

Theorem 7.2 For any $\mathbf{x}_0 \in \mathbb{R}^\infty$, there exists a function $f \in \mathcal{S}_{\mu,L}^{\infty,1}(\mathbb{R}^\infty)$ 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 7.1, but with infinite dimension and $\mathbf{e}_1 \in \mathbb{R}^\infty$ 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, $f(\mathbf{x}) \in \mathcal{S}_{\mu,L}^{\infty,1}(\mathbb{R}^\infty)$, due to Corollary 6.21.

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 $\nabla f_{\mu,L}(\mathbf{x}_0) = -\frac{\mu(L/\mu - 1)}{4} \mathbf{e}_1$, 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 6.17. ■

8 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 8.1 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 6.13.

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

Now

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

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 6.7, and the Cauchy-Schwarz inequality,

$$\Delta_k = f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \langle \nabla \mathbf{f}(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle \leq \|\nabla \mathbf{f}(\mathbf{x}_k)\|_2 r_k \leq \|\nabla \mathbf{f}(\mathbf{x}_k)\|_2 r_0. \tag{11}$$

Combining (10) and (11),

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

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 8.2 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 8.3 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 6.13 and 6.22, and it proves the first two inequalities.

Now, for $h = 2/(L + \mu)$ and again from Theorem 6.13,

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

■

Theorem 8.4 (Yuan 2010) ² 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 5.12, was only a local result. Theorems 8.1 and 8.3 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 8.1, Corollary 8.2, and 8.3, respectively) with their lower complexity bounds (Theorems 7.1 and 7.2, respectively), we possible have a huge gap.

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

8.1 Exercises

1. Prove Corollary 8.2.
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) ~~Show that a Q-linear converging sequence is a Q-sublinear converging sequence.~~
- (b) Give an example of a Q-sublinear converging sequence which is not Q-linear converging sequence.
- (c) Give an example of a R-sublinear converging sequence which is not R-linear converging sequence.

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

This algorithm was proposed for the first time by Nesterov³ 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.

Definition 9.1 A pair of sequences $\{\phi_k(\mathbf{x})\}_{k=0}^\infty$ and $\{\lambda_k\}_{k=0}^\infty$ with $\lambda_k \geq 0$ is called an *estimate sequence* of the function $f(\mathbf{x})$ if

$$\lambda_k \rightarrow 0,$$

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

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

Lemma 9.2 Given an estimate sequence $\{\phi_k(\mathbf{x})\}_{k=0}^\infty$, $\{\lambda_k\}_{k=0}^\infty$, and if for some sequence $\{\mathbf{x}_k\}_{k=0}^\infty$ we have

$$f(\mathbf{x}_k) \leq \phi_k^* := \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x})$$

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

Proof:

It follows from the definition. ■

Lemma 9.3 Assume that

1. $f \in \mathcal{S}_\mu^1(\mathbb{R}^n)$, possible with $\mu = 0$ (which means that $f \in \mathcal{F}^1(\mathbb{R}^n)$).
2. $\phi_0(\mathbf{x})$ is an arbitrary function on \mathbb{R}^n .
3. $\{\mathbf{y}_k\}_{k=0}^\infty$ is an arbitrary sequence in \mathbb{R}^n .

³Y. 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.