

RIEMANNIAN HAMILTONIAN METHODS FOR MIN-MAX OPTIMIZATION ON MANIFOLDS

ANDI HAN*, BAMDEV MISHRA†, PRATIK JAWANPURIA†,
PAWAN KUMAR‡, AND JUNBIN GAO*

Abstract. In this paper, we study min-max optimization problems on Riemannian manifolds. We introduce a Riemannian Hamiltonian function, minimization of which serves as a proxy for solving the original min-max problems. Under the Riemannian Polyak–Lojasiewicz condition on the Hamiltonian function, its minimizer corresponds to the desired min-max saddle point. We also provide cases where this condition is satisfied. For geodesic-bilinear optimization in particular, solving the proxy problem leads to the correct search direction towards global optimality, which becomes challenging with the min-max formulation. To minimize the Hamiltonian function, we propose Riemannian Hamiltonian methods (RHM) and present their convergence analyses. We extend RHM to include consensus regularization and to the stochastic setting. We illustrate the efficacy of the proposed RHM in applications such as subspace robust Wasserstein distance, robust training of neural networks, and generative adversarial networks.

Key words. Riemannian optimization, saddle point, consensus optimization, Hamiltonian gradient descent, Polyak–Lojasiewicz, geodesic-bilinear, geodesic convex concave.

AMS subject classifications. 65K05, 90C30, 90C22, 90C25, 90C26, 90C27, 90C46, 58C05, 49M15

1. Introduction. In this paper, we consider the Riemannian manifold constrained min-max problem

$$(1.1) \quad \min_{x \in \mathcal{M}_x} \max_{y \in \mathcal{M}_y} f(x, y),$$

where $\mathcal{M}_x, \mathcal{M}_y$ are complete Riemannian manifolds and $f : \mathcal{M}_x \times \mathcal{M}_y \rightarrow \mathbb{R}$ is a jointly smooth real-valued function. The aim is to find a global saddle point (x^*, y^*) that satisfies for all $(x, y) \in \mathcal{M}_x \times \mathcal{M}_y$,

$$(1.2) \quad f(x^*, y) \leq f(x^*, y^*) \leq f(x, y^*).$$

Examples of Riemannian manifolds of interest include the sphere manifold, the Stiefel manifold, the manifold of orthogonal matrices, the manifold of doubly stochastic matrices, and the symmetric positive definite manifold, to name a few [3, 14, 80, 12].

When both $\mathcal{M}_x, \mathcal{M}_y$ are the Euclidean space, problem (1.1) reduces to the classical min-max problem, which has been widely studied for applications including adversarial training [49], robust learning [21], non-linear feature learning [71, 5, 36, 37], generative adversarial networks [24, 7, 79], constrained optimization [11], multi-task learning [35, 33], and fair statistical inference [48], among others. When f is convex in x and concave in y (convex-concave), the existence of a global saddle point is guaranteed by the well-established minimax theorem [62, 82]. Algorithms converging to such saddle points include the optimistic gradient descent ascent (OGDA) algorithm [70] and the extra-gradient algorithm (EG) [23], which have been analyzed in [61, 57, 58, 56]. For the general nonconvex-nonconcave setting, however, the saddle point, be it local or global, may not exist [39], and it remains challenging to establish convergence for both OGDA and EG.

*University of Sydney (andi.han@sydney.edu.au, junbin.gao@sydney.edu.au).

†Microsoft India (bamdevm@microsoft.com, pratik.jawanpuria@microsoft.com).

‡IIIT Hyderabad India (pawan.kumar@iiit.ac.in).

On Riemannian manifolds, there exist cases where many nonconvex (or nonconcave) functions turn out to be geodesic convex (or concave), a generalized notion of convexity on Riemannian manifolds [84]. This ensures the existence of a global saddle point on manifolds under the generalized min-max theorem [85, 92]. Furthermore, there is a growing interest in the Riemannian min-max problem (1.1) with applications such as low-rank tensor learning [34, 63], orthonormal generative adversarial networks [60, 17], subspace robust Wasserstein distances [65, 46, 30, 38], and adversarial neural network training [28]. It is, therefore, motivating to study the min-max problem on manifolds.

Nevertheless, existing works that systematically study the Riemannian min-max problem are sparse. In [28], a Riemannian gradient descent ascent (RGDA) method has been proposed, yet the analysis is restricted to \mathcal{M}_y being a convex subset of the Euclidean space and $f(x, y)$ being strongly concave in y . A recent paper [92] has formally characterized the optimality conditions of the Riemannian min-max problem for geodesic convex geodesic concave functions. A Riemannian corrected extra-gradient (RCEG) algorithm has been proposed and analyzed. A follow-up work [40] completes the analysis of RGDA and RCEG under geodesic (strongly) convex (strongly) concave settings.

Contributions. In this paper, we propose a class of methods for solving the min-max problem (1.1) on Riemannian manifolds, which we call Riemannian Hamiltonian methods (RHM). The idea is to minimize the squared norm of the Riemannian gradient of (1.1), known as the Riemannian Hamiltonian. Minimizing the Hamiltonian function serves as a good proxy for solving problem (1.1). Under the Riemannian Polyak–Lojasiewicz (PL) condition [91] on the Hamiltonian function, its minimizer recovers the desired saddle point. A key motivation to consider the proxy problem instead of the original min-max problem is for geodesic-bilinear problems, where solving the proxy problem leads to the correct direction towards global optimality while existing methods either cycle or converge extremely slowly (discussed in Section 3.3). In addition, the Hamiltonian gradient methods have been considered for solving min-max problems in the Euclidean space, which show great promise in accelerating and stabilizing the convergence to saddle points [1, 8, 53, 47]. This paper generalizes many of those analysis to Riemannian manifolds.

It should be emphasized that the proposed generalization to manifolds is nontrivial as the analysis for the Euclidean counterparts, such as in [1], rely heavily on the matrix properties of the Jacobian. Generalization to Riemannian manifolds require adherence to Riemannian operations independent of the matrix structure. Another challenge is to deal with the varying inner product (Riemannian metric) structure on manifolds. We handle the above by devising novel proof strategies and proposing a metric-aware Riemannian Hamiltonian function that respects the manifold geometry.

In particular, we show global linear convergence of any Riemannian solver to saddle points of problem (1.1) as long as the Riemannian Hamiltonian of f satisfies the Riemannian PL condition [91]. We show this occurs when f is geodesic strongly convex geodesic strongly concave, and also for some nonconvex functions with sufficient geodesic linearity. We additionally extend the proposed RHM to incorporate a consensus regularization and to the stochastic setting, and prove their convergence. Existing Riemannian algorithms for solving (1.1) such as [92] make use of the exponential map to update the iterates on the manifolds. In this work, we discuss convergence results with exponential as well as general retraction maps on manifolds.

We empirically show the convergence of our proposed RHM algorithms for dif-

ferent min-max functions and compare them with existing baselines. We further demonstrate the usefulness of RHM algorithms in various applications such as learning subspace robust Wasserstein distance, robust training of neural networks and training of generative adversarial networks.

Organizations. The rest of the paper is organized as follows. Section 2 reviews the preliminary knowledge on Riemannian geometry and Riemannian optimization as well as introduces various functions classes on Riemannian manifolds. We also briefly discuss the existing literature on mix-max optimization in the Euclidean space and on Riemannian manifolds. In Section 3, we propose the Riemannian Hamiltonian function and RHM algorithms, as well as analyze their convergence under the Riemannian PL condition. We provide three cases when such condition is satisfied. Section 4 introduces and analyzes the Riemannian Hamiltonian consensus method. Sections 5 and 6 extend the proposed methods to stochastic settings and to the case of retraction. In Section 7, we empirically compare our algorithms with different baselines on various applications. Section 8 concludes the paper.

2. Preliminaries. In this section, we give a brief overview of Riemannian geometry and relevant ingredients required for Riemannian optimization. For a more complete treatment of the topic, see [3, 14]. We also briefly discuss some of the existing works on min-max optimization.

2.1. Riemannian geometry and optimization.

Basic Riemannian geometry. Riemannian manifold \mathcal{M} is a manifold with a Riemannian metric, which is a smooth, symmetric positive definite function $g : T_p\mathcal{M} \times T_p\mathcal{M} \rightarrow \mathbb{R}$ on every tangent space $T_p\mathcal{M}$, with $p \in \mathcal{M}$. It is usually written as an inner product $\langle \cdot, \cdot \rangle_p$. The metric structure induces a norm for any tangent vector $\xi \in T_p\mathcal{M}$, which is $\|\xi\|_p := \sqrt{\langle \xi, \xi \rangle_p}$. For a linear operator on the tangent space $H : T_p\mathcal{M} \rightarrow T_p\mathcal{M}$, its operator norm is defined as $\|H\|_p := \max_{\xi \in T_p\mathcal{M} : \|\xi\|_p=1} \|H[\xi]\|_p$.

A geodesic on the manifold $\gamma : [0, 1] \rightarrow \mathcal{M}$ is the locally shortest curve with zero acceleration. The exponential map at p , $\text{Exp}_p : T_p\mathcal{M} \rightarrow \mathcal{M}$ is defined as the end point of a geodesic along the initial velocity. That is, $\text{Exp}_p(\xi) = \gamma(1)$ where $\gamma'(0) = \xi$, $\gamma(0) = p$ for any $\xi \in T_p\mathcal{M}$. Riemannian distance is computed as $d(p, q) = \int_0^1 \|\gamma'(t)\|_{\gamma(t)} dt$ where $\gamma(t)$ is the distance minimizing geodesic connecting $p, q \in \mathcal{M}$. In a totally normal neighbourhood Ω where there exists a unique geodesic between any $p, q \in \Omega$, the exponential map has a well-defined inverse $\text{Exp}_p^{-1} : \mathcal{M} \rightarrow T_p\mathcal{M}$ and the Riemannian distance can be written as $d(p, q) = \|\text{Exp}_p^{-1}(q)\|_p = \|\text{Exp}_q^{-1}(p)\|_q$. Parallel transport $\Gamma_p^q : T_p\mathcal{M} \rightarrow T_q\mathcal{M}$ transports tangent vector along the geodesic while being isometric, i.e., $\langle \xi, \zeta \rangle_p = \langle \Gamma_p^q \xi, \Gamma_p^q \zeta \rangle_q$ for any $\xi, \zeta \in T_p\mathcal{M}$.

Riemannian product manifolds. The product of Riemannian manifolds $\mathcal{M} = \mathcal{M}_x \times \mathcal{M}_y$ is a Riemannian manifold with the Riemannian metric defined as, for any $p = (x, y) \in \mathcal{M}$, and $(u, u'), (v, v') \in T_p\mathcal{M}$, $\langle (u, u'), (v, v') \rangle_p = \langle u, v \rangle_x^{\mathcal{M}_x} + \langle u', v' \rangle_y^{\mathcal{M}_y}$, where $\langle \cdot, \cdot \rangle_x^{\mathcal{M}_x}$, $\langle \cdot, \cdot \rangle_y^{\mathcal{M}_y}$ are Riemannian metrics on $\mathcal{M}_x, \mathcal{M}_y$ respectively. From the metric, one can derive the geodesic, the exponential map, parallel transport, Riemannian distance, which also admit a product structure. See more details in [14].

Riemannian optimization ingredients. Riemannian optimization treats the constrained problem as an unconstrained problem on manifold by generalizing the notions of gradient and Hessian. For a differentiable function $h : \mathcal{M} \rightarrow \mathbb{R}$, the Riemannian gradient at p , $\text{grad}h(p)$ is a tangent vector that satisfies $\langle \text{grad}h(p), \xi \rangle_p =$

$Dh(p)[\xi]$ for any $\xi \in T_p\mathcal{M}$ where $Dh(p)[\xi]$ is the directional derivative of h along ξ . The Riemannian Hessian of h , $\text{Hess}h(p) : T_p\mathcal{M} \rightarrow T_p\mathcal{M}$ is a symmetric linear operator, defined as the covariant derivative of the Riemannian gradient. For a bi-function $f : \mathcal{M}_x \times \mathcal{M}_y \rightarrow \mathbb{R}$, we can similarly define Riemannian partial gradient $\text{grad}_x f(x, y)$, $\text{grad}_y f(x, y)$ as Riemannian gradient for x, y , holding the other variable constant. The Riemannian cross derivative $\text{grad}_{xy}^2 f(x, y) : T_x\mathcal{M}_x \rightarrow T_y\mathcal{M}_y$ is defined as $\text{grad}_{xy}^2 f(x, y)[u] := D_x \text{grad}_y f(x, y)[u]$ and similarly for $\text{grad}_{yx}^2 f(x, y)$.

