

## Bachelorarbeit

# The Riemannian BFGS Method and its Implementation in Julia

Abschlussarbeit zur Erlangung des akademischen Grades

Bachelor of Science

vorgelegt von  
Matrikelnummer  
1. Prüfer  
2. Prüfer

Tom-Christian Riemer  
444019  
Referee A  
Referee B

Abgabedatum

In the Future, last compiled 23. September 2020

# INHALTSVERZEICHNIS

1	THE BFGS-METHOD FOR THE EUCLIDEAN CASE	3
1.1	Preliminaries . . . . .	3
1.2	Quasi-Newton Methods . . . . .	9
1.3	The Broyden-Fletcher-Goldfarb-Shanno Formula . . . . .	13
1.4	The BFGS-Method . . . . .	18
1.5	A Cautious BFGS-Method . . . . .	21
1.6	Limited-Memory BFGS-Method . . . . .	23
2	RIEMANNIAN MANIFOLDS	29
3	THE BFGS-METHOD FOR RIEMANNIAN MANIFOLDS	32
3.1	Preliminaries . . . . .	32
3.2	Quasi-Newton Methods For Riemannian Manifolds . . . . .	36
3.3	The BFGS Formula For Riemannian Manifolds . . . . .	42
3.4	The BFGS Method On Riemannian Manifolds . . . . .	46
3.5	Cautious BFGS-Method On Riemannian Manifolds . . . . .	49
3.6	Limited-Memory BFGS-Method On Riemannian Manifolds . . . . .	51
4	NUMERICS	54
4.1	Realizing the Update-Formula . . . . .	54
5	CONCLUSION	56

# 1 THE BFGS-METHOD FOR THE EUCLIDEAN CASE

## 1.1 PRELIMINARIES

$$\min f(x), \quad x \in \mathbb{R}^n \quad (1.1.1)$$

---

**Algorithm 1** General descent method

---

```
1:  $x_0 \in \mathbb{R}^n$ ,  $k = 0$ 
2: while  $x_k$  does not satisfy any stopping criterion do
3:   Determine a descent direction  $d_k$  of  $f$  in  $x_k$ .
4:   Determine a stepsize  $\alpha_k > 0$  with  $f(x_k + \alpha_k d_k) < f(x_k)$ .
5:   Set  $x_{k+1} = x_k + \alpha_k d_k$  and  $k = k + 1$ .
6: end while
7: return  $x_k$ 
```

---

---

**Algorithm 2** Local Newton's method

---

```
1:  $x_0 \in \mathbb{R}^n$ ,  $0 \leq \epsilon < 1$ ,  $k = 0$ 
2: while  $\|\nabla f(x_k)\| > \epsilon$  do
3:   Determine  $d_k \in \mathbb{R}^n$  by solving
```

$$\nabla^2 f(x_k) d = -\nabla f(x_k).$$

```
4:   Set  $x_{k+1} = x_k + d_k$  and  $k = k + 1$ .
5: end while
6: return  $x_k$ 
```

---

### Wolfe conditions:

A popular inexact line search condition is the so called Armijo condition or Armijo rule. It stipulates a stepsize  $\alpha_k$  that should first of all give a sufficient decrease in the objective function  $f$ , as measured by the following inequality:

$$f(x_k + \alpha_k d_k) \leq f(x_k) + c_1 \alpha_k \nabla f(x_k)^T d_k \quad (1.1.2)$$

for some constant  $c_1 \in (0, 1)$ . In other words, the reduction in  $f$  should be proportional to the stepsize  $\alpha_k$  and the directional derivative  $\nabla f(x_k)^T d_k$ . Eq. (1.1.2) is often called Armijo condition. Eq. (1.1.2) is not enough by itself to ensure that the algorithm makes reasonable progress because it is

satisfied for all sufficiently small values of  $\alpha_k$ . To rule out unacceptably short steps we introduce a second requirement, which requires  $\alpha_k$  to satisfy

$$\nabla f(x_k + \alpha_k d_k)^T d_k \geq c_2 \nabla f(x_k)^T d_k \quad (1.1.3)$$

for some constant  $c_2 \in (c_1, 1)$ . Eq. (1.1.3) is often referred to in the literature as “curvature condition”. This can lead to confusion in relation to quasi-Newton methods. The geometric interpretation of Eq. (1.1.3) is that the slope  $\nabla f(x_k + \alpha_k d_k)^T d_k$ , which is simply the derivative of  $f(x_k + \alpha_k d_k)$ , at the acceptable point must be greater than or equal to some multiple  $c_2$  of the initial slope  $\nabla f(x_k)^T d_k$ . This makes sense because if the slope  $\nabla f(x_k + \alpha_k d_k)^T d_k$  is strongly negative, we have an indication that we can reduce  $f$  significantly by moving further along the chosen direction. On the other hand, if  $\nabla f(x_k + \alpha_k d_k)^T d_k$  is only slightly negative or even positive, it is a sign that we cannot expect much more decrease in  $f$  in this direction, so it makes sense to terminate the line search.

Eq. (1.1.2) and Eq. (1.1.3) are known as the Wolfe-Powell inexact line search rule, Wolfe-Powell rule or just Wolfe conditions with  $0 < c_1 < c_2 < 1$ . The requirement  $0 < c_1 < c_2 < 1$  is necessary, such that there exists stepsize  $\alpha_k$  satisfying the Wolfe-Powell rule (see Sun, Yuan, 2006, p. 104-105).

It should point out that Eq. (1.1.3) is an approximation of the orthogonal condition  $\nabla f(x_{k+1})^T d_k = 0$ . However, unfortunately, one possible disadvantage of Eq. (1.1.3) is that it does not reduce to an exact line search in the limit  $c_2 \rightarrow 0$ . In addition, a stepsize may satisfy the Wolfe conditions without being particularly close to a minimizer of  $f(x_k + \alpha d_k)$ . We can, however, modify Eq. (1.1.3) to force  $\alpha_k$  to lie in at least a broad neighborhood of a local minimizer or stationary point of  $f(x_k + \alpha d_k)$ . The strong Wolfe conditions (or strong Wolfe-Powell rule) require  $\alpha_k$  to satisfy Eq. (1.1.2) and

$$|\nabla f(x_k + \alpha_k d_k)^T d_k| \geq c_2 |\nabla f(x_k)^T d_k| \quad (1.1.4)$$

with  $0 < c_1 < c_2 < 1$ . The only difference with the Wolfe conditions is that we no longer allow the derivative  $\nabla f(x_k + \alpha_k d_k)^T d_k$  to be too positive. Hence, we exclude points that are far from stationary points of  $f(x_k + \alpha d_k)$ .

In general, the smaller the value  $c_2$ , the more exact the line search. Normally, taking  $c_2 = 0.1$  gives a fairly accurate line search, whereas the value  $c_2 = 0.9$  gives a weak line search. However, taking too small  $c_2$  may be unwise, because the smaller the value  $c_2$ , the more expensive the computing effort. Usually,  $c_1 = 0.1$  and  $c_2 = 0.4$  are suitable, and it depends on the specific problem.

It is not difficult to prove that there exist stepsizes that satisfy the Wolfe conditions for every function  $f$  that is smooth and bounded below.

**Lemma 1.1.1** (Nocedal, Wright, 2006, Lemma 3.1). *Suppose that  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable. Let  $d_k$  be a descent direction at  $x_k$ , and assume that  $f$  is bounded below along the ray  $\{x_k + \alpha d_k | \alpha > 0\}$ . Then if  $0 < c_1 < c_2 < 1$ , there exist intervals of stepsizes satisfying the Wolfe conditions and the strong Wolfe conditions.*

The Wolfe conditions are scale-invariant in a broad sense: Multiplying the objective function by a constant or making an affine change of variables does not alter them. They can be used in most line search methods, and are particularly important in the implementation of quasi-Newton methods

Nocedal, Wright, 2006, p. 33-35.

The question now arises of how the stepsize is determined. We therefore introduce the concept of a stepsize strategy:

**Definition 1.1.2** (Geiger, Kanzow, 1999, p. 27). *A map  $T$  of  $\mathbb{R}^n \times \mathbb{R}^n$  into the power set of positive real numbers  $\mathcal{P}((0, \infty))$ , i.e. a map that assigns a subset  $T(x, d)$  of  $\mathbb{R}$  to each pair  $(x, d)$ , is called stepsize strategy or stepsize rule. We call such a stepsize rule (under certain conditions) well-defined, if (under these conditions) the set  $T(x, d)$  for each pair  $(x, d)$  with  $\nabla f(x)^T d < 0$  is not empty.*

Among this set of illustrations there is a certain subset, which includes one that meets the Wolfe conditions. We talk about the efficient stepsize strategies introduced by Warth, Werner, 1977, Definition 0.1. But we use the following definition for reasons of comprehensibility:

**Definition 1.1.3** (Geiger, Kanzow, 1999, Definition 4.5). *Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be continuously differentiable,  $x \in \mathbb{R}^n$  and  $d \in \mathbb{R}^n$  a descent direction from  $f$  in  $x$ . A stepsize strategy  $T$  is called efficient if there is a constant  $\theta > 0$  independent of  $x$  and  $d$  with*

$$f(x + \alpha d) \leq f(x) - \theta \left( \frac{\nabla f(x)^T d}{\|d\|} \right)^2$$

for all  $\alpha \in T(x, d)$ .

The stepsize strategy introduced in Powell, 1976 (see. Werner, 1978/79, Definition 2.2.) meets the Wolfe conditions. It can be shown that this strategy is efficient. Because of the introduction in Powell, 1976, it is often called Wolfe-Powell stepsize strategy, in brief, the Wolfe-Powell rule.

Stepsize strategies are mainly used in theory. In practical applications, algorithms are used to obtain stepsizes that meet the respective conditions. There are different algorithms that find a stepsize that meets the (strong) Wolfe conditions. In practice (mainly because of convergence reasons) the stepsize  $\alpha_k = 1$  is preferred if it fits.

In contrast to these stepsize strategies, there is the so called exact line search or optimal line search where a  $\alpha_k > 0$  is chosen such that the objective function in the direction  $d_k$  is minimized, i.e.,  $f(x_k + \alpha_k d_k) = \min_{\alpha > 0} f(x_k + \alpha d_k)$  Sun, Yuan, 2006, p. 71.

Since, in practical computation, an exact optimal stepsize generally cannot be found, and it is also expensive to find almost an exact stepsize, therefore the inexact line search with less computation load is highly popular Sun, Yuan, 2006, p. 72.

For many optimization methods, for example, Newton method and quasi-Newton method, their convergence rate does not depend on the exact line search. Therefore, as long as there is an acceptable stepsize strategy which ensures that the objective function has sufficient descent, the exact line search can be avoided and the computing efforts will be decreased greatly Sun, Yuan, 2006, p. 102.

The line search, which should satisfy either the Wolfe conditions or the strong Wolfe conditions, should always try the stepsize  $\alpha_k = 1$  first, because this stepsize will eventually always be accepted (under certain conditions), thereby producing superlinear convergence of the overall algorithm. Computational observations strongly suggest that it is more economical, in terms of function evaluations, to perform

a fairly inaccurate line search. The values  $c_1 = 10^{-4}$  and  $c_2 = 0.9$  are commonly used for the Wolfe conditions [Nocedal, Wright, 2006](#), p. 142.

**Satz 1.1.4** ([Sun, Yuan, 2006](#), Theorem 1.2.15). *Let  $A \in \mathbb{R}^{n \times n}$  be nonsingular and  $u, v \in \mathbb{R}^n$  be arbitrary. If*

$$1 + v^T A^{-1} u \neq 0,$$

*then the rank-one update  $A + uv^T$  of  $A$  is nonsingular, and its inverse is represented by*

$$(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u}.$$

**Satz 1.1.5** (Sherman-Morrison-Woodbury Theorem [Sun, Yuan, 2006](#), Theorem 1.2.16). *Let  $A \in \mathbb{K}^{n \times n}$  be a nonsingular matrix,  $U, V \in \mathbb{K}^{n \times m}$ . If  $I + V^* A^{-1} U$  is invertible, then  $A + UV^*$  is invertible and*

$$(A + UV^*)^{-1} = A^{-1} - A^{-1}U(I + V^* A^{-1}U)^{-1}V^* A^{-1}.$$

Note: The Sherman-Morrison-Woodbury formula can be extended to rank  $R$  modifications [Ulbrich, Ulbrich, 2012](#), p. 70.

**Satz 1.1.6** ([Nocedal, Wright, 2006](#), Theorem 3.2.). *Consider any iteration of the form  $x_{k+1} = x_k + \alpha_k d_k$ , where  $d_k$  is a descent direction and  $\alpha_k$  satisfies the Wolfe conditions [Eq. \(1.1.2\)](#) and [Eq. \(1.1.3\)](#). Suppose that  $f$  is bounded below in  $\mathbb{R}^n$  and that  $f$  is continuously differentiable in an open set  $\mathcal{N}$  containing the level set  $\mathcal{L} = \{x \in \mathbb{R}^n : f(x) \leq f(x_0)\}$ , where  $x_0$  is the starting point of the iteration. Assume also that the gradient  $\nabla f$  is Lipschitz continuous on  $\mathcal{N}$ , that is, there exists a constant  $L > 0$  such that*

$$\|\nabla f(x) - \nabla f(\tilde{x})\| \leq L\|x - \tilde{x}\|, \quad \text{for all } x, \tilde{x} \in \mathcal{N}.$$

*Then*

$$\sum_{k \geq 0} \cos(\theta_k)^2 \|\nabla f(x_k)\|^2 < \infty,$$

*where*

$$\cos(\theta_k) = -\frac{\nabla f(x_k)^T d_k}{\|\nabla f(x_k)\| \|d_k\|}.$$

[Theorem 1.1.6](#) quantifies the effect of properly chosen stepsizes  $\alpha_k$ . It describes how far  $d_k$  can deviate from the steepest descent direction  $\nabla f(x_k)$  and still produce a globally convergent iteration. Various line search termination conditions can be used to establish this result, but for concreteness we will consider only the Wolfe conditions [Eq. \(1.1.2\)](#) and [Eq. \(1.1.3\)](#), [Nocedal, Wright, 2006](#), p. 38. The theorem [Theorem 1.1.6](#) implies

$$\lim_{k \rightarrow \infty} \cos(\theta_k)^2 \|\nabla f(x_k)\|^2 = 0. \quad (1.1.5)$$

This limit [Eq. \(1.1.5\)](#) can be used in turn to derive global convergence results for line search algorithms. If the method for choosing the search direction  $d_k$  ensures that the angle  $\theta_k$  is bounded away from  $90^\circ$ , there is a positive constant  $\delta$  such that

$$\cos(\theta_k) \geq \delta > 0.$$

It follows immediately from [Eq. \(1.1.5\)](#) that

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

In other words, the gradient norms  $\|\nabla f(x_k)\|$  converge to zero, provided that the search directions  $d_k$  are never too close to orthogonality with the gradient.

In some cases (including the global convergence analysis of the BFGS method), only the weaker result

$$\liminf_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0. \quad (1.1.6)$$

can be shown using [Eq. \(1.1.5\)](#). That means, just a subsequence of the gradient norms  $\|\nabla f(x_{k_j})\|$  converges to zero, rather than the whole sequence. This can be shown by a contradiction. Suppose that [Eq. \(1.1.6\)](#) does not hold, so that the gradients remain bounded away from zero, that is, there exists  $\gamma > 0$  such that

$$\|\nabla f(x_k)\| \geq \gamma, \quad \text{for all } k \text{ sufficiently large.}$$

Then from [Eq. \(1.1.5\)](#) we conclude that

$$\lim_{k \rightarrow \infty} \cos(\theta_k) = 0 \quad (1.1.7)$$

that is, the entire sequence  $\{\cos(\theta_k)\}_k$  converges to 0. To establish [Eq. \(1.1.6\)](#), therefore, it is enough to

show that a subsequence  $\{\cos(\theta_{k_j})\}_j$  is bounded away from zero.

**Satz 1.1.7** (Nocedal, Wright, 2006, Theorem 3.6). *Suppose that  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is twice continuously differentiable. Consider the iteration  $x_{k+1} = x_k + \alpha_k d_k$ , where  $d_k$  is a descent direction and  $\alpha_k$  satisfies the Wolfe conditions with  $c_1 \leq \frac{1}{2}$ . If the sequence  $\{x_k\}_k$  converges to a point  $x^*$  such that  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive definite, and if the search direction satisfies*

$$\lim_{k \rightarrow \infty} \frac{\|\nabla f(x_k) + \nabla^2 f(x_k) d_k\|}{\|d_k\|} = 0 \quad (1.1.8)$$

then

- (i) the stepsize  $\alpha_k = 1$  is admissible for all  $k$  greater than a certain index  $k_0$  and
- (ii) if  $\alpha_k = 1$  for all  $k > k_0$ ,  $\{x_k\}_k$  converges to  $x^*$  superlinearly.

It is easy to see that if  $c_1 > \frac{1}{2}$ , then the line search would exclude the minimizer of a quadratic, and unit stepsizes may not be admissible.

If  $d_k$  is a quasi-Newton search direction of the form  $d_k = -H_k^{-1} \nabla f(x_k)$ , where the symmetric and positive definite matrix  $H_k$  is updated at every iteration by a quasi-Newton updating formula, then Eq. (1.1.8) is equivalent to

$$\lim_{k \rightarrow \infty} \frac{\|(H_k - \nabla^2 f(x^*)) d_k\|}{\|d_k\|} = 0 \quad (1.1.9)$$

Hence, we have the surprising (and delightful) result that a superlinear convergence rate can be attained even if the sequence of quasi-Newton matrices  $H_k$  does not converge to  $\nabla^2 f(x^*)$ ; it suffices that the  $H_k$  become increasingly accurate approximations to  $\nabla^2 f(x^*)$  along the search directions  $d_k$ . Importantly, condition Eq. (1.1.9) is both necessary and sufficient for the superlinear convergence of quasi-Newton methods Nocedal, Wright, 2006, p. 47.

**Folgerung 1.1.8** (Geiger, Kanzow, 1999, Lemma 7.9). *Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be twice continuously differentiable,  $\{H_k\}_k$  a sequence of regular matrices in  $\mathbb{R}^{n \times n}$ ,  $x_0 \in \mathbb{R}^n$  and  $\{x_k\}_k \subseteq \mathbb{R}^n$  a sequence defined by*

$$x_{k+1} = x_k - H_k^{-1} \nabla f(x_k), \quad k = 0, 1, \dots$$

with the limit  $\lim_{n \rightarrow \infty} x_k = x^*$ ,  $x_k \neq x^*$  for all  $k \in \mathbb{N}$  and  $\nabla^2 f(x^*)$  regular. Then the following statements are equivalent

- (i)  $\{x_k\}_k \rightarrow x^*$  superlinear and  $\nabla f(x^*) = 0$ .
- (ii)  $\|(\nabla^2 f(x_k) - H_k)(x_{k+1} - x_k)\| = o(\|x_{k+1} - x_k\|)$



$$(iii) \quad \|(\nabla^2 f(x^*) - H_k)(x_{k+1} - x_k)\| = o(\|x_{k+1} - x_k\|)$$

These conditions are also called Dennis-Moré conditions. They show that for superlinear convergence it is only important that  $\nabla^2 f(x_k)(x_{k+1} - x_k)$  and  $H_k(x_{k+1} - x_k)$  match sufficiently well. It is therefore not necessary that  $H_k$  approximates the entire Hessian matrix  $\nabla^2 f(x_k)$  well.

## 1.2 QUASI-NEWTON METHODS

Quasi-Newton methods are a class of numerical methods for solving nonlinear minimization problems. As the name suggests, these are based on the Newton method, but attempt to minimize the computational effort. The class goes back to the physicist William Davidon of the Argonne National Laboratory, who developed the first algorithm in the mid 1950s.

For the Newton method, both the gradient and the Hessian are calculated in every iteration. Of course, we get useful information about curvature of our function from the Hessian, get local at least superlinear convergence and if we add a method for determining stepsizes, we even get global convergence. But there are arguments against the Newton method, mainly related to the calculation of the Hessian. For example the calculation could be too costly or not possible at all (which includes the case that the Hessian does not exist). Quasi-Newton methods follow the strategy of not calculating and instead approximating it. Henceforth we call the approximation of the Hessian matrix  $\nabla^2 f(x_k)$  used in each iteration  $H_k$ .

It is expected that the sequence  $\{H_k\}_k$  should possess positive-definiteness,  $d_k = -H_k^{-1}\nabla f(x_k)$  should be a descent direction and the resulting method should behave like Newton's method in terms of convergence. Of course the calculation should cost less.

Let  $f: D \rightarrow \mathbb{R}$  be twice continuously differentiable on an open subset  $D \subset \mathbb{R}^n$ . We consider the quadratic Taylor-approximation of  $f$  at  $x_{k+1}$ :

$$f(x) \approx m_{k+1}(x) = f(x_{k+1}) + g_{k+1}^T(x - x_{k+1}) + \frac{1}{2}(x - x_{k+1})^T G_{k+1}(x - x_{k+1})$$

where  $g_{k+1} \triangleq \nabla f(x_{k+1})$  and  $G_{k+1} \triangleq \nabla^2 f(x_{k+1})$ . For the gradient we obtain

$$\nabla f(x) \approx \nabla m_{k+1}(x) = g_{k+1} + G_{k+1}(x - x_{k+1}).$$

Setting  $x = x_k$ ,  $s_k = x_{k+1} - x_k$  and  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k) = g_{k+1} - g_k$ , we get

$$G_{k+1}^{-1}y_k \approx s_k.$$

This holds with equality, if  $f$  is a quadratic function.

Now one is interested in the fact that an approximation of the Hessian inverse  $B_{k+1}$  satisfies this relation for the quasi-Newton method, i.e.,

$$B_{k+1}y_k = s_k \quad (1.2.1)$$

which is called the quasi-Newton equation, quasi-Newton condition or secant equation. A method that uses this condition to generate its symmetric Hessian (inverse) approximations is called a quasi-Newton method.

