

Chapter 8 Gradient Methods

1. Steepest Descent
2. Analysis of Gradient Methods



Optimization Problem

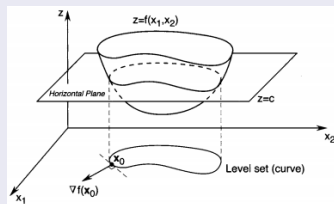
unconstrained optimization problem: $\min\{f(\mathbf{x}) \mid \mathbf{x} \in \mathbb{R}^n\}$

- The function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is multi-dimensional objective function.
- The vector $\mathbf{x} \in \mathbb{R}^n$ is n -dimensional decision variables.

Definition (level set, normal vector, and tangent vector)

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be function, and $c \in \mathbb{R}$ be a constant.

- level set of f w.r.t constant c :
 $S_c(f) = \{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) = c\}$.
 $\implies \mathbf{x}_0 \in S_c(f)$ if $f(\mathbf{x}_0) = c$.
- tangent of f at \mathbf{x}_0 :
 $l(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^\top (\mathbf{x} - \mathbf{x}_0)$.
- normal of f at \mathbf{x}_0 : $\nabla f(\mathbf{x}_0)$.



★ Given small displacement, f increases more in the direction of $\nabla f(\mathbf{x}_0)$ than in any other direction.

proof. Taylor series and the rate of increase $\langle \nabla f(\mathbf{x}_0), \mathbf{d} \rangle$ with $\|\mathbf{d}\| = 1$ and Cauchy inequality



Gradient Descent Methods

motivation for $\min\{f(\mathbf{x}) \mid \mathbf{x} \in \mathbb{R}^n\}$

Let $\mathbf{x}^{(0)}$ be an initial point, and a new point $\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})$.

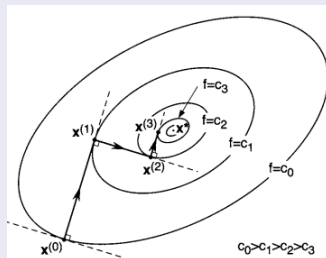
Taylor's series $\Rightarrow f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})) = f(\mathbf{x}^{(0)}) - \alpha \|\nabla f(\mathbf{x}^{(0)})\|^2 + o(\alpha)$.

If $\nabla f(\mathbf{x}^{(0)}) \neq 0 \Rightarrow f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})) < f(\mathbf{x}^{(0)})$, $\forall \alpha > 0$.

$\Rightarrow \mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})$ is an improvement over $\mathbf{x}^{(0)}$ as a minimizer.

gradient descent method for solving $\min\{f(\mathbf{x}) \mid \mathbf{x} \in \mathbb{R}^n\}$

- gradient descent method:
 - S0: given $\mathbf{x}^{(0)}$ and step size $\alpha_k > 0$,
 - S1: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$.
- steepest descent method:
 - S0: given $\mathbf{x}^{(0)}$,
 - S1: $\alpha_k = \arg \min_{\alpha > 0} f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$.
 - S2: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$.



Gradient Descent Methods

Theorem (property of steepest descent method)

If $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$ is a steepest descent sequence, then $\nabla f(\mathbf{x}^{(k+1)}) \perp \nabla f(\mathbf{x}^{(k)})$ or $(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) \perp (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$.

proof. let $\varphi(\alpha) = f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$.

By definition of $\alpha_k \implies \varphi'(\alpha_k) = 0 \implies \langle \nabla f(\mathbf{x}^{(k+1)}), \nabla f(\mathbf{x}^{(k)}) \rangle = 0$.

By iterative scheme $\mathbf{x}^{(k)} \implies \langle \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}, \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \rangle = 0$

Theorem (property of steepest descent method)

If $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$ is a steepest descent sequence, and if $\nabla f(\mathbf{x}^{(k)}) \neq 0$, then $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$.

proof. By the definition of φ and Taylor series

$$\iff \varphi(\alpha_k) = \varphi(\alpha_0) + \varphi'(\alpha_0)\alpha_k + o(\alpha_k)$$

$$\iff f(\mathbf{x}^{(k+1)}) = f(\mathbf{x}^{(k)}) - \alpha_k \|\nabla f(\mathbf{x}^{(k)})\|^2 + o(\alpha_k).$$

If $\nabla f(\mathbf{x}^{(k)}) \neq 0$, then $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$.



