

3 GRADIENT METHODS

3.1 INTRODUCTION

- $\langle \nabla f(\mathbf{x}), \mathbf{d} \rangle$ with $\|\mathbf{d}\| = 1$, is the rate of increase of f in the direction \mathbf{d} at the point \mathbf{x} . By the Cauchy-Schwarz inequality,

$$\langle \nabla f(\mathbf{x}), \mathbf{d} \rangle \leq \|\nabla f(\mathbf{x})\| \|\mathbf{d}\| = \|\nabla f(\mathbf{x})\|$$

because $\|\mathbf{d}\| = 1$.

- The direction of maximum rate of **increase** of f at \mathbf{x} is $\mathbf{d} = \nabla f(\mathbf{x}) / \|\nabla f(\mathbf{x})\|$ because

$$\left\langle \nabla f(\mathbf{x}), \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|} \right\rangle = \frac{(\nabla f(\mathbf{x}))^\top \nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|} = \frac{\|\nabla f(\mathbf{x})\|^2}{\|\nabla f(\mathbf{x})\|} = \|\nabla f(\mathbf{x})\|.$$

- The direction of **maximum rate of decrease** is $-\nabla f(\mathbf{x}) / \|\nabla f(\mathbf{x})\|$.
- Let $\mathbf{x}^{(0)}$ be a starting point and consider the point $\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})$. The first-order Taylor expansion of a function f around a point $\mathbf{x}^{(0)}$ is given by

$$f(\mathbf{x}) = f(\mathbf{x}^{(0)}) + \nabla f(\mathbf{x}^{(0)})^\top (\mathbf{x} - \mathbf{x}^{(0)}) + o(\mathbf{x} - \mathbf{x}^{(0)})$$

- Then, we have

$$\begin{aligned} f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})) &= f(\mathbf{x}^{(0)}) + \nabla f(\mathbf{x}^{(0)})^\top (\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)}) - \mathbf{x}^{(0)}) \\ &\quad + o(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)}) - \mathbf{x}^{(0)}) \\ &= f(\mathbf{x}^{(0)}) - \alpha \nabla f(\mathbf{x}^{(0)})^\top \nabla f(\mathbf{x}^{(0)}) + o(-\alpha \nabla f(\mathbf{x}^{(0)})) \\ &= f(\mathbf{x}^{(0)}) - \alpha \|\nabla f(\mathbf{x}^{(0)})\|^2 + o(\alpha). \end{aligned}$$

- If $\nabla f(\mathbf{x}^{(0)}) \neq 0$, then for sufficiently small $\alpha > 0$, we have

$$f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})) < f(\mathbf{x}^{(0)}).$$

- This means that the point $\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})$ is an improvement over the point $\mathbf{x}^{(0)}$ if we are searching for a minimizer.
- Suppose that we are given a point $\mathbf{x}^{(k)}$. To find the next point $\mathbf{x}^{(k+1)}$ we start at $\mathbf{x}^{(k)}$ and move by an amount $-\alpha_k \nabla f(\mathbf{x}^{(k)})$, where α_k is a positive scalar called the **step size**. This procedure leads to the gradient descent algorithm:

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

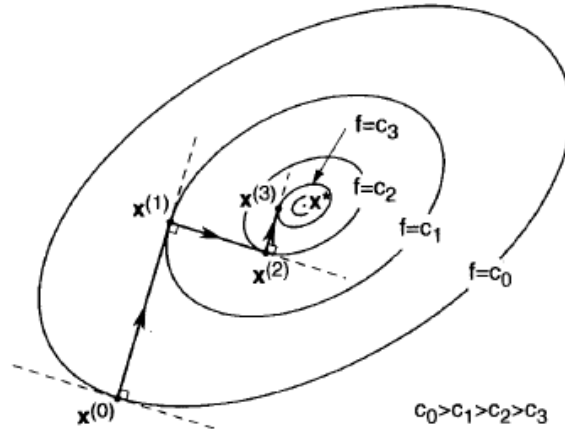
- $\{\mathbf{x}^{(k)}\}_{k=0}^\infty$ is a **steepest descent sequence** if $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$ where $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$.

