

BACKTRACKING LINE SEARCH

MARKUS GRASMAIR

1. LINE SEARCH METHODS

In the following, we will consider so called *line search methods* for the numerical solution of a free optimisation problem

$$(P) \quad \min_{x \in \mathbb{R}^d} f(x).$$

These algorithms form the largest and best known (and analysed) class of numerical solution methods for this type of problems. The main alternative are *trust region methods*, which will be discussed later after we have formulated some theoretical results concerning constrained optimisation.

1.1. General set-up. Line search methods are iterative algorithms, where one starts with an initial guess x_0 of a solution of (P), which is then successively improved. This improvement is carried out in two separate steps:

- (1) First we choose a *search direction* $p_k \in \mathbb{R}^d$.
- (2) Then we choose a *step length* $\alpha_k > 0$.

Having found p_k and α_k , we then define the next iterate as

$$x_{k+1} = x_k + \alpha_k p_k.$$

Our goal is the minimisation of the function f . Thus it makes sense to choose the search direction p_k in each step in such a way that the function values actually become smaller. Since the negative gradient $-\nabla f(x_k)$ is the direction in which the function f has the steepest downwards slope at the point x_k , it further sounds reasonable to choose $p_k = -\nabla f(x_k)$. Doing so, we then end up with the *gradient descent* or *steepest descent* method. An alternative is *Newton's method*, where one chooses $p_k = -H_f(x_k)^{-1} \nabla f(x_k)$, given that the Hessian $H_f(x_k)$ is invertible.

For the choice of the step length, we start with a trivial example: the function $f: \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = \frac{1}{2}x^2$. Here one step of the gradient descent method with step length $\alpha_k > 0$ reads

$$x_{k+1} = x_k - \alpha_k f'(x_k) = x_k - \alpha_k x_k = (1 - \alpha_k)x_k.$$

Let us now for simplicity assume that we always choose the same step length $\alpha_k = \alpha$. Then we simply have that

$$x_k = (1 - \alpha)^k x_0,$$

and the iteration converges to the unique global and local solution $x = 0$ if and only if we have chosen $\alpha < 2$.

These considerations show that we will need to bound the step length somehow in order to obtain a convergent method. The question, though, is how we choose that bound. One easily sees that we will have to do so depending on the function we want to minimise, and probably even depending on the current iterate: If we consider instead the function $f(x) = 10x^2$, then the gradient descent method with constant step length $\alpha > 0$ yields the iterates

$$x_{k+1} = (1 - 20\alpha)^k x_0.$$

This sequence converges to 0 if and only if $\alpha < 1/10$; with a step length of $\alpha = 1$ (which would have led us to the analytic solution in a single step in the first example) the same iterations with this function would rapidly diverge. Conversely, if one were to use a step length $\alpha = 1/20$ for the first function, then the iterations would converge, but would do so very slowly.

From these observations we see that the reasonable step lengths depend strongly on the function that is optimised and that it is not possible to formulate a general step length rule that does not take into account the cost function. In some situations, where one has a very good knowledge about the cost function including uniform upper and lower bounds for its curvature it can be possible to formulate a good choice of the step length a-priori. In most practical situations, however, it will be necessary to choose the step length adaptively in each iteration of the line search.

The first step towards such an adaptive step length choice is the formulation of a criterion that tells us whether a step length is acceptable or not. Since we intend to minimise the function f , it makes sense to require at a minimum that the function values actually decrease in each step. That is, one would require as a first idea a step length to satisfy the condition

$$(1) \quad f(x_k + \alpha p_k) < f(x_k)$$

for it to be acceptable. It turns out, though, that this condition is not enough to guarantee convergence of the algorithm.

As an example, consider the function

$$f(x) = x^2 - \frac{1}{10}(x^2 - 1)^2.$$

If one uses the gradient descent method with step length $\alpha = 1$ for the minimisation of this function, and starts the iteration with a value x_0 that is slightly larger than 1 (e.g. $x_0 = 1.1$), then the iterates x_k will satisfy

$$\begin{aligned} x_{2k} &\rightarrow +1 && \text{as } k \rightarrow \infty, \\ x_{2k+1} &\rightarrow -1 && \text{as } k \rightarrow \infty, \\ |f'(x_k)| &\rightarrow 2 && \text{as } k \rightarrow \infty, \\ f(x_{k+1}) &< f(x_k) && \text{for all } k. \end{aligned}$$

Thus the decrease condition (1) is satisfied in every single gradient descent step, but still the iterates do not converge.

