

MA 5037: Optimization Methods and Applications

The Gradient Method



Suh-Yuh Yang (楊肅煜)

Department of Mathematics, National Central University
Jhongli District, Taoyuan City 320317, Taiwan

First version: May 21, 2018/Last updated: June 15, 2025

Descent direction methods

We consider the *unconstrained minimization problem*:

$$\min\{f(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n\},$$

where *the objective function f is continuously differentiable over \mathbb{R}^n* . We will consider an iterative algorithm for finding *stationary points of f* . The iterative algorithm takes the form

$$\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k, \quad k = 0, 1, \dots,$$

where \mathbf{d}_k is the direction and t_k is the stepsize.

Definition: Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function. A vector $\mathbf{0} \neq \mathbf{d} \in \mathbb{R}^n$ is called a *descent direction of f at \mathbf{x}* if the directional derivative $f'(\mathbf{x}; \mathbf{d}) < 0$. (Note that $f'(\mathbf{x}; \mathbf{d}) = \nabla f(\mathbf{x})^\top \mathbf{d}$)

Descent property: If \mathbf{d} is a descent direction of f at \mathbf{x} , then $\exists \varepsilon > 0$ such that $f(\mathbf{x} + t\mathbf{d}) < f(\mathbf{x})$ for any $t \in (0, \varepsilon]$. \square

Taking small enough steps along these descent directions lead to a decrease of the objective function.

Schematic descent direction method

Initialization: Pick $\mathbf{x}_0 \in \mathbb{R}^n$.

General step: For any $k = 0, 1, \dots$, set

- (a) Pick a descent direction \mathbf{d}_k .
 - (b) Find a stepsize t_k satisfying $f(\mathbf{x}_k + t_k \mathbf{d}_k) < f(\mathbf{x}_k)$.
 - (c) Set $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$
 - (d) If a stopping criterion is satisfied then stop, \mathbf{x}_{k+1} is the output.
-

The descent direction method remains “*conceptual*” and cannot be implemented. Many details are missing in the above description:

- *What is the starting point \mathbf{x}_0 ?*
- *How to choose the descent direction \mathbf{d}_k ?*
- *What stepsize should be taken t_k ?*
- *What is the stopping criterion?*

Three popular choices of stepsize t_k

The process of finding t_k is called *line search*, since it is essentially a minimization procedure on the 1-D function $g(t) := f(\mathbf{x}_k + t\mathbf{d}_k)$.

- **constant stepsize:** $t_k = \bar{t}$ for any k .
- **exact line search:** t_k is a minimizer of f along the ray $\mathbf{x}_k + t\mathbf{d}_k$:

$$t_k \in \arg \min_{t \geq 0} f(\mathbf{x}_k + t\mathbf{d}_k).$$

- **backtracking:** The method requires three parameters: $s > 0$ (*not too small*), $\alpha, \beta \in (0, 1)$.

```
set  $t_k \leftarrow s$ 
while  $f(\mathbf{x}_k) - f(\mathbf{x}_k + t_k\mathbf{d}_k) < -\alpha t_k \overbrace{\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k}^{f'(\mathbf{x}_k; \mathbf{d}_k)}$  do
    set  $t_k \leftarrow \beta t_k$ 
```

Therefore, the stepsize is chosen as $t_k = s\beta^{i_k}$, where i_k is the smallest nonnegative integer for which (\star) is satisfied:

$$f(\mathbf{x}_k) - f(\mathbf{x}_k + s\beta^{i_k}\mathbf{d}_k) \geq -\alpha s\beta^{i_k} \nabla f(\mathbf{x}_k)^\top \mathbf{d}_k. \quad (\star)$$

The third option is in a sense a compromise between the other twos.

Validity of the sufficient decrease condition (★)

Theorem: Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function and $\mathbf{x} \in \mathbb{R}^n$. Assume that $\mathbf{0} \neq \mathbf{d} \in \mathbb{R}^n$ is a descent direction of f at \mathbf{x} and let $\alpha \in (0, 1)$. Then $\exists \varepsilon > 0$ such that for all $t \in [0, \varepsilon]$, we have

$$f(\mathbf{x}) - f(\mathbf{x} + t\mathbf{d}) \geq -\alpha t \nabla f(\mathbf{x})^\top \mathbf{d}.$$

Proof: Since f is continuously differentiable it follows that

$$f(\mathbf{x} + t\mathbf{d}) = f(\mathbf{x}) + t \nabla f(\mathbf{x})^\top \mathbf{d} + o(t\|\mathbf{d}\|),$$

and hence

$$f(\mathbf{x}) - f(\mathbf{x} + t\mathbf{d}) = -\alpha t \nabla f(\mathbf{x})^\top \mathbf{d} - (1 - \alpha)t \nabla f(\mathbf{x})^\top \mathbf{d} - o(t\|\mathbf{d}\|).$$

Since \mathbf{d} is a descent direction of f at \mathbf{x} , we have

$$\lim_{t \rightarrow 0^+} \frac{-(1 - \alpha)t \nabla f(\mathbf{x})^\top \mathbf{d} - o(t\|\mathbf{d}\|)}{t} = -(1 - \alpha) \nabla f(\mathbf{x})^\top \mathbf{d} > 0.$$

