

## TMA947 / MMG621 — Nonlinear optimization

**Lecture 5 — Unconstrained optimization algorithms**

Emil Gustavsson, Zuzana Nedělková

November 6, 2017

[Minor revision: Axel Ringh - August, 2023]

Consider the unconstrained optimization problem to

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x), \quad (1)$$

where  $f \in C^0$  on  $\mathbb{R}^n$  ( $f$  is continuous). Mostly, we assume that  $f \in C^1$  holds ( $f$  is continuously differentiable), sometimes even  $C^2$ . The choice of the algorithm depends on the size of the problem, availability of  $\nabla f(x)$  and  $\nabla^2 f(x)$ , convexity of  $f$  and if the goal is to find a local or the global minimum.

Most algorithms for unconstrained optimization problems are what we call *line search type algorithms*.

**Definition.** *Line search type algorithm*

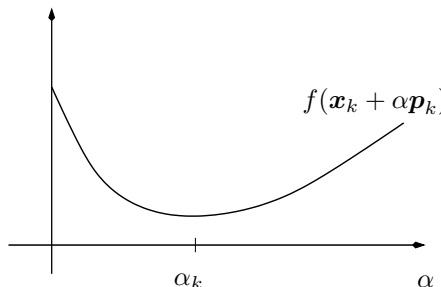
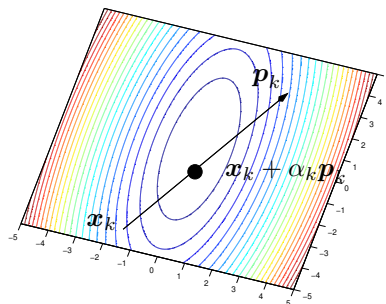
**Step 0:** **Starting point**  $x_0 \in \mathbb{R}^n$ . Let  $k := 0$ .

**Step 1:** Find **search direction**  $p_k \in \mathbb{R}^n$

**Step 2:** Perform **line search**, i.e., find  $\alpha_k > 0$  such that  $f(x_k + \alpha_k p_k) < f(x_k)$

**Step 3:** Let  $x_{k+1} := x_k + \alpha_k p_k$ .

**Step 4:** If **termination criteria** is fulfilled then stop! Otherwise, let  $k := k + 1$  and go to Step 1.



Most algorithms we consider are inherently *local*, meaning that the search direction  $p_k$  is only based on the information at the current point  $x_k$ , that is,  $f(x_k)$ ,  $\nabla f(x_k)$ , and  $\nabla^2 f(x_k)$ .

Think of a near-sighted mountain climber. The climber is in a deep fog and can only check his or her barometer for the height and feel the steepness of the slope under her feet.

## Step 1: Search directions

Vector  $\mathbf{p}_k$  is a descent direction at  $\mathbf{x}_k$  if  $f(\mathbf{x}_k + \alpha \mathbf{p}_k) < f(\mathbf{x}_k)$  for all  $\alpha \in (0, \delta]$  for some  $\delta > 0$ .

Let  $f \in C^1$  in some neighborhood of  $\mathbf{x}_k \in \mathbb{R}^n$ , if  $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$ , then  $\mathbf{p}_k = -\nabla f(\mathbf{x}_k)$  is a descent direction for  $f$  at  $\mathbf{x}_k$  (follows from optimality conditions). This search step is called *steepest descent direction* because it solves the problem to

$$\underset{\mathbf{p} \in \mathbb{R}^n: \|\mathbf{p}\|=1}{\text{minimize}} \quad \nabla f(\mathbf{x}_k)^T \mathbf{p}.$$

Let  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  be an arbitrary symmetric, positive definite matrix. Then  $\mathbf{p}_k = -\mathbf{Q} \nabla f(\mathbf{x}_k)$  is a descent direction for  $f$  at  $\mathbf{x}_k$ , because

$$\nabla f(\mathbf{x}_k)^T \mathbf{p}_k = -\nabla f(\mathbf{x}_k)^T \mathbf{Q} \nabla f(\mathbf{x}_k) < 0,$$

due to the positive definiteness of  $\mathbf{Q}$ .

Examples:

- *Steepest descent*:  $\mathbf{Q} = \mathbf{I}$ ,
- *Newton's method*:  $\mathbf{Q} = [\nabla^2 f(\mathbf{x}_k)]^{-1}$ .