3.2 THE METHOD OF STEEPEST DESCENT

- The method of steepest descent is a gradient algorithm where the step size α_k is chosen to achieve the **maximum amount of decrease** of the objective function at each individual step. Specifically, α_k is chosen to minimize $\phi_k(\alpha) \triangleq f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$. In other words,

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)})).$$

- Proposition 3.2.1** If $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$ is a steepest descent sequence for a given function $f: \mathbb{R}^n \rightarrow \mathbb{R}$, then for each k the vector $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$ is orthogonal to the vector $\mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)}$.



- Proposition 3.2.2** If $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$ is a steepest descent sequence for $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and if $\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$, then $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$.
- If for some k , we have $\nabla f(\mathbf{x}^{(k)}) = \mathbf{0}$, then the point $\mathbf{x}^{(k)}$ satisfies the FONC. The condition $\nabla f(\mathbf{x}^{(k)}) = \mathbf{0}$, however, is not directly suitable as a practical stopping criterion, because the numerical computation of the gradient will rarely be identically equal to zero.
- A practical stopping criterion is to check if, given a prespecified threshold $\varepsilon > 0$,
 - $\|\nabla f(\mathbf{x}^{(k)})\| < \varepsilon$;
 - $|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})| < \varepsilon$;
 - $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < \varepsilon$;
 - $|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})| / |f(\mathbf{x}^{(k)})| < \varepsilon$;
 - $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| / \|\mathbf{x}^{(k)}\| < \varepsilon$.
- The criteria 4 and 5 above are preferable to the criteria 1, 2 and 3 because the relative criteria are “**scale-independent**.” To avoid dividing by very small numbers in criteria 4 and 5, we can modify these stopping criteria as follows:

$$\frac{|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})|}{\max\{1, |f(\mathbf{x}^{(k)})|\}} < \varepsilon \quad \text{or} \quad \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|}{\max\{1, \|\mathbf{x}^{(k)}\|\}} < \varepsilon.$$

Example 3.2.1. Perform three iterations of the method of steepest descent to find the minimizer of $f(x_1, x_2, x_3) = (x_1 - 4)^4 + (x_2 - 3)^2 + 4(x_3 + 5)^4$. The initial point is $\mathbf{x}^{(0)} = [4, 2, -1]^\top$.

Note that

$$\nabla f(\mathbf{x}) = [4(x_1 - 4)^3, 2(x_2 - 3), 16(x_3 + 5)^3]^\top$$

$$\mathbf{x} - \alpha \nabla f(\mathbf{x}) = [x_1 - 4\alpha(x_1 - 4)^3, x_2 - 2\alpha(x_2 - 3), x_3 - 16\alpha(x_3 + 5)^3]^\top$$

$$f(\mathbf{x} - \alpha \nabla f(\mathbf{x})) = (x_1 - 4\alpha(x_1 - 4)^3 - 4)^4 + (x_2 - 2\alpha(x_2 - 3) - 3)^2 + 4(x_3 - 16\alpha(x_3 + 5)^3 + 5)^4$$

and $f(\mathbf{x}^{(0)}) = 1025$.

