

## Chapter Seven

---

### Trust-Region Methods

The plain Newton method discussed in Chapter 6 was shown to be locally convergent to any critical point of the cost function. The method does not distinguish among local minima, saddle points, and local maxima: all (nondegenerate) critical points are asymptotically stable fixed points of the Newton iteration. Moreover, it is possible to construct cost functions and initial conditions for which the Newton sequence does not converge. There even exist examples where the set of nonconverging initial conditions contains an open subset of search space.

To exploit the desirable superlinear local convergence properties of the Newton algorithm in the context of global optimization, it is necessary to embed the Newton update in some form of descent method. In Chapter 6 we briefly outlined how the Newton equation can be used to generate a descent direction that is used in a line-search algorithm. Such an approach requires modification of the Newton equation to ensure that the resulting sequence of search directions is gradient-related and an implementation of a standard line-search iteration. The resulting algorithm will converge to critical points of the cost function for *all* initial points. Moreover, saddle points and local maxima are rendered unstable, thus favoring convergence to local minimizers.

Trust-region methods form an alternative class of algorithms that combine desirable global convergence properties with a local superlinear rate of convergence. In addition to providing good global convergence, trust-region methods also provide a framework to relax the computational burden of the plain Newton method when the iterates are too far away from the solution for fast local convergence to set in. This is particularly important in the development of optimization algorithms on matrix manifolds, where the inverse Hessian computation can involve solving relatively complex matrix equations.

Trust-region methods can be understood as an enhancement of Newton's method. To this end, however, we need to consider this method from another viewpoint: instead of looking for an update vector along which the derivative of  $\text{grad } f$  is equal to  $-\text{grad } f(x_k)$ , it is equivalent to think of Newton's method (in  $\mathbb{R}^n$ ) as the algorithm that selects the new iterate  $x_{k+1}$  to be the critical point of the quadratic Taylor expansion of the cost function  $f$  about  $x_k$ .

To this end, the chapter begins with a discussion of generalized quadratic models on manifolds (Section 7.1). Here again, a key role is played by the concept of retraction, which provides a way to pull back the cost function on

the manifold to a cost function on the tangent space. It is therefore sufficient to define quadratic models on abstract vector spaces and to understand how these models correspond to the real-valued function on the manifold  $\mathcal{M}$ .

Once the notion of a quadratic model is established, a trust-region algorithm can be defined on a manifold (Section 7.2). It is less straightforward to show that all the desirable convergence properties of classical trust-region methods in  $\mathbb{R}^n$  still hold, *mutatis mutandis*, for their manifold generalizations. The difficulty comes from the fact that trust-region methods on manifolds do not work with a single cost function but rather with a succession of cost functions whose domains are different tangent spaces. The issue of computing an (approximate but sufficiently accurate) solution of the trust-region subproblems is discussed in Section 7.3. The convergence analysis is carried out in Section 7.4. The chapter is concluded in Section 7.5 with a “checklist” of steps one has to go through in order to turn the abstract geometric trust-region schemes into practical numerical algorithms on a given manifold for a given cost function; this checklist is illustrated for several examples related to Rayleigh quotient minimization.

## 7.1 MODELS

Several classical optimization schemes rely on successive local minimization of quadratic models of the cost function. In this section, we review the notion of quadratic models in  $\mathbb{R}^n$  and in general vector spaces. Then, making use of retractions, we extend the concept to Riemannian manifolds.

### 7.1.1 Models in $\mathbb{R}^n$

The fundamental mathematical tool that justifies the use of local models is Taylor’s theorem (see Appendix A.6). In particular, we have the following results.

**Proposition 7.1.1** *Let  $f$  be a smooth real-valued function on  $\mathbb{R}^n$ ,  $x \in \mathbb{R}^n$ ,  $\mathcal{U}$  a bounded neighborhood of  $x$ , and  $H$  any symmetric matrix. Then there exists  $c > 0$  such that, for all  $(x + h) \in \mathcal{U}$ ,*

$$\|f(x + h) - (f(x) + \partial f(x)h + \frac{1}{2}h^T Hh)\| \leq c\|h\|^2,$$

where  $\partial f(x) := (\partial_1 f(x), \dots, \partial_n f(x))$ . If, moreover,  $H_{i,j} = \partial_i \partial_j f(x)$ , then there exists  $c > 0$  such that, for all  $(x + h) \in \mathcal{U}$ ,

$$\|f(x + h) - (f(x) + \partial f(x)h + \frac{1}{2}h^T Hh)\| \leq c\|h\|^3.$$

### 7.1.2 Models in general Euclidean spaces

The first step towards defining quadratic models on Riemannian manifolds is to generalize the above results to (abstract) Euclidean spaces, i.e., finite-dimensional vector spaces endowed with an inner product  $\langle \cdot, \cdot \rangle$ . This is readily done using the results in Appendix A.6. (Note that Euclidean spaces

are naturally finite-dimensional normed vector spaces for the induced norm  $\|\xi\| := \sqrt{\langle \xi, \xi \rangle}$ .

**Proposition 7.1.2** *Let  $f$  be a smooth real-valued function on a Euclidean space  $\mathcal{E}$ ,  $x \in \mathcal{E}$ ,  $\mathcal{U}$  a bounded neighborhood of  $x$ , and  $H : \mathcal{E} \rightarrow \mathcal{E}$  any symmetric operator. Then there exists  $c > 0$  such that, for all  $(x+h) \in \mathcal{U}$ ,*

$$\|f(x+h) - (f(x) + \langle \text{grad } f(x), h \rangle + \frac{1}{2} \langle H[h], h \rangle)\| \leq c \|h\|^2.$$

*If, moreover,  $H = \text{Hess } f(x) : h \mapsto \text{D}(\text{grad } f)(x)[h]$ , then there exists  $c > 0$  such that, for all  $(x+h) \in \mathcal{U}$ ,*

$$\|f(x+h) - (f(x) + \langle \text{grad } f(x), h \rangle + \frac{1}{2} \langle H[h], h \rangle)\| \leq c \|h\|^3.$$

### 7.1.3 Models on Riemannian manifolds

Let  $f$  be a real-valued function on a Riemannian manifold  $\mathcal{M}$  and let  $x \in \mathcal{M}$ . A model  $m_x$  of  $f$  around  $x$  is a real-valued function defined on a neighborhood of  $x$  such that (i)  $m_x$  is a “sufficiently good” approximation of  $f$  and (ii)  $m_x$  has a “simple” form that makes it easy to tackle with optimization algorithms.

The quality of the model  $m_x$  is assessed by evaluating how the discrepancy between  $m_x(y)$  and  $f(y)$  evolves as a function of the Riemannian distance  $\text{dist}(x, y)$  between  $x$  and  $y$ . The model  $m_x$  is an *order- $q$  model*,  $q > 0$ , if there exists a neighborhood  $\mathcal{U}$  of  $x$  in  $\mathcal{M}$  and a constant  $c > 0$  such that

$$|f(y) - m_x(y)| \leq c (\text{dist}(x, y))^{q+1} \quad \text{for all } y \in \mathcal{U}.$$

Note that an order- $q$  model automatically satisfies  $m_x(x) = f(x)$ . The following result shows that the order of a model can be assessed using any retraction  $R$  on  $\mathcal{M}$ .

**Proposition 7.1.3** *Let  $f$  be a real-valued function on a Riemannian manifold  $\mathcal{M}$  and let  $x \in \mathcal{M}$ . A model  $m_x$  of  $f$  is order- $q$  if and only if there exists a neighborhood  $\mathcal{U}$  of  $x$  and a constant  $c > 0$  such that*

$$|f(y) - m_x(y)| \leq c \|R_x^{-1}(y)\|^{q+1} \quad \text{for all } y \in \mathcal{U}.$$

*Proof.* In view of the local rigidity property of retractions, it follows that  $\text{D}(\text{Exp}^{-1} \circ R_x)(0_x) = \text{id}_{T_x \mathcal{M}}$ , hence  $\|R_x^{-1}(y)\| = \|(\text{Exp}_x^{-1} \circ R_x^{-1})(\text{Exp}_x y)\| = \Omega(\|\text{Exp}_x y\|) = \Omega(\text{dist}(x, y))$ ; see Section A.4 for the definition of the asymptotic notation  $\Omega$ . In other words, there is a neighborhood  $\mathcal{U}$  of  $x$  and constants  $c_1$  and  $c_2$  such that

$$c_1 \text{dist}(x, y) \leq \|R_x^{-1}(y)\| \leq c_2 \text{dist}(x, y) \quad \text{for all } y \in \mathcal{U}.$$

□

Proposition 7.1.3 yields a conceptually simple way to build an order- $q$  model of  $f$  around  $x$ . Pick a retraction on  $\mathcal{M}$ . Consider

$$\widehat{f}_x := f \circ R_x,$$

the pullback of  $f$  to  $T_x\mathcal{M}$  through  $R_x$ . Let  $\widehat{m}_x$  be the order- $q$  Taylor expansion of  $\widehat{f}_x$  around the origin  $0_x$  of  $T_x\mathcal{M}$ . Finally, push  $\widehat{m}_x$  forward through  $R_x$  to obtain  $m_x = \widehat{m}_x \circ R_x^{-1}$ . This model is an order- $q$  model of  $f$  around  $x$ . In particular, the obvious order-2 model to choose is

$$\begin{aligned}\widehat{m}_x(\xi) &= \widehat{f}_x(0_x) + D\widehat{f}_x(0_x)[\xi] + \frac{1}{2}D^2\widehat{f}_x(0_x)[\xi, \xi] \\ &= f(x) + \langle \text{grad } f(x), \xi \rangle + \frac{1}{2}\langle \text{Hess } \widehat{f}_x(0_x)[\xi], \xi \rangle.\end{aligned}$$

Note that, since  $DR_x(0_x) = \text{id}_{T_x\mathcal{M}}$  (see Definition 4.1.1), it follows that  $Df(x) = D\widehat{f}_x(0_x)$ , hence  $\text{grad } \widehat{f}_x(0_x) = \text{grad } f(x)$ , where  $\text{grad } f(x)$  denotes the (Riemannian) gradient of  $f$  (see Section 3.6).

In practice, the second-order differentials of the function  $\widehat{f}_x$  may be difficult to compute. (The first-order differential is always straightforward since the rigidity condition of the retraction ensures  $D\widehat{f}_x(0_x) = Df(x)$ .) On the other hand, the Riemannian connection  $\nabla$  admits nice formulations on Riemannian submanifolds and Riemannian quotient manifolds (see Section 5.3). This suggests the model

$$m_x = \widehat{m}_x \circ R_x^{-1}$$

with

$$\widehat{m}_x(\xi) = f(x) + \langle \text{grad } f(x), \xi \rangle + \frac{1}{2}\langle \text{Hess } f(x)[\xi], \xi \rangle, \quad \xi \in T_x\mathcal{M}, \quad (7.1)$$

where the quadratic term is given by the Riemannian Hessian

$$\text{Hess } f(x)[\xi] = \nabla_\xi \text{grad } f(x). \quad (7.2)$$

In general, this model  $m_x$  is only order 1 because  $\text{Hess } f(x) \neq \text{Hess } \widehat{f}_x(0_x)$ . However, if  $R$  is a second-order retraction, then  $\text{Hess } f(x) = \text{Hess } \widehat{f}_x(0_x)$  (Proposition 5.5.5) and  $m_x$  is order 2. More importantly, for any retraction,  $\text{Hess } f(x_*) = \text{Hess } \widehat{f}_{x_*}(0_{x_*})$  when  $x_*$  is a critical point of  $f$  (Proposition 5.5.6).

The quadratic model (7.1) has a close connection to the Newton algorithm. Assuming that the Hessian is nonsingular at  $x$ , the critical point  $\xi_*$  of  $\widehat{m}_x$  satisfies the Newton equation

$$\text{Hess } f(x)[\eta_*] + \text{grad } f(x) = 0.$$

It follows that the geometric Newton method (Algorithm 5), with retraction  $R$  and affine connection  $\nabla$ , defines its next iterate as the critical point of the quadratic model (7.1).

We point out that all the models we have considered so far assume a quadratic form on  $T_x\mathcal{M}$ , i.e., there is a symmetric operator  $H : T_x\mathcal{M} \rightarrow T_x\mathcal{M}$  such that

$$\widehat{m}_x(\xi) = f(x) + \langle \text{grad } f(x), \xi \rangle + \frac{1}{2}\langle H[\xi], \xi \rangle.$$

Quadratic models are particularly interesting because the problem of minimizing a quadratic function under trust-region constraints is well understood and several algorithms are available (see Section 7.3).

## 7.2 TRUST-REGION METHODS

We first briefly review the principles of trust-region methods in  $\mathbb{R}^n$ . Extending the concept to manifolds is straightforward given the material developed in the preceding section.

### 7.2.1 Trust-region methods in $\mathbb{R}^n$

The basic trust-region method in  $\mathbb{R}^n$  for a cost function  $f$  consists of adding to the current iterate  $x \in \mathbb{R}^n$  the update vector  $\eta \in \mathbb{R}^n$ , solving (up to some approximation) the *trust-region subproblem*

$$\min_{\eta \in \mathbb{R}^n} m(\eta) = f(x) + \partial f(x)\eta + \frac{1}{2}\eta^T H\eta, \quad \|\eta\| \leq \Delta, \quad (7.3)$$

where  $H$  is some symmetric matrix and  $\Delta$  is the trust-region radius. Clearly, a possible choice for  $H$  is the Hessian matrix  $H_{i,j} = \partial_i \partial_j f(x)$ ; classical convergence results guarantee superlinear convergence if the chosen  $H$  is “sufficiently close” to the Hessian matrix. The algorithm used to compute an approximate minimizer  $\eta$  of the model within the trust region is termed the *inner iteration*. Once an  $\eta$  has been returned by the inner iteration, the quality of the model  $m$  is assessed by forming the quotient

$$\rho = \frac{f(x) - f(x + \eta)}{m(0) - m(\eta)}. \quad (7.4)$$

Depending on the value of  $\rho$ , the new iterate  $x + \eta$  is accepted or discarded and the trust-region radius  $\Delta$  is updated. A specific procedure (in the Riemannian setting) is given in Algorithm 10, or see the textbooks mentioned in Notes and References.