For quasi-Newton methods we replace the Hessian of our objective function  $\nabla^2 f(x_{k+1})$  in the model by an approximation  $H_{k+1}$ :

$$m_{k+1}(x) = f(x_{k+1}) + g_{k+1}^T(x - x_{k+1}) + \frac{1}{2}(x - x_{k+1})^T H_{k+1}(x - x_{k+1})$$

which satisfies the interpolation conditions:

$$m_{k+1}(x_{k+1}) = f(x_{k+1}) \quad \text{and} \quad \nabla m_{k+1}(x_{k+1}) = \nabla f(x_{k+1}).$$

Unlike the normal Newton method, in which we require that  $\nabla^2 m(x_{k+1}) = G_{k+1}$ , we want the model to satisfy

$$\nabla m_{k+1}(x_k) = g_k$$

from which follows

$$g_k = g_{k+1} + H_{k+1}(x_k - x_{k+1}).$$

So we have

$$H_{k+1}(x_{k+1} - x_k) = g_{k+1} - g_k \quad \Leftrightarrow \quad H_{k+1}s_k = y_k. \quad (1.2.2)$$

This is also called the quasi-Newton equation, quasi-Newton condition or secant equation but now expressed with the approximation of the Hessian.

Of course, we see immediately that the following relationship holds

$$H_k = B_k^{-1} \quad \text{for all } k \in \mathbb{N}_0 \quad (1.2.3)$$

and vice versa [Sun, Yuan, 2006](#).

We now come to the definition of a single iteration. A model at the current iteration  $x_k$  is defined:

$$f(x) \approx m_k(x) = f(x_k) + g_k^T(x - x_k) + \frac{1}{2}(x - x_k)^T H_k(x - x_k)$$

Assuming that  $H_k$  is positive definite, we get a quadratic convex model  $m_k$ . The minimizer  $d_k$  of it, we can write explicitly as

$$d_k = -H_k^{-1}g_k = -B_k g_k \quad (1.2.4)$$

is used as the search direction, and the new iterate is

$$x_{k+1} = x_k + \alpha_k d_k \quad (1.2.5)$$

where the stepsize  $\alpha_k$  is chosen to satisfy the Wolfe conditions. This iteration is quite similar to the line search Newton method, also called globalized Newton method. The key difference is that the approximate Hessian  $H_k$  is used in place of the true Hessian  $G_k = \nabla^2 f(x_k)$ , as already mentioned. The quasi-Newton equation requires that the symmetric positive definite matrix  $H_{k+1}$  maps  $s_k$  to  $y_k$ . This will be possible only if  $s_k$  and  $y_k$  satisfy the

$$s_k^T y_k > 0. \quad (1.2.6)$$

In the following we will refer to Eq. (1.2.6) as curvature condition. This follows from multiplying the quasi-Newton equation by  $s_k^T$  from the left, because we assume that  $H_{k+1}$  is positive definite. If the function  $f$  is strongly convex, then this inequality will be satisfied for any two points  $x_k$  and  $x_{k+1}$ . For nonconvex functions this condition will not always hold. In this scenario we have to impose restrictions on the line search procedure that chooses the stepsize  $\alpha_k$ . The curvature condition holds if we impose the (strong) Wolfe conditions on the line search. Setting  $s_k = x_{k+1} - x_k = \alpha_k d_k$  and using Eq. (1.1.3) leads to:

$$y_k^T s_k \geq (c_2 - 1)\alpha_k g_k^T d_k. \quad (1.2.7)$$

Since  $c_2 < 1$  and  $d_k$  is a descent direction, the right side is positive and the curvature condition holds. This also shows us the reason why Eq. (1.1.3) is also called curvature condition, since it guarantees Eq. (1.2.6). When the curvature condition is satisfied, the quasi-Newton equation has always a solution  $H_{k+1}$ . In fact, it admits an infinite number of solutions, since the  $n(n+1)/2$  degrees of freedom in a symmetric positive definite matrix exceed the conditions imposed by the quasi-Newton equation. The requirement of positive-definiteness imposes  $n$  additional inequalities - all principal minors must be positive - but these conditions do not absorb the remaining degrees of freedom Nocedal, Wright, 2006. Another indication that the matrices  $H_{k+1}$  satisfy the quasi-Newton equation is given by Corollary 1.1.8. Necessary and sufficient for the superlinear convergence of the sequence  $\{x_k\}_k$  to a minimizer  $x^*$  is the condition:

quasi-Newton method	Newton method
Only need the function values and gradients $\{H_k\}_k$ maintains positive definite for several updates	Need the function values, gradients and Hessians $\{G_k\}_k$ is not sure to be positive definite
Need $O(n^2)$ multiplications in each iteration	Need $O(n^3)$ multiplications in each iteration

Tabelle 1.2.1: Comparison

$$\|(G_k - H_k)(x_{k+1} - x_k)\| = o(\|x_{k+1} - x_k\|). \quad (1.2.8)$$

It can be shown that (Eq. (1.2.8)) is equivalent to

$$\|g_{k+1} - g_k - H_k(x_{k+1} - x_k)\| = o(\|x_{k+1} - x_k\|).$$

This motivates the following requirement on  $H_{k+1}$ :

$$H_{k+1}(x_{k+1} - x_k) = g_{k+1} - g_k. \quad (1.2.9)$$

We see immediately that this is the quasi-Newton equation (Eq. (1.2.2)) Geiger, Kanzow, 1999. The current theory is nevertheless sufficient to formulate a general algorithm.

---

**Algorithm 3** General Quasi-Newton Method

---

```

1:  $x_0 \in \mathbb{R}^n$ ,  $B_0 \in \mathbb{R}^{n \times n}$  approximation of  $\nabla^2 f(x_0)^{-1}$ ,  $0 \leq \epsilon < 1$ ,  $k = 0$ 
2: while  $\|\nabla f(x_k)\| > \epsilon$  do
3:   Compute  $d_k = -B_k \nabla f(x_k)$ .
4:   Determine the stepsize  $\alpha_k > 0$  by line search.
5:   Set  $x_{k+1} = x_k + \alpha_k d_k$ .
6:   Update  $B_k$  into  $B_{k+1}$  such that the quasi-Newton equation holds.
7:   Set  $k = k + 1$ .
8: end while
9: return  $x_k$ 

```

---

One commonly starts the algorithm with  $B_0 = I$ , the identity matrix or set  $B_0$  to be a finite-difference approximation to the inverse Hessian  $\nabla^2 f(x_0)^{-1}$ . If  $B_0 = I$ , the first iteration is just a steepest descent iteration. In some cases one uses the direct approximation  $H_k$  of the Hessian. In this case we need to solve a system of equations in step 3 to get  $d_k$  and we need to update  $H_k$  instead of  $B_k$ . However, since one generally wants to do without solving a system of equations, this variant is not recommended. The resulting advantages of the quasi-Newton method over the ordinary Newton method are shown in the Table Table 1.2.1.

As Newton's method is a steepest descent method under the norm  $\|\cdot\|_{G_k}$ , the quasi-Newton method is a steepest descent method under the norm  $\|\cdot\|_{H_k}$ . In fact,  $d_k$  is the solution of the minimization problem

$$\begin{aligned} \min \quad & g_k^T d \\ \text{s.t.} \quad & \|d\|_{H_k} \leq 1. \end{aligned} \tag{1.2.10}$$

It follows from

$$(g_k^T d)^2 \leq (g_k^T H_k^{-1} g_k)(d^T H_k d)$$

that the solution of (Eq. (1.2.10)) is

$$d_k = -H_k^{-1} g_k = -B_k g_k,$$

and  $g_k^T d_k$  is the smallest value. By the way, since the metric matrices  $H_k$  are positive definite and always changed from iteration to iteration, the method is also called the variable metric method [Sun, Yuan, 2006](#).

### 1.3 THE BROYDEN-FLETCHER-GOLDFARB-SHANNO FORMULA

We have seen that the search direction in a quasi-Newton method is given by

$$d_k = -B_k g_k = -H_k^{-1} g_k$$

and the new iterate is

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

This iteration is quite similar to the one of Newton's method. The key difference is that the approximate Hessian  $H_k$  is used in place of the true Hessian  $\nabla^2 f(x_k)$ . Instead of computing  $H_k$  afresh at every iteration, Davidon proposed to update it in a simple manner to account for the curvature measured during the most recent step [Nocedal, Wright, 2006](#), p. 137. The question now is how the matrix  $H_{k+1}$  (or  $B_{k+1}$ ) should be constructed from  $H_k$  (or  $B_k$ ) and other information. Various formulae have been developed for this, some of which are interrelated. In this thesis the main focus is on the Broyden-Fletcher-Goldfarb-Shanno formula or short BFGS formula, which has proven to be the most efficient quasi-Newton method in practice [Ulbrich, Ulbrich, 2012](#), p. 69. However, all approaches follow the following three important guidelines to create  $H_{k+1}$ :

- (i)  $H_{k+1}$  satisfies the quasi-Newton equation [Eq. \(1.2.2\)](#).
- (ii)  $H_{k+1}$  is symmetric and positive definite.

(iii)  $H_{k+1}$  is “near”  $H_k$ .

Of course these three characteristics should also hold for the approximation of the inverse  $B_{k+1}$ . In the previous subsection was shown that  $H_{k+1}$  should satisfy the quasi-Newton equation [Eq. \(1.2.2\)](#). The strongest motivation comes from the fact that we approximate our objective function local by a quadratic model and the Hessian of a quadratic function always satisfies the quasi-Newton equation. The fact that the distance between  $H_{k+1}$  and  $H_k$  should not be too large will be related to the rate of convergence of the resulting method and the uniqueness of the formula. It's obvious that the matrix  $H_{k+1}$  should be symmetric, since we want to approximate the Hessian and the Hessian is always symmetric in the case of a twice continuously differentiable function  $f \in C^2$ . We need positive-definiteness for efficiency, numerical stability and global convergence. If the Hessian  $\nabla^2 f(x^*)$  is positive definite, the stationary point  $x^*$  is a strong minimizer. Hence, we hope the Hessian approximations  $\{H_k\}_k$  (or inverse Hessian approximations  $\{B_k\}_k$ ) are positive definite. In addition, if  $H_k$  (or  $B_k$ ) is positive definite, the local quadratic model of  $f$  has a unique local minimizer, and the direction  $d_k$  is a descent direction [Sun, Yuan, 2006](#), p. 212.

Before we get to the BFGS formula, also called BFGS update, let us first look at another one. By exchanging variables, we then get the BFGS formula. It's the so called DFP update, proposed by Davidon in [Davidon, 1959](#) and developed later by Fletcher and Powell in [Fletcher, Powell, 1963](#). We assume that the matrix  $B_k$  approximates  $\nabla^2 f(x_k)^{-1}$  sufficiently well. Let us consider a symmetric rank-two update, that means we add two symmetric rank-one matrices to the current matrix

$$B_{k+1} = B_k + auu^T + bvv^T$$

where  $u, v \in \mathbb{R}^n$ ,  $a, b \in \mathbb{R}$  are to be determined. From the quasi-Newton equation follows

$$B_{k+1}y_k = B_k y_k + auu^T y_k + bvv^T y_k = s_k.$$

Clearly,  $u$  and  $v$  can not uniquely be determined. One possible choice is

$$u = s_k, \quad v = B_k y_k.$$

Hence we obtain

$$a = \frac{1}{u^T y_k} = \frac{1}{s_k^T y_k}, \quad b = -\frac{1}{v^T y_k} = -\frac{1}{y_k^T B_k y_k}.$$

Therefore

$$B_{k+1}^{DFP} = B_k^{DFP} + \frac{s_k s_k^T}{s_k^T y_k} - \frac{B_k^{DFP} y_k y_k^T B_k^{DFP}}{y_k^T B_k^{DFP} y_k}.$$

This is the DFP update, which approximates the inverse of the Hessian  $\nabla^2 f(x_k)^{-1}$  in every iteration [Sun, Yuan, 2006](#), p. 210.

The last two terms in the right-hand-side are symmetric rank-one matrices. This is the fundamental idea of quasi-Newton updating: Instead of recomputing the approximate Hessian (or inverse Hessian) from scratch at every iteration, we apply a simple modification that combines the most recently observed information about the objective function with the existing knowledge embedded in our current Hessian approximation [Nocedal, Wright, 2006](#), p. 139.

The BFGS formula can be obtained by simple trick: for  $H_{k+1}^{BFGS}$  replace the triple  $(B_k^{DFP}, s_k, y_k)$  in  $B_{k+1}^{DFP}$  by  $(H_k^{BFGS}, y_k, s_k)$ . Thus, BFGS update is also said to be a complement DFP update. The result is

$$H_{k+1}^{BFGS} = H_k^{BFGS} + \frac{y_k y_k^T}{s_k^T y_k} - \frac{H_k^{BFGS} s_k s_k^T H_k^{BFGS}}{s_k^T H_k^{BFGS} s_k}. \quad (1.3.1)$$

[Sun, Yuan, 2006](#)

This formula was discovered independently by Broyden in [Broyden, 1967](#), by Fletcher in [Fletcher, 1970](#), by Goldfarb in [Goldfarb, 1970](#) and by Shanno in [Shanno, 1970](#), which is the reason for the name. All four authors derive the BFGS formula in a slightly different way, which can be seen as a reason why it is superior to the other updating formulae in practice [Geiger, Kanzow, 1999](#), p. 136. It is presently considered to be the most effective of all quasi-Newton updating formulae [Nocedal, Wright, 2006](#), p. 139. The DFP update is quite effective, but it was soon superseded by the BFGS formula, which has all good properties of the DFP update [Sun, Yuan, 2006](#), p. 219.

Since  $H_k s_k = -\alpha_k g_k$  and  $H_k d_k = -g_k$ , this formula can also be written as

$$H_{k+1}^{BFGS} = H_k^{BFGS} + \frac{g_k g_k^T}{g_k^T d_k} + \frac{y_k y_k^T}{\alpha_k y_k^T d_k}.$$

By applying the Sherman–Morrison–Woodbury, [Theorem 1.1.5](#), formula twice to [Eq. \(1.3.1\)](#), we obtain

$$\begin{aligned} B_{k+1}^{BFGS} &= B_k^{BFGS} + \frac{(s_k - B_k^{BFGS} y_k) s_k^T + s_k (s_k - B_k^{BFGS} y_k)^T}{s_k^T y_k} - \frac{(s_k - B_k^{BFGS} y_k)^T y_k s_k s_k^T}{(s_k^T y_k)^2} \\ &= B_k^{BFGS} + \left( I + \frac{y_k^T B_k^{BFGS} y_k}{s_k^T y_k} \right) \frac{s_k s_k^T}{s_k^T y_k} - \frac{s_k y_k^T B_k^{BFGS} + B_k^{BFGS} y_k s_k^T}{s_k^T y_k} \\ &= \left( I - \frac{s_k y_k^T}{s_k^T y_k} \right) B_k^{BFGS} \left( I - \frac{y_k s_k^T}{s_k^T y_k} \right) + \frac{s_k s_k^T}{s_k^T y_k}. \end{aligned} \quad (1.3.2)$$

These are different formulae for the approximation of the Hessian inverse. Furthermore, reference is also made to this with “BFGS formula”. It is easy to see that [Eq. \(1.3.2\)](#) is also a rank-two modification of  $B_k^{BFGS}$ . One can easily show that

$$H_{k+1}^{BFGS} B_{k+1}^{BFGS} = B_{k+1}^{BFGS} H_{k+1}^{BFGS} = I.$$

Replacing the triple  $(B_k^{BFGS}, s_k, y_k)$  in Eq. (1.3.2) by  $(H_k^{DFP}, y_k, s_k)$ , one would get a formula for  $H_{k+1}^{DFP}$ , the direct DFP update. This describes a method for finding its dual update from a given update [Sun, Yuan, 2006](#), p.218. For this reason, the DFP and BFGS formulae are sometimes referred to as “dual” updating formulae.

Now it is checked that  $H_{k+1}^{BFGS}$  meets the given characteristics. For 1. and 2. there is the following statement:

**Satz 1.3.1** ([Ulbrich, Ulbrich, 2012](#), Theorem 13.4). (i) If  $y_k^T s_k \neq 0$  and  $s_k^T H_k^{BFGS} s_k \neq 0$  holds, the matrices  $H_{k+1}^{BFGS} \in \mathbb{R}^{n \times n}$  are well defined, symmetric and satisfy the quasi-Newton equation (Eq. (1.2.2)). (ii) If  $H_k^{BFGS}$  is positive definite and  $y_k^T s_k > 0$ , then  $H_{k+1}^{BFGS}$  is positive definite.

Such an update is also called positive definite update. The same holds of course for the approximation of the inverse  $B_{k+1}^{BFGS}$ .

In the previous subsection was shown that the curvature condition, Eq. (1.2.6), must hold. This was achieved by imposing restrictions on the line search method, see Eq. (1.2.7). So the positive-definiteness can be guaranteed just by a Wolfe line search strategy.

The statement was actually made for Broyden class matrices (the matrices are a convex combination of DFP and BFGS matrices) in [Ulbrich, Ulbrich, 2012](#), p. 69, this means that it can be transferred one-to-one to the DFP update  $H_{k+1}^{DFP}$ .

The last characteristic, that  $H_{k+1}$  should be “near”  $H_k$ , has a far more powerful meaning than the other two. Many authors use only this to define the BFGS formula which of course is perfectly legitimate. As already mentioned, this property leads to the fact that the formula can be considered as unique and it has something to do with the rate of convergence. The two go hand in hand.

One wants the BFGS method to be similar to Newton’s method in terms of convergence. This means that it should converge superlinearly. This can be proven by the Dennis-Moré condition. The connection between this condition and the characteristic is shown by

**Lemma 1.3.2** ([Ulbrich, Ulbrich, 2012](#), Lemma 13.2).  $x^*$  fulfills the sufficient condition of second order. If Algorithm 3 with  $\alpha_k = 1$  for all  $k \in \mathbb{N}$  generates a sequence  $\{x_k\}_k$  convergent to  $x^*$  and also holds

$$\lim_{k \rightarrow \infty} \|H_{k+1} - H_k\| = 0,$$

then  $H_k$  satisfies the Dennis-Moré condition, [Corollary 1.1.8](#), and  $\{x_k\}_k$  converges superlinear to  $x^*$ .

Therefore one looks for quasi-Newton updates for which  $H_{k+1}$  is close to  $H_k$  in each iteration, so that the distance between them converges towards zero.

We would now like to consider the third property ( $H_{k+1}$  is “near”  $H_k$ ) from the point of view of the uniqueness of the formula. One obtains this by considering the formula for  $B_{k+1}^{BFGS}$  as the solution to



an optimization problem. The following two statements provide us with the uniqueness of the BFGS formula:

**Lemma 1.3.3** (Geiger, Kanzow, 1999, Lemma 11.7). *Let be  $s \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^n$  with  $y \neq 0$  and a symmetric matrix  $B \in \mathbb{R}^{n \times n}$  given. Furthermore let  $W \in \mathbb{R}^{n \times n}$  be symmetric and positive definite. Then the unique solution of the inverse weighted problem*

$$\begin{aligned} \min_{B_+} \quad & \|W(B_+ - B)W\|_F^2 \\ \text{s.t.} \quad & B_+ = B_+^T, \quad B_+ y = s \end{aligned} \quad (1.3.3)$$

is given by

$$B_+^W = B + \frac{(s - By)(W^{-2}y)^T + W^{-2}y(s - By)^T}{(W^{-2}y)^T y} - \frac{y^T(s - By)W^{-2}y(W^{-2}y)^T}{((W^{-2}y)^T y)^2}.$$

In order to be able to convert this to the BFGS formula one is looking for the so-called weighting matrix  $W$ :

**Satz 1.3.4** (Geiger, Kanzow, 1999, Theorem 11.8). *Let  $B \in \mathbb{R}^{n \times n}$  be symmetric and positive definite and  $s, y \in \mathbb{R}^n$  with  $s^T y > 0$ . Let  $Q \in \mathbb{R}^{n \times n}$  be a symmetric and positive definite matrix with  $Qs = y$ , and let  $W = Q^{\frac{1}{2}}$  be a square root of  $Q$ . Then the unique solution of the inverse weighted problem (Eq. (1.3.3)) with the weighted  $W$  is given by*

$$B_+^{BFGS} = B + \frac{(s - By)s^T + s(s - By)^T}{y^T s} - \frac{(s - By)^T y s s^T}{(y^T s)^2}. \quad (1.3.4)$$

One can choose  $Q = W^2 = \tilde{G}_k$  and  $\tilde{G}_k$  is the average Hessian, i.e.

$$\tilde{G}_k = \int_0^1 \nabla^2 f(x_k + \tau s_k) d\tau = \int_0^1 \nabla^2 f(x_k + \tau \alpha_k d_k) d\tau, \quad (1.3.5)$$

which is positive definite for a strong convex function. The property

$$y_k = \nabla f(x_{k+1}) - \nabla f(x_k) = \int_0^1 \nabla^2 f(x_k + \tau s_k) s_k d\tau = \tilde{G}_k s_k$$

follows from Taylor's formula [Sun, Yuan, 2006](#), p. 214. With this choice of weighting matrix  $W$ , the norm

$$\|A\|_{W^2} = \|WAW\|_F$$

is non-dimensional, which is a desirable property, since we do not wish the solution of [Eq. \(1.3.3\)](#) to depend on the units of the problem [Nocedal, Wright, 2006](#), p. 139. The existence of  $Q$  is shown in [Geiger, Kanzow, 1999](#), Lemma 11.5. and since it is symmetric and positive definite, the existence of a symmetric positive definite matrix  $W$  follows from it, see [Geiger, Kanzow, 1999](#), Satz B.6. The specified minimum characteristic with respect to the weighted norms mentioned in the theorem automatically ensures the invariance of the BFGS method under affine-linear variable transformations. This important characteristic is also present in the Newton method [Ulbrich, Ulbrich, 2012](#), p. 69. The initial approximation  $B_0^{BFGS}$  must still be discussed. Unfortunately, there is no perfect strategy for this yet. One possibility is to use information about the problem and approximate the Hessian inverse by finite differences at  $x_0$ . One could also use a multiple of the identity matrix  $\beta I$ , where  $\beta$  is a scaling factor for the variables. But to determine this factor is problematic. If  $\beta$  is too large, so that the first step  $d_0 = -\beta g_0$  is too long, many function evaluations may be required to find a suitable value for the stepsize  $\alpha_0$ . A quite effective heuristic is to scale the starting matrix after the first step, but before the first BFGS update is performed. The provisional value  $B_0^{BFGS} = I$  is changed by setting

$$B_0^{BFGS} = \frac{y_1^T s_1}{y_1^T y_1} I$$

before applying the update to obtain  $B_1^{BFGS}$ . This formula attempts to make the size of  $B_0^{BFGS}$  similar to that of  $\nabla^2 f(x_0)^{-1}$  [Nocedal, Wright, 2006](#), p. 140.

