

Analysis of Steepest decent

Oscar Dalmau
dalmau@cimat.mx

Centro de Investigación en Matemáticas
CIMAT A.C. Mexico

February 2018

Outline

- ① Steepest descent Method
- ② Global Convergence
- ③ Rate of convergence

Steepest descent Method with exact line search

- The method of *steepest descent with exact step size or with exact line search* is a gradient algorithm where the step size is obtained by solving

$$\alpha_k = \arg \min_{\alpha > 0} \phi(\alpha)$$

with $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$, $\mathbf{d}_k = -\mathbf{g}_k$.

- The **previous method moves in orthogonal steps**, see next proposition.

Steepest descent Method with exact line search

Proposition 1.1

(Orthogonality of directions) If $\{\mathbf{x}_k\}_{k=1}^{\infty}$ 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 $\mathbf{x}_{k+2} - \mathbf{x}_{k+1}$

Remark

- 1 Note that $\alpha_k \mathbf{d}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ and $\alpha_{k+1} \mathbf{d}_{k+1} = \mathbf{x}_{k+2} - \mathbf{x}_{k+1}$. Therefore, the proposition states that two consecutive directions are orthogonal, i.e. $\mathbf{d}_k \perp \mathbf{d}_{k+1}$, where $\mathbf{d}_k = -\mathbf{g}_k$ and $\mathbf{d}_{k+1} = -\mathbf{g}_{k+1}$ then $\mathbf{g}_k \perp \mathbf{g}_{k+1}$.
- 2 The solution trajectory of the steepest-descent method with exact line search follows a zig-zag pattern.

Steepest descent Method with exact line search

Proposition

If $\{\mathbf{x}_k\}_{k=1}^{\infty}$ 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 $\mathbf{x}_{k+2} - \mathbf{x}_{k+1}$

Proof.

As α_k minimizes $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ then $\phi'(\alpha_k) = 0$. Using $\mathbf{d}_k = -\mathbf{g}_k = -\nabla f(\mathbf{x}_k)$,

$$\phi'(\alpha_k) = \nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k = \nabla f(\mathbf{x}_{k+1})^T \mathbf{d}_k = -\mathbf{g}_{k+1}^T \mathbf{g}_k$$

then $\mathbf{g}_k \perp \mathbf{g}_{k+1}$.



Steepest descent Method with exact line search

Proposition

If $\{\mathbf{x}_k\}_{k=1}^{\infty}$ 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 $\mathbf{x}_{k+2} - \mathbf{x}_{k+1}$

Corollary

\mathbf{g}_k is parallel to the tangent plane to the level set $\{\mathbf{x} \mid f(\mathbf{x}) = f(\mathbf{x}_{k+1})\}$ at \mathbf{x}_{k+1} .

Steepest descent Method with exact line search

Proposition 1.2

If $\{\mathbf{x}_k\}_{k=1}^{\infty}$ is a steepest descent sequence for a given function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and if $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$ then $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$.

Steepest descent Method with exact line search

Proposition

If $\{\mathbf{x}_k\}_{k=1}^{\infty}$ is a steepest descent sequence for a given function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and if $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$ then $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$.

Proof.

As α_k minimizes $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ then $\phi(\alpha_k) \leq \phi(\alpha)$ for all α . On the other hand, $\phi'(0) = \nabla f(\mathbf{x}_k)^T \mathbf{g}_k = -\|\mathbf{g}_k\|^2 < 0$ due to $\mathbf{g}_k \neq \mathbf{0}$.

Therefore, there exists $\hat{\alpha}$ (by the sign preserving theorem) such that $\phi'(\alpha) < 0$ for $\alpha \in (0, \hat{\alpha})$. Using Taylor (or the mean value theorem), there exists $\bar{\alpha}$ such that $\phi(\alpha) - \phi(0) = \phi'(\bar{\alpha})\alpha$ with $\bar{\alpha} \in (0, \alpha)$ then $\phi'(\bar{\alpha})\alpha < 0$ and hence $\phi(\alpha) < \phi(0)$ for $\alpha \in (0, \hat{\alpha})$. Then $\phi(\alpha_k) < \phi(0)$ for $\alpha \in (0, \hat{\alpha})$, i.e., $f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) < f(\mathbf{x}_k)$ or $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$.



