

Optimization

Hassan OMRAN

Lecture 3: Multi-Dimensional Search Methods - Part I

Télécom Physique Strasbourg
Université de Strasbourg



Outline of the talk

1. Gradient methods

2. Newton's method

1. Gradient methods

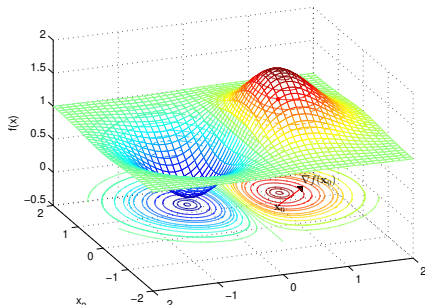
2. Newton's method

Gradient methods

Consider the continuously differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. We study the following problem:

$$\begin{array}{ll} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in \mathbb{R}^n \end{array}$$

- ◇ Remember that the gradient $\nabla f(\mathbf{x}_0)$ indicate the direction of the maximum increase of $f(\cdot)$ at \mathbf{x}_0
- ◇ The vector $-\nabla f(\mathbf{x}_0)$ indicates the direction of the **maximum decrease** of $f(\cdot)$ at \mathbf{x}_0



Gradient methods

Starting from a point \mathbf{x}_0 , we look for a new point \mathbf{x}_1 in the direction $-\nabla f(\mathbf{x}_0)$. It is easy to show that

$$f(\mathbf{x}_0 - \alpha \nabla f(\mathbf{x}_0)) = f(\mathbf{x}_0) - \alpha \|\nabla f(\mathbf{x}_0)\|^2 + o(\alpha)$$

Clearly when $\nabla f(\mathbf{x}_0) \neq \mathbf{0}$ and for sufficiently small $\alpha > 0$ we have

$$f(\mathbf{x}_0 - \alpha \nabla f(\mathbf{x}_0)) < f(\mathbf{x}_0)$$

Algorithm:

Starting from a point \mathbf{x}_k we calculate the gradient $\nabla f(\mathbf{x}_k)$, and for a suitable step size α_k we define the next point

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \quad (1)$$

The question is how to choose the step size α_k ?

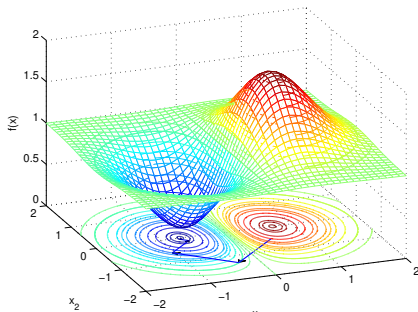
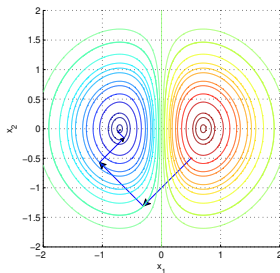
- ◇ Gradient methods with a fixed step size $\alpha_k = \alpha > 0$
- ◇ Steepest descent

Steepest descent

The steepest descent is a gradient method where the step α_k is chosen in order to have the maximum amount of decrease of the objective function $f(\cdot)$ at each iteration. That is, at each iteration we define the function $h_k(\alpha) = f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k))$ and we solve

$$\alpha_k = \arg \min_{\alpha \geq 0} h_k(\alpha) = \arg \min_{\alpha \geq 0} f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)) \quad (2)$$

Note that (2) is a one-dimensional optimization problem.



Sequence resulting from the steepest descent method.

Steepest descent: orthogonal generated directions

The next theorem shows that the directions generated by the steepest descent method are orthogonal.

Theorem 1

Let $\{\mathbf{x}_k\}_{k=0}^{\infty}$ be the sequence from the steepest descent for a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, then $\forall k \geq 0$ the vector $\mathbf{x}_{k+2} - \mathbf{x}_{k+1}$ is orthogonal to the vector $\mathbf{x}_{k+1} - \mathbf{x}_k$.

Proof.

Consider the function

$h_k : \mathbb{R} \rightarrow \mathbb{R}$ defined by $h_k(\alpha) = f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k))$. From the FONC and the chain rule

$$0 = \dot{h}_k(\alpha)|_{\alpha=\alpha_k} = (\nabla f(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)))^T (-\nabla f(\mathbf{x}_k)) = -\langle \nabla f(\mathbf{x}_{k+1}), \nabla f(\mathbf{x}_k) \rangle \quad (3)$$

Finally, since $\mathbf{x}_{k+2} = \mathbf{x}_{k+1} - \alpha_{k+1} \nabla f(\mathbf{x}_{k+1})$ and $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$

$$\langle \mathbf{x}_{k+2} - \mathbf{x}_{k+1}, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle = \alpha_{k+1} \alpha_k \langle \nabla f(\mathbf{x}_{k+1}), \nabla f(\mathbf{x}_k) \rangle \quad (4)$$

From (3) and (4) we have that $\langle \mathbf{x}_{k+2} - \mathbf{x}_{k+1}, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle = 0$.



