

---

# Optimization on Smooth Manifolds Crash Course

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Running examples . . . . .	2
1.2	Notation . . . . .	3
1.3	Terminology . . . . .	3
<b>2</b>	<b>First-order geometry</b>	<b>3</b>
2.1	First-order geometry: material not covered . . . . .	11
2.2	First-order geometry: extending beyond the embedded case . . . . .	12
<b>3</b>	<b>Second-order geometry</b>	<b>12</b>
3.1	Second-order geometry: material not covered . . . . .	16
3.2	Second-order geometry: extending beyond the embedded case . . . . .	16
<b>4</b>	<b>Retractions, Taylor expansions, and optimality conditions</b>	<b>17</b>
<b>5</b>	<b>Quotient manifolds</b>	<b>23</b>
5.1	Quotient manifolds: material not covered . . . . .	30
<b>6</b>	<b>Riemannian gradient descent</b>	<b>30</b>
<b>7</b>	<b>Riemannian second-order methods</b>	<b>30</b>
<b>A</b>	<b>Examples of manifolds</b>	<b>30</b>
<b>B</b>	<b>Regarding product manifolds</b>	<b>30</b>

# 1 Introduction

This is a crash course on optimization on smooth manifolds. Basically all of the definitions and theorems are taken from [Bou20], and I've done my best to provide the exact places in the book where you can find them. Unfortunately [Bou20] is a preprint, and newer versions have been released since I made these notes. The version I used to create these notes was compiled on November 9, 2020.<sup>1</sup> Thus, be warned that theorem/definition numbering in the most recent version of [Bou20] may not always exactly match what is given here (but sections should stay the same).

An important note about [Bou20]: They structure their textbook by first defining and proving everything for *embedded submanifolds* (manifolds that are subsets of a Euclidean space). It isn't until Chapter 8 that they go back and define things more generally. As embedded submanifolds are much simpler to work with and make up most applications, this write-up focuses on them nearly exclusively. When it makes sense to do so though I will try to provide intuition about how definitions generalize beyond the embedded case, since all of the generalizations are very natural and not something to be afraid of.

## 1.1 Running examples

We will use the following optimization problem as a running example:

$$\min_{y \in \mathcal{M}_0} f(y)$$

where

$$\mathcal{M}_0 = \{y \in \mathbb{R}^d \mid h(y) = 0\}.$$

Here  $f : \mathcal{M}_0 \rightarrow \mathbb{R}$  and  $h : \mathbb{R}^d \rightarrow \mathbb{R}^m$ . We let  $h_i : \mathbb{R}^d \rightarrow \mathbb{R}$  denote the  $i$ th component function of  $h$  for  $i \in [m]$ . The properties of this kind of manifold (i.e., a manifold defined by the zero set of a function  $h$ ) are covered in “Section 7.7: Manifolds defined by  $h(x) = 0$ ” in [Bou20], so we will borrow heavily from that section in examples.

To illustrate the particularization of a manifold defined by the zero set of a function to an application, we will use the Burer-Monteiro optimization problem:

$$\min_{Y \in \mathcal{M}_p} \langle C, YY^T \rangle$$

where

$$\mathcal{M}_p = \{Y \in \mathbb{R}^{n \times p} \mid \mathcal{A}(YY^T) = b\}$$

Here  $\mathcal{A} : \mathbb{S}^{n \times n} \rightarrow \mathbb{R}^m$  is linear,  $C \in \mathbb{S}^{n \times n}$ , and  $b \in \mathbb{R}^m$ . We can break  $\mathcal{A}$  into coordinate functions via  $A_1, \dots, A_m \in \mathbb{S}^{n \times n}$  by letting  $\mathcal{A}(YY^T)_i = \langle A_i, YY^T \rangle$ . The Burer-Monteiro problem can be viewed as a particularization of the more general problem above by setting  $f(Y) = \langle C, YY^T \rangle$  and  $h(Y) = \mathcal{A}(YY^T) - b$ , or equivalently,  $h_i(Y) = \langle A_i, YY^T \rangle - b_i$ . We define  $f(Y)$  and  $h(Y)$  in this way, where the case of  $Y$  makes it clear whether we are referring to the more general case or the Burer-Monteiro problem.

---

<sup>1</sup>You can see the compilation date on the first page of the PDF of the textbook.

## 1.2 Notation

$\mathbb{S}^{n \times n}$  denotes the space of  $n \times n$  symmetric matrices. We use  $\nabla f, \nabla^2 f$  to denote the Euclidean gradient and Euclidean Hessian respectively,<sup>2</sup> whereas  $\text{grad} f, \text{Hess} f$  denote the Riemannian gradient and Riemannian Hessian. We typically use  $\langle \cdot, \cdot \rangle$  to denote the inner product, but  $A \bullet B$  may be used as well on occasion to denote the trace inner product between matrices ( $A$  and  $B$  in this case).  $I_p$  denotes the  $p \times p$  identity matrix.  $\mathcal{M}$  is used to denote a general or unspecified manifold.

## 1.3 Terminology

[Bou20] is careful to use the term “embedded submanifold” when they specifically mean a submanifold of a linear space, as opposed to a submanifold of a more general manifold (see Section 8.14 in that reference for the latter). Unfortunately these notes may not always adhere to that convention, although generally when we say “submanifold” we mean an embedded submanifold in particular. When in doubt, we advise referring to [Bou20].

## 2 First-order geometry

In this section we build up to defining the Riemannian gradient. First, it is necessary to define the *differential* of a function between linear spaces. Later we will extend this definition to a function between manifolds.

**Definition 2.1 (Differential of a function between linear spaces)** *Let  $F : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  be smooth (infinitely differentiable). The differential<sup>3</sup> of  $F$  at  $x$  is the linear map  $DF(x) : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  defined as:*

$$DF(x)[v] = \lim_{t \rightarrow 0} \frac{F(x + tv) - F(x)}{t} = \left. \frac{d}{dt} F(x + tv) \right|_{t=0}.$$

Intuitively, just think of  $d' \leftarrow 1$ , in which case the differential eats a vector  $v$  and spits out the directional derivative of  $F$  in the direction  $v$ . Thus, it is just the linear map defined by  $\nabla F(x)$ . More generally, the differential tells you “how  $F$  changes” as you start at  $x$  and move along  $v$ . (Specifically, it is the **total derivative** or equivalently the Jacobian viewed as a linear map.) As we will see, this intuition generalizes to the definition of the differential of a map between manifolds.

We can already define what a smooth embedded submanifold (of a linear space) is:

**Definition 2.2 (Smooth embedded submanifold of a linear space)** *Let  $\mathcal{M}$  be a subset of  $\mathbb{R}^d$ .  $\mathcal{M}$  is a smooth embedded submanifold<sup>4</sup> of  $\mathbb{R}^d$  if either:*

---

<sup>2</sup>Technically this isn’t necessary since the Euclidean derivatives coincide with the Riemannian ones in linear spaces, but we will do it for illustrative purposes. [Bou20] uses  $\text{grad} f, \text{Hess} f$  to denote the Euclidean gradient and Euclidean Hessian because  $\nabla$  is used as the symbol for a *connection* on a manifold. However, these notes only briefly cover connections, and I think the clarity obtained by using the usual  $\nabla$  for the Euclidean derivatives outweighs the potential confusion arising from it also being the symbol for a connection.

<sup>3</sup>(3.12) in [Bou20].

<sup>4</sup>Definition 3.6 in [Bou20]. I should mention that Boumal also includes a rather nice illustration to accompany this definition on the same page.

1.  $\mathcal{M}$  is an open subset of  $\mathbb{R}^d$ , in which case it is called an open submanifold.
2. For a fixed integer  $k \geq 1$  and for each  $x \in \mathcal{M}$  there exists a neighborhood  $U$  of  $x$  in  $\mathbb{R}^d$  and a smooth function  $h : U \rightarrow \mathbb{R}^k$  such that
  - (a) If  $y \in U$ , then  $h(y) = 0$  if and only if  $y \in \mathcal{M}$ ; and
  - (b)  $\text{rankDh}(x) = k$ . (In other words,  $\text{Dh}(x)$  is full rank.)

$h$  is called a local defining function for  $\mathcal{M}$  at  $x$ .

For now on we may drop the word “smooth,” but know that when we refer to an “embedded submanifold,” we really always mean a smooth embedded submanifold. (The same omission is made in [Bou20].)

**Example 2.3 ( $\mathcal{M}_0$  is a smooth embedded submanifold if LICQ holds)** *In this case the whole manifold is defined by a single defining function:  $h$ . Condition (a) under Definition 2.2 is trivially satisfied. As for Condition (b), we derive*

$$\text{Dh}(y)[v] = (\langle \nabla h_1(y), v \rangle, \dots, \langle \nabla h_m(y), v \rangle)$$

*Then  $\text{Dh}(Y)$  is full rank if and only if  $\nabla h_1(y), \dots, \nabla h_m(y)$  are linearly independent. In other words, if  $\nabla h_1(y), \dots, \nabla h_m(y)$  are linearly independent at all  $y \in \mathcal{M}_p$ , then  $\mathcal{M}_0$  is indeed a smooth embedded submanifold. Note that this is equivalent to the constraint qualification **LICQ** holding at every point. Or using the nonlinear programming language of [LY16], it is equivalent to saying that every  $y \in \mathcal{M}_0$  is a regular point. (See page 325 in that reference.)*

**Example 2.4 ( $\mathcal{M}_p$  is a smooth embedded submanifold if LICQ holds)** *We have*

$$\begin{aligned} \text{Dh}(Y)[V] &= (\langle \nabla_Y(A_1 \bullet YY^T - b_1), V \rangle, \dots, \langle \nabla_Y(A_m \bullet YY^T - b_m), V \rangle) \\ &= 2(\langle A_1 Y, V \rangle, \dots, \langle A_m Y, V \rangle) \end{aligned}$$

*Then  $\text{Dh}(Y)$  is full rank if and only if  $A_1 Y, \dots, A_m Y$  are linearly independent. In other words, if  $A_1 Y, \dots, A_m Y$  are linearly independent at all  $Y \in \mathcal{M}_p$ , then  $\mathcal{M}_p$  is indeed a smooth embedded submanifold.*

Regarding Definition 2.2: It turns out that it is necessary for  $\text{Dh}(x)$  to be full rank on  $x$ —it is *not* sufficient for  $\text{Dh}(x)$  to just have constant rank at all  $x \in \mathcal{M}_p$ . However, it *is* sufficient for  $\text{Dh}(x)$  to have constant rank if it has constant rank not just on the zero set of  $h$  but on a neighborhood of the zero set as well.<sup>5</sup> This is exactly Case (b) in Assumption 1.1 in [BVB18]! (Case (a) in Assumption 1.1 is of course just Condition (b) in Definition 2.2.) This is all to say that the manifold approach can give you a way to deal with “redundant constraints” which the usual nonlinear programming approach cannot (as far as we know).

Finally, it is also worth noting that Condition (b) in Definition 2.2 can be shown to be equivalent to Definition 1 in [WW20]. This is to say that everything here is standard even if it might appear in slightly different forms in the literature.

---

<sup>5</sup>See Sections 3.10 and 8.14 in [Bou20].

**Definition 2.5 (Velocity of a curve)** Let  $c : I \subseteq \mathbb{R} \rightarrow \mathbb{R}^d$  be a curve. Its velocity<sup>6</sup> at time  $t$  is  $c'(t) = \frac{d}{dt}c(t)$ . (This is just the ordinary Euclidean derivative.)

**Definition 2.6 (Tangent space, tangent vector)** Let  $\mathcal{M}$  be a smooth embedded submanifold of  $\mathbb{R}^d$ . For all  $x \in \mathcal{M}$ , define

$$T_x\mathcal{M} = \{c'(0) \mid c : I \subseteq \mathbb{R} \rightarrow \mathcal{M} \text{ is smooth around } 0 \text{ and } c(0) = x\}.$$

In other words,  $v \in T_x\mathcal{M}$  if and only if there exists a smooth curve on  $\mathcal{M}$  passing through  $x$  with velocity  $v$ .<sup>7</sup>  $T_x\mathcal{M}$  is called the tangent space to  $\mathcal{M}$  at  $x$ .  $v \in T_x\mathcal{M}$  is called a tangent vector to  $\mathcal{M}$  at  $x$ .<sup>8</sup>

It can be shown that the tangent space is a linear space (aka a vector space).

The fact that the tangent space is a linear space is key because it allows us to easily extend many definitions from Euclidean spaces to manifolds.

**Theorem 2.7 (Expression for the tangent space of a smooth embedded submanifold)** Let  $\mathcal{M} \subseteq \mathbb{R}^d$  be an embedded submanifold. Then for  $x \in \mathcal{M}$ , if  $M$  is an open submanifold (first case in Definition 2.2), then  $T_x\mathcal{M} = \mathbb{R}^d$ . Otherwise,  $T_x\mathcal{M} = \ker Dh(x)$  where  $h$  is any local defining function at  $x$ .<sup>9</sup>

Note that this is exactly the definition of the “tangent plane” given in the nonlinear programming textbook [LY16]—see page 325! Thus, we see that notions in smooth manifold optimization (smooth embedded submanifold, tangent space) reduce to notions in nonlinear programming (LICQ, tangent plane), although the manifold approach is more general and gives us better geometric intuition.