Gradient Descent Methods

- ★ steepest descent method render the searching direction in “zigzag” phenomenon.
- ★ the gradient descent method render the objective function value monotonically decrease.

stopping criterion: let $\varepsilon > 0$ is a prespecified threshold (e.g., 10^{-6})

- 1st-order necessary condition: $\|\nabla f(\mathbf{x}^{(k)})\| < \varepsilon$.
- difference of successive objective function values:
 $|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})| < \varepsilon, \frac{|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})|}{|f(\mathbf{x}^{(k)})|} < \varepsilon, \frac{|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})|}{\max\{1, |f(\mathbf{x}^{(k)})|\}} < \varepsilon$.
- difference of successive iterates:
 $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < \varepsilon, \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|}{\|\mathbf{x}^{(k)}\|} < \varepsilon, \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|}{\max\{1, \|\mathbf{x}^{(k)}\|\}} < \varepsilon$.



Gradient Descent Methods

Example $(\min_{\mathbf{x} \in \mathbb{R}^3} f(\mathbf{x}) = (x_1 - 4)^4 + (x_2 - 3)^2 + 4(x_3 + 5)^4)$

Minimize f by steepest descent method with $\mathbf{x}^{(0)} = [4, 2, -1]^\top$.

Ans: The gradient of f is $\nabla f(\mathbf{x}) = [4(x_1 - 4)^3, 2(x_2 - 3), 16(x_3 + 5)^3]^\top$.

S1: $\because \nabla f(\mathbf{x}^{(0)}) = [0, -2, 1024]^\top$ and exact stepsize

$$\alpha_0 = \arg \min_{\alpha > 0} f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})) \xrightarrow{\text{secant method}} \alpha_0 = 3.967 \times 10^{-3}.$$

The new iterate $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \alpha_0 \nabla f(\mathbf{x}^{(0)}) = [4.000, 2.008, -5.062]^\top$.

S2: $\because \nabla f(\mathbf{x}^{(1)}) = [0.000, -1.984, -0.003875]^\top$ and exact stepsize

$$\alpha_1 = \arg \min_{\alpha > 0} f(\mathbf{x}^{(1)} - \alpha \nabla f(\mathbf{x}^{(1)})) \xrightarrow{\text{secant method}} \alpha_1 = 0.500.$$

The new iterate $\mathbf{x}^{(2)} = \mathbf{x}^{(1)} - \alpha_1 \nabla f(\mathbf{x}^{(1)}) = [4.000, 3.000, -5.060]^\top$.

S3: $\because \nabla f(\mathbf{x}^{(2)}) = [0.000, 0.000, -0.003525]^\top$ and exact stepsize

$$\alpha_2 = \arg \min_{\alpha > 0} f(\mathbf{x}^{(2)} - \alpha \nabla f(\mathbf{x}^{(2)})) \xrightarrow{\text{secant method}} \alpha_2 = 16.29.$$

The new iterate $\mathbf{x}^{(3)} = \mathbf{x}^{(2)} - \alpha_2 \nabla f(\mathbf{x}^{(2)}) = [4.000, 3.000, -5.002]^\top$.

★ The minimizer of $\mathbf{x}^* = [4, 3, -5]^\top$. It can be implemented by C/C++/matlab/... programs.



Steepest Descent Methods on Quadratic Function

Definition (quadratic function)

$f(x) = \frac{1}{2}x^\top Qx - b^\top x$, where $Q \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix, $b \in \mathbb{R}^n$, and $x \in \mathbb{R}^n$.

By 1st-order necessary condition: $0 = \nabla f(x) = Qx - b, \implies$ unique minimizer: $x^* = Q^{-1}b$.