1.2. Armijo condition and backtracking. The main idea of the Armijo condition is to replace the decrease condition (1) by the condition that the function values decrease *sufficiently*. To that end, we first note that the Taylor series expansion

$$f(x_k + \alpha p_k) = f(x_k) + \alpha \langle \nabla f(x_k), p_k \rangle + o(\alpha)$$

indicates that we should *expect* a decrease of approximately $\alpha \langle \nabla f(x_k), p_k \rangle$ for small step lengths $\alpha > 0$. Now we say that a step length α is acceptable, if the *actual* decrease $f(x_k) - f(x_k + \alpha p_k)$ is at least a fraction $0 < c_1 < 1$ of that value, that is, if the *Armijo condition* (or *sufficient decrease condition*)

$$(2) \quad f(x_k + \alpha p_k) \leq f(x_k) + c_1 \alpha \langle \nabla f(x_k), p_k \rangle$$

holds; else the step length α is deemed to large (see also Figure 1).

Now the question is, how we can find a good step length that satisfies the Armijo condition (2). One simple possibility is the idea of backtracking: Starting in each step with a *fixed* step length $\hat{\alpha} > 0$, one reduces the step length by a constant factor $\rho > 0$ until the Armijo condition is satisfied. The resulting method is summarised in Algorithm 1.

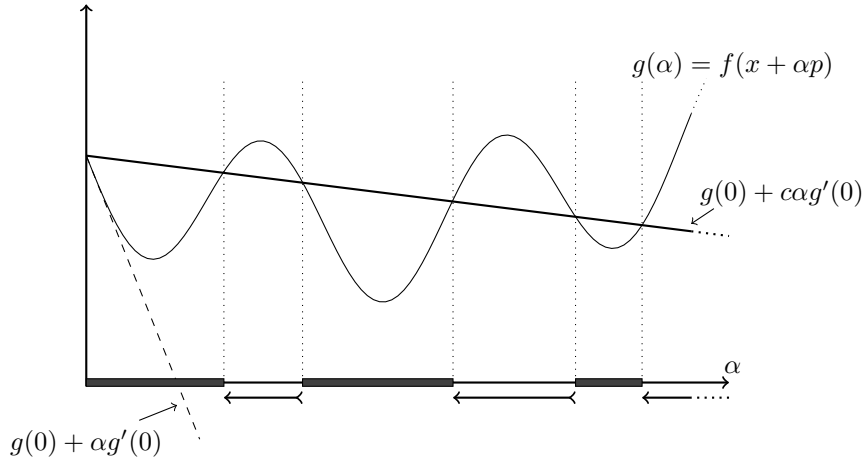


FIGURE 1. Armijo condition. The dark rectangles indicate the ranges of step lengths that are acceptable. All values outside these ranges are deemed “too large” by an algorithmic approach (indicated by the arrows below the α -axis).

Data: an objective function $f: \mathbb{R}^d \rightarrow \mathbb{R}$;
 an initial guess x_{init} ;
Result: $x^* \in \mathbb{R}^d$;

Initialisation: choose $0 < c < 1$ (sufficient decrease parameter);
 choose $\hat{\alpha} > 0$ (initial step length);
 choose $0 < \rho < 1$ (contraction factor);
 set $x_0 := x_{\text{init}}$, $k = 0$;

while convergence criterion not yet satisfied **do**
 | choose a search direction p_k with $\langle \nabla f(x_k), p_k \rangle < 0$;
 | set $\alpha = \hat{\alpha}$;
 | **while** $f(x_k + \alpha p_k) > f(x_k) + c\alpha \langle \nabla f(x_k), p_k \rangle$ **do**
 | | $\alpha \leftarrow \rho \alpha$;
 | **end**
 | define $x_{k+1} := x_k + \alpha p_k$;
 | $k \leftarrow k + 1$;
end
 define $x^* := x^{(k)}$;

Algorithm 1: Backtracking line search algorithm.