Iteration 1:

To compute $\mathbf{x}^{(1)}$, we need

$$\alpha_0 = \operatorname{argmin}_{\alpha \geq 0} f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)}))$$

Let $\mathbf{h}(\alpha) = \mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})$ and $g(\alpha) = f(\mathbf{h}(\alpha))$. Using the secant method, we obtain $\alpha_0 = 3.967 \times 10^{-3}$:

$$\alpha^{(k+1)} = \alpha^{(k)} - \frac{\alpha^{(k)} - \alpha^{(k-1)}}{g'(\alpha^{(k)}) - g'(\alpha^{(k-1)})} g'(\alpha^{(k)})$$

where $g'(\alpha) = (\nabla f(\mathbf{h}(\alpha)))^\top \begin{bmatrix} h'_1(\alpha) \\ \vdots \\ h'_1(\alpha) \end{bmatrix} = -(\nabla f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})))^\top \nabla f(\mathbf{x}^{(0)})$. Thus,

$$\begin{aligned} \mathbf{x}^{(1)} &= \mathbf{x}^{(0)} - \alpha_0 \nabla f(\mathbf{x}^{(0)}) \\ &= [4, 2, -1]^\top - 3.967 \times 10^{-3} [4(4 - 4)^3, 2(2 - 3), 16(-1 + 5)^3]^\top \\ &= [4, 2, -1]^\top - 3.967 \times 10^{-3} [0, -2, 1024]^\top \\ &\approx [4.000, 2.007934, -5.061568]^\top \\ &\approx [4.000, 2.008, -5.062]^\top, \end{aligned}$$

and $f(\mathbf{x}^{(1)}) \approx 0.984$.

Iteration 2:

We find $\alpha_1 \approx 0.5000$, where $\alpha_1 = \operatorname{argmin}_{\alpha \geq 0} f(\mathbf{x}^{(1)} - \alpha \nabla f(\mathbf{x}^{(1)}))$. Thus,

$$\begin{aligned} \mathbf{x}^{(2)} &= \mathbf{x}^{(1)} - \alpha_1 \nabla f(\mathbf{x}^{(1)}) \\ &= [4.000, 2.008, -5.062]^\top - 0.500 [4(4 - 4)^3, 2(2.008 - 3), 16(-5.062 + 5)^3]^\top \\ &\approx [4.000, 2.008, -5.062]^\top - 0.500 [0, -1.9841, -0.00877]^\top \\ &\approx [4.000, 3.000, -5.060]^\top, \end{aligned}$$

and $f(\mathbf{x}^{(2)}) = 5.326 \times 10^{-5}$.

Iteration 3:

We find $\alpha_2 \approx 16.28$, where $\alpha_2 = \underset{\alpha \geq 0}{\operatorname{argmin}} f(\mathbf{x}^{(2)} - \alpha \nabla f(\mathbf{x}^{(2)}))$. Thus,

$$\begin{aligned}\mathbf{x}^{(3)} &= \mathbf{x}^{(2)} - \alpha_2 \nabla f(\mathbf{x}^{(2)}) \\ &= [4.000, 3.000, -5.060]^\top - 16.28[4(4-4)^3, 2(3-3), 16(-5.060+5)^3]^\top \\ &\approx [4.000, 3.000, -5.060]^\top - 16.28[0.000, 0.000, 0.004]^\top \\ &\approx [4.000, 3.000, -5.003]^\top,\end{aligned}$$

and $f(\mathbf{x}^{(3)}) = 1.215 \times 10^{-8}$.

Example 3.2.2. Consider an objection function which is in quadratic form

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{b}^\top \mathbf{x},$$

where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix, $\mathbf{b} \in \mathbb{R}^n$, and $\mathbf{x} \in \mathbb{R}^n$. Show that the method of steepest descent takes the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}} \mathbf{g}^{(k)},$$

where $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) = \mathbf{Q} \mathbf{x}^{(k)} - \mathbf{b}$.

There is no loss of generality in assuming \mathbf{Q} to be a symmetric matrix. For if we are given a quadratic form $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ and $\mathbf{A} \neq \mathbf{A}^\top$, then because the transposition of a scalar equals itself, we obtain

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = (\mathbf{x}^\top \mathbf{A} \mathbf{x})^\top = \mathbf{x}^\top \mathbf{A}^\top \mathbf{x}.$$

Hence,

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \mathbf{A}^\top \mathbf{x} = \frac{1}{2} \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top) \mathbf{x} \triangleq \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x}.$$

Note that $(\mathbf{A} + \mathbf{A}^\top)^\top = \mathbf{Q}^\top = (\mathbf{A} + \mathbf{A}^\top) = \mathbf{Q}$.

