

If we consider very large problems where we can not afford  $n$  number of iterations, the above theorem says that:

- The function value can be expected to decrease fast.
- The convergence to the optimal solution  $\mathbf{x}^*$  can be arbitrarily slow.

## 6.2 Lower Complexity Bound for the class $\mathcal{S}_{\mu,L}^{\infty,1}(\ell^2)$

**Gradient Based Method:** Iterative method  $\mathcal{M}$  generated by a sequence such that

$$\mathbf{x}_k \in \mathbf{x}_0 + \text{span}\{\nabla f(\mathbf{x}_0), \nabla f(\mathbf{x}_1), \dots, \nabla f(\mathbf{x}_{k-1})\}, \quad k \geq 1.$$

Let us define

$$\ell^2 := \left\{ \{x_i\}_{i=1}^{\infty} \mid \sum_{i=1}^{\infty} x_i^2 < \infty \right\}.$$

Consider the problem class as follows

<b>Model:</b>	$\min_{\mathbf{x} \in \ell^2} f(\mathbf{x})$
<b>Oracle:</b>	$f \in \mathcal{S}_{\mu,L}^{\infty,1}(\ell^2)$ Only function and gradient values are available
<b>Approximate solution:</b>	Find $\bar{\mathbf{x}} \in \mathbb{R}^n$ such that $\begin{cases} f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) < \varepsilon \\ \ \bar{\mathbf{x}} - \mathbf{x}^*\ _2^2 < \varepsilon \end{cases}$

**Theorem 6.2** For any  $\mathbf{x}_0 \in \ell^2$ , there exists a function  $f \in \mathcal{S}_{\mu,L}^{\infty,1}(\ell^2)$  such that for any gradient based method of type  $\mathcal{M}$ , we have

$$\begin{aligned} f(\mathbf{x}_k) - f(\mathbf{x}^*) &\geq \frac{\mu}{2} \left( \frac{\sqrt{L/\mu} - 1}{\sqrt{L/\mu} + 1} \right)^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2, \\ \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 &\geq \left( \frac{\sqrt{L/\mu} - 1}{\sqrt{L/\mu} + 1} \right)^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2, \end{aligned}$$

where  $\mathbf{x}^*$  is the minimum of  $f(\mathbf{x})$ .

*Proof:*

This type of methods are invariant with respect to a simultaneous shift of all objects in the space of variables. Therefore, we can assume that  $\mathbf{x}_0 = \{0\}_{i=1}^{\infty}$ .

Consider the following quadratic function

$$f_{\mu,L}(\mathbf{x}) = \frac{\mu(L/\mu - 1)}{8} \left\{ [\mathbf{x}]_1^2 + \sum_{i=1}^{\infty} ([\mathbf{x}]_i - [\mathbf{x}]_{i+1})^2 - 2[\mathbf{x}]_1 \right\} + \frac{\mu}{2} \|\mathbf{x}\|_2^2.$$

Then

$$\nabla f_{\mu,L}(\mathbf{x}) = \left( \frac{\mu(L/\mu - 1)}{4} \mathbf{A} + \mu \mathbf{I} \right) \mathbf{x} - \frac{\mu(L/\mu - 1)}{4} \mathbf{e}_1,$$

where  $\mathbf{A}$  is the same tridiagonal matrix defined in Theorem 6.1, but with infinite dimension and  $\mathbf{e}_1 \in \ell^2$  is a vector where only the first element is one.

After some calculations, we can show that  $\mu \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}$  and therefore,  $f(\mathbf{x}) \in \mathcal{S}_{\mu,L}^{\infty,1}(\ell^2)$ , due to Corollary 5.22.

The minimal optimal solution of this function is:

$$[\mathbf{x}^*]_i := q^i = \left( \frac{\sqrt{L/\mu} - 1}{\sqrt{L/\mu} + 1} \right)^i, \quad i = 1, 2, \dots$$

Then

$$\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 = \sum_{i=1}^{\infty} [\mathbf{x}^*]_i^2 = \sum_{i=1}^{\infty} q^{2i} = \frac{q^2}{1 - q^2}.$$

Now, since  $\nabla \mathbf{f}_{\mu,L}(\mathbf{x}_0) = -\frac{\mu(L/\mu-1)}{4}\mathbf{e}_1$ , and  $\mathbf{A}$  is a tridiagonal matrix,  $[\mathbf{x}_k]_i = 0$  for  $i = k+1, k+2, \dots$ , and

$$\|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \geq \sum_{i=k+1}^{\infty} [\mathbf{x}^*]_i^2 = \sum_{i=k+1}^{\infty} q^{2i} = \frac{q^{2(k+1)}}{1 - q^2} = q^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$

Finally, the first inequality follows from Corollary 5.17. ■

## 7 The Steepest Descent Method for Differentiable Convex and Differentiable Strongly Convex Functions with Lipschitz Continuous Gradients

Let us consider the steepest descent method with constant step  $h$ .