Notations: let gradient of f at $x^{(k)}$ be $g^{(k)} = \nabla f(x^{(k)})$ and Hessian of f be $F(x) = \nabla^2 f(x^{(k)}) = Q$.

steepest descent method for quadratic function

$$\alpha_k = \frac{(g^{(k)})^\top g^{(k)}}{(g^{(k)})^\top Q g^{(k)}} = \frac{\|g^{(k)}\|_2^2}{\|g^{(k)}\|_Q^2},$$
$$x^{(k+1)} = x^{(k)} - \alpha_k g^{(k)}.$$

Question: why $\alpha_k = \frac{(g^{(k)})^\top g^{(k)}}{(g^{(k)})^\top Q g^{(k)}}$?



Steepest Descent Methods on Quadratic Function

$$\text{exact stepsize: } \alpha_k = \frac{(\mathbf{g}^{(k)})^\top \mathbf{g}^{(k)}}{(\mathbf{g}^{(k)})^\top \mathbf{Q} \mathbf{g}^{(k)}}$$

$$\because \alpha_k = \arg \min_{\alpha > 0} \varphi(\alpha) := f(\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}),$$

By optimality condition of minimizer α_k , $\therefore \varphi'(\alpha) = 0$, i.e.,

$$\begin{aligned} 0 &= \varphi'(\alpha_k) = (\mathbf{g}^{(k)})^\top \nabla f(\mathbf{x}^{(k)} - \alpha \nabla \mathbf{g}^{(k)}) = (\mathbf{g}^{(k)})^\top (\mathbf{Q} \mathbf{x}^{(k)} - \alpha \mathbf{Q} \mathbf{g}^{(k)} - \mathbf{b}) \\ &= (\mathbf{g}^{(k)})^\top (\mathbf{g}^{(k)} - \alpha \mathbf{Q} \mathbf{g}^{(k)}) = (\mathbf{g}^{(k)})^\top \mathbf{g}^{(k)} - \alpha (\mathbf{g}^{(k)})^\top \mathbf{Q} \mathbf{g}^{(k)} \end{aligned}$$

$$\text{Example } (\min_{\mathbf{x} \in \mathbb{R}^2} f(\mathbf{x}) = x_1^2 + x_2^2)$$

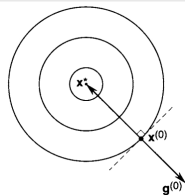
Minimize f by steepest descent method with any initial point $\mathbf{x}^{(0)} \in \mathbb{R}^2$.

Ans: \because gradient of f is $\nabla f(\mathbf{x}) = 2\mathbf{x}$,

Hessian of f is $\nabla^2 f(\mathbf{x}) = \mathbf{Q} = 2\mathbf{I}$.

$$\because \mathbf{g}^{(0)} = 2\mathbf{x}^{(0)} \implies \alpha_0 = \frac{(\mathbf{g}^{(0)})^\top \mathbf{g}^{(0)}}{(\mathbf{g}^{(0)})^\top \mathbf{Q} \mathbf{g}^{(0)}} = \frac{1}{2}$$

$$\therefore \text{new iterate: } \mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \alpha_0 \mathbf{g}^{(0)} = [0, 0]^\top.$$



Steepest Descent Methods on Quadratic Function

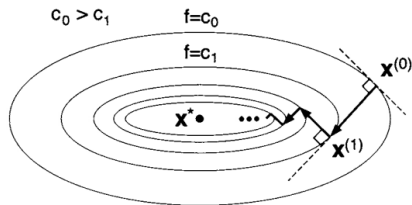
Example $(\min_{\mathbf{x} \in \mathbb{R}^2} f(\mathbf{x}) = \frac{x_1^2}{5} + x_2^2)$

Minimize f by steepest descent method with any initial point $\mathbf{x}^{(0)} \in \mathbb{R}^2$.

Ans: with any initial point $\mathbf{x}^{(0)}$,
implement steepest descent method

$$\begin{cases} \alpha_k = \frac{(\mathbf{g}^{(k)})^\top \mathbf{g}^{(k)}}{(\mathbf{g}^{(k)})^\top \mathbf{Q} \mathbf{g}^{(k)}} = \frac{\|\mathbf{g}^{(k)}\|_2^2}{\|\mathbf{g}^{(k)}\|_Q^2}, \\ \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}, \end{cases}$$

(e.g., in Matlab/C/C++)



Properties of steepest descent methods on quadratic function

- i) function values are monotonically descent, i.e., $f(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)})$;
- ii) globally convergent, $\mathbf{x}^{(k)} \xrightarrow{k \rightarrow \infty} \mathbf{x}^*$ a minimizer of f .