The unique minimizer of f can be found by setting the gradient of f to zero:

$$\nabla f(\mathbf{x}) = \frac{1}{2} (\mathbf{Q} + \mathbf{Q}^\top)^\top \mathbf{x} - \mathbf{b} = \mathbf{Q} \mathbf{x} - \mathbf{b} = \mathbf{0}.$$

The Hessian of f is $\mathbf{F}(\mathbf{x}) = \mathbf{Q} = \mathbf{Q}^\top > 0$. Let $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$. The steepest descent algorithm for the quadratic function is

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)},$$

where

$$\begin{aligned}\alpha_k &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}) \\ &= \arg \min_{\alpha \geq 0} \left(\frac{1}{2} (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)})^\top \mathbf{Q} (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}) - \mathbf{b}^\top (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}) \right).\end{aligned}$$

Because $\alpha_k \geq 0$ is a minimizer of $\phi_k(\alpha) = f(\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)})$, we apply the FONC to $\phi_k(\alpha)$ to obtain

$$\begin{aligned} 0 &= \phi'_k(\alpha) \\ 0 &= (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)})^\top \mathbf{Q}(-\mathbf{g}^{(k)}) - \mathbf{b}^\top (-\mathbf{g}^{(k)}) \\ 0 &= -\mathbf{x}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)} + \alpha \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)} + \mathbf{b}^\top \mathbf{g}^{(k)} \\ \mathbf{x}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)} - \mathbf{b}^\top \mathbf{g}^{(k)} &= \alpha \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)} \\ (\mathbf{x}^{(k)\top} \mathbf{Q} - \mathbf{b}^\top) \mathbf{g}^{(k)} &= \alpha \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)} \\ \mathbf{g}^{(k)\top} \mathbf{g}^{(k)} &= \alpha \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)} \\ \alpha &= \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}}. \end{aligned}$$

Therefore, the steepest descent algorithm for the quadratic function is

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}} \mathbf{g}^{(k)},$$

where $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) = \mathbf{Q} \mathbf{x}^{(k)} - \mathbf{b}$.

Example 3.2.3. State the steepest descent algorithms for each of the following objection functions:

(a) $f(x_1, x_2) = x_1^2 + x_2^2$.

(b) $f(x_1, x_2) = x_1^2/5 + x_2^2$.

(a) $f(x_1, x_2) = x_1^2 + x_2^2 = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{x}$

The steepest descent algorithm for the function is

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}} \mathbf{g}^{(k)},$$

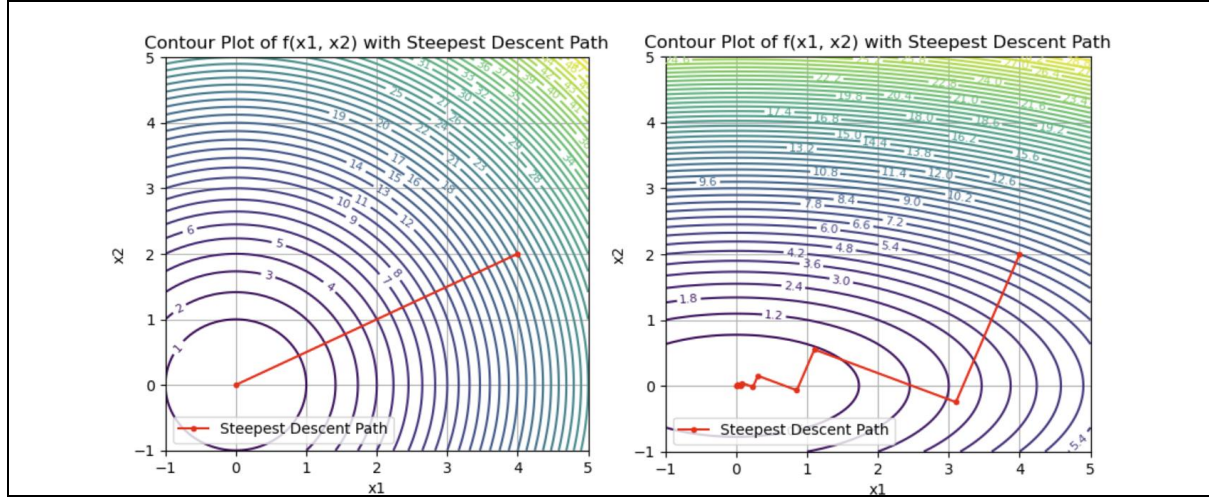
where $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) = \mathbf{Q} \mathbf{x}^{(k)}$ and $\mathbf{Q} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$.