Hence, $\exists \varepsilon > 0$ such that for all $t \in (0, \varepsilon]$, we have

$$-(1 - \alpha)t \nabla f(\mathbf{x})^\top \mathbf{d} - o(t\|\mathbf{d}\|) > 0,$$

which implies the desired result. \square

Example: exact line search for quadratic functions

Let $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} + 2\mathbf{b}^\top \mathbf{x} + c$, where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix, $\mathbf{b} \in \mathbb{R}^n$, and $c \in \mathbb{R}$. Let $\mathbf{x} \in \mathbb{R}^n$ and let $\mathbf{d} \in \mathbb{R}^n$ be a descent direction of f at \mathbf{x} . The exact line search for the stepsize can be obtained by considering

$$\min_{t \geq 0} \{g(t) := f(\mathbf{x} + t\mathbf{d})\}.$$

By a direct computation, we have

$$\begin{aligned} g(t) = f(\mathbf{x} + t\mathbf{d}) &= (\mathbf{d}^\top \mathbf{A} \mathbf{d})t^2 + 2(\mathbf{d}^\top \mathbf{A} \mathbf{x} + \mathbf{d}^\top \mathbf{b})t + \mathbf{x}^\top \mathbf{A} \mathbf{x} + 2\mathbf{b}^\top \mathbf{x} + c \\ &= (\mathbf{d}^\top \mathbf{A} \mathbf{d})t^2 + 2(\mathbf{d}^\top \mathbf{A} \mathbf{x} + \mathbf{d}^\top \mathbf{b})t + f(\mathbf{x}). \end{aligned}$$

Since $g'(t) = 2(\mathbf{d}^\top \mathbf{A} \mathbf{d})t + 2\mathbf{d}^\top (\mathbf{A} \mathbf{x} + \mathbf{b})$ and $\nabla f(\mathbf{x}) = 2(\mathbf{A} \mathbf{x} + \mathbf{b})$, it follows that $g'(t) = 0$ if and only if

$$t = t^* := -\frac{\mathbf{d}^\top \nabla f(\mathbf{x})}{2\mathbf{d}^\top \mathbf{A} \mathbf{d}} > 0,$$

where since \mathbf{d} is a descent direction of f at \mathbf{x} , $f'(\mathbf{x}; \mathbf{d}) = \mathbf{d}^\top \nabla f(\mathbf{x}) < 0$.

In what direction f decreases most rapidly?

- Making an observation, for $n = 2$, we have

$$f'(x_k; d_k) = \langle \nabla f(x_k), d_k \rangle = \|\nabla f(x_k)\| \|d_k\| \cos \theta_k,$$

where θ_k is the angle between the vectors $\nabla f(x_k)$ and d_k . Therefore, f decreases most rapidly when $\theta_k = \pi$, i.e., in the direction of $-\nabla f(x_k)$ whenever $\nabla f(x_k) \neq \mathbf{0}$.

- *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function and let $x \in \mathbb{R}^n$ be a nonstationary point, $\nabla f(x) \neq \mathbf{0}$. Then an optimal solution of $\min_{d \in \mathbb{R}^n} \{f'(x; d) : \|d\| = 1\}$ is $d = -\frac{\nabla f(x)}{\|\nabla f(x)\|}$.*

Proof: By the Cauchy-Schwarz inequality, for $\|d\| = 1$, we have

$$f'(x; d) = \nabla f(x)^\top d \geq -\|\nabla f(x)\| \|d\| = -\|\nabla f(x)\| \leftarrow \text{a lower bound}$$

Taking $d = -\frac{\nabla f(x)}{\|\nabla f(x)\|}$, we attain the lower bound. \square

The gradient method

- In the gradient method, we take $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$, provided $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$.

$$f'(\mathbf{x}_k; -\nabla f(\mathbf{x}_k)) = -\nabla f(\mathbf{x}_k)^\top \nabla f(\mathbf{x}_k) = -\|\nabla f(\mathbf{x}_k)\|^2 < 0.$$

- The gradient method

Input: Tolerance parameter $\varepsilon > 0$.

Initialization: Pick $\mathbf{x}_0 \in \mathbb{R}^n$ arbitrarily.

General step: For any $k = 0, 1, \dots$, execute

- (a) Pick a stepsize t_k by a line search procedure on the function

$$g(t) := f(\mathbf{x}_k - t\nabla f(\mathbf{x}_k)).$$

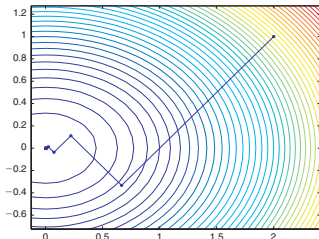
- (b) Set $\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)$.
- (c) If $\|\nabla f(\mathbf{x}_{k+1})\| \leq \varepsilon$ then stop and \mathbf{x}_{k+1} is the output.

Example

Consider the 2-D minimization problem $\min_{x,y} (x^2 + 2y^2)$ whose optimal solution is $(x,y) = (0,0)$ with corresponding optimal value 0.

- MATLAB function: `gradient_method_quadratic(...)`

For solving $\min_{x \in \mathbb{R}^n} \{x^\top A x + 2b^\top x\}$, $A \succ 0$, exact line search.



- MATLAB function: `gradient_method_constant(...)`
- MATLAB function: `gradient_method_backtracking(...)`

In computational experience, backtracking does not have real disadvantages in comparison to exact line search!

The gradient method: zig-zag effect

The zig-zag effect: *Let $\{x_k\}$ be the sequence generated by the gradient method with exact line search for solving a problem of minimizing a continuously differentiable function f . Then for any $k = 0, 1, 2, \dots$*

$$(x_{k+2} - x_{k+1})^\top (x_{k+1} - x_k) = 0.$$

Proof: By the definition of the gradient method, we have

$$x_{k+1} - x_k = -t_k \nabla f(x_k), \quad x_{k+2} - x_{k+1} = -t_{k+1} \nabla f(x_{k+1}).$$

Therefore, we wish to prove that $\nabla f(x_k)^\top \nabla f(x_{k+1}) = 0$. Since

$$g(t) := f(x_k - t \nabla f(x_k)),$$

we have

$$0 = g'(t_k) = -\nabla f(x_k)^\top \nabla f(x_k - t_k \nabla f(x_k)).$$

That is,

$$-\nabla f(x_k)^\top \nabla f(x_{k+1}) = 0,$$

which is the desired result. \square *(see the figure on page 9)*

A quadratic minimization problem

Consider the simple quadratic minimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x}) := \mathbf{x}^\top \mathbf{A} \mathbf{x}\},$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{A} \succ \mathbf{0}$. The optimal solution is obviously $\mathbf{x}^* = \mathbf{0}$.

