



## Lecture 4 : Lower-bound of first-order methods and Nesterov optimal algorithm

### 1 Introduction

In previous lectures, we already discussed about one of the most well-known optimization algorithms: *gradient descent* and its theoretical guarantee (on the convergence of iterates and of the objective functions). A natural question is: is gradient descent the *best* optimization that we can use when dealing with the class of  $L$ -smooth (and convex) functions? And what is the limit of an algorithm optimizing a function of this class? These questions will be (partly) revealed in this lecture.

### 2 First-order methods and their fundamental limits

In this section, we consider the following question:

**Question 2.1.** Given an optimization algorithm  $\mathcal{A}$  and a function class  $\mathcal{F}$ , investigate the worst objective gap at the  $k$ th iteration, i.e.,

$$\ell(\mathcal{F}, \mathcal{A}) := \sup_{f \in \mathcal{F}} f(x_k) - f^*,$$

where  $x_k$  is the  $k$ th iteration generated by  $\mathcal{A}$  and  $f^* = \inf_x f(x)$ .

The sup operator in (2.1) implies that we would like find the *worst* function making the algorithm  $\mathcal{A}$  suffered the most. This gives us a lower-bound on the performance of the algorithm  $\mathcal{A}$  when optimizing a function of  $f$ .

Of course, we expect that  $\ell(\mathcal{F}, \mathcal{A})$  depends on the initialization  $x_0$ , the number of iterations  $k$  and the properties of the function class  $\mathcal{F}$ .

To this end, we consider  $\mathcal{A}$  and  $\mathcal{F}$  as follows:

1.  $\mathcal{A}$  is a first-order algorithm, i.e., an iterative method whose sequence of iterates  $x_k$  satisfies:

$$x_k \in x_0 + \text{Lin}\{\nabla f(x_0), \dots, \nabla f(x_{k-1})\}, k \geq 1.$$

In particular, gradient descent is a first-order method because:

$$\begin{aligned} x_k &= x_{k-1} - \alpha \nabla f(x_{k-1}) \\ &= x_{k-2} - \alpha \nabla f(x_{k-2}) - \alpha \nabla f(x_{k-1}) \\ &= \dots \\ &= x_0 - \alpha(\nabla f(x_0) + \dots + \nabla f(x_{k-1})). \end{aligned}$$

Attention: in general, first-order methods are much bigger since they refer to those that use only gradient.

2.  $\mathcal{F}$  is the set of convex and  $L$ -smooth functions, i.e.,  $f$  is convex and its gradient is  $L$ -Lipschitz.

Throughout this section, we will prove the following result:

**Theorem 2.2** (Lower-bound of first-order methods performance). *For any  $x_0 \in \mathbb{R}^n$  and  $1 \leq k \leq \frac{1}{2}(n - 1)$ , there exists a function  $f \in \mathcal{F}$  such that for any first-order algorithm  $\mathcal{A}$ , we have:*

$$\begin{aligned} f(x_k) - f^* &\geq \frac{3L\|x_0 - x^*\|^2}{32(k+1)^2} \\ \|x_k - x^*\|^2 &\geq \frac{1}{8}\|x_0 - x^*\|^2. \end{aligned}$$

where  $x^* \in \operatorname{argmin} f$  and  $f^* = f(x^*)$ .

The proof consists multiple steps, which are detailed below:

**Construction of the functions  $f$**  We consider a family of  $k$  quadratic functions ( $k = 1, \dots, n$ ) given by:

$$f_k(x) = \frac{L}{4} \left\{ \frac{1}{2} \left[ [x]_1^2 + \sum_{i=1}^{k-1} ([x]_i - [x]_{i+1})^2 + [x]_k^2 \right] - [x]_1 \right\}$$

where  $[x]_i$  is the  $i$ th coordinate of the vector  $x$ . We prove that  $f_k \in \mathcal{F}$ . Given a vector  $s \in \mathbb{R}^n$ , we have:

$$\begin{aligned} f_k(x+s) - f_k(x) &= \underbrace{\frac{L}{4} \left[ [s]_1[x]_1 + \sum_{i=1}^{k-1} ([s]_i - [s]_{i+1})^\top ([x]_i - [x]_{i+1}) + [x]_k[s]_k - [s]_1 \right]}_{\text{first-order term}} \\ &\quad + \underbrace{\frac{1}{2} \cdot \frac{L}{4} \left( [s]_1^2 + \sum_{i=1}^{k-1} ([s]_i - [s]_{i+1})^2 + [s]_k^2 \right)}_{\text{second-order term}} \end{aligned}$$

Therefore, we conclude that:

$$0 \leq s^\top \nabla^2 f(x)s = \frac{L}{4} \left( [s]_1^2 + \sum_{i=1}^{k-1} ([s]_i - [s]_{i+1})^2 + [s]_k^2 \right) \leq \frac{L}{4} \left( [s]_1^2 + 2 \sum_{i=1}^{k-1} ([s]_i^2 + [s]_{i+1}^2) + [s]_k^2 \right) \leq L\|s\|_2^2.$$

Therefore,  $0 \preceq \nabla^2 f(x) \preceq L\mathbf{I}$ . Therefore,  $f$  is convex and  $L$ -smooth (prove this is an exercise).

