



DEPARTMENT OF MATHEMATICAL SCIENCES

TMA4500 - INDUSTRIAL MATHEMATICS, SPECIALIZATION
PROJECT

Riemannian Optimization using second order information on the Symplectic Stiefel manifold

Author:
Ole Gunnar Røsholt Hovland

What date should I use?

Table of Contents

List of Figures	i
List of Tables	i
1 Introduction	1
2 Introductory Theory	1
2.1 Foundational definitions	1
2.2 The Symplectic group	3
2.3 The Symplectic Stiefel manifold	3
3 Right-Invariant Framework on the Symplectic Stiefel Manifold	4
3.1 Right-invariant metric	4
3.2 Geodesics	5
3.3 Riemannian gradient of the Symplectic Stiefel manifold	5
3.4 Riemannian Hessian of the Symplectic Stiefel manifold	6
3.5 Moving from Theory to Application	7
4 Algorithms	7
4.1 Riemannian gradient descent	7
4.2 Riemannian trust-region method	8
4.3 Implementation	9
5 Numerical Experiments	9
5.1 Nearest symplectic matrix problem	10
6 Conclusion	11
Bibliography	13

List of Figures

List of Tables

1	Nearest symplectic matrix problem solved by CG, TR-1 and TR-2. The table summarizes time to converge for all the algorithms, with $n = 100$ for $k = 5, 10, 20$. . .	12
---	---	----

1 Introduction

2 Introductory Theory

2.1 Foundational definitions

This section is designed to be a reference work to set notation, and to ensure that the reader has the necessary background to understand the optimization algorithms we will be studying.

Note that even though the elements of the symplectic group and symplectic stiefel manifold are matrices, the word "point" will be used to refer to a specific matrix, as it is a point on the matrix manifold.

Definition 1 (General Linear group). *The real General Linear group is defined as the set of all invertible matrices in $\mathbb{R}^{n \times n}$, denoted by $\text{GL}(n)$. [4, Example 9.11]*

Definition 2 (Orthogonal group). *The real Orthogonal group is defined as the set of all orthogonal matrices in $\mathbb{R}^{n \times n}$, denoted by $\text{O}(n)$. [5, p. 3]*

Definition 3 (Quotient manifold). *We define the definition of quotient manifold as in [1, p. 27]. Let \mathcal{M} be a manifold equipped with the operation \sim called the equivalence relation. The equivalence relation has the following properties:*

1. (reflexive) $p \sim p$ for all $p \in \mathcal{M}$,
2. (symmetric) $p \sim q$ if and only if $q \sim p$ for all $q, p \in \mathcal{M}$, and
3. (transitive) given $p \sim q$ and $q \sim r$ this implies that $p \sim r$ for all $p, q, r \in \mathcal{M}$.

Given the set $[p] := \{q \in \mathcal{M} : q \sim p\}$ called the equivalence class of all points equivalent to p , the set

$$\mathcal{M}/\sim := \{[p] \mid p \in \mathcal{M}\}$$

is called the quotient of \mathcal{M} by \sim . It is the set of all equivalence classes of \sim in \mathcal{M} . The mapping $\pi : \mathcal{M} \rightarrow \mathcal{M}/\sim$ called the natural- or canonical projection, defined by $p \mapsto [p]$.

Definition 4 (Tangent Space). *Following [4, Def. 8.33], for a point p on a smooth manifold \mathcal{M} , denote the set of smooth curves [4, Def. 8.5] passing through p at $t = 0$ as C_p . This means that $\alpha(0) = p$ for all $\alpha \in C_p$. For $\alpha, \beta \in C_p$ we say that they are equivalent if*

$$(\phi \circ \alpha)'(0) = (\phi \circ \beta)'(0),$$

meaning their derivatives match in a coordinate chart (defined as in [9, p. 4]) if their derivatives in the coordinate chart at zero are equal. Denote this equivalence relation as $\alpha \sim \beta$. It has analogous properties to the equivalence relation in Definition 3. The equivalence class is defined as $[\alpha] = \{\beta \in C_p \mid \alpha \sim \beta\}$. Every equivalence class is called a tangent vector to \mathcal{M} at p . The tangent space at p is the quotient set

$$T_p\mathcal{M} = C_p/\sim = \{[\alpha] \mid \alpha \in C_p\}.$$

Definition 5 (Riemannian manifold). *As defined in [3, def 2.6, p. 179]: a smooth manifold \mathcal{M} , as defined in [9, p. 13], is a Riemannian manifold if we can define a field of symmetric, positive definite, bilinear forms $g(\cdot, \cdot)$, called the Riemannian metric. By field we mean that g_p is defined on the tangent space $T_p\mathcal{M}$ at each point $p \in \mathcal{M}$ [3, def 2.1, p. 178]. We will assume that g is smooth, meaning that it is of class \mathcal{C}^∞ .*

Definition 6 (Vector field on Riemannian manifold). *Following Appendix A of [8], a smooth vector field $\mathcal{X} : \mathcal{M} \rightarrow T\mathcal{M}$, $p \mapsto \mathcal{X}(p) \in T_p\mathcal{M}$ on a Riemannian manifold \mathcal{M} can be expressed through local coordinates as*

$$\mathcal{X}(p) = \sum_{i=1}^n \alpha_i \partial_i =: \boldsymbol{\alpha}^T \boldsymbol{\partial},$$

where $\boldsymbol{\alpha} \in \mathbb{R}^n$, and $\boldsymbol{\partial}$ is the canonical basis of $T_p\mathcal{M}$.

Definition 7 (Horizontal & Vertical Space). Using Definition 3, given a Riemannian manifold $\overline{\mathcal{M}}$ with Riemannian metric \overline{g} , denote a quotient manifold of $\overline{\mathcal{M}}$ as $\mathcal{M} = \overline{\mathcal{M}} / \sim$. Following the definitions in Absil et al. [1, p. 43], for a point $p \in \mathcal{M}$, the equivalence class $[p] = \pi^{-1}(p)$ induces an embedded submanifold of $\overline{\mathcal{M}}$ (see Definition 3), hence it admits a tangent space,

$$\mathcal{V}_{\overline{p}} = T_{\overline{p}}(\pi^{-1}(p)) = \ker(D\pi(p))$$

[8, p. 4] named the vertical space at \overline{p} . Note that Canonically chosen as the orthogonal complement of $\mathcal{V}_{\overline{p}}$ in $T_{\overline{p}}\overline{\mathcal{M}}$, the horizontal space [1, p. 48] is defined as

$$\mathcal{H}_{\overline{p}} := \mathcal{V}_{\overline{p}}^\perp = \{Y_{\overline{p}} \in T_{\overline{p}}\overline{\mathcal{M}} \mid \overline{g}(Y_{\overline{p}}, Z_{\overline{p}}) = 0 \quad \forall \quad Z_{\overline{p}} \in \mathcal{V}_{\overline{p}}\}.$$

The horizontal lift at $\overline{p} \in \pi^{-1}(p)$ of a tangent vector $X_p \in T_p\mathcal{M}$ is the unique tangent vector $X_{\overline{p}} \in \mathcal{H}_{\overline{p}}$ that satisfies $D\pi(\overline{p})[X_{\overline{p}}] = X_p$. Note that given the horizontal space on $\overline{\mathcal{M}}$, $\mathcal{H}_{\overline{p}} \oplus \mathcal{V}_{\overline{p}} = T_{\overline{p}}\overline{\mathcal{M}}$, where \oplus denotes the Whitney sum.