**Example 2.8 (Tangent space to  $\mathcal{M}_0$ )** Recall from Example 2.3 that we have

$$Dh(y)[v] = (\langle \nabla h_1(y), v \rangle, \dots, \langle \nabla h_m(y), v \rangle)$$

Theorem 2.7 instantly gives

$$T_y\mathcal{M}_0 = \left\{ \dot{v} \in \mathbb{R}^d \mid \langle \nabla h_i(y), \dot{v} \rangle = 0 \text{ for } i \in [m] \right\}.$$

**Example 2.9 (Tangent space to  $\mathcal{M}_p$ )** Recall from Example 2.4 that we have

$$Dh(Y)[V] = 2(\langle A_1Y, V \rangle, \dots, \langle A_mY, V \rangle).$$

Theorem 2.7 instantly gives

$$T_Y\mathcal{M}_p = \left\{ \dot{V} \in \mathbb{R}^{n \times p} \mid \langle A_iY, \dot{V} \rangle = 0 \text{ for } i \in [m] \right\}.$$

As a check, note that this is exactly how they define the tangent space in [BVB18] (Lemma 2.1). It is also exactly how the tangent space is defined in [WW20] (page 5), although they express it in a more complicated way.

<sup>6</sup>(3.12) in [Bou20].

<sup>7</sup>Definition 3.7 in [Bou20].

<sup>8</sup>Definition 3.10 in [Bou20].

<sup>9</sup>Theorem 3.8 in [Bou20].

**Definition 2.10 (Normal space)** Let  $\mathcal{M} \subseteq \mathbb{R}^d$  be an embedded submanifold. The normal space<sup>10</sup> at  $x$ , denoted  $N_x\mathcal{M}$ , is simply the orthogonal complement of  $T_x\mathcal{M}$  in  $\mathbb{R}^d$ .

**Example 2.11 (Normal space to  $\mathcal{M}_0$ )** Referencing Example 2.8, it is easy to see that

$$N_y\mathcal{M}_0 = \text{span}\{\nabla h_1(y), \dots, \nabla h_m(y)\}.$$

**Example 2.12 (Normal space to  $\mathcal{M}_p$ )** Referencing Example 2.9, it is easy to see that

$$N_Y\mathcal{M}_p = \text{span}\{A_1Y, \dots, A_mY\}.$$

**Definition 2.13 (Dimension of a manifold)** The dimension of a manifold  $\mathcal{M} \subseteq \mathbb{R}^d$ , denoted  $\dim \mathcal{M}$ , is just the dimension of its tangent space.<sup>11</sup> Note that the tangent space is linear (aka a vector space) and the dimension of  $T_x\mathcal{M}$  does not depend on  $x$ , so this is well-defined. In particular, this means that  $\dim \mathcal{M} = d$  if  $\mathcal{M}$  is an open embedded submanifold, and  $\dim \mathcal{M} = \dim \ker \text{Dh}(x) = d - \text{rankDh}(x)$  otherwise (due to Theorem 2.7).

**Example 2.14 (Dimension of  $\mathcal{M}_0$ )** Remember that we need  $\text{Dh}(y)$  to be full rank for  $\mathcal{M}_0$  to be a smooth embedded submanifold in the first place (see Example 2.3). Since  $\text{Dh}(y) : \mathbb{R}^d \rightarrow \mathbb{R}^m$ , we have that  $\text{rankDh}(y) = m$ . So  $\dim \mathcal{M}_0 = d - m$ .

**Example 2.15 (Dimension of  $\mathcal{M}_p$ )** It immediately follows from Example 2.14 that  $\dim \mathcal{M}_p = np - m$ .

As a check, this is confirmed by [BVB18] (Proposition 1.2) and [WW20] (page 5). (Although [BVB18] makes things slightly more complicated because they allow  $\text{Dh}(Y)$  to not be full rank as long as it has constant rank on a neighborhood of  $\mathcal{M}_p$  [see the discussion of this below Example 2.4 above]. This affects  $\dim \mathcal{M}_p$ , although it still just ends up being  $d - \text{rankDh}(x)$ .)

We now define what it means for a map between two manifolds to be smooth, which we will do by smoothly extending the function in question from the manifold to the linear embedding space (also called the ambient space). In fact, the notion of a smooth extension will end up being the more important part of the definition below for our purposes!

**Definition 2.16 (Smooth map between manifolds, smooth extension)** Let  $\mathcal{M} \subseteq \mathbb{R}^d$  and  $\mathcal{M}' \subseteq \mathbb{R}^{d'}$  be embedded submanifolds. A map  $F : \mathcal{M} \rightarrow \mathcal{M}'$  is smooth if there exists a function  $\bar{F} : U \rightarrow \mathbb{R}^{d'}$  which is smooth (in the usual sense, aka infinitely differentiable) on a neighborhood  $U$  of  $\mathcal{M}$  in  $\mathbb{R}^d$ . Furthermore, this map  $\bar{F}$  must be such that  $F$  and  $\bar{F}$  coincide on  $\mathcal{M} \cap U$ , i.e.,  $F(y) = \bar{F}(y)$  for all  $y \in \mathcal{M} \cap U$ . (Equivalently,  $F$  is the restriction of  $\bar{F}$  to  $\mathcal{M}$ :  $F = \bar{F}|_{\mathcal{M}}$ .) We call  $\bar{F}$  a smooth extension of  $F$ .<sup>12</sup>

To gain intuition for the above definition, it is best to think of  $\mathcal{M}'$  as just being  $\mathbb{R}$ . (Note that  $\mathbb{R}$  is an open embedded submanifold per Definition 2.2.)

<sup>10</sup>Defined in the paragraph above (5.15) in [Bou20].

<sup>11</sup>Theorem 3.10 in [Bou20].

<sup>12</sup>Proposition 3.24 in [Bou20]. Note that Boumal first defines what it means for  $F$  to be smooth at a point in Definition 3.23, but I have omitted this for simplicity.

**Example 2.17 (Smooth extension of our objective function on  $\mathcal{M}_p$ )** Clearly a smooth extension of  $f : \mathcal{M}_p \rightarrow \mathbb{R}$  defined as  $f(Y) = \langle C, YY^T \rangle$  is  $\bar{f} : \mathbb{R}^{np} \rightarrow \mathbb{R}$  defined as  $\bar{f}(Y) = \langle C, YY^T \rangle$ . In other words, the only thing that changes when we take a smooth extension is the domain! Thus, we can really view ourselves as starting with a smooth extension in this application, and then obtaining the objective function on the manifold by restricting the original function to the manifold.

We are now ready to extend Definition 2.1 to manifolds.

**Definition 2.18 (Differential of a map between manifolds)** The differential of  $F : \mathcal{M} \subseteq \mathbb{R}^d \rightarrow \mathcal{M}' \subseteq \mathbb{R}^{d'}$  at  $x$  is a linear operator  $DF(x) : T_x\mathcal{M} \rightarrow T_{F(x)}\mathcal{M}'$  defined by

$$DF(x)[v] = \left. \frac{d}{dt} F(c(t)) \right|_{t=0},$$

where  $c$  is a smooth curve on  $\mathcal{M}$  passing through  $x$  at  $t = 0$  with velocity  $v$ .<sup>13</sup>

Intuitively,  $DF(x)$  is just a map which tells you how  $F$  changes as you start at  $x$  and move along the manifold  $\mathcal{M}$  in the direction  $v \in T_x\mathcal{M}$ . Although we can't actually move along  $v$ —we have to move along a curve with velocity  $v$  called  $c(t)$ . Then,  $F(c(t))$  is just a curve on  $\mathcal{M}$ , and we would like to return how  $F$  changes as we move along this curve. Well, let's just return the velocity of this curve at  $t = 0$ —this is exactly  $DF(x)[v]$ .

It is worth mentioning that the differential of Theorem 2.18 satisfies all of the properties we would expect it to satisfy (the chain rule, the product rule, etc.). See Section 4.7 in [Bou20] for details.

The next theorem provides even more intuition:

**Theorem 2.19 (Differential of a map between manifolds via a smooth extension)** With the same notation as in 2.18, let  $\bar{F} : U \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  be a smooth extension of  $F$ . (We know such an extension exists since  $F$  is smooth.) Then  $DF(x)$  is the restriction of  $D\bar{F}(x)$  to  $T_x\mathcal{M}$ , i.e.,

$$DF(x) = D\bar{F}(x)|_{T_x\mathcal{M}}.$$

Note that  $D\bar{F}(x)$  is a differential between linear spaces, so it is defined via Definition 2.1.<sup>14</sup>

In a sense, this says that we could equivalently define the differential between manifolds via smooth extensions. This also helps us gain much better intuition as to what  $DF(x)$  is if  $F$  is a map between manifolds. Let  $d' \leftarrow 1$  so that  $F : \mathcal{M} \subseteq \mathbb{R}^d \rightarrow \mathcal{M}' \subseteq \mathbb{R}$ . Let  $\bar{F}$  be a smooth extension of  $F$ . Recall from the discussion below Definition 2.1 that  $D\bar{F}(x)$  eats a vector  $v$  and spits out the directional derivative of  $\bar{F}$  in the direction  $v$ . Then, Theorem 2.19 basically says that  $DF(x)$  does the exact same thing, except you are only allowed to feed in vectors in  $T_x\mathcal{M}$ . This restriction makes sense because intuitively, we are only allowed to move in directions that are in  $T_x\mathcal{M}$ , so  $DF(x)$  is essentially a fancy way of forcing us to only pay attention to how  $F$  changes in directions which we actually care about!

<sup>13</sup>Definition 3.27 in [Bou20]. Once again, there is an accompanying illustration which is excellent!

<sup>14</sup>Proposition 3.28 in [Bou20].

**Example 2.20 (Differential of the objective function over  $\mathcal{M}_p$ )** Define  $\bar{f}$  as in Example 2.17. Using Definition 2.1 (and particularly the discussion below it), we have that for  $Y \in \mathbb{R}^{np}$ ,  $D\bar{f}(Y)[V] : \mathbb{R}^{np} \rightarrow \mathbb{R}$  is defined via

$$\begin{aligned} D\bar{f}(Y)[V] &= \langle \nabla_Y \bar{f}(Y), V \rangle \\ &= 2 \langle CY, V \rangle \end{aligned}$$

Then using Theorem 2.19, we have that  $Df(Y) : T_Y \mathcal{M}_p \rightarrow \mathbb{R}$  is

$$Df(Y) = 2 \langle CY, V \rangle|_{T_Y \mathcal{M}_p}$$

where  $T_Y \mathcal{M}_p$  is defined in Example 2.9. We will obtain a more explicit expression for  $Df(Y)$  after we define the Riemannian gradient!

We need only a few more definitions before we can define the Riemannian gradient!

**Definition 2.21 (Tangent bundle)** The tangent bundle<sup>15</sup>  $T\mathcal{M}$  of a manifold  $\mathcal{M}$  is the disjoint union of the tangent spaces of  $\mathcal{M}$ :

$$T\mathcal{M} = \{(x, v) \mid x \in \mathcal{M} \text{ and } v \in T_x \mathcal{M}\}.$$

Following the conventions of [Bou20], we may abuse notation at times and conflate  $v$  and  $(x, v)$  for a tangent vector  $v \in T_x \mathcal{M}$ . We may write  $(x, v) \in T_x \mathcal{M}$ , or even  $v \in T\mathcal{M}$  if it is clear from context that the foot or base of  $v$  is  $x$ .

**Theorem 2.22** If  $\mathcal{M}$  is an embedded submanifold of  $\mathbb{R}^d$ , then the tangent bundle  $T\mathcal{M}$  is an embedded submanifold of  $\mathbb{R}^d \times \mathbb{R}^d$ .<sup>16</sup> (See Appendix B for more information about product manifolds.)

We don't really need Theorem 2.22 in these notes beyond one use case: In multiple places (e.g., Definition 2.23, Definition 4.1), we will need to introduce a map from or to  $T\mathcal{M}$ , and we will want such a map to be smooth per Definition 2.16. Since Definition 2.16 only applies to maps between embedded submanifolds, we need to recognize that  $T\mathcal{M}$  is itself an embedded submanifold.

**Definition 2.23 (Vector field on a manifold)** A vector field<sup>17</sup> on a manifold  $\mathcal{M}$  is a map  $V : \mathcal{M} \rightarrow T\mathcal{M}$  such that  $V(x) \in T_x \mathcal{M}$  for all  $x \in \mathcal{M}$ . If  $V$  is a smooth map, we say it is a smooth vector field.

In other words, a vector field is just a formal way of assigning to each  $x \in \mathcal{M}$  a tangent vector at that point.

---

<sup>15</sup>Definition 3.35 in [Bou20].

<sup>16</sup>Theorem 3.36 in [Bou20].

<sup>17</sup>Definition 3.37 in [Bou20].