Convergence Analysis

- **Goal:** convergence analysis of steepest descent method on $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$, where $\mathbf{Q} \succ 0$, $\mathbf{b} \in \mathbb{R}^n$, and $\mathbf{x} \in \mathbb{R}^n$.
- **Facts:** \mathbf{x}^* is minimizer (i.e., $\mathbf{Q} \mathbf{x}^* = \mathbf{b}$) $\Rightarrow f(\mathbf{x}^*) = -\frac{1}{2} (\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^*$.
- **Notations:** $V(\mathbf{x}) = f(\mathbf{x}) - f(\mathbf{x}^*) = \frac{1}{2} \|\mathbf{x} - \mathbf{x}^*\|_{\mathbf{Q}}^2 = \frac{1}{2} \|\mathbf{g}\|_{\mathbf{Q}^{-1}}^2$

Lemma (steepest descent method $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}$)

Let $\{\mathbf{x}^{(k)}\}_{k=0}^\infty$ be a sequence generated by steepest descent method. Then

$$V(\mathbf{x}^{(k+1)}) = (1 - \gamma_k) V(\mathbf{x}^{(k)}),$$

where $\mathbf{g}^{(k)} = \mathbf{Q} \mathbf{x}^{(k)} - \mathbf{b}$ and the parameter

$$\gamma_k = \begin{cases} 1, & \text{if } \mathbf{g}^{(k)} = 0, \\ \frac{\alpha_k \|\mathbf{g}^{(k)}\|_{\mathbf{Q}}^2}{\|\mathbf{g}^{(k)}\|_{\mathbf{Q}^{-1}}^2} \left(2 \frac{\|\mathbf{g}^{(k)}\|_{\mathbf{Q}}^2}{\|\mathbf{g}^{(k)}\|_{\mathbf{Q}}^2} - \alpha_k \right), & \text{if } \mathbf{g}^{(k)} \neq 0. \end{cases}$$

Convergence Analysis

proof. assume $\mathbf{g}^{(k)} \neq 0$. Therefore,

$$\begin{aligned} & V(\mathbf{x}^{(k+1)}) - V(\mathbf{x}^{(k)}) \\ &= \frac{1}{2} \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_Q^2 - \frac{1}{2} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_Q^2 \\ &= \frac{1}{2} \|\mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)} - \mathbf{x}^*\|_Q^2 - \frac{1}{2} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_Q^2 \quad [\because \text{iterative scheme}] \\ &= -\alpha_k (\mathbf{x}^{(k)} - \mathbf{x}^*)^\top Q \mathbf{g}^{(k)} + \frac{\alpha_k^2}{2} \|\mathbf{g}^{(k)}\|_Q^2 \\ &= -\alpha_k (\mathbf{g}^{(k)})^\top \mathbf{g}^{(k)} + \frac{\alpha_k^2}{2} \|\mathbf{g}^{(k)}\|_Q^2 \\ &= \underbrace{\frac{\|\mathbf{g}^{(k)}\|_{Q^{-1}}^2}{2}}_{V(\mathbf{x}^{(k)})} \underbrace{\frac{2}{\|\mathbf{g}^{(k)}\|_{Q^{-1}}^2} \left(-\alpha_k \|\mathbf{g}^{(k)}\|_2^2 + \frac{\alpha_k^2}{2} \|\mathbf{g}^{(k)}\|_Q^2 \right)}_{\gamma_k} \end{aligned}$$

$$\therefore V(\mathbf{x}^{(k+1)}) - V(\mathbf{x}^{(k)}) = \gamma_k V(\mathbf{x}^{(k)}) \implies V(\mathbf{x}^{(k+1)}) = (1 - \gamma_k) V(\mathbf{x}^{(k)}).$$



Convergence Analysis

Question: how to guarantee $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$?