Riemannian geodesic convex optimization. Geodesic convexity [89, 14] generalizes the notion of convexity to Riemannian manifold. A *geodesic convex set* $\Omega \subseteq \mathcal{M}$ requires for any two points in the set, there exist a geodesic (on \mathcal{M}) connecting them that lies entirely in the set. From this definition, any connected, complete Riemannian manifold is geodesic convex itself. A function $h : \Omega \rightarrow \mathbb{R}$ is *geodesic convex* if for any $p, q \in \Omega$, it satisfies that $h(\gamma(t)) \leq (1-t)h(p) + th(q)$ for $t \in [0, 1]$ and γ is a geodesic connecting p, q . A function is *geodesic linear* if it is both geodesic convex and geodesic concave. A twice differentiable function h is *geodesic μ -strongly convex* if $\frac{d^2 h(\gamma(t))}{dt^2} \geq \mu$. We call a function $h(p)$ *g-(strongly)-convex* if it is geodesic (strongly) convex. Similarly, we call a function $f(x, y)$ *g-(strongly)-convex-concave* if it is geodesic (strongly) convex in x and geodesic (strongly) concave in y .

Next, we define the spectrum of a linear operator on the tangent space, which is used to analyze the Riemannian Hessian as well as the Riemannian cross derivatives in the subsequent sections.

DEFINITION 2.1 (Spectrum of a linear operator). *Consider a linear operator $T : V \rightarrow W$ where V, W are two inner product spaces. If $V = W$, and T is symmetric, i.e., $T = T^*$, where T^* is the adjoint operator of T , then we say (λ, v) is an eigenpair of T if $T[v] = \lambda v$. In general, when $V \neq W$, the singular value σ of T is the square root of the eigenvalues of $T^* \circ T$.*

We use $\lambda_{\min}/\lambda_{\max}$ and $\sigma_{\min}/\sigma_{\max}$ to represent the smallest/largest eigenvalues and singular values, respectively. We also use $\lambda_{|\min|}$ to denote the minimum eigenvalue in magnitude. Below, we introduce several function classes on manifolds, generalizing the Lipschitz continuity as well as the Polyak–Łojasiewicz condition from the Euclidean space [74, 69]

DEFINITION 2.2 (Lipschitz continuity [14]). *Let $L_0, L_1, L_2 > 0$.*

- (1). *A real-valued function $h : \mathcal{M} \rightarrow \mathbb{R}$ is L_0 -Lipschitz continuous if for all $p \in \mathcal{M}$, $\|\text{grad}h(p)\|_p \leq L_0$.*
- (2). *A vector field $V \in \mathfrak{X}(\mathcal{M})$ is L_1 -Lipschitz continuous if for all $p \in \mathcal{M}$ and $s \in T_p\mathcal{M}$ such that $q = \text{Exp}_p(s) \in \Omega$, a totally normal neighbourhood of p , it satisfies $\|\Gamma_q^p V(q) - V(p)\|_x \leq L_1 \|s\|_p$.*
- (3). *A linear operator $H(p) : T_p\mathcal{M} \rightarrow T_p\mathcal{M}$ is L_2 -Lipschitz continuous if for all $p \in \mathcal{M}$ and $q = \text{Exp}_p(s) \in \Omega$, it satisfies $\|\Gamma_q^p \circ H(q) \circ \Gamma_p^q - H(p)\|_p \leq L_2 \|s\|_p$.*

DEFINITION 2.3 (Polyak–Łojasiewicz (PL) condition on Riemannian manifold [91, 42, 25]). *A function $h : \mathcal{M} \rightarrow \mathbb{R}$ satisfies the PL condition on Riemannian manifold if for any $p \in \mathcal{M}$, there exists $\delta > 0$ such that $\frac{1}{2} \|\text{grad}h(p)\|_p^2 \geq \delta(h(p) - h(p^*))$, where $p^* = \arg \min_{p \in \mathcal{M}} h(p)$ is the global minimizer of h .*

The following lemma shows the connection between smoothness of a function on manifold and its Lipschitz Riemannian gradient, which is fundamental for convergence analysis.

LEMMA 2.4 (Lipschitz Riemannian gradient and smoothness [14]). *For a func-*

tion $h : \mathcal{M} \rightarrow \mathbb{R}$, its Riemannian gradient is L_1 -Lipschitz continuous if and only if $\|\text{Hess}h(p)\|_p \leq L_1$ for all $p \in \mathcal{M}$. Suppose h has L_1 -Lipschitz Riemannian gradient, then h is L_1 -smooth on \mathcal{M} with $|h(q) - h(p) - \langle \text{grad}h(p), s \rangle_p| \leq \frac{L_1}{2} \|s\|_p^2$, for all $q = \text{Exp}_p(s) \in \Omega$ and $p \in \mathcal{M}$.

Notations. Here, we summarize the main notations used in the paper. We use $\nabla, \nabla^2, \text{grad}$, and Hess to represent the Euclidean gradient, Euclidean Hessian, Riemannian gradient, and Riemannian Hessian respectively. The boldface ∇ is used to denote the Riemannian connection. For a bi-function $f(x, y)$, we denote $\nabla_x f(x, y), \nabla_y f(x, y)$ as the partial Euclidean derivative with respect to x, y , respectively, if $x, y \in \mathbb{R}^d$. Similarly for $x, y \in \mathcal{M}$, $\text{grad}_x f(x, y), \text{grad}_y f(x, y)$ denote the partial Riemannian gradients. We also make use of $\text{grad}_{xy}^2 f(x, y), \text{grad}_{yx}^2 f(x, y)$ to represent the Riemannian cross derivatives. We use $\langle \cdot, \cdot \rangle_p^{\mathcal{M}}$ to represent the Riemannian metric at $p \in \mathcal{M}$. When the manifold considered is clear, we omit the superscript for clarity. Furthermore, we use $\langle \cdot, \cdot \rangle_2$ to denote the Euclidean inner product.

2.2. Min-max optimization. Here we discuss related works on min-max optimization both in the Euclidean space and on Riemannian manifolds.

In Euclidean space. In the Euclidean space (i.e., \mathbb{R}^n), the standard gradient descent ascent (GDA) that follows the min-max gradient is known to cycle or diverge for simple convex-concave objectives [52]. To address the cycling issue, the optimistic gradient descent ascent algorithm (OGDA) [70] modifies the GDA update to include an additional gradient momentum. On the other hand, the extra-gradient algorithm (EG) [23] employs an additional min-max gradient step at every iteration. As shown in [55, 56], both OGDA and EG methods approximate the proximal point method [73] and converge sublinearly under convex-concave settings [61, 57] and linearly under strongly-convex-strongly-concave settings [87, 55].

However, for the more general nonconvex-nonconcave settings, finding a global saddle point satisfying (1.2) is difficult and several existing works [18, 4, 50, 79] aim to find a local saddle point that satisfies (1.2) in a local neighbourhood. It should be noted that when the function is convex-concave, all local saddle points are global.

A necessary set of conditions for the saddle points is that they satisfy the first-order stationarity, i.e., the gradients with respect to x and y vanish. This motivates the Euclidean Hamiltonian gradient descent (EHGD) [53, 8, 1, 47] approach for solving the min-max problem, which minimizes the sum of the squares of the gradient norms with respect to x and y . It should be noted that EHGD works under the assumption that all such stationary points are global min-max saddle points [1, 47]. Cases are discussed where this assumption is satisfied, which allows EHGD to converge to a global min-max saddle point of the original min-max problem [1, 47]. Further, studies [53, 8, 1, 47] demonstrate good empirical performance of EHGD in a variety of applications.

It should be noted that EHGD approaches have only been studied for unconstrained problems in the Euclidean space. Challenges in the constrained settings appear with definition of the Hamiltonian and subsequent analysis.

On Riemannian manifolds. There is a growing theoretical and empirical interest in solving min-max problems under Riemannian optimization framework [46, 30, 28, 92]. An extension of the GDA algorithm to manifolds, named RGDA, has been proposed in [28]. However, [28] considers a min-max setting in which the minimization problem (in x) is on a manifold, but the maximization problem (in y) is on a convex set. In addition, it analyzes the convergence when the maximization problem over y is strongly concave. Hence, [28] does not study the general Riemannian min-max

problem (1.1). It discusses the convergence of their algorithm to first-order stationary points of the min-max problem. Additionally, they propose different stochastic extensions of their algorithm and analyze their convergence.

Recently, [92] has proposed a Riemannian corrected extra-gradient algorithm (RCEG) for the Riemannian min-max problems (1.1), which contains two steps. First, RCEG takes a step similar to the RGDA update. Then, starting from the newly obtained point, RCEG combines the RGDA direction with the direction of the first step. In the g-convex-concave settings, this correction allows [92] to prove (local) convergence of RCEG to global min-max saddle points of (1.1). The convergence is however analyzed only for averaged iterates. After the submission of this work, we notice that a recent paper [40] proves both last-iterate and average-iterate convergence of RCEG to saddle points under g-convex-concave and g-strongly-convex-concave settings. They also discuss average-iterate and last-iterate convergence of RGDA under g-convex-concave and g-strongly-convex-concave settings, respectively. Nevertheless, the convergence analysis requires a bounded domain (and curvature) and a carefully chosen stepsize that depends on the curvature and diameter bound of the domain. In contrast, we have shown in this work global convergence to saddle points with stepsize that only depends on the Lipschitz constants of the objective.

More details on the RGDA and RCEG algorithms as well as the comparisons on the convergence analysis are in Appendix C.

3. Riemannian Hamiltonian gradient methods. As mentioned earlier, the Euclidean Hamiltonian approach [53, 8, 1, 47] is a popular approach to tackle the min-max problem (1.1) when \mathcal{M}_x and \mathcal{M}_y are restricted to the Euclidean space. Specifically, the Euclidean Hamiltonian function \mathcal{E} is defined as,

$$(3.1) \quad \mathcal{E}(x, y) := \frac{1}{2} \|\nabla_x f(x, y)\|_2^2 + \frac{1}{2} \|\nabla_y f(x, y)\|_2^2,$$

where $\nabla_x f(x, y)$ and $\nabla_y f(x, y)$ are the partial derivatives of f with respect to x and y , respectively. Here, $\|\cdot\|_2$ denotes the Frobenius norm. The global minimum of the function \mathcal{E} is attained when $\mathcal{E}(x, y) = 0$, i.e., $\nabla_x f(x, y) = 0$ and $\nabla_y f(x, y) = 0$. This corresponds to a first-order stationary point of the function f . Hence, minimization of \mathcal{E} in (3.1), becomes a good proxy to solve the original min-max problem.

Building on the Euclidean Hamiltonian approach, generalization to the Riemannian min-max problem (1.1) requires understanding of first-order stationary points on manifolds \mathcal{M}_x and \mathcal{M}_y . These are necessarily identified with the points where the Riemannian gradient of f vanishes. This leads to our proposed Riemannian Hamiltonian function as

$$(3.2) \quad \mathcal{H}(x, y) := \frac{1}{2} \|\text{grad}_x f(x, y)\|_x^2 + \frac{1}{2} \|\text{grad}_y f(x, y)\|_y^2,$$

where $\text{grad}_x f(x, y)$ and $\text{grad}_y f(x, y)$ are the Riemannian partial gradients of f with respect to x and y respectively. Here, $\|\text{grad}_x f(x, y)\|_x^2 = \langle \text{grad}_x f(x, y), \text{grad}_x f(x, y) \rangle_x^{\mathcal{M}_x}$ is the square of the gradient norm in the Riemannian metric sense on \mathcal{M}_x . Similarly, $\|\text{grad}_y f(x, y)\|_y^2 = \langle \text{grad}_y f(x, y), \text{grad}_y f(x, y) \rangle_y^{\mathcal{M}_y}$ is the square of the norm on \mathcal{M}_y .