Stopping criteria

- 1 The previous proposition states that the steepest descent has the property: $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ if $\mathbf{g}_k \neq \mathbf{0}$.
- 2 If for some k , it holds $\nabla f(\mathbf{x}_k) = \mathbf{0}$ then $\mathbf{x}_{k+1} = \mathbf{x}_k$. We can use this as a stopping criterion for the algorithm, however, the gradient will rarely be identically equal to zero.
- 3 A practical stopping criterion is to check if the norm $\|\nabla f(\mathbf{x}_k)\|$ of the gradient is less than a threshold, i.e. $\|\nabla f(\mathbf{x}_k)\| \leq \tau$.

Stopping criteria

Other stopping criteria

$$\begin{aligned} |f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)| &\leq \tau \\ \|\mathbf{x}_{k+1} - \mathbf{x}_k\| &\leq \tau \end{aligned}$$

We may check the relative values of the above quantities

$$\begin{aligned} \frac{|f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)|}{|f(\mathbf{x}_k)|} &\leq \tau \\ \frac{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|}{\|\mathbf{x}_k\|} &\leq \tau \end{aligned}$$

The above two (relative) stopping criteria are preferable to the previous (absolute) criteria because the relative criteria are scale-independent.

Stopping criteria

To avoid dividing by a small number we may use the following modifications

$$\boxed{\begin{aligned}\frac{|f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)|}{\max\{1, |f(\mathbf{x}_k)|\}} &\leq \tau_f \\ \frac{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|}{\max\{1, \|\mathbf{x}_k\|\}} &\leq \tau_x\end{aligned}}$$

Stopping criteria

To avoid dividing by a small number we may use the following modifications

$$\boxed{\begin{aligned}\frac{|f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)|}{\max\{1, |f(\mathbf{x}_k)|\}} &\leq \tau_f \\ \frac{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|}{\max\{1, \|\mathbf{x}_k\|\}} &\leq \tau_x\end{aligned}}$$

$$\boxed{\begin{aligned}\|\nabla f(\mathbf{x}_k)\| &\leq \tau_g \\ k &> K_{\max}\end{aligned}}$$

Exact Steepest descent for a quadratic function

Let $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q}\mathbf{x} - \mathbf{b}^T \mathbf{x}$, with \mathbf{Q} positive definite. As

$$\alpha_k = \arg \max_{\alpha > 0} \phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

Exact Step for a quadratic function

$$\phi(\alpha) = \frac{1}{2}(\mathbf{x}_k + \alpha \mathbf{d}_k)^T \mathbf{Q}(\mathbf{x}_k + \alpha \mathbf{d}_k) - \mathbf{b}^T(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

then, from $\phi'(\alpha) = 0$ and $\mathbf{g}_k = \mathbf{Q}\mathbf{x}_k - \mathbf{b} = -\mathbf{d}_k$

$$\begin{aligned}(\mathbf{x}_k + \alpha \mathbf{d}_k)^T \mathbf{Q} \mathbf{d}_k - \mathbf{b}^T \mathbf{d}_k &= 0 \\ \mathbf{d}_k^T \mathbf{Q} \mathbf{d}_k \alpha &= -(\mathbf{Q}\mathbf{x}_k - \mathbf{b})^T \mathbf{d}_k \\ \alpha_k &= \frac{-\mathbf{g}_k^T \mathbf{d}_k}{\mathbf{d}_k^T \mathbf{Q} \mathbf{d}_k} = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}\end{aligned}$$