Definition 8 (Riemannian connection). The Riemannian connection, also known as the Levi-Civita connection, is the unique affine connection which is torsion free, and metric compatible [11, Def. 6.4]. In Appendix A of [8], denoting $\mathfrak{X}(\mathcal{M})$ as the space of smooth vector fields on \mathcal{M} , it is defined as the unique \mathbb{R} -bilinear smooth map on \mathcal{M} with Riemannian metric $\langle \cdot, \cdot \rangle_p$

$$\nabla: \mathfrak{X}(\mathcal{M}) \times \mathfrak{X}(\mathcal{M}) \rightarrow \mathfrak{X}(\mathcal{M}), \quad (X, Y) \mapsto \nabla_X Y,$$

such that the following properties hold. Given $X, Y, Z \in \mathfrak{X}(\mathcal{M})$, and $f \in \mathcal{C}^\infty(\mathcal{M})$, $\nabla_X Y$ has the following properties:

1. (first argument linearity) $\nabla_{fX} Y = f \nabla_X Y$,
2. (Leibnitz) $\nabla_X (fY) = (Xf)Y + f \nabla_X Y$,
3. (torsion free) $\nabla_X Y - \nabla_Y X = [X, Y]$, where $[\cdot, \cdot]$ is the Lie bracket, and
4. (metric compatibility) $Z\langle X, Y \rangle = \langle \nabla_Z X, Y \rangle + \langle X, \nabla_Z Y \rangle$.

Definition 9 (Christoffel symbols). The method we will employ to completely describe a connection (as defined in Definition 8) locally is to describe them through Christoffel symbols. Following the definition of [11, p. 100], let ∇ be an affine connection on \mathcal{M} . Denote a coordinate vector field on the coordinate open set $(U, p^1, \dots, p^n) \subseteq \mathcal{M}$ by $\partial_i := \partial/\partial p^i$. In this coordinate frame there exist the numbers called Christoffel symbols, Γ_{ij}^k , defined through the following

$$\nabla_{\partial_i} \partial_j = \sum_{k=1}^n \Gamma_{ij}^k \partial_k =: \Gamma_{ij}^T \partial.$$

Definition 10 (Retraction). Following [4, Def. 3.47], a retraction on a smooth manifold \mathcal{M} is a smooth map,

$$\mathcal{R}: T\mathcal{M} \rightarrow \mathcal{M}, \quad (p, X) \mapsto \mathcal{R}_p(X)$$

such that every curve generated from $c(t) = \mathcal{R}_p(tX)$ satisfies $c(0) = p$ and $\dot{c}(0) = X$. Equivalently the conditions can be stated as in [4, p. 40] without the use of curves. For all $p \in \mathcal{M}$, $\mathcal{R}_p(0) = p$, and $D\mathcal{R}_p(0): T_p\mathcal{M} \rightarrow T_p\mathcal{M}$, $D\mathcal{R}_p(0)[X] = X$ is the identify map.

Definition 11.

Definition 12.

For the rest of this paper we denote \mathcal{M} as being a Riemannian manifold.

After the basic definitions, talk about how the rest of the theory is a highlighted summary through BZ and JZ. The goal is to look at the findings in JZ, however it relies heavily on theory derived in BZ. It will be mentioned of some parts are from other works, or if they are original work.

define or ref
w sum? i
wrote it in
obsidian in
the H and V
space file.

In JZ this
was stated
wrongly
citation

do i need
a source to
define this?

2.2 The Symplectic group

The real *symplectic group* is the space overarching the symplectic Stiefel manifold, and we will look at this space first. To be able to define the symplectic group we first need some preliminary definitions. Define the *symplectic identity* as the following block matrix,

$$J_{2n} := \begin{bmatrix} 0_n & I_n \\ -I_n & 0_n \end{bmatrix},$$

where I_n denotes the $n \times n$ identity matrix. J_{2n} has some properties we will take advantage of frequently:

$$J_{2n}^T = -J_{2n} = J_{2n}^{-1} \quad (2.1)$$

The symplectic group is defined as the set of matrices which define the symplectic structure in the following sense. We define the real symplectic group as

$$\mathrm{Sp}(2n) := \{p \in \mathbb{R}^{2n \times 2n} \mid p^+ p = I_{2n}\}, \quad (2.2)$$

where $^+$ is defined as the *symplectic inverse* of any matrix $q \in \mathbb{R}^{2n \times 2k}$ such that

$$q^+ := J_{2k}^T q^T J_{2n}. \quad (2.3)$$

The Lie algebra of $\mathrm{Sp}(2n)$ is the symplectic groups' tangent space at the identity. It is given by

$$\mathfrak{sp}(2n) := \{\Omega \in \mathbb{R}^{2n \times 2n} \mid \Omega^+ = -\Omega\}, \quad (2.4)$$

where Ω is called the Hamiltonian matrix [ref ??](#). Now we can define the tangent space of $\mathrm{Sp}(2n)$ at any point p as

$$T_p \mathrm{Sp}(2n) = \{\Omega p \in \mathbb{R}^{2n \times 2n} \mid \Omega \in \mathfrak{sp}(2n)\}. \quad (2.5)$$

2.3 The Symplectic Stiefel manifold

Now that we have defined the symplectic group, we introduce the manifold of interest, the real symplectic Stiefel manifold. It is defined as

$$\mathrm{SpSt}(2n, 2k) := \{p \in \mathbb{R}^{2n \times 2k} \mid p^+ p = I_{2k}\}, \quad (2.6)$$

where p^+ is as in (2.3). Following [2, Prop. 3.1], it is explicitly connected to the symplectic group in the sense that $\mathrm{SpSt}(2n, 2k)$ is diffeomorphic to the following quotient manifold of $\mathrm{Sp}(2n)$:

$$\mathrm{SpSt}(2n, 2k) \cong \mathrm{Sp}(2n) / \mathrm{Sp}(2(n - k)),$$

where the notion of quotient manifold is as in Definition 3. It has dimension $\dim(\mathrm{SpSt}(2n, 2k)) = (2n - 2k + 1)k$.