Remark 3.1. The proposed Riemannian Hamiltonian function (3.2) generalizes the Euclidean Hamiltonian function (3.1) in two different ways:

- 1) Equation (3.2) implicitly embeds the manifold geometry of $\mathcal{M}_x, \mathcal{M}_y$ into the Hamiltonian function.

- 2) Equation (3.2) generalizes the Euclidean metric considered in (3.1) to a Riemannian metric. This generalization allows to use other varying metrics for min-max problems *in the Euclidean space*, e.g., the Fisher information metric [20] or real-projective space metrics [3, Chapter 2].

It should be noted that the Riemannian Hamiltonian (3.2) can be viewed on the product manifold $\mathcal{M} = \mathcal{M}_x \times \mathcal{M}_y$, i.e., for $p = (x, y) \in \mathcal{M}$, the Riemannian gradient is $\text{grad}_p f(p) = (\text{grad}_x f(x, y), \text{grad}_y f(x, y))$, and therefore, $\mathcal{H}(x, y) = \|\text{grad}_p f(p)\|_p^2$. Hence, we propose to solve the following problem on the product manifold as

$$(3.3) \quad \min_{p \in \mathcal{M}} \left\{ \mathcal{H}(p) = \frac{1}{2} \|\text{grad} f(p)\|_p^2 \right\}.$$

Similar to the EHGD approaches [1, 47], we work with the following assumption.

ASSUMPTION 1. *The objective f admits at least one stationary point and all stationary points are global min-max saddle points.*

It is worth noticing that under Assumption 1, solving (3.3) is equivalent to solving (1.1). On Riemannian manifolds, Assumption 1 holds when f is g-convex-concave.

We now show that the Riemannian gradient of the Riemannian Hamiltonian $\mathcal{H}(p)$ admits a simple expression.

PROPOSITION 3.2. *Riemannian gradient of \mathcal{H} is $\text{grad}\mathcal{H}(p) = \text{Hess}f(p)[\text{grad}f(p)]$.*

Proof. First, we see that \mathcal{H} is a smooth function on the manifold due to the smoothness of f and its Riemannian gradient (formally characterized later in Proposition 3.6). For any smooth vector field $U : \mathcal{M} \rightarrow T\mathcal{M}$, denoted as $U \in \mathfrak{X}(\mathcal{M})$, we have $U\mathcal{H} = \langle \text{grad}\mathcal{H}, U \rangle$, where $\langle \cdot, \cdot \rangle$ is the Riemannian metric (on any tangent space). Let ∇ be the Riemannian connection (or the Levi-Civita connection) of \mathcal{M} , which provides a way to differentiate vector fields on manifolds. By definition, the Riemannian connection satisfies the metric compatibility property [3, 14], i.e., $U\langle V, W \rangle = \langle \nabla_U V, W \rangle + \langle V, \nabla_U W \rangle$ for any vector fields U, V, W . Also, by definition, application of the Riemannian Hessian of $f : \mathcal{M} \rightarrow \mathbb{R}$ along a vector field U is $\text{Hess}f[U] = \nabla_U \text{grad}f$. Based on these claims, we show

$$\begin{aligned} U\mathcal{H} &= \frac{1}{2} U\langle \text{grad}f, \text{grad}f \rangle = \langle \nabla_U \text{grad}f, \text{grad}f \rangle = \langle \text{Hess}f[U], \text{grad}f \rangle \\ &= \langle \text{Hess}f[\text{grad}f], U \rangle, \end{aligned}$$

where the last equality follows from the self-adjoint property of the Riemannian Hessian. The proof is complete by noticing $\langle \text{Hess}f[\text{grad}f], U \rangle = \langle \text{grad}\mathcal{H}, U \rangle$ for any U . \square

Remark 3.3. The importance of the varying metric in the proposed Riemannian Hamiltonian (3.2), can be observed in Proposition 3.2, where we obtain a simple expression for the Riemannian gradient of \mathcal{H} . This allows to connect the properties of \mathcal{H} with that of the min-max objective f , discussed in detail later in Section 3.2.

Remark 3.4. It should be noted that for the Euclidean case when $x \in \mathbb{R}^m, y \in \mathbb{R}^n$, existing works [8, 1, 47] analyze the Hamiltonian methods in the form of $\mathbf{J}^\top v$, where \mathbf{J} is an asymmetric Jacobian matrix and v is the min-max gradient $(\nabla_x f(x, y), -\nabla_y f(x, y))$. For the same setting, however, Proposition 3.2 obtains the Hamiltonian gradient as $\mathbf{H}\nabla f$, where \mathbf{H} and ∇f are the (Euclidean) Hessian matrix and gradient vector $\nabla f = (\nabla_x f(x, y), \nabla_y f(x, y))$, respectively. This is not surprising as $\mathbf{J}^\top v = \mathbf{H}\nabla f$. Proposition 3.2 allows to analyze the performance of the Riemannian

Algorithm 3.1 Riemannian Hamiltonian methods (RHM)

- 1: Initialize $p_0 = (x_0, y_0) \in \mathcal{M}$.
 - 2: **for** $t = 0, \dots, T$ **do**
 - 3: Compute the step $\xi(p_t)$ from the gradient $\text{grad}\mathcal{H}(p_t) = \text{Hess}f(p_t)[\text{grad}f(p_t)]$.
 - 4: Update $p_{t+1} = \text{Exp}_{p_t}(\xi(p_t))$.
 - 5: **end for**
 - 6: **Output:** p_T .
-

Hamiltonian approach in terms of the symmetric Riemannian Hessian operator. The analysis in [1, 47] heavily rely on the matrix structure of \mathbf{J} and makes use of the linear algebraic properties of the Jacobian. Our approach, thanks to Proposition 3.2, adheres to general Riemannian manifolds as we directly deal with the operator, which is independent of the matrix structure. Hence, many of the subsequent analysis in this paper differ from [1, 47].

To minimize the Riemannian Hamiltonian (3.3), one can apply first-order Riemannian solvers including Riemannian steepest descent [88], Riemannian conjugate gradient [72], or second-order solvers, such as Riemannian trust-regions [2, 13], provided the Hessian (or approximated Hessian) of the Hamiltonian is available. We refer to such class of methods for solving min-max problems on manifolds collectively as Riemannian Hamiltonian methods (RHM). Its procedures are outlined in Algorithm 3.1, where the step $\xi(p_t)$ is computed depending on the selected solver.

Remark 3.5. We remark that Algorithm 3.1 aims to solve a proxy optimization problem (3.3) where we only require the first-order information, i.e., $\text{grad}\mathcal{H}(p)$. Although the Hessian of f , i.e., $\text{Hess}f(p)[\text{grad}\mathcal{H}(p)]$ is used in Algorithm 3.1, this essentially corresponds to the gradient information of the proxy problem. Furthermore, from the computational perspective, Algorithm 3.1 only requires one evaluation of Hessian-vector product per iteration. This is much more efficient than second-order methods, such as Riemannian trust region [2, 13] or cubic regularized Newton methods [6] that require at least several oracles to such Hessian vector product each iteration. Finally, when Hessian of f is unavailable, we find finite difference approximation is sufficient to achieve convergence in practice.

We analyze the performance of the proposed RHM. In particular, we aim to obtain the global minimizer p^* of \mathcal{H} , which satisfies $\mathcal{H}(p^*) = 0$ with RHM. However, this may not always be numerically tractable without additional structures on the Riemannian Hamiltonian. One such structure is assuming the Riemannian Hamiltonian is g -convex, for which RHM converges to the optimal p^* (g -convexity guarantees convergence to global optimality). This, however, may not lead to interesting problem classes for f . Moreover, there is no guarantee that \mathcal{H} is a g -convex even when f is g -convex-concave.

Another interesting structure is the Polyak–Łojasiewicz (PL) condition. The PL condition [69] amounts to a sufficient condition to establish linear convergence for gradient-based methods to global optimality [41]. The Riemannian version of the PL condition (Definition 2.3) has been studied in [91, 42, 93, 25]. In Section 3.1, we impose the Riemannian PL condition on the Hamiltonian \mathcal{H} as it allows convergence of RHM to global optimality. It should be noted that functions satisfying the Riemannian PL condition subsume g -(strongly)-convex functions. In Section 3.2, we discuss many interesting function classes of f that allow the Hamiltonian \mathcal{H} to satisfy the condition.

3.1. Convergence analysis. To analyze the convergence of RHM, we focus on the Riemannian steepest descent direction in the main text, i.e., $\xi(p_t) = -\eta_t \text{grad}\mathcal{H}(p_t)$ with either fixed stepsize or variable stepsize computed from backtracking line-search [15, 14]. We include the details of implementing the Riemannian conjugate gradient and Riemannian trust-region methods together with their convergence analysis in Appendix F. We make the following standard assumption [3, 14, 91, 78, 92] throughout the rest of the paper. We assume that our manifolds \mathcal{M}_x and \mathcal{M}_y are complete (and so is $\mathcal{M} = \mathcal{M}_x \times \mathcal{M}_y$).

ASSUMPTION 2. *The objective f , its Riemannian gradient, and its Riemannian Hessian are L_0, L_1, L_2 -Lipschitz continuous, respectively.*

In the next proposition, we show that the Riemannian Hamiltonian \mathcal{H} is L -smooth.

PROPOSITION 3.6 (Smoothness of Riemannian Hamiltonian). *Under Assumption 2, the Riemannian Hamiltonian is L -smooth with $L = L_0L_2 + L_1^2$, i.e., for any $p \in \mathcal{M}, q = \text{Exp}_p(\xi)$, it satisfies $\mathcal{H}(q) \leq \mathcal{H}(p) + \langle \text{grad}\mathcal{H}(p), \xi \rangle_p + \frac{L}{2} \|\xi\|_p^2$.*

Proof. According to Lemma 2.4, it is sufficient to show that the Riemannian gradient of \mathcal{H} is L -Lipschitz. From Proposition 3.2 and Assumption 2, we have for any $p \in \mathcal{M}, q = \text{Exp}_p(s) \in \Omega$, the domain of exponential map around p ,

$$\begin{aligned} \|\Gamma_p^q \text{grad}\mathcal{H}(p) - \text{grad}\mathcal{H}(q)\|_q &= \|\Gamma_p^q \text{Hess}f(p)[\text{grad}f(p)] - \text{Hess}f(q)[\text{grad}f(q)]\|_q \\ &\leq \|\Gamma_p^q \text{Hess}f(p)[\text{grad}f(p)] - \text{Hess}f(q)[\Gamma_p^q \text{grad}f(p)]\|_q \\ &\quad + \|\text{Hess}f(q)[\Gamma_p^q \text{grad}f(p)] - \text{Hess}f(q)[\text{grad}f(q)]\|_q \\ &= \|\text{Hess}f(p)[\text{grad}f(p)] - \Gamma_q^p \text{Hess}f(q)[\Gamma_p^q \text{grad}f(p)]\|_p \\ &\quad + \|\text{Hess}f(q)[\Gamma_p^q \text{grad}f(p)] - \text{Hess}f(q)[\text{grad}f(q)]\|_q \\ &\leq L_2 \|\text{grad}f(p)\|_p \|s\|_p + L_1 \|\Gamma_p^q \text{grad}f(p) - \text{grad}f(q)\|_q \\ &\leq (L_0L_2 + L_1^2) \|s\|_p, \end{aligned}$$

where we apply the triangle inequality and the isometry property of parallel transport. \square

If the Hamiltonian \mathcal{H} satisfies the Riemannian PL condition, then we show that Algorithm 3.1 with the steepest descent update (RHM-SD) converges linearly to the global minimizer of \mathcal{H} .

