Scheme for the Optimal Gradient Method" is an optimal method in terms of complexity for the dominant term $\ln(\varepsilon^{-1})$.

Remark 9.8 Many times, you will find in articles that a method has "optimal rate of convergence". In our case, if we apply the "General Scheme for the Optimal Gradient Method" to $\min_{\boldsymbol{x}\in\mathbb{R}^n} f(\boldsymbol{x})$, the number of iterations of this method to obtain $f(\boldsymbol{x}_k) - f(\boldsymbol{x}^*) < \varepsilon$ is $k = k(L, \boldsymbol{x}_0, \boldsymbol{x}^*, \varepsilon) = \mathcal{O}\left(\sqrt{\frac{L\|\boldsymbol{x}_0-\boldsymbol{x}^*\|_2^2}{\varepsilon}}\right)$ and $k = k(L, \mu, \boldsymbol{x}_0, \boldsymbol{x}^*, \varepsilon) = \mathcal{O}\left(\sqrt{\frac{L}{\mu} \ln \frac{L\|\boldsymbol{x}_0-\boldsymbol{x}^*\|_2^2}{\varepsilon}}\right)$ for $f(\boldsymbol{x}) \in \mathcal{F}_L^{1,1}(\mathbb{R}^n)$ and $\mathcal{S}_{L,\mu}^{1,1}(\mathbb{R}^n)$, respectively.

It is <u>extremely important</u> to note that this value is the maximum number of iterations in the worse case scenario. To obtain the total complexity of the method, you need to <u>multiply</u> the above number by the number of floating-point operations per iteration. This value also vary according to the method.

Now, instead of doing line search at Step 4 of the General Scheme for the Optimal Gradient Method, let us consider the constant step size iteration $\mathbf{x}_{k+1} := \mathbf{y}_k - \frac{1}{L} \nabla \mathbf{f}(\mathbf{y}_k)$ (see proof of Theorem 9.5). From the calculations given at Exercise 1, we arrive to the following simplified scheme. Hereafter, we assume that $L > \mu$ to exclude the trivial case $L = \mu$ with finished in one iteration.

Observe that the sequences $\{\boldsymbol{x}_k\}_{k=0}^{\infty}$ and $\{\boldsymbol{y}_k\}_{k=0}^{\infty}$ generated by the "General Scheme" and the "Constant Step Scheme for the Optimal Gradient Methods" are exactly the same⁴ if we choose $\boldsymbol{x}_{k+1} := \boldsymbol{y}_k - \frac{1}{L} \nabla \boldsymbol{f}(\boldsymbol{y}_k)$ in the former method. Therefore, the result of Theorem 9.6 is still valid for $\gamma_0 := \alpha_0 (\alpha_0 L - \mu)/(1 - \alpha_0)$.

Also, if we further impose $\gamma_0 = \alpha_0(\alpha_0 L - \mu)/(1 - \alpha_0) = L$, we will have the rate of convergence of Theorem 9.7.

9.1 Discussion on Particular Cases

9.1.1 Accelerated Gradient Method for Smooth (Differentiable) Strongly Convex Functions

In this case, we have $\mu > 0$ and choosing $\gamma_0 := \alpha_0(\alpha_0 L - \mu)/(1 - \alpha_0) = \mu$, we can have further simplifications:

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

⁴ strictly speaking, there is a one index difference between y_k 's on these two methods due to the order y_k is defined in the loop.

Accelerated Gradient Method for Smooth Strongly Convex FunctionStep 0:Choose $x_0 \in \mathbb{R}^n$, set $y_0 := x_0$ and k := 0.Step 1:Compute $\nabla f(y_k)$.Step 2:Set $x_{k+1} := y_k - \frac{1}{L} \nabla f(y_k)$.Step 3:Set $y_{k+1} := x_{k+1} + \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} (x_{k+1} - x_k), k := k + 1$ and go to Step 1.