Lemma

Let $\{\mathbf{x}^k\}_{k=0}^{\infty}$ be generated by steepest descent method. Then, $\{\mathbf{x}^k\}_{k=0}^{\infty}$ converges to \mathbf{x}^* for any initial point $\mathbf{x}^{(0)}$ if and only if $\sum_{k=0}^{\infty} \gamma_k = \infty$.

proof. $\because V(\mathbf{x}^{(k+1)}) = (1 - \gamma_k)V(\mathbf{x}^{(k)}) \implies V(\mathbf{x}^{(k+1)}) = \left(\prod_{i=0}^k (1 - \gamma_i) \right) V(\mathbf{x}^{(0)})$

$$\begin{aligned} \therefore \boxed{\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*} &\iff \boxed{V(\mathbf{x}^{(k+1)}) \rightarrow 0} \iff \boxed{\prod_{i=0}^k (1 - \gamma_i) \rightarrow 0} \\ &\iff \boxed{\sum_{i=0}^{\infty} \gamma_i = \infty} \iff \boxed{\sum_{i=0}^{\infty} -\log(1 - \gamma_i) = \infty} \end{aligned}$$

★ The last “ \iff ” is due to $-2x \leq \log(1 - x) \leq -x$ in calculus.



Convergence Analysis

★ **Question:** how to guarantee $\sum_{k=0}^{\infty} \gamma_k = \infty$?

Lemma (Rayleigh's inequality)

Let $\lambda_{\max}(\mathbf{Q})$ and $\lambda_{\min}(\mathbf{Q})$ be the maximal and minimal eigenvalues of $\mathbf{Q} \succ 0$.

$$\text{Then, } \begin{cases} \lambda_{\min}(\mathbf{Q}) \leq \frac{\|\mathbf{x}\|_{\mathbf{Q}}^2}{\|\mathbf{x}\|_2^2} \leq \lambda_{\max}(\mathbf{Q}), \\ \frac{1}{\lambda_{\max}(\mathbf{Q})} \leq \frac{\|\mathbf{x}\|_{\mathbf{Q}^{-1}}^2}{\|\mathbf{x}\|_2^2} \leq \frac{1}{\lambda_{\min}(\mathbf{Q})}, \end{cases} \quad \forall \mathbf{x} \neq 0.$$

(*proof: quadratic form in linear algebra.*)

Lemma

$$\text{For any } \mathbf{Q} \succ 0 \text{ and } \mathbf{x} \neq 0, \quad \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{Q})} \leq \frac{\|\mathbf{x}\|_{\mathbf{Q}}^2 \|\mathbf{x}\|_{\mathbf{Q}^{-1}}^2}{\|\mathbf{x}\|_2^4} \leq \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})}.$$

★ (stringent version) By Cauchy and Kantorovich inequalities:

$$\text{For any } \mathbf{Q} \succ 0 \text{ and } \mathbf{x} \neq 0, \quad 1 \leq \frac{\|\mathbf{x}\|_{\mathbf{Q}}^2 \|\mathbf{x}\|_{\mathbf{Q}^{-1}}^2}{\|\mathbf{x}\|_2^4} \leq \frac{[\lambda_{\max}(\mathbf{Q}) + \lambda_{\min}(\mathbf{Q})]^2}{4\lambda_{\max}(\mathbf{Q})\lambda_{\min}(\mathbf{Q})}.$$



Convergence Analysis

Theorem

The sequence $\{\mathbf{x}^{(k)}\}$ generated by steepest descent method on minimizing quadratic function converges a minimizer \mathbf{x}^ for any $\mathbf{x}^{(0)}$.*

proof. assume $\mathbf{g}^{(k)} \neq 0$.