2. CONVERGENCE ANALYSIS

We now formulate and prove the main convergence theorem for backtracking line search methods. In order to be able to apply this result for a large class of different search directions, we do this in a quite general setting. Later, we will discuss specifically the gradient descent and Newton methods.

Theorem 2.1. Assume that $f \in C^2(\mathbb{R}^d)$, that $x_0 \in \mathbb{R}^d$, and that the set

$$S := L_f(f(x_0)) = \{x \in \mathbb{R}^d : f(x) \leq f(x_0)\}$$

is bounded. Consider the iteration

$$x_{k+1} = x_k + \alpha_k p_k.$$

Assume that the following hold:

- The search direction $p_k \in \mathbb{R}^d$ is a descent direction for all $k \in \mathbb{N}$ in the sense that

$$\langle p_k, \nabla f(x_k) \rangle < 0.$$

- There exist constants $C_0 > 0$ and $C_1 > 0$ such that

$$(3) \quad |\langle \nabla f(x_k), p_k \rangle| \geq C_0 \|\nabla f(x_k)\| \|p_k\| \quad \text{and} \quad \|p_k\| \geq C_1 \|\nabla f(x_k)\|$$

for all $k \in \mathbb{N}$.

- The step length α_k is chosen by backtracking Armijo line search, that is, Algorithm 1.

Then the sequence $(x_k)_{k \in \mathbb{N}}$ is bounded,

$$\lim_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0,$$

and every accumulation point of this sequence is a stationary point of f .

Proof. By construction, the iteration will only terminate if x_k is a stationary point, that is, if $\nabla f(x_k) = 0$. Since we have assumed that $\langle \nabla f(x_k), p_k \rangle < 0$ for all k , this can never be the case.

Since f is continuous, it follows that $S = L_f(f(x_0))$ is closed. Because it is by assumption bounded and non-empty (at least the point x_0 is contained in it), it follows that f admits a global minimum on S , and consequently on \mathbb{R}^d . In particular, the function f is bounded below, say

$$f(x) \geq K \quad \text{for all } x \in \mathbb{R}^d$$

for some constant $K \in \mathbb{R}$.

Now the Armijo condition and the assumption that p_k is a descent direction imply that

$$f(x_{k+1}) = f(x_k + \alpha_k p_k) \leq f(x_k) + c_1 \alpha_k \langle \nabla f(x_k), p_k \rangle < f(x_k)$$

for all k . In particular, $f(x_k) < f(x_0)$ for all k , which implies that $x_k \in S$ for all k . Thus the sequence $(x_k)_{k \in \mathbb{N}}$ is bounded. Moreover,

$$\sum_{k=0}^{\infty} (f(x_k) - f(x_{k+1})) = f(x_0) - \lim_{k \rightarrow \infty} f(x_k) \leq f(x_0) - K < \infty,$$

which implies that

$$\lim_{k \rightarrow \infty} (f(x_k) - f(x_{k+1})) = 0.$$

Now the Armijo condition implies that

$$c_1 \alpha_k |\langle \nabla f(x_k), p_k \rangle| = -c_1 \alpha_k \langle \nabla f(x_k), p_k \rangle \leq f(x_k) - f(x_{k+1}) \rightarrow 0 \quad \text{as } k \rightarrow \infty.$$

From (3) we thus obtain that

$$(4) \quad \alpha_k \|p_k\| \|\nabla f(x_k)\| \rightarrow 0 \quad \text{as } k \rightarrow \infty.$$

Let us now assume to the contrary of the claim of the theorem that the sequence $(\|\nabla f(x_k)\|)_{k \in \mathbb{N}}$ does *not* converge to 0. Then there exist $\varepsilon > 0$ and a subsequence $(x_{k'})_{k'}$ such that

$$\|\nabla f(x_{k'})\| \geq \varepsilon \quad \text{for all } k'.$$

As a consequence, (4) implies that, necessarily,

$$\lim_{k'} \alpha_{k'} \|p_{k'}\| = 0.$$

Since $\|p_{k'}\| \geq C_1 \|\nabla f(x_{k'})\| \geq C_1 \varepsilon > 0$, this implies that also

$$\lim_{k'} \alpha_{k'} = 0.$$

In particular, we have for all sufficiently large k' that $\alpha_{k'} < \hat{\alpha}$. That is, for all these k' it was necessary to perform at least one contraction step in the backtracking algorithm.