**Definition 2.24 (Riemannian gradient)** Let  $f : \mathcal{M} \rightarrow \mathbb{R}$  be a smooth function on a Riemannian manifold<sup>18</sup>  $\mathcal{M}$ . The Riemannian gradient<sup>19</sup> of  $f$  is the vector field  $\text{grad}f$  on  $\mathcal{M}$  uniquely defined by the following: For all  $(x, v) \in \text{T}\mathcal{M}$ ,

$$\text{D}f(x)[v] = \langle \text{grad}f(x), v \rangle,$$

where  $\text{D}f(x)$  is as defined in 2.18 (or Theorem 2.19).

Note that this definition critically implies that  $\text{grad}f(x) \in \text{T}_x\mathcal{M}$ , i.e., it is a tangent vector.

**Definition 2.25 (Orthogonal projector to the tangent space)** Let  $\mathcal{M} \subseteq \mathbb{R}^d$  be an embedded submanifold. We denote the orthogonal projector<sup>20</sup> from  $\mathbb{R}^d$  to  $\text{T}_x\mathcal{M}$  as

$$\text{Proj}_x : \mathbb{R}^d \rightarrow \text{T}_x\mathcal{M} \subseteq \mathbb{R}^d.$$

Being a projector,  $\text{Proj}_x$  is a linear operator such that  $\text{Proj}_x \circ \text{Proj}_x = \text{Proj}_x$ . Being orthogonal, the matrix corresponding to  $\text{Proj}_x$  is symmetric.

**Example 2.26 (Deriving  $\text{Proj}_y$  for  $\mathcal{M}_0$ )** Recall from Theorem 2.7 that

$$\text{T}_y\mathcal{M}_0 = \ker \text{D}h(y).$$

Also, *recall* that the orthogonal projector onto the kernel of a matrix  $A$  can be represented as  $I - A^+A$ , where  $A^+$  denotes the Moore-Penrose inverse of  $A$  (aka the pseudoinverse). Thus, one expression for  $\text{Proj}_y$  is

$$\text{Proj}_y(V) = V - \text{D}h(y)^+[\text{D}h(y)[v]]. \quad (1)$$

We will also derive a second expression for  $\text{Proj}_y$ . Note that equivalently,  $\text{T}_y\mathcal{M}_p = \ker((\text{D}h(y)^T)^T)$ . *Recall* that the orthogonal projector onto the kernel of the transpose of a matrix,  $A^T$ , can be expressed as  $I - AA^+$ . Setting  $A \leftarrow \text{D}h(y)^T$ , this gives

$$\text{Proj}_y(v) = v - \text{D}h(y)^T [(\text{D}h(y)^T)^+[v]]. \quad (2)$$

(I am giving both of these expressions because although (1) is simpler, Boumal actually only gives (2) in Section 7.7 of [Bou20]. It might be because typically  $\text{D}h(V)^T$  is typically simpler to work with, although I'm not sure.)

**Example 2.27 (Deriving  $\text{Proj}_Y$  for  $\mathcal{M}_p$ )** Particularizing (1) or (2) to the Burer-Monteiro problem just involves plugging in what  $\text{D}h(Y)$  is in that context—see Example 2.4.

It is worth mentioning that at first glance, neither (1) nor (2) look equivalent to the expression for  $\text{Proj}_Y$  given in Lemma 2.2 in [BVB18]. Indeed, [BVB18] uses a more unorthodox method to derive their expression for  $\text{Proj}_Y$ . However, their expression for  $\text{Proj}_Y$  can be shown to be equivalent to (2) using *these* properties of the Moore-Penrose inverse.

<sup>18</sup>See Section 2.1 for what a Riemannian manifold is.

<sup>19</sup>Definition 3.50 in [Bou20].

<sup>20</sup>Definition 3.52 in [Bou20].

And finally:

**Theorem 2.28 (Connecting the Riemannian gradient and Euclidean gradient)** *Let  $\mathcal{M} \subseteq \mathbb{R}^d$  be a Riemannian submanifold<sup>21</sup> of  $\mathbb{R}^d$ . The Riemannian gradient of  $f$  is given by*

$$\operatorname{grad}f(x) = \operatorname{Proj}_x(\nabla\bar{f}(x)),$$

where  $\bar{f}$  is any smooth extension of  $f$  to a neighborhood of  $\mathcal{M}$  in  $\mathbb{R}^d$ .<sup>22</sup>

It is not hard to see that Theorem 2.28 follows quickly from Definition 2.24 and Theorem 2.19.

**Example 2.29 (Riemannian gradient in the optimization problem over  $M_0$ )** *Let  $\bar{f} : \mathbb{R}^d \rightarrow \mathbb{R}$  denote a smooth extension of  $f : \mathcal{M}_0 \rightarrow \mathbb{R}$ . Theorem 2.28 yields*

$$\operatorname{grad}f(y) = \nabla\bar{f}(y) - Dh(y)^T [(Dh(y)^T)^+[\nabla\bar{f}(y)]]$$

To really see what is going on here, it is helpful to reexpress this as follows:

$$\operatorname{grad}f(y) = \nabla\bar{f}(y) - \sum_{i=1}^m \lambda_i(y) \nabla h_i(y) \tag{3}$$

where

$$\lambda(y) := (Dh(y)^T)^+[\nabla\bar{f}(y)].$$

(Recall from Example 2.3 that the  $i$ th column of  $Dh(y)^T$  is simply  $\nabla h_i(y)$ .)

How should we think of the “multipliers” that make up  $\lambda(Y)$ ? Well, remember how we derived (3) in the first place via Theorem 2.28—we took  $\nabla\bar{f}(Y)$  and deleted its normal component. Thus, you should just think of the  $\lambda_i(y)$ ’s as being the multipliers which delete the component of  $\nabla\bar{f}(y)$  that lies in

$$N_y\mathcal{M}_0 = \operatorname{span}\{\nabla h_1(y), \dots, \nabla h_m(y)\}.$$

Indeed, this is exactly what the pseudoinverse does. And this makes complete sense because  $\operatorname{grad}f(Y)$  must be a tangent vector by definition, so the  $\lambda_i(y)$ ’s couldn’t be anything else.

It is also basically immediately apparent at this point from (3) that the first-order criticality condition from nonlinear programming is equivalent to  $\operatorname{grad}f(y) = 0$ . More on this later.

**Example 2.30 (Riemannian gradient in the Burer-Monteiro problem)** *Define  $\bar{f}$  as in Example 2.17. Theorem 2.28 yields*

$$\begin{aligned} \operatorname{grad}f(Y) &= \nabla\bar{f}(Y) - Dh(Y)^T [(Dh(Y)^T)^+[\nabla\bar{f}(Y)]] \\ &= 2CY - Dh(Y)^T [(Dh(Y)^T)^+[2CY]]. \end{aligned}$$

<sup>21</sup>See Section 2.1 for what a Riemannian manifold is.

<sup>22</sup>Proposition 3.53 in [Bou20].

This is not particularly informative, but following the same path as in Example 2.29, we can reexpress this as

$$\begin{aligned}\operatorname{grad}f(Y) &= 2 \left[ CY - \sum_{i=1}^m \lambda_i(Y) A_i Y \right] \\ &= 2 \left[ \left( C - \sum_{i=1}^m \lambda_i(Y) A_i \right) Y \right]\end{aligned}$$

where

$$\lambda(Y) := (Dh(Y)^T)^+ [2CY].$$

As a sanity check, this is exactly the expression derived in [BVB18]—see (2.6) and (2.5) in that paper. Although note that due to their unorthodox expression for  $\operatorname{Proj}_Y$  (see the discussion of this in Example 2.27), their expression for  $\lambda(Y)$  (which they call  $\mu$ ) is different (though equivalent).

## 2.1 First-order geometry: material not covered

This section covered a large portion of Chapter 3 from [Bou20] as well as parts of Section 7.7 from [Bou20] through the examples. That said, we still made some omissions. Probably the most important omission has to do with the notation of a “Riemannian submanifold,” which was mentioned in several places (e.g., Definition 2.24, Theorem 2.28). Basically, [Bou20] stresses that you can equip manifolds with different inner products. More accurately, you equip each tangent space with an inner product (since tangent spaces are vector spaces), and the inner product on  $T_x\mathcal{M}$  can even depend on  $x$ .<sup>23</sup> If the inner product on  $T_x\mathcal{M}$  varies smoothly with  $x$ , then we call it a *Riemannian metric*.<sup>24</sup> A manifold equipped with a Riemannian metric is a *Riemannian submanifold*.<sup>25</sup> It turns out that if you take an embedded submanifold of  $\mathbb{R}^d$  and define the inner products on the tangent spaces of the embedded submanifold to just be the restrictions of an inner product on  $\mathbb{R}^d$  to those spaces, you obtain a Riemannian submanifold. We refer to this particular type of Riemannian submanifold, where the inner product on  $T_x\mathcal{M}$  doesn’t vary with  $x$  and is inherited from the inner product on  $\mathbb{R}^d$ , a *Riemannian submanifold of  $\mathbb{R}^d$* .<sup>26</sup> Because, as noted in [Bou20], this is the most common type of Riemannian submanifold in practice and it is very natural to work with, we are only covering it briefly here.

It is important to make clear that in these notes, if a theorem/definition in [Bou20] involves a Riemannian submanifold, our “version” of the theorem/definition will basically always additionally assume that it is a Riemannian submanifold of  $\mathbb{R}^d$ . This way we don’t need to add the point where you are taking the inner product as a subscript to each inner product. For example, compare Definition 2.24 to Definition 3.50 in [Bou20]: The latter writes  $\langle v, \operatorname{grad}f(x) \rangle_x$  whereas we omit the subscript  $x$ . Once again, we do this because Riemannian submanifolds of linear spaces are the most common use case, and it also helps to think about theorems in terms of the inner products we are all used to (i.e., inner products inherited from linear spaces).

<sup>23</sup>Definition 3.43 in [Bou20].

<sup>24</sup>Definition 3.44 in [Bou20].

<sup>25</sup>Definition 3.45 in [Bou20].

<sup>26</sup>Definition 3.47 in [Bou20].

We also didn't cover retractions (Section 3.6 in [Bou20])—this will wait until a later section of these notes.

## 2.2 First-order geometry: extending beyond the embedded case

To be done.

## 3 Second-order geometry

In this section we build up to defining the Riemannian Hessian, although we will not end up defining it in full generality. This is due to the fact that the general definition takes a lot of work—you need to build up certain mathematical objects which satisfy nice properties and then prove uniqueness, etc. As expected however, working with the Riemannian Hessian on embedded submanifolds is much easier than in the general case, so we rarely need this machinery in practice. Still, we have sought to give a flavor of how it is actually defined.

Toward doing this, it will be helpful to remember how we should think about the regular Euclidean Hessian,  $\nabla^2 f(x)$  for  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ . It can be thought of as an  $n \times n$  matrix of second-order partial derivatives, but just as it is typically better to think of  $\nabla f(x)$  in terms of the linear map it induces, it is better to think of  $\nabla^2 f(x)$  as a linear map. Indeed, we can express  $\nabla^2 f(x)$  as a linear map,  $\nabla^2 f(x) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , where  $\nabla^2 f(x)[v]$  is precisely the directional derivative of  $\nabla f(x)$  along  $v$ . (Of course, behind the scenes  $\nabla^2 f(x)[v]$  just corresponds to multiplying the  $n \times n$  matrix “version” of  $\nabla^2 f(x)$  by  $v \in \mathbb{R}^d$ . For a more in-depth discussion of all of this, see the beginning of Chapter 5 in [Bou20].) Likewise, the Riemannian Hessian will end up being a linear map,  $\text{Hess}f(x) : T_x\mathcal{M} \rightarrow T_x\mathcal{M}$ , where  $\text{Hess}f(x)[v]$  should somehow be the directional derivative (whatever that means for manifolds) of the Riemannian gradient  $\text{grad}f(x)$  along  $v$ .

A good first try for defining the Riemannian Hessian is to take the differential, per Definition 2.18, of the vector field  $\text{grad}f(x)$ . After all, remember, per Definitions 2.24 and 2.23, that  $\text{grad}f(x)$  is really a map:  $\text{grad}f(x) : \mathcal{M} \rightarrow T\mathcal{M}$ . In other words, it is a map between two embedded submanifolds,<sup>27</sup> and Definition 2.18 tells us exactly how to take the differential of a map between manifolds. However, there is a problem. Note that  $D(\text{grad}f)(x) : T_x\mathcal{M} \rightarrow T_{\text{grad}f(x)}T\mathcal{M}$  per Definition 2.18. So if we define  $\text{Hess}f(x)$  as  $D(\text{grad}f)(x)$ , the domain will indeed be  $T_x\mathcal{M}$  as desired, but it is not clear that the codomain,  $T_{\text{grad}f(x)}T\mathcal{M}$ , is equivalent to  $T_x\mathcal{M}$ , the latter being what we want the codomain to be. In fact they aren't equivalent, as proven through a simple example in Section 5.1 of [Bou20].

So to recap, we have a vector field  $\text{grad}f(x)$ , and we want to take the derivative of it in a way such that the codomain is  $T_x\mathcal{M}$ . The differential does not do this, so we need a new kind of derivative, which we will call a *connection*.