### 7.2.2 Trust-region methods on Riemannian manifolds

We can now lay out the structure of a trust-region method on a Riemannian manifold  $(\mathcal{M}, g)$  with retraction  $R$ . Given a cost function  $f : \mathcal{M} \rightarrow \mathbb{R}$  and a current iterate  $x_k \in \mathcal{M}$ , we use  $R_{x_k}$  to locally map the minimization problem for  $f$  on  $\mathcal{M}$  into a minimization problem for the *pullback* of  $f$  under  $R_{x_k}$ ,

$$\widehat{f}_{x_k} : T_{x_k}\mathcal{M} \rightarrow \mathbb{R} : \xi \mapsto f(R_{x_k}\xi). \quad (7.5)$$

The Riemannian metric  $g$  turns  $T_{x_k}\mathcal{M}$  into a Euclidean space endowed with the inner product  $g_{x_k}(\cdot, \cdot)$ —which we usually denote by  $\langle \cdot, \cdot \rangle_{x_k}$ —and the trust-region subproblem on  $T_{x_k}\mathcal{M}$  reads

$$\begin{aligned} \min_{\eta \in T_{x_k}\mathcal{M}} \widehat{m}_{x_k}(\eta) &= f(x_k) + \langle \text{grad } f(x_k), \eta \rangle + \frac{1}{2} \langle H_k[\eta], \eta \rangle, \\ \text{subject to } \langle \eta, \eta \rangle_{x_k} &\leq \Delta_k^2, \end{aligned} \quad (7.6)$$

where  $H_k$  is some symmetric operator on  $T_{x_k}\mathcal{M}$ . (A possible choice is  $H_k := \text{Hess } f(x_k)$ , the Riemannian Hessian (7.2).) This is called the *trust-region subproblem*.

Next, an (approximate) solution  $\eta_k$  of the Euclidean trust-region subproblem (7.6) is computed using any available method (inner iteration). The candidate for the new iterate is then given by  $R_{x_k}(\eta_k)$ . Notice that the inner iteration has to operate on the Euclidean space  $T_{x_k}\mathcal{M}$  which, in the case of embedded submanifolds and quotient manifolds of a matrix space  $\mathbb{R}^{n \times p}$ , is represented as a linear subspace of  $\mathbb{R}^{n \times p}$ . Fortunately, most (or even all?) classical algorithms for the trust-region subproblem are readily adapted to Euclidean matrix spaces (see Section 7.3 for details).

The decisions on accepting or rejecting the candidate  $R_{x_k}(\eta_k)$  and on selecting the new trust-region radius  $\Delta_{k+1}$  are based on the quotient

$$\rho_k = \frac{f(x_k) - f(R_{x_k}(\eta_k))}{\widehat{m}_{x_k}(0_{x_k}) - \widehat{m}_{x_k}(\eta_k)} = \frac{\widehat{f}_{x_k}(0_{x_k}) - \widehat{f}_{x_k}(\eta_k)}{\widehat{m}_{x_k}(0_{x_k}) - \widehat{m}_{x_k}(\eta_k)}. \quad (7.7)$$

If  $\rho_k$  is exceedingly small, then the model is very inaccurate: the step must be rejected, and the trust-region radius must be reduced. If  $\rho_k$  is small but less dramatically so, then the step is accepted but the trust-region radius is reduced. If  $\rho_k$  is close to 1, then there is a good agreement between the model and the function over the step, and the trust-region radius can be expanded. If  $\rho_k \gg 1$ , then the model is inaccurate, but the overall optimization iteration is producing a significant decrease in the cost. In this situation a possible strategy is to increase the trust region in the hope that your luck will hold and that bigger steps will result in a further decrease in the cost, regardless of the quality of the model approximation. This procedure is formalized in Algorithm 10.

Later in the chapter we sometimes drop the subscript  $k$  and denote  $x_{k+1}$  by  $x_+$ .

In general, there is no assumption on the operator  $H_k$  in (7.6) other than being a symmetric linear operator. Consequently, the choice of the retraction  $R$  does not impose any constraint on  $\widehat{m}_{x_k}$ . In order to achieve superlinear convergence, however,  $H_k$  must approximate the Hessian (Theorem 7.4.11). The issue of obtaining an approximate Hessian in practice is addressed in Section 7.5.1.

### 7.3 COMPUTING A TRUST-REGION STEP

Step 2 in Algorithm 10 computes an (approximate) solution of the trust-region subproblem (7.6),

$$\begin{aligned} \min_{\eta \in T_{x_k}\mathcal{M}} \widehat{m}_{x_k}(\eta) &= f(x_k) + \langle \text{grad } f(x_k), \eta \rangle + \frac{1}{2} \langle H_k[\eta], \eta \rangle, \\ \text{subject to } \langle \eta, \eta \rangle_{x_k} &\leq \Delta_k^2. \end{aligned}$$

Methods for solving trust-region subproblems in the  $\mathbb{R}^n$  case can be roughly classified into two broad classes: (i) methods (based on the work of Moré and Sorensen) that compute nearly exact solutions of the subproblem; (ii) methods that compute an approximate solution to the trust-region subproblem

**Algorithm 10** Riemannian trust-region (RTR) meta-algorithm

**Require:** Riemannian manifold  $(\mathcal{M}, g)$ ; scalar field  $f$  on  $\mathcal{M}$ ; retraction  $R$  from  $T\mathcal{M}$  to  $\mathcal{M}$  as in Definition 4.1.1.

Parameters:  $\bar{\Delta} > 0$ ,  $\Delta_0 \in (0, \bar{\Delta})$ , and  $\rho' \in [0, \frac{1}{4})$ .

**Input:** Initial iterate  $x_0 \in \mathcal{M}$ .

**Output:** Sequence of iterates  $\{x_k\}$ .

```

1: for  $k = 0, 1, 2, \dots$  do
2:   Obtain  $\eta_k$  by (approximately) solving (7.6);
3:   Evaluate  $\rho_k$  from (7.7);
4:   if  $\rho_k < \frac{1}{4}$  then
5:      $\Delta_{k+1} = \frac{1}{4}\Delta_k$ ;
6:   else if  $\rho_k > \frac{3}{4}$  and  $\|\eta_k\| = \Delta_k$  then
7:      $\Delta_{k+1} = \min(2\Delta_k, \bar{\Delta})$ ;
8:   else
9:      $\Delta_{k+1} = \Delta_k$ ;
10:  end if
11:  if  $\rho_k > \rho'$  then
12:     $x_{k+1} = R_x \eta_k$ ;
13:  else
14:     $x_{k+1} = x_k$ ;
15:  end if
16: end for

```

using a computationally simple iteration that achieves at least the decrease in cost obtained for the Cauchy point. For the trust-region subproblem (7.6), assuming that  $\text{grad } f(x_k) \neq 0$ , we define the *Cauchy point* as the solution  $\eta_k^C$  of the one-dimensional problem

$$\eta_k^C = \arg \min_{\eta} \{ \widehat{m}_{x_k}(\eta) : \eta = -\tau \text{grad } f(x_k), \tau > 0, \|\eta\| \leq \Delta_k \}, \quad (7.8)$$

which reduces to the classical definition of the Cauchy point when  $\mathcal{M} = \mathbb{R}^n$ . The *Cauchy decrease* is given by  $\widehat{m}_{x_k}(0) - \widehat{m}_{x_k}(\eta_k^C)$ . We present a brief outline of the first method before providing a more detailed development of a conjugate gradient algorithm for the second approach that we prefer for the later developments.

### 7.3.1 Computing a nearly exact solution

The following statement is a straightforward adaptation of a result of Moré and Sorensen to the case of the trust-region subproblem on  $T_{x_k}\mathcal{M}$  as expressed in (7.6).

**Proposition 7.3.1** *The vector  $\eta^*$  is a global solution of the trust-region subproblem (7.6) if and only if there exists a scalar  $\mu \geq 0$  such that the*

following conditions are satisfied:

$$(H_k + \mu \text{id})\eta^* = -\text{grad } f(x_k), \quad (7.9a)$$

$$\mu(\Delta_k - \|\eta^*\|) = 0, \quad (7.9b)$$

$$(H_k + \mu \text{id}) \text{ is positive-semidefinite}, \quad (7.9c)$$

$$\|\eta^*\| \leq \Delta_k. \quad (7.9d)$$

This result suggests a strategy for computing the solution of the subproblem (7.6). Either solve (7.9) with  $\mu = 0$  or define

$$\eta(\mu) := -(H_k + \mu \text{id})^{-1} \text{grad } f(x_k)$$

and adjust  $\mu$  to achieve  $\|\eta(\mu)\| = \Delta_k$ . Several algorithms have been proposed to perform this task; see Notes and References.

### 7.3.2 Improving on the Cauchy point

In many applications, the dimension  $d$  of the manifold  $\mathcal{M}$  is extremely large (see, for example, computation of an invariant subspace of a large matrix  $A$  in Section 7.5). In such cases, solving the linear system (7.9a) of size  $d$  or checking the positive-definiteness of a  $d \times d$  matrix (7.9c) is unfeasible. Many algorithms exist that scale down the precision of the solution of the trust-region subproblem (7.6) and lighten the numerical burden.

A number of these methods start by computing the Cauchy point (7.8), and then attempt to improve on it. The improvement strategy is often designed so that, when  $H_k$  is positive-definite, and given sufficient iterations, the estimate eventually reaches the minimizer  $\eta_k^N = (H_k)^{-1} \text{grad } f(x_k)$  provided that the minimizer lies within the trust region. Among these strategies, the *truncated conjugate-gradient method* is one of the most popular. Algorithm 11 is a straightforward adaptation of the truncated CG method in  $\mathbb{R}^n$  to the trust-region subproblem (7.6) in  $T_{x_k}\mathcal{M}$ . Note that we use superscripts to denote the evolution of  $\eta$  within the inner iteration, while subscripts are used in the outer iteration.

Several comments about Algorithm 11 are in order.

The following result will be useful in the convergence analysis of trust-region methods.

**Proposition 7.3.2** *Let  $\eta^i$ ,  $i = 0, \dots, j$ , be the iterates generated by Algorithm 11 (truncated CG method). Then  $\widehat{m}_{x_k}(\eta^i)$  is strictly decreasing and  $\widehat{m}_{x_k}(\eta_k) \leq \widehat{m}_{x_k}(\eta^i)$ ,  $i = 0, \dots, j$ . Further,  $\|\eta^i\|$  is strictly increasing and  $\|\eta_k\| > \|\eta^i\|$ ,  $i = 0, \dots, j$ .*

The simplest stopping criterion to use in Step 14 of Algorithm 11 is to truncate after a fixed number of iterations. In order to achieve superlinear convergence (see Section 7.4.2), one may take the stopping criterion

$$\|r_{j+1}\| \leq \|r_0\| \min(\|r_0\|^\theta, \kappa), \quad (7.10)$$

where  $\theta > 0$  is a real parameter chosen in advance.



---

**Algorithm 11** Truncated CG (tCG) method for the trust-region subproblem

---

**Goal:** This algorithm handles Step 2 of Algorithm 10.

```

1: Set  $\eta^0 = 0$ ,  $r_0 = \text{grad } f(x_k)$ ,  $\delta_0 = -r_0$ ;  $j = 0$ ;
2: loop
3:   if  $\langle \delta_j, H_k \delta_j \rangle_{x_k} \leq 0$  then
4:     Compute  $\tau$  such that  $\eta = \eta^j + \tau \delta_j$  minimizes  $\widehat{m}_{x_k}(\eta)$  in (7.6) and
       satisfies  $\|\eta\|_{g_x} = \Delta$ ;
5:     return  $\eta_k := \eta$ ;
6:   end if
7:   Set  $\alpha_j = \langle r_j, r_j \rangle_{x_k} / \langle \delta_j, H_k \delta_j \rangle_{x_k}$ ;
8:   Set  $\eta^{j+1} = \eta^j + \alpha_j \delta_j$ ;
9:   if  $\|\eta^{j+1}\|_{g_x} \geq \Delta$  then
10:    Compute  $\tau \geq 0$  such that  $\eta = \eta^j + \tau \delta_j$  satisfies  $\|\eta\|_{g_x} = \Delta$ ;
11:    return  $\eta_k := \eta$ ;
12:   end if
13:   Set  $r_{j+1} = r_j + \alpha_j H_k \delta_j$ ;
14:   if a stopping criterion is satisfied then
15:     return  $\eta_k := \eta^{j+1}$ ;
16:   end if
17:   Set  $\beta_{j+1} = \langle r_{j+1}, r_{j+1} \rangle_{x_k} / \langle r_j, r_j \rangle_{x_k}$ ;
18:   Set  $\delta_{j+1} = -r_{j+1} + \beta_{j+1} \delta_j$ ;
19:   Set  $j = j + 1$ ;
20: end loop

```

---

In Steps 4 and 10 of Algorithm 11,  $\tau$  is found by computing the positive root of the quadratic equation

$$\tau^2 \langle \delta_j, \delta_j \rangle_{x_k} + 2\tau \langle \eta^j, \delta_j \rangle_{x_k} = \Delta_k^2 - \langle \eta^j, \eta^j \rangle_{x_k}.$$

Notice that the truncated CG algorithm is “inverse-free”, as it uses  $H_k$  only in the computation of  $H_k[\delta_j]$ .

Practical implementations of Algorithm 11 usually include several additional features to reduce the numerical burden and improve the robustness to numerical errors. For example, the value of  $\langle r_{j+1}, r_{j+1} \rangle_{x_k}$  can be stored since it will be needed at the next iteration. Because the Hessian operator  $H_k$  is an operator on a vector space of dimension  $d$ , where  $d$  may be very large, it is important to implement an efficient routine for computing  $H_k \delta$ . In many practical cases, the tangent space  $T_x \mathcal{M}$  to which the quantities  $\eta$ ,  $r$ , and  $\delta$  belong will be represented as a linear subspace of a higher-dimensional Euclidean space; to prevent numerical errors it may be useful from time to time to reproject the above quantities onto the linear subspace.