Now recall that we stop the backtracking with the *first* value of α for which the Armijo condition is satisfied. In case we have to perform at least one contraction, this means that the previous value of α did *not* satisfy the Armijo condition. That is, for all sufficiently large k' the step length $\alpha_{k'}/\rho$ did not satisfy the Armijo condition, implying that

$$(5) \quad f\left(x_{k'} + \frac{\alpha_{k'}}{\rho} p_{k'}\right) > f(x_{k'}) + c_1 \frac{\alpha_{k'}}{\rho} \langle \nabla f(x_{k'}), p_{k'} \rangle$$

for all sufficiently large k' . Now we use a Taylor series expansion of f around $x_{k'}$ to obtain that

$$f\left(x_{k'} + \frac{\alpha_{k'}}{\rho} p_{k'}\right) = f(x_{k'}) + \frac{\alpha_{k'}}{\rho} \langle \nabla f(x_{k'}), p_{k'} \rangle + \frac{\alpha_{k'}^2}{2\rho^2} \langle p_{k'}, H_f(y_{k'}) p_{k'} \rangle,$$

where the point $y_{k'}$ lies on the line segment between $x_{k'}$ and $x_{k'} + \frac{\alpha_{k'}}{\rho} p_{k'}$. Inserting this in (5) and dividing by $\alpha_{k'}/\rho > 0$, we obtain that

$$\langle \nabla f(x_{k'}), p_{k'} \rangle + \frac{\alpha_{k'}}{2\rho} \langle p_{k'}, H_f(y_{k'}) p_{k'} \rangle > c_1 \langle \nabla f(x_{k'}), p_{k'} \rangle,$$

or

$$(1 - c_1) \langle \nabla f(x_{k'}), p_{k'} \rangle > -\frac{\alpha_{k'}}{2\rho} \langle p_{k'}, H_f(y_{k'}) p_{k'} \rangle$$

for all k' . Using again (3) and the assumption that $p_{k'}$ is a descent direction, this implies that

$$(1 - c_1) C_0 \|\nabla f(x_{k'})\| \|p_{k'}\| < \frac{\alpha_{k'}}{2\rho} \langle p_{k'}, H_f(y_{k'}) p_{k'} \rangle \leq \frac{\alpha_{k'}}{2\rho} \|p_{k'}\|^2 \|H_f(y_{k'})\|_2$$

and thus, after dividing by $\|p_{k'}\|$ and using the estimate $\|\nabla f(x_{k'})\| \geq \varepsilon$,

$$(6) \quad (1 - c_1) C_0 \varepsilon < \frac{\alpha_{k'} \|p_{k'}\|}{2\rho} \|H_f(y_{k'})\|_2$$

for all k' . Now recall that $\alpha_{k'} \|p_{k'}\| \rightarrow 0$ as $k' \rightarrow \infty$. Thus (6) implies that

$$(7) \quad \|H_f(y_{k'})\|_2 \rightarrow \infty \text{ as } k' \rightarrow \infty.$$

However, the points $y_{k'}$ lie on the line segments between $x_{k'}$ and $x_{k'} + \frac{\alpha_{k'}}{\rho} p_{k'}$. Since the points $x_{k'}$ lie in the bounded set S and $\alpha_{k'} \|p_{k'}\| \rightarrow 0$ as $k' \rightarrow \infty$, it follows that the points $y_{k'}$ lie in a bounded set as well. Now f is by assumption twice continuously differentiable on \mathbb{R}^d and thus the Hessian is continuous. Therefore its norm is bounded on the bounded set $\{y_{k'}\}_{k'}$, which contradicts (7).

Therefore our initial assumption that $\|\nabla f(x_k)\|$ does not converge to 0 was false, which implies that, in fact, it does.