We begin with the convergence result for RHM-SD with fixed stepsize.

THEOREM 3.7 (Linear convergence of RHM-SD with fixed stepsize). *Let f satisfy Assumption 2 and \mathcal{H} satisfy the Riemannian PL condition, i.e., $\frac{1}{2} \|\text{grad}\mathcal{H}(p)\|_p^2 \geq \delta \mathcal{H}(p)$ (with $\mathcal{H}(p^*) = 0$). Consider Algorithm 3.1 using steepest descent direction with fixed stepsize $\eta_t = \eta = 1/L$, where $L = L_0L_2 + L_1^2$. Then, the iterates p_t satisfy $\|\text{grad}f(p_t)\|_{p_t}^2 \leq (1 - \frac{\delta}{L})^t \|\text{grad}f(p_0)\|_{p_0}^2$.*

Proof. From the smoothness of the Riemannian Hamiltonian \mathcal{H} (Proposition 3.6, Lemma 2.4) and the gradient update in Algorithm 3.1, we have

$$\begin{aligned} \mathcal{H}(p_{t+1}) - \mathcal{H}(p_t) &\leq -\eta \|\text{grad}\mathcal{H}(p_t)\|_{p_t}^2 + \frac{\eta^2 L}{2} \|\text{grad}\mathcal{H}(p_t)\|_{p_t}^2 \\ &= -\frac{1}{2L} \|\text{grad}\mathcal{H}(p_t)\|_{p_t}^2 \leq -\frac{\delta}{L} \mathcal{H}(p_t), \end{aligned}$$

where the last inequality employs the Riemannian PL condition. This leads to $\mathcal{H}(p_{t+1}) \leq (1 - \frac{\delta}{L}) \mathcal{H}(p_t)$. Applying this result recursively completes the proof. \square

Line-search methods are practically favourable because they adapt the stepsize without requiring the knowledge of the Lipschitz constant L . Here, we consider the backtracking line-search for choosing stepsize η_t for Riemannian steepest descent, which is commonly used in practice. Given an initial stepsize $\bar{\eta}$, the backtracking line-search iteratively decreases the stepsize by a factor of $\varrho \in (0, 1)$ until the Armijo-type sufficient decrease condition is satisfied, i.e.,

$$(3.4) \quad \mathcal{H}(p_t) - \mathcal{H}(\text{Exp}_{p_t}(\eta_t \zeta(p_t))) \geq r_1 \eta_t \langle -\text{grad}\mathcal{H}(p_t), \zeta(p_t) \rangle_{p_t},$$

for some update direction $\zeta(p_t)$. The complete procedure is included in Appendix B. We next present the convergence for RHM-SD with backtracking linesearch.

THEOREM 3.8 (Linear convergence of RHM-SD with backtracking line-search). *Under the same setting as in Theorem 3.7, consider Algorithm 3.1 using the steepest descent direction with backtracking line-search, parameters $r_1, \rho \in (0, 1)$, and an initial stepsize $\bar{\eta}$. Then, the iterates p_t satisfy*

$$\|\text{grad}f(p_t)\|_{p_t}^2 \leq \left(1 - 2 \min\left\{\bar{\eta}r_1, \frac{2\varrho(1-r_1)r_1}{L}\right\} \delta\right)^t \|\text{grad}f(p_0)\|_{p_0}^2.$$

Proof. Given \mathcal{H} is L -smooth, the proof follows from [14, Lemma 4.12] and the Riemannian PL condition. \square

Remark 3.9. It should be highlighted that the convergence rates to global saddle points obtained in Theorems 3.7, 3.8 are independent of the manifold curvature (which we achieve via solving a proxy Hamiltonian problem (3.3)). In contrast, the linear convergence rates shown in [92, 40] are curvature dependent.

3.2. Important problem classes for RHM. We now discuss the instances of f where the Riemannian Hamiltonian satisfies the Riemannian PL condition (Definition 2.3). This allows RHM (Algorithm 3.1) to converge to global min-max saddle points of (1.1).

From the expression of $\text{grad}\mathcal{H}(p)$ in Proposition 3.2, we observe that if all eigenvalues of $\text{Hess}f(p)$ are lower bounded in magnitude (i.e., $|\lambda| \geq \alpha > 0$), then the Riemannian Hamiltonian \mathcal{H} satisfies the Riemannian PL condition with $\delta = \alpha^2$. This is because

$$(3.5) \quad \underbrace{\frac{1}{2} \|\text{grad}\mathcal{H}(p)\|_p^2 \geq \alpha^2 \mathcal{H}(p)}_{\text{Riemannian PL condition}} \Leftrightarrow \underbrace{\frac{1}{2} \|\text{Hess}f(p)[\text{grad}f(p)]\|_p^2 \geq \frac{\alpha^2}{2} \|\text{grad}f(p)\|_p^2}_{\text{Required eigenvalue bound on Hess}f(p)}.$$

Our aim, therefore, is to identify classes of f that satisfy the right hand side of (3.5). We provide three cases where the Riemannian PL condition is naturally satisfied on the Riemannian Hamiltonian \mathcal{H} , which generalize the results in [1] to Riemannian manifolds. These include the cases when the objective f is g -strongly-convex-concave and when f is smooth with sufficient geodesic linearity.

In order to analyze function classes of f that lead to (3.5), we require the following results on the Riemannian Hessian $\text{Hess}f(p)$ of the product manifold \mathcal{M} (which are of independent interest as well).

- 1) Decomposition of the Riemannian Hessian $\text{Hess}f(p)$ and adjoint property of the cross derivatives. This is shown in Appendix D.
- 2) We establish general lower bounds on the eigenvalue magnitude of the Riemannian Hessian, which we include in Appendix E.

The above results help to bound the eigenvalues of $\text{Hess}f(p)$ in terms of the spectrum of $\text{Hess}_x f(x, y)$, $\text{Hess}_y f(x, y)$, and the cross derivatives $\text{grad}_{xy}^2 f(x, y)$, $\text{grad}_{yx}^2 f(x, y)$. We now present the main results below.

PROPOSITION 3.10 (Geodesic strongly convex strongly concave). *Let $f(x, y)$ be geodesic strongly convex in x and geodesic strongly concave in y with parameter $\mu > 0$. Then, \mathcal{H} satisfies the Riemannian PL condition (3.5) with $\delta = \mu^2$.*

Proof. We show that if there exists an eigenpair (λ, ξ) of $\text{Hess}f(p)$ such that $|\lambda| < \mu$ with $p = (x, y)$, $\xi = (u, v)$, then it leads to a contradiction. From the expression of the Riemannian Hessian in Proposition D.1, we have

$$\begin{aligned}\text{Hess}_x f(x, y)[u] + \text{grad}_{yx}^2 f(x, y)[v] &= \lambda u \\ \text{Hess}_y f(x, y)[v] + \text{grad}_{xy}^2 f(x, y)[u] &= \lambda v.\end{aligned}$$

This can be equivalently written as

$$(3.6) \quad \langle \text{Hess}_x f(x, y)[u], u \rangle_x + \langle \text{grad}_{yx}^2 f(x, y)[v], u \rangle_x = \lambda \|u\|_x^2$$

$$(3.7) \quad \langle \text{Hess}_y f(x, y)[v], v \rangle_y + \langle \text{grad}_{xy}^2 f(x, y)[u], v \rangle_y = \lambda \|v\|_y^2.$$

From (3.6), we obtain

$$(3.8) \quad \langle \text{grad}_{yx}^2 f(x, y)[v], u \rangle_x = -\langle u, (\text{Hess}_x f(x, y) - \lambda \text{id})[u] \rangle_x,$$

where id is the identity operator. From the symmetry of the Riemannian cross derivatives (Proposition D.2), we can substitute (3.8) into (3.7), which gives

$$(3.9) \quad \langle \text{Hess}_y f(x, y)[v], v \rangle_y - \langle u, (\text{Hess}_x f(x, y) - \lambda \text{id})[u] \rangle_x = \lambda \|v\|_y^2.$$

The geodesic strong convexity in x and geodesic strong concavity in y leads to $\text{Hess}_x f(x, y) \succeq \mu \text{id}$ and $\text{Hess}_y f(x, y) \preceq -\mu \text{id}$ respectively. Thus, the LHS of (3.9) is smaller than $-\mu$, which contradicts $|\lambda| < \mu$. Thus, all eigenvalues of $\text{Hess}f(p)$ satisfies $|\lambda| \geq \mu$. \square

PROPOSITION 3.11 (Smooth and geodesic linear). *Let $\sigma_{\min}(\text{grad}_{xy}^2 f(x, y)) \geq \tau > 0$ and let $f(x, y)$ be geodesic linear in one variable and has L_1 -Lipschitz Riemannian gradient in another variable. Then, \mathcal{H} satisfies the Riemannian PL condition (3.5) with $\delta = \frac{\tau^4}{2\tau^2 + L_1^2}$.*

Proof. Without loss of generality, we assume $f(x, y)$ has L_1 -Lipschitz gradient in x and geodesic linear in y . The geodesic linearity in y implies that $\text{Hess}_y f(x, y) = 0$, and therefore, we can apply Lemma E.1, which shows

$$\lambda_{|\min|}^2(\text{Hess}f(p)) \geq \frac{\sigma_{\min}^4(\text{grad}_{xy}^2 f(x, y))}{2\sigma_{\min}^2(\text{grad}_{xy}^2 f(x, y)) + \|\text{Hess}_x f(x, y)\|_x^2}.$$

Also, from Lemma 2.4, we have $\|\text{Hess}_x f(x, y)\|_x^2 \leq L_1^2$. Finally, from the assumption $\sigma_{\min}(\text{grad}_{xy}^2 f(x, y)) \geq \tau$, the proof is complete. \square

PROPOSITION 3.12 (Smooth and sufficiently geodesic-bilinear). *Let $0 < \tau \leq \sigma(\text{grad}_{xy}^2 f(x, y)) \leq \Upsilon$ and let $f(x, y)$ has L_1 -Lipschitz Riemannian gradient for both x and y . Define $\mu = \lambda_{|\min|}(\text{Hess}_x f(x, y))$, $\rho = \lambda_{|\min|}(\text{Hess}_y f(x, y))$ and let the sufficient geodesic-bilinearity condition holds: $(\tau^2 + \mu^2)(\tau^2 + \rho^2) - 4L_1^2\Upsilon^2 > 0$. Then, \mathcal{H} satisfies the Riemannian PL condition (3.5) with $\delta = \frac{(\tau^2 + \mu^2)(\tau^2 + \rho^2) - 4L_1^2\Upsilon^2}{2\tau^2 + \rho^2 + \mu^2}$.*

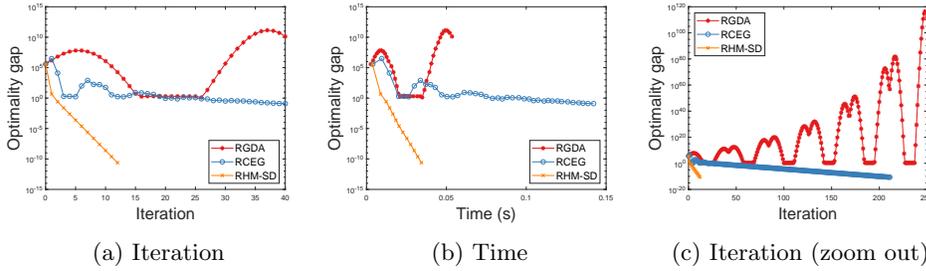


Fig. 1: RGDA [28] fails to converge on the geodesic bilinear problem $f(\mathbf{X}, \mathbf{Y}) = \log \det(\mathbf{X}) \log \det(\mathbf{Y})$. In particular, RGDA suffers from cyclic behaviour. RCEG [92, 40] converges very slowly. In contrast to RGDA and RCEG, the proposed Riemannian gradient descent method on the Riemannian Hamiltonian function (RHM-SD) quickly converges on such challenging bilinear problems. Notably, the proposed RHM-SD achieves an optimality gap lower than 10^{-10} in just 12 iterations while RCEG takes 220 iterations.