Exact Steepest descent for a quadratic function

The update formula for the Steepest descent with Exact step size for the quadratic function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{b}^T \mathbf{x}$, with \mathbf{Q} positive definite, is:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$$

with

$$\begin{aligned} \mathbf{d}_k &= -\mathbf{g}_k \\ \alpha_k &= \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} \end{aligned}$$

Steepest descent: elimination of line search

If $f(x)$ is not quadratic but the **Hessian \mathbf{H}_k is available**: then, we can approximate

$$\phi(\alpha) = f(x_k + \alpha \mathbf{d}_k) \approx f(x_k) + \alpha \mathbf{g}_k^T \mathbf{d}_k + \frac{1}{2} \alpha^2 \mathbf{d}_k^T \mathbf{H}_k \mathbf{d}_k$$

from $\phi'(\alpha) = 0$ and $\mathbf{d}_k = -\mathbf{g}_k$

$$\alpha = \alpha_k \approx \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{H}_k \mathbf{g}_k}$$

then

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{H}_k \mathbf{g}_k} \mathbf{g}_k$$

Steepest descent: elimination of line search

If the Hessian \mathbf{H}_k is not available: suppose we have an estimate $\hat{\alpha}$ of α_k . Then $\hat{f} = f(x_k + \hat{\alpha} \mathbf{d}_k)$

$$\phi(\hat{\alpha}) = f(x_k + \hat{\alpha} \mathbf{d}_k) \approx f(x_k) + \hat{\alpha} \mathbf{g}_k^T \mathbf{d}_k + \frac{1}{2} \hat{\alpha}^2 \mathbf{d}_k^T \mathbf{H}_k \mathbf{d}_k$$

$$\mathbf{g}_k^T \mathbf{H}_k \mathbf{g}_k \approx \frac{2\hat{f} - f_k + \hat{\alpha} \mathbf{g}_k^T \mathbf{g}_k}{\hat{\alpha}^2}$$

from $\phi'(\alpha) = 0$ and $\mathbf{d}_k = -\mathbf{g}_k$

$$\alpha_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{H}_k \mathbf{g}_k} \approx \frac{\mathbf{g}_k^T \mathbf{g}_k \hat{\alpha}^2}{2(\hat{f} - f_k + \hat{\alpha} \mathbf{g}_k^T \mathbf{g}_k)}$$

then

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{\mathbf{g}_k^T \mathbf{g}_k \hat{\alpha}^2}{2(\hat{f} - f_k + \hat{\alpha} \mathbf{g}_k^T \mathbf{g}_k)} \mathbf{g}_k$$

we can use for example: $\hat{\alpha} = \alpha_{k-1}$.

Steepest descent: with fixed step size

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{g}_k$$

with $\alpha_k = \alpha$ where $\alpha > 0$ is a small value provided by the user, for example, $\alpha = 1e-2$

Global Convergence

- 1 An iterative algorithm is *globally convergent* if for any arbitrary starting point the algorithm is guaranteed to generate a sequence of points converging to a point that satisfies the *first order necessary condition* for a minimizer.
- 2 When the algorithm is not globally convergent, it may still generate a sequence that converges to a point satisfying the *first order necessary condition*, if the initial point is sufficiently close to the point. In this case, we say that the algorithm is *locally convergent*.
- 3 Another issue of interest, for both locally or globally convergent algorithms, is the rate of convergence; i.e., how fast the algorithm converges to a solution point.

Quadratic case

Lets start by the Quadratic case.

Let $f(x) = \frac{1}{2}x^T Qx - b^T x$, with Q is a symmetric positive definite matrix.

Note that there is not loss of generality in considering Q to be a symmetric matrix, due to, if A is not symmetric

$$x^T A x = \frac{1}{2}(x^T A x + x^T A^T x) = x^T \frac{A + A^T}{2} x := x^T Q x$$

where $Q = \frac{A + A^T}{2}$ is symmetric!.