Steepest descent: descent property

At each iteration, if $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$, then the value of the function will strictly decrease

Theorem 2

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

Proof

Consider the functions $h_k : \mathbb{R} \rightarrow \mathbb{R}$ defined by $h_k(\alpha) = f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k))$. Remember from (2) that

$$h_k(\alpha_k) \leq h_k(\alpha), \quad \forall \alpha \geq 0 \quad (5)$$

Also

$$\dot{h}_k(0) = \left(\nabla f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)) \right)^T \left(-\nabla f(\mathbf{x}_k) \right) \Big|_{\alpha=0} = -\|\nabla f(\mathbf{x}_k)\|^2 < 0$$

Thus, from Taylor theorem we have

$$h_k(\alpha) = h_k(0) + \dot{h}_k(0)\alpha + o(\alpha)$$

Steepest descent: descent property

Proof (Cont.)

then for sufficiently small $\bar{\alpha} > 0$ we have

$$h_k(\bar{\alpha}) < h_k(0) \quad (6)$$

From (5) and (6)

$$\begin{aligned} h_k(\alpha_k) &\leq h_k(\bar{\alpha}) < h_k(0) \\ f(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)) &< f(\mathbf{x}_k) \\ f(\mathbf{x}_{k+1}) &< f(\mathbf{x}_k) \end{aligned}$$



Steepest descent: stopping criteria

If at a certain iteration we have that $\nabla f(\mathbf{x}_k) = \mathbf{0}$, then the point satisfies the FONC, and the algorithm can be stopped.

Practically, it is not always possible to obtain $\nabla f(\mathbf{x}_k) = \mathbf{0}$. One of the following criteria can be used as a stopping condition (with a small $\epsilon > 0$):

$$\diamond \quad \|\nabla f(\mathbf{x}_k)\| \leq \epsilon$$

$$\diamond \quad |f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)| < \epsilon$$

$$\diamond \quad \|\mathbf{x}_{k+1} - \mathbf{x}_k\| < \epsilon$$

$$\diamond \quad \frac{|f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)|}{|f(\mathbf{x}_k)|} < \epsilon$$

$$\diamond \quad \frac{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|}{\|\mathbf{x}_k\|} < \epsilon$$

The last two criteria are scale-independent. They are to be used with caution as dividing by very small numbers may occur.

Steepest descent: the case of quadratic function

The case of quadratic function

Consider the following problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} + \mathbf{q}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{x} \in \mathbb{R}^n \end{aligned} \tag{7}$$

for a given matrix $0 < Q = Q^T \in \mathbb{R}^{n \times n}$ and a vector $\mathbf{q} \in \mathbb{R}^n$. Note that $\nabla f(\mathbf{x}) = Q\mathbf{x} + \mathbf{q}$ and $D^2 f(\mathbf{x}) = Q > 0$.

The case of quadratic function

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{q}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{x} \in \mathbb{R}^n \end{aligned} \quad (7)$$

From SOSC with $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $D^2f(\mathbf{x}) = \mathbf{Q} > 0$ we find that the solution is $\mathbf{x}^* = -\mathbf{Q}^{-1}\mathbf{q}$. In this case, we will show that we can provide the expression of α_k .

$$\begin{aligned}\alpha_k &= \arg \min_{\alpha \geq 0} h_k(\alpha) = \arg \min_{\alpha \geq 0} f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)) \\ &= \arg \min_{\alpha \geq 0} \frac{1}{2} (\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k))^T Q (\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)) + \mathbf{q}^T (\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k))\end{aligned}$$