$$\therefore \gamma_k = \frac{\alpha_k \|\mathbf{g}^{(k)}\|_Q^2}{\|\mathbf{g}^{(k)}\|_{Q^{-1}}^2} \left(2 \frac{\|\mathbf{g}^{(k)}\|_2^2}{\|\mathbf{g}^{(k)}\|_Q^2} - \alpha_k \right)$$

$$\xrightarrow{\alpha_k = \frac{\|\mathbf{g}^{(k)}\|_2^2}{\|\mathbf{g}^{(k)}\|_Q^2}} \gamma_k = \frac{\|\mathbf{g}^{(k)}\|_2^4}{\|\mathbf{g}^{(k)}\|_Q^2 \|\mathbf{g}^{(k)}\|_{Q^{-1}}^2}.$$

$$\therefore \frac{4\lambda_{\max}(\mathbf{Q})\lambda_{\min}(\mathbf{Q})}{[\lambda_{\max}(\mathbf{Q}) + \lambda_{\min}(\mathbf{Q})]^2} \leq \gamma_k \leq 1$$

$$\iff \sum_{i=0}^{\infty} \gamma_k = \infty$$

$$\iff V(\mathbf{x}^{(k)}) \rightarrow 0$$

$$\iff \mathbf{x}^{(k)} \rightarrow \mathbf{x}^*.$$



Convergence Analysis

Notation: GD=gradient descent method.

Corollary (GD with fixed stepsize: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}$)

The sequence $\{\mathbf{x}^{(k)}\}$ produced by GD with fixed stepsize on minimizing quadratic function converges to \mathbf{x}^ for any $\mathbf{x}^{(0)}$ and $\alpha \in (0, \frac{2}{\lambda_{\max}(\mathbf{Q})})$.*

proof. $\gamma_k = \frac{\alpha_k \|\mathbf{g}^{(k)}\|_Q^2}{\|\mathbf{g}^{(k)}\|_{Q^{-1}}^2} \left(2 \frac{\|\mathbf{g}^{(k)}\|_2^2}{\|\mathbf{g}^{(k)}\|_Q^2} - \alpha_k \right) \xrightarrow{\alpha_k \equiv \alpha} \gamma_k = \alpha \frac{\|\mathbf{g}^{(k)}\|_Q^2}{\|\mathbf{g}^{(k)}\|_{Q^{-1}}^2} \left(2 \frac{\|\mathbf{g}^{(k)}\|_2^2}{\|\mathbf{g}^{(k)}\|_Q^2} - \alpha \right).$

$$\because \frac{\|\mathbf{g}^{(k)}\|_Q^2}{\|\mathbf{g}^{(k)}\|_{Q^{-1}}^2} = \frac{\|\mathbf{g}^{(k)}\|_Q^2}{\|\mathbf{g}^{(k)}\|_2^2} \frac{\|\mathbf{g}^{(k)}\|_2^2}{\|\mathbf{g}^{(k)}\|_{Q^{-1}}^2} \xrightarrow[\text{inequality}]{\text{Rayleigh's}} \frac{\|\mathbf{g}^{(k)}\|_Q^2}{\|\mathbf{g}^{(k)}\|_{Q^{-1}}^2} \geq [\lambda_{\min}(\mathbf{Q})]^2$$

$$\because \frac{\|\mathbf{g}^{(k)}\|_2^2}{\|\mathbf{g}^{(k)}\|_Q^2} \geq \frac{1}{\lambda_{\max}(\mathbf{Q})}.$$

$$\therefore \gamma_k \geq \alpha [\lambda_{\min}(\mathbf{Q})]^2 \left(\frac{2}{\lambda_{\max}(\mathbf{Q})} - \alpha \right) \xrightarrow{\alpha \in (0, \frac{2}{\lambda_{\max}(\mathbf{Q})})} \text{Const} > 0$$

$$\therefore \sum_{i=0}^{\infty} \gamma_k = \infty \iff V(\mathbf{x}^{(k)}) \rightarrow 0 \iff \mathbf{x}^{(k)} \rightarrow \mathbf{x}^*.$$

★ $\alpha \in (0, \frac{2}{\lambda_{\max}(\mathbf{Q})})$ is necessary and sufficient condition, i.e., any α beyond this interval can't guarantee the convergence (contradictory proof in textbook).