## 7.4 CONVERGENCE ANALYSIS

In this section, we study the global convergence properties of the RTR scheme (Algorithm 10) without any assumption on the way the trust-region subproblems are solved (Step 2), except that the approximate solution  $\eta_k$  must produce a decrease in the model that is at least a fixed fraction of the Cauchy decrease. Under mild additional assumptions on the retraction and on the cost function, it is shown that the sequences  $\{x_k\}$  produced by Algorithm 10 converge to the set of critical points of the cost function. This result is well known in the  $\mathbb{R}^n$  case; in the case of manifolds, the convergence analysis has to address the fact that a different lifted cost function  $\hat{f}_{x_k}$  is considered at each iterate  $x_k$ .

In the second part of the section we analyze the local convergence of Algorithm 10-11 around nondegenerate local minima. Algorithm 10-11 refers to the RTR framework where the trust-region subproblems are approximately solved using the truncated CG algorithm with stopping criterion (7.10). It is shown that the iterates of the algorithm converge to nondegenerate critical points with an order of convergence of at least  $\min\{\theta + 1, 2\}$ , where  $\theta$  is the parameter chosen for the stopping condition (7.10).

### 7.4.1 Global convergence

The objective of this section is to show that, under appropriate assumptions, the sequence  $\{x_k\}$  generated by Algorithm 10 converges to the critical set of the cost function; this generalizes a classical convergence property of trust-region methods in  $\mathbb{R}^n$ . In what follows,  $(\mathcal{M}, g)$  is a Riemannian manifold of dimension  $d$  and  $R$  is a retraction on  $\mathcal{M}$  (Definition 4.1.1). We define the pullback cost

$$\hat{f} : T\mathcal{M} \mapsto \mathbb{R} : \xi \mapsto f(R\xi) \quad (7.11)$$

and, in accordance with (7.5),  $\hat{f}_x$  denotes the restriction of  $\hat{f}$  to  $T_x\mathcal{M}$ . We denote by  $B_\delta(0_x) = \{\xi \in T_x\mathcal{M} : \|\xi\| < \delta\}$  the open ball in  $T_x\mathcal{M}$  of radius  $\delta$  centered at  $0_x$ , and  $B_\delta(x)$  stands for the set  $\{y \in \mathcal{M} : \text{dist}(x, y) < \delta\}$ , where  $\text{dist}$  denotes the Riemannian distance (i.e., the distance defined in terms of the Riemannian metric; see Section 3.6). We denote by  $P_\gamma^{t \leftarrow t_0} v$  the vector of  $T_{\gamma(t)}\mathcal{M}$  obtained by parallel translation (with respect to the Riemannian connection) of the vector  $v \in T_{\gamma(t_0)}\mathcal{M}$  along a curve  $\gamma$ .

As in the classical  $\mathbb{R}^n$  proof, we first show that at least one accumulation point of  $\{x_k\}$  is a critical point of  $f$ . The convergence result requires that  $\hat{m}_{x_k}(\eta_k)$  be a sufficiently good approximation of  $\hat{f}_{x_k}(\eta_k)$ . In classical proofs, this is often guaranteed by the assumption that the Hessian of the cost function is bounded. It is, however, possible to weaken this assumption, which leads us to consider the following definition.

**Definition 7.4.1 (radially L- $C^1$  function)** *Let  $\hat{f} : T\mathcal{M} \rightarrow \mathbb{R}$  be defined as in (7.11). We say that  $\hat{f}$  is radially Lipschitz continuously differentiable*

if there exist reals  $\beta_{RL} > 0$  and  $\delta_{RL} > 0$  such that, for all  $x \in \mathcal{M}$ , for all  $\xi \in T_x \mathcal{M}$  with  $\|\xi\| = 1$ , and for all  $t < \delta_{RL}$ , it holds that

$$\left| \frac{d}{d\tau} \widehat{f}_x(\tau\xi)|_{\tau=t} - \frac{d}{d\tau} \widehat{f}_x(\tau\xi)|_{\tau=0} \right| \leq \beta_{RL} t. \quad (7.12)$$

For the purpose of Algorithm 10, which is a descent algorithm, this condition needs only to be imposed for all  $x$  in the level set

$$\{x \in M : f(x) \leq f(x_0)\}. \quad (7.13)$$

A key assumption in the classical global convergence result in  $\mathbb{R}^n$  is that the approximate solution  $\eta_k$  of the trust-region subproblem (7.6) produces at least as much decrease in the model function as a fixed fraction of the Cauchy decrease. The definition (7.8) of the Cauchy point is equivalent to the closed-form definition  $\eta_k^C = -\tau_k \text{grad } f(x_k)$  with

$$\tau_k = \begin{cases} \frac{\Delta_k}{\|\text{grad } f(x_k)\|} & \text{if } \langle H_k[\text{grad } f(x_k)], \text{grad } f(x_k) \rangle_{x_k} \leq 0, \\ \frac{\Delta_k}{\|\text{grad } f(x_k)\|} \min \left( \frac{\|\text{grad } f(x_k)\|^3}{\Delta_k \langle H_k[\text{grad } f(x_k)], \text{grad } f(x_k) \rangle_{x_k}}, 1 \right) & \text{otherwise.} \end{cases}$$

(Note that the definition of the Cauchy point excludes the case  $\text{grad } f(x_k) = 0$ , for which convergence to a critical point becomes trivial.) The assumption on the decrease in  $f$  then becomes

$$\widehat{m}_{x_k}(0) - \widehat{m}_{x_k}(\eta_k) \geq c_1 \|\text{grad } f(x_k)\| \min \left( \Delta_k, \frac{\|\text{grad } f(x_k)\|}{\|H_k\|} \right), \quad (7.14)$$

for some constant  $c_1 > 0$ , where  $\|H_k\|$  is defined as

$$\|H_k\| := \sup\{\|H_k \zeta\| : \zeta \in T_{x_k} \mathcal{M}, \|\zeta\| = 1\}. \quad (7.15)$$

In particular, the Cauchy point satisfies (7.14) with  $c_1 = \frac{1}{2}$ . Hence the tangent vector  $\eta_k$  returned by the truncated CG method (Algorithm 11) satisfies (7.14) with  $c_1 = \frac{1}{2}$  since the truncated CG method first computes the Cauchy point and then attempts to improve the model decrease.

With these ideas in place, we can state and prove the first global convergence result. Note that this theorem is presented under weak assumptions; stronger but arguably easier to check assumptions are given in Proposition 7.4.5.

**Theorem 7.4.2** *Let  $\{x_k\}$  be a sequence of iterates generated by Algorithm 10 with  $\rho' \in [0, \frac{1}{4})$ . Suppose that  $f$  is  $C^1$  and bounded below on the level set (7.13), that  $\widehat{f}$  is radially  $L$ - $C^1$  (Definition 7.4.1), and that there is a constant  $\beta$  such that  $\|H_k\| \leq \beta$  for all  $k$ . Further suppose that all  $\eta_k$ 's obtained in Step 2 of Algorithm 10 satisfy the Cauchy decrease inequality (7.14) for some positive constant  $c_1$ . We then have*

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

*Proof.* From the definition of the ratio  $\rho_k$  in (7.7), we have

$$|\rho_k - 1| = \left| \frac{\widehat{m}_{x_k}(\eta_k) - \widehat{f}_{x_k}(\eta_k)}{\widehat{m}_{x_k}(0) - \widehat{m}_{x_k}(\eta_k)} \right|. \quad (7.16)$$

Proposition A.6.1 (Taylor) applied to the function  $t \mapsto \widehat{f}_{x_k}(t \frac{\eta_k}{\|\eta_k\|})$  yields

$$\begin{aligned} \widehat{f}_{x_k}(\eta_k) &= \widehat{f}_{x_k}(0_{x_k}) + \|\eta_k\| \left. \frac{d}{d\tau} \widehat{f}_{x_k}\left(\tau \frac{\eta_k}{\|\eta_k\|}\right) \right|_{\tau=0} + \epsilon' \\ &= f(x_k) + \langle \text{grad } f(x_k), \eta_k \rangle_{x_k} + \epsilon', \end{aligned}$$

where  $|\epsilon'| = \frac{1}{2}\beta_{RL}\|\eta_k\|^2$  whenever  $\|\eta_k\| < \delta_{RL}$  and  $\beta_{RL}$  and  $\delta_{RL}$  are the constants in the radially L- $C^1$  property (7.12). Therefore, it follows from the definition (7.6) of  $\widehat{m}_{x_k}$  that

$$\begin{aligned} |\widehat{m}_{x_k}(\eta_k) - \widehat{f}_{x_k}(\eta_k)| &= \left| \frac{1}{2} \langle H_k \eta_k, \eta_k \rangle_{x_k} - \epsilon' \right| \\ &\leq \frac{1}{2} \beta \|\eta_k\|^2 + \frac{1}{2} \beta_{RL} \|\eta_k\|^2 \leq \beta' \|\eta_k\|^2 \end{aligned} \quad (7.17)$$

whenever  $\|\eta_k\| < \delta_{RL}$ , where  $\beta' = \max(\beta, \beta_{RL})$ .

Assume for contradiction that the claim does not hold; i.e., assume there exist  $\epsilon > 0$  and a positive index  $K$  such that

$$\|\text{grad } f(x_k)\| \geq \epsilon \quad \text{for all } k \geq K. \quad (7.18)$$

From (7.14), for  $k \geq K$ , we have

$$\begin{aligned} \widehat{m}_{x_k}(0) - \widehat{m}_{x_k}(\eta_k) &\geq c_1 \|\text{grad } f(x_k)\| \min\left(\Delta_k, \frac{\|\text{grad } f(x_k)\|}{\|H_k\|}\right) \\ &\geq c_1 \epsilon \min\left(\Delta_k, \frac{\epsilon}{\beta'}\right). \end{aligned} \quad (7.19)$$

Substituting (7.17) and (7.19) into (7.16), we have that

$$|\rho_k - 1| \leq \frac{\beta' \|\eta_k\|^2}{c_1 \epsilon \min\left(\Delta_k, \frac{\epsilon}{\beta'}\right)} \leq \frac{\beta' \Delta_k^2}{c_1 \epsilon \min\left(\Delta_k, \frac{\epsilon}{\beta'}\right)} \quad (7.20)$$

whenever  $\|\eta_k\| < \delta_{RL}$ . Let  $\widehat{\Delta}$  be defined as  $\widehat{\Delta} = \min\left(\frac{c_1 \epsilon}{2\beta'}, \frac{\epsilon}{\beta'}, \delta_{RL}\right)$ . If  $\Delta_k \leq \widehat{\Delta}$ , then  $\min\left(\Delta_k, \frac{\epsilon}{\beta'}\right) = \Delta_k$  and (7.20) becomes

$$|\rho_k - 1| \leq \frac{\beta' \widehat{\Delta} \Delta_k}{c_1 \epsilon \min\left(\Delta_k, \frac{\epsilon}{\beta'}\right)} \leq \frac{\Delta_k}{2 \min\left(\Delta_k, \frac{\epsilon}{\beta'}\right)} = \frac{1}{2}.$$

Therefore,  $\rho_k \geq \frac{1}{2} > \frac{1}{4}$  whenever  $\Delta_k \leq \widehat{\Delta}$ , so that by the workings of Algorithm 10, it follows (from the argument above) that  $\Delta_{k+1} \geq \Delta_k$  whenever  $\Delta_k \leq \widehat{\Delta}$ . It follows that a reduction in  $\Delta_k$  (by a factor of  $\frac{1}{4}$ ) can occur in Algorithm 10 only when  $\Delta_k > \widehat{\Delta}$ . Therefore, we conclude that

$$\Delta_k \geq \min\left(\Delta_K, \widehat{\Delta}/4\right) \quad \text{for all } k \geq K. \quad (7.21)$$

Suppose now that there is an infinite subsequence  $\mathcal{K}$  such that  $\rho_k \geq \frac{1}{4} > \rho'$  for  $k \in \mathcal{K}$ . If  $k \in \mathcal{K}$  and  $k \geq K$ , we have from (7.19) that

$$\begin{aligned} f(x_k) - f(x_{k+1}) &= f_{x_k} - \hat{f}_{x_k}(\eta_k) \\ &\geq \frac{1}{4}(\hat{m}_{x_k}(0) - \hat{m}_{x_k}(\eta_k)) \geq \frac{1}{4}c_1\epsilon \min\left(\Delta_k, \frac{\epsilon}{\beta'}\right). \end{aligned} \quad (7.22)$$

Since  $f$  is bounded below on the level set containing these iterates, it follows from this inequality that  $\lim_{k \in \mathcal{K}, k \rightarrow \infty} \Delta_k = 0$ , clearly contradicting (7.21). Then such an infinite subsequence as  $\mathcal{K}$  cannot exist. It follows that we must have  $\rho_k < \frac{1}{4}$  for all  $k$  sufficiently large so that  $\Delta_k$  will be reduced by a factor of  $\frac{1}{4}$  on every iteration. Then we have  $\lim_{k \rightarrow \infty} \Delta_k = 0$ , which again contradicts (7.21). Hence our assumption (7.18) is false, and the proof is complete.  $\square$

To further show that all accumulation points of  $\{x_k\}$  are critical points, we need to make an additional regularity assumption on the cost function  $f$ . The convergence result in  $\mathbb{R}^n$  requires that  $f$  be Lipschitz continuously differentiable. That is, for any  $x, y \in \mathbb{R}^n$ ,

$$\|\text{grad } f(y) - \text{grad } f(x)\| \leq \beta_1 \|y - x\|. \quad (7.23)$$

A key to obtaining a Riemannian counterpart of this global convergence result is to adapt the notion of being Lipschitz continuously differentiable to the Riemannian manifold  $(\mathcal{M}, g)$ . The expression  $\|x - y\|$  on the right-hand side of (7.23) naturally becomes the Riemannian distance  $\text{dist}(x, y)$ . For the left-hand side of (7.23), observe that the operation  $\text{grad } f(x) - \text{grad } f(y)$  is not well defined in general on a Riemannian manifold since  $\text{grad } f(x)$  and  $\text{grad } f(y)$  belong to two different tangent spaces, namely,  $T_x\mathcal{M}$  and  $T_y\mathcal{M}$ . However, if  $y$  belongs to a normal neighborhood of  $x$ , then there is a unique geodesic  $\alpha(t) = \text{Exp}_x(t \text{Exp}_x^{-1} y)$  in this neighborhood such that  $\alpha(0) = x$  and  $\alpha(1) = y$ , and we can parallel-translate  $\text{grad } f(y)$  along  $\alpha$  to obtain the vector  $P_\alpha^{0 \leftarrow 1} \text{grad } f(y)$  in  $T_x\mathcal{M}$ . A lower bound on the size of the normal neighborhoods is given by the *injectivity radius*, defined as

$$i(\mathcal{M}) := \inf_{x \in \mathcal{M}} i_x,$$

where

$$i_x := \sup\{\epsilon > 0 : \text{Exp}_x|_{B_\epsilon(0_x)} \text{ is a diffeomorphism for all } x \in \mathcal{M}\}.$$

This yields the following definition.