**Theorem 7.1** Let  $f \in \mathcal{F}_L^{1,1}(\mathbb{R}^n)$ , and  $0 < h < \frac{2}{L}$ . The steepest descent method with constant step generates a sequence which converges as follows:

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{2(f(\mathbf{x}_0) - f(\mathbf{x}^*))\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 + kh(2 - Lh)(f(\mathbf{x}_0) - f(\mathbf{x}^*))}.$$

*Proof:*

Denote  $r_k = \|\mathbf{x}_k - \mathbf{x}^*\|_2$ . Then

$$\begin{aligned} r_{k+1}^2 &= \|\mathbf{x}_k - \mathbf{x}^* - h\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2 \\ &= r_k^2 - 2h\langle \nabla \mathbf{f}(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle + h^2\|\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2 \\ &= r_k^2 - 2h\langle \nabla \mathbf{f}(\mathbf{x}_k) - \nabla \mathbf{f}(\mathbf{x}^*), \mathbf{x}_k - \mathbf{x}^* \rangle + h^2\|\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2 \\ &\leq r_k^2 - h\left(\frac{2}{L} - h\right)\|\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2, \end{aligned}$$

where the last inequality follows from Theorem 5.13.

Therefore, since  $0 < h < \frac{2}{L}$ ,  $r_{k+1} < r_k < \dots < r_0$ .

Now

$$\begin{aligned} f(\mathbf{x}_{k+1}) &\leq f(\mathbf{x}_k) + \langle \nabla \mathbf{f}(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{L}{2}\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \\ &= f(\mathbf{x}_k) - h\|\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2 + \frac{L}{2}\| -h\nabla \mathbf{f}(\mathbf{x}_k) \|^2 \end{aligned} \tag{10}$$

$$= f(\mathbf{x}_k) - \omega\|\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2 < f(\mathbf{x}_k), \tag{11}$$

where  $\omega = h(1 - \frac{L}{2}h)$ . Denoting by  $\Delta_k = f(\mathbf{x}_k) - f(\mathbf{x}^*)$ , from the convexity of  $f(\mathbf{x})$ , Theorem 5.7, and the Cauchy-Schwarz inequality,

$$\Delta_k = f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \langle \nabla \mathbf{f}(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle \leq \|\nabla \mathbf{f}(\mathbf{x}_k)\|_2 r_k \leq \|\nabla \mathbf{f}(\mathbf{x}_k)\|_2 r_0. \tag{12}$$

Combining (11) and (12),

$$\Delta_{k+1} \leq \Delta_k - \frac{\omega}{r_0^2} \Delta_k^2.$$

Thus dividing by  $\Delta_k \Delta_{k+1}$ ,

$$\frac{1}{\Delta_{k+1}} \geq \frac{1}{\Delta_k} + \frac{\omega}{r_0^2} \frac{\Delta_k}{\Delta_{k+1}} \geq \frac{1}{\Delta_k} + \frac{\omega}{r_0^2}.$$

since  $\frac{\Delta_k}{\Delta_{k+1}} \geq 1$ . Summing up these inequalities we get

$$\frac{1}{\Delta_{k+1}} \geq \frac{1}{\Delta_0} + \frac{\omega}{r_0^2} (k+1).$$

■

To obtain the optimal step size, it is sufficient to find the maximum of the function  $\omega := \omega(h) = h(1 - \frac{L}{2}h)$  which is  $h^* := 1/L$ .

**Corollary 7.2** If  $f \in \mathcal{F}_L^{1,1}(\mathbb{R}^n)$ , the steepest descent method with constant step  $h = 1/L$  yields

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{2L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{k+4}.$$

That is,  $\{f(\mathbf{x}_k)\}_{k=0}^\infty$  converges  $R$ -sublinearly to  $f(\mathbf{x}^*)$ .

*Proof:*

Left for exercise. ■

**Theorem 7.3** Let  $f \in \mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$ , and  $0 < h \leq \frac{2}{\mu+L}$ . The steepest descent method with constant step generates a sequence which converges as follows:

$$\begin{aligned} \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 &\leq \left(1 - \frac{2h\mu L}{\mu + L}\right)^k \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2, \\ f(\mathbf{x}_k) - f(\mathbf{x}^*) &\leq \frac{L}{2} \left(1 - \frac{2h\mu L}{\mu + L}\right)^k \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2. \end{aligned}$$

If  $h = \frac{2}{\mu+L}$ , then

$$\begin{aligned} f(\mathbf{x}_k) - f(\mathbf{x}^*) &\leq \frac{L}{2} \left(\frac{L/\mu - 1}{L/\mu + 1}\right)^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2, \\ \|\mathbf{x}_k - \mathbf{x}^*\|_2 &\leq \left(\frac{L/\mu - 1}{L/\mu + 1}\right)^k \|\mathbf{x}_0 - \mathbf{x}^*\|_2. \end{aligned}$$

That is,  $\{\mathbf{x}_k\}_{k=0}^\infty$  and  $\{f(\mathbf{x}_k)\}_{k=0}^\infty$  converges  $R$ -linearly to  $\mathbf{x}^*$  and  $f(\mathbf{x}^*)$ , respectively.