The following piece of insight can give some further intuition on what the Symplectic Stiefel manifold is. We note that the Stiefel manifold is a quotient space, as defined in Definition 3, of the orthogonal group as defined in Definition 2 such that $\mathrm{St}(2n, 2k) = \mathrm{O}(n) / \mathrm{O}(n - k)$.

The expression for the tangent space follows straightforwardly from the definition of $\mathrm{SpSt}(2n, 2k)$. Assume we have a curve, $c(t) \in \mathrm{SpSt}(2n, 2k)$, s.t. $c(0) = p$ and $\dot{c}(0) := \frac{d}{dt}c(t)|_{t=0} = X$. Since $c(t)$ is a curve in $\mathrm{SpSt}(2n, 2k)$, by (2.6) it must satisfy the following condition:

$$c(t)^T J_{2n} c(t) = J_{2k}. \quad (2.7)$$

Taking the derivative of (2.7) with respect to t at $t = 0$ we get

$$\dot{c}^T(t) J_{2n} c(t) + c^T(t) J_{2n} \dot{c}|_{t=0} = X^T J_{2n} p + p^T J_{2n} X = 0_{2k}.$$

After moving the first term over to the left hand side, and multiplying with J_{2n} from the left, we get

$$p^+ X = -X^+ p.$$

We recognize this condition as $p^+ X \in \mathfrak{sp}(2k)$ as defined in (2.4). This means that for a point p ,

$$T_p \mathrm{SpSt}(2n, 2k) = \{X \in \mathbb{R}^{2n \times 2k} \mid p^+ X \in \mathfrak{sp}(2k)\}. \quad (2.8)$$

Why do we need sp? We will use quotient properties to map stuff to Spst

Do I want to find a historical reference?

maybe explain a little more how we get this tangent space

Is a version of this info too much of a detour?

3 Right-Invariant Framework on the Symplectic Stiefel Manifold

One of the key insights of Bendokat and Zimmermann [2, p. 11] is that using a right invariant framework one is able to construct geodesics on $\text{SpSt}(2n, 2k)$. Following in their— footsteps, we will first define a right-invariant metric on $\text{Sp}(2n)$ and its corresponding geodesics, then transport this metric to $\text{SpSt}(2n, 2k)$. This will allow us to define geodesics on $\text{SpSt}(2n, 2k)$. Finally, through this framework we will be able to define the Riemannian gradient, hessian, and other tools necessary for the implementation of the optimization algorithms on the $\text{SpSt}(2n, 2k)$.

3.1 Right-invariant metric

We begin by defining the point-wise right-invariant metric on $\text{Sp}(2n)$ as the mapping $g_p^{\text{Sp}}: T_p\text{Sp}(2n) \times T_p\text{Sp}(2n) \rightarrow \mathbb{R}$,

$$g_p^{\text{Sp}}(X_1, X_2) := \frac{1}{2} \text{tr}((X_1 p^+)^T X_2 p^+), \quad X_1, X_2 \in T_p\text{Sp}(2n). \quad (3.1)$$

It is right-invariant in the sense that $g_{pq}^{\text{Sp}}(X_1 q, X_2 q) = \frac{1}{2} \text{tr}((X_1 q q^+ p^+)^T X_2 q q^+ p^+) = g_p^{\text{Sp}}(X_1, X_2)$ for all $p \in \text{Sp}(2n)$. Now that we have defined g_p^{Sp} , we want to, in a sense, transport it to $\text{SpSt}(2n, 2k)$ in a way that preserves the right-invariance. To achieve this we will use a *horizontal lift* to define a metric on $\text{SpSt}(2n, 2k)$ through 3.1. Split $T_p\text{Sp}(2n)$ into two parts: the horizontal- and vertical part, with respect to g_p^{Sp} and π :

rewrite these two sentences

$$T_p\text{Sp}(2n) = \text{Ver}_p \text{Sp}(2n) \oplus \text{Hor}_p \text{Sp}(2n). \quad (3.2)$$

Noting that $\mathfrak{sp}(2n) = T_I\text{Sp}(2n)$, we can express these spaces through $\mathfrak{sp}(2n)$ as

$$\begin{aligned} \text{Ver}_p \text{Sp}(2n) &:= \{\Omega p \mid \Omega \in \text{Ver } \mathfrak{sp}(2n)\}, \\ \text{Hor}_p \text{Sp}(2n) &:= \{\Omega p \mid \Omega \in \text{Hor } \mathfrak{sp}(2n)\}, \end{aligned}$$

where horizontal- and vertical space is defined as in Definition 7. Explicit expressions of $\text{Ver } \mathfrak{sp}(2n)$ and $\text{Hor } \mathfrak{sp}(2n)$ is given in [2, p. 11]. We will only need the horizontal space to define our desired geodesic, so we will only explicitly state $\text{Hor } \text{Sp}(2n)$. It is expressed as

$$\text{Hor}_p \text{Sp}(2n) = \{\bar{\Omega} p \mid \Omega r + r \Omega - r \Omega r = \bar{\Omega} \in \text{Hor } \mathfrak{sp}(2n)\},$$

where $r = J_{2n}^T \pi(p) \pi(p)^+ J_{2n}$. From this quotient perspective, for $p = \pi(q)$ where $q \in \text{Sp}(2n)$, we can express vectors from $T_p\text{SpSt}(2n, 2k)$ through $\text{Hor}_p \text{Sp}(2n)$. Following [8, p. 5] any $X \in T_p\text{SpSt}(2n, 2k)$ and a specific $p = \pi(q)$, there exists a unique horizontal lift (see Definition ??)

$$\mathfrak{h}_q(X) = \bar{\Omega}(X) q, \quad (3.3)$$

where

$$\bar{\Omega}(X) = X(p^T p)^{-1} p^T + J_{2n} p (p^T p)^{-1} X^T (I_{2n} - J_{2n}^T p (p^T p)^{-1} p^T J_{2n}) J_{2n}. \quad (3.4)$$

Look at 20, Theorem 2.28 in BZ to maybe skip the derivation up to gp

The point-wise right-invariant Riemannian metric on $\text{SpSt}(2n, 2k)$ can then be defined as the mapping $g_p: T_p\text{SpSt}(2n, 2k) \times T_p\text{SpSt}(2n, 2k) \rightarrow \mathbb{R}$, $g_p(X_1, X_2) := g_p^{\text{Sp}}(\mathfrak{h}_q(X_1), \mathfrak{h}_q(X_2))$. More explicitly

$$g_p(X_1, X_2) = \text{tr} \left(X_1^T \left(I_{2n} - \frac{1}{2} J_{2n}^T p (p^T p)^{-1} p^T J_{2n} \right) X_2 (p^T p)^{-1} \right), \quad (3.5)$$

for $X_1, X_2 \in T_p\text{SpSt}(2n, 2k)$.