From FONC $\dot{h}_k(\alpha)|_{\alpha=\alpha_k} = 0$ we find that

Steepest descent: the case of quadratic function

$$\begin{aligned}
 0 &= \dot{h}_k(\alpha) \Big|_{\alpha=\alpha_k} = \frac{d}{d\alpha} f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)) \Big|_{\alpha=\alpha_k} \\
 &= \left(\nabla f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)) \right)^T \left(-\nabla f(\mathbf{x}_k) \right) \Big|_{\alpha=\alpha_k} \\
 &= \left(Q(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)) + \mathbf{q} \right)^T \left(-\nabla f(\mathbf{x}_k) \right) \Big|_{\alpha=\alpha_k} \\
 &= \left((Q\mathbf{x}_k + \mathbf{q}) - \alpha_k Q \nabla f(\mathbf{x}_k) \right)^T \left(-\nabla f(\mathbf{x}_k) \right)
 \end{aligned}$$

we suppose that $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$, since if it is the case, then $\mathbf{x}_k = \mathbf{x}^*$ and the algorithm stops

$$\Rightarrow \alpha_k = \frac{(Q\mathbf{x}_k + \mathbf{q})^T \nabla f(\mathbf{x}_k)}{\nabla f(\mathbf{x}_k)^T Q \nabla f(\mathbf{x}_k)} = \frac{\nabla f(\mathbf{x}_k)^T \nabla f(\mathbf{x}_k)}{\nabla f(\mathbf{x}_k)^T Q \nabla f(\mathbf{x}_k)} \quad (8)$$

finally

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)^T \nabla f(\mathbf{x}_k)}{\nabla f(\mathbf{x}_k)^T Q \nabla f(\mathbf{x}_k)} \nabla f(\mathbf{x}_k) \quad (9)$$

In the case of quadratic function, it is shown that the steepest descent always converges.

The descent method in the case of quadratic function in (7) converges for any initial condition. That is $\mathbf{x}_k \rightarrow \mathbf{x}^* \forall \mathbf{x}_0 \in \mathbb{R}^n$.

Also, in this case the convergence rate is linear.

The descent method in the case of quadratic function in (7) has a linear convergence

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \|\mathbf{x}_k - \mathbf{x}^*\| \sqrt{\frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)} - 1}$$

the ratio $\frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)} \geq 1$ is called the **condition number** of Q . The smaller this number, the faster \mathbf{x}_k converges to \mathbf{x}^* .

Steepest descent: the case of quadratic function

Proof

Consider the error $e_k = \mathbf{x}_k - \mathbf{x}^*$, and the function $V(e_k) = e_k^T Q e_k$. For simplicity, we will use the notation \mathbf{g}_k for $\nabla f(\mathbf{x}_k)$.

$$\begin{aligned} V(e_{k+1}) &= (\mathbf{x}_{k+1} - \mathbf{x}^*)^T Q (\mathbf{x}_{k+1} - \mathbf{x}^*) \\ &= (\mathbf{x}_k - \mathbf{x}^* - \alpha_k \mathbf{g}_k)^T Q (\mathbf{x}_k - \mathbf{x}^* - \alpha_k \mathbf{g}_k) \\ &= (\mathbf{x}_k - \mathbf{x}^*)^T Q (\mathbf{x}_k - \mathbf{x}^*) - \alpha_k \mathbf{g}_k^T Q (\mathbf{x}_k - \mathbf{x}^*) - \alpha_k (\mathbf{x}_k - \mathbf{x}^*)^T Q \mathbf{g}_k + \alpha_k^2 \mathbf{g}_k^T Q \mathbf{g}_k \\ &= V(e_k) - 2\alpha_k \underbrace{\mathbf{g}_k^T Q (\mathbf{x}_k - \mathbf{x}^*)}_{Q\mathbf{x}_k + \mathbf{q} = \mathbf{g}_k} + \alpha_k^2 \mathbf{g}_k^T Q \mathbf{g}_k \end{aligned}$$

From (8) we have

$$\begin{aligned} V(e_{k+1}) &= V(e_k) - 2\left(\frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T Q \mathbf{g}_k}\right) \mathbf{g}_k^T \mathbf{g}_k + \left(\frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T Q \mathbf{g}_k}\right)^2 \mathbf{g}_k^T Q \mathbf{g}_k \\ &= V(e_k) - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T Q \mathbf{g}_k} \end{aligned}$$

Steepest descent: the case of quadratic function

Proof (Cont.)

$$\frac{V(e_{k+1})}{V(e_k)} = 1 - \frac{(g_k^T g_k)^2}{g_k^T Q g_k} \frac{1}{V(e_k)}$$

Note that

$$\begin{aligned} V(e_k) &= e_k^T Q e_k \\ &= (\mathbf{x}_k - \mathbf{x}^*)^T Q (\mathbf{x}_k - \mathbf{x}^*) \\ &= (Q \mathbf{x}_k - Q \mathbf{x}^*)^T Q^{-1} (Q \mathbf{x}_k - Q \mathbf{x}^*) \\ &= g_k^T Q^{-1} g_k \end{aligned}$$

Thus

$$\frac{V(e_{k+1})}{V(e_k)} = 1 - \frac{(g_k^T g_k)^2}{(g_k^T Q g_k)(g_k^T Q^{-1} g_k)}$$

Steepest descent: the case of quadratic function

Proof (Cont.)