**Definition 7.4.3 (Lipschitz continuously differentiable)** *Assume that  $(\mathcal{M}, g)$  has a positive injectivity radius. A real function  $f$  on  $\mathcal{M}$  is Lipschitz continuously differentiable if it is differentiable and if there exists  $\beta_1$  such that, for all  $x, y$  in  $\mathcal{M}$  with  $\text{dist}(x, y) < i(\mathcal{M})$ , it holds that*

$$\|P_\alpha^{0 \leftarrow 1} \text{grad } f(y) - \text{grad } f(x)\| \leq \beta_1 \text{dist}(y, x), \quad (7.24)$$

where  $\alpha$  is the unique minimizing geodesic with  $\alpha(0) = x$  and  $\alpha(1) = y$ .

Note that (7.24) is symmetric in  $x$  and  $y$ ; indeed, since the parallel transport is an isometry, it follows that

$$\|P_\alpha^{0\leftarrow 1} \text{grad } f(y) - \text{grad } f(x)\| = \|\text{grad } f(y) - P_\alpha^{1\leftarrow 0} \text{grad } f(x)\|.$$

Moreover, we place one additional requirement on the retraction  $R$ , that there exist  $\mu > 0$  and  $\delta_\mu > 0$  such that

$$\|\xi\| \geq \mu \text{dist}(x, R_x \xi) \quad \text{for all } x \in \mathcal{M}, \text{ for all } \xi \in T_x \mathcal{M}, \|\xi\| \leq \delta_\mu. \quad (7.25)$$

Note that for the exponential retraction, (7.25) is satisfied as an equality with  $\mu = 1$ . The bound is also satisfied when  $\mathcal{M}$  is compact (Corollary 7.4.6).

We are now ready to show that under some additional assumptions, the gradient of the cost function converges to zero on the whole sequence of iterates. Here again we refer to Proposition 7.4.5 for a simpler (but slightly stronger) set of assumptions that yield the same result.

**Theorem 7.4.4** *Let  $\{x_k\}$  be a sequence of iterates generated by Algorithm 10. Suppose that all the assumptions of Theorem 7.4.2 are satisfied. Further suppose that  $\rho' \in (0, \frac{1}{4})$ , that  $f$  is Lipschitz continuously differentiable (Definition 7.4.3), and that (7.25) is satisfied for some  $\mu > 0$ ,  $\delta_\mu > 0$ . It then follows that*

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

*Proof.* Consider any index  $m$  such that  $\text{grad } f(x_m) \neq 0$ . Define the scalars

$$\epsilon = \frac{1}{2} \|\text{grad } f(x_m)\|, \quad r = \min \left( \frac{\|\text{grad } f(x_m)\|}{2\beta_1}, i(M) \right) = \min \left( \frac{\epsilon}{\beta_1}, i(M) \right).$$

In view of the Lipschitz property (7.24), we have for all  $x \in B_r(x_m)$ ,

$$\begin{aligned} \|\text{grad } f(x)\| &= \|P_\alpha^{0\leftarrow 1} \text{grad } f(x)\| \\ &= \|P_\alpha^{0\leftarrow 1} \text{grad } f(x) + \text{grad } f(x_m) - \text{grad } f(x_m)\| \\ &\geq \|\text{grad } f(x_m)\| - \|P_\alpha^{0\leftarrow 1} \text{grad } f(x) - \text{grad } f(x_m)\| \\ &\geq 2\epsilon - \beta_1 \text{dist}(x, x_m) \\ &> 2\epsilon - \beta_1 \min \left( \frac{\|\text{grad } f(x_m)\|}{2\beta_1}, i(M) \right) \\ &\geq 2\epsilon - \frac{1}{2} \|\text{grad } f(x_m)\| \\ &= \epsilon. \end{aligned}$$

If the entire sequence  $\{x_k\}_{k \geq m}$  stays inside the ball  $B_r(x_m)$ , then we have  $\|\text{grad } f(x_k)\| > \epsilon$  for all  $k \geq m$ , a contradiction to Theorem 7.4.2. Thus the sequence eventually leaves the ball  $B_r(x_m)$ . Let the index  $l \geq m$  be such that  $x_{l+1}$  is the first iterate after  $x_m$  outside  $B_r(x_m)$ . Since  $\|\text{grad } f(x_k)\| > \epsilon$  for

$k = m, m+1, \dots, l$ , we have, in view of the Cauchy decrease condition (7.14),

$$\begin{aligned}
 f(x_m) - f(x_{l+1}) &= \sum_{k=m}^l f(x_k) - f(x_{k+1}) \\
 &\geq \sum_{k=m, x_k \neq x_{k+1}}^l \rho' (\widehat{m}_{x_k}(0) - \widehat{m}_{x_k}(\eta_k)) \\
 &\geq \sum_{k=m, x_k \neq x_{k+1}}^l \rho' c_1 \|\text{grad } f(x_k)\| \min \left( \Delta_k, \frac{\|\text{grad } f(x_k)\|}{\|B_k\|} \right) \\
 &\geq \sum_{k=m, x_k \neq x_{k+1}}^l \rho' c_1 \epsilon \min \left( \Delta_k, \frac{\epsilon}{\beta} \right).
 \end{aligned}$$

We distinguish two cases. If  $\Delta_k > \epsilon/\beta$  in at least one of the terms of the sum, then

$$f(x_m) - f(x_{l+1}) \geq \rho' c_1 \epsilon \frac{\epsilon}{\beta}. \quad (7.26)$$

In the other case, we have

$$f(x_m) - f(x_{l+1}) \geq \rho' c_1 \epsilon \sum_{k=m, x_k \neq x_{k+1}}^l \Delta_k \geq \rho' c_1 \epsilon \sum_{k=m, x_k \neq x_{k+1}}^l \|\eta_k\|. \quad (7.27)$$

If  $\|\eta_k\| > \delta_\mu$  in at least one term in the sum, then

$$f(x_m) - f(x_{l+1}) \geq \rho' c_1 \epsilon \delta_\mu. \quad (7.28)$$

Otherwise, (7.27) yields

$$\begin{aligned}
 f(x_m) - f(x_{l+1}) &\geq \rho' c_1 \epsilon \sum_{k=m, x_k \neq x_{k+1}}^l \mu \text{dist}(x_k, R_{x_k} \eta_k) \\
 &= \rho' c_1 \epsilon \mu \sum_{k=m, x_k \neq x_{k+1}}^l \text{dist}(x_k, x_{k+1}) \\
 &\geq \rho' c_1 \epsilon \mu r = \rho' c_1 \epsilon \mu \min \left( \frac{\epsilon}{\beta_1}, i(M) \right).
 \end{aligned} \quad (7.29)$$

It follows from (7.26), (7.28), and (7.29) that

$$f(x_m) - f(x_{l+1}) \geq \rho' c_1 \epsilon \min \left( \frac{\epsilon}{\beta}, \delta_\mu, \frac{\epsilon \mu}{\beta_1}, i(M) \mu \right). \quad (7.30)$$

Because  $\{f(x_k)\}_{k=0}^\infty$  is decreasing and bounded below, we have

$$f(x_k) \downarrow f^* \quad (7.31)$$

for some  $f^* > -\infty$ . It then follows from (7.30) that

$$\begin{aligned} f(x_m) - f^* &\geq f(x_m) - f(x_{l+1}) \\ &\geq \rho' c_1 \epsilon \min \left( \frac{\epsilon}{\beta}, \delta_\mu, \frac{\epsilon \mu}{\beta_1}, i(M) \mu \right) \\ &= \frac{1}{2} \rho' c_1 \|\text{grad } f(x_m)\| \\ &\quad \min \left( \frac{\|\text{grad } f(x_m)\|}{2\beta}, \delta_\mu, \frac{\|\text{grad } f(x_m)\| \mu}{2\beta_1}, i(M) \mu \right). \end{aligned}$$

Taking  $m \rightarrow \infty$  in the latter expression yields  $\lim_{m \rightarrow \infty} \|\text{grad } f(x_m)\| = 0$ .  $\square$

Note that this theorem reduces gracefully to the classical  $\mathbb{R}^n$  case, taking  $\mathcal{M} = \mathbb{R}^n$  endowed with the classical inner product and  $R_x \xi := x + \xi$ . Then  $i(M) = +\infty > 0$ ,  $R$  satisfies (7.25), and the Lipschitz condition (7.24) reduces to the classical expression, which subsumes the radially L- $C^1$  condition.

The following proposition shows that the regularity conditions on  $f$  and  $\widehat{f}$  required in the previous theorems are satisfied under stronger but possibly easier to check conditions. These conditions impose a bound on the Hessian of  $f$  and on the “acceleration” along curves  $t \mapsto R(t\xi)$ . Note also that all these conditions need only be checked on the level set  $\{x \in \mathcal{M} : f(x) \leq f(x_0)\}$ .

**Proposition 7.4.5** *Suppose that  $\|\text{grad } f(x)\| \leq \beta_g$  and that  $\|\text{Hess } f(x)\| \leq \beta_H$  for some constants  $\beta_g, \beta_H$ , and all  $x \in \mathcal{M}$ . Moreover, suppose that*

$$\left\| \frac{D}{dt} \frac{d}{dt} R(t\xi) \right\| \leq \beta_D \tag{7.32}$$

for some constant  $\beta_D$ , for all  $\xi \in T\mathcal{M}$  with  $\|\xi\| = 1$  and all  $t < \delta_D$ , where  $\frac{D}{dt}$  denotes the covariant derivative along the curve  $t \mapsto R(t\xi)$ . Then the Lipschitz- $C^1$  condition on  $f$  (Definition 7.4.3) is satisfied with  $\beta_L = \beta_H$ ; the radially Lipschitz- $C^1$  condition on  $\widehat{f}$  (Definition 7.4.1) is satisfied for  $\delta_{RL} < \delta_D$  and  $\beta_{RL} = \beta_H(1 + \beta_D \delta_D) + \beta_g \beta_D$ ; and the condition (7.25) on  $R$  is satisfied for values of  $\mu$  and  $\delta_\mu$  satisfying  $\delta_\mu < \delta_D$  and  $\frac{1}{2} \beta_D \delta_\mu < \frac{1}{\mu} - 1$ .

*Proof.* By a standard Taylor argument (see Lemma 7.4.7), boundedness of the Hessian of  $f$  implies the Lipschitz- $C^1$  property of  $f$ .

For (7.25), define  $u(t) = R(t\xi)$  and observe that

$$\text{dist}(x, R(t\xi)) \leq \int_0^t \|u'(\tau)\| d\tau$$

where  $\int_0^t \|u'(\tau)\| d\tau$  is the length of the curve  $u$  between 0 and  $t$ . Using the Cauchy-Schwarz inequality and the invariance of the metric by the connection, we have

$$\begin{aligned} \left| \frac{d}{d\tau} \|u'(\tau)\| \right| &= \left| \frac{d}{d\tau} \sqrt{\langle u'(\tau), u'(\tau) \rangle_{u(\tau)}} \right| = \left| \frac{\langle \frac{D}{dt} u'(\tau), u'(\tau) \rangle_{u(\tau)}}{\|u'(\tau)\|} \right| \\ &\leq \frac{\beta_D \|u'(\tau)\|}{\|u'(\tau)\|} \leq \beta_D \end{aligned}$$



for all  $t < \delta_D$ . Therefore

$$\int_0^t \|u'(\tau)\| \, d\tau \leq \int_0^t \|u'(0)\| + \beta_D \tau \, d\tau = \|\xi\|t + \frac{1}{2}\beta_D t^2 = t + \frac{1}{2}\beta_D t^2,$$

which is smaller than  $\frac{t}{\mu}$  if  $\frac{1}{2}\beta_D t < \frac{1}{\mu} - 1$ .

For the radially Lipschitz- $C^1$  condition, let  $u(t) = R(t\xi)$  and  $h(t) = f(u(t)) = \hat{f}(t\xi)$  with  $\xi \in T_x\mathcal{M}$ ,  $\|\xi\| = 1$ . Then

$$h'(t) = \langle \text{grad } f(u(t)), u'(t) \rangle_{u(t)}$$

and

$$\begin{aligned} h''(t) &= \frac{D}{dt} \langle \text{grad } f(u(t)), u'(t) \rangle_{u(t)} \\ &= \left\langle \frac{D}{dt} \text{grad } f(u(t)), u'(t) \right\rangle_{u(t)} + \langle \text{grad } f(u(t)), \frac{D}{dt} u'(t) \rangle_{u(t)}. \end{aligned}$$

Now,  $\frac{D}{dt} \text{grad } f(u(t)) = \nabla_{u'(t)} \text{grad } f(u(t)) = \text{Hess } f(u(t))[u'(t)]$ . It follows that  $|h''(t)|$  is bounded on  $t \in [0, \delta_D)$  by the constant  $\beta_{RL} = \beta_H(1 + \beta_D\delta_D) + \beta_g\beta_D$ . Then

$$|h'(t) - h'(0)| \leq \int_0^t |h''(\tau)| \, d\tau \leq t\beta_{RL}.$$

□

**Corollary 7.4.6 (smoothness and compactness)** *If the cost function  $f$  is smooth and the Riemannian manifold  $\mathcal{M}$  is compact, then all the conditions in Proposition 7.4.5 are satisfied.*

The major manifolds considered in this book (the Grassmann manifold and the Stiefel manifold) are compact, and the cost functions based on the Rayleigh quotient are smooth.