We will now derive *Newton's method*. First assume that  $\nabla^2 f(\mathbf{x})$  is positive definite. A second-order Taylor approximation is then:

$$f(\mathbf{x}_k + \mathbf{p}) - f(\mathbf{x}_k) \approx \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p} =: \varphi_{\mathbf{x}_k}(\mathbf{p})$$

We now try to minimize this approximation by setting the gradient of  $\varphi_{\mathbf{x}_k}(\mathbf{p})$  to zero:

$$\nabla_{\mathbf{p}} \varphi_{\mathbf{x}_k}(\mathbf{p}) = \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k) \mathbf{p} = \mathbf{0} \quad \Leftrightarrow \quad \nabla^2 f(\mathbf{x}_k) \mathbf{p} = -\nabla f(\mathbf{x}_k)$$

Now by choosing the vector fulfilling this we obtain  $\mathbf{p}_k = -[\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$  as the search direction. When  $n = 1$ , we get that  $p_k = -f'(x_k)/f''(x_k)$ .

When the Hessian  $\nabla^2 f(\mathbf{x}_k)$  is positive definite this search direction is a descent direction. But when  $\nabla^2 f(\mathbf{x}_k)$  is negative definite (may be also non invertible), the search direction is an ascent direction, meaning that Newton's method does differentiate between minimization and maximization problem. The solution to this problem is to modify  $\nabla^2 f(\mathbf{x}_k)$  by adding a diagonal matrix  $\gamma \mathbf{I}$  such that  $(\nabla^2 f(\mathbf{x}_k) + \gamma \mathbf{I})$  is positive definite (this can always be done, why?). This method is called the *Levenberg-Marquardt modification*. We thus take as search direction

$$\mathbf{p}_k = -[\nabla^2 f(\mathbf{x}_k) + \gamma \mathbf{I}]^{-1} \nabla f(\mathbf{x}_k).$$

Note that

- *Steepest descent*:  $\gamma = \infty$ ,
- *Newton's method*:  $\gamma = 0$ .

What happens when we can not compute  $\nabla^2 f(\mathbf{x}_k)$ ? Try to approximate the Hessian in some way choosing approximate matrix  $\mathbf{B}_k$ . From Taylor expansion for  $\nabla f(\mathbf{x}_k)$  we have that

$$\nabla^2 f(\mathbf{x}_k)(\mathbf{x}_k - \mathbf{x}_{k-1}) \approx \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1})$$

so the approximate matrix  $\mathbf{B}_k$  has to fulfill

$$\mathbf{B}_k(\mathbf{x}_k - \mathbf{x}_{k-1}) = \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1}).$$

Many different choices of  $\mathbf{B}_k$  exist, and they lead to what is called *quasi-Newton* methods.

To summarize:

Steepest descent:	$\mathbf{p}_k = -\nabla f(\mathbf{x}_k)$
Newton's method:	$\nabla^2 f(\mathbf{x}_k)\mathbf{p}_k = -\nabla f(\mathbf{x}_k)$
Levenberg-Marquardt:	$(\nabla^2 f(\mathbf{x}_k) + \gamma \mathbf{I})\mathbf{p}_k = -\nabla f(\mathbf{x}_k)$
Quasi-Newton:	$\mathbf{B}_k\mathbf{p}_k = -\nabla f(\mathbf{x}_k)$ .

## Step 2: Line search

In each iteration one would like to solve

$$\underset{\alpha \geq 0}{\text{minimize}} \varphi(\alpha) := f(\mathbf{x}_k + \alpha \mathbf{p}_k).$$

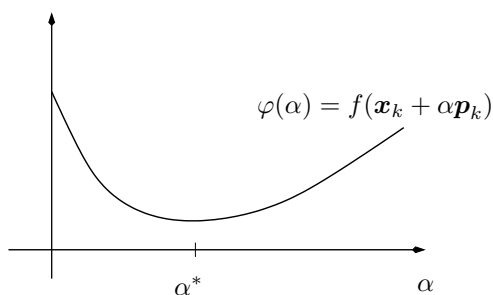
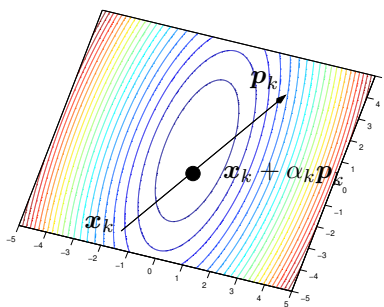
The optimality conditions for the problem are

$$\begin{aligned} \varphi'(\alpha^*) &\geq 0, \\ \alpha^* \varphi'(\alpha^*) &= 0, \\ \alpha^* &\geq 0. \end{aligned}$$