(b) $f(x_1, x_2) = \frac{x_1^2}{5} + x_2^2 = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 2/5 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{x}$

The steepest descent algorithm for the function is

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}} \mathbf{g}^{(k)},$$

where $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) = \mathbf{Q} \mathbf{x}^{(k)}$ and $\mathbf{Q} = \begin{bmatrix} 2/5 & 0 \\ 0 & 2 \end{bmatrix}$.



3.3 THE CONDITION NUMBER

- Consider the quadratic minimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^\top \mathbf{A} \mathbf{x}$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix. The optimal solution is obviously $\mathbf{x}^* = \mathbf{0}$. The steepest descent algorithm is

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)},$$

where $\mathbf{g}^{(k)} = 2\mathbf{A}\mathbf{x}^{(k)}$ is the gradient of $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$ at $\mathbf{x}^{(k)}$ and the stepsize α_k is chosen by the minimization rule

$$\alpha_k = \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{2\mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)}}.$$

- Therefore,

$$\begin{aligned} f(\mathbf{x}^{(k+1)}) &= \mathbf{x}^{(k+1)\top} \mathbf{A} \mathbf{x}^{(k+1)} = (\mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)})^\top \mathbf{A} (\mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}) \\ &= \mathbf{x}^{(k)\top} \mathbf{A} \mathbf{x}^{(k)} - 2\alpha_k \mathbf{g}^{(k)\top} \mathbf{A} \mathbf{x}^{(k)} + \alpha_k^2 \mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)} \\ &= \mathbf{x}^{(k)\top} \mathbf{A} \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)\top} \mathbf{g}^{(k)} + \alpha_k^2 \mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)} \\ &= \mathbf{x}^{(k)\top} \mathbf{A} \mathbf{x}^{(k)} - \left(\frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{2\mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)}} \right) \mathbf{g}^{(k)\top} \mathbf{g}^{(k)} + \left(\frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{2\mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)}} \right)^2 \mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)} \\ &= \mathbf{x}^{(k)\top} \mathbf{A} \mathbf{x}^{(k)} - \frac{1}{2} \frac{(\mathbf{g}^{(k)\top} \mathbf{g}^{(k)})^2}{\mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)}} + \frac{1}{4} \frac{(\mathbf{g}^{(k)\top} \mathbf{g}^{(k)})^2}{\mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)}} \\ &= \mathbf{x}^{(k)\top} \mathbf{A} \mathbf{x}^{(k)} - \frac{1}{4} \frac{(\mathbf{g}^{(k)\top} \mathbf{g}^{(k)})^2}{\mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)}} \frac{\mathbf{x}^{(k)\top} \mathbf{A} \mathbf{x}^{(k)}}{\mathbf{x}^{(k)\top} \mathbf{A} \mathbf{x}^{(k)}} \\ &= \mathbf{x}^{(k)\top} \mathbf{A} \mathbf{x}^{(k)} \left(1 - \frac{1}{4} \frac{(\mathbf{g}^{(k)\top} \mathbf{g}^{(k)})^2}{(\mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)}) (\mathbf{x}^{(k)\top} \mathbf{A} \mathbf{x}^{(k)})} \right) \\ &= f(\mathbf{x}^{(k)}) \left(1 - \frac{(\mathbf{g}^{(k)\top} \mathbf{g}^{(k)})^2}{(\mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)}) (2\mathbf{x}^{(k)\top} \mathbf{A} \mathbf{A}^{-1} 2\mathbf{A} \mathbf{x}^{(k)})} \right) \end{aligned}$$