It remains to show that every accumulation point of the sequence $\{x_k\}_{k \in \mathbb{N}}$ is a stationary point of f . Assume therefore that \hat{x} is an accumulation point of this sequence. Then there exists a subsequence $\{x_{k''}\}_{k''}$ with $\lim_{k''} x_{k''} = \hat{x}$. Now the continuity of ∇f implies that $\nabla f(\hat{x}) = \lim_{k''} \nabla f(x_{k''}) = 0$, which concludes the proof. \square

Remark 2.2. The second condition in (3) prevents the search directions from being “too short” compared to the gradient. The first condition in (3) prevents the search conditions from being almost orthogonal to the gradient.

Remark 2.3. From Theorem 2.1, it does not follow that the sequence $\{x_k\}_{k \in \mathbb{N}}$ necessarily converges. However, the situations where this can happen are rather limited:

First we recall (?) that a bounded sequence $\{x_k\}_{k \in \mathbb{N}}$ that has only a single accumulation point necessarily converges to that point.¹ Next, we note that the sequence $\{f(x_k)\}_{k \in \mathbb{N}}$ is decreasing and bounded below, and thus it converges to some value $v := \lim_{k \rightarrow \infty} f(x_k)$.

Now assume that the sequence $\{x_k\}_{k \in \mathbb{N}}$ in Theorem 2.1 does *not* converge. Since the sequence is bounded, it has at least one accumulation point. Since it does not converge, we can even conclude that it has at least *two* accumulation points, say $y_1, y_2 \in \mathbb{R}^d$. Moreover, by Theorem 2.1, we have that $\nabla f(y_1) = \nabla f(y_2) = 0$, and from the convergence of $\{f(x_k)\}_{k \in \mathbb{N}}$ we obtain that $f(y_1) = f(y_2) = v$.

Thus a necessary condition for the sequence $\{x_k\}_{k \in \mathbb{N}}$ *not* to converge, is that there exist different critical points of f with the same function value. Put differently, if all the critical points of f have a different function value, then the sequence $\{x_k\}_{k \in \mathbb{N}}$ necessarily converges to one of these points.

3. APPLICATION TO DIFFERENT METHODS

3.1. Backtracking gradient descent. In the gradient descent method, we choose the search direction as $p_k = -\nabla f(x_k)$. In this case, the estimates in (3) are trivially satisfied with $C_0 = C_1 = 1$ (and equalities instead of inequalities). Thus Theorem 2.1 implies that the backtracking gradient descent method converges in the sense that $\nabla f(x_k) \rightarrow 0$, provided that the function f is twice continuously differentiable and the lower level-set $L_f(f(x_0))$ is bounded.

In addition, we can estimate

$$\|x_{k+1} - x_k\| = \|\alpha_k p_k\| = \|\alpha_k \nabla f(x_k)\| \leq \hat{\alpha} \|\nabla f(x_k)\|.$$

Since $\lim_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0$, we can conclude that also

$$\lim_{k \rightarrow \infty} \|x_{k+1} - x_k\| = 0.$$

That is, the difference between consecutive iterates in the algorithm tends to zero.

3.2. General search directions. In several important cases (in particular Newton's method and quasi-Newton methods), the search direction is defined as

$$p_k = -B_k^{-1} \nabla f(x_k)$$

where $B_k \in \mathbf{Sym}_d$ is a positive definite, symmetric matrix depending on the current iterate x_k , the function f , and possibly also earlier iterates. In this case, one can show that the estimates in (3) hold, if there exists a constant $C > 0$ such that

$$\|B_k\|_2 \leq C \quad \text{and} \quad \|B_k^{-1}\|_2 \leq C$$

for each $k \in \mathbb{N}$.

First we note that, in this case,

$$\|\nabla f(x_k)\| = \|B_k p_k\| \leq \|B_k\|_2 \|p_k\| \leq C \|p_k\|,$$

and thus the second estimate in (3) holds with $C_1 = 1/C$.

¹Here is a sketch of the proof: Denote by \hat{x} the single accumulation points of the sequence, but assume that the sequence does not converge to \hat{x} . Then there exist a subsequence $\{x_{k'}\}_{k'}$ and $\varepsilon > 0$ with $\|x_{k'} - \hat{x}\| > \varepsilon$ for all k' . Now the sequence $\{x_{k'}\}_{k'}$ is still bounded and thus has a convergent subsequence, say $\{x_{k''}\}_{k''}$, with $y = \lim_{k'' \rightarrow \infty} x_{k''}$. By construction $\|y - \hat{x}\| \geq \varepsilon > 0$, but at the same time y is an accumulation point of $\{x_{k'}\}_{k' \in \mathbb{N}}$, which contradicts the uniqueness of \hat{x} .