**Definition 3.1 (Connection, informal)** *A connection on a manifold  $\mathcal{M}$  is an operator*

$$\nabla : T\mathcal{M} \times \mathcal{X}(\mathcal{M}) \rightarrow T\mathcal{M},$$

---

<sup>27</sup> $T\mathcal{M}$  is indeed an embedded submanifold—see Theorem 3.36 in [Bou20].

where  $\mathcal{X}(\mathcal{M})$  is the set of all smooth vector fields on  $\mathcal{M}$ .<sup>28</sup> So that we can talk about specific inputs and outputs, we will say that it maps  $(u, V) \in \mathbb{T}\mathcal{M} \times \mathcal{X}(\mathcal{M})$  to  $\nabla_u V \in \mathbb{T}\mathcal{M}$ . It must satisfy the property that if  $u \in \mathbb{T}_x\mathcal{M}$ , then  $\nabla_u V \in \mathbb{T}_x\mathcal{M}$ —in other words, the tangent vector that it spits out must lie in the same tangent space as the tangent vector given in the input. Furthermore, a connection needs to satisfy a list of other properties (smoothness, linearity, etc.) so that it really behaves like a derivative. See Definition 5.1 in [Bou20] for the full definition of a connection, which includes a list of these properties.

So we have a notion of the derivative of a vector field which enforces our rule that it should map vectors in  $\mathbb{T}_x\mathcal{M}$  to vectors in  $\mathbb{T}_x\mathcal{M}$ . Indeed, for a vector field  $V$  over a manifold  $\mathcal{M}$ , this is obtained by fixing the second argument in the connection to be  $V$ , so that for a given  $x \in \mathcal{M}$  we obtain a map from  $\mathbb{T}_x\mathcal{M}$  to  $\mathbb{T}_x\mathcal{M}$ . We are only a few steps away from being able to define the Riemannian Hessian. It turns out that the list of properties a connection must satisfy (see Definition 5.1 in [Bou20]) is not expansive enough for our purposes. In particular, we can only define the Riemannian Hessian (and the Riemannian gradient for that matter) after introducing an inner product which makes our manifold into a Riemannian manifold.<sup>29</sup> We would like the connection we use to define the Riemannian Hessian to play nicely with this inner product/metric. This is the subject of the next theorem:

**Theorem 3.2 (Fundamental theorem of Riemannian geometry, informal)** *On a Riemannian manifold  $\mathcal{M}$ , there exists a unique connection which “plays nicely” with the metric on your manifold. This connection is called the Levi-Civita or Riemannian connection. See Theorem 5.5 in [Bou20] for a formal statement.*

Now we can define the Riemannian Hessian:

**Definition 3.3 (Riemannian Hessian)** *Let  $\mathcal{M}$  be a Riemannian manifold with its Riemannian connection  $\nabla$ . Let  $f : \mathcal{M} \rightarrow \mathbb{R}$  be a smooth function. The Riemannian Hessian<sup>30</sup> of  $f$  at  $x \in \mathcal{M}$  is a linear operator  $\text{Hess}f(x) : \mathbb{T}_x\mathcal{M} \rightarrow \mathbb{T}_x\mathcal{M}$  defined as follows*

$$\text{Hess}f(x)[u] = \nabla_u \text{grad}f.$$

Recall from Definition 3.1 that  $\nabla_u \text{grad}f$  denotes the result you get from feeding  $u \in \mathbb{T}_x\mathcal{M}$  and  $\text{grad}f \in \mathcal{X}(\mathcal{M})$  into the connection  $\nabla$  as its input.

The following intuitively had to be true, but let’s make it explicit:

**Theorem 3.4 (The Riemannian Hessian is self-adjoint)** *The Riemannian Hessian is self-adjoint.<sup>31</sup> In particular, when we are working with real numbers, the matrix corresponding to the Riemannian Hessian is symmetric.*

<sup>28</sup>[Bou20] uses similar notation for the set of all vector fields on  $\mathcal{M}$ —see Definition 3.37.

<sup>29</sup>See Section 2.1.

<sup>30</sup>Definition 5.13 in [Bou20].

<sup>31</sup>Proposition 5.16 in [Bou20].

Now we begin the process of simplifying all we have done above for the case where  $\mathcal{M}$  is a Riemannian submanifold of a Euclidean space. First, we obtain an explicit expression for the Riemannian connection:

**Theorem 3.5 (Riemannian connection of a Riemannian submanifold of  $\mathbb{R}^d$ )** *Let  $\mathcal{M}$  be a Riemannian submanifold of a Euclidean space (i.e.,  $\mathbb{R}^d$ ). Then the Riemannian connection is precisely<sup>32</sup>*

$$\nabla_u V = \text{Proj}_x(D\bar{V}(x)[u]),$$

where  $\bar{V}$  is any smooth extension of  $V$ . (Here  $V$  is a smooth vector field on  $\mathcal{M}$  and  $u \in T\mathcal{M}$ .)

This is incredibly intuitive based on the discussion we had toward the beginning of this chapter. Remember that the whole problem with using the differential to obtain the derivative of a vector field was that the result might not be in  $T_x\mathcal{M}$ . Well, Theorem 3.5 is basically saying that the “correct” way to fix this for Riemannian submanifolds of Euclidean spaces is to just force the output to lie in  $T_x\mathcal{M}$  through orthogonal projection!

Definition 3.3 and Theorem 3.5 combine to instantly yield the following corollary:

**Corollary 3.6 (Riemannian Hessian on a Riemannian submanifold of  $\mathbb{R}^d$ )** *Let  $\mathcal{M}$  be a Riemannian submanifold of the Euclidean space  $\mathbb{R}^d$ . Let  $f : \mathcal{M} \rightarrow \mathbb{R}$  be a smooth function. Let  $\overline{\text{grad}f} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  denote a smooth extension of  $\text{grad}f$ . Then<sup>33</sup>*

$$\text{Hess}f(x)[u] = \text{Proj}_x(D(\overline{\text{grad}f})(x)[u]).$$

In other words, for a Riemannian submanifold of  $\mathbb{R}^d$ , you can compute the Riemannian Hessian by differentiating *in the classical sense* (i.e., using Definition 2.1) a smooth extension of the Riemannian gradient vector field, and then orthogonally projecting the result to  $T_x\mathcal{M}$ .

Note that  $D(\overline{\text{grad}f})(x)[u]$  is *not* equivalent to  $\nabla^2 f(x)[u]$ , meaning Corollary 3.6 is *not* saying that  $\text{Hess}f(x)[u]$  is just the projection of the Euclidean Hessian to  $T_x\mathcal{M}$ , aka  $\text{Proj}_x(\nabla^2 f(x)[u])$ . In particular,  $D(\overline{\text{grad}f})(x)[u]$  and  $\nabla^2 f(x)[u]$  are not equivalent because  $\overline{\text{grad}f}$  and  $\nabla f$  are not equivalent. The latter are not equivalent because they need not coincide on the manifold  $\mathcal{M}$  itself. Indeed, by definition  $\overline{\text{grad}f}$  and  $\text{grad}f$  must coincide on  $\mathcal{M}$ , and we know from Theorem 3.5 that  $\text{grad}f$  and  $\nabla f$  need not coincide on  $\mathcal{M}$ . (After all, what then would even be the point of defining the Riemannian gradient, as it would always be equivalent to the Euclidean gradient?)

However, it is still possible to directly relate the Riemannian Hessian to the Euclidean Hessian. We only provide an informal statement here:

**Theorem 3.7 (Relating the Riemannian Hessian to the Euclidean Hessian, informal)** *Let  $\mathcal{M}$  be a Riemannian submanifold of  $\mathbb{R}^d$ . Let  $f : \mathcal{M} \rightarrow \mathbb{R}$  be smooth with smooth extension  $\bar{f} : \mathbb{R}^d \rightarrow \mathbb{R}$ . The Riemannian Hessian of  $f$  is given by:*

$$\text{Hess}f(x)[u] = \text{Proj}_x(\nabla^2 \bar{f}(x)[u]) + \text{correction}(\text{Proj}_x^\perp(\nabla \bar{f}(x))),$$

---

<sup>32</sup>Theorem 5.8 in [Bou20].

<sup>33</sup>Corollary 5.14 in [Bou20].

where  $\text{Proj}_x^\perp = \text{Id} - \text{Proj}_x$  denotes the orthogonal projector onto the normal space  $N_x\mathcal{M}$ , and correction  $: \mathbb{R}^d \rightarrow T_x\mathcal{M}$  is some “correctional function” which we do not specify here. (Id is the identity map.)

See Corollary 5.17 in [Bou20] for the formal statement.

In other words, this says that for Riemannian submanifolds of Euclidean spaces, the Riemannian Hessian is the projected Euclidean Hessian plus a correctional term which depends only on the normal part of the Euclidean gradient.

**Example 3.8 (Riemannian Hessian in the optimization problem over  $\mathcal{M}_0$ )** Using the notation of Example 2.29, recall that

$$\text{grad}f(y) = \nabla\bar{f}(y) - \sum_{i=1}^m \lambda_i(y) \nabla h_i(y) \quad (4)$$

where

$$\lambda(y) := (Dh(y)^T)^+ [\nabla\bar{f}(y)]. \quad (5)$$

We now wish to apply Corollary 3.6 to obtain an expression for  $\text{Hess}f(x)$ . To do so we need to obtain a smooth extension of  $\text{grad}f$ , which we will denote  $\overline{\text{grad}f}$ . Well, the smooth extensions of  $\nabla\bar{f}(y)$  and  $\nabla h_i(y)$  can just be  $\nabla\bar{f}(y)$  and  $\nabla h_i(y)$  respectively.<sup>34</sup> The only question remaining is how to smoothly extend  $\lambda(y)$ . Well, from (5) we can view  $\lambda(y)$  as a smooth function on the open subset of  $\mathbb{R}^d$  consisting of all points  $y$  where  $Dh(y)$  has full rank. This is an open neighborhood of  $\mathcal{M}_0$ , so we have our smooth extension. In other words,

$$\overline{\text{grad}f(x)} = \nabla\bar{f}(y) - \sum_{i=1}^m \lambda_i(y) \nabla h_i(y), \quad (6)$$

where we have replaced  $\nabla\bar{f}(y)$ ,  $\lambda_i(y)$ ,  $\nabla h_i(y)$  with smooth extensions of themselves.

Then we can differentiate (6) normally (i.e., using Definition 2.1) to get

$$D(\overline{\text{grad}f(x)})(x)[u] = \nabla^2\bar{f}(y)[u] - \sum_{i=1}^m D\lambda_i(x)[u] \cdot \nabla h_i(x) - \sum_{i=1}^m \lambda_i(y) \nabla^2 h_i(y)[u].$$

Applying Corollary 3.6 yields

$$\begin{aligned} \text{Hess}f(y)[u] &= \text{Proj}_x \left( \nabla^2\bar{f}(y)[u] - \sum_{i=1}^m D\lambda_i(x)[u] \cdot \nabla h_i(x) - \sum_{i=1}^m \lambda_i(y) \nabla^2 h_i(y)[u] \right) \\ &= \text{Proj}_x \left( \nabla^2\bar{f}(y)[u] - \sum_{i=1}^m \lambda_i(y) \nabla^2 h_i(y)[u] \right), \end{aligned} \quad (7)$$

<sup>34</sup>Formally, we are extending the domain of  $\nabla\bar{f}(y)$  and  $\nabla h_i(y)$  as they appear in (4) from  $\mathcal{M}$  to  $\mathbb{R}^d$ . Although in other sense, the  $\nabla\bar{f}(y)$  and  $\nabla h_i(y)$  which appear in (4) are really restrictions of the “original”  $\nabla\bar{f}(y)$  and  $\nabla h_i(y)$  from  $\mathbb{R}^d$  to  $\mathcal{M}$ , so it is more like we are lifting this restriction.

where the second equality follows because

$$\text{Proj}_x \left( \sum_{i=1}^m D\lambda_i(x)[u] \cdot \nabla h_i(x) \right) = 0$$

since each  $\nabla h_i(x) \in N_y \mathcal{M}$  (see Example 2.11).

Note that (7) looks super similar to the second-order criticality condition in nonlinear programming (in other words, the Hessian of the Lagrangian)—more on this later.

**Example 3.9 (Riemannian Hessian in the Burer-Monteiro problem)** Applying (7), we obtain

$$\begin{aligned} \text{Hess}f(Y)[U] &= 2 \cdot \text{Proj}_Y \left( CU - \sum_{i=1}^m \lambda_i(Y) A_i U \right) \\ &= 2 \cdot \text{Proj}_Y \left( \begin{bmatrix} C - \sum_{i=1}^m \lambda_i(Y) A_i \\ U \end{bmatrix} \right). \end{aligned}$$

where  $\lambda(Y)$  is as in Example 2.30 and  $\text{Proj}_Y$  is as in Example 2.27.

As a sanity check, this is precisely what the authors of [BVB18] derive—see (2.7) in that paper.