3.2 Geodesics

To be able to define what it means to travel a (local) path of minimal length on any manifold, we need to define geodesics. They are the generalization of straight lines in Euclidean space, and since most optimization algorithms travel along paths, they are essential in transferring Euclidean optimization algorithms to the Riemannian domain. In this section we will begin by defining geodesics on $\text{Sp}(2n)$. Through the lens of quotient manifolds, and our right invariant framework, we will define geodesics on $\text{SpSt}(2n, 2k)$ from the geodesics on $\text{Sp}(2n)$.

Following [2, Prop. 2.1], given $p \in \text{Sp}(2n)$, $X \in T_p \text{Sp}(2n)$ and the right-invariant Riemannian metric (3.1), the respective geodesic $\gamma(t)$ is defined as

$$\gamma(t) := \exp(t(Xp^+ - (Xp^+)^T)) \exp(t(Xp^+)^T)p,$$

where $\gamma(0) = p$, $\dot{\gamma}(0) = X$ and $^+$ is the symplectic inverse (as in (2.3)). Here, \exp denotes the matrix exponential.

To be able define geodesics from the right invariant metric on $\text{SpSt}(2n, 2k)$ we need the following result. According to cite[Cor. 7.46]???, if g_p^{SpSt} has a horizontal tangent vector at every point, it projects to a Riemannian geodesic on $\text{SpSt}(2n, 2k)$. ref

The last preliminary definition we need is the following. In [BZ] Bendokat and Zimmermann proves that if we, for a point $p \in \text{Sp}(2n)$, define a geodesic $\gamma(t)$ through a horizontal tangent vector $X \in \text{Hor}_p \text{Sp}(2n)$, then $\dot{\gamma}(t) \in \text{Hor}_{\gamma(t)} \text{Sp}(2n)$. We call such a geodesic a *horizontal geodesic*.

Following [BZ], we can now define a geodesic on $\text{SpSt}(2n, 2k)$ through a horizontal geodesic $\gamma(t) \subseteq \text{Sp}(2n)$. Let $q \in \text{SpSt}(2n, 2k)$, $Y \in T_q \text{SpSt}(2n, 2k)$, and $p \in \pi^{-1}(q) \subset \text{Sp}(2n)$. Then the geodesic from p in direction Y is

$$\varphi(t) = \pi(\gamma(t)) = \exp(t(\bar{\Omega}(Y) - \bar{\Omega}(Y)^T)) \exp(t\bar{\Omega}(Y)^T)p. \quad (3.6)$$

3.3 Riemannian gradient of the Symplectic Stiefel manifold

Another component to many optimization algorithms is the gradient. In the Euclidean setting the gradient of the cost function is a vector that points in the direction of the steepest ascent. Analogously, the Riemannian gradient points in the direction of the steepest ascent on the manifold. A useful analogy of what this means is to consider the sphere, and some cost function defining a surface on the sphere. The riemannian gradient would then point in the direction towards the closest highest point on the sphere. In this section we will state the Riemannian gradient explicitly, and prove that it satisfies necessary properties.

Proposition 1. *Given a function $f: \text{SpSt}(2n, 2k) \rightarrow \mathbb{R}$, the Riemannian gradient with respect to g_p is given by* add general def of R grad.?

$$\text{grad } f(p) = \nabla f(p)p^T p + J_{2n}p(\nabla f(p))^T J_{2n}p, \quad (3.7)$$

where $\nabla f(p)$ is the Euclidean gradient of a smooth extension around $p \in \text{SpSt}(2n, 2k)$ in $\mathbb{R}^{2n \times 2k}$ at p .

Proof. We can see that this is the Riemannian gradient by the following two observations stated in [2, p. 12], which we verify ourselves below.

Firstly, gradient must be in $T_p \text{SpSt}(2n, 2k)$, which means by (2.8) that $0 = p^+ \text{grad } f(p) + (\text{grad } f(p))^+ p$. Computing this we get

$$p^T J \nabla f(p) p^T p + p^T J J p (\nabla f(p))^T J p + p^T p (\nabla f(p))^T J p + p^T J^T \nabla f(p) p^T J^T J p = 0,$$

where we have used $JJ = -J^T J = -I_{2n}$ and (2.1).

Secondly, the gradient also has to satisfy $g_p(\text{grad } f(p), X) = df_p(X) = \text{tr}((\nabla f(p))^T X)$ for all $X \in T_p \text{SpSt}(2n, 2k)$:

$$g_p(\text{grad } f(p), X) = \text{tr}((p^T p (\nabla f(p))^T + p^T J^T \nabla f(p) p^T J^T)(I_{2n} - \frac{1}{2}G)X(p^T p)^{-1}),$$

where $G := J^T p(p^T p)^{-1} p^T J$. Expanding this expression we obtain

$$\begin{aligned} &= \text{tr}(p^T p(\nabla f(p))^T X(p^T p)^{-1}) - \frac{1}{2} \text{tr}(p^T p(\nabla f(p))^T G X(p^T p)^{-1}) \\ &+ \text{tr}(p^T J^T \nabla f(p) p^T J^T X(p^T p)^{-1}) - \frac{1}{2} \text{tr}(p^T J^T \nabla f(p) p^T J^T G X(p^T p)^{-1}), \end{aligned}$$

where the cancellations used the fact that the trace is invariant under circular shifts. Noting that the first term is by definition $d f_p(X)$, and inserting the definition of G , the expression becomes

$$\begin{aligned} &= d f_p(X) - \frac{1}{2} \text{tr}((\nabla f(p))^T J^T p(p^T p)^{-1} p^T J X) \\ &+ \text{tr}(p^T J^T \nabla f(p) p^T J^T X(p^T p)^{-1}) \\ &- \frac{1}{2} \text{tr}(p^T J^T \nabla f(p) p^T J^T J^T p(p^T p)^{-1} p^T J X(p^T p)^{-1}). \end{aligned}$$