Proof. We can directly apply Lemma E.2 and set $a = 2\tau^2 + \rho^2 + \mu^2$ and $b = (\tau^2 + \mu^2)(\tau^2 + \rho^2) - 4L_1^2\Upsilon^2 > 0$ by assumption. \square

It is worth noticing that the sufficient geodesic-bilinearity condition in Proposition 3.12 can be interpreted as requiring a sufficiently large weight on the geodesic-bilinear component in the objective function f . To see this, suppose $f(x, y) = c_1 f_0(x, y) + f_1(x) + f_2(y)$ where f_0 is geodesic linear in each x and y (i.e. bilinear) with the weight $c_1 > 0$ and f_1, f_2 have L_1 -Lipschitz Riemannian gradient. Because by definition, Riemannian Hessian of a geodesic linear function is zero, f has $2L_1$ -Lipschitz Riemannian gradient (by Lemma 2.4). Let τ_0, Υ_0 be the minimum and maximum singular values of $\text{grad}_{xy}^2 f_0(x, y)$. Then, $\tau = c_1 \tau_0, \Upsilon = c_1 \Upsilon_0$. The sufficient geodesic bilinearity condition is satisfied for $c_1 \geq 4L_1 \Upsilon_0 / \tau_0^2$. This is because $(\tau^2 + \mu^2)(\tau^2 + \rho^2) > \tau^4 = c_1^4 \tau_0^4 \geq 16L_1^2 c_1^2 \Upsilon_0^2 = 16L_1^2 \Upsilon^2$.

Remark 3.13. When $f_1(x) = f_2(y) = 0$, it should be noted that $f(x, y) = c_1 f_0(x, y)$ is geodesic bilinear. Additionally, \mathcal{H} satisfies the Riemannian PL condition with $\delta = \frac{c_1^2 \tau_0^2}{2}$.

3.3. The geodesic-bilinear example. Here, we give a motivating example to show how the Riemannian Hamiltonian approach achieves convergence to global saddle points. To this end, we consider the problem $f(\mathbf{X}, \mathbf{Y}) = \log \det(\mathbf{X}) \log \det(\mathbf{Y})$ where $\mathbf{X}, \mathbf{Y} \in \mathbb{S}_{++}^d$, the set of $d \times d$ symmetric positive definite (SPD) matrices. When endowed with the affine-invariant metric [12], i.e., $\langle \mathbf{U}, \mathbf{V} \rangle_{\mathbf{X}} = \text{tr}(\mathbf{X}^{-1} \mathbf{U} \mathbf{X}^{-1} \mathbf{V})$ for any $\mathbf{U}, \mathbf{V} \in T_{\mathbf{X}} \mathbb{S}_{++}^d$, the set becomes a Riemannian manifold. Under this metric, one can show that the function is g-bilinear, i.e., geodesic linear in both \mathbf{X}, \mathbf{Y} , but not g-strongly-convex-concave (Proposition 7.1). However, the Riemannian Hamiltonian of the objective, i.e., $\mathcal{H}(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} (\|\text{grad}_{\mathbf{X}} f(\mathbf{X}, \mathbf{Y})\|_{\mathbf{X}}^2 + \|\text{grad}_{\mathbf{Y}} f(\mathbf{X}, \mathbf{Y})\|_{\mathbf{Y}}^2)$ satisfies the Riemannian PL condition (Proposition 7.2). This suggests that the vanilla Riemannian gradient descent method for minimizing the Riemannian Hamiltonian (3.3) converges to the global saddle points of f . In Appendix G, we show that the geodesic-bilinear function $f(\mathbf{X}, \mathbf{Y})$ does not satisfy the min-max Riemannian PL condition on

the original problem. This further justifies the merit of the proposed Hamiltonian proxy problem (3.3) of $f(\mathbf{X}, \mathbf{Y})$ that satisfies the Riemannian PL condition.

On the other hand, the RGDA algorithm [28] follows the negative of the min-max Riemannian gradient of f . Specifically, let $\mathcal{M} = \mathbb{S}_{++}^d \times \mathbb{S}_{++}^d$ and $P = (\mathbf{X}, \mathbf{Y}) \in \mathcal{M}$ be the product manifold and its elements. The min-max gradient of f is derived as $G(P) = (\text{grad}_{\mathbf{X}} f(\mathbf{X}, \mathbf{Y}), -\text{grad}_{\mathbf{Y}} f(\mathbf{X}, \mathbf{Y})) = (\mathbf{X} \log \det(\mathbf{Y}), -\mathbf{Y} \log \det(\mathbf{X}))$. We compare this expression with the gradient of the Riemannian Hamiltonian, which is $\text{grad}\mathcal{H}(P) = (\mathbf{X} \log \det(\mathbf{X}), \mathbf{Y} \log \det(\mathbf{Y}))$. We observe that $\langle G(P), \text{grad}\mathcal{H}(P) \rangle_P = 0$, which implies that the min-max gradient of f is always orthogonal to the gradient of its Riemannian Hamiltonian. In fact, such orthogonality holds for any g-bilinear objective (see Proposition G.1 in Appendix G). Given that the negative Hamiltonian gradient $-\text{grad}\mathcal{H}(P)$ points to the global saddle points, the orthogonality of its direction implies that RGDA provably cycles around the saddle points. For the RCEG algorithm [92, 40], since it also makes use of the RGDA-style updates, we expect its slow convergence for g-bilinear problems.

We illustrate the above findings in Figure 1, where we compare the Riemannian Hamiltonian steepest descent method (RHM-SD) against both RGDA [28] and RCEG [92, 40] with properly tuned stepsize. The convergence is measured in optimality gap, i.e., $|\det(\mathbf{X}) - 1| + |\det(\mathbf{Y}) - 1|$ given that the global saddle points of f satisfy $\det(\mathbf{X}^*) = \det(\mathbf{Y}^*) = 1$ (Proposition 7.2). We observe that the proposed RHM-SD takes only 12 iterations to obtain an optimality gap of lower than 10^{-10} for this challenging setup. However, RGDA experiences cyclic behaviour, which matches our analysis above. While RCEG incorporates a correction step for RGDA-style updates to address the cycling issue, it still exhibits a slight cyclic behavior during the initial phase but eventually converges (Figure 1(c)). Overall, RCEG takes 220 iterations for the same optimality gap. From Figures 1(a) and 1(b), we also observe that per iteration runtime cost of RHM-SD is similar to RCEG.

More details and discussions can be found in Section 7.1, where we generalize the findings to include quadratic terms.

4. Riemannian Hamiltonian consensus method. In the Euclidean setting, [53] proposes the consensus method for solving min-max problems in the Euclidean space. The consensus method has also been viewed as a perturbation of the Euclidean Hamiltonian method [1]. In this section, we propose an extension of RHM with steepest descent update, namely the Riemannian Hamiltonian consensus method (RHM-CON), by combining the Hamiltonian gradient direction with the min-max gradient direction. In practice, particularly for some deep learning applications, Assumption 1 may not satisfy. Thus solving the Hamiltonian proxy problem (3.3) may lead to undesired stationary points that are not saddle points. The consensus direction provides a regularization and is usually practically favourable for general nonconvex-nonconcave min-max problem. We show such an example in Section 7.5.

The update of RHM-CON is given by

$$p_{t+1} = \text{Exp}_{p_t}(-\eta_t \zeta(p_t)) = \text{Exp}_{p_t}\left(-\eta_t (\gamma v(p_t) + \text{grad}\mathcal{H}(p_t))\right),$$

with $\gamma \geq 0$ and $v(p_t) := (\text{grad}_x f(x_t, y_t), -\text{grad}_y f(x_t, y_t))$ is the min-max gradient. When $\gamma = 0$, this reduces to RHM-SD. The RHM-CON method is formalized in Algorithm 4.1. Below, we provide the convergence result for RHM-CON.

THEOREM 4.1 (Linear convergence of RHM-CON). *Under Assumption 2 with $L = L_0 L_2 + L_1^2$, suppose that the Riemannian Hamiltonian \mathcal{H} satisfies the PL condi-*

Algorithm 4.1 Riemannian Hamiltonian consensus (RHM-CON) method

-
- 1: **Input:** Stepsize η and regularization parameter γ .
 - 2: Initialize $p_0 = (x_0, y_0) \in \mathcal{M}$.
 - 3: **for** $t = 0, \dots, T$ **do**
 - 4: Compute the min-max gradient $v(p_t) = (\text{grad}_x f(x_t, y_t), -\text{grad}_y f(x_t, y_t))$.
 - 5: Compute the update direction $\zeta(p_t) = \gamma v(p_t) + \text{Hess}f(p_t)[\text{grad}f(p_t)]$.
 - 6: Update $p_{t+1} = \text{Exp}_{p_t}(-\eta_t \zeta(p_t))$.
 - 7: **end for**
 - 8: **Output:** p_T .
-

tion. Let $c > 0$ such that $\|\zeta(p_t)\|^2 = \|\gamma v(p_t) + \text{grad}\mathcal{H}(p_t)\|^2 \geq c\|\text{grad}\mathcal{H}(p_t)\|^2$ for all the iterates p_t . Set $\gamma < \sqrt{\delta}$, $\eta_t = \eta \leq \frac{1}{L}$, then Algorithm 4.1 converges with

$$\|\text{grad}f(p_t)\|_{p_t}^2 \leq (1 - \nu)^t \|\text{grad}f(p_0)\|_{p_0}^2,$$

where $\nu = (c\delta + \delta - \gamma^2)\eta - Lc\delta\eta^2 > 0$.

Proof. First, we highlight that

$$\frac{1}{2}\|v(p)\|_p^2 = \frac{1}{2}\|\text{grad}f(p)\|_p^2 = \mathcal{H}(p).$$

From the smoothness of Riemannian Hamiltonian (Proposition 3.6, Lemma 2.4) and the update in Algorithm 4.1, we have

$$\begin{aligned} & \mathcal{H}(p_{t+1}) - \mathcal{H}(p_t) \\ & \leq -\eta \langle \text{grad}\mathcal{H}(p_t), \zeta(p_t) \rangle_{p_t} + \frac{\eta^2 L}{2} \|\zeta(p_t)\|_{p_t}^2 \\ & = -\frac{\eta}{2} \|\text{grad}\mathcal{H}(p_t)\|_{p_t}^2 + \frac{\eta}{2} \|\zeta(p_t) - \text{grad}\mathcal{H}(p_t)\|_{p_t}^2 - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2}\right) \|\zeta(p_t)\|_{p_t}^2 \\ & \leq \left(-\frac{\eta}{2} - \frac{\eta c}{2} + \frac{\eta^2 L c}{2}\right) \|\text{grad}\mathcal{H}(p_t)\|_{p_t}^2 + \frac{\eta \gamma^2}{2} \|v(p_t)\|_{p_t}^2 \\ & \leq (-\eta - \eta c + \eta^2 L c) \delta \mathcal{H}(p_t) + \eta \gamma^2 \mathcal{H}(p_t) \\ & = (Lc\delta\eta^2 - c\delta\eta - \delta\eta + \eta\gamma^2) \mathcal{H}(p_t), \end{aligned}$$

where the second inequality follows from $\eta \leq \frac{1}{L}$ (which gives $\frac{\eta}{2} - \frac{\eta^2 L}{2} \geq 0$) and the lower bound on $\|\zeta(p_t)\|_{p_t}^2$. The last inequality uses the PL condition and $\eta \leq \frac{1}{L} < \frac{1+c}{Lc}$, which ensures $-\frac{\eta}{2} - \frac{\eta c}{2} + \frac{\eta^2 L c}{2} < 0$. From the choice of η and γ as well as the definition of ν , we have $\nu > 0$. This is because $\nu = \eta(c\delta + \delta - \gamma^2 - Lc\delta\eta) \geq \eta(\delta - \gamma^2) > 0$. Thus, $\mathcal{H}(p_{t+1}) = (1 - \nu)\mathcal{H}(p_t)$ ensuring linear convergence. Applying this result recursively completes the proof. \square

From Theorem 4.1, we see that linear convergence is achieved provided that the weight γ on min-max gradient direction is sufficiently small. Also, we highlight that a uniform parameter $c > 0$ always exists in a compact set as long as $\gamma v(p_t) \neq -\text{grad}\mathcal{H}(p_t)$. This can be ensured by choosing a small value for γ .