9.1.2 Accelerated Gradient Method for Smooth (Differentiable) Convex Functions

In the case $\mu = 0$, there are much simpler variation of the method⁵.

Nesterov's Accelerated Gradient Method for Smooth Convex Function Step 0: Choose $\boldsymbol{x}_0 \in \mathbb{R}^n$, set $\boldsymbol{y}_0 := \boldsymbol{x}_0$, $t_0 := 1$, and k := 0. Step 1: Compute $\nabla \boldsymbol{f}(\boldsymbol{y}_k)$. Step 2: Set $\boldsymbol{x}_{k+1} := \boldsymbol{y}_k - \frac{1}{L} \nabla \boldsymbol{f}(\boldsymbol{y}_k)$. Step 3: $t_{i+1} := \frac{1 + \sqrt{1 + 4t_i^2}}{2}$. Step 4: Set $\boldsymbol{y}_{k+1} := \boldsymbol{x}_{k+1} + \frac{t_i - 1}{t_{i+1}} (\boldsymbol{x}_{k+1} - \boldsymbol{x}_k)$, k := k + 1 and go to Step 1.

Moreover, this is equivalent to the following update as well.

Nesterov's Accelerated Gradient Method for Smooth Convex Function Step 0: Choose $x_0 \in \mathbb{R}^n$, set $y_0 := x_0$ and k := 1. Step 1: Compute $\nabla f(y_{k-1})$. Step 2: Set $x_k := y_{k-1} - \frac{1}{L} \nabla f(y_{k-1})$. Step 3: Set $y_k := x_k + \frac{k-1}{k+2}(x_k - x_{k-1}), k := k+1$ and go to Step 1.

The Nesterov's Accelerated Gradient Method for $f \in \mathcal{F}_L^{1,1}(\mathbb{R}^n)$ generates a sequence $\{x_k\}_{k=0}^{\infty}$ such that

$$f(\boldsymbol{x}_k) - f(\boldsymbol{x}^*) \le \frac{2L \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2^2}{(k+1)^2}.$$

Recently, it was shown that an extension of this method guarantee a $o(k^{-2})$ convergence for $f(\boldsymbol{x}_k) - f(\boldsymbol{x}^*)$ by Attouch and Peypouquet⁶.

9.2 Exercises

1. We want to justify the Constant Step Scheme of the Optimal Gradient Method. This is a particular case of the General Scheme for the Optimal Gradient Method for the following choice:

$$\begin{split} \gamma_{k+1} &:= L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu \\ \boldsymbol{y}_k &= \frac{\alpha_k\gamma_k\boldsymbol{v}_k + \gamma_{k+1}\boldsymbol{x}_k}{\gamma_k + \alpha_k\mu} \\ \boldsymbol{x}_{k+1} &= \boldsymbol{y}_k - \frac{1}{L}\boldsymbol{\nabla}\boldsymbol{f}(\boldsymbol{y}_k) \\ \boldsymbol{v}_{k+1} &= \frac{(1 - \alpha_k)\gamma_k\boldsymbol{v}_k + \alpha_k\mu\boldsymbol{y}_k - \alpha_k\boldsymbol{\nabla}\boldsymbol{f}(\boldsymbol{y}_k)}{\gamma_{k+1}} \end{split}$$

⁵Y. 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. It also has a scheme to estimate L in the case this constant in unknown.

⁶Hedy Attouch and Juan Peypouquet, "The rate of convergence of Nesterovs accelerated forward-backward method is actually faster than $1/k^2$," SIAM Journal on Optimization **26** (2016), pp. 1824-1834.

(a) Show that v_{k+1} = x_k + ¹/_{α_k}(x_{k+1} - x_k).
(b) Show that y_{k+1} = x_{k+1} + β_k(x_{k+1} - x_k) for β_k = ^{α_{k+1}γ_{k+1}(1-α_k)}/_{α_k(γ_{k+1}+α_{k+1}μ)}.
(c) Show that β_k = ^{α_k(1-α_k)}/_{α_k²+α_{k+1}}.
(d) Explain why α²_{k+1} = (1 - α_{k+1})α²_k + ^μ/_Lα_{k+1}.

