

General Scheme for the “Optimal” Gradient Method	
Step 0:	Choose $\mathbf{x}_0 \in \mathbb{R}^n$, let $\gamma_0 > 0$ such that $L \geq \gamma_0 \geq \mu \geq 0$. Set $\mathbf{v}_0 := \mathbf{x}_0$ and $k := 0$.
Step 1:	Compute $\alpha_k \in (0, 1]$ from the equation $L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu$.
Step 2:	Set $\gamma_{k+1} := (1 - \alpha_k)\gamma_k + \alpha_k\mu$, $\mathbf{y}_k := \frac{\alpha_k\gamma_k\mathbf{v}_k + \gamma_{k+1}\mathbf{x}_k}{\gamma_k + \alpha_k\mu}$.
Step 3:	Compute $f(\mathbf{y}_k)$ and $f'(\mathbf{y}_k)$.
Step 4:	Find \mathbf{x}_{k+1} such that $f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}_k) - \frac{1}{2L}\ f'(\mathbf{y}_k)\ _2^2$ using “line search”.
Step 5:	Set $\mathbf{v}_{k+1} := \frac{(1-\alpha_k)\gamma_k\mathbf{v}_k + \alpha_k\mu\mathbf{y}_k - \alpha_k f'(\mathbf{y}_k)}{\gamma_{k+1}}$, $k := k + 1$ and go to Step 1.

Theorem 9.6 Consider $f \in \mathcal{S}_{\mu, L}^{1,1}(\mathbb{R}^n)$, possible with $\mu = 0$ (which means that $f \in \mathcal{F}_L^{1,1}(\mathbb{R}^n)$). The general scheme of the “optimal” gradient method generates a sequence $\{\mathbf{x}_k\}_{k=0}^\infty$ such that

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \lambda_k \left[f(\mathbf{x}_0) + \frac{\gamma_0}{2} \|\mathbf{x}^* - \mathbf{x}_0\|_2^2 - f(\mathbf{x}^*) \right],$$

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

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

Proof:

The first part is obvious from the definition and Lemma 9.2.

We already know that $\alpha_k \geq \sqrt{\frac{\mu}{L}}$, therefore,

$$\lambda_k = \prod_{i=-1}^{k-1} (1 - \alpha_i) = \prod_{i=0}^{k-1} (1 - \alpha_i) \leq \left(1 - \sqrt{\frac{\mu}{L}} \right)^k,$$

which only has a meaning if $\mu > 0$. For the case $\mu = 0$, let us prove first that $\gamma_k = \gamma_0 \lambda_k$. Obviously $\gamma_0 = \gamma_0 \lambda_0$, and assuming the induction hypothesis,

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

Therefore, $L\alpha_k^2 = \gamma_{k+1} = \gamma_0 \lambda_{k+1}$. Since λ_k is a decreasing sequence

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

Thus

$$\frac{1}{\sqrt{\lambda_k}} \geq 1 + \frac{k}{2} \sqrt{\frac{\gamma_0}{L}}$$

and we have the result. ■

Theorem 9.7 Consider $f \in \mathcal{S}_{\mu, L}^{1,1}(\mathbb{R}^n)$, possible with $\mu = 0$ (which means that $f \in \mathcal{F}_L^{1,1}(\mathbb{R}^n)$). If we take $\gamma_0 = L$, the general scheme of the “optimal” gradient method generates a sequence $\{\mathbf{x}_k\}_{k=0}^\infty$ such that

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq L \min \left\{ \left(1 - \sqrt{\frac{\mu}{L}} \right)^k, \frac{4}{(k+2)^2} \right\} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$

This means that it is “optimal” for the class of functions from $\mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$ with $\mu > 0$, or $\mathcal{F}_L^{1,1}(\mathbb{R}^n)$.

In the particular case of $\mu > 0$, we have the following inequality for k sufficiently large:

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

Proof:

The first inequality follows from the previous theorem, $f(\mathbf{x}_0) - f(\mathbf{x}^*) \leq \langle f'(\mathbf{x}^*), \mathbf{x}_0 - \mathbf{x}^* \rangle + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$, and the fact that $f'(\mathbf{x}^*) = \mathbf{0}$.

Let us analyze first the case when $\mu > 0$. From Theorem 7.2, we know that we can find functions such that

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \geq \frac{\mu}{2} \left(\frac{\sqrt{L/\mu} - 1}{\sqrt{L/\mu} + 1} \right)^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \geq \frac{\mu}{2} \exp \left(-\frac{4k}{\sqrt{L/\mu} - 1} \right) \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2,$$

where the second inequality follows from $\ln(\frac{a-1}{a+1}) = -\ln(\frac{a+1}{a-1}) \geq 1 - \frac{a+1}{a-1} \geq -\frac{2}{a-1}$, for $a \in (1, +\infty)$. Therefore, the worst case bound to find \mathbf{x}_k such that $f(\mathbf{x}_k) - f(\mathbf{x}^*) < \varepsilon$ can not be better than

$$k > \frac{\sqrt{L/\mu} - 1}{4} \left(\ln \frac{1}{\varepsilon} + \ln \frac{\mu}{2} + 2 \ln \|\mathbf{x}_0 - \mathbf{x}^*\|_2 \right).$$

On the other hand, from the above result

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

where the second inequality follows from $\ln(1-a) \leq -a$, $a < 1$. Therefore, we can guarantee that $k > \sqrt{L/\mu} (\ln \frac{1}{\varepsilon} + \ln L + 2 \ln \|\mathbf{x}_0 - \mathbf{x}^*\|_2)$.

For the case $\mu = 0$, the conclusion is obvious from Theorem 7.1.

Finally, for $\mu > 0$, since $\frac{\mu}{2} \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 + f(\mathbf{x}^*) \leq f(\mathbf{x}_k)$ from the definition, we have the second inequality. ■

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 $\mathbf{x}_{k+1} := \mathbf{y}_k - \frac{1}{L} f'(\mathbf{y}_k)$ (See proof of Theorem 9.5). From the calculations given at Exercise 2, 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.

Constant Step Scheme I for the “Optimal” Gradient Method	
Step 0:	Choose $\mathbf{x}_0 \in \mathbb{R}^n$, $\alpha_0 \in (0, 1)$ such that $\mu \leq \frac{\alpha_0(\alpha_0 L - \mu)}{1 - \alpha_0} \leq L$, set $\mathbf{y}_0 := \mathbf{x}_0$ and $k := 0$.
Step 1:	Compute $f'(\mathbf{y}_k)$.
Step 2:	Set $\mathbf{x}_{k+1} := \mathbf{y}_k - \frac{1}{L} f'(\mathbf{y}_k)$.
Step 3:	Compute $\alpha_{k+1} \in (0, 1)$ from the equation $\alpha_{k+1}^2 = (1 - \alpha_{k+1})\alpha_k^2 + \frac{\mu}{L}\alpha_{k+1}$.
Step 4:	Set $\beta_k := \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}$.
Step 5:	Set $\mathbf{y}_{k+1} := \mathbf{x}_{k+1} + \beta_k(\mathbf{x}_{k+1} - \mathbf{x}_k)$, $k := k + 1$ and go to Step 1.

Observe that the sequences generated by the General Scheme and the Constant Step Scheme I for the “Optimal” Gradient Methods are different. However, the rate of convergence of the above method is similar to Theorem 9.6 for $\gamma_0 := \alpha_0(\alpha_0 L - \mu)/(1 - \alpha_0)$. 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:

$$f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*) \leq L \min \left\{ \left(1 - \sqrt{\frac{\mu}{L}}\right)^k, \frac{4}{(k+2)^2} \right\} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$