5. Stochastic min-max optimization. Applications such as domain generalization, robust training, and generative adversarial networks yield a min-max problem

with a stochastic function f , e.g., with a finite sum structure of the function [47]. Under the stochastic setting, the objective function in (1.1) can be expressed as an expectation, i.e.,

$$\min_{x \in \mathcal{M}_x} \max_{y \in \mathcal{M}_y} \left\{ f(x, y) = \mathbb{E}_\omega [f(x, y; \omega)] \right\},$$

where $\omega \in \mathcal{D}$ is a random variable following a certain distribution \mathcal{D} . This implies an expectation structure on the Riemannian Hamiltonian as

$$\mathcal{H}(p) = \frac{1}{2} \left\| \mathbb{E}_\omega [\text{grad} f(p; \omega)] \right\|_p^2 = \frac{1}{2} \mathbb{E}_\omega \mathbb{E}_\varphi \langle \text{grad} f(p; \omega), \text{grad} f(p; \varphi) \rangle_p,$$

for $\omega, \varphi \in \mathcal{D}$. Modifying Proposition 3.2 for the stochastic setting leads to

$$\text{grad} \mathcal{H}(p) = \frac{1}{2} \mathbb{E}_{\omega, \varphi} \left[\text{Hess} f(p; \omega) [\text{grad} f(p; \varphi)] + \text{Hess} f(p; \varphi) [\text{grad} f(p; \omega)] \right].$$

Let $\text{grad} \mathcal{H}_{\omega, \varphi}(p) := \frac{1}{2} \text{Hess} f(p; \omega) [\text{grad} f(p; \varphi)] + \frac{1}{2} \text{Hess} f(p; \varphi) [\text{grad} f(p; \omega)]$. We can modify RHM-SD by replacing the gradient of Hamiltonian with its stochastic version (which we call RHM-SGD) as

$$(5.1) \quad \text{grad} \mathcal{H}_{\mathcal{S}, \mathcal{S}'}(p_t) := \frac{1}{|\mathcal{S}| |\mathcal{S}'|} \sum_{\omega \in \mathcal{S}, \varphi \in \mathcal{S}'} \text{grad} \mathcal{H}_{\omega, \varphi}(p),$$

where $\mathcal{S} = \{\omega_1, \dots, \omega_{|\mathcal{S}|}\}$, $\mathcal{S}' = \{\varphi_1, \dots, \varphi_{|\mathcal{S}'|}\}$ are randomly selected subsets with $\omega_i, \varphi_j \in \mathcal{D}$. The stochastic Hamiltonian gradient provides an unbiased estimate of the full gradient, i.e., $\mathbb{E}_{\mathcal{S}, \mathcal{S}'} [\text{grad} \mathcal{H}_{\mathcal{S}, \mathcal{S}'}(p)] = \text{grad} \mathcal{H}(p)$. We now show the convergence result of RHM-SGD.

THEOREM 5.1 (Convergence of RHM-SGD with fixed and decaying stepsize).

Let Assumption 2 hold with $L = L_0 L_2 + L_1^2$, and let the Riemannian Hamiltonian \mathcal{H} satisfy the PL condition with parameter δ . Assume also that the variance of the stochastic gradient is bounded, i.e., $\mathbb{E}_{\omega, \varphi} \|\text{grad} \mathcal{H}_{\omega, \varphi}(p_t)\|_{p_t}^2 \leq G$. Then, RHM-SGD with fixed stepsize $\eta_t = \eta < \frac{1}{2\delta}$ converges with $\mathbb{E} \|\text{grad} f(p_t)\|_{p_t}^2 \leq (1 - 2\eta\delta)^t \mathbb{E} \|\text{grad} \mathcal{H}(p_0)\|_{p_0}^2 + \frac{\eta L G}{4}$. Also, RHM-SGD with decaying stepsize $\eta_t = \frac{2t+1}{2\delta(t+1)^2}$, converges with $\mathbb{E} \|\text{grad} f(p_t)\|_{p_t}^2 \leq \frac{L G}{2\delta^2 t}$.

Proof. The proof follows from [41, Theorem 4] and can be easily adapted to the Riemannian manifold setting, and therefore, is omitted. \square

We can similarly consider the stochastic version of RHM-CON, which we denote as RHM-SCON, with the update step as

$$\zeta_{\mathcal{S}, \mathcal{S}'}(p_t) = \gamma(v_{\mathcal{S}}(p_t) + v_{\mathcal{S}'}(p_t))/2 + \text{grad} \mathcal{H}_{\mathcal{S}, \mathcal{S}'}(p_t),$$

where $v_{\mathcal{S}}(p_t)$ is the stochastic min-max gradient on sample set \mathcal{S} . Theorem 5.1 can be adapted to prove the convergence of RHM-SCON following similar assumptions and analysis in Theorem 4.1.

6. Convergence under retraction. Existing algorithms for solving (1.1), such as RCEG [92], employs the exponential map to update iterates on the manifolds. However, in many cases, the computational cost of implementing the exponential map for many Riemannian manifolds is prohibitive. An alternative is to consider the

more general retraction operation [3, Chapter 4]. In this section, we show that the use of retraction (instead of the exponential map) in RHM algorithms guarantees similar convergence under an additional mild assumption.

Retraction $R_p : T_p\mathcal{M} \rightarrow \mathcal{M}$ is a map that satisfies for all $p \in \mathcal{M}$, (1) $R_p(0) = p$ (2) $DR_p(0)[\xi] = \xi$ for all $\xi \in T_p\mathcal{M}$. From the definition, we observe that the exponential map is a special case of retraction. In practice, when an efficient retraction is available, the Hamiltonian gradient update can be performed via retraction, i.e., $p_{t+1} = R_{p_t}(-\eta \text{grad}\mathcal{H}(p_t))$. To analyze the convergence, we make the following additional assumption that bound the differential operator of the retraction map.

ASSUMPTION 3. *There exists constants $\theta_1, \theta_2 > 0$ such that the retraction curve $c(t) := R_p(t\xi)$ with $\|\xi\|_p = 1$ satisfies $\|c'(t)\|_{c(t)} \leq \theta_1$ and $\|c''(t)\|_{c(t)} \leq \theta_2$ for all t where $c(t) \in \mathcal{U}$, where \mathcal{U} is a compact subset of \mathcal{M} .*

This assumption is always satisfied for a compact manifold \mathcal{M} . The compactness appears to be necessary for retraction-based analysis for first-order algorithms [25, 78, 42, 14]. We remark that for the case of the exponential map, the retraction curve coincides with the geodesic curve. Then, $\theta_1 = 1$ because $\|c'(t)\|_{c(t)} = \|\Gamma_p^{c(t)}\xi\|_{c(t)} = 1$ by isometric property of parallel transport. Also, $\theta_2 = 0$ from the definition of the geodesic.

PROPOSITION 6.1. *Under Assumptions 2 and 3, the Riemannian Hamiltonian \mathcal{H} is retraction L_R -smooth with $L_R = \theta_1^2 L + \theta_2 L_1 L_0$, i.e., for any $p \in \mathcal{M}$, $q = R_p(\xi) \in \mathcal{U}$, we have $\mathcal{H}(q) \leq \mathcal{H}(p) + \langle \text{grad}\mathcal{H}(p), \xi \rangle_p + \frac{L_R}{2} \|\xi\|_p^2$.*

Proof. For any retraction curve $c(t) = R_p(t\xi)$ with $\|\xi\|_p = 1$ and $t \geq 0$ such that $c(t) \in \mathcal{U}$, we obtain

$$\begin{aligned} \frac{d^2}{dt^2} \mathcal{H}(c(t)) &= \langle \text{Hess}\mathcal{H}(c(t))[c'(t)], c'(t) \rangle_{c(t)} + \langle \text{grad}\mathcal{H}(c(t)), c''(t) \rangle_{c(t)} \\ &\leq L\theta_1^2 + \theta_2 \|\text{Hess}f(c(t))[\text{grad}f(c(t))]\|_{c(t)} \\ (6.1) \quad &\leq L\theta_1^2 + \theta_2 L_1 L_0 = L_R, \end{aligned}$$

where the second inequality applies the gradient of Hamiltonian is L -Lipschitz (Proposition 3.6, Lemma 2.4) and Assumption 3. The last inequality follows from Assumption 2. The proof from (6.1) to L_R -smoothness of \mathcal{H} is due to [31, Lemma 3.2], which we include here for completeness.

For any $\xi \in T_p\mathcal{M}$ such that $R_p(\xi) \in \mathcal{U}$, let $\alpha = \|\xi\|_p$, $\zeta = \xi/\|\xi\|_p$ and hence $\xi = \alpha\zeta$ with $\|\zeta\|_p = 1$. Applying Taylor's Theorem on $\mathcal{H} \circ R_p$ gives

$$\begin{aligned} \mathcal{H}(R_p(\xi)) - \mathcal{H}(p) &= \mathcal{H}(R_p(\alpha\zeta)) - \mathcal{H}(R_p(0)) \\ &= \alpha \frac{d}{dt} \mathcal{H}(R_p(t\zeta)) \Big|_{t=0} + \frac{\alpha^2}{2} \frac{d^2}{dt^2} \mathcal{H}(R_p(t\zeta)) \Big|_{t=\tilde{t}} \\ &\leq \alpha \langle \text{grad}\mathcal{H}(p), \zeta \rangle_p + \frac{\alpha^2 L_R}{2} \\ &= \langle \text{grad}\mathcal{H}(p), \xi \rangle + \frac{L_R}{2} \|\xi\|_p^2, \end{aligned}$$

where $\tilde{t} \in [0, \alpha]$. Thus, the proof is complete. \square

Using Proposition 6.1, we show below that RHM-SD attains a linear convergence rate with retraction.

THEOREM 6.2 (Linear convergence of RHM-SD under retraction). *Under same settings as in Theorem 3.7, suppose Assumption 3 holds, and the iterates stay in the compact set \mathcal{U} . Then, RHM-SD with retraction and $\eta = 1/L_R$ converges with $\|\text{grad}f(p_t)\|_{p_t}^2 \leq (1 - \frac{\delta}{L_R})^t \|\text{grad}f(p_0)\|_{p_0}^2$.*

The proof is similar to the proof of Theorem 3.7 and is omitted. A similar analysis with the retraction operation can be performed for other variants of RHM including RHM-CON, RHM-SGD, and RHM-SCON.

7. Experiments. In this section, we discuss empirical performance of the proposed Riemannian Hamiltonian methods for various min-max optimization problems on manifolds. The algorithms are implemented in Matlab using the Manopt package [16] except for Section 7.4, 7.5 where we use Pytorch with the Geoopt package [43]. We highlight that there exist many other manifold optimization packages, such as ROPTLIB [32], Manopt.jl [10], Pymnapt [86], McTorch [51], and RiemOpt [83], where RHM can also be implemented efficiently. We use the following acronyms for the various RHM algorithms considered in this section.

- RHM-SD-F: RHM with steepest descent direction with fixed stepsize.
- RHM-SD: RHM with steepest descent direction with backtracking line search.
- RHM-CON: RHM consensus method with fixed stepsize (Section 4).
- RHM-CG: RHM with the conjugate gradient method.
- RHM-TR: RHM with the trust-region method where we use Hessian approximation with finite differentiation [13].
- RHM-SGD: RHM with stochastic gradient (Section 5).
- RHM-SCON: RHM with stochastic consensus method (Section 5).

We compare the proposed Riemannian Hamiltonian methods with the Riemannian gradient descent ascent (RGDA) [28] and the Riemannian corrected extra-gradient (RCEG) [92]. As discussed previously, RGDA has not been studied and analyzed for solving the general min-max problem (1.1), but when \mathcal{M}_y is a convex subset of the Euclidean space [28]. In our experiments, however, we extend RGDA to solve (1.1).