10 Extension of the Optimal Gradient Method (First-Order Method, Accelerated Gradient Method, Fast Gradient Method) for the Min-Max Problems over Simple Closed Convex Sets

Suppose we are given Q a <u>closed convex</u> subset of \mathbb{R}^n , <u>simple enough</u> to have an easy projection onto it. *E.g.*, positive orthant, *n*-dimensional box, simplex, Euclidean ball, ellipsoids, *etc.*

Given $f_i \in \mathcal{S}_{\mu,L}^{1,1}(Q)$ (i = 1, 2, ..., m), we define the following function $f: Q \to \mathbb{R}$,

$$f(\boldsymbol{x}) := \max_{1 \le i \le m} f_i(\boldsymbol{x}) \quad \text{for} \quad \boldsymbol{x} \in Q.$$
(15)

This function is non-differentiable in general, but convex (see Theorem 6.6). We will see that the method discussed so far can be easily adapted for the following min-max-type convex optimization problem.

$$\begin{cases} \text{minimize} & f(\boldsymbol{x}) \\ \text{subject to} & \boldsymbol{x} \in Q, \end{cases}$$
(16)

where Q is a closed convex set with a simple structure, and f(x) is defined as above.

For a given $\bar{x} \in Q$, let us define the following linearization of f(x) at \bar{x} .

$$f(\bar{\boldsymbol{x}}; \boldsymbol{x}) := \max_{1 \le i \le m} \left[f_i(\bar{\boldsymbol{x}}) + \langle \boldsymbol{\nabla} \boldsymbol{f}_i(\bar{\boldsymbol{x}}), \boldsymbol{x} - \bar{\boldsymbol{x}} \rangle
ight], \quad \text{for } \boldsymbol{x} \in \mathbb{R}^n.$$

Lemma 10.1 Let $f_i \in S^{1,1}_{\mu,L}(Q)$ (i = 1, 2, ..., m). For $x \in Q$, we have

$$egin{aligned} f(m{x}) &\geq f(m{x};m{x}) + rac{\mu}{2} \|m{x} - m{x}\|_2^2, \ f(m{x}) &\leq f(m{x};m{x}) + rac{L}{2} \|m{x} - m{x}\|_2^2. \end{aligned}$$

Proof:

It follows from the properties of $f_i \in \mathcal{S}_{\mu,L}^{1,1}(Q)$.

Theorem 10.2 A point $x^* \in Q$ is an optimal solution of (16) with $f_i \in S^{1,1}_{\mu,L}(Q)$ (i = 1, 2, ..., m) if and only if

$$f(\boldsymbol{x}^*; \boldsymbol{x}) \geq f(\boldsymbol{x}^*; \boldsymbol{x}^*) = f(\boldsymbol{x}^*), \quad \forall \boldsymbol{x} \in Q.$$

Proof:

Indeed, if the inequality is true, it follows from Lemma 10.1 that

$$f(\boldsymbol{x}) \ge f(\boldsymbol{x}^*; \boldsymbol{x}) + \frac{\mu}{2} \|\boldsymbol{x} - \boldsymbol{x}^*\|_2^2 \ge f(\boldsymbol{x}^*) + \frac{\mu}{2} \|\boldsymbol{x} - \boldsymbol{x}^*\|_2^2 \ge f(\boldsymbol{x}^*), \quad \forall \boldsymbol{x} \in Q.$$

For the converse, let \mathbf{x}^* be an optimal solution of the minimization problem (16). Assume by contradiction that there is a $\mathbf{x} \in Q$ such that $f(\mathbf{x}^*; \mathbf{x}) < f(\mathbf{x}^*)$.