## 7.4.2 Local convergence

We now state local convergence properties of Algorithm 10-11: local convergence to local minimizers (Theorem 7.4.10) and superlinear convergence (Theorem 7.4.11). We begin with a few preparatory lemmas.

As before,  $(\mathcal{M}, g)$  is a Riemannian manifold of dimension  $d$  and  $R$  is a retraction on  $\mathcal{M}$  (Definition 4.1.1). The first lemma is a first-order Taylor formula for tangent vector fields.

**Lemma 7.4.7 (Taylor's theorem)** *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  $y \in \mathcal{V}$ ,*

$$P_\gamma^{0 \leftarrow 1} \zeta_y = \zeta_x + \nabla_\xi \zeta + \int_0^1 (P_\gamma^{0 \leftarrow \tau} \nabla_{\gamma'(\tau)} \zeta - \nabla_\xi \zeta) \, d\tau, \quad (7.33)$$

where  $\gamma$  is the unique minimizing geodesic satisfying  $\gamma(0) = x$  and  $\gamma(1) = y$ , and  $\xi = \text{Exp}_x^{-1} y = \gamma'(0)$ .

*Proof.* Start from

$$P_\gamma^{0 \leftarrow 1} \zeta_y = \zeta_x + \int_0^1 \frac{d}{d\tau} P_\gamma^{0 \leftarrow \tau} \zeta \, d\tau = \zeta_x + \nabla_\xi \zeta + \int_0^1 \left( \frac{d}{d\tau} P_\gamma^{0 \leftarrow \tau} \zeta - \nabla_\xi \zeta \right) d\tau$$

and use the formula for the connection in terms of the parallel transport (see [dC92, Ch. 2, Ex. 2]), to obtain

$$\frac{d}{d\tau} P_\gamma^{0 \leftarrow \tau} \zeta = \frac{d}{d\epsilon} P_\gamma^{0 \leftarrow \tau} P_\gamma^{\tau \leftarrow \tau + \epsilon} \zeta \Big|_{\epsilon=0} = P_\gamma^{0 \leftarrow \tau} \nabla_{\gamma'} \zeta.$$

□

We use this lemma to show that in some neighborhood of a nondegenerate local minimizer  $v$  of  $f$ , the norm of the gradient of  $f$  can be taken as a measure of the Riemannian distance to  $v$ .

**Lemma 7.4.8** *Let  $v \in \mathcal{M}$  and let  $f$  be a  $C^2$  cost function such that  $\text{grad } f(v) = 0$  and  $\text{Hess } f(v)$  is positive-definite with maximal and minimal eigenvalues  $\lambda_{\max}$  and  $\lambda_{\min}$ . Then, given  $c_0 < \lambda_{\min}$  and  $c_1 > \lambda_{\max}$ , there exists a neighborhood  $\mathcal{V}$  of  $v$  such that, for all  $x \in \mathcal{V}$ , it holds that*

$$c_0 \text{dist}(v, x) \leq \|\text{grad } f(x)\| \leq c_1 \text{dist}(v, x). \quad (7.34)$$

*Proof.* From Taylor (Lemma 7.4.7), it follows that

$$\begin{aligned} P_\gamma^{0 \leftarrow 1} \text{grad } f(v) &= \text{Hess } f(v)[\gamma'(0)] \\ &+ \int_0^1 (P_\gamma^{0 \leftarrow \tau} \text{Hess } f(\gamma(\tau))[\gamma'(\tau)] - \text{Hess } f(v)[\gamma'(0)]) \, d\tau. \end{aligned} \quad (7.35)$$

Since  $f$  is  $C^2$  and since  $\|\gamma'(\tau)\| = \text{dist}(v, x)$  for all  $\tau \in [0, 1]$ , we have the following bound for the integral in (7.35):

$$\begin{aligned} &\left\| \int_0^1 P_\gamma^{0 \leftarrow \tau} \text{Hess } f(\gamma(\tau))[\gamma'(\tau)] - \text{Hess } f(v)[\gamma'(0)] \, d\tau \right\| \\ &= \left\| \int_0^1 (P_\gamma^{0 \leftarrow \tau} \circ \text{Hess } f(\gamma(\tau)) \circ P_\gamma^{\tau \leftarrow 0} - \text{Hess } f(v)) [\gamma'(0)] \, d\tau \right\| \\ &\leq \epsilon(\text{dist}(v, x)) \text{dist}(v, x), \end{aligned}$$

where  $\lim_{t \rightarrow 0} \epsilon(t) = 0$ . Since  $\text{Hess } f(v)$  is nonsingular, it follows that  $|\lambda_{\min}| > 0$ . Take  $\mathcal{V}$  sufficiently small so that  $\lambda_{\min} - \epsilon(\text{dist}(v, x)) > c_0$  and  $\lambda_{\max} + \epsilon(\text{dist}(v, x)) < c_1$  for all  $x$  in  $\mathcal{V}$ . Then, using the fact that the parallel translation is an isometry, (7.34) follows from (7.35). □

We need a relation between the gradient of  $f$  at  $R_x(\xi)$  and the gradient of  $\hat{f}_x$  at  $\xi$ .

**Lemma 7.4.9** *Let  $R$  be a retraction on  $\mathcal{M}$  and let  $f$  be a  $C^1$  cost function on  $\mathcal{M}$ . Then, given  $v \in \mathcal{M}$  and  $c_5 > 1$ , there exists a neighborhood  $\mathcal{V}$  of  $v$  and  $\delta > 0$  such that*

$$\|\text{grad } f(R\xi)\| \leq c_5 \|\text{grad } \hat{f}(\xi)\|$$

for all  $x \in \mathcal{V}$  and all  $\xi \in T_x \mathcal{M}$  with  $\|\xi\| \leq \delta$ , where  $\hat{f}$  is as in (7.11).

*Proof.* Consider a parameterization of  $\mathcal{M}$  at  $v$  and consider the corresponding parameterization of  $T\mathcal{M}$  (see Section 3.5.3). We have

$$\partial_i \widehat{f}_x(\xi) = \sum_j \partial_j f(R\xi) A_i^j(\xi),$$

where  $A(\xi)$  stands for the differential of  $R_x$  at  $\xi \in T_x\mathcal{M}$ . Then,

$$\begin{aligned} \|\text{grad } \widehat{f}_x(\xi)\|^2 &= \sum_{i,j} \partial_i \widehat{f}_x(\xi) g^{ij}(x) \partial_j \widehat{f}_x(\xi) \\ &= \sum_{i,j,k,\ell} \partial_k f(R_x \xi) A_i^k(\xi) g^{ij}(x) A_j^\ell(\xi) \partial_\ell f(R_x \xi) \end{aligned}$$

and

$$\|\text{grad } f(R_x \xi)\|^2 = \sum_{i,j} \partial_i f(R_x \xi) g^{ij}(R_x \xi) \partial_j f(R_x \xi).$$

The conclusion follows by a real analysis argument, invoking the smoothness properties of  $R$  and  $g$ , and the compactness of the set  $\{(x, \xi) : x \in \mathcal{V}, \xi \in T_x\mathcal{M}, \|\xi\| \leq \delta\}$  and using  $A(0_x) = \text{id}$ .  $\square$

Finally, we will make use of Lemma 5.5.6 stating that the Hessians of  $f$  and  $\widehat{f}$  coincide at critical points.

We now state and prove the local convergence results. The first result states that the nondegenerate local minima are attractors of Algorithm 10-11.

**Theorem 7.4.10 (local convergence to local minima)** *Consider Algorithm 10-11—i.e., the Riemannian trust-region algorithm where the trust-region subproblems (7.6) are solved using the truncated CG algorithm with stopping criterion (7.10)—with all the assumptions of Theorem 7.4.2. Let  $v$  be a nondegenerate local minimizer of  $f$ , i.e.,  $\text{grad } f(v) = 0$  and  $\text{Hess } f(v)$  is positive-definite. Assume that  $x \mapsto \|H_x^{-1}\|$  is bounded on a neighborhood of  $v$  and that (7.25) holds for some  $\mu > 0$  and  $\delta_\mu > 0$ . Then there exists a neighborhood  $\mathcal{V}$  of  $v$  such that, for all  $x_0 \in \mathcal{V}$ , the sequence  $\{x_k\}$  generated by Algorithm 10-11 converges to  $v$ .*

*Proof.* Take  $\delta_1 > 0$  with  $\delta_1 < \delta_\mu$  such that  $\|H_x^{-1}\|$  is bounded on  $B_{\delta_1}(v)$ , that  $B_{\delta_1}(v)$  contains only  $v$  as critical point, and that  $f(x) > f(v)$  for all  $x \in \bar{B}_{\delta_1}(v)$ . (In view of the assumptions, such a  $\delta_1$  exists.) Take  $\delta_2$  small enough that, for all  $x \in B_{\delta_2}(v)$ , it holds that  $\|\eta^*(x)\| \leq \mu(\delta_1 - \delta_2)$ , where  $\eta^*$  is the (unique) solution of  $H_x \eta^* = -\text{grad } f(x)$ ; such a  $\delta_2$  exists because of Lemma 7.4.8 and the bound on  $\|H_x^{-1}\|$ . Consider a level set  $\mathcal{L}$  of  $f$  such that  $\mathcal{V} := \mathcal{L} \cap B_{\delta_1}(v)$  is a subset of  $B_{\delta_2}(v)$ ; invoke that  $f \in C^1$  to show that such a level set exists. Let  $\eta^{tCG}(x, \Delta)$  denote the tangent vector  $\eta_k$  returned by the truncated CG algorithm (Algorithm 11) when  $x_k = x$  and  $\Delta_k = \Delta$ . Then,  $\mathcal{V}$  is a neighborhood of  $v$ , and for all  $x \in \mathcal{V}$  and all  $\Delta > 0$ , we have

$$\text{dist}(x, x_+) \leq \frac{1}{\mu} \|\eta^{tCG}(x, \Delta)\| \leq \frac{1}{\mu} \|\eta^*\| \leq (\delta_1 - \delta_2),$$

where we used the fact that  $\|\eta\|$  is increasing along the truncated CG process (Proposition 7.3.2). It follows from the equation above that  $x_+$  is in  $B_{\delta_1}(v)$ . Moreover, since  $f(x_+) \leq f(x)$ , it follows that  $x_+ \in \mathcal{V}$ . Thus  $\mathcal{V}$  is invariant. But the only critical point of  $f$  in  $\mathcal{V}$  is  $v$ , so  $\{x_k\}$  goes to  $v$  whenever  $x_0$  is in  $\mathcal{V}$ .  $\square$

Now we study the order of convergence of sequences that converge to a nondegenerate local minimizer.

**Theorem 7.4.11 (order of convergence)** *Consider Algorithm 10-11 with stopping criterion (7.10). Suppose that  $f$  is a  $C^2$  cost function on  $\mathcal{M}$  and that*

$$\|H_k - \text{Hess } \widehat{f}_{x_k}(0_k)\| \leq \beta_{\mathcal{H}} \|\text{grad } f(x_k)\|, \quad (7.36)$$

*i.e.,  $H_k$  is a sufficiently good approximation of  $\text{Hess } \widehat{f}_{x_k}(0_{x_k})$ . Let  $v \in \mathcal{M}$  be a nondegenerate local minimizer of  $f$  (i.e.,  $\text{grad } f(v) = 0$  and  $\text{Hess } f(v)$  is positive-definite). Further assume that  $\text{Hess } \widehat{f}_x$  is Lipschitz-continuous at  $0_x$  uniformly in  $x$  in a neighborhood of  $v$ , i.e., there exist  $\beta_{L2} > 0$ ,  $\delta_1 > 0$ , and  $\delta_2 > 0$  such that, for all  $x \in B_{\delta_1}(v)$  and all  $\xi \in B_{\delta_2}(0_x)$ , it holds that*

$$\|\text{Hess } \widehat{f}_x(\xi) - \text{Hess } \widehat{f}_x(0_x)\| \leq \beta_{L2} \|\xi\|, \quad (7.37)$$

*where  $\|\cdot\|$  on the left-hand side denotes the operator norm in  $T_x\mathcal{M}$  defined as in (7.15). Then there exists  $c > 0$  such that, for all sequences  $\{x_k\}$  generated by the algorithm converging to  $v$ , there exists  $K > 0$  such that for all  $k > K$ ,*

$$\text{dist}(x_{k+1}, v) \leq c (\text{dist}(x_k, v))^{\min\{\theta+1, 2\}} \quad (7.38)$$

*with  $\theta > 0$  as in (7.10).*

*Proof.* The proof relies on a set of bounds which are justified after the main result is proved. Assume that there exist  $\tilde{\Delta}, c_0, c_1, c_2, c_3, c'_3, c_4, c_5$  such that, for all sequences  $\{x_k\}$  satisfying the conditions asserted, all  $x \in \mathcal{M}$ , all  $\xi$  with  $\|\xi\| < \tilde{\Delta}$ , and all  $k$  greater than some  $K$ , it holds that

$$c_0 \text{dist}(v, x_k) \leq \|\text{grad } f(x_k)\| \leq c_1 \text{dist}(v, x_k), \quad (7.39)$$

$$\|\eta_k\| \leq c_4 \|\text{grad } \widehat{m}_{x_k}(0)\| \leq \tilde{\Delta}, \quad (7.40)$$

$$\rho_k > \rho', \quad (7.41)$$

$$\|\text{grad } f(R_{x_k}\xi)\| \leq c_5 \|\text{grad } \widehat{f}_{x_k}(\xi)\|, \quad (7.42)$$

$$\|\text{grad } \widehat{m}_{x_k}(\xi) - \text{grad } \widehat{f}_{x_k}(\xi)\| \leq c_3 \|\xi\|^2 + c'_3 \|\text{grad } f(x_k)\| \|\xi\|, \quad (7.43)$$

$$\|\text{grad } \widehat{m}_{x_k}(\eta_k)\| \leq c_2 \|\text{grad } \widehat{m}_{x_k}(0)\|^{\theta+1}, \quad (7.44)$$

where  $\{\eta_k\}$  is the sequence of update vectors corresponding to  $\{x_k\}$ .