The gradient method with exact line search takes the form

$$\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k, \quad \mathbf{d}_k = -\nabla f(\mathbf{x}_k) = -2\mathbf{A} \mathbf{x}_k, \quad t_k = \frac{\mathbf{d}_k^\top \mathbf{d}_k}{2\mathbf{d}_k^\top \mathbf{A} \mathbf{d}_k}.$$

Then we have

$$\begin{aligned} f(\mathbf{x}_{k+1}) &= \mathbf{x}_{k+1}^\top \mathbf{A} \mathbf{x}_{k+1} = (\mathbf{x}_k + t_k \mathbf{d}_k)^\top \mathbf{A} (\mathbf{x}_k + t_k \mathbf{d}_k) \\ &= \mathbf{x}_k^\top \mathbf{A} \mathbf{x}_k + 2t_k \mathbf{d}_k^\top \mathbf{A} \mathbf{x}_k + t_k^2 \mathbf{d}_k^\top \mathbf{A} \mathbf{d}_k \\ &= \mathbf{x}_k^\top \mathbf{A} \mathbf{x}_k - t_k \mathbf{d}_k^\top \mathbf{d}_k + t_k^2 \mathbf{d}_k^\top \mathbf{A} \mathbf{d}_k \\ &= \mathbf{x}_k^\top \mathbf{A} \mathbf{x}_k - \frac{1}{4} \frac{(\mathbf{d}_k^\top \mathbf{d}_k)^2}{\mathbf{d}_k^\top \mathbf{A} \mathbf{d}_k} \\ &= \mathbf{x}_k^\top \mathbf{A} \mathbf{x}_k \left(1 - \frac{1}{4} \frac{(\mathbf{d}_k^\top \mathbf{d}_k)^2}{(\mathbf{d}_k^\top \mathbf{A} \mathbf{d}_k)(\mathbf{x}_k^\top \mathbf{A} \mathbf{A}^{-1} \mathbf{A} \mathbf{x}_k)} \right). \end{aligned}$$

Kantorovich inequality

Since $d_k = -2Ax_k$, we have

$$f(x_{k+1}) = \left(1 - \frac{(d_k^\top d_k)^2}{(d_k^\top A d_k)(d_k^\top A^{-1} d_k)}\right) f(x_k).$$

Kantorovich inequality: Let $A \in \mathbb{R}^{n \times n}$ and $A \succ 0$. Then $\forall 0 \neq x \in \mathbb{R}^n$,

$$\frac{(x^\top x)^2}{(x^\top A x)(x^\top A^{-1} x)} \geq \frac{4\lambda_{\max}(A)\lambda_{\min}(A)}{(\lambda_{\max}(A) + \lambda_{\min}(A))^2}.$$

Proof: Let $m := \lambda_{\min}(A) > 0$ and $M := \lambda_{\max}(A) > 0$. Then the eigenvalues of $A + MmA^{-1}$ are $\lambda_i(A) + \frac{Mm}{\lambda_i(A)}$, $i = 1, 2, \dots, n$. The maximum value of the 1-D function $\varphi(t) = t + \frac{Mm}{t}$ on $[m, M]$ can be attained at $t = m$ and $t = M$ and the value is $M + m$. Therefore, the eigenvalues of $A + MmA^{-1}$ are smaller than $M + m$. Thus

$$A + MmA^{-1} \preceq (M + m)I,$$

which implies that

$$x^\top A x + Mm(x^\top A^{-1} x) \leq (M + m)(x^\top x).$$

Using the inequality $\alpha\beta \leq \frac{1}{4}(\alpha + \beta)^2$, we obtain the desired result

$$(x^\top A x)(Mm(x^\top A^{-1} x)) \leq \frac{1}{4}(x^\top A x + Mm(x^\top A^{-1} x))^2 \leq \frac{(M + m)^2}{4}(x^\top x)^2. \quad \square$$

Convergence rate analysis

Returning to the convergence rate analysis of the gradient method for the quadratic minimization problem, we have

$$\begin{aligned} f(\mathbf{x}_{k+1}) &= \left(1 - \frac{(\mathbf{d}_k^\top \mathbf{d}_k)^2}{(\mathbf{d}_k^\top \mathbf{A} \mathbf{d}_k)(\mathbf{d}_k^\top \mathbf{A}^{-1} \mathbf{d}_k)}\right) f(\mathbf{x}_k) \\ &\leq \left(1 - \frac{4Mm}{(M+m)^2}\right) f(\mathbf{x}_k) = \left(\frac{M-m}{M+m}\right)^2 f(\mathbf{x}_k), \end{aligned}$$

which implies a *linear rate* to the optimal value,

$$|f(\mathbf{x}_{k+1}) - 0| = f(\mathbf{x}_{k+1}) \leq c f(\mathbf{x}_k) = c |f(\mathbf{x}_k) - 0| \quad \text{and} \quad f(\mathbf{x}_k) \leq c^k f(\mathbf{x}_0),$$

$$c := \left(\frac{M-m}{M+m}\right)^2 = \left(\frac{\chi-1}{\chi+1}\right)^2 < 1, \quad \chi := \frac{M}{m} = \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})}.$$

Definition: $\chi(\mathbf{A}) := \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})}$ is called the *condition number* of \mathbf{A} .

Note: Although the condition number can be defined for general matrices, here we restrict ourselves to SPD real matrices.