From Rayleigh's inequality we have

$$\begin{aligned} \frac{V(e_{k+1})}{V(e_k)} &= 1 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{(\mathbf{g}_k^T Q \mathbf{g}_k)(\mathbf{g}_k^T Q^{-1} \mathbf{g}_k)} \leq 1 - \frac{1}{\lambda_{\max}(Q)\lambda_{\max}(Q^{-1})} \\ \frac{\lambda_{\min}(Q)\|\mathbf{e}_{k+1}\|^2}{\lambda_{\max}(Q)\|\mathbf{e}_k\|^2} &\leq \frac{V(e_{k+1})}{V(e_k)} \leq 1 - \frac{\lambda_{\min}(Q)}{\lambda_{\max}(Q)} \end{aligned}$$

finally

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2}{\|\mathbf{x}_k - \mathbf{x}^*\|^2} \leq \frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)} - 1$$



Steepest descent: the case of quadratic function

Examples

Apply the steepest descent method to the problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} + \mathbf{q}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{x} \in \mathbb{R}^n \end{aligned}$$

in the following cases

◇

$$Q = \lambda I, \lambda > 0, \quad \forall \mathbf{q} \in \mathbb{R}^2 \quad \forall \mathbf{x}_0 \in \mathbb{R}^2$$

◇

$$Q = \begin{bmatrix} 10 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{q} = \begin{bmatrix} -3 \\ -3 \end{bmatrix}, \quad \mathbf{x}_0 = \begin{bmatrix} -2 \\ -7 \end{bmatrix}$$

◇

$$Q = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{q} = \begin{bmatrix} -3 \\ -3 \end{bmatrix}, \quad \mathbf{x}_0 = \begin{bmatrix} -2 \\ -7 \end{bmatrix}$$

Steepest descent: remarks

- ◇ *Gradient methods are simple, and have (often) good performance*
- ◇ *They require the calculation of the gradient, which may not be always possible*
- ◇ *They have good convergence properties for convex QP problems*
- ◇ *They may have bad performance for some functions*

We will see how, for the case of Rosenbock's function, they have a slow zig-zag convergence near the optimal point, and require a big number of function evaluations.

1. Gradient methods

2. Newton's method

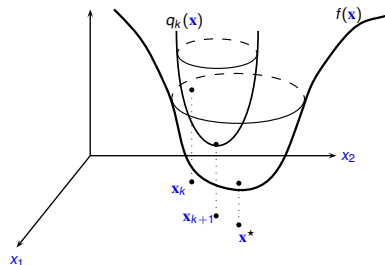
Newton's method

Consider the **twice** continuously differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, we study the following problem:

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{x} \in \mathbb{R}^n \end{aligned} \tag{10}$$

The idea is to minimize a second-order approximation of $f(\cdot)$ at each iteration

- ◇ Start with an initial value \mathbf{x}_0
- ◇ At each iteration minimize second-order approximation of $f(\cdot)$ at \mathbf{x}_k .
- ◇ Use the minimizer of the quadratic approximate as a starting point for the next iteration



The approximation of the function $f(\cdot)$ by a quadratic function $q_k(\cdot)$.

Newton's method

At iteration k , using Taylor series expansion of $f(\cdot)$ at \mathbf{x}_k , we find the approximate $q_k(\cdot)$

$$q_k(\mathbf{x}) = f(\mathbf{x}_k) + (\mathbf{x} - \mathbf{x}_k)^T \nabla f(\mathbf{x}_k) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_k)^T D^2 f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) \quad (11)$$

Note that

$$\nabla q_k(\mathbf{x}) = \nabla f(\mathbf{x}_k) + D^2 f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k), \quad D^2 q_k(\mathbf{x}) = D^2 f(\mathbf{x}_k)$$

If $D^2 f(\mathbf{x}_k) > 0$ then $q_k(\cdot)$ achieves a minimum for $\nabla q_k(\mathbf{x}_{k+1}) = \mathbf{0}$, that is

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (D^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k) \quad (12)$$

Denoting $A = D^2 f(\mathbf{x}_k)$ and $b = \nabla f(\mathbf{x}_k)$, the recursive formula is composed of two steps:

- solving $Ay = b$
- $\mathbf{x}_{k+1} = \mathbf{x}_k - y$

Newton's method: convergence rate

The case of quadratic function

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} + \mathbf{q}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{x} \in \mathbb{R}^n \end{aligned} \tag{13}$$

For a given matrix $0 < Q = Q^T \in \mathbb{R}^{n \times n}$ and a vector $\mathbf{q} \in \mathbb{R}^n$. Note that $\nabla f(\mathbf{x}) = Q\mathbf{x} + \mathbf{q}$ and $D^2 f(\mathbf{x}^*) = Q > 0$.