*Proof:*

Denote  $r_k = \|\mathbf{x}_k - \mathbf{x}^*\|_2$ . Then

$$\begin{aligned} r_{k+1}^2 &= \|\mathbf{x}_k - \mathbf{x}^* - h \nabla f(\mathbf{x}_k)\|_2^2 \\ &= r_k^2 - 2h \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle + h^2 \|\nabla f(\mathbf{x}_k)\|_2^2 \\ &= r_k^2 - 2h \langle \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*), \mathbf{x}_k - \mathbf{x}^* \rangle + h^2 \|\nabla f(\mathbf{x}_k)\|_2^2 \\ &\leq r_k^2 - 2h \left( \frac{\mu L}{\mu + L} r_k^2 + \frac{1}{\mu + L} \|\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*)\|_2^2 \right) + h^2 \|\nabla f(\mathbf{x}_k)\|_2^2 \\ &= \left(1 - \frac{2h\mu L}{\mu + L}\right) r_k^2 + h \left(h - \frac{2}{\mu + L}\right) \|\nabla f(\mathbf{x}_k)\|_2^2 \end{aligned}$$

from Theorems 5.13 and 5.23, and it proves the first two inequalities.

Now, for  $h = 2/(L + \mu)$  and again from Theorem 5.13,

$$\begin{aligned} f(\mathbf{x}_k) - f(\mathbf{x}^*) - \langle \nabla f(\mathbf{x}^*), \mathbf{x}_k - \mathbf{x}^* \rangle &\leq \frac{L}{2} \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \\ &\leq \frac{L}{2} \left( \frac{L/\mu - 1}{L/\mu + 1} \right)^{2k} r_0^2. \end{aligned}$$

■

**Theorem 7.4 (Yuan 2010)** <sup>2</sup> In the special case of a strongly convex quadratic function  $f(\mathbf{x}) = \frac{1}{2} \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle + \langle \mathbf{a}, \mathbf{x} \rangle + \alpha$  with  $\lambda_1(\mathbf{A}) = L \geq \lambda_n(\mathbf{A}) = \mu > 0$ , we can obtain

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

for the steepest descent method with “exact line search”.

- Note that the previous result for the steepest descent method, Theorem 4.18, was only a local result. Theorems 7.1 and 7.3 guarantee that the steepest descent method converges for any starting point  $\mathbf{x}_0 \in \mathbb{R}^n$  (due to convexity).
- Comparing the rate of convergence of the steepest descent method for the classes  $\mathcal{F}_L^{1,1}(\mathbb{R}^n)$  and  $\mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$  (Theorems 7.1, Corollary 7.2, and 7.3, respectively) with their lower complexity bounds (Theorems 6.1 and 6.2, respectively), we possibly have a huge gap.

## 7.1 Exercises

1. Prove Corollary 7.2.
2. Consider a sequence  $\{\beta_k\}_{k=0}^\infty$  which converges to zero.

The sequence is said to converge *Q-sublinearly* if

$$\lim_{k \rightarrow \infty} \sup \left| \frac{\beta_{k+1}}{\beta_k} \right| = 1.$$

A zero converging sequence  $\{\beta_k\}_{k=0}^\infty$  is said to converge *R-sublinearly* if it is dominated by a Q-sublinearly converging sequence. That is, if there is a Q-sublinearly converging sequence  $\{\hat{\beta}_k\}_{k=0}^\infty$  such that  $0 \leq |\beta_k| \leq \hat{\beta}_k$ .

- (a) Give an example of a Q-sublinear converging sequence which is not Q-linear converging sequence.
- (b) Give an example of a R-sublinear converging sequence which is not R-linear converging sequence.

## 8 The Optimal Gradient Method (First-Order Method, Accelerated Gradient Method, Fast Gradient Method)

This algorithm was proposed for the first time by Nesterov<sup>3</sup> in 1983. In [Nesterov03], he gives a reinterpretation of the algorithm and provides another justification of it which attains the same complexity bound of the original article.

<sup>2</sup>Y.-X. Yuan, “A short note on the Q-linear convergence of the steepest descent method”, *Mathematical Programming* **123** (2010), pp. 339–343.

<sup>3</sup>Y. Nesterov, “A method for solving the convex programming problem with convergence rate  $\mathcal{O}(1/k^2)$ ,” *Dokl. Akad. Nauk SSSR* **269** (1983), pp. 543–547.