Convergence Analysis

Example (GD with fixed stepsize: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}$)

$$\text{Minimize } f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 4 & 2\sqrt{2} \\ 0 & 5 \end{bmatrix} \mathbf{x} + \mathbf{x}^\top \begin{bmatrix} 3 \\ 6 \end{bmatrix} + 24.$$

Ans: $\mathbf{Q} = \frac{1}{2} \left(\begin{bmatrix} 4 & 2\sqrt{2} \\ 0 & 5 \end{bmatrix} + \begin{bmatrix} 4 & 2\sqrt{2} \\ 0 & 5 \end{bmatrix}^\top \right) = \begin{bmatrix} 4 & \sqrt{2} \\ \sqrt{2} & 5 \end{bmatrix}$

$$\implies \lambda_{\max}(\mathbf{Q}) = 12$$

\therefore GD with fixed stepsize $\alpha \in (0, \frac{1}{6})$ converges to \mathbf{x}^* .



Convergence Analysis

Theorem (convergence rate of steepest descent method)

If steepest descent method minimizes quadratic function, then

$$V(\mathbf{x}^{(k+1)}) \leq \left(\frac{\kappa-1}{\kappa+1}\right)^2 V(\mathbf{x}^{(k)}).$$

proof. $\because V(\mathbf{x}^{(k+1)}) = (1 - \gamma_k) V(\mathbf{x}^{(k)})$ and $\frac{4\lambda_{\max}(\mathbf{Q})\lambda_{\min}(\mathbf{Q})}{[\lambda_{\max}(\mathbf{Q}) + \lambda_{\min}(\mathbf{Q})]^2} \leq \gamma_k$.

$$\begin{aligned} V(\mathbf{x}^{(k+1)}) &\leq \frac{[\lambda_{\max}(\mathbf{Q}) - \lambda_{\min}(\mathbf{Q})]^2}{[\lambda_{\max}(\mathbf{Q}) + \lambda_{\min}(\mathbf{Q})]^2} V(\mathbf{x}^{(k)}) \\ &\xrightarrow{\kappa = \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})}} V(\mathbf{x}^{(k+1)}) \leq \frac{[\kappa-1]^2}{[\kappa+1]^2} V(\mathbf{x}^{(k)}). \end{aligned}$$

★ $c = \frac{[\kappa-1]^2}{[\kappa+1]^2}$ is called convergent rate. [converge fast if $c \downarrow 0$; and vice verse]

★ one step convergent (i.e., $c = 0$) $\iff \lambda_{\max}(\mathbf{Q}) = \lambda_{\min}(\mathbf{Q})$.



Convergence Rate

Definition (Let $\{\mathbf{x}^k\}_{k=0}^{\infty}$ be a sequence converging to \mathbf{x}^*)

The order of convergence of the sequence $\{\mathbf{x}^k\}_{k=0}^{\infty}$ is defined by

$$\begin{cases} p, & \text{if } \lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_2}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^p} \in (0, \infty), \\ \infty, & \text{if } \lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_2}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^p} = 0, \forall p > 0 \end{cases}$$

★ The larger the p is, the faster the convergence of $\{\mathbf{x}^k\}_{k=0}^{\infty}$

particular cases: let $\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_2}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^p} = \mu$

- sublinear: $p = 1$ and $\mu = 1$;
- linear: $p = 1$ and $\mu \in (0, 1)$;
- superlinear: $p = 1$ and $\mu = 0$ (also $p > 1$);
- quadratic: $p = 2$.

$$x^{(k)} = \frac{1}{k}$$

$$\therefore x^* = \lim_{k \rightarrow \infty} x^{(k)} = 0$$

$$\therefore \lim_{k \rightarrow \infty} \frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^p} = \lim_{k \rightarrow \infty} \frac{k^p}{k+1} = \begin{cases} 0, & p < 1, \\ 1, & p = 1, \\ \infty & p > 1, \end{cases} \Rightarrow \text{1-order and sublinear}$$

