of Nesterov's acceleration. The understanding has also deepened.

- **ODE perspective:** Differential Equation[Su et al., 2014], Integral Quadratic Constraints [Lessard et al., 2016], Variational Perspective[Wibisono et al., 2016].
- **Geometry perspective:** Linear Coupling[Allen-Zhu and Orecchia, 2014], Geometric interpretation loosely inspired by the ellipsoid method[Bubeck et al., 2015]
- **Game theory perspective:** view this method as a natural iterative buyer-supplier game[Lan and Zhou, 2018]

Estimate Sequences

The Nesterov acceleration schemes and efficiency bounds of optimal methods are based on the notion of estimating sequences.

Definition 6 (Definition 2.2.1 of Nesterov [2013])

A pair of sequences $\{\phi_k(\mathbf{x})\}_{k=0}^{\infty}$ and $\{\lambda_k\}_{k=0}^{\infty}$, $\lambda_k \geq 0$, are called the **estimating sequences** of function of $f(\mathbf{x})$ if

$$\lambda_k \rightarrow 0,$$

and for any $x \in \mathbb{R}^n$ and all $k \geq 0$ we have

$$\phi_k(\mathbf{x}) \leq (1 - \lambda_k)f(\mathbf{x}) + \lambda_k\phi_0(\mathbf{x}). \quad (5)$$

Estimate Sequences

Lemma 7 (Lemma 2.2.1)

If for some sequence of points $\{\mathbf{x}_k\}$, we have

$$f(\mathbf{x}_k) \leq \phi_k^* \equiv \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x}), \quad (6)$$

Then $f(\mathbf{x}_k) - f^* \leq \lambda_k [\phi_0(\mathbf{x}^*) - f^*] \rightarrow 0$.

Proof. Indeed,

$$\begin{aligned} f(\mathbf{x}_k) &\stackrel{(6)}{\leq} \phi_k^* = \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x}) \stackrel{(5)}{\leq} \min_{\mathbf{x} \in \mathbb{R}^n} [(1 - \lambda_k)f(\mathbf{x}) + \lambda_k \phi_0(\mathbf{x})] \\ &\leq (1 - \lambda_k)f(\mathbf{x}^*) + \lambda_k \phi_0(\mathbf{x}^*). \end{aligned}$$



Estimate Sequences

$$f(\mathbf{x}_k) \leq \phi_k^*, \quad \phi_\infty(\mathbf{x}) \leq f(\mathbf{x}),$$

Thus, for any sequence $\{\mathbf{x}_k\}$, satisfying (6), we can derive its rate of convergence directly from the convergence rate of the sequence $\{\lambda_k\}$.

However, at this moment we have two serious questions.

- Firstly, we do not know how to form the estimating sequences.
- Secondly, we do not know how to satisfy inequalities (6).

Estimate Sequences

The first question is simpler.

Lemma 8 (Lemma 2.2.2 of Nesterov [2013])

Assume that

- 1 $f \in \mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$,
- 2 $\phi_0(\mathbf{x})$ is an arbitrary convex function on \mathbb{R}^n ,
- 3 $\{\mathbf{y}_k\}_{k=0}^\infty$ is an arbitrary sequence of points in \mathbb{R}^n ,
- 4 the coefficients $\{\alpha_k\}_{k=0}^\infty$ satisfy conditions $\alpha_k \in (0, 1)$, $\sum_{k=0}^\infty \alpha_k = \infty$, ,
- 5 we choose $\lambda_0 = 1$.

To be continued...

Estimate Sequences

Lemma 8

(Continued.) Then, the pair of sequences $\{\phi_k(\mathbf{x})\}_{k=0}^{\infty}$ and $\{\lambda_k\}_{k=0}^{\infty}$ defined recursively by the relations:

$$\begin{aligned}\lambda_{k+1} &= (1 - \alpha_k)\lambda_k, \\ \phi_{k+1}(\mathbf{x}) &= (1 - \alpha_k)\phi_k(\mathbf{x}) \\ &\quad + \alpha_k[f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|^2],\end{aligned}\tag{7}$$

are estimating sequences.

注： ϕ_k 是 ϕ_0 和一系列下界的平均。

Estimate Sequences

Proof. Indeed, $\phi_0(\mathbf{x}) \leq (1 - \lambda_0)f(\mathbf{x}) + \lambda_0\phi_0(\mathbf{x}) \equiv \phi_0(\mathbf{x})$. Further, let (5) hold for some $k \geq 0$. Then

$$\begin{aligned}
 \phi_{k+1}(\mathbf{x}) &\leq (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k f(\mathbf{x}) \text{ (strongly convex)} \\
 &= (1 - (1 - \alpha_k)\lambda_k)f(\mathbf{x}) + (1 - \alpha_k)\underbrace{(\phi_k(\mathbf{x}) - (1 - \lambda_k)f(\mathbf{x}))}_{\leq \lambda_k\phi_0(\mathbf{x})} \\
 &\leq (1 - (1 - \alpha_k)\lambda_k)f(\mathbf{x}) + (1 - \alpha_k)\lambda_k\phi_0(\mathbf{x}) \text{ (definition)} \\
 &= (1 - \lambda_{k+1})f(\mathbf{x}) + \lambda_{k+1}\phi_0(\mathbf{x}).
 \end{aligned}$$

Thus, we arrive at

$$\phi_k(\mathbf{x}) \leq (1 - \lambda_k)f(\mathbf{x}) + \lambda_k\phi_0(\mathbf{x}).$$

Estimate Sequences

Proof. (Continued.) It remains to note that condition 4 ensures $\lambda_k \rightarrow 0$.