Nonquadratic objective functions

- Matrices with large condition number are called *ill-conditioned*. Matrices with small condition number are called *well-conditioned*.
- The entire discussion until now was on the restrictive class of quadratic objective functions, where the Hessian matrix is constant, *but the notion of condition number also appears in the context of nonquadratic objective functions. In that case, it is well known that the rate of convergence of \mathbf{x}_k to a given stationary point \mathbf{x}^* depends on the condition number of $\chi(\nabla^2 f(\mathbf{x}^*))$.*
- We will not focus on these theoretical results, but will illustrate it on a well-known ill-conditioned problem, *the Rosenbrock function*, see next page.

The Rosenbrock function (control theory)

The Rosenbrock function is $f(x_1, x_2) := 100(x_2 - x_1^2)^2 + (1 - x_1)^2$.

- The optimal solution (global minimum) is $(x_1, x_2) = (1, 1)$ with corresponding optimal value 0.
- The gradient and Hessian of f are respectively

$$\begin{aligned}\nabla f(x_1, x_2) &= \begin{bmatrix} -400x_1(x_2 - x_1^2) - 2(1 - x_1) \\ 200(x_2 - x_1^2) \end{bmatrix}, \\ \nabla^2 f(x_1, x_2) &= \begin{bmatrix} -400x_2 + 1200x_1^2 + 2 & -400x_1 \\ -400x_1 & 200 \end{bmatrix}.\end{aligned}$$

- $(x_1, x_2) = (1, 1)$ is the unique stationary point and

$$\nabla^2 f(1, 1) = \begin{bmatrix} 802 & -400 \\ -400 & 200 \end{bmatrix}.$$

```
>> A = [802, -400; -400, 200];  
>> cond(A)  
ans = 2.5080e+003
```

A condition number of more than 2500 (ill-conditioned) should have severe effects on the convergence speed of the gradient method.

Sensitivity of solutions to linear systems

- We are given a linear system $Ax = b$, where $A \in \mathbb{R}^{n \times n}$ is symmetric, $A \succ 0$ and $b \in \mathbb{R}^n$. Then the solution is $x = A^{-1}b$.
- We consider a perturbation $b + \Delta b$ in the RHS. The new solution is denoted by $x + \Delta x$, i.e., $A(x + \Delta x) = b + \Delta b$. We have $x + \Delta x = A^{-1}(b + \Delta b) = x + A^{-1}\Delta b$. Then

$$\begin{aligned}\frac{\|\Delta x\|}{\|x\|} &= \frac{\|A^{-1}\Delta b\|}{\|x\|} \leq \frac{\|A^{-1}\| \|\Delta b\|}{\|x\|} = \frac{\lambda_{\max}(A^{-1}) \|\Delta b\|}{\|x\|} \\ &= \frac{1}{\lambda_{\min}(A)} \frac{\|\Delta b\|}{\|x\|} = \frac{1}{\lambda_{\min}(A)} \frac{\|\Delta b\|}{\|A^{-1}b\|} \leq \frac{1}{\lambda_{\min}(A)} \frac{\|\Delta b\|}{\lambda_{\min}(A^{-1}) \|b\|} \\ &= \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)} \frac{\|\Delta b\|}{\|b\|} = \chi(A) \frac{\|\Delta b\|}{\|b\|}, \quad \text{where we have used}\end{aligned}$$

$$\|A^{-1}b\| = \sqrt{b^\top A^{-2}b} \geq \sqrt{\lambda_{\min}(A^{-2}) \|b\|^2} = \lambda_{\min}(A^{-1}) \|b\|.$$

- *We can therefore deduce that the sensitivity of the solution of the linear system to right-hand-side perturbations depends on the condition number of the coefficients matrix.*

Scaling for ill-conditioned problems

We consider the unconstrained minimization problem:

$$\min\{f(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n\}.$$

Let $\mathbf{S} \in \mathbb{R}^{n \times n}$ be a nonsingular matrix. Let $\mathbf{y} := \mathbf{S}^{-1}\mathbf{x}$. Then $\mathbf{x} = \mathbf{S}\mathbf{y}$ and we obtain the equivalent problem:

$$\min\{g(\mathbf{y}) := f(\mathbf{S}\mathbf{y}) : \mathbf{y} \in \mathbb{R}^n\}.$$

Since $\nabla_{\mathbf{y}}g(\mathbf{y}) = \mathbf{S}^\top \nabla f(\mathbf{S}\mathbf{y}) = \mathbf{S}^\top \nabla f(\mathbf{x})$, the gradient method for solving $\min_{\mathbf{y} \in \mathbb{R}^n} g(\mathbf{y})$ takes the form:

$$\mathbf{y}_{k+1} = \mathbf{y}_k - t_k \mathbf{S}^\top \nabla f(\mathbf{S}\mathbf{y}_k).$$

Multiplying \mathbf{S} and letting $\mathbf{x}_k := \mathbf{S}\mathbf{y}_k$, we have

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \underbrace{\mathbf{S}\mathbf{S}^\top}_{:=\mathbf{D}} \nabla f(\mathbf{x}_k) := \mathbf{x}_k - t_k \mathbf{D} \nabla f(\mathbf{x}_k).$$

Then we obtain *the scaled gradient method* with scaling matrix \mathbf{D} .