## 1.4 THE BFGS-METHOD

We present a globalized BFGS method. For that, it must be ensured that the curvature condition, [Eq. \(1.2.6\)](#), is fulfilled by choosing the stepsize accordingly. The globalization of the BFGS method is similar to the globalization of Newton's method. In contrast the Wolfe-Powell stepsize strategy, which ensures  $s_k^T y_k > 0$  for all  $k \in \mathbb{N}$ , is chosen and not the Armijo rule (a stepsize strategy which determines a stepsize  $\alpha_k$  satisfying [Eq. \(1.1.2\)](#)). This type of line search is called inexact line search, approximate line search or acceptable line search. The stepsize  $\alpha_k > 0$  is chosen such that the objective function has an acceptable descent amount, i.e., such that the descent  $f(x_k) - f(x_k + \alpha_k d_k) > 0$  is acceptable [Sun, Yuan, 2006](#), p. 71.

We call the following algorithm "Inverse Global BFGS Method" because the updating formula for the approximation of the inverse of the Hessian ( $B_k^{BFGS} \mapsto B_{k+1}^{BFGS}$ ) is used, since we are spared the solving of a system of equations and we only have to work with matrix-vector-multiplications. The algorithm could be formulated with the approximation of the actual Hessian  $H_k^{BFGS}$ , but that would increase the effort again to  $O(n^3)$ , which is not desirable [Nocedal, Wright, 2006](#), p. 141. In practice, it must be

decided whether solving a system of equations or matrix-vector-multiplication is more advantageous for the underlying problem. We assume in this thesis that the latter is the better choice.

---

**Algorithm 4** Inverse Global BFGS Method

---

```

1: Given starting point  $x_0 \in \mathbb{R}^n$ , convergence tolerance  $\epsilon > 0$ , an initial symmetric and positive
   definite matrix  $B_0^{BFGS} \in \mathbb{R}^{n \times n}$ ,  $k = 0$ .
2: while  $\|\nabla f(x_k)\| > \epsilon$  do
3:   Compute search direction  $d_k = -B_k^{BFGS} \nabla f(x_k)$ .
4:   Find a stepsize  $\alpha_k$  that satisfies the (strong) Wolfe conditions.
5:   Set  $x_{k+1} = x_k + \alpha_k d_k$ ,  $s_k = x_{k+1} - x_k$ ,  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$ .
6:   Compute  $B_{k+1}^{BFGS}$  by means of Eq. (1.3.2).
7:   Set  $k = k + 1$ .
8: end while
9: return  $x_k$ 

```

---

It can be shown that Algorithm 4 is well defined:

**Satz 1.4.1** (Geiger, Kanzow, 1999, Theorem 11.37). *Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be continuously differentiable and bounded from below. Then for the globalized BFGS method Algorithm 4:*

- (i)  $s_k^T y_k > 0$  for all  $k \in \mathbb{N}$ .
- (ii) The matrices  $B_{k+1}^{BFGS}$  are symmetric and positive definite for all  $k \in \mathbb{N}$ .
- (iii) The method is well defined.

The derivation to this theorem in Geiger, Kanzow, 1999 shows that finding a stepsize  $\alpha_k$ , which satisfies the Wolfe conditions (Eq. (1.1.2) and Eq. (1.1.3)) or strong Wolfe conditions (Eq. (1.1.2) and Eq. (1.1.4)), is crucial Geiger, Kanzow, 1999, p. 166. Eq. (1.2.7) shows that the stepsize ensures the curvature condition  $s_k^T y_k > 0$ . This in turn ensures that the positive-definiteness of the matrix  $B_k^{BFGS}$  is passed to  $B_{k+1}^{BFGS}$ , see Theorem 1.3.1.

Let us now turn to the convergence analysis of Algorithm 4. It is desirable that each limit point of a sequence  $\{x_k\}_k$  generated by Algorithm 4 is a stationary point of  $f$  and that we get locally superlinear convergence. Unfortunately, neither of these statements is true in general Geiger, Kanzow, 1999, p. 167. Let us first deal with the global convergence. The difficulty and importance of the convergence problem of whether the BFGS method with the Wolfe line search converges globally for general functions has been addressed in many situations. But recent studies provide a negative answer to it for nonconvex functions (see e.g. Dai, 2002 or Mascarenhas, 2004) Dai, 2012, p. 3.

We present a statement about global convergence from Nocedal, Wright, 2006, which is based on results of Powell, 1976 and is also very common. For that we require that the objective function is convex. To be more precise, the following assumptions must be made for a reasonable convergence statement:

**Annahme 1.4.2** (Nocedal, Wright, 2006, Assumption 6.1.).

- (i) The objective function  $f$  is twice continuously differentiable.

(ii) The level set  $\mathcal{L} = \{x \in \mathbb{R}^n : f(x) \leq f(x_0)\}$  is convex, and there exist positive constants  $m$  and  $M$  such that

$$m\|z\|^2 \leq z^T \nabla^2 f(x) z \leq M\|z\|^2$$

for all  $z \in \mathbb{R}^n$  and  $x \in \mathcal{L}$ .

The second part of [Annahme 1.4.2](#) implies that  $\nabla^2 f(x)$  is positive definite on  $\mathcal{L}$  and that  $f$  has an unique minimizer  $x^*$  in  $\mathcal{L}$  [Nocedal, Wright, 2006](#), p. 153.

**Satz 1.4.3** ([Nocedal, Wright, 2006](#), Theorem 6.5.). Let  $H_0$  be any symmetric positive definite initial matrix, and let  $x_0$  be a starting point for which [Annahme 1.4.2](#) is satisfied. Then the sequence  $\{x_k\}_k$  generated by [Algorithm 4](#) (with  $\epsilon = 0$ ) converges to the minimizer  $x^*$  of  $f$ .

The proof can be found in [Nocedal, Wright, 2006](#), p. 154. Nevertheless, it should be noted that [Theorem 1.1.6](#), the Zoutendijk condition, is crucial to prove global convergence.

We see that the BFGS Method with the Wolfe-Powell stepsize strategy applied to a smooth convex function  $f$  from an arbitrary starting point  $x_0 \in \mathbb{R}^n$  and from any initial approximation  $B_0^{BFGS} \in \mathbb{R}^{n \times n}$ , that is symmetric and positive definite, is globally convergent. This is a very strong convergence result for the BFGS method, and it is currently not known whether this also applies to the DFP method [Nocedal, Wright, 2006](#), p. 156.

It can also be shown that [Algorithm 4](#) not only converges globally by using the Wolfe-Powell stepsize strategy but also by using a great number of inexact, efficient stepsize strategies (see [Definition 1.1.3](#)) used in practice for uniformly convex objective functions, which can be seen as an indication of the numerical stability of this method [Werner, 1978/79](#), p. 327.

The local BFGS method (a variant where the stepsize is always equal one, i.e.  $\alpha_k = 1$ ,  $\forall k$ ) achieves superlinear convergence to a stationary point second order [Geiger, Kanzow, 1999](#), Satz 11.33. The crux, however, is that this stepsize must satisfy the Wolfe conditions to be accepted and thus we get superlinear convergence for [Algorithm 4](#).

In practical implementations [Algorithm 4](#) this unit stepsize is usually used as the first trial stepsize. Under suitable assumptions on  $f$ , this stepsize will be accepted by the line search as the iterates tend to the solution and will enable superlinear convergence [Dai, 2012](#), p. 6. In [Dennis, Moré, 1974](#) this idea is pursued and a detailed derivation is given. But we want to present a very well known result from [Nocedal, Wright, 2006](#), which applies to general objective functions and is also based on the results of [Dennis, Moré, 1974](#).

**Annahme 1.4.4** ([Nocedal, Wright, 2006](#), Assumption 6.2.). The Hessian matrix is Lipschitz continuous at  $x^*$ , that is,

$$\|\nabla^2 f(x) - \nabla^2 f(x^*)\| \leq L\|x - x^*\|,$$

for all  $x$  near  $x^*$ , where  $L$  is a positive constant.

**Satz 1.4.5** ([Nocedal, Wright, 2006](#), Theorem 6.6.). Suppose that  $f$  is twice continuously differentiable

and that the iterates  $\{x_k\}_k$  generated by [Algorithm 4](#) converge to a minimizer  $x^*$  at which [Annahme 1.4.4](#) holds. Suppose also that

$$\sum_{k=0}^{\infty} \|x - x^*\| < \infty \quad (1.4.1)$$

holds. Then  $\{x_k\}_k$  converges to  $x^*$  at a superlinear rate.

The analysis of the theorem's proof makes use of the [Theorem 1.1.7](#) characterization of superlinear convergence.

In practice it is often observed that the sequence of matrices  $\{H_k^{BFGS}\}_k$  generated by [Algorithm 4](#) using the direct formula [Eq. \(1.3.1\)](#) converges to  $\nabla^2 f(x^*)$ . In this case, the superlinear convergence follows directly from [Eq. \(1.1.9\)](#).

## 1.5 A CAUTIOUS BFGS-METHOD

In [Dai, 2012](#) a four-dimensional example where the objective function is smooth (polynomial) and nonconvex was presented such that the globalized BFGS method does not converge. From this it can be concluded that [Algorithm 4](#) unfortunately does not converge in general. In this subsection we present a simple modification of the globalized BFGS method, which allows us to establish a global convergence theorem for nonconvex problems.

We present the method from [Li, Fukushima, 2001](#), which modifies [Algorithm 4](#) by a so-called cautious update. It can be shown that this method with a stepsize strategy, which satisfies the Wolfe conditions, converges globally if the objective function has Lipschitz continuous gradients. Moreover, under appropriate conditions, it can be shown that the cautious update eventually reduces to the ordinary update, i.e.  $B_k^{BFGS}$  is updated using [Eq. \(1.3.2\)](#).

The idea is to add a simple “check” which decides whether the matrix  $B_{k+1}^{CBFGS}$  is updated. To be precise, the following update is now used in each iteration:

$$B_{k+1}^{CBFGS} = \begin{cases} \text{using Eq. (1.3.2)} & \frac{y_k^T s_k}{\|s_k\|^2} \geq \mu \|\nabla f(x_k)\|^\lambda \\ B_k^{CBFGS} & \text{otherwise.} \end{cases} \quad (1.5.1)$$

where  $\mu$  and  $\lambda$  are positive constants.

The only requirement for the matrix  $B_{k+1}^{CBFGS}$  to be updated is that the inner product  $y_k^T s_k$  must be greater than or equal to a positive value depending on the gradient. That means, we now demand that the iteration fulfills a “little bit” more than just the curvature condition, [Eq. \(1.2.6\)](#).

This modification is motivated by the fact that it can be ensured that even stepsize strategies which do not fulfil [Eq. \(1.1.3\)](#) (or [Eq. \(1.1.4\)](#)) generate positive definite updates, since it is ensured that  $B_{k+1}^{CBFGS}$  is updated if [Eq. \(1.2.6\)](#) holds, which in turn implies that the positive-definiteness is inherited (see

[Theorem 1.3.1](#)), and that convenient convergence results can now also be established for nonconvex functions. We are interested in the latter because we continue to assume that a stepsize strategy is used that fulfills [Eq. \(1.1.3\)](#) (or [Eq. \(1.1.4\)](#)) and therefore [Eq. \(1.2.6\)](#) is always fulfilled. We get the following algorithm:

---

**Algorithm 5** Cautious BFGS Algorithm

---

```

1: Given starting point  $x_0 \in \mathbb{R}^n$ , convergence tolerance  $\epsilon > 0$ , an initial symmetric and positive
   definite matrix  $B_0^{CBFGS} \in \mathbb{R}^{n \times n}$ , choose constants  $\lambda > 0$  and  $\mu > 0$ ,  $k = 0$ .
2: while  $\|\text{grad } f(x_k)\| > \epsilon$  do
3:   Compute search direction  $d_k = -B_k^{CBFGS} \nabla f(x_k)$ .
4:   Find a stepsize  $\alpha_k$  that satisfies the Wolfe conditions.
5:   Set  $x_{k+1} = x_k + \alpha_k d_k$ ,  $s_k = x_{k+1} - x_k$ ,  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$ .
6:   Compute  $B_{k+1}^{CBFGS}$  by means of Eq. \(1.5.1\).
7:   Set  $k = k + 1$ .
8: end while
9: return  $x_k$ 

```

---

By this simple extension of [Eq. \(1.5.1\)](#) for [Algorithm 4](#), one gets of course other convergence results. Fortunately, these now concern not only convex functions and therefore have weaker assumptions:

**Annahme 1.5.1** (Li, Fukushima, 2001, Assumption A). *The level set*

$$\mathcal{L} = \{x \in \mathbb{R}^n : f(x) \leq f(x_0)\}$$

*is bounded, the function  $f$  is continuously differentiable on  $\mathcal{L}$ , and there exists a constant  $L > 0$  such that*

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathcal{L}. \quad (1.5.2)$$

We see that there is no longer a requirement that the sublevel set  $\mathcal{L}$  must be convex. Instead the assumption, that the function  $f$  is Lipschitz continuously differentiable, was added. Since  $\{f(x_k)\}_k$  is a decreasing sequence, it follows that the sequence  $\{x_k\}_k$  generated by [Algorithm 5](#) is contained in  $\mathcal{L}$ . We have the following statement on global convergence:

**Satz 1.5.2** (Li, Fukushima, 2001, Theorem 3.3.). *Let [Annahme 1.5.1](#) hold and  $\{x_k\}_k$  be generated by [Algorithm 5](#). Then*

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

*holds.*

**Theorem 1.5.2** shows that there exists a subsequence of  $\{x_k\}_k$  converging to a stationary point  $x^*$  of Eq. (1.1.1). If  $f$  is convex, then  $x^*$  is a global minimum of  $f$ . Since the sequence  $\{f(x_k)\}_k$  converges, it is clear that every cluster point of  $\{x_k\}_k$  is a global optimal solution of Eq. (1.1.1).

**Folgerung 1.5.3** (Li, Fukushima, 2001, Corollary 3.4.). *Let Annahme 1.5.1 hold and  $\{x_k\}_k$  be generated by Algorithm 5. If  $f$  is convex, then the whole sequence  $\{\nabla f(x_k)\}_k$  converges to zero. Consequently, every accumulation point of  $\{x_k\}_k$  is a global optimal solution.*

This can be generalized to nonconvex functions as follows:

**Satz 1.5.4** (Li, Fukushima, 2001, Theorem 3.5.). *Let  $f$  be twice continuously differentiable. Suppose that  $\lim_{k \rightarrow \infty} s_k = 0$ . If there exists an accumulation point  $x^*$  of  $\{x_k\}_k$  at which  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive definite, then the whole sequence  $\{x_k\}_k$  converges to  $x^*$ . If in addition,  $\nabla^2 f(\cdot)$  is Hölder continuous and  $c_1 \in (0, 0.5)$  holds, then the convergence rate is superlinear.*

The assumptions imply that  $x^*$  is a strict local optimal solution. In the proof is shown that Algorithm 5 is moving to Algorithm 4, that means when  $k$  is sufficiently large, the condition  $y_k^T s_k / \|s_k\|^2 \geq \mu \|\nabla f(x_k)\|^\lambda$  is always satisfied, which implies that the algorithm reduces to the “ordinary” BFGS method. The superlinear convergence then follows from the results of the convergence analysis of Algorithm 4.

In Li, Fukushima, 2001 it is also noted that the parameters need not be constant while the method is running. By choosing  $\lambda = 0.01$  if  $\|\nabla f(x_k)\| \geq 1$  and  $\lambda = 3$  if  $\|\nabla f(x_k)\| < 1$  one tries to make the cautious update Eq. (1.5.1) closer to the original BFGS update Eq. (1.3.2). Another option is that  $\lambda$  can also be within an interval  $[a, b]$  with  $a > 0$ . More generally, the value  $\mu \|\nabla f(x_k)\|^\lambda$  can be replaced by a general forcing function  $\theta(\|\nabla f(x_k)\|)$ , which is strictly monotone with  $\theta(0) = 0$ . For all these variants of parameter selection or adjustments of the cautious trigger, the convergence results hold.

## 1.6 LIMITED-MEMORY BFGS-METHOD

One of the disadvantages of the Quasi-Newton methods is that a  $n \times n$  matrix (namely  $B_{k+1}^{BFGS}$ ) must be stored in each iteration. Even when using the symmetry of this matrix, a memory requirement of  $n(n+1)/2$  matrix entries remains. For large-scale optimization problems is this not feasible Geiger, Kanzow, 1999, p. 197.

Limited-memory quasi-Newton methods, also called variable-storage quasi-Newton methods, are useful for solving large problems whose Hessian matrices cannot be computed at a reasonable cost or are not sparse. The methods save only a few  $n$ -dimensional vectors, instead of storing and computing fully dense  $n \times n$  approximations of the Hessian. The main idea is to use the curvature information (the information about the curvature of the model, which is obtained by approximating the Hessian) from only the most recent iterations to construct the Hessian approximation. Curvature information from earlier iterations, which is less likely to be relevant to the actual behavior of the Hessian at the current iteration, is discarded in the interest of saving storage. Their rate of convergence is often acceptable (albeit linear), despite these modest storage requirements Nocedal, Wright, 2006.



Due to the outstanding importance of the BFGS-method in the class of quasi-Newton methods [Geiger, Kanzow, 1999](#), p. 197, it is also predominantly used as a limited-memory variant, called L-BFGS. But there are also limited-memory versions of other quasi-Newton methods such as the Symmetric Rank-One (SR1) method [Nocedal, Wright, 2006](#), p. 177.

In subsection 1.3 three different formulae for the approximation of the Hessian inverse were introduced, see [Eq. \(1.3.2\)](#). We start from the last one

$$B_{k+1}^{BFGS} = \left( I - \frac{s_k y_k^T}{s_k^T y_k} \right) B_k^{BFGS} \left( I - \frac{y_k s_k^T}{s_k^T y_k} \right) + \frac{s_k s_k^T}{s_k^T y_k}.$$

For given vectors  $s_k, y_k \in \mathbb{R}^n$  with  $s_k^T y_k > 0$  one sets

$$\rho_k = \frac{1}{s_k^T y_k}, \quad V_k = I - \rho_k y_k s_k^T, \quad (1.6.1)$$

obtaining

$$B_{k+1}^{BFGS} = V_k^T B_k^{BFGS} V_k + \rho_k s_k s_k^T. \quad (1.6.2)$$

The matrix  $B_{k+1}^{BFGS}$  is obtained by updating  $B_k^{BFGS}$  using the pair  $\{s_k, y_k\}$  [Sun, Yuan, 2006](#), p. 293. Since the inverse Hessian approximation  $B_k^{BFGS}$  will generally be dense, the cost of storing and manipulating it is prohibitive when the number of variables is large. To circumvent this problem, one stores a modified version of  $B_k^{BFGS}$  implicitly, by storing a certain number (say,  $m$ ) of the vector pairs  $\{s_i, y_i\}$  [Nocedal, Wright, 2006](#), p. 177. After the new iterate  $x_{k+1}$  is computed, the oldest vector pair in the set  $\{s_i, y_i\}_{i=k-m}^{k-1}$  (namely  $\{s_{k-m}, y_{k-m}\}$ ) is discarded and the new pair  $\{s_k, y_k\}$  obtained from the current step is added. This works according to the strategy “first in, first out”, which means that the vectors that were imported first are also discarded first. In this way, the set of vector pairs includes curvature information from the  $m$  most recent iterations. Practical experience has shown that modest values of  $m$  (between 3 and 20) often produce satisfactory results. The strategy of keeping the  $m$  most recent pairs  $\{s_i, y_i\}_{i=k-m}^{k-1}$  works well in practice. Indeed no other strategy has yet proved to be consistently better [Nocedal, Wright, 2006](#), p. 179. Normally, for large-scale problems, one takes  $m \ll n$ . In practice, the choice of  $m$  depends on the dimension of the problem and the storage of the employed computer [Sun, Yuan, 2006](#), p. 295.

The update process in detail: at iteration  $k$ , the current iterate is  $x_k$  and the set of vector pairs is given by  $\{s_i, y_i\}_{i=k-m}^{k-1}$ . At first some initial Hessian approximation  $B_k^{(0)}$  is chosen for the  $k$ -th iteration (in contrast to the standard BFGS iteration, this initial approximation is allowed to vary from iteration to iteration). The formula [Eq. \(1.6.2\)](#) is applied  $m$  times repeatedly, i.e.

$$B_k^{(j+1)} = V_{k-m+j}^T B_k^{(j)} V_{k-m+j} + \rho_{k-m+j} s_{k-m+j} s_{k-m+j}^T, \quad j = 0, 1, \dots, m-1. \quad (1.6.3)$$

The L-BFGS approximation, called  $B_k^{L-BFGS}$ , reads the following:



$$\begin{aligned}
 B_k^{L-BFGS} &= B_k^{(m)} = V_{k-1}^T B_k^{(m-1)} V_{k-1} + \rho_{k-1} s_{k-1} s_{k-1}^T = \\
 &= \dots = \\
 &= (V_{k-1}^T \dots V_{k-m}^T) B_k^{(0)} (V_{k-m} V_{k-m+1} \dots V_{k-1}) + \\
 &\quad + \rho_{k-m} (V_{k-1}^T \dots V_{k-m+1}^T) s_{k-m} s_{k-m}^T (V_{k-m+1} \dots V_{k-1}) \\
 &\quad + \rho_{k-m+1} (V_{k-1}^T \dots V_{k-m+2}^T) s_{k-m+1} s_{k-m+1}^T (V_{k-m+2} \dots V_{k-1}) \\
 &\quad + \dots \\
 &\quad + \rho_{k-1} s_{k-1} s_{k-1}^T.
 \end{aligned}$$

That means  $B_k^{L-BFGS}$  can be calculated completely from  $B_k^{(0)}$  and the vector pairs  $\{s_i, y_i\}_{i=k-m}^{k-1}$ .  $B_k^{L-BFGS}$  must be considered as an approximation of  $B_k^{BFGS}$ . Nevertheless this matrix fulfills the Quasi-Newton equation [Eq. \(1.2.2\) Geiger, Kanzow, 1999](#).