Since $\lambda_k > 0$ holds for all $k = 0, 1, \dots, k$. We have $\ln \lambda_{k+1} = \ln \{(1 - \alpha_k)\lambda_k\}$. Therefore,

$$\ln \lambda_{k+1} - \ln \lambda_k = \ln(1 - \alpha_k) \leq -\alpha_k. \quad (\text{using } \ln(1 + x) \leq x \text{ here.})$$

We telescope the above inequality and obtain

$$\ln \lambda_{k+1} - \ln \lambda_0 \leq -\sum_{i=0}^k \alpha_i.$$

With condition 4, it implies that $\ln \lambda_{k+1} \rightarrow -\infty$, when $k \rightarrow \infty$. That is $k \rightarrow \infty$, $\lambda_k \rightarrow 0$. \square

Estimate Sequence

Thus, the above statement provides us with some rules for updating the estimating sequences. Now we have two control sequences which can help us to maintain recursively the relation (6).

At this moment, we are also free in our choice of initial function $\phi_0(x)$.

Let us choose it as a simple quadratic function. Then, we can obtain a closed form recurrence for values ϕ_k^* .

$$f(\mathbf{x}_k) \leq \phi_k^* \equiv \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x}), \quad (6)$$

Estimate Sequence

Lemma 9 (Lemma 2.2.3)

Let

$$\phi_0(\mathbf{x}) = \phi_0^* + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{v}_0\|^2.$$

Then the process (7) preserves the canonical form of functions $\{\phi_k(x)\}$:

$$\phi_k(\mathbf{x}) \equiv \phi_k^* + \frac{\gamma_k}{2} \|\mathbf{x} - \mathbf{v}_k\|^2, \quad (8)$$

Where the sequences $\{\gamma_k\}$, $\{\mathbf{v}_k\}$ and $\{\phi_k^*\}$ are defined as follows:

To be continued...

Estimate Sequence

Lemma 9

(Continued.)

$$\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k\mu,$$

$$\mathbf{v}_{k+1} = \frac{1}{\gamma_{k+1}} [(1 - \alpha_k)\gamma_k \mathbf{v}_k + \alpha_k\mu \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k)],$$

$$\begin{aligned} \phi_{k+1}^* &= (1 - \alpha_k)\phi_k^* + \alpha_k f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|^2 \\ &\quad + \frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \left(\frac{\mu}{2} \|\mathbf{y}_k - \mathbf{v}_k\|^2 + \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle \right). \end{aligned}$$

Remark

- 1 First of all, we have to confirm whether the canonical form can be preserved. If possible, $\phi_k(\mathbf{x})$ is quadratic, so it can be formulated as (8);
- 2 Then, what the lemma does is deduce the definitions of $\{\gamma_k\}$, $\{\mathbf{v}_k\}$ and $\{\phi_k^*\}$.
- 3 Set γ_k first, and deduce the rest directly. If the deduced result conforms to the definitions in the lemma, then the lemma is proved.

$$\begin{aligned}\phi_k(\mathbf{x}) &\equiv \phi_k^* + \frac{\gamma_k}{2} \|\mathbf{x} - \mathbf{v}_k\|^2, \\ \gamma_{k+1} &= (1 - \alpha_k)\gamma_k + \alpha_k\mu,\end{aligned}\tag{8}$$

Estimate Sequence

$$\phi_k(\mathbf{x}) \equiv \phi_k^* + \frac{\gamma_k}{2} \|\mathbf{x} - \mathbf{v}_k\|^2. \quad (8)$$

Proof. For canonical form: (Induction) Note that $\nabla^2 \phi_0(\mathbf{x}) = \gamma_0 I_n$. Let us show that $k \geq 0$, $\nabla^2 \phi_k(\mathbf{x}) = \gamma_k I_n$ for all $k \geq 0$. Indeed, if it is true for some k , then

$$\underbrace{\nabla^2 \phi_{k+1}(\mathbf{x}) = (1 - \alpha_k) \nabla^2 \phi_k(\mathbf{x}) + \alpha_k \mu I_n}_{\text{by (7)}} = ((1 - \alpha_k) \gamma_k + \alpha_k \mu) I_n \equiv \gamma_{k+1} I_n.$$

This justifies the canonical form (7) of the function $\phi_k(\mathbf{x})$.

Estimate Sequence

Proof. (Continued.) For \mathbf{v}_{k+1} (which is the optimal value of ϕ_{k+1}): Further,

$$\begin{aligned}\phi_{k+1}(\mathbf{x}) &= (1 - \alpha_k) \left(\phi_k^* + \frac{\gamma_k}{2} \|\mathbf{x} - \mathbf{v}_k\|^2 \right) \\ &\quad + \alpha_k \left[f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|^2 \right].\end{aligned}\tag{9}$$

Therefore the equation $\nabla \phi_{k+1}(\mathbf{x}) = 0$, which is the first order optimality condition for the function $\phi_{k+1}(\mathbf{x})$, is as follows:

$$(1 - \alpha_k)\gamma_k(\mathbf{x} - \mathbf{v}_k) + \alpha_k \nabla f(\mathbf{y}_k) + \alpha_k \mu(\mathbf{x} - \mathbf{y}_k) = 0.$$

From this equation, we get a closed form expression for the point \mathbf{v}_{k+1} , the minimum of the function $\phi_{k+1}(\mathbf{x})$. This arrives at $\mathbf{v}_{k+1} = \frac{1}{\gamma_{k+1}} [(1 - \alpha_k)\gamma_k \mathbf{v}_k + \alpha_k \mu \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k)]$.