**The optimal solution of the function  $f_k$**  We compute the minimum of the function  $f_k$ . Since  $f$  is convex, it is sufficient to find  $x$  such that  $\nabla f(x) = 0$ . The gradient of  $f$  (up to a constant  $\frac{L}{4}$ ) is given by:

$$\nabla f(x) = \begin{pmatrix} 2[x]_1 - [x]_2 \\ -[x]_1 + 2[x]_2 - [x]_3 \\ \vdots \\ -[x]_{k-2} + 2[x]_{k-1} - [x]_k \\ -[x]_{k-1} + 2[x]_k \end{pmatrix} - e_1$$

Therefore, to have  $\nabla f(x) = 0$ , one deduces (by induction) that:

$$[x]_i = (k - i + 1)[x]_k,$$

and hence, we obtain:

$$[x_k^*]_i = \begin{cases} 1 - \frac{i}{k+1}, & i = 1, \dots, k \\ \mathbb{R}, & k+1 \leq i \leq n. \end{cases}$$

The optimal value of  $f_k$  is given by:

$$f_k^* = f_k(x_k^*) = -\frac{L}{8} \cdot \frac{k}{k+1} = -\frac{L}{8} \left(1 - \frac{1}{k+1}\right).$$

**The iterates generated by first-order methods  $\mathcal{A}$**  Let's fix a value  $1 \leq p \leq n$ .

**Proposition 2.3** (Properties of first-order methods). *Let  $x_0 = 0$ . For any sequence  $\{x_k\}_{k=0}^p$  generated by a first-order method  $\mathcal{A}$ , we have:*

$$x_k \in \mathbb{R}^{k,n} := \{x \in \mathbb{R}^n \mid x_i = 0, k+1 \leq i \leq n\}.$$

As a consequence, we have:  $f_p(x_k) = f_k(x_k) \geq f_k^*, \forall p = k, \dots, n$ .

*Proof.* The proof is conducted by induction and the form of the gradient of  $f$ .

Since  $f_p - f_k$  contains monomial of the form  $x_j x_l$  where  $j, l \geq k$ , we have:

$$f_p(x_k) = f_k(x_k) \geq f_k^*.$$

□

**The actual proof of Theorem 2.2** WLOG, we consider first-order methods  $\mathcal{A}$  whose initialization is  $x_0 = 0$ . Otherwise, we can use the following remark:

**Remark 2.4.** One can choose different  $x_0 \neq 0$ . However, the result remains the same by considering  $g_k(\cdot) = f_k(\cdot + x_0)$ .

We prove two claims one by one: since we assume  $k \leq \frac{1}{2}(n-1)$ , we consider the function  $f = f_{2k+1}$  and  $x^* = x_{2k+1}^*$ : We have:

$$\begin{aligned} \|x^* - x_0\| &= \sum_{i=1}^{2k+1} \left( \frac{i}{2k+2} \right)^2 = \frac{1}{(2k+2)^2} \sum_{i=1}^{2k+1} i^2 \\ &= \frac{1}{(2k+2)^2} \frac{(2k+1)(2k+2)(4k+3)}{6} \leq \frac{2k+2}{3}. \end{aligned}$$

1. **First claim:**

$$\frac{f(x_k) - f^*}{\|x_0 - x^*\|^2} \geq \frac{L}{8} \cdot \frac{\frac{1}{k+1} - \frac{1}{2k+2}}{\frac{2k+2}{3}} = \frac{3L}{32} \cdot \frac{1}{(k+1)^2}.$$

2. **Second claim:** Consider the quantity  $\|x_k - x^*\|^2$ :

$$\begin{aligned} \|x_k - x^*\|_2^2 &\geq \sum_{i=k+1}^{2k+1} ([x_{2k+1}^*]_i)^2 = \sum_{i=k+1}^{2k+1} \left(1 - \frac{i}{2k+2}\right)^2 \\ &= k+1 - \frac{1}{k+1} \sum_{i=k+1}^{2k+1} i + \frac{1}{4(k+1)^2} \sum_{i=k+1}^{2k+1} i^2 \\ &= k+1 - \frac{3k+2}{2} + \frac{(2k+1)(7k+6)}{24(k+1)} \\ &\geq \frac{2k^2 + 7k + 6}{24(k+1)} \geq \frac{2k^2 + 7k + 6}{16(k+1)^2} \|x_0 - x^*\|^2 \geq \frac{1}{8} \|x_0 - x^*\|^2. \end{aligned}$$