Quadratic case

The convergence analysis is more convenient if we consider the following function

$$\begin{aligned} E(\mathbf{x}) &= f(\mathbf{x}) + \frac{1}{2}(\mathbf{x}^*)^T \mathbf{Q} \mathbf{x}^* \\ &= \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T \mathbf{Q} (\mathbf{x} - \mathbf{x}^*) \\ &= \frac{1}{2} \|\mathbf{x} - \mathbf{x}^*\|_{\mathbf{Q}}^2 \end{aligned}$$

which differs from $f(\mathbf{x})$ in the constant $\frac{1}{2}(\mathbf{x}^*)^T \mathbf{Q} \mathbf{x}^* = \frac{1}{2} \|\mathbf{x}^*\|_{\mathbf{Q}}^2$, ie

$$E(\mathbf{x}) = f(\mathbf{x}) + \frac{1}{2} \|\mathbf{x}^*\|_{\mathbf{Q}}^2$$

Note: $\|\mathbf{x}\|_{\mathbf{Q}}^2 = \mathbf{x}^T \mathbf{Q} \mathbf{x}$ is the weighted norm.

Quadratic case

Lemma 2.1

The iterates $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{g}_k$ with $\mathbf{g}_k = \mathbf{Q}\mathbf{x}_k - \mathbf{b}$ satisfies that if $\mathbf{g}_k \neq \mathbf{0}$ then $E(\mathbf{x}_{k+1}) = (1 - \gamma_k)E(\mathbf{x}_k)$ where

$$\gamma_k = \alpha_k \frac{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} \left(2 \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} - \alpha_k \right)$$

if additionally $\alpha_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}$ then

$$\gamma_k = \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k \mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k}$$

Note: if $\mathbf{g}_k = \mathbf{0}$ we consider $\gamma_k = 1$

Quadratic case

$$\begin{aligned} E(\mathbf{x}_{k+1}) &= \frac{1}{2}(\mathbf{x}_k - \mathbf{x}^* - \alpha_k \mathbf{g}_k)^T \mathbf{Q}(\mathbf{x}_k - \mathbf{x}^* - \alpha_k \mathbf{g}_k) \\ &= \frac{1}{2}(\mathbf{x}_k - \mathbf{x}^*)^T \mathbf{Q}(\mathbf{x}_k - \mathbf{x}^*) + \frac{1}{2} \alpha_k^2 \mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k \\ &\quad - \alpha_k (\mathbf{x}_k - \mathbf{x}^*)^T \mathbf{Q} \mathbf{g}_k \end{aligned}$$

Quadratic case

On the other hand $\mathbf{g}_k = \mathbf{Q}\mathbf{x}_k - \mathbf{b} = \mathbf{Q}(\mathbf{x}_k - \mathbf{x}^*)$ then

$$\mathbf{x}_k - \mathbf{x}^* = \mathbf{Q}^{-1}\mathbf{g}_k$$

and

$$E(\mathbf{x}_k) = \frac{1}{2}(\mathbf{x}_k - \mathbf{x}^*)^T \mathbf{Q}(\mathbf{x}_k - \mathbf{x}^*) = \frac{1}{2}\mathbf{g}_k^T \mathbf{Q}^{-1}\mathbf{g}_k$$

$$\begin{aligned} E(\mathbf{x}_{k+1}) &= \frac{1}{2}(\mathbf{x}_k - \mathbf{x}^*)^T \mathbf{Q}(\mathbf{x}_k - \mathbf{x}^*) + \frac{1}{2}\alpha_k^2 \mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k - \alpha_k \mathbf{g}_k^T \mathbf{g}_k \\ &= E(\mathbf{x}_k) + \left(\frac{1}{2}\alpha_k^2 \mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k - \alpha_k \mathbf{g}_k^T \mathbf{g}_k \right) \frac{E(\mathbf{x}_k)}{E(\mathbf{x}_k)} \\ &= [1 - \gamma_k] E(\mathbf{x}_k) \end{aligned}$$