Estimate Sequence

Proof. (Continued.) For ϕ_{k+1}^* : In view of the recurrence (7) for the sequence $\{\phi_k(x)\}$, we have

$$\begin{aligned}
 \phi_{k+1}^* + \frac{\gamma_{k+1}}{2} \|\mathbf{y}_k - \mathbf{v}_{k+1}\|^2 &= \phi_{k+1}(\mathbf{y}_k) \quad (\text{Substitute } \mathbf{y}_k) \\
 &= (1 - \alpha_k) \underbrace{\left(\phi_k^* + \frac{\gamma_k}{2} \|\mathbf{y}_k - \mathbf{v}_k\|^2 \right)}_{\phi_k(\mathbf{y}_k)} + \alpha_k f(\mathbf{y}_k) \\
 &= \boxed{(1 - \alpha_k) \frac{\gamma_k}{2}} \|\mathbf{y}_k - \mathbf{v}_k\|^2 + (1 - \alpha_k) \phi_k^* + \alpha_k f(\mathbf{y}_k).
 \end{aligned}$$

Remark.

$$\begin{aligned}
 \phi_{k+1}(x) &= (1 - \alpha_k) \phi_k(x) \\
 &\quad + \alpha_k [f(y_k) + \langle \nabla f(y_k), x - y_k \rangle + \frac{\mu}{2} \|x - y_k\|^2],
 \end{aligned}$$

Estimate Sequence

Proof. (Continued.) By the recursive relation for \mathbf{v}_{k+1} , we have (两边都减去 \mathbf{y}_k),

$$\mathbf{v}_{k+1} - \mathbf{y}_k = \frac{1}{\gamma_{k+1}} [(1 - \alpha_k)\gamma_k(\mathbf{v}_k - \mathbf{y}_k) - \alpha_k \nabla f(\mathbf{y}_k)].$$

Remark.

$$\begin{aligned}\mathbf{v}_{k+1} - \mathbf{y}_k &= \frac{1}{\gamma_{k+1}} [(1 - \alpha_k)\gamma_k \mathbf{v}_k + \alpha_k \mu \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k) - \gamma_{k+1} \mathbf{y}_k] \\ &= \frac{1}{\gamma_{k+1}} [(1 - \alpha_k)\gamma_k(\mathbf{v}_k - \mathbf{y}_k) - \alpha_k \nabla f(\mathbf{y}_k)].\end{aligned}$$

The last step is obtained by expanding $\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k \mu$.

Estimate Sequence

Proof. (Continued.) Therefore,

$$\begin{aligned} \frac{\gamma_{k+1}}{2} \|\mathbf{v}_{k+1} - \mathbf{y}_k\|^2 &= \boxed{\frac{1}{2\gamma_{k+1}} [(1 - \alpha_k)^2 \gamma_k^2 \|\mathbf{v}_k - \mathbf{y}_k\|^2} \\ &\quad - 2\alpha_k(1 - \alpha_k)\gamma_k \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle} \\ &\quad + \alpha_k^2 \|\nabla f(\mathbf{y}_k)\|^2]. \end{aligned}$$

$$\mathbf{v}_{k+1} - \mathbf{y}_k = \frac{1}{\gamma_{k+1}} [(1 - \alpha_k)\gamma_k(\mathbf{v}_k - \mathbf{y}_k) - \alpha_k \nabla f(\mathbf{y}_k)].$$

Estimate Sequence

Proof. (Continued.) It remains to substitute this relation into (9), taking into account that the multiplicative factor for the term $\|\mathbf{y}_k - \mathbf{v}_{k+1}\|^2$ in the resulting expression is as follows:

$$\begin{aligned} \boxed{(1 - \alpha_k) \frac{\gamma_k}{2}} - \boxed{\frac{1}{2\gamma_{k+1}} (1 - \alpha_k)^2 \gamma_k^2} &= (1 - \alpha_k) \frac{\gamma_k}{2} \left(1 - \frac{(1 - \alpha_k) \gamma_k}{\gamma_{k+1}} \right) \\ &= (1 - \alpha_k) \frac{\gamma_k}{2} \cdot \frac{\alpha_k \mu}{\gamma_{k+1}}. \end{aligned}$$

□

$$\phi_{k+1}^* + \frac{\gamma_{k+1}}{2} \|\mathbf{y}_k - \mathbf{v}_{k+1}\|^2 = (1 - \alpha_k) \left(\phi_k^* + \frac{\gamma_k}{2} \|\mathbf{y}_k - \mathbf{v}_k\|^2 \right) + \alpha_k f(\mathbf{y}_k).$$

Constructing x_k

(Induction) Indeed, assume that we already have x_k : $\phi_k^* \geq f(x_k)$. Let's check how to make $\phi_{k+1}^* \geq f(x_{k+1})$. Then, in the view of Lemma 9,

$$\begin{aligned}\phi_{k+1}^* &\geq (1 - \alpha_k)f(x_k) + \alpha_k f(y_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(y_k)\|^2 \\ &\quad + \frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \langle \nabla f(y_k), v_k - y_k \rangle.\end{aligned}$$

Lemma 9: $\phi_{k+1}^* = (1 - \alpha_k) \underbrace{\phi_k^*}_{\phi_k^* \geq f(x_k)} + \alpha_k f(y_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(y_k)\|^2$

$$+ \frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \left(\frac{\mu}{2} \|y_k - v_k\|^2 + \langle \nabla f(y_k), v_k - y_k \rangle \right).$$

Constructing x_k

Since $f(\mathbf{x}_k) \geq f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k \rangle$, we get the following estimate:

$$\begin{aligned}\phi_{k+1}^* &\geq f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|^2 \\ &\quad + (1 - \alpha_k) \langle \nabla f(\mathbf{y}_k), \frac{\alpha_k \gamma_k}{\gamma_{k+1}} (\mathbf{v}_k - \mathbf{y}_k) + \mathbf{x}_k - \mathbf{y}_k \rangle.\end{aligned}$$

$$\begin{aligned}\phi_{k+1}^* &\geq (1 - \alpha_k) f(\mathbf{x}_k) + \alpha_k f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|^2 \\ &\quad + \frac{\alpha_k (1 - \alpha_k) \gamma_k}{\gamma_{k+1}} \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle.\end{aligned}$$