### 3.1 Second-order geometry: material not covered

We went over much of the first half of Chapter 5 of [Bou20] in this section. Obviously some formality regarding connections was left out, but these details aren't really necessary for working with Riemannian submanifolds of Euclidean spaces due to Corollary 3.6 and Theorem 3.7. It is worth mentioning that Definition 3.1 (and the corresponding Definition 5.1 in [Bou20]) is not actually the standard definition of the connection. However, it is equivalent to the standard definition (Definition 5.19 in [Bou20]), as proven in Section 5.6 of [Bou20]. Boumal chooses to start with a nonstandard definition because it makes it easier to define the Riemannian Hessian.

The second half of Chapter 5 in [Bou20] contains other important material: differentiating vector fields along curves, acceleration, geodesics, second-order Taylor expansions on curves, second-order retractions, etc. We will see some of this later in these notes.

Finally, we followed Section 7.7 of [Bou20] in Example 3.8. It is worth noting that in Section 7.8 of [Bou20], an alternative derivation of the Riemannian Hessian for this particular problem is given. The latter approach mirrors the derivation of the second-order necessary condition in nonlinear programming (see page 333 in [LY16]), and as noted in [Bou20], can be extended more easily to higher-order derivatives.

### 3.2 Second-order geometry: extending beyond the embedded case

To be done.



## 4 Retractions, Taylor expansions, and optimality conditions

We now discuss the tool used to move along manifolds in a particular direction. Formally, given a point  $x \in \mathcal{M}$  and a tangent vector  $v \in T_x\mathcal{M}$ , starting at  $x$  and moving in the direction  $v$  corresponds to picking a curve  $c : I \subseteq \mathbb{R} \rightarrow \mathcal{M}$  such that  $c(0) = x$  and  $c'(0) = v$ , and then moving along  $c$ . There may be many such curves. A *retraction* picks a particular curve for each  $(x, v) \in T\mathcal{M}$  in a smooth way.

**Definition 4.1 (Retraction)** *A retraction<sup>35</sup> on  $\mathcal{M}$  is a smooth map  $R : T\mathcal{M} \rightarrow \mathcal{M}$  with the following properties. For each  $x \in \mathcal{M}$ , let  $R_x : T_x\mathcal{M} \rightarrow \mathcal{M}$  denote the restriction of  $R$  at  $x$ , so that  $R_x(v) = R(x, v)$ . Then*

1.  $R_x(0) = x$ , and
2.  $DR_x(0) : T_x\mathcal{M} \rightarrow T_x\mathcal{M}$  is the identity map:  $DR_x(0)[v] = v$ .

*Equivalently, each curve  $c(t) = R_x(tv)$  satisfies  $c(0) = x$  and  $c'(0) = v$ .*

First of all, note that  $DR_x(0)$  is indeed a map from  $T_x\mathcal{M}$  to  $T_x\mathcal{M}$  per Definition 2.18. In particular, the domain is  $T_x\mathcal{M}$  because the tangent space to  $T_x\mathcal{M}$  is  $T_x\mathcal{M}$ , due to the fact that  $T_x\mathcal{M}$  is a linear space.

Next, let us interpret the conditions of Definition 4.1. The first condition is easy to interpret, but what of  $DR_x(0)[v] = v$ ? Well, it is helpful to go back to Definition 2.18 and the text below it. Remember that  $DR_x(0)[v]$  intuitively tells us how  $R_x$  changes when we start with  $0 \in T_x\mathcal{M}$  as the input and push the input in the direction  $v$ . So the fact that  $DR_x(0)[v] = v$  says that if we push the input along  $v$ , the output (which starts at  $R_x(0)$ ) had better also be pushed along  $v$ . (Or more precisely, e.g., the curve  $R_x(tv) : \mathbb{R} \rightarrow \mathcal{M}$  better have velocity  $v$  when  $t = 0$ , which is exactly the text at the bottom of Definition 4.1.) In other words, the  $DR_x(0)[v] = v$  condition intuitively forces the retraction to actually “send us in the direction we specify.”

**Example 4.2 (Retractions on a sphere)** *The following are natural retractions on the sphere  $S^{n-1}$ :*

$$R_x(v) = \frac{x + v}{\|x + v\|_2} = \frac{x + v}{\sqrt{1 + \|v\|^2}}$$

and

$$R_x(v) = \cos(\|v\|)x + \frac{\sin(\|v\|)}{\|v\|}v. \tag{8}$$

*The former is simply projection. The latter has you move along the great circle traced out by  $v$ . See Example 3.40 in [Bou20] for more details.*

---

<sup>35</sup>Definition 3.41 in [Bou20]

**Example 4.3 (Retraction on  $\mathcal{M}_0$ )** A natural retraction on  $\mathcal{M}_0$  is the metric projection:

$$R_x(v) = \operatorname{argmin}_{y \in \mathcal{M}_0} \|x + v - y\|_2. \quad (9)$$

In other words, we return the point on the manifold closest to  $x + v$ . This is well-defined for small  $v$ , but  $R_x(v)$  may not be uniquely defined for larger  $v$ . It may also be difficult to compute this retraction.

**Example 4.4 (Retractions on  $\mathcal{M}_p$ )** As far as I know, the particularization of (9) to  $\mathcal{M}_p$  does not yield a nice expression in general.<sup>36</sup> However, for particular instantiations of  $\mathcal{M}_p$ , (9) might yield a nice expression (which I’m guessing is why manifold-based algorithms for optimization on  $\mathcal{M}_p$  seem to be limited to these special cases) or other nice retractions might exist. See [Bou16] and [JBAS10] for examples.

The following object is key to the general method of designing manifold optimization algorithms:

**Definition 4.5 (Pullback)** Let  $f : \mathcal{M} \rightarrow \mathbb{R}$  be smooth and let  $R : T\mathcal{M} \rightarrow \mathcal{M}$  be a retraction on  $\mathcal{M}$ . The composition  $f \circ R : T\mathcal{M} \rightarrow \mathbb{R}$  is called the pullback<sup>37</sup> of  $f$  to the tangent spaces. In particular,  $f \circ R_x : T_x\mathcal{M} \rightarrow \mathbb{R}$  is the pullback of  $f$  to the tangent space at  $x$ . Importantly, this is a smooth function on a linear space.

Since typical optimization algorithms on manifolds are retraction based, meaning they iterate

$$x_{k+1} = R_{x_k}(s_k)$$

for some step  $s_k$ , the change in the cost function value from one iterate ( $x_k$ ) to the next ( $x_{k+1}$ ) can be understood through the pullback  $f \circ R_{x_k}$ . In other words, we want to pick  $s_k \in T_{x_k}\mathcal{M}$  so as to minimize  $f \circ R_{x_k}$ . However,  $f \circ R_{x_k}$  may in general be very complicated, so we seek to approximate it using Taylor expansions. We then minimize the approximation instead of  $f \circ R_{x_k}$  itself.<sup>38</sup> This is the basis for both Riemannian gradient descent and Riemannian Newton’s method—more on this later. First we need to write out the relevant Taylor expansions, which will also help us derive optimality conditions.

**Theorem 4.6 (First-order Taylor expansion on curves)** Let  $\mathcal{M}$  be a Riemannian submanifold of  $\mathbb{R}^d$ ,<sup>39</sup> and let  $f : \mathcal{M} \rightarrow \mathbb{R}$  be smooth. Let  $c : I \rightarrow \mathcal{M}$  be a smooth curve on  $\mathcal{M}$  with  $c(0) = x$  and  $c'(0) = v$ , where  $I$  is an open interval of  $\mathbb{R}$  around  $t = 0$ . Then<sup>40</sup>

$$f(c(t)) = f(x) + t \langle \operatorname{grad} f(x), v \rangle + O(t^2).$$

The big  $O$  notation here is with respect to  $t$  going to 0, meaning that the residual goes to zero at least as fast as  $t^2$  as  $t \rightarrow 0$ .

---

<sup>36</sup>I haven’t actually tried to compute this though. But I’m guessing this is the case since I have not seen such a result.

<sup>37</sup>Section 4.1 in [Bou20].

<sup>38</sup>For reference, I took this discussion from Section 6.2 of [Bou20].

<sup>39</sup>See Section 2.1 for a reminder of what this means.

<sup>40</sup>Section 4.1 in [Bou20].

**Theorem 4.7 (First-order Taylor expansion of the pullback)** *With the same setup as in Theorem 4.6, suppose we additionally have a retraction  $R : T\mathcal{M} \rightarrow \mathcal{M}$ . Then<sup>41</sup>*

$$f(R_x(s)) = f(x) + t \langle \text{grad}f(x), s \rangle + O(\|s\|^2).$$

The proofs of Theorems 4.6 and 4.7 are very short and simple. Theorem 4.6 is proven by a simple application of the chain rule. Theorem 4.7 is proven by applying Theorem 4.6 to the curve  $c(t) = R_x(tv)$  and then introducing  $s = tv$ . See Section 4.1 in [Bou20] for the details.

**Theorem 4.8 (First-order necessary condition)** *Let  $f : \mathcal{M} \rightarrow \mathbb{R}$  be a smooth function on a Riemannian manifold. If  $x$  is a local minimizer of  $f$ , then  $\text{grad}f(x) = 0$ .<sup>42</sup>*

Once again, the proof is simple. If  $\text{grad}f(x)$  were not zero, then the first-order Taylor expansion reveals that  $f$  will decrease along any curve with velocity  $-\text{grad}f(x)$ .

**Definition 4.9 (Critical/stationary point)** *Given a smooth function  $f$  on a Riemannian manifold  $\mathcal{M}$ , we call  $x \in \mathcal{M}$  a critical point or stationary point<sup>43</sup> of  $f$  if  $\text{grad}f(x) = 0$ .*

Now we move into the second-order realm where things inevitably get more complicated:

**Theorem 4.10 (Second-order Taylor expansion on curves)** *With the same setup as in Theorem 4.6, we have<sup>44</sup>*

$$f(c(t)) = f(x) + t \langle \text{grad}f(x), v \rangle + \frac{t^2}{2} \langle \text{Hess}f(x)[v], v \rangle + \frac{t^2}{2} \langle \text{grad}f(x), c''(0) \rangle + O(t^3). \quad (10)$$

Clearly this is not as clean of an expression as the first-order Taylor expansion due to the annoying  $\frac{t^2}{2} \langle \text{grad}f(x), c''(0) \rangle$  term (which in particular depends on  $c$  beyond just the velocity of  $c$  at  $t = 0$ ). In fact, we have not yet even defined what  $c''(0)$  means exactly.  $c'(t)$  can be thought of as a function which maps each  $t$  in some open interval of  $\mathbb{R}$  containing 0 to a tangent vector at  $c(t)$ . Indeed,  $c'(t) \in T_{c(t)}\mathcal{M}$  by the literal definition of the tangent space (Definition 2.6). More generally,  $c'(t)$  is an example of a *smooth vector field on a curve* which is essentially produced by taking a curve  $c$  and composing it with a smooth vector field. Formally:

**Definition 4.11 (Smooth vector field on a curve)** *Given a smooth curve  $c : I \rightarrow \mathcal{M}$  ( $I$  is open and connected in  $\mathbb{R}$ ), the map  $Z : I \rightarrow T\mathcal{M}$  is a smooth vector field on  $c$  if  $Z(t)$  is in  $T_{c(t)}\mathcal{M}$  for all  $t \in I$ , and if it is smooth as a map from  $I$  (an open submanifold of  $\mathbb{R}$ ) to  $T\mathcal{M}$ .<sup>45</sup>*

Just as we needed to define a new kind of derivative, the connection, in Section 3 so that the Hessian spits out tangent vectors, we need to define a new kind of derivative to “properly” take

---

<sup>41</sup>Section 4.1 in [Bou20].

<sup>42</sup>Proposition 4.4 in [Bou20].

<sup>43</sup>Definition 4.5 in [Bou20].

<sup>44</sup>Section 5.9 in [Bou20].

<sup>45</sup>Definition 5.27 in [Bou20].

the derivative of a smooth vector field along a curve. It is called the *induced covariant derivative*, and it is defined in a similar way to how the Riemannian connection (Theorem 3.2) is defined: You basically show that there exists a unique operator which satisfies the “nice properties” which the derivative of a vector field along a curve “should satisfy.” See Theorem 5.28 in [Bou20] for the formal statement.

Just like for the connection, one of the most important properties of the induced covariant derivative is that the result of applying it to a smooth vector field along a curve should produce a smooth vector field along a curve, meaning in particular that the result is still assigning tangent vectors to each  $t \in I$  and not arbitrary vectors in the linear embedding space. Furthermore, just like for the connection (Theorem 3.5), in the cases we care about in these notes, the induced covariant derivative just corresponds to composing the regular Euclidean derivative with an orthogonal projection onto the tangent space. (See Proposition 5.29 in [Bou20] for a formal statement.) For this reason, *the  $c''(0)$  in Theorem 4.10 is just the regular Euclidean derivative of  $c'(0)$  composed with an orthogonal projection onto the tangent space.* Formally,

$$c''(t) := \text{Proj}_{c(t)} \left( \frac{d}{dt} c'(t) \right), \quad (11)$$

where  $\frac{d}{dt} c'(t)$  is the Euclidean derivative of  $c'(t)$ . In fact,  $c'(t)$  is formally called the *acceleration*<sup>46</sup> of the curve  $c$ . The classical or *extrinsic* acceleration of  $c$  is denoted in [Bou20] (see Section 5.8) as

$$\ddot{c}(t) = \frac{d^2}{dt^2} c(t) = \frac{d}{dt} c'(t)$$

Then (11) is equivalent to

$$c''(t) := \text{Proj}_{c(t)} (\ddot{c}(t)), \quad (12)$$

which says that for a Riemannian submanifold of  $\mathbb{R}^d$ , the acceleration of a curve  $c$  is the tangential part of its extrinsic acceleration in the embedding space.