For all the experiments, we implement the algorithms with exponential map for comparability with RCEG, except for the applications of subspace robust Wasserstein distance (Section 7.3), robust training (Section 7.4) and generative adversarial networks (Section 7.5) where we implement with retraction map because the manifolds considered do not have a well-defined logarithm map. Hence, for these applications, RCEG is excluded for comparison. In robust training and generative adversarial network experiments, we also test stochastic algorithms for RGDA and RHM. The codes are available at <https://github.com/andyjm3>.

7.1. Geodesic quadratic bilinear optimization. The first example we consider is

$$(7.1) \quad f(\mathbf{X}, \mathbf{Y}) = c_q(\log \det(\mathbf{X}))^2 + c_l \log \det(\mathbf{X}) \log \det(\mathbf{Y}) - c_q(\log \det(\mathbf{Y}))^2,$$

where $\mathbf{X}, \mathbf{Y} \in \mathbb{S}_{++}^d$, the set of $d \times d$ symmetric positive definite (SPD) matrices. The weights $c_l, c_q \geq 0$ control the balance between the linear and quadratic terms.

For $\mathbf{X} \in \mathbb{S}_{++}^d$, the tangent space $T_{\mathbf{X}}\mathbb{S}_{++}^d$ is the set of symmetric matrices. When endowed with the affine-invariant (AI) metric, i.e., $\langle \mathbf{U}, \mathbf{V} \rangle_{\mathbf{X}} = \text{tr}(\mathbf{X}^{-1}\mathbf{U}\mathbf{X}^{-1}\mathbf{V})$, for any $\mathbf{U}, \mathbf{V} \in T_{\mathbf{X}}\mathbb{S}_{++}^d$, one can derive the geodesic, exponential map, and other Riemannian optimization ingredients [26, 12, 67]. We include the expressions in Appendix A. Here, we use \mathcal{M}_{SPD} to represent the SPD manifold with the AI metric. It is worth noticing that the function (7.1) is nonconvex-nonconcave in the Euclidean space (with details included in Appendix G).

However, the log-det function is geodesic linear on SPD manifold with the AI metric [84] and we show in the following proposition that $f(\mathbf{X}, \mathbf{Y})$ is g-convex-concave, although not necessarily g-strongly-convex-concave.

PROPOSITION 7.1. *The function (7.1) is g-convex-concave on \mathcal{M}_{SPD} but not g-strongly-convex-concave.*

We next prove that the Riemannian Hamiltonian \mathcal{H} of the objective (7.1) satisfies the PL condition, which allows linear convergence of the proposed RHM algorithms.

PROPOSITION 7.2. *The Riemannian Hamiltonian of (7.1) satisfies the PL condition with $\delta = (4c_q^2 + c_l^2)d^2$. A point $(\mathbf{X}^*, \mathbf{Y}^*)$ is a global saddle point of (7.1) if and only if it satisfies $\det(\mathbf{X}^*) = \det(\mathbf{Y}^*) = 1$.*

In Proposition 7.2, we see that there exist a continuum of global saddle points. Consequently, we define an optimality gap criterion as $|\det(\mathbf{X}) - 1| + |\det(\mathbf{Y}) - 1|$ for a candidate point (\mathbf{X}, \mathbf{Y}) .

Experiment settings and results. We consider $d = 30$ and discuss results on various combinations of c_q, c_l . We compare our RHM with RGDA [28] and RCEG [92]. All the choices of stepsize are tuned to reflect the best performance except for RHM-SD, RHM-CG, RHM-TR where the stepsizes are selected adaptively by the algorithms. For RHM-CON, we set $\gamma = 0.5$. Convergence of an algorithm is measured in terms of $\|\text{grad}f(p_t)\|_{p_t}$, which is equivalent to $\sqrt{2\mathcal{H}(p_t)}$. This measure of convergence has also been considered in [92] for min-max problems on manifolds. Algorithms are stopped either when gradient norm falls below 10^{-10} or the max iteration has been reached. Results are reported in Fig. 2.

From Fig. 2, we observe rapid convergence of RHM algorithms in all the settings. The convergence for RGDA varies across different choices of c_q, c_l where it converges faster when the weight on the quadratic term (c_q) is relatively higher and is not able to converge when c_l increases. We also observe convergence for RCEG in all cases but the rate is slower compared to RHM algorithms. In Fig. 2f, we further compare the optimality gap where we observe all the proposed RHM algorithms reach below 10^{-10} at a faster rate than the baselines. The slopes of RHM-SD-F and RHM-CON are steeper than that of RCEG (indicating better theoretical rates for RHM). Additional results on optimality gap comparisons are in Fig. 5 in Appendix H. Finally, Fig. 2g shows the runtime performance of various algorithms, with the markers indicating the progress of respective algorithms per iteration. We observe that the per-iteration computational cost of RHM is higher than RGDA. This is because RHM exploits second-order information of f to compute the gradient of \mathcal{H} . Also, we see that RCEG can be costly because it requires evaluation of the exponential map twice and the logarithm map once per iteration.

7.2. Robust geometry-aware PCA. Geometry-aware principal component analysis (PCA) on \mathcal{M}_{SPD} [27] concerns dimensionality reduction for SPD matrices while preserving geometric structures on the manifold. The robust PCA (or robust Fréchet mean) on SPD manifolds has been considered in [92]. For a set of SPD matrices $\mathbf{M}_i \in \mathbb{S}_{++}^d, i = 1, \dots, n$, the aim is to find the Fréchet mean $\mathbf{M} \in \mathbb{S}_{++}^d$ that is bounded away from zero, i.e.,

$$(7.2) \quad \min_{\mathbf{M} \in \mathcal{M}_{\text{SPD}}} \max_{\mathbf{x} \in \mathcal{S}^d} \mathbf{x}^\top \mathbf{M} \mathbf{x} + \frac{\alpha}{n} \sum_{i=1}^n \text{dist}^2(\mathbf{M}, \mathbf{M}_i),$$

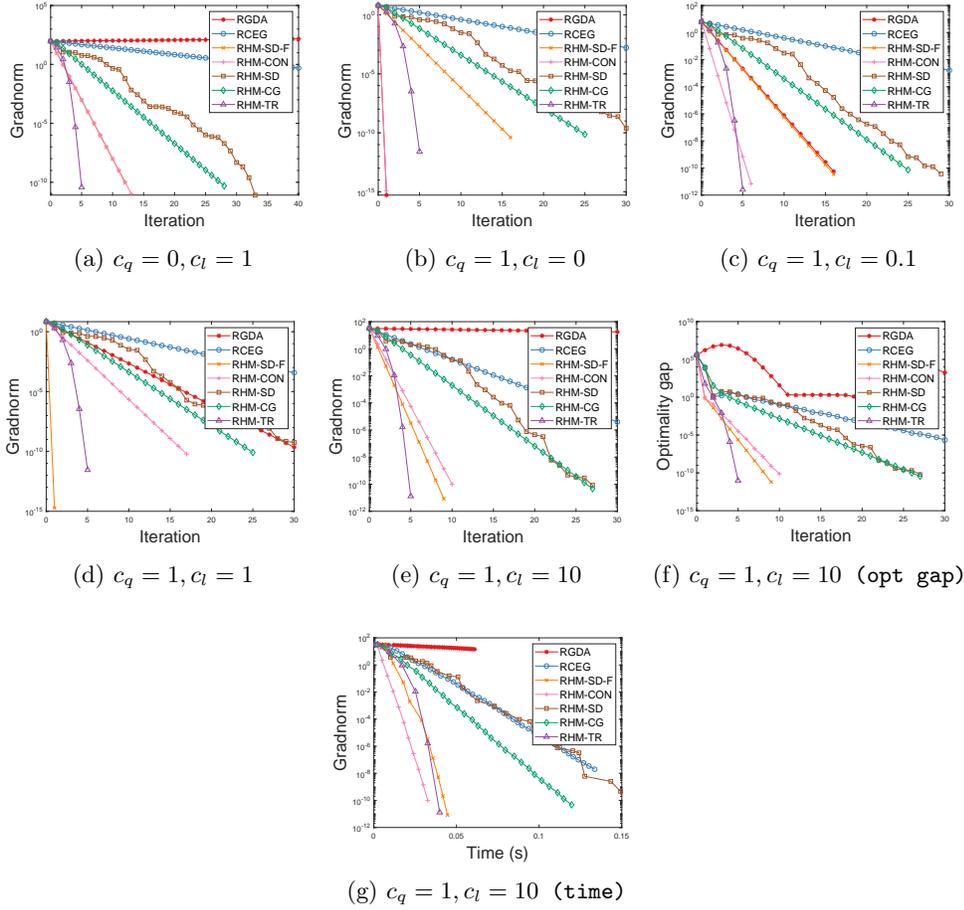


Fig. 2: Experiments on the geodesic quadratic bilinear problem (7.1) with $d = 30$, under varying weights c_q, c_l . We observe that our RHM algorithms converge quickly in all settings while baselines such as RGDA [28] and RCEG [92]. The performance of RGDA varies greatly with the settings where it converges only for a few settings and for the others RGDA fails to converge. RCEG presents a relatively more stable convergence behavior than RGDA but with a rate that is slower than our proposed RHM algorithms.

where $\alpha > 0$ and $\mathcal{S}^{d-1} := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$ denotes the sphere manifold and $\text{dist} : \mathcal{S}_{++}^d \times \mathcal{S}_{++}^d$ is the Riemannian distance on \mathcal{M}_{SPD} .

We first note that the function in (7.2) is geodesic strongly convex in \mathbf{M} and geodesic nonconcave in \mathbf{x} . Also, it is difficult to verify the Riemannian PL condition on the Hamiltonian of (7.2). Hence, this is a challenging problem instance as it does not fall into the studied settings of the existing works [28, 92] including ours.

Experiment settings and results. For this problem, we follow the same settings as discussed in [92] for generating the SPD matrices \mathbf{M}_i with the eigenvalues bounded in $[\mu_0, \mu_1]$. Following [92], we choose $d = 50$, $n = 40$, $\mu_0 = 0.2$, and $\mu_1 = 4.5$.

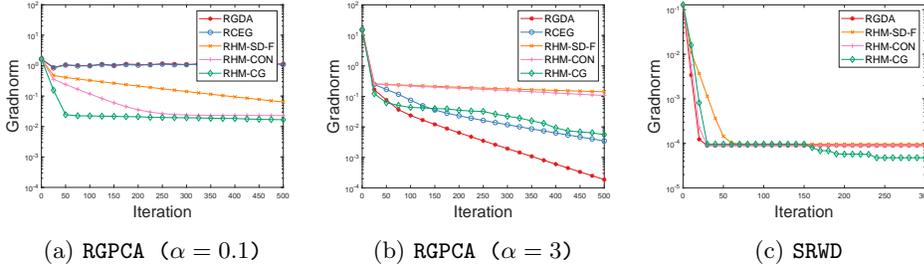


Fig. 3: Convergence on the robust geometry-aware PCA (RGPCA) problem (7.2) with $d = 50, n = 40, \mu = 0.2, L = 4.5$, and subspace robust Wasserstein distance (SRWD) problem on the example of fragmented hypercube [46]. We observe that the baselines RGDA and RCEG fail to converge for $\alpha = 0.1$ (approximately bilinear setting), whereas the proposed RHM algorithms show convergence for $\alpha = 0.1$ and $\alpha = 3$.

The convergence results are presented in Figs. 3a and 3b, where we only include RHM-SD-F, RHM-CON ($\gamma = 0.5$), and RHM-CG for clarity (RHM-TR performs similar to RHM-CG). We observe that although RGDA and RCEG converge faster than RHM when $\alpha = 3$, they fail to converge when $\alpha = 0.1$. The latter finding is not surprising as both RGDA and RCEG seem to perform poorly on approximately bilinear problems (as also observed in Section 7.1). In contrast, we observe that RHM algorithms converge in both the settings, which is also validated by our analysis in Section 3.2. It is known that the conjugate gradient based methods outperforms steepest descent methods on more challenging optimization problems. This explains the faster convergence of RHM-CG over RHM-SD-F and RHM-CON. Overall, the results in Fig. 3 show the benefit of the Riemannian Hamiltonian modeling in non standard settings.