After using $J^T J^T = -I_{2n}$ and (2.1) on the last term, we notice that we can cancel $p^T p(p^T p)^{-1}$, making it equal to the second to last term. Now focusing on the second term: for the first equality we use the fact that for any matrix, A , $\text{tr}(A) = \text{tr}(A^T)$, and for the second equality we utilize the cyclic property of the trace, and (2.1),

$$\begin{aligned} \frac{1}{2} \text{tr}((\nabla f(p))^T J^T p(p^T p)^{-1} p^T J X) &= \frac{1}{2} \text{tr}(X^T J^T p(p^T p)^{-1} p^T J \nabla f(p)) \\ &= -\frac{1}{2} \text{tr}(p^T J^T \nabla f(p) X^T J^T p(p^T p)^{-1}) \end{aligned} \quad (3.8)$$

Inserting (3.8) into our expression we end up with:

$$d f_p(X) = d f_p(X) + \frac{1}{2} \text{tr}(p^T J \nabla f(p) X^T J^T p(p^T p)^{-1}) + \frac{1}{2} \text{tr}(p^T J^T \nabla f(p) p^T J^T X(p^T p)^{-1}),$$

where the last two terms cancel after applying (2.1), and the tangent space condition (2.8), $p^T J X = -X^T J p$. \square

3.4 Riemannian Hessian of the Symplectic Stiefel manifold

Define the remaining preliminaries for Hessian. Make an intro to section. loosely following Appendix A in Jz

If one knows more about the manifold, one can make more assumptions for an optimization algorithm. The hope would be that these additional assumptions would make the algorithm more efficient. In the Euclidean setting, the Hessian gives us curvature information about the cost function. Generalizing the Hessian to the Riemannian setting it gives us something similar, and we can generalize many popular algorithms. In this section we will therefore define the Riemannian Hessian of the symplectic manifold. We begin by defining what it means to take a second derivative on a manifold. After defining the necessary tools, we will give an analytical expression for the Riemannian Hessian.

For two smooth vector fields, $\mathcal{X}(p) = \alpha^T \boldsymbol{\theta}$ and $\mathcal{Y}(p) = \beta^T \boldsymbol{\theta}$ defined as in Definition 6, the covariant derivative (defined through the Riemannian connection by Definition 8) written in local coordinates is

$$\nabla_{\mathcal{X}} \mathcal{Y} = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \partial_i (\beta_j) \partial_j + \alpha_i \beta_j \sum_{k=1}^n \Gamma_{ij}^k \partial_k.$$

In preparation for the Hessian, we include the description of the covariant derivative from [11, p. 96] to restrict the covariant derivative further. We want to define the covariant derivative of a vector field along a curve $c(t)$. $c(t)$ is the smooth curve, $c: I \rightarrow \mathcal{M}$, $t \mapsto (\gamma_1(t), \dots, \gamma_n(t))$, where $I := [a, b] \subseteq \mathbb{R}$. Since we have an affine connection on \mathcal{M} , the following unique map exists:

$$\frac{D}{dt}: \Gamma(T\mathcal{M}|_{c(t)}) \rightarrow \Gamma(T\mathcal{M}|_{c(t)}),$$

where $\Gamma(T\mathcal{M}|_{c(t)})$ denotes the the vector space of all smooth vector fields along $c(t)$. If $V \in \Gamma(T\mathcal{M}|_{c(t)})$ is induced by \mathcal{X} , meaning $V(t) = \mathcal{X}|_{c(t)}$, then

$$\frac{DV}{dt}(t) = \nabla_{\dot{c}(t)} \mathcal{X} = \dot{\alpha}(t) + \Gamma(\alpha(t), \dot{\gamma}(t)), \quad \Gamma(u, v) = \begin{bmatrix} u^T \Gamma^1 v \\ \vdots \\ u^T \Gamma^n v \end{bmatrix}.$$

$\Gamma(u, v)$ is called the *Christoffel function*. If $\dot{c}(t)$ is a geodesic, the expression above reduces to

$$\ddot{\gamma}(t) = -\Gamma(\dot{\gamma}(t), \dot{\gamma}(t)), \quad (3.9)$$

since by definition geodesics must satisfy $\frac{D}{Dt}\dot{\gamma}(t) = 0$ and $\dot{\alpha}(t) = \dot{\gamma}(t)$. Importantly, once we have found the Christoffel symbols through (3.9), we can still use them for curves that are not geodesics. This is because the Christoffel symbols only depend on the Riemannian metric, and the local coordinates. To do this, we recover the Christoffel function for two different inputs through polarization [5, p. 312]

$$\Gamma(X, Y) = \frac{1}{4}(\Gamma(X + Y, X + Y) - \Gamma(X - Y, X - Y)),$$

where $X, Y \in \Gamma(T\mathcal{M}|_{c(t)})$.

i do not understand polarization.

To find the Christoffel symbols for $\text{SpSt}(2n, 2k)$ with respect to the right invariant metric g defined in (3.5), we differentiate the geodesic formula from ref, and use (3.9) to achieve the following formula,

$$\Gamma(X, X) = -\ddot{\gamma}(0) = -(\overline{\Omega}(X) - \overline{\Omega}(X)^T)(X + \overline{\Omega}(X)^T p) - (\overline{\Omega}(X)^T)^2 p.$$

Here $X = \dot{\gamma}(0) \in T_p \text{SpSt}(2n, 2k)$, $p \in \text{SpSt}(2n, 2k)$, and $\overline{\Omega}(X)$ is as in ref. With our metric g , the Hessian at p of a smooth function $f: \text{SpSt}(2n, 2k) \rightarrow \mathbb{R}$ is the endomorphism

$$\begin{aligned} \text{Hess } f(p): T_p \text{SpSt}(2n, 2k) &\rightarrow T_p \text{SpSt}(2n, 2k), \\ \text{Hess } f(p)[X] &= \left. \frac{d}{dt} \text{grad } f(c(t)) \right|_{t=0} + \Gamma(\text{grad } f(p), X), \end{aligned}$$

where $\text{grad } f(\cdot)$ is as in 3.7. $c(t) \in \text{SpSt}(2n, 2k)$ is an arbitrary curve such that $c(0) = p$ and $c'(0) = X$

3.5 Moving from Theory to Application

Additional stuff to bridge the gap between theory and application. I.e. cayley transformation (p.20 in BZ, and prop. 5.2). Approx to the Riemannian hessian

4 Algorithms

In this section, we introduce the algorithms of interest in this Specialization project. They are used by Jensen and Zimmermann for optimization on $\text{SpSt}(2n, 2k)$. The goal is to study the feasibility of these algorithms by attempting to reproduce some of the findings of Jensen and Zimmermann.