This is all great, but we are still left with this ugly  $\frac{t^2}{2} \langle \text{grad} f(x), c''(0) \rangle$  term in Theorem 4.10. Well, fortunately it turns out that in most of the cases we care about, this term will get zeroed out. The first case where this will happen is when we are at a stationary point, since then  $\text{grad} f(x) = 0$ . This yields the following:

**Theorem 4.12 (Second-order necessary condition)** *Consider a smooth function  $f : \mathcal{M} \rightarrow \mathbb{R}$  defined on a Riemannian manifold  $\mathcal{M}$ . If  $x$  is a local minimizer of  $f$ , then  $\text{grad} f(x) = 0$  and  $\text{Hess} f(x) \succeq 0$ .<sup>47</sup>*

The proof is simple. We already know from Theorem 4.8 that  $x$  being a local minimum implies  $\text{grad} f(x) = 0$ , so (10) becomes

$$f(c(t)) = f(x) + \frac{t^2}{2} \langle \text{Hess} f(x)[v], v \rangle + O(t^3). \quad (13)$$

Then if Hess weren't positive semidefinite, (13) implies that  $f$  decreases along a curve whose velocity is an eigenvector of Hess with a negative corresponding eigenvalue.

<sup>46</sup>Definition 5.35 in [Bou20].

<sup>47</sup>Proposition 6.1 in [Bou20].

**Definition 4.13 (Second-order critical/stationary point)** Given a smooth function  $f$  on a Riemannian manifold  $\mathcal{M}$ , we call  $x \in \mathcal{M}$  a second-order critical point or a second-order stationary point<sup>48</sup> of  $f$  if  $\text{grad}f(x) = 0$  and  $\text{Hess}f(x) \succeq 0$ .

**Example 4.14 (Necessary conditions for optimization over  $\mathcal{M}_0$ )** We derive the criticality conditions for optimization over  $\mathcal{M}_0$ . Combining (3) and Theorem 4.8/Definition 4.9, we have that  $y \in \mathcal{M}_0$  is a first-order critical point if and only if

$$\text{grad}f(y) = \nabla \bar{f}(y) - \sum_{i=1}^m \lambda_i(y) \nabla h_i(y) = 0 \quad (14)$$

where

$$\lambda(y) := (Dh(y)^T)^+ [\nabla \bar{f}(y)].$$

Note that this is equivalent to the first-order nonlinear programming criticality condition with respect to the optimization of  $\bar{f}$  over  $\mathcal{M}_0$ . Indeed, we say that  $y$  is a first-order critical point with respect to the latter problem if  $Dh(y)$  has full rank (i.e., LICQ holds) and there exists unique multipliers  $\nu \in \mathbb{R}^m$  such that

$$\nabla \bar{f}(y) - \sum_{i=1}^m \nu_i \nabla h_i(y) = 0$$

The fact that  $\nabla h_1(y), \dots, \nabla h_m(y)$  are linearly independent (see Example 2.3) implies that  $\lambda(y) = \nu$  at a first-order critical point.

From (7) and Theorem 4.12/Definition 4.13, we have that  $y \in \mathcal{M}_0$  is a second-order critical point if and only if (14) holds and

$$\text{Proj}_x \circ \left( \nabla^2 \bar{f}(y) - \sum_{i=1}^m \lambda_i(y) \nabla^2 h_i(y) \right) \circ \text{Proj}_x \succeq 0. \quad (15)$$

(Note that we can precompose (7) with  $\text{Proj}_x$  to get the equivalent expression (15) since the input to the Riemannian Hessian lies in  $\mathbb{T}_x \mathcal{M}$  anyway by the definition of the Riemannian Hessian. We follow [Bou20]—see (7.79) in that reference—and make this change since (15) is “pleasantly symmetric.”)

It is immediately apparent that this is equivalent to the second-order criticality condition in nonlinear programming. See Exercise 7.8 in [Bou20] for the details.

**Example 4.15 (Necessary conditions for the Burer-Monteiro problem)** Combining Example 2.30 with Theorem 4.8/Definition 4.9, we have that  $Y \in \mathcal{M}_p$  is a first-order critical point if and only if

$$\text{grad}f(Y) = 2 \left[ \left( C - \sum_{i=1}^m \lambda_i(Y) A_i \right) Y \right] = 0 \quad (16)$$

---

<sup>48</sup>Definition 6.2 in [Bou20].

where

$$\lambda(Y) := (Dh(Y)^T)^+[2CY].$$

From applying (15) to the Burer-Monteiro problem, we get that  $Y \in \mathcal{M}_p$  is a second-order critical point if and only if (16) holds and

$$2 \left[ C - \sum_{i=1}^m \lambda_i(Y) A_i \right] \bullet UU^T \geq 0$$

for all  $U \in T_Y \mathcal{M}_p$ . (See Example 2.9 for an expression for  $T_Y \mathcal{M}_p$ .) Clearly these are equivalent to the nonlinear programming criticality conditions for the Burer-Monteiro problem. (See also Example 4.14 above.)

As a sanity check, these are indeed equivalent to the conditions given in [BVB18]—see Definition 2.3 in that reference.

It turns out that even if we aren't at a critical point, there is a way to remove the ugly  $\frac{t^2}{2} \langle \text{grad} f(x), c''(0) \rangle$  term. If we can't make  $\text{grad} f(x)$  equal to zero, we just need to make the acceleration  $c''(0)$  equal to zero! This leads to a special kind of curve:

**Definition 4.16 (Geodesic)** *On a Riemannian manifold  $\mathcal{M}$ , a geodesic<sup>49</sup> is a smooth curve  $c : I \rightarrow \mathcal{M}$  such that  $c''(t) = 0$  for all  $t \in I$ , where  $I$  is an open interval of  $\mathbb{R}$ .*

Due to (12), a curve  $c$  on a Riemannian submanifold  $\mathcal{M}$  is a geodesic if and only if its classical (or extrinsic) acceleration  $\ddot{c}$  is everywhere normal to  $\mathcal{M}$ .

We can also use this idea to define a special kind of retraction:

**Definition 4.17 (Second-order retraction)** *A second-order retraction<sup>50</sup>  $R$  on a Riemannian manifold  $\mathcal{M}$  is a retraction such that, for all  $x \in \mathcal{M}$  and all  $v \in T_x \mathcal{M}$ , the curve  $c(t) = R_x(tv)$  has zero acceleration at  $t = 0$ , that is,  $c''(0) = 0$ .*

What can we say about the existence of geodesics and second-order retractions? Well, it turns out that geodesics exist between any two sufficiently close points on an arbitrary Riemannian manifold.<sup>51</sup> If the Riemannian manifold is additionally complete<sup>52</sup> (which is satisfied in particular by compact manifolds<sup>53</sup> and finite-dimensional Euclidean spaces<sup>54</sup>), then any two points  $x, y$  in the same connected component can be connected by a geodesic.<sup>55</sup> It turns out that every Riemannian manifold admits a special second-order retraction called the exponential map.<sup>56</sup> In general, the

<sup>49</sup>Definition 5.37 in [Bou20].

<sup>50</sup>Definition 5.41 in [Bou20].

<sup>51</sup>See the note to the right of the proof of Proposition 6.3 in [Bou20].

<sup>52</sup>See Definitions 10.5 and 10.6 as well as Theorem 10.7 in [Bou20].

<sup>53</sup>Example 10.9 in [Bou20].

<sup>54</sup>Example 10.10 in [Bou20].

<sup>55</sup>Theorem 10.8 in [Bou20].

<sup>56</sup>Definition 10.13 in [Bou20].

exponential map may only be defined on an open subset of  $\text{TM}$ , but if the Riemannian manifold is additionally complete, then it is defined on all of  $\text{TM}$ .<sup>57</sup> As an example, (8) is the exponential map on  $S^{n-1}$ .<sup>58</sup> The curve  $c(t) = R_x(tv)$  defined by (8) is a geodesic.<sup>59</sup>

Using a second-order retraction, we can obtain a second-order version of Theorem 4.7:

**Theorem 4.18 (Second-order Taylor expansion of the pullback)** *Consider a Riemannian manifold  $\mathcal{M}$  of  $\mathbb{R}^d$  equipped with a second-order retraction  $R$ . Then for all points  $x \in \mathcal{M}$  we have<sup>60</sup>*

$$f(R_x(s)) = f(x) + \langle \text{grad}f(x), s \rangle + \frac{1}{2} \langle \text{Hess}f(x)[s], s \rangle + O(\|s\|^3).$$

Finally, we give a second-order sufficient condition:

**Theorem 4.19 (Second-order sufficient condition)** *Consider a smooth function  $f : \mathcal{M} \rightarrow \mathbb{R}$  defined on a Riemannian manifold  $\mathcal{M}$ . If  $\text{grad}f(x) = 0$  and  $\text{Hess}f(x) \succ 0$ , then  $x$  is a strict local minimizer of  $f$ .<sup>61</sup>*

The fact that there exists a geodesic between any two sufficiently close points on a Riemannian manifold is useful for proving this, since it lets us remove dependence on the acceleration  $c''$  from our Taylor expansions (in a neighborhood at least).

## 5 Quotient manifolds

Quotient manifolds are a powerful tool to gain further geometric insight into manifolds which can be decomposed into equivalence classes. Very informally speaking, a quotient manifold is the manifold formed by the equivalence classes of another manifold. The latter is known as the *total space*. As each element of a quotient manifold is an equivalence class, quotient manifolds are rather abstract objects. As a result, we seek to understand properties of and operators on quotient manifolds in terms of the corresponding properties of and operators on the total space. Indeed, this is much of what Chapter 9 of [Bou20] focuses on.

One might worry that the abstractness of quotient manifolds would hinder the development of effective algorithms on the quotient manifold. Fortunately this isn't the case. Running Riemannian gradient descent on the quotient manifold is actually equivalent to running Riemannian gradient descent on the total space. So in practice, we never actually run Riemannian gradient descent on the quotient manifold itself but instead run it on the simpler total space. Still, identifying Riemannian gradient descent as being equivalent on both can lead to better theoretical guarantees,<sup>62</sup> since while the Riemannian Hessian on the total space will never be positive definite, it may be positive definite

<sup>57</sup>See the text below Definition 10.13 in [Bou20].

<sup>58</sup>Example 10.18 in [Bou20].

<sup>59</sup>Example 5.36 in [Bou20].

<sup>60</sup>Proposition 5.43 in [Bou20].

<sup>61</sup>Proposition 6.3 in [Bou20].

<sup>62</sup>See Section 9.9 in [Bou20] for a discussion of this.

on the quotient manifold. Running Riemannian Newton’s method on the total space *isn’t* equivalent to running it on the quotient manifold, although the two return numerically close solutions.<sup>63</sup>

To get started, let  $\sim$  be an equivalence relation on a manifold  $\overline{\mathcal{M}}$  with equivalence classes

$$[x] = \{y \in \overline{\mathcal{M}} \mid x \sim y\},$$

and let

$$\mathcal{M} = \overline{\mathcal{M}} / \sim = \{[x] : x \in \overline{\mathcal{M}}\}$$

denote the resulting quotient set.  $\overline{\mathcal{M}}$  is known as the *total space*,<sup>64</sup> and the following is the key operator which links the total space  $\overline{\mathcal{M}}$  to the quotient set  $\mathcal{M}$ :

**Definition 5.1 (Canonical/natural projection)** *The canonical projection or natural projection  $\pi : \overline{\mathcal{M}} \rightarrow \mathcal{M}$  links the total space  $\overline{\mathcal{M}}$  to the quotient  $\mathcal{M}$  in the following way:*

$$\pi : x \mapsto \pi(x) = [x].$$

*In other words,  $\pi$  sends  $x \in \overline{\mathcal{M}}$  to its equivalence class.*

So far we have been careful to refer to  $\mathcal{M}$  as a quotient set, but we would like to know when it is also a quotient manifold. Unfortunately, fully understanding the definition of a quotient manifold requires one to also understand the definition of a general manifold (as opposed to just the special case of an embedded submanifold of a linear space—Definition 2.2—which we have covered in these notes). This is because even if the total space  $\overline{\mathcal{M}}$  is an embedded submanifold of  $\mathbb{R}^d$ , the quotient set  $\mathcal{M}$  may not be identifiable as an embedded submanifold of a linear space. (At least, I think this is the case due to the fact that it is really a set of equivalence classes.) Indeed, this is why the chapter on quotient manifolds (Chapter 9) comes after the chapter on general manifolds (Chapter 8) in [Bou20]. Still, as we will see, we can get very far in our practical understanding of quotient manifolds if we just accept the definition of a general manifold to be a black box.