In fact, there is no need to compute and save  $B_k^{L-BFGS}$  explicitly. Instead, one only saves the pairs  $\{s_i, y_i\}_{i=k-m}^{k-1}$  and computes  $B_k^{L-BFGS} g_k = B_k^{L-BFGS} \nabla f(x_k)$  [Sun, Yuan, 2006](#). The product can be obtained by performing a sequence of inner products and vector summations involving  $g_k$  and the pairs  $\{s_i, y_i\}_{i=k-m}^{k-1}$  [Nocedal, Wright, 2006](#). So, we have

$$\begin{aligned}
 B_k^{L-BFGS} g_k &= (V_{k-1}^T \dots V_{k-m}^T) B_k^{(0)} (V_{k-m} V_{k-m+1} \dots V_{k-1}) g_k + \\
 &\quad + \rho_{k-m} (V_{k-1}^T \dots V_{k-m+1}^T) s_{k-m} s_{k-m}^T (V_{k-m+1} \dots V_{k-1}) g_k \\
 &\quad + \rho_{k-m+1} (V_{k-1}^T \dots V_{k-m+2}^T) s_{k-m+1} s_{k-m+1}^T (V_{k-m+2} \dots V_{k-1}) g_k \\
 &\quad + \dots \\
 &\quad + \rho_{k-1} s_{k-1} s_{k-1}^T g_k.
 \end{aligned}$$

Since  $V_i g_k = (I - \rho_i y_i s_i^T) g_k$  for  $i = k-1, \dots, k-m$ , one can derive a recursive method to compute the product efficiently [Sun, Yuan, 2006](#), p. 293. See the L-BFGS two-loop recursion for  $B_k^{L-BFGS} g_k$ , [Algorithm 6](#).

Without considering the multiplication  $B_k^{(0)} q$ , the L-BFGS two-loop recursion requires  $4mn$  multiplications. If  $B_k^{(0)}$  is diagonal, then  $n$  additional multiplications are needed. Apart from being inexpensive, this recursion has the advantage that the multiplication by the initial matrix  $B_k^{(0)}$  is isolated from the rest of the computations, allowing this matrix to be chosen freely and to vary between iterations. One may even use an implicit choice of  $B_k^{(0)}$  by defining some initial approximation  $H_k^{(0)}$  to the Hessian (not its inverse) and obtaining  $r$  by solving the system  $H_k^{(0)} r = q$  [Nocedal, Wright, 2006](#), p. 178.

$B_k^{(0)}$  can be an arbitrarily, symmetrical and positive definite matrix. In general  $B_k^{(0)}$  will be a multiple of the identity matrix, so that it can be stored very easily [Geiger, Kanzow, 1999](#), p. 198. A method for choosing  $B_k^{(0)}$  that has proven effective in practice is to set  $B_k^{(0)} = \gamma_k I$ , where

**Algorithm 6** L-BFGS two-loop recursion for  $B_k^{L-BFGS} g_k$ 


---

```

1:  $q = g_k$ 
2: for  $i = k - 1, k - 2, \dots, k - m$  do
3:    $\rho_i = \frac{1}{s_i^T y_i}$ 
4:    $\alpha_i = \rho_i s_i^T q$ 
5:    $q = q - \alpha_i y_i$ 
6: end for
7:  $r = B_k^{(0)} q$ 
8: for  $i = k - m, k - m + 1, \dots, k - 1$  do
9:    $\beta = \rho_i y_i^T r$ 
10:   $r = r + s_i(\alpha_i - \beta)$ 
11: end for
12: stop with result  $B_k^{L-BFGS} g_k = r$ 

```

---

$$\gamma_k = \frac{s_{k-1}^T y_{k-1}}{y_{k-1}^T y_{k-1}}. \quad (1.6.4)$$

$\gamma_k$  is the scaling factor that attempts to estimate the size of the true Hessian matrix along the most recent search direction. This choice helps to ensure that the search direction  $d_k$  is well scaled, and as a result the stepsize  $\alpha_k = 1$  is accepted in most iterations. It is important that the line search is based on the (strong) Wolfe conditions, so that the BFGS updating is stable [Nocedal, Wright, 2006](#), p. 178-179. With the previous theory, the following algorithm can be created for the L-BFGS-method:

**Algorithm 7** L-BFGS-Method

---

```

1:  $x_0 \in \mathbb{R}^n$ ,  $B_0^{L-BFGS} \in \mathbb{R}^{n \times n}$  spd,  $0 \leq \epsilon < 1$ ,  $\sigma \in (0, \frac{1}{2})$ ,  $\rho \in (\sigma, 1)$ ,  $m \in \mathbb{N}$ , set  $k = 0$ 
2: while  $\|\nabla f(x_k)\| > \epsilon$  do
3:   Choose  $B_k^{(0)}$  (e.g.  $B_k^{(0)} = \gamma_k I$  from Eq. \(3.6.2\)).
4:   Compute  $d_k = -B_k^{L-BFGS} \nabla f(x_k)$  with  $B_k^{L-BFGS} \nabla f(x_k)$  from Algorithm 6.
5:   Find a stepsize  $\alpha_k = \alpha(\sigma, \rho)$  that satisfies the (strong) Wolfe conditions.
6:   Set  $x_{k+1} = x_k + \alpha_k d_k$ .
7:   if  $k > m$  then
8:     Discard the vector pairs  $\{s_{k-m}, y_{k-m}\}$  from storage.
9:   end if
10:  Compute and save  $s_k = x_{k+1} - x_k$ ,  $y_k = g_{k+1} - g_k$ .
11:  Set  $k = k + 1$ .
12: end while
13: return  $x_k$ 

```

---

Unlike the conventional quasi-Newton methods, this one follows a different algorithmic sequence. Here, first the approximation of the Hessian inverse  $B_k^{L-BFGS}$  is calculated for the current iteration  $x_k$  and then the next iteration  $x_{k+1}$  is calculated. For example with the (globalized) BFGS-method, the new iteration  $x_{k+1}$  is calculated first and then the approximation of the Hessian inverse for the new iteration  $B_{k+1}^{BFGS}$  is calculated. This is the core idea of this method. Instead of passing the completely calculated matrix, only the vector pairs  $\{s_i, y_i\}_{i=k-m}^{k-1}$  are passed in each iteration and at the beginning

the approximation  $B_k^{L-BFGS}$  is created from these. The matrices  $B_k^{L-BFGS}$  are not explicitly stored, but only the vector pairs needed for the calculation and the start matrix  $B_k^{(0)}$ . For small values of  $m$  and larger dimensions  $n$ , the memory requirement for the L-BFGS-method is thus considerably lower than for the (globalized) BFGS-method itself, namely  $O(mn)$  instead of  $O(n^2)$ , which is due to the fact that the stepsize  $d_k$  can be obtained with  $O(mn)$  operations Geiger, Kanzow, 1999, p. 200-201.

In practical applications of the L-BFGS-method, the strong Powell-Wolfe stepsize strategy is used. This is because the L-BFGS-method seems to depend more on the choice of a “good” stepsize  $\alpha_k > 0$  than, for example, the (globalized) BFGS-method, and because the “optimal” stepsize can be better approximated by means of the strong Powell-Wolfe rule Geiger, Kanzow, 1999, p. 212-213.

During its first  $m - 1$  iterations, the L-BFGS-method (Algorithm 7) is equivalent to the inverse global BFGS-method (Algorithm 4) if the initial matrix is the same in both methods ( $B_0^{L-BFGS} = B_0^{BFGS}$ ), and if L-BFGS chooses  $B_k^{(0)} = B_0^{L-BFGS}$  at each iteration Nocedal, Wright, 2006, p. 179.

Before discussing the convergence properties, it must first be ensured that the L-BFGS-method (Algorithm 7) is well defined:

**Satz 1.6.1** (Geiger, Kanzow, 1999, Note 12.3). *If  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable and bounded below, then  $s_k^T y_k > 0$  holds for the sequences  $\{s_k\}_k$  and  $\{y_k\}_k$  generated by the L-BFGS-method (Algorithm 7). Furthermore, the matrices of the sequence  $\{B_k^{L-BFGS}\}_k$  are symmetric and positive definite and the L-BFGS-method (Algorithm 7) is well defined.*

The following statement can be made about global convergence:

**Satz 1.6.2** (Sun, Yuan, 2006, Theorem 5.7.4). *Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be a twice continuously differentiable and uniformly convex function. Then the iterative sequence  $\{x_k\}_k$  generated by the L-BFGS-method (Algorithm 7) converges to the unique minimizer  $x^*$  of  $f$ .*

Thus sequences  $\{x_k\}_k$  generated by the L-BFGS-method, like by the inverse global BFGS-method, converge globally for twice continuously differentiable, uniformly convex functions. The following can be said about their rate of convergence:

**Satz 1.6.3** (Sun, Yuan, 2006, Theorem 5.7.7). *Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be a twice continuously differentiable and uniformly convex function. Assume that the iterative sequence  $\{x_k\}_k$  generated by the L-BFGS-method (Algorithm 7) converges to the unique minimizer  $x^*$  of  $f$ . Then the rate of convergence is at least  $R$ -linear.*

This theorem indicates that the L-BFGS-method often converges slowly, which leads to a relatively large number of function evaluations. Also, it is inefficient on highly ill-conditioned optimization problems. Though there are some weaknesses, L-BFGS-method is a main choice for large-scale problems in which the true Hessian is not sparse, because, in this case, it may outperform other rival algorithms Sun, Yuan, 2006.

The memoryless BFGS-method should also be mentioned. Here,  $B_k^{BFGS} = I$  is inserted into the formula Eq. (1.6.2). It reads

$$B_{k+1}^{BFGS} = V_k^T V_k + \rho_k s_k s_k^T.$$

This formula satisfies the quasi-Newton equation [Eq. \(1.2.2\)](#), is positive definite and is called the memoryless BFGS formula. Obviously, if  $m = 1$  and  $B_k^{(0)} = I$  for all  $k \in \mathbb{N}$ , the limited-memory BFGS-method is just the memoryless BFGS-method. The idea is to use only the information from the previous iteration [Sun, Yuan, 2006](#).

## 2 RIEMANNIAN MANIFOLDS

- Riemannian Manifold
- Tangent Space with the defined operations
- Self-adjoint Operators on the tangent space
- Matrix Representation of Operators
- Gradient
- Hessian
- Retraction
- Exponential
- Logarithm
- Vector Transport
- Isometric Vector Transport
- Parallel Transport/Translation
- Vector transport by differentiated retraction
- Flat and sharp operations
- Tangent Vector Field
- Normal Neighborhood
- Levi-Civita connection

**Definition 2.0.1** (Absil, Mahony, Sepulchre, 2008, Definition 4.1.1). A retraction on a manifold  $\mathcal{M}$  is a smooth mapping  $R$  from the tangent bundle  $\mathcal{T}\mathcal{M}$  onto  $\mathcal{M}$  with the following properties. Let  $R_x$  denote the restriction of  $R$  to  $\mathcal{T}_x\mathcal{M}$ .

- (i)  $R_x(0_x) = x$ , where  $0_x$  denotes the zero element of  $\mathcal{T}_x\mathcal{M}$ .
- (ii) With the canonical identification  $\mathcal{T}_{0_x}\mathcal{T}_x\mathcal{M} \simeq \mathcal{T}_x\mathcal{M}$ ,  $R_x$  satisfies

$$DR_x(0_x) = \text{id}_{\mathcal{T}_x\mathcal{M}},$$

where  $\text{id}_{\mathcal{T}_x\mathcal{M}}$  denotes the identity mapping on  $\mathcal{T}_x\mathcal{M}$ .

**Definition 2.0.2** (Absil, Mahony, Sepulchre, 2008, Definition 8.1.1). A vector transport on a manifold  $\mathcal{M}$  is a smooth mapping

$$T: \mathcal{TM} \oplus \mathcal{TM} \rightarrow \mathcal{TM}$$

$$(\eta_x, \xi_x) \mapsto T_{\eta_x}(\xi_x)$$

satisfying the following properties for all  $x \in \mathcal{M}$ :

- (i) (Associated retraction) There exists a retraction  $R$ , called the retraction associated with  $T$ , such that the following diagram commutes

$$\begin{array}{ccc} (\eta_x, \xi_x) & \xrightarrow{T} & T_{\eta_x}(\xi_x) \\ \downarrow & & \downarrow \pi \\ \eta_x & \xrightarrow{R} & \pi(T_{\eta_x}(\xi_x)) \end{array}$$

where  $\pi(T_{\eta_x}(\xi_x))$  denotes the foot of the tangent vector  $T_{\eta_x}(\xi_x)$ .

- (ii) (Consistency)  $T_{0_x}(\xi_x) = \xi_x$  for all  $\xi_x \in \mathcal{T}_x\mathcal{M}$ ;  
 (iii) (Linearity)  $T_{\eta_x}(a\xi_x + b\xi_x) = aT_{\eta_x}(\xi_x) + bT_{\eta_x}(\xi_x)$ .

**Definition 2.0.3** (Huang, 2013, p. 10). A vector transport  $T: \mathcal{TM} \oplus \mathcal{TM} \rightarrow \mathcal{TM}$  with associated retraction  $R$  is called isometric if it satisfies for all  $x \in \mathcal{M}$

$$g_{R(\eta_x)}(T_{\eta_x}(\xi_x), T_{\xi_x}(\xi_x)) = g_x(\eta_x, \xi_x)$$

for all  $(\eta_x, \xi_x) \in \mathcal{T}_x\mathcal{M} \oplus \mathcal{T}_x\mathcal{M}$ .

**Definition 2.0.4** (Absil, Mahony, Sepulchre, 2008, p. 172). A vector transport  $T: \mathcal{TM} \oplus \mathcal{TM} \rightarrow \mathcal{TM}$  by differentiated retraction is a vector transport given by

$$T_{\eta_x}(\xi_x) = DR_x(\eta_x)[\xi_x]$$

i.e.,

$$T_{\eta_x}(\xi_x) = \frac{d}{dt} R_x(\eta_x + t\xi_x)|_{t=0}$$

where  $R_x: \mathcal{T}_x\mathcal{M} \rightarrow \mathcal{M}$  is a retraction.

**Definition 2.0.5** (Qi, 2011, Definition 2.2.1). Let  $(\mathcal{M}, g)$  be a Riemannian manifold and let  $X = X^i \partial_i$  be a vector field on  $\mathcal{M}$ , where  $\{\partial_i\}$  is a local frame for the tangent bundle  $\mathcal{TM}$ . The flat of  $X$  is defined by  $X^\flat = g_{ij} X^i dx^j = X_j dx^j$  where  $\{dx^i\}$  is the dual coframe and the metric  $g$  is defined locally, using Einstein notation, as  $g = g_{ij} dx^i \otimes dx^j$ . Equivalently, we have  $X^\flat(Y) = g(X, Y)$  for all vectors  $X$  and  $Y$ .

**Definition 2.0.6** (Absil, Mahony, Sepulchre, 2008, p. 192). Let  $X$  be a topological space. A neighborhood of a point  $x \in X$  is a subset of  $X$  that includes an open set containing  $x$ .

## 3 THE BFGS-METHOD FOR RIEMANNIAN MANIFOLDS

### 3.1 PRELIMINARIES

- Riemannian Newton-Method
- (strong) retraction-convexity
- Table of differences to the Euclidean setting

#### Generalization of the Wolfe conditions to Riemannian manifolds:

In order to select a suitable stepsize, a generalization of the Wolfe conditions to a Riemannian manifold is required. The generalized Wolfe conditions on  $\mathcal{M}$  are

$$f(\text{retr}_{x_k} \alpha_k \eta_k) \leq f(x_k) + c_1 \alpha_k g_{x_k}(\text{grad } f(x_k), \eta_k) \quad (3.1.1)$$

$$\frac{d}{dt} f(\text{retr}_{x_k} t \eta_k)|_{t=\alpha_k} \geq c_2 \frac{d}{dt} f(\text{retr}_{x_k} t \eta_k)|_{t=0} \quad (3.1.2)$$

with  $0 < c_1 < c_2 < 1$ . Condition [Eq. \(3.1.1\)](#) is the generalized Armijo condition, [Eq. \(1.1.2\)](#), and [Eq. \(3.1.2\)](#) is the generalized “curvature condition”, [Eq. \(1.1.3\)](#).

Other generalizations of the Wolfe conditions are possible (see e.g. [Ring, Wirth, 2012](#)). If the vector transport  $T$  is isometric, then [Eq. \(3.1.2\)](#) can be replaced by:

$$g_{x_k}(T_{x_k, \alpha_k \eta_k}^{-1}(\text{grad } f(\text{retr}_{x_k} \alpha_k \eta_k)), \eta_k) \geq c_2 g_{x_k}(\text{grad } f(x_k), \eta_k) \quad (3.1.3)$$

[Eq. \(3.1.3\)](#) transports the tangent vector that is in the tangent space of the potential next iterate  $\text{retr}_{x_k} \alpha_k \eta_k$  to  $\mathcal{T}_{x_k} \mathcal{M}$  and applies [Eq. \(1.1.3\)](#). For parallel transport,  $P$ , and the exponential map,  $\exp$ , as the retraction, conditions [Eq. \(3.1.2\)](#) and [Eq. \(3.1.3\)](#) are identical [Qi, 2011](#), p. 12-13.

A generalization of the strong Wolfe conditions consists of [Eq. \(3.1.1\)](#) and

$$|g_{\text{retr}_{x_k} \alpha_k \eta_k}(\text{grad } f(\text{retr}_{x_k} \alpha_k \eta_k), \text{Dretr}_{x_k} \alpha_k \eta_k[\eta_k])| \leq c_2 |g_{x_k}(\text{grad } f(x_k), \eta_k)| \quad (3.1.4)$$

where  $0 < c_1 < c_2 < 1$ .



The existence of a stepsize  $\alpha_k$  satisfying Eq. (3.1.1) and Eq. (3.1.4) can be shown by an almost literal repetition of that for the strong Wolfe conditions in the Euclidean case (see Lemma 1.1.1).

**Proposition 3.1.1** (Sato, Iwai, 2015, Proposition 2.1.). *Let  $\mathcal{M}$  be a Riemannian manifold with a retraction  $\text{retr}$ . If a smooth objective function  $f$  on  $\mathcal{M}$  is bounded below on  $\{\text{retr}_{x_k} \alpha_k \eta_k | \alpha > 0\}$  for  $x_k \in \mathcal{M}$  and for a descent direction  $\eta_k \in \mathcal{T}_{x_k} \mathcal{M}$ , and if constants  $c_1$  and  $c_2$  satisfy  $0 < c_1 < c_2 < 1$ , then there exists a stepsize  $\alpha_k$  which satisfies the strong Wolfe conditions Eq. (3.1.1) and Eq. (3.1.4).*

Using a vector transport by differentiated retraction (see Definition 2.0.4), then Eq. (3.1.4) can be expressed as

$$|g_{\text{retr}_{x_k} \alpha_k \eta_k}(\text{grad } f(\text{retr}_{x_k} \alpha_k \eta_k), T_{x_k, \alpha_k \eta_k}(\eta_k))| \leq c_2 |g_{x_k}(\text{grad } f(x_k), \eta_k)| \quad (3.1.5)$$

Sato, Iwai, 2015.

**Definition 3.1.2** (Huang, 2013, Definition 4.3.1). *For a function  $f: \mathcal{M} \rightarrow \mathbb{R}: x \mapsto f(x)$  on a Riemannian manifold  $\mathcal{M}$  with retraction  $R$  define  $\tilde{m}_{x,\eta}(t) = f(R_x(t\eta))$  for  $x \in \mathcal{M}$  and  $\eta \in \mathcal{T}_x \mathcal{M}$ . The function  $f$  is retractionconvex with respect to the retraction  $R$  in a set  $\mathcal{S}$  if for all  $x \in \mathcal{S}$ , all  $\eta \in \mathcal{T}_x \mathcal{M}$  and  $\|\eta\| = 1$ ,  $\tilde{m}_{x,\eta}(t)$  is convex for all  $t$  which satisfies  $R_x(t\eta) \in \mathcal{S}$ . Moreover,  $f$  is strongly retraction-convex in  $\mathcal{S}$  if  $\tilde{m}_{x,\eta}(t)$  is strongly convex for all  $x \in \mathcal{S}$  and all  $\|\eta\| = 1$  such that  $R_x(\eta) \in \mathcal{S}$ .*

**Definition 3.1.3** (Cruz Neto, Melo, Sousa, 2017, p. 5). *A function  $f: \mathcal{M} \rightarrow \mathbb{R}$  on a Riemannian manifold  $\mathcal{M}$  is convex if its restriction to every geodesic in  $\mathcal{M}$  is a convex function along the geodesic, i.e., if for every geodesic segment  $\gamma: [a, b] \rightarrow \mathcal{M}$  and every  $t \in [0, 1]$ ,*

$$f(\gamma((1-t)a + tb)) \leq (1-t)f(\gamma(a)) + tf(\gamma(b)).$$

*A convex function  $f$  is strictly convex if this inequality is strict whenever  $t \in (0, 1)$ . A convex function is always continuous. If  $f$  is smooth, it is known that  $f$  is (strictly) convex provided its Hessian is positive (definite) semidefinite, or equivalently if  $(f \circ \gamma)'' \geq 0$  ( $> 0$ ) for every geodesic  $\gamma: I \subset \mathbb{R} \rightarrow \mathcal{M}$ .*

**Definition 3.1.4.** *Let  $H$  be a Hilbert space with inner product  $\langle \cdot, \cdot \rangle$  and  $\mathcal{L}(H)$  the set of linear and continuous operators on  $H$ . An operator  $A \in \mathcal{L}(H)$  is called positive definite if*

$$\langle Ax, x \rangle \geq 0$$

*holds for all  $x \in H$ .*

**Definition 3.1.5.** *Let  $H$  be a Hilbert space with inner product  $\langle \cdot, \cdot \rangle$  and  $\mathcal{L}(H)$  the set of linear and continuous operators on  $H$ . An operator  $A \in \mathcal{L}(H)$  is called self-adjoint if  $A = A^*$  holds, i.e.*

$$\langle Ax, y \rangle = \langle x, Ay \rangle$$

holds for all  $x, y \in H$ .