4.1 Riemannian gradient descent

The first algorithm we will define is Riemannian gradient descent (GD). Although seemingly simple, the Euclidean gradient descent (E-GD) is a powerful algorithm because of its simplicity. The lack of assumptions makes the algorithm easily applicable to a wide range of problems, and its simplicity also makes it computationally efficient. Initially following [4, p. 56], the E-GD is intuitively transferred to the Riemannian framework. For the sequence $\{p_k\}$ where $k \in \mathbb{N}$, we have the point $p_k \in \text{SpSt}(2n, 2k)$, where the next point in the sequence is computed as

$$p_{k+1} = \mathcal{R}_{p_k}(-t_k \text{grad } f(p_k)).$$

Here $t_k > 0$ is some step-size to be determined.

i felt JZ was unclear. Is this correct?

refs

ref geod.

define this somewhere

To ensure that each step sufficiently decreases the cost function, we will use the Riemannian version of the Armijo condition to compute t_k . In [6, p. 17] the Armijo condition for each iteration k given for $\beta \in (0, 1)$, and a search direction X_k as

$$f(\mathcal{R}_{p_k}(t_k X_k)) \leq f(p_k) + \beta t_k \langle \text{grad } f(p_k), X_k \rangle_{p_k}. \quad (4.1)$$

The step-size t_k is calculated as $t_k = \gamma_k \delta^h$, where $\gamma \in (0, 1)$ is the backtracking parameter and h is the smallest integer such that (4.1) is satisfied.

Algorithm 1 Riemannian Gradient descent

Input: Initial point $p_0 \in \text{SpSt}(2n, 2k)$, objective function $f: \text{SpSt}(2n, 2k) \rightarrow \mathbb{R}$, retraction \mathcal{R} , parameters $\beta, \gamma \in (0, 1)$, steplength range $0 < \gamma_{\min} < \gamma_{\max}$, initial step size $\gamma_0 = f(p_0)$, maximum number of iterations $N \in \mathbb{N}$, step parameters $h_{\min} < h_{\max} \in \mathbb{Z}$, tolerance parameters $\epsilon, \epsilon_x, \epsilon_f > 0$, Riemannian metric $\langle \cdot, \cdot \rangle_p$, with gradient grad_f where $\langle \cdot, \cdot \rangle$ denotes the Euclidean inner product.

```

1: for  $0 \leq k \leq N$  do
2:    $X_k = -\text{grad } f(p_k)$ 
3:    $\gamma_k = \max(\gamma_{\min}, (\gamma_k, \gamma_{\max}))$ 
4:   Find  $t_k$  through solving (4.1)
5:    $p_{k+1} = \mathcal{R}_{p_k}(t_k X_k)$ 
6:   if  $\|\text{grad } f(p_k)\|_F < \epsilon$  then
7:     if  $\frac{|f(p_k) - f(p_{k+1})|}{|f(p_k) + 1|} < \epsilon_f$  and  $\frac{\|p_k - p_{k+1}\|_F}{\sqrt{2n}} < \epsilon_x$  then
8:       Break
9:     end if
10:  end if
11: end for

```

Output: Iterates $\{p_k\}$

For robustness it is important to know if this algorithm behaves properly, in the sense that we want it to reliably converge to critical points of f . It is shown in [6, Cor. 5.8] that for $\text{SpSt}(2n, 2k)$ Algorithm 2 generates an infinite sequence of iterates $\{p_k\}$ [6, Prop. 5.6]. It can be shown that every accumulation point p^* of $\{p_k\}$ is a critical point of f , meaning $\text{grad } f(p^*) = 0$.

4.2 Riemannian trust-region method

Since the gradient descent method only uses first order information, a natural question would be to investigate how a method that utilizes second order information would perform. The criterion for this second order method to be considered better, for a specific problem, than GD would be that it despite the (expected) increased computing time per step, it would converge with few enough steps as to still beat CG.

In choosing a second order method, the simplest choice would be a Riemannian version of Newton's method (NM). Despite its simple design, NM is known to be quite unstable, and need to be sufficiently close to a local minima to guarantee convergence [4, p. 122]. The trust-region method (TR) addresses several of the problems with NM, while still having comparatively fast convergence properties as NM locally.

Following [4, p. 131], the goal of each step k of TR is for a point p_k , is to approximate the pullback $f \circ \mathcal{R}_{p_k}(X)$ through an approximation of $T_{p_k} \text{SpSt}(2n, 2k)$,

$$f(\mathcal{R}_{p_k}(X)) \approx m_{p_k}(X) = f(p_k) + \langle \text{grad } f(p_k), X \rangle_{p_k} + \frac{1}{2} \langle H_{p_k}(X), X \rangle_{p_k}.$$

H_{p_k} can be any self-adjoint linear map on $T_{p_k} \text{SpSt}(2n, 2k)$, but in our experiments H_{p_k} will be either $\text{Hess } f(p_k)[X]$ as in ??, or $\pi_{p_k}(\text{D grad } f(p_k)[X])$ as in ??. From [4, Prop. 5.44] the convergence rate of this method is essentially the same as $\text{Hess } f(p)[X]$.

To choose a step size for TR we demand that the step must reduce the value of $m_{p_k}(X)$. Since our model is an approximation of the tangent space, we construct a trust region which is the region

around p_k where we assume that the error in our approximation is negligible. We solve the TR subproblem

$$\min_{X \in T_{p_k} \text{SpSt}(2n, 2k)} m_k(X) \quad \text{subject to} \quad \|X\|_{p_k} \leq \Delta_k \quad (4.2)$$

to find the candidate step X_k , where candidate for next iterate then is $\hat{p}_k = \mathcal{R}_{p_k}(X_k)$. Δ_k denotes the radius of the trust region at that iterate. Depending on how the new step performs (see line ? of Algorithm ??) it is either accepted or rejected. Finally the trust region radius is evaluated to see if it needs to be modified. The procedure is codified in Algorithm ??, which is adapted from [4, Algorithm 3.3]. Further reading can be found in [4, p. 131].