Starting from any initial point \mathbf{x}_0 , we have

$$\begin{aligned} \mathbf{x}_1 &= \mathbf{x}_0 - Q^{-1}(Q\mathbf{x}_0 + \mathbf{q}) \\ &= -Q^{-1}\mathbf{q} \\ &= \mathbf{x}^* \end{aligned}$$

In this case, the Newton's method find the solution in one iteration, starting from any initial condition.

Newton's method: convergence rate

In the general case, the method has good convergence properties when starting close the solution.

Theorem 5

Consider the problem in (10) with three times continuously differentiable $f(\cdot)$, and a point \mathbf{x}^ with $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and invertible $D^2f(\mathbf{x}^*)$. Then, Newton's method is well defined $\forall k \geq 0$, and converges to \mathbf{x}^* if the starting point \mathbf{x}_0 is sufficiently close to \mathbf{x}^* .*

Moreover, it has a quadratic convergence

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^2} \leq c, \quad c > 0.$$

*Note that theorem above shows the **local** convergence property of the method.*

If the Hessian $D^2f(\mathbf{x}_k)$ is not positive definite, the method may not have the descent property.

Newton's method: convergence rate

The following theorem shows that the algorithm generates decreasing directions when $D^2f(\mathbf{x}_k) > 0$.

Theorem 6

If the Hessian $D^2f(\mathbf{x}_k) > 0$ and $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$, then the direction generated by Newton's method

$$\mathbf{d} = \mathbf{x}_{k+1} - \mathbf{x}_k = -(D^2f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k)$$

is a descent direction.

Proof

Consider the function $h_k(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d})$. Using Taylor's theorem we have

$$h_k(\alpha) = h_k(0) + \underbrace{\dot{h}_k(0)}_{\nabla f(\mathbf{x}_k)^T \mathbf{d} = \mathbf{d}^T \nabla f(\mathbf{x}_k) = -\mathbf{d}^T (D^2f(\mathbf{x}_k)) \mathbf{d}} \alpha + o(\alpha)$$

Newton's method: convergence rate

Proof (Cont.)

$$f(\mathbf{x}_k + \alpha \mathbf{d}) = f(\mathbf{x}_k) - \alpha \mathbf{d}^T (D^2 f(\mathbf{x}_k)) \mathbf{d} + o(\alpha)$$

Since $\mathbf{d} \neq \mathbf{0}$, and $D^2 f(\mathbf{x}_k) > 0$, then $\exists \bar{\alpha} > 0$ such that

$$f(\mathbf{x}_k + \alpha \mathbf{d}) < f(\mathbf{x}_k), \quad \forall \alpha \in (0, \bar{\alpha})$$

□

Newton's method: modifications

By Theorem 6 we see that instead of $\mathbf{x}_{k+1} = \mathbf{x}_k - (D^2f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k)$, we can introduce the following one-dimensional optimization

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k (D^2f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k)$$

with

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}_k - \alpha (D^2f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k)) \quad (14)$$

Note that by Theorem 6 the method in (14) guarantees the descent property $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ (when $D^2f(\mathbf{x}_k) > 0$ and $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$).

Finally, the **Levenberg–Marquardt** method is a modification of Newton's method which overcomes the problem caused by non definite positive Hessians

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k (\mu_k I + D^2f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k), \quad \text{with } \mu_k \geq 0$$

This can be seen as a combination of Newton's method when $\mu_k \rightarrow 0$, and gradient methods $\mu_k \rightarrow \infty$.

Newton's method: remarks

- ◇ The method is also called Newton-Raphson method
- ◇ Can also be seen as a method to find $\nabla f(\mathbf{x}) = \mathbf{0}$
- ◇ The quadratic approximate matches the first and second derivatives of $f(\cdot)$ at \mathbf{x}_k , thus more information are used than the gradient methods
- ◇ The method performs better than the steepest descent which uses only the first derivative (if the initial point is close the the minimizer)
- ◇ Drawbacks:
 - The evaluation of the Hessian can be computationally difficult
 - Requires solving a set of linear equations, which is of $O(n^3)$ complexity
 - Only local convergence is guaranteed
 - Modifications are needed if the Hessian is not definite positive