**Definition 3.1.6.** An iterative update scheme of an algorithm on a Riemannian manifold  $\mathcal{M}$  is defined as: Starting with  $x_0 \in \mathcal{M}$  (an initial guess) the algorithm computes

$$x_{k+1} = R_{x_k}(\alpha_k \eta_k), \quad k = 1, 2, \dots,$$

where  $\eta_k \in \mathcal{T}_{x_k} \mathcal{M}$  is a tangent vector and  $\alpha_k > 0$  is a stepsize, which are determined in each iteration and  $R_{x_k}: \mathcal{T}_{x_k} \mathcal{M} \rightarrow \mathcal{M}$  a retraction, depending on the iterate  $x_k \in \mathcal{M}$ .

**Definition 3.1.7** (Zhang, Sra, 2016, Definition 2). A differentiable function  $f: \mathcal{M} \rightarrow \mathbb{R}$  is said to be geodesically  $\mu$ -strongly convex if for any  $x, y \in \mathcal{M}$ ,

$$f(y) \geq f(x) + g_x(\nabla f(x), \log_x y) + \frac{\mu}{2} d(x, y)^2$$

or, equivalently, for any geodesic  $\gamma$  such that  $\gamma(0) = x$ ,  $\gamma(1) = y$  and  $t \in [0, 1]$ ,

$$f(\gamma(t)) \leq (1-t)f(x) + tf(y) - \frac{\mu}{2} t(1-t)d(x, y)^2.$$

**Satz 3.1.8** (Deng, 2011, Theorem 1.1). Let  $H$  and  $K$  be Hilbert spaces over the same field. Let  $A \in \mathcal{L}(H)$  and  $G \in \mathcal{L}(K)$  both be invertible, and  $Y, Z \in \mathcal{L}(K, H)$ . Then  $A + YGZ^*$  is invertible iff  $G^{-1} + Z^*A^{-1}Y$  is invertible. In which case,

$$(A + YGZ^*)^{-1} = A^{-1} - A^{-1}Y(G^{-1} + Z^*A^{-1}Y)^{-1}Z^*A^{-1}. \quad (3.1.6)$$

**Satz 3.1.9** (Qi, 2011, Theorem 2.4.1). Consider any iteration of form  $x_{k+1} = \exp_{x_k} \alpha_k \eta_k$ , where  $\eta_k$  is a descent direction and  $\alpha_k$  satisfies the Wolfe conditions Eq. (3.1.1) and Eq. (3.1.2). Suppose that  $f$  is bounded below on  $\mathcal{M}$  and that  $f$  is continuously differentiable in an open set  $\mathcal{N}$  containing the level set  $\mathcal{L} = \{x: f(x) \leq f(x_0)\}$ , where  $x_0$  is the starting point of the iteration. Assume also that the gradient  $\text{grad } f$  is Lipschitz continuous on  $\mathcal{N}$ , then

$$\sum_{k \geq 0} \cos(\theta_k)^2 \|\text{grad } f(x_k)\|^2 < \infty,$$

where

$$\cos(\theta_k) = -\frac{g_{x_k}(\text{grad } f(x_k), \eta_k)}{\|\text{grad } f(x_k)\| \|\eta_k\|}.$$

We can conclude from [Theorem 3.1.9](#), if  $\lim_{k \rightarrow \infty} \|\text{grad } f(x_k)\| = 0$  holds, that the search directions are never too close to orthogonality with the gradient, i.e.  $\cos(\theta_k)^2$  stays away from 0. It follows that the algorithm would achieve global convergence to a set of stationary points. In practice, given the instability of an iteration at stationary points, such an algorithm is often effective at converging to an isolated minimizer when starting close enough, i.e., the global convergence result is used in a local manner [Qi, 2011](#), p. 25.

**Satz 3.1.10** ([Qi, 2011](#), Theorem 2.3.1). *Let  $\mathcal{M}$  be a manifold endowed with a  $C^2$  vector transport  $T_{\leftarrow}^{\text{retr}}$  and an associated retraction  $\text{retr}$ . Let  $F$  be a  $C^2$  tangent vector field on  $\mathcal{M}$ . Also let  $\mathcal{M}$  be endowed with an affine connection  $\nabla$ . Let  $\mathbb{D}F(x)$  denote the linear transformation of  $\mathcal{T}_x \mathcal{M}$  defined by  $\mathbb{D}F(x)[\xi_x] = \nabla_{\xi_x} F$  for all tangent vectors  $\xi_x$  to  $\mathcal{M}$  at  $x$ . Let  $\{\mathcal{B}_k\}_k$  be a sequence of bounded nonsingular linear transformations of  $\mathcal{T}_{x_k} \mathcal{M}$ , where  $k = 0, 1, \dots$ ,  $x_{k+1} = \text{retr}_{x_k} \eta_k$ , and  $\eta_k = -\mathcal{B}_k^{-1}[F(x_k)]$ . Assume that  $\mathbb{D}F(x^*)$  is nonsingular,  $x_k \neq x^*$ ,  $\forall k$ , and  $\lim_{k \rightarrow \infty} x_k = x^*$ . Then  $\{x_k\}_k$  converges superlinearly to  $x^*$  and  $F(x^*) = 0$  if and only if*

$$\lim_{k \rightarrow \infty} \frac{\|\mathcal{B}_k[\eta_k] - T_{\xi_k \leftarrow x^*}^{\text{retr}} \circ \mathbb{D}F(x^*) \circ (T_{\xi_k \leftarrow x^*}^{\text{retr}})^{-1}[\eta_k]\|}{\|\eta_k\|} = 0,$$

where  $\xi_k \in \mathcal{T}_{x^*} \mathcal{M}$  is defined by  $\xi_k = \text{retr}_{x^*}^{-1} x_k$ , i.e.  $\text{retr}_{x^*} \xi_k = x_k$ .

[Theorem 3.1.10](#) gives a necessary and sufficient condition for a Riemannian quasi-Newton algorithm that defines its search direction based on vector transport, its associated retraction and the transport of a linear transformation to achieve superlinear convergence [Qi, 2011](#), p. 89.

**Definition 3.1.11** ([Huang, Absil, Gallivan, 2018](#), Definition 4.1). *Let  $T_{\leftarrow(\cdot)}[\text{retr}]$  be a vector transport associated with a retraction  $\text{retr}$ . A function  $f$  on  $\mathcal{M}$  is said to be Lipschitz continuously differentiable with respect to  $T_{\leftarrow(\cdot)}[\text{retr}]$  on  $\mathcal{U} \subset \mathcal{M}$  if there exists  $L_1 > 0$  such that*

$$\|T_{x,\eta}^{\text{retr}}(\text{grad } f(x)) - \text{grad } f(\text{retr}_x \eta)\| \leq L_1 \|\eta\|$$

for all  $x \in \mathcal{U}$ ,  $\eta \in \mathcal{T}_x \mathcal{M}$  such that  $\text{retr}_x \eta \in \mathcal{U}$ .

**Definition 3.1.12** ([Huang, Absil, Gallivan, 2018](#), Definition 4.1). *Let  $T_{\leftarrow(\cdot)}[\text{retr}]$  be a vector transport associated with a retraction  $\text{retr}$ . A function  $f$  on  $\mathcal{M}$  is said to be Lipschitz continuously differentiable with respect to  $T_{\leftarrow(\cdot)}[\text{retr}]$  on  $\mathcal{U} \subset \mathcal{M}$  if there exists  $L_1 > 0$  such that*

$$\|T_{x,\eta}^{\text{retr}}(\text{grad } f(x)) - \text{grad } f(\text{retr}_x \eta)\| \leq L_1 \|\eta\|$$

for all  $x \in \mathcal{U}$ ,  $\eta \in \mathcal{T}_x \mathcal{M}$  such that  $\text{retr}_x \eta \in \mathcal{U}$ .

### 3.2 QUASI-NEWTON METHODS FOR RIEMANNIAN MANIFOLDS

This subsection creates a foundation for quasi-Newton methods on Riemannian manifolds. We try to name the most important points of the general theory, so that the derivation of the BFGS method on this structure is reasonable and makes sense. Quasi-Newton methods on Riemannian manifolds are often obtained by generalizing their Euclidean counterparts. This will also happen in some aspects in this paper but the purpose is that we want to start ... and build the theory

In subsection 1.2 the following properties have proved to be important for the approximating matrices  $H_k$  and  $B_k$  in the Euclidean case: positive-definiteness, symmetry and satisfying the quasi-Newton equation. We will now investigate the Riemannian analogue. Let us begin with the quasi-Newton equation, the core of the theory. The equation reads for the Euclidean case:

$$H_{k+1}(x_{k+1} - x_k) = \nabla f(x_{k+1}) - \nabla f(x_k) \quad \text{or} \quad H_{k+1}s_k = y_k.$$

This cannot be transferred one-to-one to manifolds. For the generalization to the Riemannian metric one idea would be to use  $\text{retr}^{-1}$  as a minus operator, since it returns a tangent vector which defines a curve connecting the two points  $x_k$  and  $x_{k+1}$ , and to compare the vectors, one could map  $\text{grad } f(x_k)$  to the tangent space of  $x_{k+1}$  using a vector transport. And since we now have an equation between tangent vectors, the matrix  $H_{k+1}$  becomes a linear operator  $\mathcal{H}_{k+1}$  on the tangent space  $\mathcal{T}_{x_{k+1}} \mathcal{M}$ . This leads to a naive quasi-Newton equation:

$$\mathcal{H}_{k+1}[T_{x_k, \eta_k}^{\text{retr}}(\eta_k)] = \text{grad } f(x_{k+1}) - T_{x_k, \eta_k}^{\text{retr}}(\text{grad } f(x_k)) \quad (3.2.1)$$

where  $\eta_k \in \mathcal{T}_{x_k} \mathcal{M}$  is the update vector at the iterate  $x_k$ , i.e.  $\text{retr}_{x_k} \eta_k = x_{k+1}$  [Absil, Mahony, Sepulchre, 2008](#), p. 179.

But here the first problems arise, namely the determination of the vector transport  $T_{\cdot, \cdot}(\cdot)$  and the retraction  $\text{retr}_{\cdot}$ . In Chapter 2 the two concepts were presented and kept as general as possible. So there is no unique vector transport  $T_{\cdot, \cdot}(\cdot)$  and no unique retraction  $\text{retr}_{\cdot}$  in general. The quasi-Newton equation and the resulting quasi-Newton method depend crucially on the choice of these two maps. The above equation is not wrong a priori but to get a well-defined method and clear results regarding its convergence, it is useful to take a closer look at these functions.

We now consider the most natural generalization of the Euclidean quasi-Newton equation [Eq. \(1.2.2\)](#). It lends itself to parallel transport  $P_{\leftarrow}(\cdot)$  and the exponential map  $\exp_{\cdot}$  or its inverse, the logarithmic map  $\log_{\cdot}$ . In subsection 1.2 Taylor's theorem was used, among other things, to deduce the secant equation. Similarly, there is a version of Taylor's Theorem for a vector field on a manifold [Huang, 2013](#). The following Lemma is a first-order Taylor formula for tangent vector fields.

**Lemma 3.2.1** (Absil, Mahony, Sepulchre, 2008, Lemma 7.4.7). *Let  $x \in \mathcal{M}$ , let  $\mathcal{V}$  be a normal neighborhood of  $x$ , and let  $\zeta$  be a  $C^1$  tangent vector field on  $\mathcal{M}$ . Then, for all  $y \in \mathcal{V}$ ,*

$$P_{x \leftarrow y}(\zeta_y) = \zeta_x + \nabla_{\xi} \zeta + \int_0^1 (P_{x \leftarrow \gamma(\tau)}(\nabla_{\gamma'(\tau)} \zeta) - \nabla_{\xi} \zeta) d\tau, \quad (3.2.2)$$

where  $\gamma$  is the unique minimizing geodesic satisfying  $\gamma(0) = x$  and  $\gamma(1) = y$ , and  $\xi = \log_x y = \gamma'(0)$ .

Applying Taylor's Theorem on the gradient of  $f$  at  $x_{k+1}$ , one obtains

$$P_{x_k \leftarrow x_{k+1}}(\text{grad } f(x_{k+1})) = \text{grad } f(x_k) + \nabla_{\xi} \text{grad } f(x_k) + \int_0^1 (P_{x_k \leftarrow \gamma(\tau)}(\nabla_{\gamma'(\tau)} \text{grad } f(x_k)) - \nabla_{\xi} \text{grad } f(x_k)) d\tau,$$

where  $\gamma_k$  is the unique minimizing geodesic satisfying  $\gamma_k(0) = x_k$  and  $\gamma_k(1) = x_{k+1}$ , and  $\xi = \log_{x_k} x_{k+1} = \gamma'_k(0)$ .

Ignoring the integral remainder term and rearranging yields

$$P_{x_k \leftarrow x_{k+1}}(\text{grad } f(x_{k+1})) - \text{grad } f(x_k) \approx \nabla_{\xi} \text{grad } f(x_k) = \text{Hess } f(x_k)[\log_{x_k} x_{k+1}]$$

This formula is very similar to the Euclidean quasi-Newton equation. It is defined on  $\mathcal{T}_{x_k} \mathcal{M}$ , the desired approximation  $\mathcal{H}_{k+1}$  of the Hessian  $\text{Hess } f(x_{k+1})$  must be an operator on  $\mathcal{T}_{x_{k+1}} \mathcal{M}$ . Applying parallel transport, yields

$$\text{grad } f(x_{k+1}) - P_{x_{k+1} \leftarrow x_k}(\text{grad } f(x_k)) = \mathcal{H}_{k+1}[P_{x_{k+1} \leftarrow x_k}(\log_{x_k} x_{k+1})] \quad (3.2.3)$$

This is one Riemannian version of the quasi-Newton equation Eq. (1.2.2). There are several possible generalizations of the Euclidean secant condition to a Riemannian manifold Huang, 2013, p. 17. In the following, the term Riemannian quasi-Newton equation is used to refer to Eq. (3.2.3). We introduce  $s_k = P_{x_{k+1} \leftarrow x_k}(\log_{x_k} x_{k+1})$  and  $y_k = \text{grad } f(x_{k+1}) - P_{x_{k+1} \leftarrow x_k}(\text{grad } f(x_k))$  to shorten Eq. (3.2.3) - as in the Euclidean case - further to

$$y_k = \mathcal{H}_{k+1}[s_k] \quad \text{or equivalently} \quad \mathcal{B}_{k+1}[y_k] = s_k,$$

where  $\mathcal{B}_{k+1} = \mathcal{H}_{k+1}^{-1}$ .

From the definition of  $s_k$  we can conclude that we use exponential map  $\exp \cdot$  as retraction  $\text{retr} \cdot$  in this variant, since  $\log \cdot = \exp \cdot^{-1}$ . This means the iterative-scheme has the form

$$x_{k+1} = \exp_{x_k} \alpha_k \eta_k,$$

where  $\alpha_k > 0$  is a stepsize to be determined and  $\eta_k \in \mathcal{T}_{x_k} \mathcal{M}$  is an update direction to be determined. The equation Eq. (3.2.3) explicitly uses the exponential map and parallel transport, due to its origins in Taylor's Theorem, but there's no need for that. As already mentioned, the Riemannian quasi-Newton method depends on the choice of these two maps. Alternate forms of this equation can be derived by using a different retraction or a different vector transport [Huang, 2013](#). If these change, the work will draw attention to this and discuss the resulting effects on the method. For now we work with the equation Eq. (3.2.3), i.e. we get the next iteration by using the exponential map  $\exp_{\cdot}$  as the retraction  $\text{retr}_{\cdot}$  and use the parallel transport  $P_{\leftarrow}(\cdot)$  as vector transport  $T_{\cdot}(\cdot)$  to compare tangent vectors. Now we come to the calculation of the update direction  $\eta_k$ . From the equation Eq. (3.2.3) it can be concluded that the operator  $\mathcal{H}_k$  approximates the Hessian operator  $\text{Hess } f(x_k)$  of  $f$  at  $x_k$ . We consider the pullback  $\hat{f}_{x_k} = f \circ \exp_{x_k} : \mathcal{T}_{x_k} \mathcal{M} \rightarrow \mathbb{R}$  at the current iterate  $x_k$ , because we want to build a quadratic model of  $f$  around  $x_k$ . Let  $m_k$  be the order-2 Taylor expansion of  $\hat{f}_{x_k}$  around the origin  $0_{x_k}$  of  $\mathcal{T}_{x_k} \mathcal{M}$ , i.e.

$$\begin{aligned} \hat{f}_{x_k}(\eta) &\approx m_k(\eta) = \hat{f}_{x_k}(0_{x_k}) + D\hat{f}_{x_k}(0_{x_k})[\eta] + \frac{1}{2}D^2\hat{f}_{x_k}(0_{x_k})[\eta, \eta] \\ &= f(x_k) + g_{x_k}(\eta, \text{grad } \hat{f}_{x_k}(0_{x_k})) + \frac{1}{2}g_{x_k}(\eta, \text{Hess } \hat{f}_{x_k}(0_{x_k})[\eta]) \\ &= f(x_k) + g_{x_k}(\eta, \text{grad } f(x_k)) + \frac{1}{2}g_{x_k}(\eta, \text{Hess } f(x_k)[\eta]) \end{aligned}$$

since  $D \text{retr}_{x_k} 0_{x_k} = \text{id}_{\mathcal{T}_{x_k} \mathcal{M}}$  (see ??), it follows that  $D\hat{f}_{x_k}(0_{x_k}) = Df(x_k)$ , hence  $\text{grad } \hat{f}_{x_k}(0_{x_k}) = \text{grad } f(x_k)$ . Since the exponential map  $\exp_{\cdot}$  is a second-order retraction,  $\text{Hess } \hat{f}_{x_k}(0_{x_k}) = \text{Hess } f(x_k)$  holds and  $m_k$  is in fact order 2 [Absil, Mahony, Sepulchre, 2008](#), p. 139. Since the Hessian of  $f$  at the iterate  $x_k$  should be approximated by the operator  $\mathcal{H}_k$ , we get the following model

$$m_k(\eta) = f(x_k) + g_{x_k}(\eta, \text{grad } f(x_k)) + \frac{1}{2}g_{x_k}(\eta, \mathcal{H}_k[\eta]).$$

If we take the derivative  $\frac{d}{d\eta}$  of the model, set it to zero and solve the equation, we get

$$\eta_k = -\mathcal{H}_k^{-1}[\text{grad } f(x_k)] = -\mathcal{B}_k[\text{grad } f(x_k)] \in \mathcal{T}_{x_k} \mathcal{M}.$$

It remains to be clarified whether this is a descent direction and thus can be used as an update vector. Before discussing this, we consider a property which is required for the approximation in the Euclidean case. In this case we demand that  $H_k$  or  $B_k$  are symmetric for all  $k$ . Since the Hessian matrix  $\nabla^2 f(x_k)$  is for twice continuously differentiable functions always symmetric and since  $H_k$  or  $B_k$  are approximating

it, it makes sense that this should be required. The Riemannian Hessian has a similar characteristic:

**Proposition 3.2.2** (Absil, Mahony, Sepulchre, 2008, Proposition 5.5.3). *The Riemannian Hessian is symmetric (in the sense of the Riemannian metric). That is,*

$$g(\text{Hess } f[\eta], \xi) = g(\eta, \text{Hess } f[\xi])$$

for all  $\eta, \xi \in \mathfrak{X}(M)$ .

We see immediately that symmetry can be generalized almost one to one for tangent spaces. But the term “symmetrical” could be misleading in some places. Since the Hessian can be seen as an operator on a tangent space, we call it self-adjoint instead. It follows: On a Riemannian manifold the Hessian  $\text{Hess } f(\cdot)$  of a function  $f$  is always a self-adjoint operator. The natural consequence is that it is required that the approximations  $\mathcal{H}_k$  and  $\mathcal{B}_k$  should also be self-adjoint Huang, 2013, p. 20.

Finally, it remains to be clarified whether the  $\eta_k$ , as determined above, is really a descent direction. In the Euclidean case, a descent direction is defined by the fact that the inner product of the vector indicating this direction and the gradient is negative. Consequently it is required that for quasi-Newton methods the approximations  $H_k$  and  $B_k$  are positive definite in each iteration to ensure a continuous descent. This idea can also be adopted here. In the Riemannian case a descent direction of a function  $f: \mathcal{M} \rightarrow \mathbb{R}$  at a point  $x \in \mathcal{M}$  denotes a tangent vector  $\eta \in \mathcal{T}_x \mathcal{M}$  with  $g_x(\text{grad } f(x), \eta) < 0$ . This property ensures that the objective function  $f$  indeed decreases along the search direction Ring, Wirth, 2012, p. 5. In our case this means

$$g_{x_k}(\text{grad } f(x_k), -\alpha_k \mathcal{B}_k[\text{grad } f(x_k)]) < 0.$$

Using the linearity of  $g_{x_k}$  and since  $\alpha_k > 0$ , this implies that  $g_{x_k}(\text{grad } f(x_k), \mathcal{B}_k[\text{grad } f(x_k)]) = g_{x_k}(\mathcal{B}_k[\text{grad } f(x_k)], \text{grad } f(x_k)) > 0$  must hold in every iteration. So it makes sense that we require that the  $\mathcal{B}_k$  (and thus of course also  $\mathcal{H}_k$ ) is a positive definite operator in every iteration, so that the  $\eta_k$  is a descent direction. The goal is to find updates hence all operators in the sequence  $\{\mathcal{B}_k\}_k$  are positive definite and self-adjoint to create a continuous descent.