Now denote for $A \in \mathbf{Sym}_d$ by $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ the largest and smallest eigenvalue of A , respectively. Recall moreover that, for positive definite symmetric matrices, we have the equalities

$$\|A\|_2 = \lambda_{\max}(A) \quad \text{and} \quad \|A^{-1}\|_2 = \frac{1}{\lambda_{\min}(A)}.$$

Moreover, if $A \in \mathbf{Sym}_d$ is positive definite, we can estimate

$$\langle y, Ay \rangle \geq \lambda_{\min}(A) \|y\|^2$$

for all $y \in \mathbb{R}^d$. Applying this estimate with $A = B_k$, we obtain that

$$\begin{aligned} |\langle p_k, \nabla f(x_k) \rangle| &= \langle p_k, B_k p_k \rangle \geq \lambda_{\min}(B_k) \|p_k\|^2 \\ &= \frac{1}{\|B_k^{-1}\|} \|p_k\|^2 \geq \frac{1}{C} \|p_k\|^2 \geq \frac{1}{C^2} \|p_k\| \|\nabla f(x_k)\|. \end{aligned}$$

In the last step, we have here used the estimate $\|p_k\| \geq C \|\nabla f(x_k)\|$ shown previously. Thus the first estimate in (3) is also satisfied with $C_0 = 1/C^2$.

3.3. Backtracking (damped) Newton's method. Finally, we consider Newton's method, where the search directions are defined as

$$p_k = -H_f(x_k)^{-1} \nabla f(x_k).$$

Before we starting the discussion of the convergence of this method, we first have to address its well-posedness. Indeed, the definition of the search direction in Newton's method only makes sense, if the Hessian $H_f(x_k)$ is non-singular; else p_k is not well-defined. In addition, even if $H_f(x_k)$ is invertible, it is not guaranteed that p_k actually is a descent direction and thus useful as a search direction. If, for instance, $H_f(x_k) = -\text{Id}$, then we would obtain the search direction $p_k = \nabla f(x_k)$, and the function f would *increase* along p_k . We can, however, guarantee that $p_k = -H_f(x_k)^{-1} \nabla f(x_k)$ is a search direction in the case where $H_f(x_k)$ is positive definite: Since the inverse of a positive definite symmetric matrix is positive definite as well, we obtain in this case that

$$\langle \nabla f(x_k), p_k \rangle = -\langle \nabla f(x_k), H_f(x_k)^{-1} \nabla f(x_k) \rangle < 0$$

as long as $\nabla f(x_k) \neq 0$. Since the positive definiteness of the Hessian H_f is strongly linked to the convexity of f , we will therefore discuss first the case where f is a convex function. In fact, we will require a somewhat stronger notion:

Definition 3.1. A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is *strongly convex*, if there exists $C > 0$ such that the function $g: \mathbb{R}^d \rightarrow \mathbb{R}$, $g(x) := f(x) - C\|x\|^2$, is convex.²

One can show the following characterisation of twice continuously differentiable strongly convex functions:

Proposition 3.2. Assume that $f \in C^2(\mathbb{R}^d)$. Then f is strongly convex, if and only if $H_f(x)$ is non-singular for all $x \in \mathbb{R}^d$ and there exists $C > 0$ such that

$$\|H_f(x)^{-1}\|_2 \leq C \quad \text{for all } x \in \mathbb{R}^d.$$

It is easy to see that strongly convex functions are in particular strictly convex. In addition, it is possible to show that every strongly convex function is coercive. Thus it follows that every strongly convex function has a unique minimiser.

Corollary 3.3. Assume that $f \in C^2(\mathbb{R}^d)$ is strongly convex and the iterates $x_{k+1} = x_k + \alpha_k p_k$ are generated with Newton's method with backtracking Armijo line search. Then the sequence $\{x_k\}_{k \in \mathbb{N}}$ converges to the unique minimiser of f .