Constructing x_k

Let us look at this inequality. We want to have $\phi_{k+1}^* \geq f(\mathbf{x}_{k+1})$. Recall that we can ensure the inequality

$$f(\mathbf{y}_k) - \frac{1}{2L} \|\nabla f(\mathbf{y}_k)\|^2 \geq f(\mathbf{x}_{k+1}),$$

in many different ways. The simplest one is just to take the gradient step

$$\boxed{\mathbf{x}_{k+1} = \mathbf{y}_k - h_k \nabla f(\mathbf{y}_k)},$$

with $h_k = \frac{1}{L}$. (此时, 我们还有 y_k 和 α_k 两个自由度可以调)

$$0 \leq f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

Constructing x_k

adjust \mathbf{y}_k and α_k such that

$$\begin{aligned} \phi_{k+1}^* \geq f(\mathbf{y}_k) - \underbrace{\frac{\alpha_k^2}{2\gamma_{k+1}}}_{=1/2L} \|\nabla f(\mathbf{y}_k)\|^2 \\ + (1 - \alpha_k) \langle \nabla f(\mathbf{y}_k), \underbrace{\frac{\alpha_k \gamma_k}{\gamma_{k+1}} (\mathbf{v}_k - \mathbf{y}_k) + \mathbf{x}_k - \mathbf{y}_k}_{=0} \rangle. \end{aligned}$$

Constructing x_k

Let us define α_k as a positive root of the quadratic equation

$$\boxed{L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu \quad (= \gamma_{k+1})}.$$

Then $\frac{\alpha_k^2}{2\gamma_{k+1}} = \frac{1}{2L}$, and we can replace the previous inequality by the following one:

$$\phi_{k+1}^* \geq f(\mathbf{x}_{k+1}) + (1 - \alpha_k)\langle \nabla f(\mathbf{y}_k), \underbrace{\frac{\alpha_k\gamma_k}{\gamma_{k+1}}(\mathbf{v}_k - \mathbf{y}_k) + \mathbf{x}_k - \mathbf{y}_k}_{=0} \rangle.$$

Constructing x_k

Let us now use our freedom in the choice of y_k . It can be found from the equation:

$$\frac{\alpha_k \gamma_k}{\gamma_{k+1}} (\mathbf{v}_k - \mathbf{y}_k) + \mathbf{x}_k - \mathbf{y}_k = 0$$

This is

$$\boxed{\mathbf{y}_k = \frac{\alpha_k \gamma_k \mathbf{v}_k + \gamma_{k+1} \mathbf{x}_k}{\gamma_k + \alpha_k \mu}},$$

and we come to the following methods, which are often addressed as **accelerated gradient method**.

Optimal Scheme

General Scheme of Optimal Method

- 0. Choose the point $\mathbf{x}_0 \in \mathbb{R}^n$, some $\gamma_0 > 0$, and set $\mathbf{v}_0 = \mathbf{x}_0$.
- 1. k -th iteration ($k \geq 0$).
 - (a) Compute $\alpha_k \in (0, 1)$ from the equation $L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu$. Set $\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k\mu$.

(b) Choose

$$\mathbf{y}_k = \frac{\alpha_k \gamma_k \mathbf{v}_k + \gamma_{k+1} \mathbf{x}_k}{\gamma_k + \alpha_k \mu}.$$

(10)

Compute $f(\mathbf{y}_k)$ and $\nabla f(\mathbf{y}_k)$.

- (c) Find \mathbf{x}_{k+1} such that $f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}_k) - \frac{1}{2L} \|\nabla f(\mathbf{y}_k)\|^2$
- (d) Set $\mathbf{v}_{k+1} = \frac{(1 - \alpha_k)\gamma_k \mathbf{v}_k + \alpha_k \mu \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k)}{\gamma_{k+1}}.$

Optimal Scheme

Note that in Step (c) of this scheme we can choose an arbitrary \mathbf{x}_{k+1} satisfying the inequality

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}_k) - \frac{\omega}{2} \|\nabla f(\mathbf{y}_k)\|^2,$$

with some $\omega > 0$. Then the constant $\frac{1}{\omega}$ replaces L in the equation of step (a).

Part II

Theoretically Analysis for AGD

Analysis of Optimal Scheme

Theorem 10 (Theorem 2.2.1)

Scheme (10) generates a sequence of points $\{\mathbf{x}_k\}_{k=0}^{\infty}$ such that

$$f(\mathbf{x}_k) - f^* \leq \lambda_k [f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2],$$

where $\lambda_0 = 1$ and $\lambda_k = \prod_{i=0}^{k-1} (1 - \alpha_i)$.

Proof. Indeed, let us choose $\phi_0(\mathbf{x}) = f(\mathbf{x}_0) + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{x}_0\|^2$. Then $f(\mathbf{x}_0) = \phi_0^*$ and we get $f(\mathbf{x}_k) \leq \phi_k^*$ by the rules of the scheme. It remains to use Lemma 7.

(From Lemma 7, we have $f(\mathbf{x}_k) - f^* \leq \lambda_k [\phi_0(\mathbf{x}^*) - f^*] \rightarrow 0$. 代入 $\phi_0(\mathbf{x}^*)$ 即可) 。 \square

Analysis of Optimal Scheme

Lemma 11 (Lemma 2.2.4)

If in the method (10) we choose $\gamma_0 \geq \mu$, then

$$\lambda_k \leq \min \left\{ \left(1 - \sqrt{\frac{\mu}{L}} \right)^k, \frac{4L}{(2\sqrt{L} + k\sqrt{\gamma_0})^2} \right\}. \quad (11)$$

Proof. Indeed, if $\gamma_k \geq \mu$, then $\gamma_{k+1} = L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu \geq \mu$. Since $\gamma_0 \geq \mu$, we conclude that this inequality is valid for all γ_k . Hence $\alpha_k \geq \sqrt{\frac{\mu}{L}}$ and we have proved the first inequality in (11).

Remark. $\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k\mu > (1 - \alpha_k)\mu + \alpha_k\mu = \mu$.