The positive-definiteness of the approximating operators  $\mathcal{H}_k$  and  $\mathcal{B}_k$  has a consequence that is similar to one we have already seen in the Euclidean case: the curvature condition Eq. (1.2.6). Much as in the Euclidean case, it is essential that

$$g_{x_{k+1}}(s_k, y_k) > 0 \tag{3.2.4}$$

holds, otherwise the secant condition  $\mathcal{H}_{k+1}[s_k] = y_k$  cannot hold with  $\mathcal{H}_{k+1}$  positive definite, whereas positive definiteness of the operators is the key to guarantee that the search directions  $\eta_k$  are descent directions Huang, 2013, p. 54. In the Euclidean setting, the inequality Eq. (1.2.6) holds for any two points  $x_k$  and  $x_{k+1}$ , if the objective  $f$  is strongly convex. This can also be adopted almost one-to-one.

But we need a different characterization of geodesically  $\mu$ -strongly convex functions on Riemannian manifolds.

**Satz 3.2.3.** *A differentiable function  $f: \mathcal{M} \rightarrow \mathbb{R}$  is geodesically  $\mu$ -strongly convex if and only if*

$$g_x(P_{y \leftarrow x}(\text{grad } f(y)) - \text{grad } f(x), \log_x y) \geq \mu d(x, y)^2$$

*holds for any  $x, y \in \mathcal{M}$ .*

*Beweis.* From Definition 3.1.7 we have

$$g_x(\text{grad } f(x), \log_x y) \leq f(y) - f(x) - \frac{\mu}{2} d(x, y)^2 \quad \text{A}$$

$$g_y(\text{grad } f(y), \log_y x) \leq f(x) - f(y) - \frac{\mu}{2} d(x, y)^2. \quad \text{B}$$

Now we draw  $-1$  twice from the inner product  $g_y(\cdot, \cdot)$  of the left side of inequality B, i.e.

$$g_y(\text{grad } f(y), \log_y x) = (-1) \cdot g_y(\text{grad } f(y), -\log_y x) = g_y(-\text{grad } f(y), -\log_y x).$$

Then we use, that  $-\log_y x = P_{y \leftarrow x}(\log_x y)$  and apply the parallel transport  $P_{x \leftarrow y}(\cdot)$  on both arguments

$$g_y(-\text{grad } f(y), -\log_y x) = g_x(-P_{x \leftarrow y}(\text{grad } f(y)), \log_x y).$$

That means

$$g_x(-P_{x \leftarrow y}(\text{grad } f(y)), \log_x y) \leq f(x) - f(y) - \frac{\mu}{2} d(x, y)^2. \quad \text{B}$$

holds. Now adding the inequalities A and B leads to

$$g_x(\text{grad } f(x) - P_{x \leftarrow y}(\text{grad } f(y)), \log_x y) \leq -\mu d(x, y)^2$$

Multiplying both sides with  $-1$  gives us the inequality. □



If  $y$  is set equal to  $x_{k+1}$  and  $x$  equal to  $x_k$ , only the parallel transport  $P_{x_{k+1} \leftarrow x_k}(\cdot)$  has to be applied to both arguments and we see that the Riemannian curvature condition Eq. (3.2.4) holds for geodesically  $\mu$ -strongly convex functions for all points  $x_k, x_{k+1} \in \mathcal{M}$ . If the function is not geodesically  $\mu$ -strongly convex, it cannot be guaranteed that the condition always holds. This is where determining the stepsize comes into play. Assuming that the stepsize  $\alpha_k$  meets the Wolfe conditions Eq. (3.1.3) or strong Wolfe conditions Eq. (3.1.4), we get

$$\begin{aligned}
g_{x_{k+1}}(s_k, y_k) &= g_{x_{k+1}}(P_{x_{k+1} \leftarrow x_k}(\alpha_k \eta_k), \text{grad } f(x_{k+1}) - P_{x_{k+1} \leftarrow x_k}(\text{grad } f(x_k))) \\
&= g_{x_{k+1}}(P_{x_{k+1} \leftarrow x_k}(\alpha_k \eta_k), \text{grad } f(x_{k+1})) - g_{x_k}(\alpha_k \eta_k, \text{grad } f(x_k)) \\
&= \alpha_k g_{x_{k+1}}(P_{x_{k+1} \leftarrow x_k}(\eta_k), \text{grad } f(x_{k+1})) - \alpha_k g_{x_k}(\eta_k, \text{grad } f(x_k)) \\
&= \alpha_k g_{x_k}(P_{x_k \leftarrow x_{k+1}}(\text{grad } f(x_{k+1})), \eta_k) - \alpha_k g_{x_k}(\eta_k, \text{grad } f(x_k)) \\
&\geq \alpha_k c_2 g_{x_k}(\eta_k, \text{grad } f(x_k)) - \alpha_k g_{x_k}(\eta_k, \text{grad } f(x_k)) \\
&= (c_2 - 1) \alpha_k g_{x_k}(\eta_k, \text{grad } f(x_k))
\end{aligned} \tag{3.2.5}$$

Since  $c_2 < 1$  and  $\eta_k$  is a descent direction, the right side is positive and the curvature condition holds. Using the strong Wolfe conditions Eq. (3.1.4), the inequality still remains correct.

At the end of this subsection a very general quasi-Newton method will be presented. We will now dispense with an exact definition of retraction and vector transport in order to preserve the generality. Depending on the choice of a certain combination of these two and the update formula for the approximation, this leads to different algorithms and thus to different results with respect to convergence.

---

**Algorithm 8** General Riemannian Quasi-Newton Method

---

- 1: Riemannian manifold  $\mathcal{M}$  with Riemannian metric  $g$ , vector transport  $T$  on  $\mathcal{M}$  with associated retraction  $\text{retr}$ , smooth real-valued function  $f$  on  $\mathcal{M}$ , initial iterate  $x_0 \in \mathcal{M}$ , initial Hessian approximation  $\mathcal{H}_0$ .
  - 2: **while** not converged **do**
  - 3:   Obtain  $\eta_k \in \mathcal{T}_{x_k} \mathcal{M}$  by solving  $\mathcal{H}_k[\eta_k] = -\text{grad } f(x_k)$ .
  - 4:   Determine the stepsize  $\alpha_k > 0$  by line search.
  - 5:   Set  $x_{k+1} = \text{retr}_{x_k}(\alpha_k \eta_k)$ .
  - 6:   Define  $s_k = T_{x_k, \alpha_k \eta_k}(\alpha_k \eta_k)$  and  $y_k = \text{grad } f(x_{k+1}) - T_{x_k, \alpha_k \eta_k}(\text{grad } f(x_k))$ .
  - 7:   Define the linear operator  $\mathcal{H}_{k+1}: \mathcal{T}_{x_{k+1}} \mathcal{M} \rightarrow \mathcal{T}_{x_{k+1}} \mathcal{M}$  by using  $s_k, y_k$  and  $\tilde{\mathcal{H}}_k = T_{x_k, \alpha_k \eta_k} \circ \mathcal{H}_k \circ (T_{x_k, \alpha_k \eta_k})^{-1}$ , such that
$$\mathcal{H}_{k+1}[s_k] = y_k$$
holds.
  - 8:   Set  $k = k + 1$ .
  - 9: **end while**
  - 10: **return**  $x_k$
-

## 3.3 THE BFGS FORMULA FOR RIEMANNIAN MANIFOLDS

As we have seen in Section 1.3, the success of the euclidean BFGS update was based on the properties of positive-definiteness, symmetry and closeness to the former approximation. Since we now work with linear operators, the setup has to be adjusted a little bit (matrices are only operators after all). But the theoretical basics remain the same. This means we transfer these properties to the Riemannian setup, since the Euclidean BFGS update, Eq. (1.3.1) or Eq. (1.3.2), admits several generalizations. We therefore require that the approximation  $\mathcal{H}_{k+1}$  of the Hessian operator  $\text{Hess } F(x_{k+1})$  meets the following characteristics:

- (i)  $\mathcal{H}_{k+1}$  satisfies the Riemannian quasi-Newton equation Eq. (3.2.3).
- (ii)  $\mathcal{H}_{k+1}$  is self-adjoint and positive definite.
- (iii)  $\mathcal{H}_{k+1}$  is "near"  $\mathcal{H}_k$ .

Of course, the approximation of the Hessian inverse,  $\mathcal{B}_{k+1}$ , should satisfy the same. With the fulfillment of these properties as a goal, we try to find an update formula which follows the methodology of Eq. (1.3.1). Most of the preparatory work for the Riemannian BFGS formula, short RBFGS formula, was already done in the Euclidean case. Often a derivation of this in the literature is skipped, because it seems to be clear, which is perfectly ok.

We assume that  $\mathcal{B}_k$  is a self-adjoint, positive definite operator on  $\mathcal{T}_{x_k} \mathcal{M}$ . We now want to create an operator on  $\mathcal{T}_{x_{k+1}} \mathcal{M}$  from this using the information from the tangent vectors  $y_k, s_k \in \mathcal{T}_{x_{k+1}} \mathcal{M}$ . We don't want the update formula to depend on a particular choice of retraction and vector transport (which would affect the definitions of  $y_k$  and  $s_k$ ), so we only assume that  $y_k$  and  $s_k$  meet the curvature condition, Eq. (3.2.4). As in the Euclidean case, the new operator  $\mathcal{B}_{k+1}$  should emerge from the old one  $\mathcal{B}_k$  by a rank 2 update. In Section 1.3 we saw that this is done by adding two dyadic products, which are rank one matrices, multiplied by a scalar to the old matrix  $\mathcal{B}_k$ .

This cannot be realized one-to-one in the Riemannian setup. First of all, the operators  $\mathcal{B}_{k+1}$  and  $\mathcal{B}_k$  are defined in different tangent spaces, which generally does not allow an addition. To overcome this obstacle, we introduce the following

$$\tilde{\mathcal{B}}_k = \mathcal{T}_{x_{k+1} \leftarrow x_k} \circ \mathcal{B}_k \circ \mathcal{T}_{x_k \leftarrow x_{k+1}} : \mathcal{T}_{x_{k+1}} \mathcal{M} \rightarrow \mathcal{T}_{x_{k+1}} \mathcal{M} \quad (3.3.1)$$

where  $\mathcal{T}_{x_{k+1} \leftarrow x_k} : \mathcal{T}_{x_k} \mathcal{M} \rightarrow \mathcal{T}_{x_{k+1}} \mathcal{M}$  represents at first a general vector transport and its inverse, i.e.  $\mathcal{T}_{x_k \leftarrow x_{k+1}} = \mathcal{T}_{x_{k+1} \leftarrow x_k}^{-1}$ . It shows, that the vector transport seems to be needed not only for comparing two vectors in different tangent spaces, but also for transporting operators from one tangent space to the appropriate one Huang, 2013, p. 20.

Now we have to find a similar concept for the dyadic products in the Riemannian setup. Therefore we introduce the concept of the so called index lowering function, musical isomorphism, or canonical isomorphism Definition 2.0.5. Put simply, it means: let  $\eta_x \in \mathcal{T}_x \mathcal{M}$  then  $\eta_x^\flat$  represents the flat of  $\eta_x$ , i.e.,  $\eta_x^\flat : \mathcal{T}_x \mathcal{M} \rightarrow \mathbb{R}$ ,  $\xi_x \mapsto \eta_x^\flat[\xi_x] = g_x(\eta_x, \xi_x)$ . With that we can create rank 1 operators in the tangent space  $\mathcal{T}_{x_{k+1}} \mathcal{M}$ . Overall, this enables us to write the following update:

Vector transport is used here, to transport the linear operator.

$$\mathcal{B}_{k+1}[\cdot] = \tilde{\mathcal{B}}_k[\cdot] + auu^b[\cdot] + bvv^b[\cdot]$$

where  $u, v \in \mathcal{T}_{x_{k+1}}\mathcal{M}$  and  $a, b \in \mathbb{R}$  are to be determined. One sees immediately that  $uu^b[\cdot]$  and  $vv^b[\cdot]$  are self-adjoint and positive definite operators on  $\mathcal{T}_{x_{k+1}}\mathcal{M}$  with rank 1.

At this point, a property must now be added to the vector transport T. We require that  $\mathcal{B}_{k+1}$  is a self-adjoint operator. Therefore it makes sense to require that  $\tilde{\mathcal{B}}_k$  is also a self-adjoint operator, so that this property is preserved by the update. Looking at Eq. (3.3.1) we know an isometric vector transport T guarantees that  $\tilde{\mathcal{B}}_k$  is self-adjoint if  $\mathcal{B}_k$  is self-adjoint Huang, 2013, p. 20.

With an isometric vector transport and the demand that this new operator fulfills Eq. (3.2.3), it follows:

*etwas lang. Was muss T sein um zusätzlich erfüllen? das steht da nämlich nicht. Oder brauchen wir iso? Dann reicht ein Satz damit.*

$$\mathcal{B}_{k+1}[y_k] = \tilde{\mathcal{B}}_k[y_k] + auu^b[y_k] + bvv^b[y_k] = s_k$$

Clearly,  $u, v \in \mathcal{T}_{x_{k+1}}\mathcal{M}$  <sup>are</sup> ~~can~~ not uniquely be determined. One possible choice is

$$u = s_k, \quad v = \tilde{\mathcal{B}}_k[y_k].$$

Since the vector transport is assumed to be isometric and  $\mathcal{B}_k$  is a self-adjoint operator, which implies  $\tilde{\mathcal{B}}_k^* = \tilde{\mathcal{B}}_k$ , it holds  $\xi \in \mathcal{T}_{x_{k+1}}\mathcal{M}$ :

*and hence for all*

$$v^b[\xi] = (\tilde{\mathcal{B}}_k[y_k])^b[\xi] = g_{x_{k+1}}(\tilde{\mathcal{B}}_k[y_k], \xi) = g_{x_{k+1}}(y_k, \tilde{\mathcal{B}}_k[\xi]) = y_k^b[\tilde{\mathcal{B}}_k[\xi]].$$

*we obtain*

*(V, I) bitte immer in Tafel mitschreiben, hier gerne ausschreiben*

Hence we obtain

$$a = \frac{1}{u^b[y_k]} = \frac{1}{s_k^b[y_k]}, \quad b = -\frac{1}{v^b[y_k]} = -\frac{1}{y_k^b[\tilde{\mathcal{B}}_k[y_k]]}.$$

With these constructions and updating the notation to  $\mathcal{B}_k = \mathcal{B}_k^{RDFP}$  for all  $k \in \mathbb{N}_0$  we get:

$$\mathcal{B}_{k+1}^{RDFP}[\cdot] = \tilde{\mathcal{B}}_k^{RDFP}[\cdot] + s_k \frac{s_k^b[\cdot]}{s_k^b[y_k]} - \tilde{\mathcal{B}}_k^{RDFP}[y_k] \frac{y_k^b(\tilde{\mathcal{B}}_k^{RDFP}[\cdot])}{y_k^b(\tilde{\mathcal{B}}_k^{RDFP}[y_k])}. \quad (3.3.2)$$

This is the Riemannian DFP update formula for approximating the Hessian inverse  $(\text{Hess } f(x_{k+1}))^{-1}$ . Since the direct DFP update is known for the approximation of the Hessian (see e.g. Huang, 2013, p. 19), it is left to the reader to show that this formula is its inverse.

Also in the Riemannian setup, the idea of quasi-Newton methods is very clear. Instead of calculating a complete approximation of the Hessian  $\text{Hess } f(x_{k+1})$  or its inverse at every iteration, the previous operator is simply updated using the obtained information included in  $s_k, y_k \in \mathcal{T}_{x_{k+1}}\mathcal{M}$  from the step.

*Nein. In dem paper kann man das schreiben wenn man sich sicher ist. In einer Abschlussarbeit ist das zu belegen (Ref)*

If we now follow the same strategy as in [Section 1.3](#) and replace the triple  $(\mathcal{B}_k^{DFP}, s_k, y_k)$  by  $(\mathcal{H}_k^{BFGS}, y_k, s_k)$  in [Eq. \(3.3.2\)](#), we obtain a direct update for the Hessian:

$$\mathcal{H}_{k+1}^{RBFGS}[\cdot] = \tilde{\mathcal{H}}_k^{RBFGS}[\cdot] + y_k \frac{y_k^b[\cdot]}{s_k^b[y_k]} - \tilde{\mathcal{H}}_k^{RBFGS}[s_k] \frac{s_k^b(\tilde{\mathcal{H}}_k^{RBFGS}[\cdot])}{s_k^b(\tilde{\mathcal{H}}_k^{RBFGS}[s_k])}. \quad (3.3.3)$$

This is the direct Riemannian BFGS formula for approximating the Hessian  $\text{Hess } f(x_{k+1})$ . It turns out that even in the Riemannian case we can speak of dual update formulae. We see that the computation of  $\mathcal{H}_{k+1}^{RBFGS}$  requires only first-order information, namely the gradient at  $x_k$  and  $x_{k+1}$ , which is a big advantage over the operator  $\text{Hess } f(x_{k+1})$  used in Newton's method, which involves second-order information [Gabay, 1982](#), p. 206.

It can be shown that this operator  $\mathcal{H}_{k+1}^{RBFGS}$ , generated by [Eq. \(3.3.3\)](#), fulfills the **above mentioned** characteristics. At first to the positive-definiteness and self-adjointness:

**Lemma 3.3.1** ([Qi, 2011](#), Lemma 2.4.1 + Lemma 2.4.2).

Let  $T$  be an isometric vector transport.

- (i) Let  $y_k = \text{grad } f(x_{k+1}) - T_{x_k, \alpha_k \eta_k}^{\text{etr}}(\text{grad } f(x_k)) \in \mathcal{T}_{x_{k+1}} \mathcal{M}$ ,  $s_k = T_{x_k, \alpha_k \eta_k}(\alpha_k \eta_k) \in \mathcal{T}_{x_{k+1}} \mathcal{M}$ ,  $s_k \neq 0$  and assume  $T$  represents an isometric vector transport. Let  $\{\mathcal{H}_k\}_k$  be a sequence of bounded invertible linear transformation of  $\mathcal{T}_{x_k} \mathcal{M}$ , where  $k = 0, 1, 2, \dots$ . If  $\mathcal{H}_k$  on  $\mathcal{T}_{x_k} \mathcal{M}$  is self-adjoint and positive definite with respect to the Riemannian metric then there exists an invertible linear transformation,  $\mathcal{J}_{k+1}$ , on  $\mathcal{T}_{x_{k+1}} \mathcal{M}$  such that

$$y_k = \mathcal{J}_{k+1} \circ \mathcal{J}_{k+1}^*[s_k]$$

if and only if

$$g_{x_{k+1}}(s_k, y_k) > 0.$$

- (ii) Using the notation and assumptions above, the sequence of linear transformations  $\{\mathcal{H}_k\}_k$  defined above is the same as the sequence defined by [Eq. \(3.3.3\)](#).

The operator  $\mathcal{J}_{k+1}$  is of course not unique. Since  $\mathcal{J}_{k+1}: \mathcal{T}_{x_{k+1}} \mathcal{M} \rightarrow \mathcal{T}_{x_{k+1}} \mathcal{M}$  is an invertible linear operator, the operator  $\mathcal{H}_k^{RBFGS}$  transfers the positive-definiteness and self-adjointness to the operator  $\mathcal{H}_{k+1}^{RBFGS} = \mathcal{J}_{k+1} \circ \mathcal{J}_{k+1}^*$  through the update [Eq. \(3.3.3\)](#). Therefore if the curvature condition, [Eq. \(3.2.4\)](#), holds in every iteration, the RBFGS formula, [Eq. \(1.3.1\)](#), produces a series of linear operators  $\mathcal{H}_k^{RBFGS}$  on  $\mathcal{T}_{x_k} \mathcal{M}$  that are all self-adjoint and positive definite with respect to the Riemannian metric if an isometric vector transport is used. This property is very important in order to prove convergence results later. Note that we have not bounded the condition number of  $\mathcal{H}_k^{RBFGS}$  [Qi, 2011](#), p.21-23. There is of course also a RBFGS formula for the approximation of the Hessian inverse  $(\text{Hess } f(x_{k+1}))^{-1}$ . This was usually also produced by simple generalization of the Euclidean counterpart [Eq. \(1.3.2\)](#). But one could also apply the Sherman–Morrison–Woodbury identity for operators, [Eq. \(3.1.6\)](#), to [Eq. \(3.3.3\)](#), this would result in an equivalent update:

Ref bitte  
genauer...  
von 5.42?  
Vielleicht die in eine  
Umgebung mit Nummer?

unterschiedliche  
transporte?

$$\begin{aligned}
\mathcal{B}_{k+1}^{RBFGS}[\cdot] &= \tilde{\mathcal{B}}_k^{RBFGS}[\cdot] - s_k \frac{y_k^b[\tilde{\mathcal{B}}_k^{RBFGS}[\cdot]]}{y_k^b[s_k]} - \tilde{\mathcal{B}}_k^{RBFGS}[y_k] \frac{s_k^b[\cdot]}{s_k^b[y_k]} + s_k \frac{y_k^b[\tilde{\mathcal{B}}_k^{RBFGS}[y_k]] s_k^b[\cdot]}{(y_k^b[s_k])^2} + s_k \frac{s_k^b[\cdot]}{s_k^b[y_k]} \\
&= \left( \text{id}[\cdot] - \frac{s_k y_k^b[\cdot]}{s_k^b[y_k]} \right) \tilde{\mathcal{B}}_k^{RBFGS}[\cdot] \left( \text{id}[\cdot] - \frac{y_k s_k^b[\cdot]}{s_k^b[y_k]} \right) + \frac{s_k s_k^b[\cdot]}{s_k^b[y_k]}
\end{aligned} \tag{3.3.4}$$

As the inverse of  $\mathcal{H}_{k+1}^{RBFGS}$  fulfills  $\mathcal{B}_{k+1}^{RBFGS}$ , created by Eq. (3.3.4), also the positive-definiteness, the self-adjointness and the Riemannian quasi-Newton equation, Eq. (3.2.3), for all  $k \in \mathbb{N}$ .