²Her må man være litt forsiktig med oversettelser til norsk: "strictly convex" blir oversatt som "strengt konveks," mens "strongly convex" blir "sterkt konveks."

Proof. Since f is strongly convex, it actually has a unique minimiser. Moreover, since f is differentiable and convex, this minimiser is also the unique critical point of f . In view of Remark 2.3, it is thus sufficient to show that the conditions of Theorem 2.1 are satisfied.

The boundedness of the level set $L_f(f(x_0))$ follows from the coercivity of f (which is a consequence of its strong convexity). Since $p_k = -H_f(x_k)^{-1}\nabla f(x_k)$ and f is strongly convex, it follows that p_k is a descent direction for each k . Following the discussion in Section 3.2, it thus remains to show that there exists $C > 0$ such that $\|H_f(x_k)\|_2 \leq C$ and $\|H_f(x_k)^{-1}\|_2 \leq C$ for all k . The latter inequality is an immediate consequence of Proposition 3.2. For the former inequality, we note that H_f is continuous and that the iterates $\{x_k\}_{k \in \mathbb{N}}$ are contained in the bounded set $L_f(f(x_0))$. This proves the boundedness of $\{\|H_f(x_k)\|_2\}_{k \in \mathbb{N}}$, which in turn concludes the proof of the assertion. \square

Now we will briefly discuss the case, where f is not necessarily (strongly) convex. Here it is necessary to modify Newton's method in order to obtain a well-posed (and convergent) optimisation method. For that, we discuss two relatively simple methods:

- In the first method, we change the Hessian matrix by modifying its eigenvalues if necessary. For that, we first choose a lower bound $\varepsilon > 0$ for the eigenvalues.

In each step of our algorithm, we then start by computing an orthogonal eigendecomposition of $H_f(x_k)$. That is, we compute an orthogonal matrix $U_k \in \mathbb{R}^{d \times d}$ and a diagonal matrix $\Lambda_k = \text{diag}(\lambda_1, \dots, \lambda_d)$ such that

$$H_f(x_k) = U_k \Lambda_k U_k^T.$$

Such an orthogonal eigendecomposition exists, as $H_f(x_k)$ is a symmetric matrix. Now we define $m_k := \max\{\varepsilon, \lambda_k\}$, set $M_k := \text{diag}(m_1, \dots, m_d)$, and $B_k := U_k M_k U_k^T$, and compute

$$p_k = -B_k^{-1}\nabla f(x_k) = -U_k M_k^{-1} U_k^T \nabla f(x_k).$$

Then the eigenvalues of the matrix B_k are bounded below by $\varepsilon > 0$, and thus B_k is positive definite with $\|B_k^{-1}\|_2 \leq 1/\varepsilon$. At the same time, $\|B_k\|_2 \leq \max\{\varepsilon, \|H_f(x_k)\|_2\}$, and thus $\|B_k\|_2$ will be uniformly bounded above, if f is C^2 and the sequence $\{x_k\}_{k \in \mathbb{N}}$ stays bounded. Therefore, the conditions from Section 3.2 are satisfied.

Of course, different choices of the new eigenvalues m_k are possible, provided that they are positive and bounded away from zero. One interesting alternative that is worth considering is to define $m_k := \max\{\varepsilon, |\lambda_k|\}$.

- As an alternative, we now discuss a method that does not require the computation of an eigendecomposition of the Hessian. Here, we simply switch to the gradient descent direction, whenever the Newton direction does not make sense. Again, we fix a parameter $\varepsilon > 0$ for the method.

In each step of our algorithm, we start by *trying* to compute $p_k := -H_f(x_k)^{-1}\nabla f(x_k)$. If this is not possible (because the matrix $H_f(x_k)$ is singular or ill-conditioned), we instead use the gradient descent direction $p_k := -\nabla f(x_k)$. Also, we switch to the gradient descent direction, if $\langle \nabla f(x_k), p_k \rangle < \varepsilon \|\nabla f(x_k)\| \|p_k\|$. This guarantees that the first condition in (3) is satisfied in each step. Moreover, the second condition in (3) is satisfied because of the estimate $\|\nabla f(x_k)\| \leq \max\{1, \|H_f(x_k)\|\} \|p_k\|$.