The scaled gradient

- The matrix $D = SS^\top$ is positive definite (cf. Exercise 2.6). The direction $-D\nabla f(\mathbf{x}_k)$ is a descent of f at \mathbf{x}_k when $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$ since

$$f'(\mathbf{x}_k; -D\nabla f(\mathbf{x}_k)) = -\nabla f(\mathbf{x}_k)^\top D\nabla f(\mathbf{x}_k) < 0.$$

- To summarize the above discussion, we have shown that the scaled gradient method with scaling matrix $D \succ \mathbf{0}$ is equivalent to the gradient method employed on the function

$$g(\mathbf{y}) = f(D^{1/2}\mathbf{y}),$$

where $\mathbf{y} := D^{-1/2}\mathbf{x}$ ($\iff \mathbf{x} = D^{1/2}\mathbf{y}$). We note that the gradient and Hessian of g are given by

$$\begin{aligned}\nabla_{\mathbf{y}}g(\mathbf{y}) &= D^{1/2}\nabla f(D^{1/2}\mathbf{y}) = D^{1/2}\nabla_{\mathbf{x}}f(\mathbf{x}), \\ \nabla_{\mathbf{y}}^2g(\mathbf{y}) &= D^{1/2}\nabla^2f(D^{1/2}\mathbf{y})D^{1/2} = D^{1/2}\nabla_{\mathbf{x}}^2f(\mathbf{x})D^{1/2}.\end{aligned}$$

The scaled gradient method

Input: Tolerance parameter $\varepsilon > 0$.

Initialization: Pick $\mathbf{x}_0 \in \mathbb{R}^n$ arbitrarily.

General step: For any $k = 0, 1, \dots$, execute

(a) Pick a scaling matrix $\mathbf{D}_k \succ \mathbf{0}$.

(b) Pick a stepsize t_k by a line search procedure on the function

$$h(t) := f(\mathbf{x}_k - t\mathbf{D}_k \nabla f(\mathbf{x}_k)).$$

(c) Set $\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \mathbf{D}_k \nabla f(\mathbf{x}_k)$.

(d) If $\|\nabla f(\mathbf{x}_{k+1})\| \leq \varepsilon$ then stop and \mathbf{x}_{k+1} is the output.

It is often beneficial to choose the scaling matrix differently at each iteration.

How to choose the D_k ? damped Newton's method

- To accelerate the rate of convergence of $\{x_k\}$, which depends on the condition number of *the scaled Hessian* $D_k^{1/2} \nabla^2 f(x_k) D_k^{1/2}$. The scaling matrix is often chosen to make this scaled Hessian to be as close as possible to the identity matrix.
- When $\nabla^2 f(x_k) \succ 0$, we choose $D_k = (\nabla^2 f(x_k))^{-1}$ and the scaled Hessian becomes the identity matrix. *The resulting method is the so-called damped Newton's method:*

$$x_{k+1} = x_k - t_k (\nabla^2 f(x_k))^{-1} \nabla f(x_k)$$

One difficulty associated with damped Newton's method is that it requires full knowledge of the Hessian.

- The term $(\nabla^2 f(x_k))^{-1} \nabla f(x_k)$ suggests that a linear system of the form $\nabla^2 f(x_k) d = \nabla f(x_k)$ needs to be solved at each iteration, which might be costly from a computational point of view.
- The simplest of all scaling matrices are diagonal matrices. A natural choice for diagonal elements is $D_{ii} = (\nabla^2 f(x_k))_{ii}^{-1}$.

The Gauss-Newton method

We consider the nonlinear least squares (NLS) problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \{g(\mathbf{x}) := \sum_{i=1}^m (f_i(\mathbf{x}) - c_i)^2\},$$

where f_1, f_2, \dots, f_m are continuously differentiable over \mathbb{R}^n and $c_1, c_2, \dots, c_m \in \mathbb{R}$. The problem can be reformulated as

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|F(\mathbf{x})\|^2,$$

where the vector-valued function F is given by

$$F(\mathbf{x}) := \begin{bmatrix} f_1(\mathbf{x}) - c_1 \\ f_2(\mathbf{x}) - c_2 \\ \vdots \\ f_m(\mathbf{x}) - c_m \end{bmatrix}.$$

The Gauss-Newton method (A linearization method):

Given the iterate \mathbf{x}_k , find

$$\mathbf{x}_{k+1} := \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \sum_{i=1}^m \left(f_i(\mathbf{x}_k) + \nabla f_i(\mathbf{x}_k)^\top (\mathbf{x} - \mathbf{x}_k) - c_i \right)^2 \right\}.$$

The Gauss-Newton method (cont'd)

The minimization problem is essentially a linear LS problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{A}_k \mathbf{x} - \mathbf{b}_k\|^2,$$

where

$$\mathbf{A}_k = \begin{bmatrix} \nabla f_1(\mathbf{x}_k)^\top \\ \nabla f_2(\mathbf{x}_k)^\top \\ \vdots \\ \nabla f_m(\mathbf{x}_k)^\top \end{bmatrix} := J(\mathbf{x}_k),$$

is the so-called Jacobian matrix and

$$\mathbf{b}_k = \begin{bmatrix} \nabla f_1(\mathbf{x}_k)^\top \mathbf{x}_k - f_1(\mathbf{x}_k) + c_1 \\ \nabla f_2(\mathbf{x}_k)^\top \mathbf{x}_k - f_2(\mathbf{x}_k) + c_2 \\ \vdots \\ \nabla f_m(\mathbf{x}_k)^\top \mathbf{x}_k - f_m(\mathbf{x}_k) + c_m \end{bmatrix} := J(\mathbf{x}_k) \mathbf{x}_k - F(\mathbf{x}_k).$$

The underlying assumption is that $J(\mathbf{x}_k)$ is of a full column rank; otherwise the minimization will not produce a unique minimizer.