### 3 Optimal algorithm with Nesterov acceleration

Knowing that the lower-bound is  $O(\frac{1}{k^2})$ , we wonder if there exists an optimization attaining this limit. The answer is affirmative, which is the famous Nesterov acceleration method. It is defined as follows:

$$\begin{aligned}\lambda_k &= \frac{1 + \sqrt{1 + 4\lambda_{k-1}^2}}{2} \\ \gamma_k &= \frac{1 - \lambda_k}{\lambda_{k+1}} \\ x_{k+1} &= y_k - \frac{1}{L} \nabla f(y_k) \\ y_{k+1} &= (1 - \gamma_k)x_{k+1} + \gamma_k x_k\end{aligned}\tag{Nes-Acc}$$

where  $\lambda_0 = 0$  and  $x_0 = y_0$  is a (random) initialization. The main result of this section is the following:

**Theorem 3.1** (Theoretical guarantee of (Nes-Acc)). *Given a convex and  $L$ -smooth function  $f$ , the iterates generated by (Nes-Acc) satisfy:*

$$f(x_k) - f^* \leq \frac{2L\|y_1 - x^*\|^2}{k^2}.$$

where we assume  $x^* \in \operatorname{argmin} f$  is a non-empty set and  $f^* = f(x^*)$ .

*Proof.* Since  $f$  is convex and  $L$ -smooth, we have:

$$\begin{aligned}f\left(x - \frac{1}{L} \nabla f(x)\right) - f(y) &\leq f\left(x - \frac{1}{L} \nabla f(x)\right) - f(x) + \nabla f(x)^\top (x - y) \\ &\leq -\frac{1}{2L} \|\nabla f(x)\|^2 + \nabla f(x)^\top (x - y).\end{aligned}$$

Applying the previous results for  $x = y_k$  and  $y = x_k$ , we have:

$$\begin{aligned}f(x_{k+1}) - f(x_k) &= f\left(y_k - \frac{1}{L} \nabla f(y_k)\right) - f(x_k) \\ &\leq -\frac{1}{2L} \|\nabla f(y_k)\|^2 + \nabla f(y_k)^\top (y_k - x_k) \\ &= -\frac{L}{2} \|x_{k+1} - y_k\|^2 - L(x_{k+1} - y_k)^\top (y_k - x_k).\end{aligned}\tag{1}$$

Similarly, we apply the result for  $x = y_k$  and  $y = x^*$ :

$$f(x_{k+1}) - f(x^*) \leq -\frac{L}{2} \|x_{k+1} - y_k\|^2 - L(x_{k+1} - y_k)^\top (y_k - x^*)\tag{2}$$

Multiplying (1) by  $(\lambda_k - 1)$  and adding to (2) yields:

$$\lambda_k \underbrace{(f(x_{k+1}) - f^*)}_{\delta_{k+1}} - (\lambda_k - 1) \underbrace{(f(x_k) - f^*)}_{\delta_k} \leq -\frac{L\lambda_k}{2} \|x_{k+1} - y_k\|^2 - L(x_{k+1} - y_k)^\top (\lambda_k y_k - (\lambda_k - 1)x_k - x^*)$$

By definition, we have:  $\lambda_{k-1}^2 = \lambda_k^2 - \lambda_k$ . Therefore, multiplying the previous inequality by  $\lambda_k$ , we have:

$$\begin{aligned}\lambda_k^2 \delta_{k+1} - \lambda_{k-1}^2 \delta_k &\leq -\frac{L}{2} \left( \|\lambda_k(x_{k+1} - y_k)\|^2 + 2\lambda_k(x_{k+1} - y_k)^\top (\lambda_k y_k - (\lambda_k - 1)x_k - x^*) \right) \\ &= -\frac{L}{2} \left( \|\lambda_k x_{k+1} - (\lambda_k - 1)x_k - x^*\|^2 - \|\lambda_k y_k - (\lambda_k - 1)x_k - x^*\|^2 \right) \\ &= -\frac{L}{2} \left( \|\lambda_{k+1} y_{k+1} - (\lambda_{k+1} - 1)x_{k+1} - x^*\|^2 - \|\lambda_k y_k - (\lambda_k - 1)x_k - x^*\|^2 \right)\end{aligned}$$

Telescoping, we get:

$$\lambda_{k-1}^2 \delta_k \leq \frac{L}{2} \|y_1 - x^*\|^2.$$

The proof is finished by noticing that  $\lambda_{k-1} \geq \frac{k}{2}, \forall k \geq 1$ .  $\square$

## 4 Summary on the theoretical guarantees of first-order methods

---

Hypothesis	Gradient descent	Nesterov acceleration
<ul style="list-style-type: none"> <li>• <math>f</math> is <math>L</math>-smooth.</li> <li>• An optimal solution <math>\theta^*</math> exist.</li> <li>• Step-size <math>\alpha = \frac{1}{L}</math></li> <li>• <math>f</math> is convex.</li> </ul>	$f(x_k) - f^* = O\left(\frac{1}{k}\right)$ $x_k \rightarrow x^* \in \operatorname{argmin} f$	$f(x_k) - f^* = O\left(\frac{1}{k^2}\right)$ $x_k \rightarrow x^*$

---

**Table 1:** Comparison between gradient descent and Nesterov acceleration