In the Euclidean case we have seen that the third characteristic ( $H_{k+1}$  is “near”  $H_k$ ) is important for the showing the uniqueness of Eq. (1.3.2) and determining the convergence rate, generated by a quasi-Newton method using this formula. This was achieved by showing that Eq. (1.3.2) is a solution to an optimization problem (see Theorem 1.3.4). Unfortunately, similar relationships have not yet been explored in the Riemannian setup, but this is not absolutely necessary, since the uniqueness is more or less obvious and the convergence rate is determined even without this approach. Nevertheless, a generalization of Eq. (1.3.3) has been made in Huang, 2013, p. 19 and  $\mathcal{B}_{k+1}^{RBFGS}$ , defined by Eq. (3.3.4), is the solution:

*sehr schwammig!*

$$\begin{aligned}
&\min_{\mathcal{B}: \mathcal{T}_{x_{k+1}} \mathcal{M} \rightarrow \mathcal{T}_{x_{k+1}} \mathcal{M}} \|\mathcal{B} - \tilde{\mathcal{B}}_k\|_{\mathcal{W}_{\mathcal{H}}^2} \\
&\text{s.t. } \mathcal{B} = \mathcal{B}^*, \quad \mathcal{B}[y_k] = s_k
\end{aligned} \tag{3.3.5}$$

where  $\tilde{\mathcal{B}}_k$ ,  $s_k$  and  $y_k$  are as before.  $\mathcal{W}_{\mathcal{H}}^2$  is a self-adjoint and positive definite operator on  $\mathcal{T}_{x_{k+1}} \mathcal{M}$  satisfying  $\mathcal{W}_{\mathcal{H}}^2[s_k] = y_k$  and  $\|\mathcal{A}\|_{\mathcal{W}_{\mathcal{H}}^2} = \|\hat{\mathcal{W}}_{\mathcal{H}} G^{\frac{1}{2}} \hat{\mathcal{A}} G^{-\frac{1}{2}} \hat{\mathcal{W}}_{\mathcal{H}}\|_{\mathbb{F}}$ ,  $G$  is the matrix expression of the metric and  $\hat{\mathcal{A}}$ ,  $\hat{\mathcal{W}}_{\mathcal{H}}$  denote the matrix expression of the operators  $\mathcal{A}$ ,  $\mathcal{W}_{\mathcal{H}}$ .

Again, a possible choice for  $\mathcal{W}_{\mathcal{H}}^2$  would be the average Hessian. A Riemannian generalization of Eq. (1.3.5) is

$$\tilde{\mathcal{G}}_k = \int_0^1 P_{x_k, \gamma_{\eta_k}(t)} \circ \text{Hess } f(\gamma_{\eta_k}(t)) \circ (P_{x_k, \gamma_{\eta_k}(t)})^{-1} dt. \tag{3.3.6}$$

where  $\gamma_{\eta_k}(t): t \mapsto \exp_{x_k} t \eta_k$ ,  $\dot{\gamma}_{\eta_k}(t) = P_{x_k, \gamma_{\eta_k}(t)}(\dot{\gamma}_{\eta_k}(0))$ ,  $\gamma_{\eta_k}(0) = x_k$ ,  $\gamma_{\eta_k}(1) = x_{k+1}$  and  $s_k = \dot{\gamma}_{\eta_k}(1)$ . This is a self-adjoint operator, which is generally not positive definite. If the parallel transport,  $P$ , and the exponential map,  $\exp$ , are used, it can be shown under certain assumptions that  $\tilde{\mathcal{G}}_k$  is positive definite and  $\tilde{\mathcal{G}}_k[s_k] = y_k$  holds.

The question remains open what is used as the first approximation. Since we want the positive-definiteness and self-adjointness to be passed on, the first approximation  $\mathcal{H}_0^{RBFGS}$  or  $\mathcal{B}_0^{RBFGS}$  should definitely fulfill these properties. Of course, for reasons of quick availability, the choice often falls on the identity operator  $\mathcal{B}_0^{RBFGS} = \text{id}_{\mathcal{T}_{x_0} \mathcal{M}}$  or on the multiplication of the tangent vector by a number  $\mathcal{B}_0^{RBFGS}[\eta] = \beta \cdot \eta$ . One approach to determine this factor  $\beta$  would be to transfer the formula from the

Euclidean, since it is defined only by internal products, i.e. one could choose

$$\beta = \frac{g_{x_1}(y_0, s_0)}{g_{x_1}(y_0, y_0)} \Rightarrow \mathcal{B}_0^{RBFGS} = \frac{g_{x_1}(y_0, s_0)}{g_{x_1}(y_0, y_0)} \cdot \text{id}_{\mathcal{T}_{x_0}\mathcal{M}}$$

and use this for the update to obtain  $\mathcal{B}_1^{RBFGS}$  after the search direction has first been calculated using  $\mathcal{B}_0^{RBFGS} = \text{id}_{\mathcal{T}_{x_0}\mathcal{M}}$ . However, it is not yet clear whether this has certain advantages over using just the identity.

### 3.4 THE BFGS METHOD ON RIEMANNIAN MANIFOLDS

We now discuss a specific BFGS method for Riemannian manifolds, which can be considered the most natural, since it uses parallel transport,  $P$ , as vector transport and the exponential map,  $\exp$ , as retraction. The choice of this combination is intrinsic because it accomplishes the respective tasks, namely the comparison of vectors and smooth differentiation on Riemannian manifolds, in the best possible way. We consider here the results of Qi, 2011, a generalization of the results of Gabay, 1982, which was also the first work to deal with the BFGS method on Riemannian manifolds.

Since the convergence analysis depends on the choice of retraction and vector transport, we define the algorithm in the corresponding selection of these maps. This leads to the following algorithm, a modification of Qi, 2011, Algorithm 2:

#### Algorithm 9 Inverse Global RBFGS Method

- 1: Given starting point  $x_0 \in \mathcal{M}$ , convergence tolerance  $\epsilon > 0$ , an initial self-adjoint and positive definite operator  $\mathcal{B}_0^{RBFGS}: \mathcal{T}_{x_0}\mathcal{M} \rightarrow \mathcal{T}_{x_0}\mathcal{M}$ ,  $k = 0$ .
- 2: **while**  $\|\text{grad } f(x_k)\|_{x_k} > \epsilon$  **do**
- 3:   Compute search direction  $\eta_k = -\mathcal{B}_k^{RBFGS}[\text{grad } f(x_k)]$ .
- 4:   Find a stepsize  $\alpha_k$  that satisfies the (strong) Wolfe conditions.
- 5:   Set  $x_{k+1} = \exp_{x_k}(\alpha_k \eta_k)$ ,  $s_k = P_{x_k, \alpha_k \eta_k}(\alpha_k \eta_k)$ ,  $y_k = \text{grad } f(x_{k+1}) - P_{x_k, \alpha_k \eta_k}(\text{operatorname{grad} } f(x_k))$ .
- 6:   Compute  $\mathcal{B}_{k+1}^{RBFGS}: \mathcal{T}_{x_{k+1}}\mathcal{M} \rightarrow \mathcal{T}_{x_{k+1}}\mathcal{M}$  by means of Eq. (3.3.4).
- 7:   Set  $k = k + 1$ .
- 8: **end while**
- 9: **return**  $x_k$

We see immediately that the general structure, as it was in Algorithm 4, has not changed. As with the euclidean BFGS algorithm, Algorithm 9 can also be formulated to work with the Hessian approximation, Eq. (3.3.3), rather than with the inverse Hessian approximation  $\mathcal{B}_{k+1}^{RBFGS}$ . This yields a mathematically equivalent algorithm. Again, the question arises which variant is more suitable in practice for the underlying problem. Eq. (3.3.4) makes it possible to cheaply compute an approximation of the inverse of the Hessian. This may make Algorithm 9 advantageous even in the case where we have a cheap exact formula for the Hessian but not for its inverse or when the cost of solving linear systems is unacceptably high Qi, 2011, p. 13.

We note that the selection of vector transport and retraction of course also influences the update

haben wir schon.



formula Eq. (3.3.4). On the one hand, the current operator  $\mathcal{B}_k^{RBFGS} \mapsto T_{x_k, \alpha_k \eta_k} \circ \mathcal{B}_k^{RBFGS} \circ (T_{x_k, \alpha_k \eta_k})^{-1}$ , the current search direction times stepsize  $\alpha_k \eta_k \mapsto \text{vectorTransportDir}_{x_k, \alpha_k \eta_k}(\alpha_k \eta_k)$  (in  $s_k$ ) and the current gradient  $\text{grad } f(x_k) \mapsto \text{vectorTransportDir}_{x_k, \alpha_k \eta_k}(\text{grad } f(x_k))$  (in  $y_k$ ) are mapped into the new tangent space using the selected vector transport  $T$ , and on the other hand, the new iterate  $x_{k+1}$  is determined by the selected retraction  $\text{retr}$ , and this affects also the definition of  $y_k$ , since  $\text{grad } f(x_{k+1})$  is used.

In Algorithm 9 it must be ensured that the curvature condition, Eq. (1.2.6), is fulfilled by choosing the stepsize  $\alpha_k > 0$  accordingly. As in the Euclidean case, we therefore demand that a stepsize  $\alpha_k$  is chosen that meets the Wolfe conditions Eq. (3.1.1) and Eq. (3.1.2). This is where the selection of the maps  $P$  and  $\exp$  takes on its meaning. Thereby the condition Eq. (3.1.2) is equivalent to the condition Eq. (3.1.3). This in turn implies that the curvature condition, Eq. (3.2.4), is fulfilled in each iteration (see Eq. (3.2.5)), leading to a positive definite update according to Lemma 3.3.1. The lemma is used to justify the update step of the algorithm and to show that it preserves the positive-definiteness and self-adjointness of all  $\mathcal{B}_k^{RBFGS}$  when the vector transport used is isometric Qi, 2011, p. 23. If a different vector transport and/or retraction is chosen, this most naturally generalized second Wolfe conditions, Eq. (3.1.2), does not guarantee the satisfaction of Eq. (3.2.4). In order for a different selection of maps to generate a positive definite update, they must fulfil additional properties (see Huang, Gallivan, Absil, 2015, Lemma 2.1). The main obstacle is that any method that uses the Riemannian second Wolfe condition, Eq. (3.1.2), will require at least the action of the differentiated retraction along some particular direction Huang, Absil, Gallivan, 2016, p. 2.

Let us now turn to the convergence analysis of Algorithm 9. In the Euclidean case a sufficient condition to achieve global convergence for a convex objective function and local superlinear convergence for a general objective function is preserving the symmetric positive-definiteness when updating the Hessian approximation  $H_k^{BFGS}$  (or its inverse  $B_k^{BFGS}$ ) that defines the step Qi, 2011, p. 20-21. As we have just seen, the operators  $\mathcal{B}_k^{RBFGS}$  generated by Algorithm 9 preserve positive-definiteness and self-adjointness. Now we have to clarify whether the same convergence results can be shown in the Riemannian setup. To show global convergence, the following assumptions must be made:

außer Übersetzen hastest du doch genau das vorne in 2.3 auch nutzen.  
**Annahme 3.4.1** (Qi, 2011, Assumptions 2.4.2.).

- (i) The objective function  $f$  is twice continuously differentiable.
- (ii) The level set  $\Omega = \{x \in \mathcal{M} : f(x) \leq f(x_0)\}$  is geodesically convex. Let  $(\mathcal{M}, g)$  be a Riemannian manifold. A subset  $C$  of  $\mathcal{M}$  is said to be a geodesically convex set if, given any two points in  $C$ , there is a geodesic arc contained within  $C$  that joins those two points.
- (iii) There exists positive constants  $n$  and  $N$  such that

$$ng_x(z, z) \leq g_x(G(x)[z], z) \leq Ng_x(z, z) \quad \text{for all } z \in \mathcal{T}_x \mathcal{M} \text{ and } x \in \Omega$$

where  $G(x)$  denotes the lifted Hessian  $G(x)[\xi] = \text{Hess } \hat{f}_x(\xi) = \text{Hess } f(\text{retr}_x \xi)$ .

We see immediately that Annahme 3.4.1 is a generalization for Riemannian manifolds of Annahme 1.4.2. Therefore it is not surprising that the following theorem about global convergence of Algorithm 9 can be seen as a generalization of Theorem 1.4.3.

**Satz 3.4.2** (Qi, 2011, Theorem 2.4.3.). Let  $x_0 \in \mathcal{M}$  be starting point for which Annahme 3.4.1 is satisfied and let  $\mathcal{B}_0^{RBFGS}$  be any linear transformation on  $\mathcal{T}_{x_0} \mathcal{M}$  that is self-adjoint and positive definite with

respect to the Riemannian metric  $g$ . The sequence  $\{x_k\}_k$  generated by [Algorithm 9](#) using parallel transport and the exponential map as the retraction converges to the minimizer  $x^*$  of  $f$ .

The convexity of the objective function  $f$  is ~~only~~ used to guarantee that there is a unique minimizer. One way for this to happen is if  $f$  is convex function for the entire domain of interest [Qi, 2011](#), p. 28. *zu  
informell* Briefly on the idea of the proof, which provides important insights. Like in the Euclidean case one is also interested to avoid establishing a bound on the condition number of the approximations  $\{\mathcal{B}_k^{RBFGS}\}_k$ . Instead one estimates the size of the largest and smallest eigenvalues of the operators and shows that the search directions and stepsizes satisfy the conditions of [Theorem 3.1.9](#). That means that the proof of [Theorem 3.4.2](#) depends on the use of parallel transport  $P$  in the definition of the average Riemannian Hessian, [Eq. \(3.3.6\)](#), which is used for the estimation of the eigenvalues and since the “exponential map version” of the Zoutendijk condition, [Theorem 3.1.9](#), is used, the line search is restricted to use the exponential map  $\exp$  as the retraction to define the next iterate  $x_{k+1}$  [Qi, 2011](#), p. 25.

[Theorem 3.4.2](#) can be used to justify two important conclusions for a more general nonconvex objective function  $f$ .

**Folgerung 3.4.3** ([Qi, 2011](#), Corollary 2.4.1.). Suppose  $f$  is a nonconvex objective function on  $\mathcal{M}$  and let  $x^* \in \mathcal{M}$  be a nondegenerate local minimizer of  $f$ , i.e.,  $\text{grad } f(x^*)$  and  $\text{Hess } f(x^*)$  is positive definite. Let  $x_0$  be starting point that is close enough to  $x^*$  so that it is in the neighborhood around  $x^*$  where the Hessian is positive definite, i.e., for which [Annahme 3.4.1](#) are satisfied and let  $\mathcal{B}_0^{RBFGS}$  be any linear transformation on  $\mathcal{T}_{x_0}\mathcal{M}$  that is self-adjoint and positive definite with respect to the Riemannian metric  $g$ .

The sequence  $\{x_k\}_k$  generated by [Algorithm 9](#) using parallel transport  $P$  and the exponential map  $\exp$  as the retraction converges to the minimizer  $x^*$  of  $f$ , i.e., it is locally convergent to any nondegenerate minimizer.

Additionally, if the convexity assumption is removed from [Annahme 3.4.1](#) then from any  $x_0$  the sequence  $\{x_k\}_k$  generated by [Algorithm 9](#) using parallel transport and the exponential map as the retraction converges to a set of critical points of  $f$ , i.e., there is global convergence to such a set.

In summary, [Theorem 3.4.2](#) guarantees us the following:

- (i) Global convergence of RBFGS using parallel transport to a unique minimizer when the objective function  $f$  is convex on the domain of interest.
- (ii) Global convergence of RBFGS using parallel transport to a set of stationary points when the objective function  $f$  is not convex on the domain of interest.
- (iii) Local convergence of RBFGS using parallel transport to a nondegenerate minimizer,  $x^*$ , when the objective function  $f$  is not convex on the domain of interest but the initial guess  $x_0$  is sufficiently close to  $x^*$  [Qi, 2011](#), p. 89.

We have seen that [Algorithm 9](#) converges globally under certain conditions. We are now interested in achieving acceptably rapid convergence rate for [Algorithm 9](#), e.g., superlinear, as it is guaranteed with [Algorithm 4](#) on  $\mathbb{R}^n$ . This requires the following assumption:

**Annahme 3.4.4** ([Qi, 2011](#), Assumptions 2.4.4.). Let  $x^* \in \mathcal{M}$  be a nondegenerate local minimizer of  $f$ ,



i.e.,  $\text{grad } f(x^*) = 0$  and  $\text{Hess } f(x^*)$  is positive definite. There is  $L > 0$  such that, for all  $\xi \in \mathcal{T}_{x^*} \mathcal{M}$  and all  $\eta \in \mathcal{T}_{\text{retr}_{x^*} \xi} \mathcal{M}$  small enough, we have

$$\|(P_{\gamma_\eta(t) \leftarrow \gamma_\eta(0)})^{-1} \circ \text{Hess } f(y) \circ P_{\gamma_\eta(t) \leftarrow \gamma_\eta(0)} - P_{\gamma_\xi(1) \leftarrow \gamma_\xi(0)} \circ \text{Hess } f(x^*) \circ (P_{\gamma_\xi(1) \leftarrow \gamma_\xi(0)})^{-1}\| \leq L \max\{\text{dist}(y, x^*), \text{dist}(x, x^*)\}$$

for  $0 \leq t \leq 1$  where  $x = \exp_{x^*} \xi$ ,  $y = \exp_x \eta$  and  $\gamma_\xi(t)$  and  $\gamma_\eta(t)$  are the associated geodesics.

Once again, we see that [Annahme 3.4.4](#) is a generalized version of [Annahme 1.4.4](#).

**Satz 3.4.5** (Qi, 2011, Theorem 2.4.5.). Suppose that  $f$  is twice continuously differentiable and that the iterates,  $x_k$ , generated by the RBFGS Algorithm using parallel transport and the exponential map converge to a nondegenerate minimizer  $x^* \in \mathcal{M}$  at which [Annahme 3.4.4](#) holds. If

$$\sum_{k=0}^{\infty} \text{dist}(x_k, x^*) < \infty \quad (3.4.1)$$

holds then  $x_k$  converges to  $x^*$  superlinearly.

To prove superlinear convergence for [Algorithm 4](#), [Theorem 1.1.7](#) was used. The same is done in the Riemannian setup. In [Theorem 3.1.10](#) a key requirement on the evolution of the action of  $\mathcal{B}_k^{\text{RBFGS}}$  in the direction of  $\eta_k$  relative to the action of the covariant derivative is identified. The requirement is quite general and only requires the transport be twice continuously differentiable. In order to apply [Theorem 3.1.10](#) for proving superlinear convergence, sufficient conditions on the vector transport and retraction used in [Algorithm 9](#) that guarantee the required action of  $\mathcal{B}_k^{\text{RBFGS}}$  must be identified Qi, 2011, p. 29.

Experiments provide substantial evidence that, in practice, both isometric and non-isometric vector transport achieve superlinear convergence with RBFGS. [Theorem 3.1.10](#) is probably a key part of the explanation for this behavior Qi, 2011, p. 20. That means a non-isometric vector transport can be used but it is not provably convergent in general, at least thus far Huang, 2013, p. 20.

works?

### 3.5 CAUTIOUS BFGS-METHOD ON RIEMANNIAN MANIFOLDS

As we have seen in the previous chapter and as confirmed by other sources (see e.g. Huang, Gallivan, Absil, 2015, Ring, Wirth, 2012), for global convergence the objective functions is required to satisfy a Riemannian version of convexity. Therefore, the idea of modifying the method to apply it to a larger class of functions is not out of the question. Again, it is not surprising that ideas which worked well in the Euclidean are transferred to the Riemannian setup.

We present a method found by [Huang, Absil, Gallivan, 2018](#) based on a Riemannian generalization of the cautious update by [Li, Fukushima, 2001](#) and a weak line search condition which was introduced in [Byrd, Nocedal, 1989](#). The motivation was to develop a method that does not need a vector transport by differentiated retraction, which allows the use of line search conditions other than the Wolfe conditions (see [Eq. \(3.1.2\)](#)), and convexity of the cost function.

First, we discuss the new conditions that our stepsize must meet. In the Euclidean BFGS setup the following rule was introduced in [Byrd, Nocedal, 1989](#) for the stepsize  $\alpha_k$  of the BFGS method. It was assumed that the procedure for choosing  $\alpha_k$  is such that at each iteration either

$$h_k(\alpha_k) - h_k(0) \leq -\chi_1 \frac{h'_k(0)^2}{\|\eta_k\|^2} \quad (3.5.1)$$

or

$$h_k(\alpha_k) - h_k(0) \leq \chi_2 h'_k(0) \quad (3.5.2)$$

holds, where  $\chi_1, \chi_2$  are positive constants and  $h_k(t) = f(x_k + td_k)$ . The motivation for this was that the analysis would cover a large class of line search strategies. We see immediately that [Eq. \(3.5.1\)](#) leads to an efficient stepsize strategy (see [Definition 1.1.3](#)). The line search conditions [Eq. \(3.5.1\)](#), [Eq. \(3.5.2\)](#) are weak, since it has been shown in [Byrd, Nocedal, 1989](#) and [Werner, 1978/79](#) that many search conditions, including the Wolfe conditions, imply either [Eq. \(3.5.1\)](#) or [Eq. \(3.5.2\)](#) if the gradient of the objective function is Lipschitz continuous (see [Eq. \(1.5.2\)](#)).