The Gauss-Newton method (cont'd)

We write an explicit expression for the Gauss-Newton iterates (see Chapter 3)

$$\mathbf{x}_{k+1} = (J(\mathbf{x}_k)^\top J(\mathbf{x}_k))^{-1} J(\mathbf{x}_k)^\top \mathbf{b}_k.$$

The method can also be written as

$$\begin{aligned}\mathbf{x}_{k+1} &= (J(\mathbf{x}_k)^\top J(\mathbf{x}_k))^{-1} J(\mathbf{x}_k)^\top \mathbf{b}_k \\ &= (J(\mathbf{x}_k)^\top J(\mathbf{x}_k))^{-1} J(\mathbf{x}_k)^\top (J(\mathbf{x}_k) \mathbf{x}_k - F(\mathbf{x}_k)) \\ &= \mathbf{x}_k - (J(\mathbf{x}_k)^\top J(\mathbf{x}_k))^{-1} J(\mathbf{x}_k)^\top F(\mathbf{x}_k).\end{aligned}$$

The Gauss-Newton direction is therefore

$$\mathbf{d}_k = -(J(\mathbf{x}_k)^\top J(\mathbf{x}_k))^{-1} J(\mathbf{x}_k)^\top F(\mathbf{x}_k).$$

Noting that $\nabla g(\mathbf{x}) = 2J(\mathbf{x})^\top F(\mathbf{x})$, we can conclude that

$$\mathbf{d}_k = -\frac{1}{2} (J(\mathbf{x}_k)^\top J(\mathbf{x}_k))^{-1} \nabla g(\mathbf{x}_k)$$

meaning that the Gauss-Newton method is essentially a scaled gradient method with $t_k = 1$ and the following positive definite scaling matrix:

$$\mathbf{D}_k = \frac{1}{2} (J(\mathbf{x}_k)^\top J(\mathbf{x}_k))^{-1}.$$

The damped Gauss-Newton method

The method described so far is also called the pure Gauss-Newton method since no stepsize is really involved. To transform this method into a practical algorithm, *a stepsize is introduced, leading to the damped Gauss-Newton method.*

The damped Gauss-Newton method

Input: Tolerance parameter $\varepsilon > 0$.

Initialization: Pick $\mathbf{x}_0 \in \mathbb{R}^n$ arbitrarily.

General step: For any $k = 0, 1, \dots$, execute

- (a) Set $\mathbf{d}_k = -(J(\mathbf{x}_k)^\top J(\mathbf{x}_k))^{-1} J(\mathbf{x}_k)^\top F(\mathbf{x}_k)$.
- (b) Set stepsize t_k by a line search procedure on the function

$$h(t) := g(\mathbf{x}_k + t\mathbf{d}_k).$$

- (c) Set $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$.
- (d) If $\|\nabla g(\mathbf{x}_{k+1})\| \leq \varepsilon$ then stop and \mathbf{x}_{k+1} is the output.

Lipschitz property of the gradient

We consider the following unconstrained minimization problem

$$\min\{f(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n\},$$

where the objective function f is *continuously differentiable*.

- **Definition:** ∇f is Lipschitz continuous over $\mathbb{R}^n \Leftrightarrow \exists L \geq 0$ such that $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$.
- $C_L^{1,1}(\mathbb{R}^n)$ or $C_L^{1,1}$ or $C^{1,1}(\mathbb{R}^n)$ or $C^{1,1}$: the class of functions over \mathbb{R}^n with Lipschitz gradient with constant L .
- $C_L^{1,1}(D)$: the set of all functions over $D \subseteq \mathbb{R}^n$ whose gradient satisfies the above Lipschitz condition for any $\mathbf{x}, \mathbf{y} \in D$.
- **Examples:**
 - (1) Linear functions: given $\mathbf{a} \in \mathbb{R}^n, f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x}$ is in $C_0^{1,1}$.
 - (2) Quadratic functions: let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be symmetric, $\mathbf{b} \in \mathbb{R}^n$ and $c \in \mathbb{R}$. Then $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} + 2\mathbf{b}^\top \mathbf{x} + c$ is in $C_L^{1,1}$, since

$$\begin{aligned}\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| &= 2\|(\mathbf{A}\mathbf{x} + \mathbf{b}) - (\mathbf{A}\mathbf{y} + \mathbf{b})\| \\ &\leq 2\|\mathbf{A}\|\|\mathbf{x} - \mathbf{y}\| := L\|\mathbf{x} - \mathbf{y}\|.\end{aligned}$$

The Fundamental Theorem of Calculus (FTC)

The FTC: Let $f : [a, b] \rightarrow \mathbb{R}$ be a real-valued function.

Part 1: Let $f \in \mathcal{R}[a, b]$. Define $F(x) := \int_a^x f(t)dt$, $x \in [a, b]$. Then (i) $F(x)$ is continuous on $[a, b]$; (ii) $F'(x) = f(x)$ for $x \in (a, b)$ where f is continuous.

Part 2: If $f' \in \mathcal{R}[a, b]$, then $\int_a^b f'(x)dx = f(b) - f(a)$.

Application: Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function over $D \subseteq \mathbb{R}^n$. Let $\mathbf{x}, \mathbf{y} \in D$ and $[\mathbf{x}, \mathbf{y}] \subseteq D$. Define $g(t) := f((1-t)\mathbf{x} + t\mathbf{y})$ for $t \in [0, 1]$. Using the chain rule and the FTC, we respectively obtain $g'(t) = \nabla f((1-t)\mathbf{x} + t\mathbf{y}) \cdot (\mathbf{y} - \mathbf{x})$ and

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) &= g(1) - g(0) = \int_0^1 g'(t)dt = \int_0^1 \nabla f((1-t)\mathbf{x} + t\mathbf{y}) \cdot (\mathbf{y} - \mathbf{x})dt \\ &= \int_0^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle dt. \end{aligned}$$

In addition, if f is twice continuously differentiable over $D \subseteq \mathbb{R}^n$, then

$$\begin{aligned} f_{\mathbf{x}}(\mathbf{y}) - f_{\mathbf{x}}(\mathbf{x}) &= \int_0^1 \nabla(f_{\mathbf{x}})(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) \cdot (\mathbf{y} - \mathbf{x})dt, \\ f_{\mathbf{y}}(\mathbf{y}) - f_{\mathbf{y}}(\mathbf{x}) &= \int_0^1 \nabla(f_{\mathbf{y}})(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) \cdot (\mathbf{y} - \mathbf{x})dt. \end{aligned}$$