Quadratic case

$$\begin{aligned} E(\mathbf{x}_{k+1}) &= (1 - \gamma_k)E(\mathbf{x}_k) \\ \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_{\mathbf{Q}}^2 &= (1 - \gamma_k)\|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2 \end{aligned}$$

with

$$\begin{aligned} \gamma_k &= -\frac{\frac{1}{2}\alpha_k^2 \mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k - \alpha_k \mathbf{g}_k^T \mathbf{g}_k}{\frac{1}{2} \mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} \\ &= \alpha_k \frac{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} \left(2 \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} - \alpha_k \right) \end{aligned}$$

Quadratic case

$$\text{If } \alpha_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}$$

$$\begin{aligned} \gamma_k &= \alpha_k \frac{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} \left(2 \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} - \alpha_k \right) \\ &= \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} \frac{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} \left(2 \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} - \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} \right) \\ &= \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k \mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} \end{aligned}$$

Quadratic case

Remark 2.2

Note that $\gamma_k = 1 - \frac{E(\mathbf{x}_{k+1})}{E(\mathbf{x}_k)} \geq 0$ due to $E(\mathbf{x}_{k+1}) \leq E(\mathbf{x}_k)$ and $\gamma_k \leq 1$ due to $E(\mathbf{x}) \geq 0$, then $0 \leq \gamma_k \leq 1$.

Quadratic case

Theorem 2.3

Let $\{x_k\}$ be the sequence resulting from a gradient algorithm $x_{k+1} = x_k - \alpha_k g_k$. Let γ_k be as defined in the previous Lemma, and suppose that $\gamma_k > 0$ for all k . Then, $\{x_k\}$ converges to x^ for any initial condition x_0 if and only if*

$$\sum_{k=1}^{\infty} \gamma_k = \infty$$

Quadratic case

From $E(\mathbf{x}_{k+1}) = (1 - \gamma_k)E(\mathbf{x}_k)$ we obtain

$$E(\mathbf{x}_k) = \prod_{i=0}^{k-1} (1 - \gamma_i) E(\mathbf{x}_0)$$

- 1 If $\gamma_i = 1$ the result is trivial.
- 2 Assume $\gamma_i < 1$ Then, $\mathbf{x}_k \rightarrow \mathbf{x}^*$ if and only if $E(\mathbf{x}_k) \rightarrow 0$.

Therefore $\prod_{i=0}^{\infty} (1 - \gamma_i) \rightarrow 0$, iff $-\sum_{i=0}^{\infty} \log(1 - \gamma_i) = \infty$

Quadratic case

It remains to proof that $\sum_{i=0}^{\infty} -\log(1 - \gamma_i) = \infty$ iff $\sum_{i=0}^{\infty} \gamma_i = \infty$.

(\Leftarrow) if $\sum_{i=0}^{\infty} \gamma_i = \infty$ and taking into account that

$$\log(x) \leq x - 1$$

then

$$-\log(1 - x) \geq x$$

and therefore

$$\sum_{i=0}^{\infty} -\log(1 - \gamma_i) \geq \sum_{i=0}^{\infty} \gamma_i = \infty$$

.

Quadratic case

(\Rightarrow) (by contradiction) Suppose $\sum_{i=0}^{\infty} \gamma_i < \infty$ therefore $\gamma_i \rightarrow 0$, ie

$$1 - \gamma_i \approx 1$$

for all $i \geq j$, for sufficiently large j .

As $\log(x) \geq 2(x - 1)$ for x close to 1 then $-\log(1 - x) \leq 2x$.
Therefore

$$\sum_{i=j}^{\infty} -\log(1 - \gamma_i) \leq 2 \sum_{i=j}^{\infty} \gamma_i < \infty$$

then

$$\sum_{i=0}^{\infty} -\log(1 - \gamma_i) < \infty$$

that contradicts the hypothesis that $\sum_{i=0}^{\infty} -\log(1 - \gamma_i) = \infty$.