$$= f(\mathbf{x}^{(k)}) \left(1 - \frac{(\mathbf{g}^{(k)\top} \mathbf{g}^{(k)})^2}{(\mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)}) (\mathbf{g}^{(k)\top} \mathbf{A}^{-1} \mathbf{g}^{(k)})} \right)$$

- **Theorem 3.3.1 (Kantorovich inequality)** Let \mathbf{A} be a positive definite $n \times n$ matrix. Then for any $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$,

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

- Hence, by Kantorovich inequality,

$$\frac{(\mathbf{g}^{(k)\top} \mathbf{g}^{(k)})^2}{(\mathbf{g}^{(k)\top} \mathbf{A} \mathbf{g}^{(k)}) (\mathbf{g}^{(k)\top} \mathbf{A}^\top \mathbf{g}^{(k)})} \geq \frac{4\lambda_{\max}(\mathbf{A})\lambda_{\min}(\mathbf{A})}{(\lambda_{\max}(\mathbf{A}) + \lambda_{\min}(\mathbf{A}))^2}$$

and

$$f(\mathbf{x}^{(k+1)}) \leq \left(1 - \frac{4Mm}{(M+m)^2} \right) f(\mathbf{x}^{(k)}) = \left(\frac{M-m}{M+m} \right)^2 f(\mathbf{x}^{(k)}) = \left(\frac{M/m - 1}{M/m + 1} \right)^2 f(\mathbf{x}^{(k)}),$$

where $M = \lambda_{\max}(\mathbf{A})$, $m = \lambda_{\min}(\mathbf{A})$.

- The speed of convergence depends on the ratio $\lambda_{\max}(\mathbf{A})/\lambda_{\min}(\mathbf{A})$. As $\lambda_{\max}(\mathbf{A})/\lambda_{\min}(\mathbf{A})$ gets larger, the convergence speed becomes slower. The ratio $\lambda_{\max}(\mathbf{A})/\lambda_{\min}(\mathbf{A})$ is called the **condition number**.
- Matrices with large condition number are called **ill-conditioned**, and matrices with small condition number are called **well-conditioned**.
- The discussion assumes quadratic objective functions, where the Hessian matrix is constant. In practice, the matrix \mathbf{A} is replaced by the Hessian matrix $\nabla^2 f(\mathbf{x}^*)$.

Example 3.3.1. The Rosenbrock function is defined as:

$$f(x_1, x_2) = 100(x_2 - x_1^2)^2 - (1 - x_1)^2.$$

Given that $\mathbf{x}^* = [1, 1]$ is the unique stationary point, find the condition number.