Convergence Rate

$$x^{(k)} = \gamma^k \text{ with } \gamma \in (0, 1)$$

$$\therefore x^* = \lim_{k \rightarrow \infty} x^{(k)} = 0$$

$$\therefore \lim_{k \rightarrow \infty} \frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^p} = \lim_{k \rightarrow \infty} \gamma^{k(1-p)+1} \Rightarrow \begin{cases} 0, & p < 1 \\ \gamma, & p = 1 \Rightarrow \text{1-order and linear} \\ \infty, & p > 1 \end{cases}$$

$$x^{(k)} = \gamma^{q^k} \text{ with } q > 1 \text{ and } \gamma \in (0, 1)$$

$$\therefore x^* = \lim_{k \rightarrow \infty} x^{(k)} = 0$$

$$\therefore \lim_{k \rightarrow \infty} \frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^p} = \lim_{k \rightarrow \infty} \gamma^{(q-p)q^k} \Rightarrow \begin{cases} 0, & p < q \\ 1, & p = q \Rightarrow q\text{-order} \\ \infty, & p > 1 \end{cases}$$

$$x^{(k)} = 1$$

$$\therefore x^* = \lim_{k \rightarrow \infty} x^{(k)} = 1, \therefore \lim_{k \rightarrow \infty} \frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^p} = 0, \forall p > 0 \Rightarrow \infty\text{-order}$$

Convergence Rate

Definition

$a = O(h)$: $\exists c$ such that $|a| \leq c|h|$ for sufficiently small h .

$a = o(h)$: $\frac{a}{h} \rightarrow 0$ as $h \rightarrow 0$.

Definition

If $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_2 = O(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^p)$, then $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ in $\geq p$ -order.

If $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_2 = o(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^p)$, then $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ in $> p$ -order.

Example (GD for minimizing $f(x) = x^2 - \frac{x^3}{3}$ with $\alpha_k \equiv \frac{1}{2}$, $x^{(0)} = 1$)

$$\therefore x^{(k+1)} = x^{(k)} - \frac{1}{2}f'(x^{(k)}) \implies x^{(k+1)} = \frac{1}{2}(x^{(k)})^2 \implies x^* = 0.$$

$$\therefore \lim_{k \rightarrow \infty} \frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^2} = \lim_{k \rightarrow \infty} \frac{|x^{(k+1)}|}{|x^{(k)}|^2} = \frac{1}{2} \implies 2\text{-order}$$

$$\star x^{(k+1)} = \frac{1}{2}(x^{(k)})^2 \xrightarrow{\text{why?}} x^{(k+1)} = \left(\frac{1}{2}\right)^{2^{k-1}}$$



Convergence Rate

Question: For $V(\mathbf{x}^{(k+1)}) = (1 - \gamma_k)V(\mathbf{x}^{(k)})$, when $\gamma_k = 1$?

Lemma

For steepest descent method, $\gamma_k = 1$ iff $\mathbf{g}^{(k)}$ is an eigenvector of \mathbf{Q} .

proof. (\Leftarrow) if $\mathbf{g}^{(k)}$ is an eigenvector of $\mathbf{Q} \implies \gamma_k = \frac{\|\mathbf{g}^{(k)}\|_2^4}{\|\mathbf{g}^{(k)}\|_Q^2 \|\mathbf{g}^{(k)}\|_{Q^{-1}}^2} = 1$.

(\Rightarrow) if $\gamma_k = 1 \implies V(\mathbf{x}^{(k+1)}) = 0 \implies \mathbf{x}^{(k+1)} = \mathbf{x}^*$

$$\therefore \mathbf{x}^* = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}$$

$$\therefore \mathbf{Q}\mathbf{x}^* = \mathbf{Q}\mathbf{x}^{(k)} - \alpha_k \mathbf{Q}\mathbf{g}^{(k)}, \text{ i.e.,}$$

$$\frac{1}{\alpha_k} \mathbf{g}^{(k)} = \mathbf{Q}\mathbf{g}^{(k)} \implies \mathbf{g}^{(k)} \text{ is an eigenvector of } \mathbf{Q}$$

★ (inversion of Lemma) if $\mathbf{g}^{(k)}$ isn't an eigenvector of \mathbf{Q} , then $\gamma_k < 1$.

homework

exercise in textbook: 8.4, 8.10, 8.16