Algorithm 2 Riemannian Trust-region method

Input: Initial point $p_0 \in \text{SpSt}(2n, 2k)$, objective function $f: \text{SpSt}(2n, 2k) \rightarrow \mathbb{R}$, retraction \mathcal{R} , maximum number of iterations $N \in \mathbb{N}$, maximal radius $\overline{\Delta} > 0$, initial radius $\Delta_0 \in (0, \overline{\Delta})$, ratio of model improvement threshold $\gamma_{\min} > 0$

tolerance parameters $\epsilon, \epsilon_x, \epsilon_f > 0$, Riemannian metric $\langle \cdot, \cdot \rangle_p$, with gradient grad_f where $\langle \cdot, \cdot \rangle$ denotes the Euclidean inner product.

- 1: **for** $0 \leq k \leq N$ **do**
- 2: Find X_k through solving (4.2)
- 3: $\hat{p} = \mathcal{R}_{p_k}(X_k)$
- 4: $\gamma_k = (f(p_k) - f(\hat{p})) / (m_k(0) - m_k(X_k))$
- 5: Compute new iterate

$$p_{k+1} = \begin{cases} \hat{p} & \text{if } \gamma_k > \gamma_{\min} \\ p_k & \text{otherwise} \end{cases}$$

- 6: Compute new trust-region radius

$$\Delta_{k+1} = \begin{cases} \frac{1}{4}\Delta_k & \text{if } \gamma_k < \frac{1}{4} \\ \min \left\{ 2\Delta_k, \overline{\Delta}_k \right\} & \text{if } \gamma_k > \frac{3}{4} \text{ and } \|X_k\| = \Delta_k \\ \Delta_k & \text{otherwise} \end{cases}$$

- 7: Todo: add tolerances
- 8: **end for**

Output: Iterates $\{p_k\}$

Convergence and stability of TR is presented in-depth in [4, p. 147]. However, for consistency we note here that given sufficient assumptions (which are met in our examples), by [4, Cor. 6.24] for $\{p_k\}$ generated by TR,

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

In other words, this means that for all $\epsilon > 0$ and K there exists $k \geq K$ such that $\|\text{grad } f(p_k)\|_{p_k} \leq \epsilon$.

4.3 Implementation

All files used to produce the data in this report are available in the following [github repository](#)[7].

5 Numerical Experiments

The following experiments were performed using Julia version 1.10.2 on the Laptop ACER Swift SF314-43 processor AMD Ryzen 7 5700U with Radeon Graphics 1.80GHz and 16GB of RAM.

To test and compare the feasibility of our algorithms we will, as in [8, p. 15], try to solve the following problem. For a matrix $q \in \mathbb{R}^{2n \times 2k}$ we want to find the closest symplectic matrix $p \in \text{SpSt}(2n, 2k)$. We formalize this in as the following optimization problem,

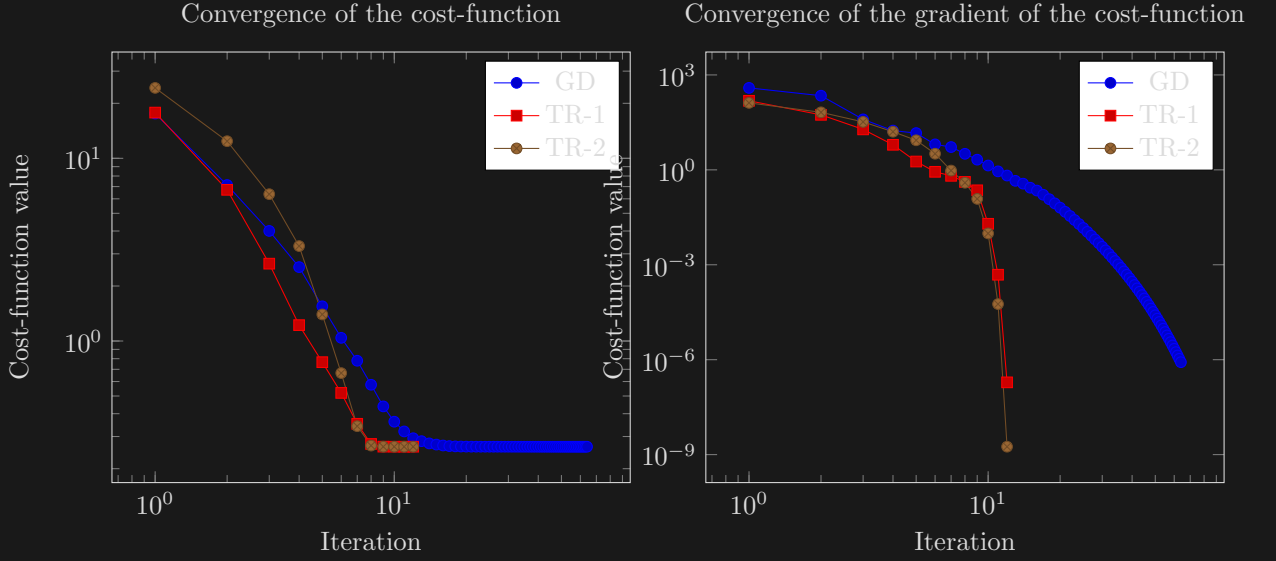
$$\min_{p \in \text{SpSt}(2n, 2k)} f(p), \quad (5.1)$$