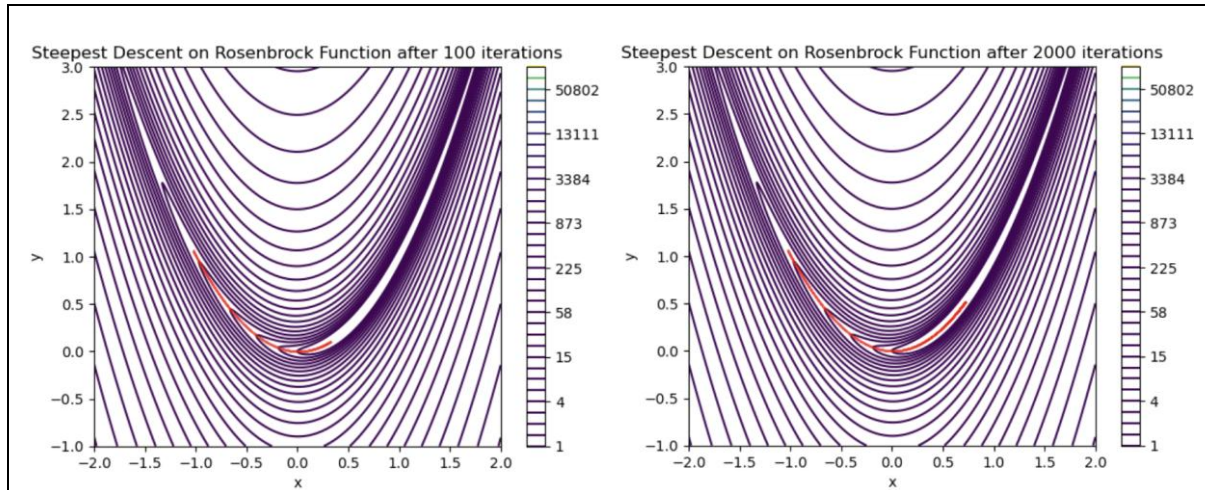
$$\nabla f(\mathbf{x}) = \begin{bmatrix} -400x_1(x_2 - x_1^2) - 2(1 - x_1) \\ 200(x_2 - x_1^2) \end{bmatrix} \quad \nabla^2 f(\mathbf{x}) = \begin{bmatrix} -400x_2 + 1200x_1^2 + 2 & -400x_1 \\ -400x_1 & 200 \end{bmatrix}$$

Hence,

$$\nabla^2 f(\mathbf{x}^*) = \begin{bmatrix} -400 + 1200 + 2 & -400 \\ -400 & 200 \end{bmatrix} = \begin{bmatrix} 802 & -400 \\ -400 & 200 \end{bmatrix}$$

Therefore, the condition number is

$$\frac{\lambda_{\max}(\nabla^2 f(\mathbf{x}^*))}{\lambda_{\min}(\nabla^2 f(\mathbf{x}^*))} \approx \frac{1001}{0.4} = 2502.5$$



3.3 THE ORDER OF CONVERGENCE

- The **order of convergence** of a sequence quantifies **how quickly** the sequence approaches its limit; a **higher order** indicates a **faster rate** of convergence.
- Let $\{\mathbf{x}^{(k)}\}$ be the sequence that converges to \mathbf{x}^* , that is, $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0$. We say that the order of convergence is p , where $p \in \mathbb{R}$, if

$$0 < \lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} < \infty.$$

If, for all $p > 0$,

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} = 0,$$

then we say that the order of convergence is ∞ . Note that in the definition above, $0/0$ should be understood to be 0.

Example 3.3.1. Determine the order of convergence of each of the following sequences:

- $x^{(k)} = 1/k$.
- $x^{(k)} = \gamma^k$, where $0 < \gamma < 1$.
- $x^{(k)} = \gamma^{(q^k)}$, where $q > 1$ and $0 < \gamma < 1$.
- $x^{(k)} = 1$ for all k .

(a) Note that $x^{(k)} \rightarrow x^* = 0$. Then,

$$\frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^p} = \frac{|x^{(k+1)}|}{|x^{(k)}|^p} = \frac{1/(k+1)}{1/k^p} = \frac{k^p}{k+1}.$$

If $p < 1$, the sequence converges to 0. If $p > 1$, it grows to ∞ . If $p = 1$, the sequence converges to 1. Hence, the order of convergence is 1.

(b) Note that $x^{(k)} \rightarrow x^* = 0$. Then,

$$\frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^p} = \frac{\gamma^{k+1}}{(\gamma^k)^p} = \gamma^{k+1-kp} = \gamma^{k(1-p)+1}.$$

If $p < 1$, the sequence converges to 0. If $p > 1$, it grows to ∞ . If $p = 1$, the sequence converges to $\gamma < \infty$. Hence, the order of convergence is 1.