Analysis of Optimal Scheme

Since

$$\lambda_k = \prod_{i=0}^{k-1} (1 - \alpha_i),$$

and for every α_i we have

$$\alpha_i \geq \sqrt{\frac{\mu}{L}}.$$

Thus, we obtain

$$\lambda_k \leq \left(1 - \sqrt{\frac{\mu}{L}}\right)^k.$$

Analysis of Optimal Scheme

Proof. (Continued.) For the second inequality, let us first prove that $\gamma_k \geq \gamma_0 \lambda_k$. Indeed, since $\gamma_0 = \gamma_0 \lambda_0$, we can use induction:

$$\gamma_{k+1} \geq (1 - \alpha_k) \gamma_k \geq (1 - \alpha_k) \gamma_0 \lambda_k = \gamma_0 \lambda_{k+1}.$$

Therefore $L \alpha_k^2 = \gamma_{k+1} \geq \gamma_0 \lambda_{k+1}$. That is

$$\sqrt{\lambda_{k+1}} \leq \alpha_k \sqrt{\frac{L}{\gamma_0}}. \quad (12)$$

$$\gamma_{k+1} = (1 - \alpha_k) \gamma_k + \alpha_k \mu.$$

Analysis of Optimal Scheme

Proof. (Continued.) Denote $b_k = \frac{1}{\sqrt{\lambda_k}}$. Since $\{\lambda_k\}$ is a decreasing sequence, we have

$$\begin{aligned} b_{k+1} - b_k &= \frac{\sqrt{\lambda_k} - \sqrt{\lambda_{k+1}}}{\sqrt{\lambda_k \lambda_{k+1}}} = \frac{\lambda_k - \lambda_{k+1}}{\sqrt{\lambda_k \lambda_{k+1}} (\sqrt{\lambda_k} + \sqrt{\lambda_{k+1}})} \\ &\geq \frac{\lambda_k - \lambda_{k+1}}{2\lambda_k \sqrt{\lambda_{k+1}}} = \frac{\lambda_k - (1 - \alpha_k)\lambda_k}{2\lambda_k \sqrt{\lambda_{k+1}}} = \frac{\alpha_k}{2\sqrt{\lambda_{k+1}}} \\ &\geq \frac{1}{2} \sqrt{\frac{\gamma_0}{L}}. (\text{Since (12)}) \end{aligned}$$

Thus, $b_k \geq 1 + \frac{k}{2} \sqrt{\frac{\gamma_0}{L}}$ and the lemma is proved. □

Remark. In Theorem 10, set $\lambda_0 = 1$, $\lambda_k = \prod_{i=0}^{k-1} (1 - \alpha_i)$.

Analysis of Optimal Scheme

Theorem 12 (Theorem 2.2.2)

Let us take in (10) $\gamma_0 = L$. Then this scheme generates a sequence $\{\mathbf{x}_k\}_{k=0}^{\infty}$ such that

$$f(\mathbf{x}_k) - f^* \leq L \min \left\{ \left(1 - \sqrt{\frac{\mu}{L}}\right)^k, \frac{4}{(k+2)^2} \right\} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

This means that (10) is optimal for unconstrained minimization of the functions from $\mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$, $\mu \geq 0$.

Proof. We get the above inequality using $f(\mathbf{x}_0) - f^* \leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$ and Theorem 10 with lemma 11.

Analysis of Optimal Scheme

Proof. (Continued.) Let $\mu > 0$. From the lower complexity bounds for the class, we have

$$f(x_k) - f^* \geq \frac{\mu}{2} \left(\frac{\sqrt{Q_f} - 1}{\sqrt{Q_f} + 1} \right)^{2k} R^2 \geq \frac{\mu}{2} \exp \left(-\frac{4k}{\sqrt{Q_f} - 1} \right) R^2,$$

where $Q_f = L/\mu$ and $R = \|\mathbf{x}_0 - \mathbf{x}^*\|$. Therefore, the worst case bound for finding x_k satisfying $f(\mathbf{x}_k) - f^* \leq \epsilon$ cannot be better than

$$k \geq \frac{\sqrt{Q_f} - 1}{4} \left[\ln \frac{1}{\epsilon} + \ln \frac{\mu}{2} + 2 \ln R \right].$$

Remark. $\ln(1 + x) \leq x$

Analysis of Optimal Scheme

Proof. (Continued.) For our scheme we have

$$f(\mathbf{x}_k) - f^* \leq LR^2 \left(1 - \sqrt{\frac{\mu}{L}}\right)^k \leq LR^2 \exp\left(-\frac{k}{\sqrt{Q_f}}\right).$$

Therefore, we guarantee that $k \leq \sqrt{Q_f} \left[\ln \frac{1}{\epsilon} + \ln L + 2 \ln R\right]$. Thus, the main term in this estimate, $\sqrt{Q_f} \ln \frac{1}{\epsilon}$, is proportional to the lower bound. The same reasoning can be used for the class $\mathcal{S}_{0,L}^{1,1}(\mathbb{R}^n)$. □

Variant of Optimal Scheme

Let us analyze a variant of the scheme (10).

Constant Step Scheme, I

- 0. Choose $\mathbf{x}_0 \in \mathbb{R}^n$ and $\gamma_0 > 0$. Set $\mathbf{v}_0 = \mathbf{x}_0$.
- 1. k -th iteration ($k \geq 0$).
 - (a) Compute $\alpha_k \in (0, 1)$ from the equation $L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu$. Set $\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k\mu$.
 - (b) Choose $\mathbf{y}_k = \frac{\alpha_k\gamma_k\mathbf{v}_k + \gamma_{k+1}\mathbf{x}_k}{\gamma_k + \alpha_k\mu}$. Compute $f(\mathbf{y}_k)$ and $\nabla f(\mathbf{y}_k)$.
 - (c) Set $\mathbf{x}_{k+1} = \mathbf{y}_k - \frac{1}{L}\nabla f(\mathbf{y}_k)$ and

$$\mathbf{v}_{k+1} = \frac{1}{\gamma_{k+1}}[(1 - \alpha_k)\gamma_k\mathbf{v}_k + \alpha_k\mu\mathbf{y}_k - \alpha_k\nabla f(\mathbf{y}_k)].$$

(13)

Variant of Optimal Scheme

Let us demonstrate that this scheme can be rewritten in a simpler form. Note that

$$\mathbf{y}_k = \frac{1}{\gamma_k + \alpha_k \mu} (\alpha_k \boxed{\gamma_k \mathbf{v}_k} + \gamma_{k+1} \mathbf{x}_k), \quad (14)$$

$$\mathbf{x}_{k+1} = \mathbf{y}_k - \frac{1}{L} \nabla f(\mathbf{y}_k), \quad (15)$$

$$\mathbf{v}_{k+1} = \frac{1}{\gamma_{k+1}} [(1 - \alpha_k) \boxed{\gamma_k \mathbf{v}_k} + \alpha_k \mu \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k)]. \quad (16)$$

Remark. Eliminate \mathbf{v}_k .