Quadratic case

Lemma 2.4

Let \mathbf{Q} be an $n \times n$ real symmetric positive definite matrix. Then, for any $\mathbf{x} \in \mathbb{R}^n$, we have

$$\frac{a}{A} \leq \frac{(\mathbf{x}^T \mathbf{x})^2}{\mathbf{x}^T \mathbf{Q} \mathbf{x} \mathbf{x}^T \mathbf{Q}^{-1} \mathbf{x}} \leq \frac{A}{a}$$

where a and A are, respectively, the smallest and largest eigenvalues of \mathbf{Q} .

Quadratic case

Proof.

Applying Rayleigh's inequality, we get

$$a \leq \frac{\mathbf{x}^T \mathbf{Q} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \leq A, \text{ and } \frac{1}{A} \leq \frac{\mathbf{x}^T \mathbf{Q}^{-1} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \leq \frac{1}{a}$$

therefore $\frac{a}{A} \leq \frac{(\mathbf{x}^T \mathbf{x})^2}{\mathbf{x}^T \mathbf{Q} \mathbf{x} \mathbf{x}^T \mathbf{Q}^{-1} \mathbf{x}} \leq \frac{A}{a}.$



Proposition 2.5

Kantorovich inequality: Let \mathbf{Q} be a positive definite symmetric $n \times n$ matrix. For any vector \mathbf{x} there holds

$$\frac{(\mathbf{x}^T \mathbf{x})^2}{\mathbf{x}^T \mathbf{Q} \mathbf{x} \mathbf{x}^T \mathbf{Q}^{-1} \mathbf{x}} \geq \frac{4aA}{(a + A)^2}$$

where a and A are, respectively, the smallest and largest eigenvalues of \mathbf{Q} .

Remark 2.6

Note that $\frac{4aA}{(a+A)^2} \geq \frac{a}{A}$ therefore

$$\frac{a}{A} \leq \frac{4aA}{(a + A)^2} \leq \frac{(\mathbf{x}^T \mathbf{x})^2}{\mathbf{x}^T \mathbf{Q} \mathbf{x} \mathbf{x}^T \mathbf{Q}^{-1} \mathbf{x}} \leq \frac{A}{a}$$

Quadratic case: Steepest decent with exact line search

Theorem 2.7

In the steepest decent algorithm $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{g}_k$, with exact line search, i.e. $\alpha_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}$, we have that $\mathbf{x}_k \rightarrow \mathbf{x}^$ for any \mathbf{x}_0 .*

Proof.

As $\gamma_k = \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k \mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} \geq \frac{a}{A} > 0$. Therefore $\sum_k^\infty \gamma_k = \infty$ and using Theorem 2.3, the sequence $\{\mathbf{x}_k\}$ converges to \mathbf{x}^* for any initial condition \mathbf{x}_0 □

Quadratic case: Gradient method with fixed step size

Theorem 2.8

In the steepest decent algorithm $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \mathbf{g}_k$, with fixed step size, i.e. $\alpha_k = \alpha$ for all k , we have that $\mathbf{x}_k \rightarrow \mathbf{x}^$ for any \mathbf{x}_0 iff $0 < \alpha < \frac{2}{A}$.*

Proof.

(\Leftarrow) By Rayleigh's inequality

$$a \mathbf{g}_k^T \mathbf{g}_k \leq \mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k \leq A \mathbf{g}_k^T \mathbf{g}_k$$

$$\mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k \leq \frac{1}{a} \mathbf{g}_k^T \mathbf{g}_k$$

As $\gamma_k = \alpha \frac{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} \left(2 \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} - \alpha \right) \geq a \alpha \left(\frac{2}{A} - \alpha \right) > 0$ then $\sum_k \gamma_k = \infty$ and using Theorem 2.3, the sequence $\mathbf{x}_k \rightarrow \mathbf{x}^*$ □

Gradient method with fixed step size

Proof.