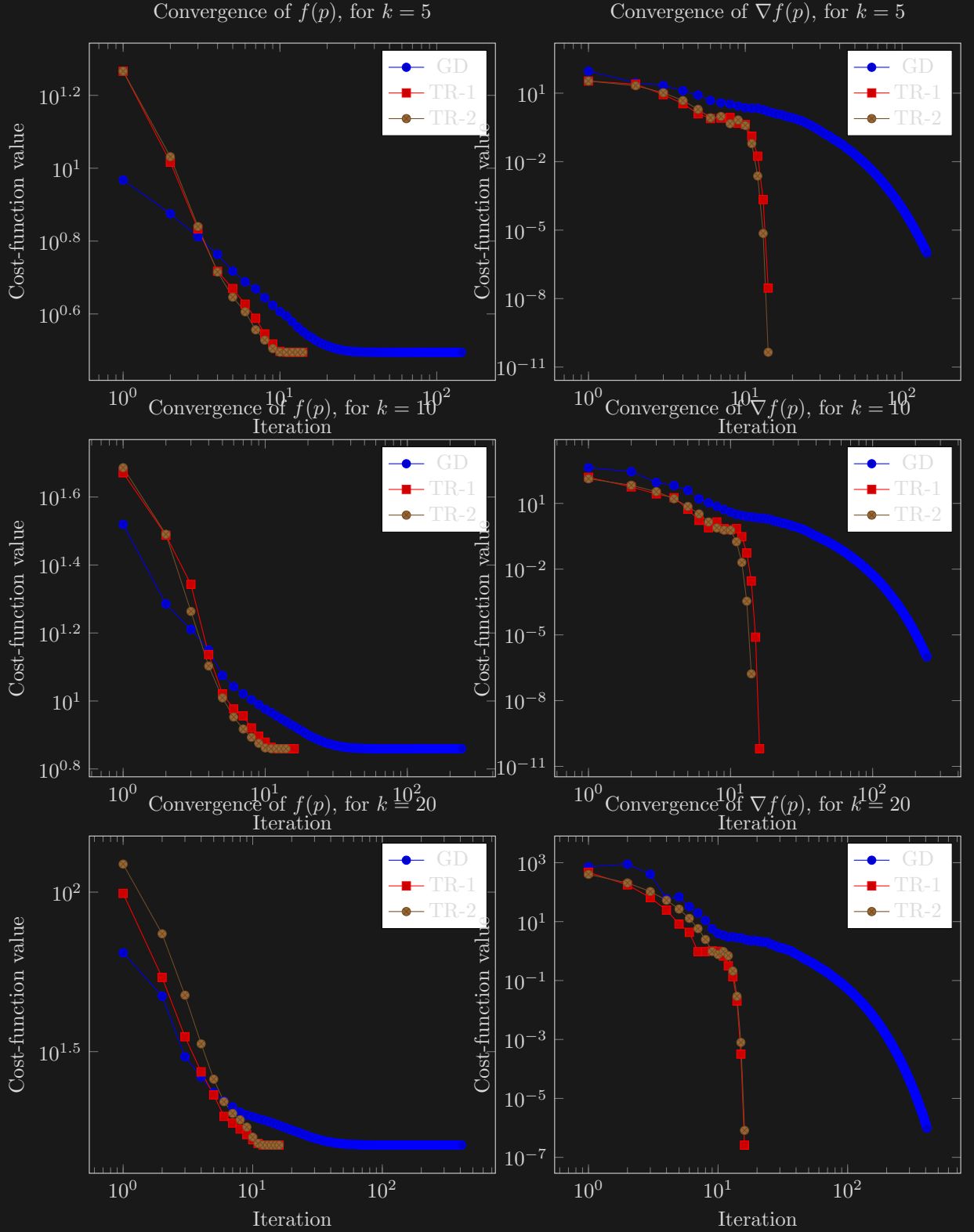
where $f(p) := \frac{1}{2} \|q - p\|_F^2$. For a point $p \in \text{SpSt}(2n, 2k)$ and $X \in T_p \text{SpSt}(2n, 2k)$, Euclidean gradient and Hessian are, respectively,

$$\nabla f(p) = p - q, \quad \nabla^2 f(p)[X] = X.$$

For the experiments we generate q randomly, and normalize it, $q \cdot \|q\|_F^{-1}$. For the optimization runs we choose $n = 1000$ and $k = \{10, 50, 100\}$. The results is displayed in Table ??? and in Figure ??.

5.1 Nearest symplectic matrix problem





6 Conclusion

Jensen and Zimmermann chose to leave out the Riemannian BFGS method from their experiments [8, p. 11]. They did this because they found it to not be competitive to the other methods they used. Despite this it could be interesting to try to validate these findings using Manopt.jl and/or

Table 1: Nearest symplectic matrix problem solved by CG, TR-1 and TR-2. The table summarizes time to converge for all the algorithms, with $n = 100$ for $k = 5, 10, 20$

	Runtime (s)		
	GD	TR-1	TR-2
$k = 5$	0.051	9.7	0.70
$k = 10$	0.23	12	5.4
$k = 20$	1.1	35	23

on problems expected to be better suited to the Riemannian BFGS method.

Bibliography

- [1] P-A Absil, R Mahony and R Sepulchre. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2008. ISBN: 9780691132983.
- [2] Thomas Bendokat and Ralf Zimmermann. *The real symplectic Stiefel and Grassmann manifolds: metrics, geodesics and applications*. 2021. arXiv: [arXiv:2108.12447v1](https://arxiv.org/abs/2108.12447).
- [3] William M. Boothby. *An Introduction to Differentiable Manifolds and Riemannian Geometry*. Vol. 63. Elsevier Science (USA), 1975. DOI: [https://doi.org/10.1016/S0079-8169\(08\)61028-4](https://doi.org/10.1016/S0079-8169(08)61028-4).
- [4] Nicolas Boumal. *An introduction to optimization on smooth manifolds*. Cambridge University Press, 2023. DOI: [10.1017/9781009166164](https://doi.org/10.1017/9781009166164). URL: <https://www.nicolasboumal.net/book>.
- [5] Alan Edelman, Tomás A. Arias and Steven T. Smith. ‘The Geometry of Algorithms with Orthogonality Constraints’. In: *SIAM Journal on Matrix Analysis and Applications* 20 (1998). DOI: [10.1137/s0895479895290954](https://doi.org/10.1137/s0895479895290954). arXiv: [arXiv:physics/9806030v1](https://arxiv.org/abs/9806030).
- [6] Bin Gao et al. ‘Riemannian optimization on the symplectic Stiefel manifold’. In: *SIAM Journal on Optimization* 31 (2021). DOI: <https://doi.org/10.1137/20M1348522>. arXiv: [arXiv:2006.15226](https://arxiv.org/abs/2006.15226).
- [7] Ole Gunnar Røsholt Hovland. *TMA4500-Project-Hovland-Symplectic-Stiefel*. 2024. URL: <https://github.com/kellertuer/TMA4500-Project-Hovland-Symplectic-Stiefel.git>.
- [8] Rasmus Jensen and Ralf Zimmermann. *Riemannian optimization on the symplectic Stiefel manifold using second-order information*. 2024. arXiv: [arXiv:2404.08463v2](https://arxiv.org/abs/2404.08463).
- [9] John M. Lee. *Introduction to Smooth Manifolds*. 2nd ed. Springer New York, NY, 2012. DOI: [10.1007/978-1-4419-9982-5](https://doi.org/10.1007/978-1-4419-9982-5).
- [10] Department of Marine Technology NTNU. *IMT Software Wiki - LaTeX*. URL: <https://www.ntnu.no/wiki/display/imtsoftware/LaTeX> (visited on 15th Sept. 2020).
- [11] Loring W. Tu. *Differential Geometry*. Springer International Publishing, 2017. DOI: [10.1007/978-3-319-55084-8](https://doi.org/10.1007/978-3-319-55084-8).