That is, we have

$$\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) = \int_0^1 \nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x})dt.$$

Boundedness of the Hessian

Theorem: *Let f be a twice continuously differentiable function over \mathbb{R}^n . Then*

$$f \in C_L^{1,1}(\mathbb{R}^n) \iff \|\nabla^2 f(x)\| \leq L, \forall x \in \mathbb{R}^n.$$

Proof: (\Leftarrow) By the fundamental theorem of calculus, $\forall x, y \in \mathbb{R}^n$, we have

$$\nabla f(y) - \nabla f(x) = \int_0^1 \nabla^2 f(x + t(y-x))(y-x) dt = \int_0^1 \nabla^2 f(x + t(y-x)) dt (y-x).$$

Thus, we have

$$\begin{aligned} \|\nabla f(y) - \nabla f(x)\| &\leq \left\| \int_0^1 \nabla^2 f(x + t(y-x)) dt \right\| \|y-x\| \\ &\leq \left(\int_0^1 \|\nabla^2 f(x + t(y-x))\| dt \right) \|y-x\| \leq L \|y-x\|. \end{aligned}$$

(\Rightarrow) By the fundamental theorem of calculus, $\forall d \in \mathbb{R}^n$ and $\alpha > 0$, we have

$$\nabla f(x + \alpha d) - \nabla f(x) = \int_0^\alpha \nabla^2 f(x + td) d dt.$$

Thus, we have

$$\left\| \left(\int_0^\alpha \nabla^2 f(x + td) dt \right) d \right\| = \|\nabla f(x + \alpha d) - \nabla f(x)\| \leq \alpha L \|d\|.$$

Dividing by α and taking the limit $\alpha \rightarrow 0^+$, we obtain $\|\nabla^2 f(x)d\| \leq L\|d\|$, *where we have used the mean value theorem for definite integrals for each matrix component of $\nabla^2 f(x + td)$.* \square

The descent lemma

The following descent lemma is fundamental in convergence proofs of gradient-based methods.

The descent lemma: *Let $D \subseteq \mathbb{R}^n$ and $f \in C_L^{1,1}(D)$ for some $L > 0$. Then for any $\mathbf{x}, \mathbf{y} \in D$ satisfying $[\mathbf{x}, \mathbf{y}] \subseteq D$ it holds that*

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

Proof: By the fundamental theorem of calculus, we have

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) &= \int_0^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle dt \\ &= \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_0^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle dt. \end{aligned}$$

Therefore, we have

$$\begin{aligned} |f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| &= \left| \int_0^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle dt \right| \\ &\leq \int_0^1 |\langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| dt \\ &\leq \int_0^1 \|\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\| \|\mathbf{y} - \mathbf{x}\| dt \\ &\leq \int_0^1 tL \|\mathbf{y} - \mathbf{x}\|^2 dt = \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2. \quad \square \end{aligned}$$

A sufficient decrease lemma

- Note that the proof of the descent lemma actually shows both upper and lower bounds on the function:

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

- A sufficient decrease lemma:** *Suppose that $f \in C_L^{1,1}(\mathbb{R}^n)$. Then $\forall \mathbf{x} \in \mathbb{R}^n$ and $t > 0$, we have*

$$f(\mathbf{x}) - f(\mathbf{x} - t\nabla f(\mathbf{x})) \geq t \left(1 - \frac{Lt}{2}\right) \|\nabla f(\mathbf{x})\|^2.$$

Proof: By the descent lemma we have

$$\begin{aligned} f(\mathbf{x} - t\nabla f(\mathbf{x})) &\leq f(\mathbf{x}) - t\|\nabla f(\mathbf{x})\|^2 + \frac{Lt^2}{2} \|\nabla f(\mathbf{x})\|^2 \\ &= f(\mathbf{x}) - t \left(1 - \frac{Lt}{2}\right) \|\nabla f(\mathbf{x})\|^2. \end{aligned}$$

The result then follows by simple rearrangement of terms. \square

Sufficient decrease of the gradient method

Theorem: Let $f \in C_L^{1,1}(\mathbb{R}^n)$. Let $\{\mathbf{x}_k\}_{k \geq 0}$ be generated by the gradient method for solving $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ with one of the following stepsize strategies:

- (1) constant stepsize $\bar{t} \in (0, \frac{2}{L})$,
- (2) exact line search,
- (3) backtracking procedure with parameters $s > 0, \alpha, \beta \in (0, 1)$.

Then we have

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq M \|\nabla f(\mathbf{x}_k)\|^2,$$

where

$$M := \begin{cases} \bar{t}(1 - \frac{\bar{t}L}{2}) & \text{constant stepsize,} \\ \frac{1}{2L} & \text{exact line search,} \\ \alpha \min\{s, \frac{2(1-\alpha)\beta}{L}\} & \text{backtracking.} \end{cases}$$

Proof: constant stepsize and exact line search

- (1) **(constant stepsize)** By the sufficient decrease lemma, we immediately have

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq \bar{t} \left(1 - \frac{L\bar{t}}{2}\right) \|\nabla f(\mathbf{x}_k)\|^2 \geq 0 \quad \text{for } \bar{t} \in (0, \frac{2}{L}). \quad \square$$

Furthermore, if we wish to obtain the largest guaranteed bound on the decrease, then we seek the maximum of

$$\bar{t} \left(1 - \frac{L\bar{t}}{2}\right), \quad \forall \bar{t} \in (0, \frac{2}{L}).$$

One can show that this maximum is attained at $\bar{t} = \frac{1}{L}$ and we have

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k) - f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right) \geq \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|^2. \quad (\star)$$

- (2) **(exact line search)** In the exact line search setting, $t_k \in \operatorname{argmin}_{t \geq 0} f(\mathbf{x}_k - t \nabla f(\mathbf{x}_k))$. By the definition of t_k we know that

$$f(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)) \leq f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right).$$