(\Rightarrow) (By contradiction) Suppose that $\alpha \leq 0$ or $\alpha \geq \frac{2}{A}$. On the other hand, we select \mathbf{x}_0 such that $\mathbf{x}_0 - \mathbf{x}^*$ is the eigenvector of \mathbf{Q} that corresponds to A .

$$\begin{aligned}\mathbf{x}_{k+1} - \mathbf{x}^* &= \mathbf{x}_k - \alpha \mathbf{g}_k - \mathbf{x}^* = \mathbf{x}_k - \alpha(\mathbf{Q}\mathbf{x}_k - \mathbf{b}) - \mathbf{x}^* \\ &= \mathbf{x}_k - \mathbf{x}^* - \alpha(\mathbf{Q}\mathbf{x}_k - \mathbf{Q}\mathbf{x}^*) = (\mathbf{I} - \alpha\mathbf{Q})(\mathbf{x}_k - \mathbf{x}^*) \\ &= (\mathbf{I} - \alpha\mathbf{Q})^{k+1}(\mathbf{x}_0 - \mathbf{x}^*) = (1 - \alpha A)^{k+1}(\mathbf{x}_0 - \mathbf{x}^*)\end{aligned}$$

then

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| = |1 - \alpha A|^{k+1} \|\mathbf{x}_0 - \mathbf{x}^*\|$$

For $\alpha \leq 0$ or $\alpha \geq \frac{2}{A}$ we have $|1 - \alpha A| \geq 1$ therefore $\|\mathbf{x}_{k+1} - \mathbf{x}^*\|$ does not converge to 0 and $\{\mathbf{x}_k\}$ does not converge to \mathbf{x}^* .

Quadratic case: Steepest decent with exact line search

Theorem 3.1

In the steepest decent algorithm with exact line search, i.e.

$\alpha_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}$, applied to the quadratic function we have

$$E(\mathbf{x}_{k+1}) \leq \left(1 - \frac{a}{A}\right) E(\mathbf{x}_k)$$

Proof.

As $\gamma_k = \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k \mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} \geq \frac{a}{A}$ then $1 - \gamma_k \leq 1 - \frac{a}{A}$. Therefore

$$E(\mathbf{x}_{k+1}) = (1 - \gamma_k) E(\mathbf{x}_k) \leq \left(1 - \frac{a}{A}\right) E(\mathbf{x}_k)$$



Quadratic case: Summary

- The previous theorem is very important for the *convergence* of the steepest decent algorithm
- The ratio $\kappa := \frac{A}{a} = \|\mathbf{Q}\|_2 \|\mathbf{Q}^{-1}\|_2 = \kappa(\mathbf{Q}) \geq 1$ is the so-called *condition number* of \mathbf{Q}
- Recall: $\|\mathbf{Q}\|_2 = \sqrt{\lambda_M(\mathbf{Q}^T \mathbf{Q})} = \sigma_M(\mathbf{Q})$. If $\mathbf{Q} \succ 0$ is symmetric then $\sigma_M(\mathbf{Q}) = \lambda_M(\mathbf{Q})$
- The term $1 - \frac{a}{A} = 1 - \frac{1}{\kappa}$ plays an important role in the convergence of the sequence $\{E(\mathbf{x}_k)\}$ (and therefore of the convergence of \mathbf{x}_k to \mathbf{x}^*).

Quadratic case: Summary

- The method of steepest descent converges linearly with a ratio **no greater than $1 - \frac{1}{\kappa}$** .
- The smaller the value of κ , the smaller the relative value of $E(\mathbf{x}_{k+1})$ with respect to $E(\mathbf{x}_k)$ and therefore $\{E(\mathbf{x}_k)\}$ converges faster to 0.
- If $\kappa = 1$, i.e., the level sets are circulars and $A = a$, the algorithm converges in one iteration to the minimizer. If κ increases then the rate of convergence decreases.