These conditions state that if  $\alpha^* > 0$ , then  $\varphi'(\alpha^*) = 0$ , which implies that

$$\nabla f(\mathbf{x}_k + \alpha^* \mathbf{p}_k)^T \mathbf{p}_k = 0,$$

meaning that the search direction  $\mathbf{p}_k$  is orthogonal to the gradient of  $f$  at  $\mathbf{x}_k + \alpha^* \mathbf{p}_k$ .



However, solving the line search problem to optimality is unnecessary. The optimal solution to the original problem lies elsewhere anyway. Examples of methods to choose step lengths  $\alpha_k$

- *Interpolation*: Use  $f(\mathbf{x}_k)$ ,  $\nabla f(\mathbf{x}_k)$ , and  $\nabla f(\mathbf{x}_k)^T \mathbf{p}_k$  to approximate  $\varphi = f(\mathbf{x}_k + \alpha \mathbf{p}_k)$  quadratically. Then minimize this approximation of  $\varphi$  analytically.
- *Newton's method*: Repeat improvements from a quadratic approximation:  $\alpha = \alpha - \varphi'(\alpha) / \varphi''(\alpha)$
- *Golden section*: Derivative-free method which shrinks an interval wherein a solution to  $\varphi'(\alpha) = 0$  lies.

We will often use what is denoted as the *Armijo rule*. The idea is to choose a step length  $\alpha$  which provides sufficient decrease in  $f$ . We have that

$$f(\mathbf{x}_k + \alpha \mathbf{p}_k) \approx f(\mathbf{x}_k) + \alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k,$$

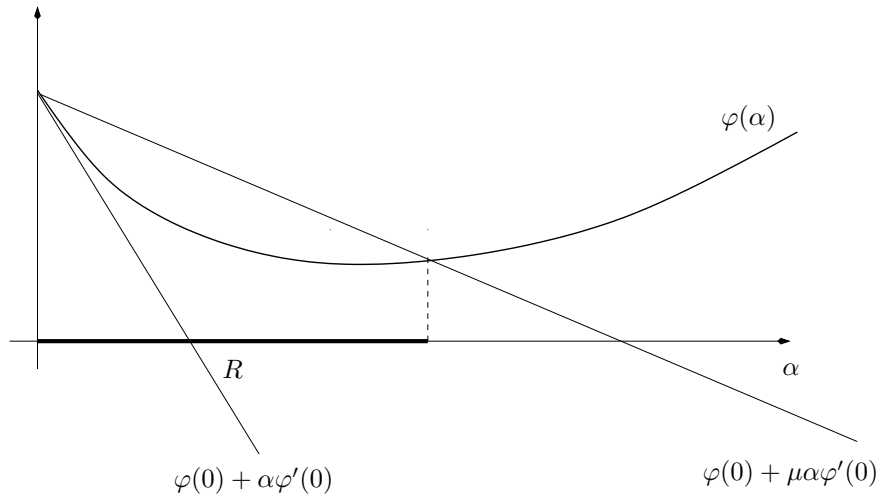
for very small values of  $\alpha > 0$ , meaning that we predict that the objective function will decrease with  $\alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k$  if we move a step length  $\alpha$  in the direction of  $\mathbf{p}_k$ . Now this might be too optimistic, and we will therefore accept the step length if the actual decrease is at least a fraction  $\mu$  ( $\mu$  is small, typically  $\mu \in [0.001, 0.01]$ ) of the predicted decrease, i.e., we will accept  $\alpha$  if

$$f(\mathbf{x}_k + \alpha \mathbf{p}_k) - f(\mathbf{x}_k) \leq \mu \alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k,$$

or equivalently, if

$$\varphi(\alpha) - \varphi(0) \leq \mu \alpha \varphi'(0).$$

We usually start with  $\alpha = 1$ . If this is not fulfilled, then choose  $\alpha := \alpha/2$ .