Given the bounds (7.39) to (7.44), the proof proceeds as follows. For all  $k > K$ , it follows from (7.39) and (7.41) that

$$c_0 \text{dist}(v, x_{k+1}) \leq \|\text{grad } f(x_{k+1})\| = \|\text{grad } f(R_{x_k}\eta_k)\|,$$

from (7.42) and (7.40) that

$$\|\text{grad } f(R_{x_k}\eta_k)\| \leq c_5 \|\text{grad } \widehat{f}_{x_k}(\eta_k)\|,$$

from (7.40) and (7.43) and (7.44) that

$$\begin{aligned} \|\text{grad } \widehat{f}_{x_k}(\eta_k)\| &\leq \|\text{grad } \widehat{m}_{x_k}(\eta_k) - \text{grad } \widehat{f}_{x_k}(\eta_k)\| + \|\text{grad } \widehat{m}_{x_k}(\eta_k)\| \\ &\leq (c_3c_4^2 + c_3'c_4)\|\text{grad } \widehat{m}_{x_k}(0)\|^2 + c_2\|\text{grad } \widehat{m}_{x_k}(0)\|^{1+\theta}, \end{aligned}$$

and from (7.39) that

$$\|\text{grad } \widehat{m}_{x_k}(0)\| = \|\text{grad } f(x_k)\| \leq c_1 \text{dist}(v, x_k).$$

Consequently, taking  $K$  larger if necessary so that  $\text{dist}(v, x_k) < 1$  for all  $k > K$ , it follows that

$$c_0 \text{dist}(v, x_{k+1}) \leq \|\text{grad } f(x_{k+1})\| \tag{7.45}$$

$$\leq c_5(c_3c_4^2 + c_3'c_4)\|\text{grad } f(x_k)\|^2 + c_5c_2\|\text{grad } f(x_k)\|^{\theta+1} \tag{7.46}$$

$$\leq c_5((c_3c_4^2 + c_3'c_4)c_1^2(\text{dist}(v, x_k))^2 + c_2c_1^{1+\theta}(\text{dist}(v, x_k))^{1+\theta})$$

$$\leq c_5((c_3c_4^2 + c_3'c_4)c_1^2 + c_2c_1^{1+\theta})(\text{dist}(v, x_k))^{\min\{2, 1+\theta\}}$$

for all  $k > K$ , which is the desired result.

It remains to show that the bounds (7.39)–(7.44) hold under the assumptions or the theorem.

Equation (7.39) comes from Lemma 7.4.8 and is due to the fact that  $v$  is a nondegenerate critical point.

We prove (7.40). Since  $\{x_k\}$  converges to the nondegenerate local minimizer  $v$  where  $\text{Hess } \widehat{f}_v(0_v) = \text{Hess } f(v)$  (in view of Lemma 5.5.6), and since  $\text{Hess } f(v)$  is positive-definite with  $f \in C^2$ , it follows from the approximation condition (7.36) and from (7.39) that there exists  $c_4 > 0$  such that, for all  $k$  greater than some  $K$ ,  $H_k$  is positive-definite and  $\|H_k^{-1}\| < c_4$ . Given a  $k > K$ , let  $\eta_k^*$  be the solution of  $H_k\eta_k^* = -\text{grad } \widehat{m}_{x_k}(0)$ . It follows that  $\|\eta_k^*\| \leq c_4\|\text{grad } \widehat{m}_{x_k}(0)\|$ . Since  $\{x_k\}$  converges to a critical point of  $f$  and since  $\text{grad } \widehat{m}_{x_k}(0) = \text{grad } f(x_k)$  in view of (7.6), we obtain that  $\|\eta_k^*\| \leq c_4\|\text{grad } \widehat{m}_{x_k}(0)\| \leq \tilde{\Delta}$  for any given  $\tilde{\Delta} > 0$  by choosing  $K$  larger if necessary. Then, since the sequence of  $\eta_k^j$ 's constructed by the truncated CG inner iteration (Algorithm 11) is strictly increasing in norm (Proposition 7.3.2) and would reach  $\eta_k^*$  at  $j = d$  in the absence of the stopping criterion, it follows that (7.40) holds.

We prove (7.41). Let  $\gamma_k$  denote  $\|\text{grad } f(x_k)\|$ . It follows from the definition of  $\rho_k$  that

$$\rho_k - 1 = \frac{\widehat{m}_{x_k}(\eta_k) - \widehat{f}_{x_k}(\eta_k)}{\widehat{m}_{x_k}(0_{x_k}) - \widehat{m}_{x_k}(\eta_k)}. \tag{7.47}$$

From Taylor's theorem, it holds that

$$\begin{aligned} \widehat{f}_{x_k}(\eta_k) &= \widehat{f}_{x_k}(0_{x_k}) + \langle \text{grad } f(x_k), \eta_k \rangle_{x_k} \\ &\quad + \int_0^1 \langle \text{Hess } \widehat{f}_{x_k}(\tau\eta_k)[\eta_k], \eta_k \rangle_{x_k} (1 - \tau) d\tau. \end{aligned}$$

It follows that

$$\begin{aligned}
& \left| \widehat{m}_{x_k}(\eta_k) - \widehat{f}_{x_k}(\eta_k) \right| \\
&= \left| \int_0^1 \left( \langle H_k[\eta_k], \eta_k \rangle_{x_k} - \langle \text{Hess } \widehat{f}_{x_k}(\tau\eta_k)[\eta_k], \eta_k \rangle_{x_k} \right) (1 - \tau) \, d\tau \right| \\
&\leq \int_0^1 \left| \langle (H_k - \text{Hess } \widehat{f}_{x_k}(0_{x_k}))[\eta_k], \eta_k \rangle_{x_k} \right| (1 - \tau) \, d\tau \\
&\quad + \int_0^1 \left| \langle (\text{Hess } \widehat{f}_{x_k}(0_{x_k}) - \text{Hess } \widehat{f}(\tau\eta_k))[\eta_k], \eta_k \rangle_{x_k} \right| (1 - \tau) \, d\tau \\
&\leq \frac{1}{2} \beta_{\mathcal{H}} \gamma_k \|\eta_k\|^2 + \frac{1}{6} \beta_{L2} \|\eta_k\|^3.
\end{aligned}$$

It then follows from (7.47), using the Cauchy bound (7.14), that

$$|\rho_k - 1| \leq \frac{(3\beta_{\mathcal{H}}\gamma_k + \beta_{L2}\|\eta_k\|)\|\eta_k\|^2}{6\gamma_k \min\{\Delta_k, \gamma_k/\beta\}},$$

where  $\beta$  is an upper bound on the norm of  $H_k$ . Since  $\|\eta_k\| \leq \Delta_k$  and  $\|\eta_k\| \leq c_4\gamma_k$ , it follows that

$$|\rho_k - 1| \leq \frac{(3\beta_{\mathcal{H}} + \beta_{L2}c_4)(\min\{\Delta_k, c_4\gamma_k\})^2}{6 \min\{\Delta_k, \gamma_k/\beta\}}. \quad (7.48)$$

Either  $\Delta_k$  is active in the denominator of (7.48), in which case we have

$$|\rho_k - 1| \leq \frac{(3\beta_{\mathcal{H}} + \beta_{L2}c_4)\Delta_k c_4 \gamma_k}{6\Delta_k} = \frac{(3\beta_{\mathcal{H}} + \beta_{L2}c_4)c_4}{6} \gamma_k,$$

or  $\gamma_k/\beta$  is active in the denominator of (7.48), in which case we have

$$|\rho_k - 1| \leq \frac{(3\beta_{\mathcal{H}} + \beta_{L2}c_4)(c_4\gamma_k)^2}{6\gamma_k/\beta} = \frac{(3\beta_{\mathcal{H}} + \beta_{L2}c_4)c_4^2\beta}{6} \gamma_k.$$

In both cases,  $\lim_{k \rightarrow \infty} \rho_k = 1$  since, in view of (7.39),  $\lim_{k \rightarrow \infty} \gamma_k = 0$ .

Equation (7.42) comes from Lemma 7.4.9.

We prove (7.43). It follows from Taylor's formula (Lemma 7.4.7, where the parallel translation becomes the identity since the domain of  $\widehat{f}_{x_k}$  is the Euclidean space  $T_{x_k}\mathcal{M}$ ) that

$$\begin{aligned}
\text{grad } \widehat{f}_{x_k}(\xi) &= \text{grad } \widehat{f}_{x_k}(0_{x_k}) + \text{Hess } \widehat{f}_{x_k}(0_{x_k})[\xi] \\
&\quad + \int_0^1 \left( \text{Hess } \widehat{f}_{x_k}(\tau\xi) - \text{Hess } \widehat{f}_{x_k}(0_{x_k}) \right) [\xi] \, d\tau.
\end{aligned}$$

The conclusion comes by the Lipschitz condition (7.37) and the approximation condition (7.36).

Finally, equation (7.44) comes from the stopping criterion (7.10) of the inner iteration. More precisely, the truncated CG loop (Algorithm 11) terminates if  $\langle \delta_j, H_k \delta_j \rangle \leq 0$  or  $\|\eta_{j+1}\| \geq \Delta$  or the criterion (7.10) is satisfied. Since  $\{x_k\}$  converges to  $v$  and  $\text{Hess } f(v)$  is positive-definite, it follows that  $H_k$  is positive-definite for all  $k$  greater than a certain  $K$ . Therefore, for all

$k > K$ , the criterion  $\langle \delta_j, H_k \delta_j \rangle \leq 0$  is never satisfied. In view of (7.40) and (7.41), it can be shown that the trust region is eventually inactive. Therefore, increasing  $K$  if necessary, the criterion  $\|\eta_{j+1}\| \geq \Delta$  is never satisfied for all  $k > K$ . In conclusion, for all  $k > K$ , the stopping criterion (7.10) is satisfied each time a computed  $\eta_k$  is returned by the truncated CG loop. Therefore, the truncated CG loop behaves as a classical linear CG method. Consequently,  $\text{grad } \widehat{m}_{x_k}(\eta_j) = r_j$  for all  $j$ . Choose  $K$  such that for all  $k > K$ ,  $\|\text{grad } f(x_k)\| = \|\text{grad } \widehat{m}_{x_k}(0)\|$  is so small—it converges to zero in view of (7.39)—that the stopping criterion (7.10) yields

$$\|\text{grad } \widehat{m}_{x_k}(\eta_j)\| = \|r_j\| \leq \|r_0\|^{1+\theta} = \|\text{grad } \widehat{m}_{x_k}(0)\|^{1+\theta}. \quad (7.49)$$

This is (7.44) with  $c_2 = 1$ .  $\square$

The constants in the proof of Theorem 7.4.11 can be chosen as  $c_0 < \lambda_{\min}$ ,  $c_1 > \lambda_{\max}$ ,  $c_4 > 1/\lambda_{\min}$ ,  $c_5 > 1$ ,  $c_3 \geq \beta_{L2}$ ,  $c'_3 \geq \beta_{\mathcal{H}}$ ,  $c_2 \geq 1$ , where  $\lambda_{\min}$  and  $\lambda_{\max}$  are the smallest and largest eigenvalue of  $\text{Hess } f(v)$ , respectively. Consequently, the constant  $c$  in the convergence bound (7.38) can be chosen as

$$c > \frac{1}{\lambda_{\min}} \left( (\beta_{L2}/\lambda_{\min}^2 + \beta_{\mathcal{H}}/\lambda_{\min}) \lambda_{\max}^2 + \lambda_{\max}^{1+\theta} \right). \quad (7.50)$$

A nicer-looking bound holds when convergence is evaluated in terms of the norm of the gradient, as expressed in the theorem below which is a direct consequence of (7.45) and (7.46).

**Theorem 7.4.12** *Under the assumptions of Theorem 7.4.11, if  $\theta + 1 < 2$ , then given  $c_g > 1$  and  $\{x_k\}$  generated by the algorithm, there exists  $K > 0$  such that*

$$\|\text{grad } f(x_{k+1})\| \leq c_g \|\text{grad } f(x_k)\|^{\theta+1}$$

for all  $k > K$ .

Nevertheless, (7.50) suggests that the algorithm may not perform well when the relative gap  $\lambda_{\max}/\lambda_{\min}$  is large. In spite of this, numerical experiments on eigenvalue problems have shown that the method tends to behave as well as, or even better than, other methods in the presence of a small relative gap.

### 7.4.3 Discussion

The main global convergence result (Theorem 7.4.4) shows that RTR-tCG method (Algorithm 10-11) converges to a set of critical points of the cost function for *all* initial conditions. This is an improvement on the pure Newton method (Algorithm 5), for which only local convergence results exist. However, the convergence theory falls short of showing that the algorithm always converges to a local minimizer. This is not surprising: since we have ruled out the possibility of checking the positive-definiteness of the Hessian of the cost function, we have no way of testing whether a critical point is

a local minimizer or not. (Note as an aside that even checking the positive-definiteness of the Hessian is not always sufficient for determining if a critical point is a local minimizer or not: if the Hessian is singular and nonnegative-definite, then no conclusion can be drawn.) In fact, for the vast majority of optimization methods, only convergence to critical points can be secured unless some specific assumptions (like convexity) are made. Nevertheless, it is observed in numerical experiments with random initial conditions that the algorithm systematically converges to a local minimizer; convergence to a saddle point is observed only on specifically crafted problems, e.g., when the iteration is started on a point that is a saddle point in computer arithmetic. This is due to the fact that the algorithm is a descent method, i.e.,  $f(x_{k+1}) < f(x_k)$  whenever  $x_{k+1} \neq x_k$ . Therefore, saddle points or local maxima are unstable fixed points of the algorithm.

There are cases where the bound (7.38) holds with order  $\min\{\theta + 1, 3\}$ ; i.e., by choosing  $\theta \geq 2$ , one obtains cubic convergence. This situation occurs when the Taylor expansion of the cost function around the limit point has no third-order contribution. Thus, the second-order approximation used in the algorithm becomes an effective third-order approximation, and the order of convergence benefits as expected. In practice, this condition holds more often than one might guess since any cost function that is symmetric around the local minimizer  $v$ , i.e.,  $f(\text{Exp}_x(\xi)) = f(\text{Exp}_x(-\xi))$ , will have only even contributions to its Taylor expansion.