Quadratic case: Steepest descent with exact line search

Lemma 3.2

In the steepest descent algorithm with exact line search, i.e.

$\alpha_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}$, if $\mathbf{g}_k \neq \mathbf{0}$ then $\gamma_k = 1$ iff \mathbf{g}_k is an eigenvector of \mathbf{Q}

Proof.

(\Leftarrow) if \mathbf{g}_k is an eigenvector of \mathbf{Q} then $\mathbf{Q} \mathbf{g}_k = \lambda \mathbf{g}_k$ and $\mathbf{Q}^{-1} \mathbf{g}_k = \lambda^{-1} \mathbf{g}_k$ therefore

$$\gamma_k = \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k \mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} = 1$$



Quadratic case: Steepest descent with exact line search

Proof.

(\Rightarrow) if $\gamma_k = 1$ then $E(\mathbf{x}_{k+1}) = \frac{1}{2}(\mathbf{x}_{k+1} - \mathbf{x}^*)^T \mathbf{Q}(\mathbf{x}_{k+1} - \mathbf{x}^*) = 0$ therefore $\mathbf{x}_{k+1} = \mathbf{x}^*$. Hence,

$$\begin{aligned}\mathbf{x}_{k+1} &= \mathbf{x}_k - \alpha_k \mathbf{g}_k \\ \mathbf{x}^* &= \mathbf{x}_k - \alpha_k \mathbf{g}_k \\ \mathbf{Q}\mathbf{x}^* &= \mathbf{Q}\mathbf{x}_k - \alpha_k \mathbf{Q}\mathbf{g}_k \\ \alpha_k \mathbf{Q}\mathbf{g}_k &= \mathbf{Q}\mathbf{x}_k - \mathbf{b} \\ \mathbf{Q}\mathbf{g}_k &= \frac{1}{\alpha_k} \mathbf{g}_k\end{aligned}$$

and then \mathbf{g}_k is an eigenvector of \mathbf{Q} .



Quadratic case: Steepest descent with exact line search

Theorem 3.3

In the steepest decent algorithm with exact line search, the error norm $E(\cdot)$ satisfies $E(\mathbf{x}_{k+1}) \leq \left(\frac{A-a}{A+a}\right)^2 E(\mathbf{x}_k)$

Proof.

Using the Kantorovich inequality: $\gamma_k = \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k \mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k} \geq \frac{4aA}{(a+A)^2}$
then $1 - \gamma_k \leq 1 - \frac{4aA}{(a+A)^2} = \left(\frac{A-a}{A+a}\right)^2$. Therefore

$$E(\mathbf{x}_{k+1}) = (1 - \gamma_k) E(\mathbf{x}_k) \leq \left(\frac{A-a}{A+a}\right)^2 E(\mathbf{x}_k)$$

Note also that: $\left(\frac{A-a}{A+a}\right)^2 = \left(\frac{\kappa-1}{\kappa+1}\right)^2$



Non-quadratic case

Theorem 3.4

Non-quadratic case Suppose f is defined on \mathbb{R}^n , has continuous second partial derivatives, and has a relative minimum at \mathbf{x}^ . Suppose further that the Hessian matrix of f , $\mathbf{H}(\mathbf{x}^*)$, has smallest eigenvalue $a > 0$ and largest eigenvalue $A > 0$. If $\{\mathbf{x}_k\}$ is a sequence generated by the method of steepest descent that converges to \mathbf{x}^* , then the sequence of objective values $\{f(\mathbf{x}_k)\}$ converges to $f(\mathbf{x}^*)$ linearly with a convergence ratio no greater than $\left(\frac{A-a}{A+a}\right)^2$, i.e., for all k sufficiently large, we have*

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(\frac{A-a}{A+a}\right)^2 [f(\mathbf{x}_k) - f(\mathbf{x}^*)]$$