(c) Note that $x^{(k)} \rightarrow x^* = 0$. Then,

$$\frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^p} = \frac{\gamma^{(q^{k+1})}}{(\gamma^{(q^k)})^p} = \gamma^{(q^{k+1}-pq^k)} = \gamma^{(q-p)q^k}.$$

If $p < q$, the sequence converges to 0. If $p > q$, it grows to ∞ . If $p = q$, the sequence converges to $1 < \infty$. Hence, the order of convergence is q .

(d) Note that $x^{(k)} \rightarrow x^* = 1$. Then,

$$\frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^p} = \frac{|x^{(k+1)} - 1|}{|x^{(k)} - 1|^p} = \frac{0}{0^p} = 0$$

for all p . Hence, the order of convergence is ∞ .

Example 3.3.2 Consider the problem of finding a minimizer of the function $f: \mathbb{R} \rightarrow \mathbb{R}$ given by

$$f(x) = x^2 - \frac{x^3}{3}.$$

Suppose that we use the algorithm $x^{(k+1)} = x^{(k)} - \alpha f'(x^{(k)})$ with step size $\alpha = 1/2$ and initial condition $x^{(0)} = 1$. Show that the order of convergence is 2.

We first show that the algorithm converges to a local minimizer of f . Indeed, we have $f'(x) = 2x - x^2$. The fixed-step-size gradient algorithm with step size $\alpha = 1/2$ is therefore given by

$$x^{(k+1)} = x^{(k)} - \alpha f'(x^{(k)}) = \frac{1}{2}(x^{(k)})^2.$$

With $x^{(0)} = 1$, we can derive the expression $x^{(k)} = (1/2)^{2^k-1}$. Hence, the algorithm converges to 0, a strict local minimizer of f .

When $p = 2$, we have

$$\frac{|x^{(k+1)}|}{|x^{(k)}|^p} = \frac{(1/2)^{2^{k+1}-1}}{((1/2)^{2^k-1})^p} = (1/2)^{2^{k+1}-1-p2^k+2p} = (1/2)^{2^{k+1}-1-2 \cdot 2^k+2} = \frac{1}{2}.$$

Therefore, the order of convergence is 2.

What happen if we try $p = 1$ or 3:

$$\frac{|x^{(k+1)}|}{|x^{(k)}|} = \frac{(1/2)^{2^{k+1}-1}}{(1/2)^{2^k-1}} = (1/2)^{2^{k+1}-1-2^k+1} = \frac{1}{2^2}.$$

$$\frac{|x^{(k+1)}|}{|x^{(k)}|^3} = \frac{(1/2)^{2^{k+1}-1}}{(1/2)^{3(2^k-1)}} = (1/2)^{2^{k+1}-1-3 \cdot 2^k+3} = (1/2)^{-2^k+2} = 2^{2^k-2} \rightarrow \infty.$$

- In the analysis, we assume the objection function is a quadratic function of the form

$$f(x) = \frac{1}{2}x^\top Qx - b^\top x,$$

where $Q \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix, $b \in \mathbb{R}^n$, and $x \in \mathbb{R}^n$.

- In example 3.2.2, the steepest descent algorithm for the quadratic function is derived as

$$x^{(k+1)} = x^{(k)} - \frac{g^{(k)\top} g^{(k)}}{g^{(k)\top} Q g^{(k)}} g^{(k)},$$

where $g^{(k)} = \nabla f(x^{(k)}) = Qx^{(k)} - b$.

- **Theorem 3.3.1** Let $\{x^{(k)}\}$ be a convergent sequence of iterates of the steepest descent algorithm applied to a function f . Then, the order of convergence of is 1 in the worst case; that is, there exist a function f and an initial condition $x^{(0)}$ such that the order of convergence of $\{x^{(k)}\}$ is 1.