Variant of Optimal Scheme

Therefore,

$$\begin{aligned} \mathbf{v}_{k+1} &= \frac{1}{\gamma_{k+1}} \left\{ \frac{1 - \alpha_k}{\alpha_k} [(\gamma_k + \alpha_k \mu) \mathbf{y}_k - \gamma_{k+1} \mathbf{x}_k] + \alpha_k \mu \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k) \right\} \\ &= \frac{1}{\gamma_{k+1}} \left\{ \frac{(1 - \alpha_k) \gamma_k}{\alpha_k} \mathbf{y}_k + \mu \mathbf{y}_k \right\} - \frac{1 - \alpha_k}{\alpha_k} \mathbf{x}_k - \frac{\alpha_k}{\gamma_{k+1}} \nabla f(\mathbf{y}_k) \\ &= \mathbf{x}_k + \frac{1}{\alpha_k} (\mathbf{y}_k - \mathbf{x}_k) - \frac{1}{\alpha_k L} \nabla f(\mathbf{y}_k) \\ &= \mathbf{x}_k + \frac{1}{\alpha_k} \left[(\mathbf{y}_k - \mathbf{x}_k) - \frac{1}{L} \nabla f(\mathbf{y}_k) \right] = \mathbf{x}_k + \frac{1}{\alpha_k} (\mathbf{x}_{k+1} - \mathbf{x}_k). \end{aligned}$$

$$\gamma_{k+1} = L \alpha_k^2 = (1 - \alpha_k) \gamma_k + \alpha_k \mu. \quad \mathbf{x}_{k+1} = \mathbf{y}_k - \frac{1}{L} \nabla f(\mathbf{y}_k).$$

Variant of Optimal Scheme

Hence,

$$\begin{aligned}\mathbf{y}_{k+1} &= \frac{1}{\gamma_{k+1} + \alpha_{k+1}\mu} (\alpha_{k+1}\gamma_{k+1}\mathbf{v}_{k+1} + \gamma_{k+2}\mathbf{x}_{k+1}) \\ &= \mathbf{x}_{k+1} + \frac{\alpha_{k+1}\gamma_{k+1}(\mathbf{v}_{k+1} - \mathbf{x}_{k+1})}{\gamma_{k+1} + \alpha_{k+1}\mu} = \mathbf{x}_{k+1} + \beta_k(\mathbf{x}_{k+1} - \mathbf{x}_k),\end{aligned}$$

where

$$\beta_k = \frac{\alpha_{k+1}\gamma_{k+1}(1 - \alpha_k)}{\alpha_k(\gamma_{k+1} + \alpha_{k+1}\mu)}.$$

Thus, \mathbf{y}_k is not involved with \mathbf{v}_k .

Variant of Optimal Scheme

Further eliminate γ_k : Let us do the same with γ_k . We have

$$\alpha_k^2 L = (1 - \alpha_k)\gamma_k + \mu\alpha_k \equiv \gamma_{k+1}.$$

Therefore,

$$\begin{aligned}\beta_k &= \frac{\alpha_{k+1}\gamma_{k+1}(1 - \alpha_k)}{\alpha_k(\gamma_{k+1} + \alpha_{k+1}\mu)} = \frac{\alpha_{k+1}\gamma_{k+1}(1 - \alpha_k)}{\alpha_k(\gamma_{k+1} + \alpha_{k+1}^2 L - (1 - \alpha_{k+1})\gamma_{k+1})} \\ &= \frac{\gamma_{k+1}(1 - \alpha_k)}{\alpha_k(\gamma_{k+1} + \alpha_{k+1}L)} = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}.\end{aligned}$$

Variant of Optimal Scheme

The iteration of α_{k+1} is processed below.

- Note that $\alpha_{k+1}^2 = (1 - \alpha_{k+1})\alpha_k^2 + q\alpha_{k+1}$ (since $L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu$ and $\gamma_{k+1} = L\alpha_k^2$), where $q = \mu/L$,
- And we also have

$$\alpha_0^2 L = (1 - \alpha_0)\gamma_0 + \mu\alpha_0.$$

The latter relation means that γ_0 can be seen as a function of α_0 . (Therefore, setting α_0 is equivalent to setting γ_0 .)

Thus, we can completely eliminate the sequence $\{\gamma_k\}$.

Variant of Optimal Scheme

Constant Step Scheme, II

- 0. Choose $\mathbf{x}_0 \in \mathbb{R}^n$ and $\alpha_0 \in (0, 1)$. Set $\mathbf{y}_0 = \mathbf{x}_0$ and $q = \frac{\mu}{L}$.
- 1. k -th iteration ($k \geq 0$).
 - (a) Compute $f(\mathbf{y}_k)$ and $\nabla f(\mathbf{y}_k)$. Set $\mathbf{x}_{k+1} = \mathbf{y}_k - \frac{1}{L} \nabla f(\mathbf{y}_k)$.
 - (b) Compute $\alpha_{k+1} \in (0, 1)$ from equation

$$\alpha_{k+1}^2 = (1 - \alpha_{k+1})\alpha_k^2 + q\alpha_{k+1}, \quad (17)$$

$$\text{and set } \beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}},$$

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

Analysis of the Variant

Theorem 13 (Theorem 2.2.3)

If in scheme (17)

$$\alpha_0 \geq \sqrt{\frac{\mu}{L}}, \quad (18)$$

then

$$\begin{aligned} f(\mathbf{x}_k) - f^* &\leq \min \left\{ \left(1 - \sqrt{\frac{\mu}{L}} \right)^k, \frac{4L}{(2\sqrt{L} + k\sqrt{\gamma_0})^2} \right\} \\ &\quad \times [f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2], \end{aligned}$$

where $\gamma_0 = \frac{\alpha_0(\alpha_0 L - \mu)}{1 - \alpha_0}$.