**Definition 5.2 (Quotient manifold)** *Suppose we endow the quotient set  $\mathcal{M} = \overline{\mathcal{M}} / \sim$  with “smooth structure” so that it is a manifold<sup>65</sup> in its own right. If additionally the canonical projection  $\pi$  is smooth and its differential  $D\pi(x) : T_x\overline{\mathcal{M}} \rightarrow T_{[x]}\mathcal{M}$  has rank  $\dim \mathcal{M}$  (i.e., it is full rank) for all  $x \in \overline{\mathcal{M}}$ , then we say that  $\mathcal{M}$  is a quotient manifold.<sup>66</sup>*

$D\pi(x)$  is a very important operator in its own right: Just as  $\pi$  provides a connection between points on  $\overline{\mathcal{M}}$  and points on  $\mathcal{M}$ ,  $D\pi(x)$  provides a connection between tangent vectors on  $\overline{\mathcal{M}}$  and tangent vectors on  $\mathcal{M}$  (or more accurately, between tangent vectors in  $T_x\overline{\mathcal{M}}$  and tangent vectors in  $T_{[x]}\mathcal{M}$ ). Indeed, we will not actually formally cover in these notes what a tangent vector on

<sup>63</sup>See Section 9.12 in [Bou20].

<sup>64</sup>Section 9.1 in [Bou20]. We follow the notational convention of [Bou20] in that the total space is always denoted  $\overline{\mathcal{M}}$  and the quotient set is always denoted  $\mathcal{M}$ .

<sup>65</sup>Roughly, you can think of a general manifold as being a space which locally resembles Euclidean space. See Definition 8.20 in [Bou20] for a formal definition.

<sup>66</sup>Definition 9.1 in [Bou20].



$\mathcal{M}$  is, since as discussed earlier,  $\mathcal{M}$  may not be an embedded submanifold of a linear space even if  $\overline{\mathcal{M}}$  is. (Definition 2.6 is only for embedded submanifolds of a linear space, although the general definition of a tangent vector, Definition 8.29 in [Bou20], is a natural extension.) That said, we can still understand tangent vectors on  $\mathcal{M}$  using  $D\pi(x)$ . In fact, even with the definition of a general tangent vector, it is still typically more helpful to understand tangent vectors on  $\mathcal{M}$  through  $D\pi(x)$ !

Typically we do not actually use Definition 5.2 to check that a quotient set is a quotient manifold; instead, we use properties of the equivalence relation  $\sim$ . In fact, there exists a characterization of which equivalence relations give rise to quotient manifolds,<sup>67</sup> but it is also unwieldy. Fortunately, there exists a class of equivalence relations defined through *group actions* on manifolds that are simple to identify and ubiquitous in practice. We will not formally go over this sufficient condition for  $\mathcal{M}$  to be a quotient manifold,<sup>68</sup> but we will provide intuition for it through an example.

First, we will informally go over a few definitions. A *Lie group* is a set that is both a group (in the abstract algebra sense) and a manifold. Examples include  $O(n)$  (the orthogonal group or the set of  $n \times n$  orthogonal matrices),  $SO(n)$  (the rotation group or the set of  $n \times n$  orthogonal matrices with determinant 1), and  $GL(n)$  (the general linear group or the set of invertible  $n \times n$  matrices.) Elements of a group can be used to transform points on a manifold—formally, these are called left group actions and right group actions. For example, under the right conditions, right-multiplying by an orthogonal matrix is a right-group action and left-multiplying by an orthogonal matrix is a left-group action. These actions naturally induce equivalence relations on the set that is getting acted on.<sup>69</sup> (See Section 9.2 in [Bou20] for a formal treatment of the contents of this paragraph.)

If a Lie group  $\mathcal{G}$  acts on a smooth manifold  $\overline{\mathcal{M}}$  in a way that satisfies certain properties, then  $\overline{\mathcal{M}}/\mathcal{G}$  is a quotient manifold. We will give the formal statement next, although we will not formally define what all of the words in the theorem mean. That said, we will get some intuition for them in Example 5.5.

**Theorem 5.3 (Quotient manifold through a group action)** *If a Lie group  $\mathcal{G}$  acts smoothly, freely, and properly on a smooth manifold  $\overline{\mathcal{M}}$ , then the quotient space  $\overline{\mathcal{M}}/\mathcal{G}$  is a quotient manifold of dimension  $\dim \overline{\mathcal{M}} - \dim \mathcal{G}$ .<sup>70</sup> For manifolds that arise this way, the quotient space  $\overline{\mathcal{M}}/\mathcal{G}$  is also known as the orbit space.<sup>71</sup>*

We will gain some intuition behind this theorem in Example 5.5, but first, a theorem that will be useful in the upcoming example:

**Theorem 5.4 (Open subsets of embedded submanifolds are embedded submanifolds)** *Let  $\mathcal{M}$  be an embedded submanifold of  $\mathbb{R}^d$ . Any open subset of  $\mathcal{M}$  is also an embedded submanifold of  $\mathbb{R}^d$ , with the same dimension and tangent spaces as  $\mathcal{M}$ .<sup>72</sup>*

**Example 5.5 ( $\mathcal{M}_p^{full}/O(p)$  is a quotient manifold)** *Let  $\mathcal{M}_p^{full}$  denote the open subset of  $\mathcal{M}_p$  that contains its rank  $p$  elements. By Theorem 5.4, we have that  $\mathcal{M}_p^{full}$  is an embedded submanifold of  $\mathbb{R}^{np}$  with the same dimension and tangent spaces as  $\mathcal{M}_p$ .*

<sup>67</sup>Proposition 3.4.2 in [AMS07].

<sup>68</sup>See Section 9.2 in [Bou20] for a formal treatment.

<sup>69</sup>Definition 9.12 in [Bou20].

<sup>70</sup>Theorem 8.17 in [Bou20].

<sup>71</sup>Definition 9.12 in [Bou20].

<sup>72</sup>Proposition 3.17 in [Bou20].

We claim that  $\mathcal{M}_p^{full}/\mathcal{O}(p)$  is a quotient manifold, where the equivalence classes are formed by  $\mathcal{O}(p)$  acting on  $\mathcal{M}_p^{full}$  through right multiplication. In other words,

$$Y \sim V \iff Y = VQ \text{ for some } Q \in \mathcal{O}(p). \quad (17)$$

To prove that  $\mathcal{M}_p^{full}/\mathcal{O}(p)$  is indeed a quotient manifold, it is sufficient to check that  $\mathcal{O}(p)$  acts smoothly, freely, and properly on  $\mathcal{M}_p^{full}$ . The group action is smooth since the map from  $\mathcal{M}_p^{full} \times \mathcal{O}(p)$  to  $\mathcal{M}_p^{full}$  defined by  $(Y, Q) \mapsto YQ$  is smooth. The group action is free because only the identity element of the group fixes any given  $Y$ . In other words,

$$YQ = Y \text{ for some } Y \in \mathcal{M}_p^{full} \Rightarrow Q = I_p. \quad (18)$$

The group action is proper because  $\mathcal{O}(p)$  is compact, and every smooth action by a compact Lie group is proper.<sup>73</sup> We conclude that  $\mathcal{M}_p^{full}/\mathcal{O}(p)$  is indeed a quotient manifold. As a sanity check,  $\mathcal{M}_p^{full}$  is stated to be a quotient manifold in [JBAS10] (see Section 4) and [WW20] (see page 18).

We can also obtain the dimension of  $\mathcal{M}_p^{full}/\mathcal{O}(p)$  using Theorem 5.3:

$$\begin{aligned} \dim(\mathcal{M}_p^{full}/\mathcal{O}(p)) &= \dim \mathcal{M}_p^{full} - \dim \mathcal{O}(p) \\ &= np - m - \frac{p(p-1)}{2}, \end{aligned}$$

where we used the fact  $\dim \mathcal{M}_p^{full} = \dim \mathcal{M}_p$  due to Theorem 5.4, and the result from Example 2.15. The expression for  $\dim \mathcal{O}(p)$  can be found in Section 7.4 in [Bou20]. To confirm our expression, see page 18 of [WW20].

A natural question is whether  $\mathcal{M}_p/\mathcal{O}(p)$  (with equivalence classes defined as in (17)) is also a quotient manifold. It turns out that if  $\mathcal{M}_p$  contains rank-deficient matrices, then  $\mathcal{M}_p/\mathcal{O}(p)$  is not a quotient manifold, as noted in the footnote at the bottom of page 18 in [WW20].<sup>74</sup>

For the next definition, it will be useful to be able to talk about the set containing all points in the total space  $\overline{\mathcal{M}}$  which map into the same equivalence class on the quotient manifold  $\mathcal{M}$ . The notation we will use for this is  $\pi^{-1}(\pi(x))$ . ( $\pi$  is defined in Definition 5.1.) In other words,

$$\pi^{-1}(\pi(x)) = \{y \in \overline{\mathcal{M}} \mid x \sim y\}.$$

**Definition 5.6 (Fiber, orbit)** Let  $\mathcal{M} = \overline{\mathcal{M}}/\sim$  be a quotient manifold. For any  $x \in \overline{\mathcal{M}}$ , the equivalence class  $\mathcal{F} = \pi^{-1}(\pi(x))$ , also called a fiber, is closed in  $\overline{\mathcal{M}}$  and it is an embedded submanifold of  $\overline{\mathcal{M}}$ . Its tangent spaces are given by

$$T_y \mathcal{F} = \ker D\pi(y) \subseteq T_y \overline{\mathcal{M}}.$$

<sup>73</sup>Proposition 9.16 in [Bou20].

<sup>74</sup>They prove that  $\mathcal{M}_p/\mathcal{O}(p)$  is not a quotient manifold by showing that it violates Condition (i) of Proposition 3.4.2 in [AMS07]. (Proposition 3.4.2 is a complete characterization of which equivalence relations lead to quotient manifolds.) I think another way to prove it would be to show that the fibers (Definition 5.6) that would arise were  $\mathcal{M}_p/\mathcal{O}(p)$  a quotient manifold could have different dimensions, which is a contradiction. This is discussed on page 203 of [Bou20]; see also Exercise 9.7 in that reference.

In particular,  $\dim \mathcal{F} = \dim \overline{\mathcal{M}} - \dim \mathcal{M}$ .<sup>75</sup>

If  $\mathcal{M}$  arises as the action of a Lie group  $\mathcal{G}$  on  $\overline{\mathcal{M}}$  as in Theorem 5.3, then the fiber  $\pi^{-1}(\pi(x))$  is also known as the orbit of  $x \in \overline{\mathcal{M}}$  under the action associated with  $\mathcal{G}$ .<sup>76</sup>

This definition (which is really part theorem) says that the equivalence classes on the total space partition the total space into many embedded submanifolds. Furthermore, we begin to see through Definition 5.6 the importance of the operator  $D\pi(y)$ , which was discussed earlier.

**Example 5.7 (Fibers of  $\mathcal{M}_p^{full}/O(p)$ )** The fiber containing  $Y \in \mathcal{M}_p^{full}$  is precisely

$$\mathcal{F} = \{Z \in \mathcal{M}_p^{full} \mid Y \sim Z\} = \{YQ \mid Q \in O(p)\}.$$

In this case, it is not hard to calculate  $T_V \mathcal{F}$  directly for some  $V \in \mathcal{F}$  using Definition 2.6. All tangent vectors in  $T_V \mathcal{F}$  take the form  $c'(0)$  for some smooth curve  $c : I \subseteq \mathbb{R} \rightarrow \mathcal{F}$  with  $c(0) = V$ . Any such curve is necessarily of the form  $c(t) = VQ(t)$  where  $Q : I \rightarrow O(p)$  is a smooth curve on the manifold  $O(p)$  with  $Q(0) = I_p$ . Then, all tangent in  $T_V \mathcal{F}$  are of the form  $VQ'(0)$ . The tangent space to  $O(p)$  at  $Q(0) = I_p$  turns out to be the set of  $p \times p$  skew-symmetric matrices,<sup>77</sup> so we have that

$$T_V \mathcal{F} = \{VS \mid S \in \mathbb{R}^{p \times p}, S + S^T = 0\}.$$

It is not hard to check that  $T_V \mathcal{F}$  is indeed a subset of  $T_V \mathcal{M}_p^{full}$ , as expected. (See Example 2.9, and recall that  $T_V \mathcal{M}_p^{full} = T_V \mathcal{M}_p$  due to Theorem 5.4.)

For a derivation of  $T_V \mathcal{F}$  involving  $\ker D\pi(V)$ , see Example 9.4 in [Bou20]. (The above also basically follows Example 9.4, although we are working with a slightly different, albeit functionally the same, manifold.)

As a sanity check, this is actually the expression obtained in [WW20]—see the bottom of page 6 in that paper. (Note that they use the term “orbit” instead of “fiber.”)