## 7.5 APPLICATIONS

In this section, we briefly review the essential “ingredients” necessary for applying the RTR-tCG method (Algorithm 10-11), and we present two examples in detail as an illustration. One of these examples is optimization of the Rayleigh quotient, leading to an algorithm for computing extreme invariant subspaces of symmetric matrices. Since trust-region algorithms with an exact Hessian can be thought of as enhanced Newton methods, and since the Newton equation for Rayleigh quotient optimization is equivalent to the Jacobi equation (see Notes and References in Chapter 6), it is not surprising that this algorithm has close links with the celebrated Jacobi-Davidson approach to the eigenproblem. The Riemannian trust-region approach sheds new light on the Jacobi-Davidson method. In particular, it yields new ways to deal with the Jacobi equation so as to reduce the computational burden while preserving the superlinear convergence inherited from the Newton approach and obtaining strong global convergence results supported by a detailed analysis.

### 7.5.1 Checklist

The following elements are required for applying the RTR method to optimizing a cost function  $f$  on a Riemannian manifold  $(\mathcal{M}, g)$ : (i) a tractable



numerical representation for points  $x$  on  $\mathcal{M}$ , for tangent spaces  $T_x\mathcal{M}$ , and for the inner products  $\langle \cdot, \cdot \rangle_x$  on  $T_x\mathcal{M}$ ; (ii) a retraction  $R_x : T_x\mathcal{M} \rightarrow \mathcal{M}$  (Definition 4.1.1); (iii) formulas for  $f(x)$ ,  $\text{grad } f(x)$  and an approximate Hessian  $H_x$  that satisfies the properties required for the convergence results in Section 7.4.

Formulas for  $\text{grad } f(x)$  and  $\text{Hess } \hat{f}_x(0_x)$  can be obtained by identification in a Taylor expansion of the lifted cost function  $\hat{f}_x$ , namely

$$\hat{f}_x(\eta) = f(x) + \langle \text{grad } f(x), \eta \rangle_x + \frac{1}{2} \langle \text{Hess } \hat{f}_x(0_x)[\eta], \eta \rangle_x + O(\|\eta\|^3),$$

where  $\text{grad } f(x) \in T_x\mathcal{M}$  and  $\text{Hess } \hat{f}_x(0_x)$  is a linear transformation of  $T_x\mathcal{M}$ . A formula for  $\text{Hess } \hat{f}_x(0_x)$  is not needed, though; the convergence theory requires only an “approximate Hessian”  $H_x$  that satisfies the approximation condition (7.36). To obtain such an  $H_x$ , one can pick  $H_x := \text{Hess}(f \circ \tilde{R}_x)(0_x)$ , where  $\tilde{R}$  is any retraction. Then, assuming sufficient smoothness of  $f$ , the bound (7.36) follows from Lemmas 7.4.8 and 5.5.6. In particular, the choice  $\tilde{R}_x = \text{Exp}_x$  yields

$$H_x := \nabla \text{grad } f(x) \quad (= \text{Hess } f(x)), \quad (7.51)$$

where  $\nabla$  denotes the Riemannian connection, and the model  $\hat{m}_x$  takes the form (7.1). If  $\mathcal{M}$  is a Riemannian submanifold or a Riemannian quotient of a Euclidean space, then  $\nabla \text{grad } f(x)$  admits a simple formula; see Section 5.3.

## 7.5.2 Symmetric eigenvalue decomposition

Let  $\mathcal{M}$  be the orthogonal group

$$\mathcal{M} = O_n = \{Q \in \mathbb{R}^{n \times n} : Q^T Q = I_n\}.$$

This manifold is an embedded submanifold of  $\mathbb{R}^{n \times n}$  (Section 3.3). The tangent spaces are given by  $T_Q O_n = \{Q\Omega : \Omega = -\Omega^T\}$  (Section 3.5.7). The canonical Euclidean metric  $g(A, B) = \text{tr}(A^T B)$  on  $\mathbb{R}^{n \times n}$  induces on  $O_n$  the metric

$$\langle Q\Omega_1, Q\Omega_2 \rangle_Q = \text{tr}(\Omega_1^T \Omega_2). \quad (7.52)$$

A retraction  $R_Q : T_Q O_n \rightarrow O_n$  must be chosen that satisfies the properties stated in Section 7.2. Several possibilities are mentioned in Section 4.1.1.

Consider the cost function

$$f(Q) = \text{tr}(Q^T A Q N),$$

where  $A$  and  $N$  are  $n \times n$  symmetric matrices. For  $N = \text{diag}(\mu_1, \dots, \mu_n)$ ,  $\mu_1 < \dots < \mu_n$ , the minimum of  $f$  is realized by the orthonormal matrices of eigenvectors of  $A$  sorted in decreasing order of corresponding eigenvalue (this is a consequence of the critical points analysis in Section 4.8). Assume that the retraction  $R$  approximates the exponential at least to order 2. With the metric  $g$  defined as in (7.52), we obtain

$$\begin{aligned} \hat{f}_Q(Q\Omega) &:= f(R_Q(Q\Omega)) \\ &= \text{tr}((I + \Omega + \frac{1}{2}\Omega^2 + O(\Omega^3))^T Q^T A Q (I + \Omega + \frac{1}{2}\Omega^2 + O(\Omega^3)) N) \\ &= f(Q) + 2\text{tr}(\Omega^T Q^T A Q N) \\ &\quad + \text{tr}(\Omega^T Q^T A Q \Omega N - \Omega^T \Omega Q^T A Q N) + O(\Omega^3), \end{aligned}$$

from which it follows that

$$\begin{aligned} D\widehat{f}_Q(0)[Q\Omega] &= 2 \operatorname{tr}(Q^T A Q \Omega N) \\ \frac{1}{2} D^2 \widehat{f}_Q(0)[Q\Omega_1, Q\Omega_2] &= \operatorname{tr}(\Omega_1^T Q^T A Q \Omega_2 N - \frac{1}{2}(\Omega_1^T \Omega_2 + \Omega_2^T \Omega_1) Q^T A Q N) \\ \operatorname{grad} \widehat{f}_Q(0) &= \operatorname{grad} f(Q) = Q[Q^T A Q, N] \end{aligned}$$

$\operatorname{Hess} \widehat{f}_Q(0)[Q\Omega] = \operatorname{Hess} f(Q)[Q\Omega] = \frac{1}{2}Q[[Q^T A Q, \Omega], N] + \frac{1}{2}Q[[N, \Omega], Q^T A Q]$ , where  $[A, B] := AB - BA$ . It is now straightforward to replace these expressions in the general formulation of Algorithm 10-11 and obtain a practical matrix algorithm.

An alternative way to obtain  $\operatorname{Hess} \widehat{f}_Q(0)$  is to exploit Proposition 5.5.5, which yields  $\operatorname{Hess} \widehat{f}_Q(0) = \nabla \operatorname{grad} f(Q)$ . Since the manifold  $\mathcal{M}$  is an embedded Riemannian submanifold of  $\mathbb{R}^{n \times p}$ , the covariant derivative  $\nabla$  is obtained by projecting the derivative in  $\mathbb{R}^{n \times p}$  onto the tangent space to  $\mathcal{M}$ ; see Section 5.3.3. We obtain  $\operatorname{Hess} f(Q)[Q\Omega] = Q \operatorname{skew}(\Omega[Q^T A Q, N] + [\Omega^T Q^T A Q + Q^T A Q \Omega, N])$ , which yields the same result as above.

All these ingredients can now be used in Algorithm 10-11 to obtain an iteration that satisfies the convergence properties proven in Section 7.4. Convergence to the critical points of the cost function means convergence to the matrices whose column vectors are the eigenvectors of  $A$ . Only the matrices containing eigenvectors in decreasing order of eigenvalue can be stable fixed points for the algorithm. They are asymptotically stable, with superlinear convergence, when all the eigenvalues are simple.

### 7.5.3 Computing an extreme eigenspace

Following up on the geometric Newton method for the generalized eigenvalue problem obtained in Section 6.4.3, we assume again that  $A$  and  $B$  are  $n \times n$  symmetric matrices with  $B$  positive-definite, and we consider the generalized eigenvalue problem

$$Av = \lambda Bv.$$

We want to compute the leftmost  $p$ -dimensional invariant subspace of the pencil  $(A, B)$ .

We consider the Rayleigh quotient function  $f$  defined by

$$f(\operatorname{span}(Y)) = \operatorname{tr}((Y^T A Y)(Y^T B Y)^{-1}), \tag{7.53}$$

where  $Y$  belongs to the set of full-rank  $n \times p$  matrices and  $\operatorname{span}(Y)$  denotes the column space of  $Y$ . The critical points of  $f$  are the invariant subspaces of the pencil  $(A, B)$ , and the minimizers of  $f$  correspond to the leftmost invariant subspaces of  $(A, B)$ .

In Section 6.4.3, we chose a noncanonical Riemannian metric (6.29) that yields a relatively short formula (6.34) for the Riemannian Hessian of  $f$ , obtained using the theory of Riemannian submersions (Proposition 5.3.4). In this section, we again use the definition

$$\mathcal{H}_Y = \{Z \in \mathbb{R}^{n \times p} : Y^T B Z = 0\}$$

for the horizontal spaces, and we take

$$R_{\mathcal{Y}}(\xi) = \text{span}(Y + \bar{\xi}_Y), \quad (7.54)$$

where  $\mathcal{Y} = \text{span}(Y)$  and  $\bar{\xi}_Y$  stands for the horizontal lift of the tangent vector  $\xi \in T_{\mathcal{Y}} \text{Grass}(p, n)$ . But now we define the Riemannian metric as

$$\langle \xi, \zeta \rangle_{\mathcal{Y}} = \text{tr} \left( (Y^T B Y)^{-1} \bar{\xi}_Y^T \bar{\zeta}_Y \right). \quad (7.55)$$

Notice that the horizontal space is not orthogonal to the vertical space with respect to the new Riemannian metric (7.55), so we have to renounce using the Riemannian submersion theory. However, we will see that with these choices for the horizontal space, the retraction, and the Riemannian metric, the second-order Taylor development of  $\widehat{f}_x := f \circ R_x$  admits quite a simple form.

For the Rayleigh cost function (7.53), using the notation

$$P_{U,V} = I - U(V^T U)^{-1} V^T \quad (7.56)$$

for the projector parallel to the span of  $U$  onto the orthogonal complement of the span of  $V$ , we obtain

$$\begin{aligned} & \widehat{f}_{\mathcal{Y}}(\xi) \\ &= f(R_{\mathcal{Y}}(\xi)) = \text{tr} \left( ((Y + \bar{\xi}_Y)^T B (Y + \bar{\xi}_Y))^{-1} ((Y + \bar{\xi}_Y)^T A (Y + \bar{\xi}_Y)) \right) \\ &= \text{tr} \left( (Y^T B Y)^{-1} Y^T A Y \right) + 2 \text{tr} \left( (Y^T B Y)^{-1} \bar{\xi}_Y^T A Y \right) \\ &\quad + \text{tr} \left( (Y^T B Y)^{-1} \bar{\xi}_Y^T (A \bar{\xi}_Y - B \bar{\xi}_Y (Y^T B Y)^{-1} (Y^T A Y)) \right) + O(\|\xi\|^3) \\ &= \text{tr} \left( (Y^T B Y)^{-1} Y^T A Y \right) + 2 \text{tr} \left( (Y^T B Y)^{-1} \bar{\xi}_Y^T P_{BY, BY} A Y \right) \\ &\quad + \text{tr} \left( (Y^T B Y)^{-1} \bar{\xi}_Y^T P_{BY, BY} (A \bar{\xi}_Y - B \bar{\xi}_Y (Y^T B Y)^{-1} (Y^T A Y)) \right) \\ &\quad + O(\|\xi\|^3), \end{aligned} \quad (7.57)$$

where the introduction of the projectors does not modify the expression since  $P_{BY, BY} \bar{\xi}_Y = \bar{\xi}_Y$ . By identification, using the noncanonical metric (7.55), we obtain

$$\overline{\text{grad } f(\mathcal{Y})}_Y = \overline{\text{grad } \widehat{f}_{\mathcal{Y}}(0)}_Y = 2P_{BY, BY} A Y \quad (7.58)$$

and

$$\overline{\text{Hess } \widehat{f}_{\mathcal{Y}}(0_{\mathcal{Y}})}[\xi]_Y = 2P_{BY, BY} (A \bar{\xi}_Y - B \bar{\xi}_Y (Y^T B Y)^{-1} (Y^T A Y)). \quad (7.59)$$

Notice that  $\text{Hess } \widehat{f}_{\mathcal{Y}}(0_{\mathcal{Y}})$  is symmetric with respect to the metric, as required.

We choose to take

$$H_{\mathcal{Y}} := \text{Hess } \widehat{f}_{\mathcal{Y}}(0_{\mathcal{Y}}). \quad (7.60)$$

Consequently, the approximation condition (7.36) is trivially satisfied. The model (7.6) is thus

$$\begin{aligned} \widehat{m}_{\mathcal{Y}}(\xi) &= f(\mathcal{Y}) + \langle \text{grad } f(\mathcal{Y}), \xi \rangle_{\mathcal{Y}} + \frac{1}{2} \langle H_{\mathcal{Y}} \xi, \xi \rangle_{\mathcal{Y}} \\ &= \text{tr} \left( (Y^T B Y)^{-1} Y^T A Y \right) + 2 \text{tr} \left( (Y^T B Y)^{-1} \bar{\xi}_Y^T A Y \right) \\ &\quad + \text{tr} \left( (Y^T B Y)^{-1} \bar{\xi}_Y^T (A \bar{\xi}_Y - B \bar{\xi}_Y (Y^T B Y)^{-1} Y^T A Y) \right). \end{aligned} \quad (7.61)$$

(Observe that for  $B = I$  and  $p = 1$ , and choosing  $y = Y$  of unit norm, the model (7.61) becomes  $\widehat{m}_Y(\xi) = \text{tr}(y^T Ay) + 2\bar{\xi}_y^T Ay + \bar{\xi}_y^T (A - y^T AyI)\bar{\xi}_y$ ,  $y^T \bar{\xi}_y = 0$ . Assuming that the Hessian is nonsingular, this model has a unique critical point that is the solution of the Newton equation (6.17) for the Rayleigh quotient on the sphere.)