Analysis of the Variant

We do not need to prove this theorem since the initial scheme is not changed. We change only notation. In Theorem13 condition (18) is equivalent to $\gamma_0 \geq \mu$.

If we choose $\alpha_0 = \sqrt{\frac{\mu}{L}}$ (this corresponds to $\gamma_0 = \mu$), Scheme (17) becomes very simple. Then

$$\alpha_k = \sqrt{\frac{\mu}{L}}, \beta_k = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}.$$

for all $k \geq 0$.

Analysis of the Variant

Thus, we come to the following process.

Constant step scheme, III

0. Choose $\mathbf{y}_0 = \mathbf{x}_0 \in \mathbb{R}^n$.

1. k -th iteration ($k \geq 0$).

$$\begin{aligned}\mathbf{x}_{k+1} &= \mathbf{y}_k - \frac{1}{L} \nabla f(\mathbf{y}_k), \\ \mathbf{y}_{k+1} &= \mathbf{x}_{k+1} + \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} (\mathbf{x}_{k+1} - \mathbf{x}_k).\end{aligned}\tag{19}$$

However, note that this process does not work for $\mu = 0$. The choice $\gamma_0 = L$ (which changes corresponding value of α_0) is safer.

Analysis of the Variant

1 [Nesterov, 1983]

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

$$\mathbf{y}_{t+1} = \mathbf{x}_{t+1} + \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} (\mathbf{x}_{t+1} - \mathbf{x}_t).$$

Analysis of the Variant

1 Heavy-ball method [Polyak, 1964]:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \alpha \boxed{\nabla f(\mathbf{x}_t)} + \beta(\mathbf{x}_t - \mathbf{x}_{t-1}). \text{ Set } \alpha = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2}, \beta = \left(\frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^2$$

$$\mathbf{x}_{t+1} = \mathbf{y}_t - \frac{4}{(\sqrt{L} + \sqrt{\mu})^2} \boxed{\nabla f(\mathbf{y}_t)}.$$

$$\mathbf{y}_{t+1} = \mathbf{x}_{t+1} + \left(\frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^2 (\mathbf{x}_{t+1} - \mathbf{x}_t).$$

For $\mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2 \leq \left(\frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^t \|\mathbf{x}_0 - \mathbf{x}^*\|_2.$$

Analysis of the Variant

1 [Nesterov and Nemirovsky, 1988]

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

$$\mathbf{y}_{t+1} = \mathbf{x}_{t+1} + \frac{t}{t+3}(\mathbf{x}_{t+1} - \mathbf{x}_t).$$

For $\mathcal{S}_{0,L}^{1,1}(\mathbb{R}^n)$,

$$f(\mathbf{x}_t) - f^* \leq 2L \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{(t+1)^2}.$$

Analysis of the Variant

1 FISTA [Beck and Teboulle, 2009]:

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

where $\lambda_0 = 0$, and $\lambda_{t+1} = \frac{1 + \sqrt{1 + 4\lambda_t^2}}{2}$ for all $t > 0$.

For $\mathcal{S}_{0,L}^{1,1}(\mathbb{R}^n)$,

$$f(\mathbf{x}_t) - f^* \leq 2L \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{t^2}.$$

Analysis of the Variant

Lemma 14

By definition of $\mathbf{x}_{t+1} = \mathbf{y}_t - \frac{1}{L} \nabla f(\mathbf{y}_t)$, we have

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}) \leq -\frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{y}_t\|_2^2 - L \langle \mathbf{x}_{t+1} - \mathbf{y}_t, \mathbf{y}_t - \mathbf{x} \rangle.$$

Proof.

We have the following fact:

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

Remark. Please refer to the course website of Niao He:

[https://github.com/niaohel/Big-Data-Optimization-Course/blob/main/lecture_scribe/IE598-lecture9-gradient-descent-and-acceleration.](https://github.com/niaohel/Big-Data-Optimization-Course/blob/main/lecture_scribe/IE598-lecture9-gradient-descent-and-acceleration.pdf)

pdf

Analysis of the Variant

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}) \leq -\frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{y}_t\|_2^2 - L \langle \mathbf{x}_{t+1} - \mathbf{y}_t, \mathbf{y}_t - \mathbf{x} \rangle.$$

Proof. (Continued.) Using the facts that $f(\mathbf{x}_{t+1}) - f(\mathbf{y}_t) \leq -\frac{1}{2L} \|\nabla f(\mathbf{y}_t)\|_2^2$ (pp. 26 of nesterov book) and $f(\mathbf{y}_t) - f(\mathbf{x}) \leq \langle \nabla f(\mathbf{y}_t), \mathbf{y}_t - \mathbf{x} \rangle$, we have

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}) \leq -\frac{1}{2L} \|\nabla f(\mathbf{y}_t)\|_2^2 + \langle \nabla f(\mathbf{y}_t), \mathbf{y}_t - \mathbf{x} \rangle.$$

Hence, by using definition of $\mathbf{x}_{t+1} = \mathbf{y}_t - \frac{1}{L} \nabla f(\mathbf{y}_t)$, we have the result. □

Analysis of the Variant

Theorem 15 (FISTA)

For $\mathcal{S}_{0,L}^{1,1}(\mathbb{R}^n)$,

$$f(\mathbf{x}_t) - f^* \leq 2L \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{t^2}.$$