Therefore, we have

$$\begin{aligned} f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) &= f(\mathbf{x}_k) - f(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)) \\ &\geq f(\mathbf{x}_k) - f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right) \geq \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|^2, \end{aligned}$$

where the last inequality comes from (\star) . \square

Proof: backtracking

- (3) **(backtracking)** In the backtracking setting we seek a small enough stepsize t_k for which we have

$$f(\mathbf{x}_k) - f(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)) \geq \alpha t_k \|\nabla f(\mathbf{x}_k)\|^2, \quad \alpha \in (0, 1).$$

We would like to find a lower bound on t_k . There are two options. Either $t_k = s$ (the initial value of the stepsize) or the stepsize t_k/β is not acceptable, i.e.,

$$f(\mathbf{x}_k) - f(\mathbf{x}_k - \frac{t_k}{\beta} \nabla f(\mathbf{x}_k)) < \alpha \frac{t_k}{\beta} \|\nabla f(\mathbf{x}_k)\|^2. \quad (\star 1)$$

By the sufficient decrease lemma with $\mathbf{x} = \mathbf{x}_k$ and $t = \frac{t_k}{\beta}$, we have

$$f(\mathbf{x}_k) - f(\mathbf{x}_k - \frac{t_k}{\beta} \nabla f(\mathbf{x}_k)) \geq \frac{t_k}{\beta} \left(1 - \frac{Lt_k}{2\beta}\right) \|\nabla f(\mathbf{x}_k)\|^2. \quad (\star 2)$$

From $(\star 1)$ and $(\star 2)$, we obtain

$$\frac{t_k}{\beta} \left(1 - \frac{Lt_k}{2\beta}\right) < \alpha \frac{t_k}{\beta} \iff t_k > \frac{2(1-\alpha)\beta}{L}.$$

Overall, we have

$$t_k \geq \min \left\{ s, \frac{2(1-\alpha)\beta}{L} \right\}.$$

Finally, we obtain

$$f(\mathbf{x}_k) - f(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)) \geq \alpha \min \left\{ s, \frac{2(1-\alpha)\beta}{L} \right\} \|\nabla f(\mathbf{x}_k)\|^2. \quad \square$$

Convergence of the gradient method

Theorem: Let $f \in C_L^{1,1}(\mathbb{R}^n)$. Let $\{\mathbf{x}_k\}_{k \geq 0}$ be generated by the gradient method for solving $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ with one of the following stepsize strategies:

- (1) constant stepsize $\bar{t} \in (0, \frac{2}{L})$,
- (2) exact line search,
- (3) backtracking procedure with parameters $s > 0, \alpha, \beta \in (0, 1)$.

Assume that f is bounded below over \mathbb{R}^n , i.e., $\exists m \in \mathbb{R}$ such that $f(\mathbf{x}) > m$ for all $\mathbf{x} \in \mathbb{R}^n$. Then we have the following:

- (a) The sequence $\{f(\mathbf{x}_k)\}_{k \geq 0}$ is nonincreasing. In addition, for any $k \geq 0$, $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ unless $\nabla f(\mathbf{x}_k) = \mathbf{0}$.
- (b) $\nabla f(\mathbf{x}_k) \rightarrow \mathbf{0}$ as $k \rightarrow \infty$.

Proof of the convergence theorem

(a) By the sufficient decrease of the gradient method, we have

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq M \|\nabla f(\mathbf{x}_k)\|^2, \quad (\star\star)$$

for some constant $M > 0$, and hence the equality $f(\mathbf{x}_k) = f(\mathbf{x}_{k+1})$ can hold only when $\nabla f(\mathbf{x}_k) = \mathbf{0}$.

(b) Since the sequence $\{f(\mathbf{x}_k)\}_{k \geq 0}$ is nonincreasing, and bounded below, it converges. Thus, in particular

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \rightarrow 0 \quad \text{as } k \rightarrow \infty,$$

which combined with $(\star\star)$ implies $\|\nabla f(\mathbf{x}_k)\| \rightarrow 0$ as $k \rightarrow \infty$, according to the squeeze theorem. Therefore, we obtain

$$\nabla f(\mathbf{x}_k) \rightarrow \mathbf{0} \quad \text{as } k \rightarrow \infty. \quad \square$$

“Rate of convergence” of gradient norms

Theorem: Under the setting of the previous theorem, let f^* be the limit of the convergent sequence $\{f(\mathbf{x}_k)\}_{k \geq 0}$. Then for any $\ell = 0, 1, 2, \dots$

$$\min_{k=0,1,\dots,\ell} \|\nabla f(\mathbf{x}_k)\| \leq \sqrt{\frac{f(\mathbf{x}_0) - f^*}{M(\ell + 1)}},$$

where

$$M = \begin{cases} \bar{t}(1 - \frac{\bar{t}L}{2}) & \text{constant stepsize,} \\ \frac{1}{2L} & \text{exact line search,} \\ \alpha \min\{s, \frac{2(1-\alpha)\beta}{L}\} & \text{backtracking.} \end{cases}$$

Proof: Summing the inequality (**) on the previous page over $k = 0, 1, \dots, \ell$, we obtain the following inequality

$$f(\mathbf{x}_0) - f(\mathbf{x}_{\ell+1}) \geq M \sum_{k=0}^{\ell} \|\nabla f(\mathbf{x}_k)\|^2.$$

Since $f(\mathbf{x}_{\ell+1}) \geq f^*$, we can conclude that

$$f(\mathbf{x}_0) - f^* \geq M \sum_{k=0}^{\ell} \|\nabla f(\mathbf{x}_k)\|^2 \geq M(\ell + 1) \min_{k=0,1,\dots,\ell} \|\nabla f(\mathbf{x}_k)\|^2,$$

implying the desired result. \square