Finally, we calculate  $\dim \mathcal{F}$  using our expression for  $T_V \mathcal{F}$  and the fact that by definition, the dimension of a manifold is just the dimension of its tangent space (Definition 2.13). This immediately implies that  $\dim \mathcal{F} = \binom{p}{2} = \frac{p(p-1)}{2}$ . As a quick check, this is indeed equal to

$$\dim \mathcal{M}_p^{full} - \dim(\mathcal{M}_p^{full}/O(p)) = (np - m) - \left( np - m - \frac{p(p-1)}{2} \right).$$

We would now like to obtain a correspondence between tangent vectors of the total space  $\overline{\mathcal{M}}$  and tangent vectors of the quotient manifold  $\mathcal{M} = \overline{\mathcal{M}}/\sim$ . The appropriate tool to do this is  $D\pi(x) : T_x \overline{\mathcal{M}} \rightarrow T_{[x]} \mathcal{M}$ , which is surjective due to Definition 5.2. However, it is not one-to-one, so to fix this we need to restrict its domain. Definition 5.6 provides a very natural way to do this—we can partition the domain of  $D\pi(x)$ , aka  $T_x \overline{\mathcal{M}}$ , into those vectors in  $T_x \mathcal{F} = \ker D\pi(x)$  and those vectors in its orthogonal complement,  $(T_x \mathcal{F})^\perp = (\ker D\pi(x))^\perp$ . Then, the restriction  $D\pi(x)$  to  $(T_x \mathcal{F})^\perp = (\ker D\pi(x))^\perp$  forms a bijection between  $(T_x \mathcal{F})^\perp = (\ker D\pi(x))^\perp$  and  $T_{[x]} \mathcal{M}$ . We formalize this below with proper terminology:

<sup>75</sup>Proposition 9.3 in [Bou20].

<sup>76</sup>Definition 9.12 in [Bou20].

<sup>77</sup>Section 7.4 in [Bou20].

**Definition 5.8 (Vertical space, horizontal space)** For a quotient manifold  $\mathcal{M} = \overline{\mathcal{M}}/\sim$ , the vertical space<sup>78</sup> at  $x \in \overline{\mathcal{M}}$  is the subspace

$$\mathbf{V}_x = \mathbf{T}_x \mathcal{F} = \ker D\pi(x)$$

where  $\mathcal{F} = \{y \in \overline{\mathcal{M}} \mid y \sim x\}$  is the fiber of  $x$ . If  $\overline{\mathcal{M}}$  is Riemannian submanifold of  $\mathbb{R}^d$  (so that we have an inner product), we call the orthogonal complement of  $\mathbf{V}_x$  the horizontal space at  $x$ :

$$\mathbf{H}_x = (\mathbf{V}_x)^\perp = \{u \in \mathbf{T}_x \overline{\mathcal{M}} \mid \langle u, v \rangle = 0 \text{ for all } v \in \mathbf{V}_x\}.$$

Intuitively, the vertical directions (aka  $\mathbf{V}_x$ ) are all of the “uninteresting” directions of  $\mathbf{T}_x \overline{\mathcal{M}}$ .

Recall from the text above Definition 5.8 that the restriction of  $D\pi(x) : \mathbf{T}_x \overline{\mathcal{M}} \rightarrow \mathbf{T}_{[x]} \mathcal{M}$  to  $\mathbf{H}_x = (\mathbf{T}_x \mathcal{F})^\perp = (\ker D\pi(x))^\perp$ , aka  $D\pi(x)|_{\mathbf{H}_x}$ , is a bijection. This bijection allows us to use “concrete” horizontal vectors in  $\mathbf{T}_x \overline{\mathcal{M}}$  to represent “abstract” vectors in  $\mathbf{T}_{[x]} \mathcal{M}$ .

**Definition 5.9 (Horizontal lift)** Consider a point  $x \in \overline{\mathcal{M}}$  and a tangent vector  $\xi \in \mathbf{T}_{[x]} \mathcal{M}$ . The horizontal lift<sup>79</sup> of  $\xi$  at  $x$  is the unique horizontal vector  $u \in \mathbf{H}_x$  such that  $D\pi(x)[u] = \xi$ . We write

$$u = (D\pi(x)|_{\mathbf{H}_x})^{-1} [\xi] = \text{lift}_x(\xi).$$

Horizontal lifts are critical for connecting objects on the quotient manifold to objects on the total space (e.g., retractions, vector fields, connections, etc.) We only include two examples in these notes:

**Theorem 5.10 (Riemannian gradient result for quotient manifolds)** The Riemannian gradient of  $f$  on a Riemannian quotient manifold is related to the Riemannian gradient of the lifted function  $\overline{f} = f \circ \pi$  on the total space via<sup>80</sup>

$$\text{lift}_x(\text{grad} f([x])) = \text{grad} \overline{f}(x)$$

for all  $x \in \overline{\mathcal{M}}$ .

We haven’t actually talked about objective functions at all in this section, so in particular we haven’t talked about what it means to lift a function. It is super intuitive though—precomposing a function  $f$  defined on the quotient manifold with the canonical projection  $\pi$  instantly gives you an “equivalent” function on the total space. And in fact, the only objective functions on the total space that are “interesting” in the first place (with respect to the quotient geometry at least) are those that are the lift of some function on the quotient manifold. If this isn’t the case, then your objective function on the total space has different values among elements in the same equivalence class, at which point the quotient geometry doesn’t help much.

Note that Theorem 5.10 implies that the Riemannian gradient on the total space, aka  $\text{grad} \overline{f}(x)$ , is always a horizontal vector (because the output of a lift is always a horizontal vector). This makes

<sup>78</sup>Definition 9.23 in [Bou20].

<sup>79</sup>Definition 9.24 [Bou20].

<sup>80</sup>Proposition 9.38 in [Bou20].

sense intuitively—the Riemannian gradient should be orthogonal to the “flat” directions, aka the vertical space  $V_x$ . (After all, the Euclidean gradient is always orthogonal to the level curves—we should expect the same of the Riemannian gradient.)

We now give a similar result for the Hessian without much commentary:

**Theorem 5.11 (Riemannian Hessian result for quotient manifolds)** *The Riemannian Hessian of  $f$  on a Riemannian quotient manifold is related to the Riemannian Hessian of the lifted function  $\bar{f} = f \circ \pi$  on the total space as*

$$\text{lift}_x(\text{Hess}f([x])[\xi]) = \text{Proj}_x^{\text{H}}(\text{Hess}\bar{f}(x)[u])$$

for all  $x \in \bar{\mathcal{M}}$  and  $\xi \in \mathbb{T}_{[x]}\mathcal{M}$ , with  $u = \text{lift}_x(\xi)$ . Here  $\text{Proj}_x^{\text{H}}$  is the orthogonal projector onto  $\mathbb{H}_x$ .

The following theorems are useful for analyzing critical points in the total space:

**Theorem 5.12 (Critical points result for quotient manifolds)** *Let  $f : \mathcal{M} \rightarrow \mathbb{R}$  be smooth on a Riemannian quotient manifold  $\mathcal{M} = \bar{\mathcal{M}}/\sim$  with canonical projection  $\pi$ , and denote the lifted function  $\bar{f} = f \circ \pi$ . From Theorem 5.10, it is clear that  $x \in \bar{\mathcal{M}}$  is a first-order critical point for  $\bar{f}$  if and only if  $[x]$  is a first-order critical point for  $f$ . Furthermore, it can be shown that  $x$  is a second-order critical point for  $\bar{f}$  if and only if  $[x]$  is a second-order critical point for  $f$ .<sup>81</sup>*

We can say more about the Riemannian Hessian at a first-order critical point:

**Theorem 5.13 (Riemannian Hessian at a critical point result for quotient manifolds)** *With the same setup as in Theorem 5.12, let  $x \in \bar{\mathcal{M}}$  be a first-order critical point. Then the eigenvalues of  $\text{Hess}\bar{f}(x)$  are exactly the eigenvalues of  $\text{Hess}f([x])$  together with a set of  $\dim V_x = \dim \bar{\mathcal{M}} - \dim \mathcal{M}$  eigenvalues equal to zero. In particular, the vertical space  $V_x$  is included in the kernel of  $\text{Hess}\bar{f}(x)$ . For every eigenvector of  $\text{Hess}f([x])$ , there is a corresponding eigenvector of  $\text{Hess}\bar{f}(x)$  which lies in  $\mathbb{H}_x$  and has the same eigenvalue.<sup>82</sup> (Indeed, the latter is just the horizontal lift of the former.)*

The above implies that if  $\dim \mathcal{M} < \dim \bar{\mathcal{M}}$ , the cost function  $\bar{f}$  on the total space  $\bar{\mathcal{M}}$  cannot admit second-order critical points where the Hessian is positive definite. Fortunately this isn’t an issue because as explained in Section 9.9 in [Bou20], running Riemannian gradient descent on the total space is equivalent to running it on the Hessian. Thus, if the Hessian is positive definite on the quotient manifold, we will get the same benefits as if it were positive definite on the total space.

It is also possible to extend Theorem 4.19 to this scenario:

**Theorem 5.14 (Second-order sufficient condition for quotient manifolds)** *With the same setup as in Theorem 5.12, let  $x \in \bar{\mathcal{M}}$  be such that  $\text{grad}\bar{f}(x) = 0$  (which is true if and only if  $\text{grad}f([x]) = 0$  due to Theorem 5.10) and  $\text{Hess}f([x]) \succ 0$ . Then  $x \in \bar{\mathcal{M}}$  is a local minimum.<sup>83</sup>*

<sup>81</sup>Exercise 9.45 in [Bou20].

<sup>82</sup>Exercise 9.45 and Lemma 9.40 in [Bou20]. Although part of this theorem I derived myself when solving Exercise 9.45.

<sup>83</sup>This result is not given in [Bou20], but I was able to prove it myself. Note that even if  $x$  is a local minimum, it will never be a strict local minimum due to the “flat” directions corresponding to the vertical space  $V_x$ .

Note that due to Theorem 5.13, we do not need to check that  $\text{Hess}f([x]) \succ 0$  directly. It is instead sufficient to show that the only eigenvectors of  $\text{Hess}\bar{f}(x)$  with 0 as their eigenvalue are those  $\dim V_x = \dim \bar{\mathcal{M}} - \dim \mathcal{M}$  of them which lie in  $V_x$ .

For example, Theorem 5.14 is used heavily in [WW20] to show that certain second-order critical points are additionally local minima. See Definition 3 and Remark 1 in that paper.

## 5.1 Quotient manifolds: material not covered

The material in this section is heavily based off of Chapter 9 of [Bou20]. That said, Chapter 9 is very long and we skipped a ton of material. Nearly all of the material skipped has to do with connecting objects on the total space to the corresponding objects on the quotient manifold (e.g., smooth maps, vector fields, retractions, connections, etc.) Section 9.13 is worth reading in particular since they talk about the specific scenario where the total space  $\bar{\mathcal{M}}$  is an embedded submanifold of  $\mathbb{R}^d$ .

## 6 Riemannian gradient descent

To be done. For now, refer to Chapter 4 in [Bou20].

## 7 Riemannian second-order methods

To be done. For now, refer to Chapter 6 in [Bou20].

## A Examples of manifolds

To be done. For now, refer to Chapter 7 in [Bou20], which is full of examples.

## B Regarding product manifolds

**Definition B.1 (Product manifold)** *Let  $\mathcal{M}, \mathcal{M}'$  be embedded submanifolds of  $\mathbb{R}^d, \mathbb{R}^{d'}$  respectively. Then  $\mathcal{M} \times \mathcal{M}'$  is an embedded submanifold of  $\mathbb{R}^d \times \mathbb{R}^{d'}$  of dimension  $\dim \mathcal{M} + \dim \mathcal{M}'$  such that*

$$T_{(x,x')}(\mathcal{M} \times \mathcal{M}') = T_x \mathcal{M} \times T_{x'} \mathcal{M}'.$$

Basically, product manifolds behave exactly as you would expect them to: “X” for a product manifold is typically just the product or concatenation of “X” over each individual manifold which makes up the product manifold.

More details/examples to be added later.

## References

- [AMS07] P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, USA, 2007.
- [Bou16] Nicolas Boumal. A riemannian low-rank method for optimization over semidefinite matrices with block-diagonal constraints, 2016.
- [Bou20] Nicolas Boumal. An introduction to optimization on smooth manifolds. Available online, Nov 2020.
- [BVB18] Nicolas Boumal, Vladislav Voroninski, and Afonso Bandeira. Deterministic guarantees for burer-monteiro factorizations of smooth semidefinite programs. *Communications on Pure and Applied Mathematics*, 73, 04 2018.
- [JBAS10] M. Journée, F. Bach, P.-A. Absil, and R. Sepulchre. Low-rank optimization on the cone of positive semidefinite matrices. *SIAM Journal on Optimization*, 20(5):2327–2351, 2010.
- [LY16] David G. Luenberger and Yinyu Ye. *Linear and nonlinear programming, fourth edition*, volume 228. New York, NY: Springer, 2016.
- [WW20] Irène Waldspurger and Alden Waters. Rank optimality for the Burer-Monteiro factorization. *SIAM J. Optim.*, 30(3):2577–2602, 2020.