Proof. According to lemma 14, and let $\mathbf{x} = \mathbf{x}_t$ and $\mathbf{x} = \mathbf{x}^*$, respectively. We have

$$\begin{aligned} [f(\mathbf{x}_{t+1}) - f(\boxed{\mathbf{x}_t})](\lambda_t - 1) &\leq \left[-\frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{y}_t\|_2^2 + \langle \mathbf{x}_{t+1} - \mathbf{y}_t, \mathbf{y}_t - \boxed{\mathbf{x}_t} \rangle \right] (\lambda_t - 1), \\ [f(\mathbf{x}_{t+1}) - f(\boxed{\mathbf{x}^*})] &\leq \left[-\frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{y}_t\|_2^2 + \langle \mathbf{x}_{t+1} - \mathbf{y}_t, \mathbf{y}_t - \boxed{\mathbf{x}^*} \rangle \right]. \end{aligned}$$

Analysis of the Variant

Proof. (Continued.) Adding above two equations, and let $\epsilon_t = f(\mathbf{x}_t) - f(\mathbf{x}^*)$, we arrive at

$$\begin{aligned}\lambda_t f(\mathbf{x}_{t+1}) - \lambda_t f(\mathbf{x}_t) + f(\mathbf{x}_t) - f(\mathbf{x}^*) &= \lambda_t \epsilon_{t+1} - (\lambda_t - 1) \epsilon_t \\ &\leq -\frac{\lambda_t L}{2} \|\mathbf{x}_{t+1} - \mathbf{y}_t\|_2^2 \\ &\quad - L \langle (\mathbf{x}_{t+1} - \mathbf{y}_t), \lambda_t \mathbf{y}_t - (\lambda_t - 1) \mathbf{x}_t - \mathbf{x}^* \rangle.\end{aligned}$$

Multiplying both sides by λ_t :

$$\lambda_t^2 \epsilon_{t+1} - \lambda_t (\lambda_t - 1) \epsilon_t \leq -\frac{L}{2} \left[\lambda_t^2 \|\mathbf{x}_{t+1} - \mathbf{y}_t\|_2^2 + 2\lambda_t \langle (\mathbf{x}_{t+1} - \mathbf{y}_t), \lambda_t \mathbf{y}_t - (\lambda_t - 1) \mathbf{x}_t - \mathbf{x}^* \rangle \right].$$

Analysis of the Variant

Proof. (Continued.) Since $\lambda_0 = 0$, and $\lambda_{t+1} = \frac{1+\sqrt{1+4\lambda_t^2}}{2}$ for all $t \geq 0$, we know the following is true,

$$\lambda_{t+1}^2 - \lambda_{t+1} = \lambda_t^2.$$

Thus, rearranging terms, we have

$$\begin{aligned} \lambda_t^2 \epsilon_{t+1} - \lambda_{t-1}^2 \epsilon_t &\leq -\frac{L}{2} \left[\|\lambda_t \mathbf{x}_{t+1} - (\lambda_t - 1) \mathbf{x}_t - \mathbf{x}^*\|_2^2 \right. \\ &\quad \left. - \|\lambda_t \mathbf{y}_t - (\lambda_t - 1) \mathbf{x}_t - \mathbf{x}^*\|_2^2 \right]. \end{aligned}$$

Invoking the definition of \mathbf{y}_{t+1} , one can show that

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

Analysis of the Variant

Proof. (Continued.) Hence, define $\mathbf{u}_0 = \mathbf{x}_0 - \mathbf{x}^*$, $\mathbf{u}_t = \lambda_t \mathbf{y}_t - (\lambda_t - 1)\mathbf{x}_t - \mathbf{x}^*$, $\forall t \geq 1$, the above equation simplifies to

$$\lambda_t^2 \epsilon_{t+1} - \lambda_{t-1}^2 \epsilon_t \leq -\frac{L}{2} \left[\|\mathbf{u}_{t+1}\|_2^2 - \|\mathbf{u}_t\|_2^2 \right].$$

By induction, we have

$$\lambda_{t-1}^2 \epsilon_t \leq \frac{L}{2} (\|\mathbf{u}_0\|_2^2 - \|\mathbf{u}_t\|_2^2) \leq \frac{L}{2} \|\mathbf{u}_0\|_2^2.$$

Since $\lambda_{t-1} > t/2$, $\forall t \geq 1$, we have

$$\epsilon_t \leq \frac{2L \|\mathbf{u}_0\|_2^2}{t^2}.$$

References I

- Yu Nesterov. A method of solving a convex programming problem with convergence rate $o(1/k^2)$. In *Sov. Math. Dokl*, volume 27, 1983.
- Yurii Nesterov and A Nemirovsky. A general approach to polynomial-time algorithms design for convex programming. *Report, Central Economical and Mathematical Institute, USSR Academy of Sciences, Moscow*, 1988.
- Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- Yu Nesterov. Gradient methods for minimizing composite objective function. Technical report, Université catholique de Louvain, Center for Operations Research, 2007.
- Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.

References II

- Paul Tseng and Sangwoon Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1-2):387–423, 2009.
- Weijie Su, Stephen Boyd, and Emmanuel Candes. A differential equation for modeling nesterov’s accelerated gradient method: Theory and insights. In *Advances in Neural Information Processing Systems*, pages 2510–2518, 2014.
- Laurent Lessard, Benjamin Recht, and Andrew Packard. Analysis and design of optimization algorithms via integral quadratic constraints. *SIAM Journal on Optimization*, 26(1):57–95, 2016.
- Andre Wibisono, Ashia C Wilson, and Michael I Jordan. A variational perspective on accelerated methods in optimization. *proceedings of the National Academy of Sciences*, 113(47):E7351–E7358, 2016.

References III

- Zeyuan Allen-Zhu and Lorenzo Orecchia. Linear coupling: An ultimate unification of gradient and mirror descent. *arXiv preprint arXiv:1407.1537*, 2014.
- Sébastien Bubeck, Yin Tat Lee, and Mohit Singh. A geometric alternative to nesterov's accelerated gradient descent. *arXiv preprint arXiv:1506.08187*, 2015.
- Guanghui Lan and Yi Zhou. An optimal randomized incremental gradient method. *Mathematical programming*, 171(1-2):167–215, 2018.
- Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- Yurii Nesterov. *Lectures on convex optimization*, volume 137. Springer, 2018.

Thank You!

Email: qianhui@zju.edu.cn