**Figure 1:** The interval ( $R$ ) accepted by the Armijo step length rule

## Convergence

In order to state a convergence result for the algorithm, we make an additional assumption for the search directions. We need the directions  $\mathbf{p}_k$  to fulfill

$$-\frac{\nabla f(\mathbf{x}_k)^T \mathbf{p}_k}{\|\nabla f(\mathbf{x}_k)\| \cdot \|\mathbf{p}_k\|} \geq s_1, \quad \|\mathbf{p}_k\| \geq s_2 \|\nabla f(\mathbf{x}_k)\|, \quad \text{and} \quad \|\mathbf{p}_k\| \leq M \quad (2)$$

for some  $s_1, s_2 > 0$ , where the first inequality makes the angle between  $\mathbf{p}_k$  and  $\nabla f(\mathbf{x}_k)$  stay between 0 and  $\pi/2$ , but not too close to  $\pi/2$ . The second inequality makes sure that the only case when  $\mathbf{p}_k$  can be zero is when the gradient is zero. These two conditions guarantee a certain descent quality.

**Theorem** (convergence of unconstrained algorithm). *Suppose  $f \in C^1$  and for the starting point  $\mathbf{x}_0$  the level set  $\{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$  is bounded. Consider the iterative algorithm described above. Suppose that for all  $k$ ,  $\mathbf{p}_k$  fulfills (2) and  $\alpha_k$  is chosen according to the Armijo rule. Then*

- a) the sequence  $\{\mathbf{x}_k\}$  is bounded,
- b) the sequence  $\{f(\mathbf{x}_k)\}$  is descending and lower bounded, and
- c) every limit point of  $\{\mathbf{x}_k\}$  is a stationary point.

*Proof.* See Theorem 11.4 in the book. □

If we add the assumption that  $f$  is a convex function, then we can show that

$$\text{optimum exists} \quad \Longleftrightarrow \quad \{\mathbf{x}_k\} \text{ converges to an optimal solution.}$$

## Step 4: Termination criteria

We can not terminate the algorithm when  $\nabla f(\mathbf{x}_k) = 0$ , since this rarely happens. We need to have some tolerance level. Three examples are

- a)  $\|\nabla f(\mathbf{x}_k)\| \leq \varepsilon_1(1 + |f(\mathbf{x}_k)|)$ , where  $\varepsilon_1 > 0$  is small.
- b)  $f(\mathbf{x}_{k-1}) - f(\mathbf{x}_k) \leq \varepsilon_2(1 + |f(\mathbf{x}_k)|)$ , where  $\varepsilon_2 > 0$  is small.
- c)  $\|\mathbf{x}_k - \mathbf{x}_{k-1}\| \leq \varepsilon_3(1 + \|\mathbf{x}_k\|)$ , where  $\varepsilon_3 > 0$  is small.

Can also use the max-norm  $\|\cdot\|_\infty$  instead.

## A note on trust region methods

Trust region methods use a quadratic approximation of the function around the current iterate  $\mathbf{x}_k$ , avoid a line search but instead bound the length of the search direction. Let

$$\varphi_{\mathbf{x}_k}(\mathbf{p}) := f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p}.$$

Since this is a local approximation, we restrict our approximation to a *trust region* in the neighborhood of  $\mathbf{x}_k$ , i.e., we trust the model in the region where  $\|\mathbf{p}\| \leq \Delta_k$ . We then solve the problem to

$$\begin{aligned} &\text{minimize } \varphi_{\mathbf{x}_k}(\mathbf{p}), \\ &\text{subject to } \|\mathbf{p}\| \leq \Delta_k. \end{aligned}$$

and let the solution be  $\mathbf{p}_k$ . Then we update our iterate as  $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k$ . We also update the trust region parameter  $\Delta_k$  depending on the progress so far (actual reduction/predicted reduction). The method is robust and possess strong convergence. More detailed information about trust region methods can be found in the book on pages 301–302.

## A note on black-box functions

In some cases the value of the objective function  $f(\mathbf{x})$  is given through some unknown simulation procedure. This implies that we do not have a clear representation of the gradient of the objective function. In some cases, we can perform *numerical differentiation* and approximate the partial derivatives as, e.g.,

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + \alpha \mathbf{e}_i) - f(\mathbf{x})}{\alpha},$$

where  $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^T$  is the unit vector in  $\mathbb{R}^n$ .

If the simulation is not accurate, we get a bas derivative information. We can use *derivative-free methods* instead. These try to build a model  $\hat{f}$  of the objective function  $f$  from evaluating the objective function at some specific test points and optimize the model  $\hat{f}$  instead of the function  $f$ .