In the Riemann case, a stepsize  $\alpha_k$  is chosen which fulfills either conditions [Eq. \(3.5.1\)](#) or condition [Eq. \(3.5.2\)](#) for the function  $h_k(t) = f(\text{retr}_{x_k} t \eta_k)$ . Since  $f \circ \text{retr}_{x_k} : \mathcal{T}_{x_k} \mathcal{M} \rightarrow \mathbb{R}$  is defined on a linear space which implies that the Euclidean results about line search are applicable. That means many search conditions, including the Wolfe conditions, imply either [Eq. \(3.5.1\)](#) or [Eq. \(3.5.2\)](#) if the gradient of the objective function satisfies the Riemannian Lipschitz continuous condition in [Definition 3.1](#):

The weak line search condition removes completely the need to consider the differentiated retraction [Huang, Absil, Gallivan, 2016](#), p. 1.

$$\mathcal{H}_{k+1}^{\text{CRBFGS}}[\cdot] = \begin{cases} \tilde{\mathcal{H}}_k^{\text{CRBFGS}}[\cdot] + y_k \frac{y_k^\flat[\cdot]}{s_k^\flat[y_k]} - \tilde{\mathcal{H}}_k^{\text{CRBFGS}}[s_k] \frac{s_k^\flat(\tilde{\mathcal{H}}_k^{\text{CRBFGS}}[\cdot])}{s_k^\flat(\tilde{\mathcal{H}}_k^{\text{CRBFGS}}[s_k])} & \frac{g_{x_{k+1}}(y_k, s_k)}{\|s_k\|_{x_{k+1}}^2} \geq \theta(\|\text{grad } f(x_k)\|) \\ \tilde{\mathcal{H}}_k^{\text{CRBFGS}} & \text{otherwise.} \end{cases} \quad (3.5.3)$$

**Annahme 3.5.1** ([Huang, Absil, Gallivan, 2018](#), Assumptions 4.1.+4.2.). (i) The level set  $\Omega = \{x \in \mathcal{M} : f(x) \leq f(x_0)\}$  is compact.

(ii) The function  $f$  is Lipschitz continuously differentiable with respect to the isometric vector transport  $T$  on  $\Omega$  (see [Definition 3.1.12](#)).

**Algorithm 10** Cautious Riemannian BFGS-Algorithm

---

```

1: Riemannian manifold  $\mathcal{M}$  with Riemannian metric  $g$ ; retraction  $\text{retr}$ ; isometric vector transport  $T_{\cdot, S}^{\text{retr}}(\cdot)$ , with  $\text{retr}$  as the associated retraction; continuously differentiable real-valued function  $f$  on  $\mathcal{M}$ , bounded below; initial iterate  $x_0 \in \mathcal{M}$ ; initial Hessian approximation  $\mathcal{H}_0^{\text{CRBFGS}}$  that is symmetric positive definite with respect to the metric  $g$ ; convergence tolerance  $\epsilon > 0$ ;  $k = 0$ .
2: while  $\|\text{grad } f(x_k)\| > \epsilon$  do
3:   Obtain  $\eta_k \in \mathcal{T}_{x_k} \mathcal{M}$  by solving  $\mathcal{H}_k^{\text{CRBFGS}}[\eta_k] = -\text{grad } f(x_k)$ .
4:   Determine the stepsize  $\alpha_k > 0$  that satisfies
5:   Set  $x_{k+1} = \text{retr}_{x_k}(\alpha_k \eta_k)$ .
6:   Define  $s_k = T_{x_k, \alpha_k \eta_k}(\alpha_k \eta_k)$  and  $y_k = \text{grad } f(x_{k+1}) - T_{x_k, \alpha_k \eta_k}(\text{grad } f(x_k))$ .
7:   Define the linear operator  $\mathcal{H}_{k+1}^{\text{CRBFGS}}: \mathcal{T}_{x_{k+1}} \mathcal{M} \rightarrow \mathcal{T}_{x_{k+1}} \mathcal{M}$  by
8:   Set  $k = k + 1$ .
9: end while
10: return  $x_k$ 

```

---

**Satz 3.5.2** (Huang, Absil, Gallivan, 2018, Theorem 4.2.). *Let  $\{x_k\}_k$  be sequences generated by Algorithm 10. If the Annahme 3.5.1 hold, then*

$$\liminf_{k \rightarrow \infty} \|\text{grad } f(x_k)\| = 0.$$

As we can see, this is the Riemannian generalization of Theorem 1.5.2. The global convergence analysis of the cautious RBFGS does not require a convexity assumption on the objective function.

**Annahme 3.5.3** (Huang, Absil, Gallivan, 2018, Assumptions 5.1.). (i) *The objective function  $f$  is twice continuously differentiable in the level set  $\Omega$ .*

(ii) *The retraction  $\text{retr}$  is twice continuously differentiable.*

(iii) *The isometric vector transport  $T_{\cdot, S}^{\text{retr}}(\cdot)$  with associated retraction  $\text{retr}$  is continuous and satisfies*

**Satz 3.5.4** (Huang, Absil, Gallivan, 2018, Theorem 5.1.). *Under the assumptions of Lemma 5.1, if  $s_k \rightarrow 0$  and  $x^*$  is an accumulation point of  $\{x_k\}_k$  generated by Algorithm 10, then the sequence  $\{x_k\}_k$  converges to  $x^*$  and Algorithm 10 reduces to the ordinary RBFGS when  $\{x_k\}_k$  is sufficiently close to  $x^*$ .*

In summary, Algorithm 10 provides a method for which a vector transport by differentiated retraction is not required, global convergence for nonconvex objective functions and superlinear convergence are guaranteed.

### 3.6 LIMITED-MEMORY BFGS-METHOD ON RIEMANNIAN MANIFOLDS

$$T_{x, S_\xi}^{\text{retr}}(\xi) = \beta T_{x, \xi}^{\text{retr}}(\xi), \quad \beta = \frac{\|\xi\|}{\|T_{x, \xi}^{\text{retr}}(\xi)\|} \quad (3.6.1)$$

$B_k^{(0)}$  can be an arbitrarily, symmetrical and positive definite matrix. In general  $B_k^{(0)}$  will be a multiple of the identity matrix, so that it can be stored very easily Geiger, Kanzow, 1999, p. 198. A method for choosing  $B_k^{(0)}$  that has proven effective in practice is to set  $B_k^{(0)} = \gamma_k I$ , where

$$\gamma_k = \frac{s_{k-1}^T y_{k-1}}{y_{k-1}^T y_{k-1}}. \quad (3.6.2)$$

$\gamma_k$  is the scaling factor that attempts to estimate the size of the true Hessian matrix along the most recent search direction. This choice helps to ensure that the search direction  $d_k$  is well scaled, and as a result the stepsize  $\alpha_k = 1$  is accepted in most iterations. It is important that the line search is based on the (strict) Wolfe conditions, so that the BFGS updating is stable Nocedal, Wright, 2006, p. 178-179.

The following algorithm is a modified version of Huang, but follows the same structures:

---

**Algorithm 11** Limited-Memory Riemannian BFGS-Algorithm

---

```

1: Riemannian manifold  $\mathcal{M}$  with Riemannian metric  $g$ , a retraction  $\text{retr}$ , isometric vector transport  $\text{T}_{\cdot, S}^{\text{retr}}(\cdot)$  that satisfies Eq. (3.6.1), smooth function  $f$  on  $\mathcal{M}$ , initial iterate  $x_0 \in \mathcal{M}$ , an integer  $m > 0$ ,  $k = 0$ ,  $\epsilon > 0$ ,  $0 < c_1 < \frac{1}{2} < c_2 < 1$ ,  $\gamma_0 = 1$ ,  $l = 0$ 
2: while  $\|\text{grad } f(x_k)\| > \epsilon$  do
3:    $\mathcal{B}_k^{(0)} = \gamma_k \text{id}$ . Obtain  $\eta_k \in \mathcal{T}_{x_k} \mathcal{M}$  by the following algorithm:
4:    $q = \text{grad } f(x_k)$ 
5:   for  $i = k - 1, k - 2, \dots, k - m$  do
6:      $\alpha_i = \rho_i s_i^T q$ 
7:      $q = q - \alpha_i y_i$ 
8:   end for
9:    $r = \mathcal{B}_k^{(0)}[q]$ 
10:  for  $i = k - m, k - m + 1, \dots, k - 1$  do
11:     $\beta = \rho_i y_i^T r$ 
12:     $r = r + s_i(\alpha_i - \beta)$ 
13:  end for
14:  Set  $\eta_k = -r$ 
15:  Find  $\alpha_k$  that satisfies Wolfe conditions Eq. (3.1.1) and Eq. (3.1.2) (or Eq. (3.1.3)).
16:  Set  $x_{k+1} = \text{retr}_{x_k} \alpha_k \eta_k$ .
17:  Set  $s_k = \text{T}_{x_k, S, \alpha_k \eta_k}^{\text{retr}}(\alpha_k \eta_k)$ ,  $\beta = \frac{\|\alpha_k \eta_k\|}{\|\text{T}_{x_k, \alpha_k \eta_k}^{\text{retr}}(\alpha_k \eta_k)\|}$ ,  $y_k = \beta^{-1} \text{grad } f(x_{k+1}) - \text{T}_{x_k, S, \alpha_k \eta_k}^{\text{retr}}(\text{grad } f(x_k))$ .
18:
19:  if  $k > m$  then
20:    Discard the vector pairs  $\{s_{k-m}, y_{k-m}\}$  from storage.
21:  end if
22:  Save  $s_k$  and  $y_k$ .
23:  Set  $k = k + 1$ .
24: end while
    
```

---

The following theorem is a generalization of Sun, Yuan, 2006, Theorem 5.7.4.

**Satz 3.6.1.** *Let  $f: \mathcal{M} \rightarrow \mathbb{R}$  be a twice continuously differentiable and uniformly convex function. Then the iterative sequence  $\{x_k\}_k$  generated by the L-RBFGS-method ([Algorithm 11](#)) converges to the unique minimizer  $x^*$  of  $f$ .*

The following theorem is a generalization of [Sun, Yuan, 2006](#), Theorem 5.7.7.

**Satz 3.6.2.** *Let  $f: \mathcal{M} \rightarrow \mathbb{R}$  be a twice continuously differentiable and uniformly convex function. Assume that the iterative sequence  $\{x_k\}_k$  generated by the L-RBFGS-method (??) converges to the unique minimizer  $x^*$  of  $f$ . Then the rate of convergence is at least  $R$ -linear.*

## 4 NUMERICS

A practical implementation of RBFSG requires the following ingredients:

- (i) an efficient numerical representation for points  $x$  on  $\mathcal{M}$ , tangent spaces  $\mathcal{T}_x\mathcal{M}$  and the inner products  $g_x(\xi_1, \xi_2)$  on  $\mathcal{T}_x\mathcal{M}$
- (ii) an implementation of the chosen retraction  $\text{retr}_x: \mathcal{T}_x\mathcal{M} \rightarrow \mathcal{M}$
- (iii) efficient formulae for  $f(x)$  and  $\text{grad } f(x)$
- (iv) an implementation of the chosen vector transport  $T_{x,\eta}$  and its inverse  $(T_{x,\eta})^{-1}$
- (v) a method for solving  $\mathcal{H}_k^{\text{RBFSG}}[\eta_k] = -\text{grad } f(x_k)$  or alternatively, a method for computing  $\eta_k = -\mathcal{H}_k^{\text{RBFSG}}[\text{grad } f(x_k)]$ .
- (vi) Locking condition

Qi, 2011

### 4.1 REALIZING THE UPDATE-FORMULA

#### Hessian Inverse Approximation

Let  $\mathcal{M}$  be a Riemannian manifold with  $\dim(\mathcal{M}) = n$ . Let  $\{e_1, \dots, e_n\}$  be a canonical basis of  $\mathcal{T}_{x_{k+1}}\mathcal{M}$ . Let  $\mathcal{B}_{k+1}^{\text{RBFSG}} = \{\mathbb{I}_1^{k+1}, \dots, \mathbb{I}_n^{k+1}\}$ . Let  $\eta \in \mathcal{T}_{x_{k+1}}\mathcal{M}$ .

$$\mathcal{B}_{k+1}^{\text{RBFSG}}[\eta] = g_{x_{k+1}}(\mathbb{I}_1^{k+1}, \eta) \cdot e_1 + \dots + g_{x_{k+1}}(\mathbb{I}_n^{k+1}, \eta) \cdot e_n.$$

$$\mathbb{I}_i^{k+1} = \mathbb{I}_i^k - \frac{g_{x_{k+1}}(e_i, s_k)}{g_{x_{k+1}}(s_k, y_k)} \cdot \tilde{\mathcal{B}}_k^{\text{RBFSG}}[y_k] - \frac{g_{x_{k+1}}(\tilde{\mathcal{B}}_k^{\text{RBFSG}}[y_k], e_i)}{g_{x_{k+1}}(s_k, y_k)} \cdot s_k + \frac{g_{x_{k+1}}(y_k, \tilde{\mathcal{B}}_k^{\text{RBFSG}}[y_k]) \cdot g_{x_{k+1}}(s_k, e_i)}{(g_{x_{k+1}}(s_k, y_k))^2} \cdot s_k + \frac{g_{x_{k+1}}(s_k, e_i)}{g_{x_{k+1}}(s_k, y_k)} \cdot s_k$$

for  $i = 1, \dots, n$ . We define

$$\tilde{\mathcal{B}}_k^{\text{RBFSG}}[y_k] = g_{x_{k+1}}(\mathbb{I}_1^k, y_k) \cdot e_1 + \dots + g_{x_{k+1}}(\mathbb{I}_n^k, y_k) \cdot e_n$$

where  $\widetilde{\lfloor}_i^k = P_{x_{k+1} \leftarrow x_k}(\lfloor_i^k)$ .

### Hessian Approximation

Let  $\mathcal{M}$  be a Riemannian manifold with  $\dim(\mathcal{M}) = n$ . Let  $\{e_1, \dots, e_n\}$  be a canonical basis of  $\mathcal{T}_{x_{k+1}}\mathcal{M}$ . Let  $\mathcal{H}_{k+1}^{RBFGS} = \{\langle_1^{k+1}, \dots, \langle_n^{k+1}\rangle$ . Let  $\eta \in \mathcal{T}_{x_{k+1}}\mathcal{M}$ .

$$\mathcal{H}_{k+1}^{RBFGS}[\eta] = g_{x_{k+1}}(\langle_1^{k+1}, \eta) \cdot e_1 + \dots + g_{x_{k+1}}(\langle_n^{k+1}, \eta) \cdot e_n.$$

$$\langle_i^{k+1} = \widetilde{\langle}_i^k - \frac{g_{x_{k+1}}(\widetilde{\langle}_i^k, s_k)}{g_{x_{k+1}}(s_k, \widetilde{\mathcal{H}}_k^{RBFGS}[s_k])} \cdot \widetilde{\mathcal{H}}_k^{RBFGS}[s_k] + \frac{g_{x_{k+1}}(e_i, y_k)}{g_{x_{k+1}}(y_k, s_k)} \cdot y_k$$

for  $i = 1, \dots, n$ . We define

$$\widetilde{\mathcal{H}}_k^{RBFGS}[s_k] = g_{x_{k+1}}(\widetilde{\langle}_1^k, s_k) \cdot e_1 + \dots + g_{x_{k+1}}(\widetilde{\langle}_n^k, s_k) \cdot e_n$$

where  $\widetilde{\langle}_i^k = P_{x_{k+1} \leftarrow x_k}(\langle_i^k)$ .

see [Qi, 2011](#), Chapter 3.3

## 5 CONCLUSION



## LITERATURE

## LITERATUR

- Absil, P.-A.; R. Mahony; R. Sepulchre (2008). *Optimization Algorithms on Matrix Manifolds*. Princeton University Press.
- Broyden, C. G. (1967). "Quasi-Newton Methods and their Application to Function Minimisation". *Mathematics of Computation* 21.99, S. 368–381.
- Byrd, R. H.; J. Nocedal (1989). "A Tool for the Analysis of Quasi-Newton Methods with Application to Unconstrained Minimization". *SIAM Journal on Numerical Analysis* 26.3, S. 727–739.
- Cruz Neto, J.; I. Melo; P. Sousa (2017). "Convexity and Some Geometric Properties". *Journal of Optimization Theory and Applications* 173. DOI: [10.1007/s10957-017-1087-2](https://doi.org/10.1007/s10957-017-1087-2).
- Dai, Y.-H. (2002). "Convergence Properties of the BFGS Algorithm". *SIAM Journal on Optimization* 13, S. 693–701. DOI: [10.1137/S1052623401383455](https://doi.org/10.1137/S1052623401383455).
- Dai, Y.-H. (2012). "A perfect example for the BFGS method". *Mathematical Programming* 138. DOI: [10.1007/s10107-012-0522-2](https://doi.org/10.1007/s10107-012-0522-2).
- Davidon, W. C. (1959). "VARIABLE METRIC METHOD FOR MINIMIZATION". DOI: [10.2172/4252678](https://doi.org/10.2172/4252678).
- Deng, C. (2011). "A generalization of the Sherman-Morrison-Woodbury formula". *Appl. Math. Lett.* 24, S. 1561–1564. DOI: [10.1016/j.aml.2011.03.046](https://doi.org/10.1016/j.aml.2011.03.046).
- Dennis, J.; J. J. Moré (1974). "Quasi-Newton Methods, Motivation and Theory". *Siam Review* 19, S. 46–89.
- Fletcher, R. (1970). "A new approach to variable metric algorithms". *The Computer Journal* 13.3, S. 317–322. DOI: [10.1093/comjnl/13.3.317](https://doi.org/10.1093/comjnl/13.3.317).
- Fletcher, R.; M. J. D. Powell (1963). "A Rapidly Convergent Descent Method for Minimization". DOI: [10.1093/comjnl/6.2.163](https://doi.org/10.1093/comjnl/6.2.163).
- Gabay, D. (1982). "Minimizing a differentiable function over a differential manifold". *Journal of Optimization Theory and Applications* 37. DOI: [10.1007/BF00934767](https://doi.org/10.1007/BF00934767).
- Geiger, C.; C. Kanzow (1999). *Numerische Verfahren zur Lösung unrestringierter Optimierungsaufgaben*. Springer. DOI: [10.1007/978-3-642-58582-1](https://doi.org/10.1007/978-3-642-58582-1).
- Goldfarb, D. (1970). "A family of variable-metric methods derived by variational means".
- Huang, W. (2013). *Optimization Algorithms on Riemannian Manifolds with Applications*. Florida State University.
- Huang, W.; P.-A. Absil; K. A. Gallivan (2018). "A Riemannian BFGS Method without Differentiated Retraction for Nonconvex Optimization Problems". *SIAM Journal on Optimization* 28.1, S. 470–495.
- Huang, W.; P.-A. Absil; K. A. Gallivan (2016). "A Riemannian BFGS Method for Nonconvex Optimization Problems". *Numerical Mathematics and Advanced Applications ENUMATH 2015*. Hrsg. von B. Karasözen; M. Manguoğlu; M. Tezer-Sezgin; S. Göktepe; Ö. Uğur. Cham: Springer International Publishing, S. 627–634.
- Huang, W.; K. A. Gallivan; P.-A. Absil (2015). "A Broyden Class of Quasi-Newton Methods for Riemannian Optimization". *SIAM Journal on Optimization* 25.3, S. 1660–1685. DOI: [10.1137/140955483](https://doi.org/10.1137/140955483).
- Li, D.-H.; M. Fukushima (2001). "On the Global Convergence of the BFGS Method for Nonconvex Unconstrained Optimization Problems". *SIAM Journal on Optimization* 11.4, S. 1054–1064. DOI: [10.1137/s1052623499354242](https://doi.org/10.1137/s1052623499354242).

- Mascarenhas, W. (2004). "The BFGS method with exact line searches fails for non-convex objective functions". *Math. Program.* 99, S. 49–61. DOI: [10.1007/s10107-003-0421-7](#).
- Nocedal, J.; S. J. Wright (2006). *Numerical Optimization*. Second Edition. Springer. DOI: [10.1007/978-0-387-40065-5](#).
- Powell, M. D. (1976). "Some global convergence properties of a variable metric algorithm for minimization without exact lin".
- Qi, C. (2011). "Numerical Optimization Methods on Riemannian Manifolds". Florida State University.
- Ring, W.; B. Wirth (2012). "Optimization Methods on Riemannian Manifolds and Their Application to Shape Space". *SIAM Journal on Optimization* 22. DOI: [10.1137/11082885X](#).
- Sato, H.; T. Iwai (2015). "A new, globally convergent Riemannian conjugate gradient method". *Optimization* 64.4, S. 1011–1031. DOI: [10.1080/02331934.2013.836650](#).
- Shanno, D. F. (1970). "Conditioning of Quasi-Newton Methods for Function Minimization". *Mathematics of Computation* 24.111, S. 647–656.
- Sun, W.; Y.-X. Yuan (2006). *Optimization Theory and Methods: Nonlinear Programming*. Bd. 1. Springer. DOI: [10.1007/b106451](#).
- Ulbrich, M.; S. Ulbrich (2012). *Nichtlineare Optimierung*. Springer. DOI: [10.1007/978-3-0346-0654-7](#).
- Warth, W.; J. Werner (1977). "Effiziente Schrittweitenfunktionen bei unrestringierten Optimierungsaufgaben". *Computing* 19, S. 59–72. DOI: [10.1007/BF02260741](#).
- Werner, J. (1978/79). "Über die globale Konvergenz von Variable-Metrik-Verfahren mit nicht-exakter Schrittweitenbestimmung." *Numerische Mathematik* 31, S. 321–334.
- Zhang, H.; S. Sra (2016). "First-order Methods for Geodesically Convex Optimization". *29th Annual Conference on Learning Theory*. Hrsg. von V. Feldman; A. Rakhlin; O. Shamir. Bd. 49. Proceedings of Machine Learning Research. Columbia University, New York, New York, USA: PMLR, S. 1617–1638.