Since the Rayleigh cost function (7.53) is smooth on  $\text{Grass}(p, n)$  and since  $\text{Grass}(p, n)$  is compact, it follows that all the assumptions involved in the convergence analysis of the general RTR-tCG algorithm (Section 7.4) are satisfied. The only complication is that we do not have a closed-form expression for the distance involved in the superlinear convergence result (7.38). But since  $B$  is fixed and positive-definite, the distances induced by the non-canonical metric (7.55) and by the canonical metric—(7.55) with  $B := I$ —are locally equivalent, and therefore for a given sequence both distances yield the same rate of convergence. (Saying that two distances  $\text{dist}_1$  and  $\text{dist}_2$  are *locally equivalent* means that given  $x \in \mathcal{M}$ , there is a neighborhood  $\mathcal{U}$  of  $x$  and constants  $c_1, c_2 > 0$  such that, for all  $y$  in  $\mathcal{U}$ , we have  $c_1 \text{dist}_1(x, y) \leq \text{dist}_2(x, y) \leq c_2 \text{dist}_1(x, y)$ .)

We have now all the required information to use the RTR-tCG method (Algorithm 10-11) for minimizing the Rayleigh cost function (7.53) on the Grassmann manifold  $\text{Grass}(p, n)$  endowed with the noncanonical metric (7.55). A matrix version of the inner iteration is displayed in Algorithm 12, where we omit the horizontal lift notation for conciseness and define

$$H_Y[Z] := P_{BY, BY}(AZ - BZ(Y^T BY)^{-1}Y^T AY). \quad (7.62)$$

Note that omission of the factor 2 in both the gradient and the Hessian does not affect the sequence  $\{\eta\}$  generated by the truncated CG algorithm.

According to the retraction formula (7.54), the returned  $\eta_k$  yields a candidate new iterate

$$Y_{k+1} = (Y_k + \eta_k)M_k,$$

where  $M_k$  is chosen such that  $Y_{k+1}^T BY_{k+1} = I$ . The candidate is accepted or rejected, and the trust-region radius is updated as prescribed in the outer RTR method (Algorithm 10), where  $\rho$  is computed using  $\widehat{m}_k$  as in (7.61) and  $\widehat{f}$  as in (7.57).

The resulting algorithm converges to the set of invariant subspaces of  $(A, B)$ —which are the critical points of the cost function (7.53)—and convergence to the leftmost invariant subspace  $\mathcal{V}$  is expected to occur in practice since the other invariant subspaces are numerically unstable. Moreover, since  $\mathcal{V}$  is a nondegenerate local minimum (under our assumption that  $\lambda_p < \lambda_{p+1}$ ), it follows that the rate of convergence is  $\min\{\theta + 1, 2\}$ , where  $\theta$  is the parameter appearing in the stopping criterion (7.10) of the inner (truncated CG) iteration.

Numerical experiments illustrating the convergence of the Riemannian trust-region algorithm for extreme invariant subspace computation are presented in Figures 7.1 and 7.2. The  $\theta$  parameter of the inner stopping criterion (7.10) was chosen equal to 1 to obtain quadratic convergence (see Theorem 7.4.11). The right-hand plot shows a case with a small eigenvalue gap,

---

**Algorithm 12** Truncated CG method for the generalized eigenvalue problem

---

**Require:** Symmetric  $n \times n$  matrices  $A$  and  $B$  with  $B$  positive-definite; a  $B$ -orthonormal full-rank  $n \times p$  matrix  $Y$  (i.e.,  $Y^T B Y = I$ ); operator  $H_Y$  defined in (7.62);  $\Delta > 0$ .

```

1: Set  $\eta^0 = 0 \in \mathbb{R}^{n \times p}$ ,  $r_0 = P_{B Y, B Y} A Y$ ,  $\delta_0 = -r_0$ ;  $j = 0$ ;
2: loop
3:   if  $\text{tr}(\delta_j^T H_Y[\delta_j]) \leq 0$  then
4:     Compute  $\tau > 0$  such that  $\eta = \eta^j + \tau \delta_j$  satisfies  $\text{tr}(\eta^T \eta) = \Delta$ ;
5:     return  $\eta_k := \eta$ ;
6:   end if
7:   Set  $\alpha_j = \text{tr}(r_j^T r_j) / \text{tr}(\delta_j^T H_Y[\delta_j])$ ;
8:   Set  $\eta^{j+1} = \eta^j + \alpha_j \delta_j$ ;
9:   if  $\text{tr}((\eta^{j+1})^T \eta^{j+1}) \geq \Delta$  then
10:    Compute  $\tau \geq 0$  such that  $\eta = \eta^j + \tau \delta_j$  satisfies  $\text{tr}(\eta^T \eta) = \Delta$ ;
11:    return  $\eta_k := \eta$ ;
12:   end if
13:   Set  $r_{j+1} = r_j + \alpha H_Y[\delta_j]$ ;
14:   Set  $\beta_{j+1} = \text{tr}(r_{j+1}^T r_{j+1}) / \text{tr}(r_j^T r_j)$ ;
15:   Set  $\delta_{j+1} = -r_{j+1} + \beta_{j+1} \delta_j$ ;
16:   if a stopping criterion is satisfied then
17:     return  $\eta_k := \eta^j$ ;
18:   end if
19: end loop

```

---

which implies that the smallest eigenvalue of the Hessian of the cost function at the solution is much smaller than its largest eigenvalue. This suggests that the multiplicative constant in the superlinear convergence bound (7.38) is large, which explains why superlinear convergence sets in less clearly than on the left-hand plot featuring a large eigenvalue gap. An experiment with a smaller machine epsilon would reveal a quadratic convergence pattern; i.e., the number of digits of accuracy eventually approximately doubles at each new iterate. (All the experiments described in this book were performed with a machine epsilon of approximately  $2 \cdot 10^{-16}$ .)

The eigenvalue algorithm resulting from the proposed approach is surprisingly competitive with other eigenvalue methods in spite of the fact that it is just a brute-force application of a general optimization scheme that does not make any attempt to exploit the specific form of the Rayleigh quotient cost function. The efficiency of the algorithm is supported by its similarity to the Jacobi-Davidson eigenvalue method, in particular to the JDCG method of Notay. Nevertheless, the algorithm admits several enhancements, including subspace acceleration techniques and the possible monitoring of the  $\rho$  ratio within the inner iteration at a low computational cost. These enhancements, as well as comparisons with state-of-the-art eigenvalue methods, are

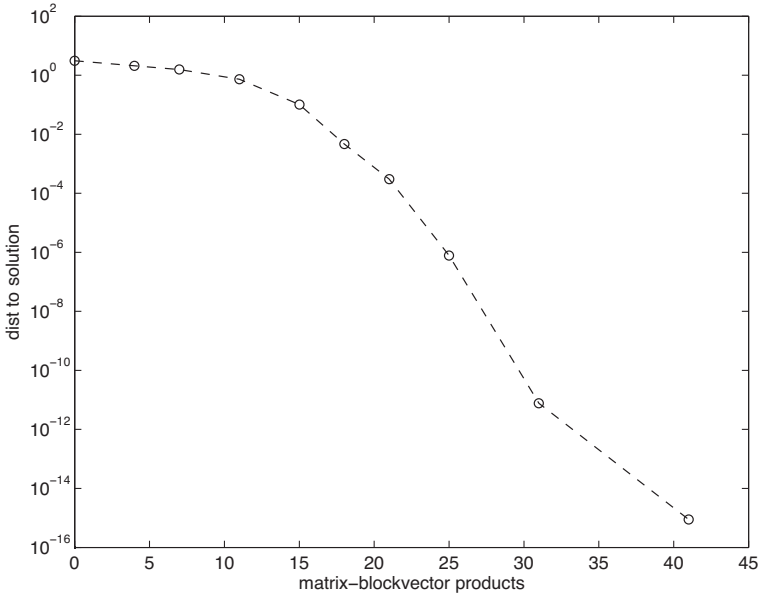


Figure 7.1 Numerical experiments on a trust-region algorithm for minimizing the Rayleigh cost function (7.53) on the Grassmann manifold  $\text{Grass}(p, n)$ , with  $n = 100$  and  $p = 5$ .  $B = I$ , and  $A$  is chosen with  $p$  eigenvalues evenly spaced on the interval  $[1, 2]$  and the other  $(n - p)$  eigenvalues evenly spaced on the interval  $[10, 11]$ ; this is a problem with a large eigenvalue gap. The horizontal axis gives the number of multiplications of  $A$  times a block of  $p$  vectors. The vertical axis gives the distance to the solution, defined as the square root of the sum of the canonical angles between the current subspace and the leftmost  $p$ -dimensional invariant subspace of  $A$ . (This distance corresponds to the geodesic distance on the Grassmann manifold endowed with its canonical metric (3.44).)

presented in articles mentioned in Notes and References.

## 7.6 NOTES AND REFERENCES

The Riemannian trust-region approach was first proposed in [ABG04]. Most of the material in this section comes from [ABG07].

For more information on trust-region methods in  $\mathbb{R}^n$ , we refer the reader to Conn *et al.* [CGT00]. Trust-region methods are also discussed in textbooks on numerical optimization such as Nocedal and Wright [NW99]; see also Hei [Hei03], Gould *et al.* [GOST05], and Walmag and Delhez [WD05] for recent developments. Algorithm 10 reduces to [NW99, Alg. 4.1] in the classical  $\mathbb{R}^n$  case; variants can be found in Conn *et al.* [CGT00, Ch. 10].

The method for computing an accurate solution of the trust-region sub-

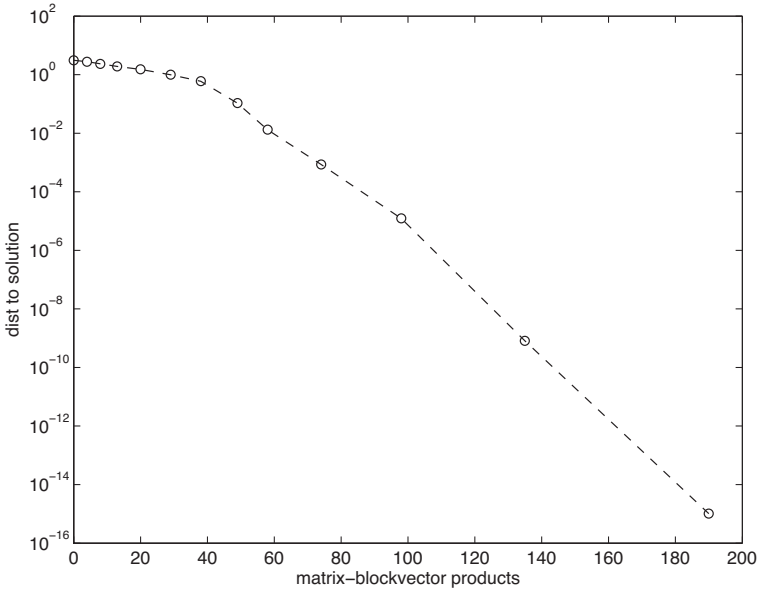


Figure 7.2 Same situation as in Figure 7.1 but now with  $B = I$  and  $A = \text{diag}(1, \dots, n)$ .

problem is due to Moré and Sorensen [MS83]. Proposition 7.3.1 is a straightforward transcription of [CGT00, Th. 7.4.1], which itself generalizes results from [MS83] (or see [NW99, Th. 4.3]) to general norms.

The truncated CG method presented in Algorithm 11 closely follows the algorithm proposed by Steihaug [Ste83]; see also the work of Toint [Toi81]. Proposition 7.3.2 is due to Steihaug [Ste83, Th. 2.1]. The reader interested in the underlying principles of the Steihaug-Toint truncated CG method should refer to [Ste83], [NW99], or [CGT00].

Besides the truncated CG method, available algorithms for (approximately) solving trust-region subproblems include the dogleg method of Powell [Pow70], the double-dogleg method of Dennis and Mei [DM79], the method of Moré and Sorensen [MS83], the two-dimensional subspace minimization strategy of Byrd *et al.* [BSS88], the method based on the difference of convex functions proposed by Pham Dinh Tao and Le Thi Hoai An [TA98], the truncated Lanczos approach of Gould *et al.* [GLRT99], the matrix-free eigenproblem-based algorithm of Rojas *et al.* [RSS00], and the sequential subspace method of Hager [Hag01, HP05]. These and other methods are discussed in Conn *et al.* [CGT00, §7.5.4].

The classical global convergence results for trust-region methods in  $\mathbb{R}^n$  can be found in Nocedal and Wright [NW99] (see in particular Theorem 4.8) and Conn *et al.* [CGT00]. A Taylor development similar to Lemma 7.4.7 can be found in Smith [Smi94]. The principle of the argument of Theorem 7.4.10 is closely related to the capture theorem, see Bertsekas [Ber95, Th 1.2.5]. A

discussion on cubic versus quadratic convergence can be found in Dehaene and Vandewalle [DV00]. A proof of Corollary 7.4.6 is given in [ABG06a].

Experiments, comparisons, and further developments are presented in [ABG06b, ABGS05, BAG06] for the Riemannian trust-region approach to extreme invariant subspace computation (Section 7.5.3). Reference [ABG06b] works out a brute-force application of the Riemannian trust-region method to the optimization of the Rayleigh quotient cost function on the sphere: comparisons are made with other eigenvalue methods, in particular the JDCG algorithm of Notay [Not02] and the Tracemin algorithm of Sameh, Wisniewski, and Tong [SW82, ST00], and numerical results are presented. Reference [ABGS05] proposes a two-phase method that combines the advantages of the (unshifted) Tracemin method and of the RTR-tCG method with an order-2 model; it allows one to make efficient use of a preconditioner in the first iterations by relaxing the trust-region constraint. The implicit RTR method proposed in Baker *et al.* [BAG06] makes use of the particular structure of the eigenvalue problem to monitor the value of the ratio  $\rho$  in the course of the inner iteration with little computational overhead, thereby avoiding the rejection of iterates because of poor model quality; for some problems, this technique considerably speeds up the iteration while the iterates are still far away from the solution, especially when a good preconditioner is available.