7.3. Subspace robust Wasserstein distance. We next consider the problem of learning subspace robust Wasserstein distance [66, 46, 30], where the aim is to compute the Wasserstein distance over the worst-case optimal transport cost on a low-dimensional space. Given two discrete measures on \mathbb{R}^d , $\mu = \sum_{i=1}^m a_i \delta_{\mathbf{x}_i}, \nu = \sum_{j=1}^n b_j \delta_{\mathbf{y}_j}$ where $\delta_{\mathbf{x}}$ is the Dirac at location \mathbf{x} . The weights a_i, b_j belong to the probability simplex, i.e., $\sum_i a_i = \sum_j b_j = 1$. The objective (with entropy regularization) is then given as

$$(7.3) \quad \min_{\Gamma \in \Pi(\mu, \nu)} \max_{\mathbf{U}: \mathbf{U} \in \text{St}(d, r)} \sum_{i,j} \left(\Gamma_{i,j} \|\mathbf{U}^\top \mathbf{x}_i - \mathbf{U}^\top \mathbf{y}_j\|_2^2 + \epsilon \pi_{i,j} (\log(\pi_{i,j}) - 1) \right),$$

where $\text{St}(d, r) := \{\mathbf{U} \in \mathbb{R}^{d \times r} : \mathbf{U}^\top \mathbf{U} = \mathbf{I}\}$ is the set of column orthonormal matrices ($d \geq r$), known as the Stiefel manifold. $\Pi(\mu, \nu) := \{\Gamma \in \mathbb{R}^{m \times n} : \Gamma_{i,j} > 0, \sum_i \Gamma_{i,j} = b_j, \sum_j \Gamma_{i,j} = a_i, \forall i, j\}$ is the set of couplings, which forms the so-called doubly stochastic manifold (or coupling manifold) [20, 80, 54].

Experiment settings and results. We follow the same experiment settings as in [46, 30] and consider a uniform distribution over hypercube $[-1, 1]^d$ and a push-forward map defined as $T(\mathbf{x}) = \mathbf{x} + 2 \text{sign}(\mathbf{x}) \odot (\sum_{i=1}^k \mathbf{e}_i)$, where $\text{sign}(\mathbf{x})$ extracts the sign of \mathbf{x} elementwise and $\{\mathbf{e}_i\}_{i=1}^d$ are the canonical basis of \mathbb{R}^d .

We choose $d = 30, r = 5, k = 2, n = 100, \epsilon = 0.2$ and compare the proposed RHM-SD-F, RHM-CON ($\gamma = 0.5$), RHM-CG with RGDA in Fig. 3c. RCEG cannot

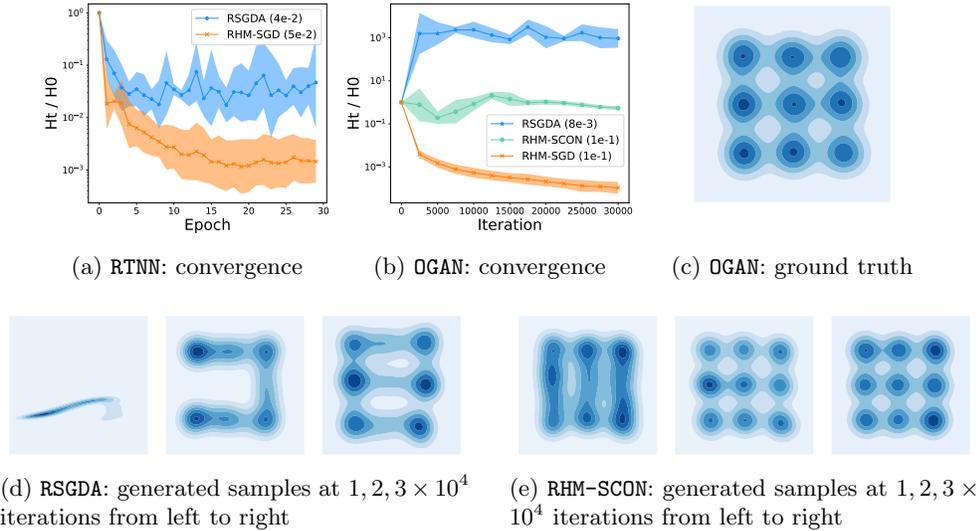


Fig. 4: (4a): Convergence on adversarial robust training of neural network (RTNN). (4b): generative adversarial networks with orthonormal weights (OGAN). (4c): Ground truth distribution. (4d), (4e): generated samples from RSGDA, RHM-SCON respectively where we see RHM-SCON quickly converge to the ground truth distribution while RSGDA suffers from mode collapse. The numbers in the parentheses indicate the best tuned stepsizes for different algorithms.

be implemented to solve (7.3) because the doubly stochastic manifold does not have a well-defined logarithm map. From the results, we see similar convergence speed of all methods while due to the inbuilt line-search algorithm of RHM-CG, it converges to a point with a smaller gradient norm.

7.4. Robust training of neural networks with orthonormal weights. We next consider adversarial robust training of deep neural networks with orthonormal weights [28]. Adversarial training of neural networks provide robust prediction against small data perturbations. Orthonormality on parameters has shown to improve generalization accuracy as well as accelerate and stabilize convergence of neural network models [9, 19, 90, 29]. This corresponds to optimization over the Stiefel manifold.

In particular, we consider the adversarial training to defend against a universal perturbation \mathbf{p} proposed in [59]. The perturbation set we consider is the sphere manifold $\mathcal{S}^{d-1}(r) := \{\mathbf{p} \in \mathbb{R}^d : \|\mathbf{p}\|_2 = r\}$ with radius r . This requires the perturbed samples to stay a certain distance away from the original ones, a strategy also applied in [45]. Given a set of data-target pairs $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^d$ are the feature vectors. The objective of adversarial training is

$$\min_{\{\mathbf{W}_\ell\}_{\ell=1}^L: \mathbf{W}_\ell \in \text{St}(d_\ell, d_{\ell+1})} \max_{\mathbf{p} \in \mathcal{S}^{d-1}(r)} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(h(\mathbf{x}_i + \mathbf{p}; \{\mathbf{W}_\ell\}_{\ell=1}^L), y_i),$$

where $\mathcal{L}(\cdot, \cdot)$ is a loss function and $h(\cdot)$ represents the forward function of a neural network.

Experiment settings and results. The adversarial training is implemented for classification tasks on MNIST images [44] where we include two hidden layers of size 16 with the orthonormality constraint. We compare the proposed stochastic version of RHM (RHM-SGD), detailed in Section 5, with Riemannian stochastic gradient descent ascent (RSGDA) algorithm [28]. We highlight that RHM-SCON performs similarly to RHM-SGD, and thus, we exclude its result for clarity. Because we require dual sampling per-iteration to compute the stochastic Hamiltonian gradient $\text{grad}\mathcal{H}_{\mathcal{S},\mathcal{S}'}(p_t) = \frac{1}{|\mathcal{S}||\mathcal{S}'|} \sum_{\omega \in \mathcal{S}, \varphi \in \mathcal{S}'} \text{grad}\mathcal{H}_{\omega,\varphi}(p)$, we choose the batch size to be 32 for both $\mathcal{S}, \mathcal{S}'$ and 64 for RSGDA. Hence, the per-iteration sampling cost is identical. We measure convergence in terms of the relative Hamiltonian $\mathcal{H}(p_t)/\mathcal{H}(p_0)$, where the Hamiltonian is evaluated on the full training set. The stepsize is fixed for both the algorithms.

We plot the convergence results (with the best tuned stepsize) in Fig. 4a, which are averaged over five different runs. We see a clear advantage of RHM-SGD compared to RSGDA with faster and more stable convergence.

7.5. Orthonormal generative adversarial networks. Generative adversarial networks (GAN) [24, 8] are popular in generating synthetic samples by optimizing a min-max game between a generator and a discriminator. The orthonormality constraint on weight parameters of the discriminator has shown to benefit the training of GANs [17, 60]. In particular, given samples $\{\mathbf{x}_i\}_{i=1}^n$ we consider the following min-max problem

$$\min_{\{\mathbf{W}_\ell^D\}} \max_{\{\mathbf{W}_\ell^G\}: \mathbf{W}_\ell^D \in \text{St}(d_\ell, d_{\ell+1})} \frac{1}{n} \sum_{i=1}^n \left(\log(\sigma(D(\mathbf{x}_i))) + \log(1 - \sigma(D(G(\mathbf{z}_i)))) \right),$$

where $D(\cdot), G(\cdot)$ represent the discriminator and generator with $\{\mathbf{W}_\ell^D\}, \{\mathbf{W}_\ell^G\}$ denoting their network weight parameters respectively. Here, $\sigma(\cdot)$ is the sigmoid function and the prior \mathbf{z}_i is sampled from the standard normal distribution.

Experiment settings and results. Following [8], we train the GAN model on 2-d samples from a multimodal mixture of Gaussian distribution. The ground truth is shown in Fig. 4c. Both the generator and discriminator have 5 hidden layers with 128 units and ReLU activation. The dimension of the prior \mathbf{z}_i is 64. For simplicity, we add the orthonormal constraint only for the penultimate layer of the discriminator model. For this experiment, we apply RHM-SCON with $\gamma = 0.5$ and compare against RSGDA, both with fixed stepsize. The batch size is chosen to be 128 for RHM-SCON and 256 for RSGDA. Similarly, the best choices of stepsize are reported, and the results are averaged over five different runs.

The convergence in terms of the relative Hamiltonian are shown in Fig. 4b, where we see RSGDA diverges while RHM-SCON is more stable. We also examine the solution quality by providing the generated samples from both algorithms at iteration $10^4, 2 \times 10^4$, and 3×10^4 in Figs 4d and 4e respectively. We note that RSGDA results in undesired mode collapse, an observation also made in [8] for training SGDA on the Euclidean space. In contrast, RHM-SCON quickly converges and recovers the ground truth distribution. Even though RHM-SGD converges to a lower Hamiltonian value, its performance in recovery of the ground truth is poor, as shown in Fig. 6 in Appendix H where the generated samples collapse to a single point. It indicates that RHM-SGD converges to a stationary point which is not a saddle point (not surprising as Assumption 1 may not be satisfied). This also highlights the practical benefit of consensus regularization for RHM (Section 4), as evidenced in the good performance

of RHM-SCON.

8. Concluding remarks. Building on the success of the Hamiltonian methods for solving min-max problems in the Euclidean space, we have considered a more general problem on manifolds, and proposed a Riemannian Hamiltonian function \mathcal{H} that respects the manifold geometry. This leads to a gradient expression (in Proposition 3.2) that allows simple analysis for the resulting optimization methods. Adapting the proofs from the Euclidean space to Riemannian manifolds requires to forgo the matrix structure of the ingredients, which includes addressing a varying inner product (Riemannian metric). The proposed Riemannian Hamiltonian methods (RHM) come with convergence guarantees and various extensions. The experiments validate the good performance of RHM in different applications. As future work, one direction is to explore the utility of RHM for more general nonconvex nonconcave problems without the Riemannian PL assumption. In addition, the current convergence analysis is measured in the Riemannian Hamiltonian, which is the gradient norm squared of the original objective f . It remains a question whether linear convergence can be maintained in terms of the optimality gap on function value of f .

Appendix A. Riemannian geometries of the considered manifolds.

In this section, we review the Riemannian optimization-related ingredients of several manifolds that are considered in the experiments section. The expressions are from the works [3, 14, 84, 80, 20, 54].

A.1. Symmetric positive definite manifold. Consider the set of the symmetric positive definite matrices of size $d \times d$, $\mathbb{S}_{++}^d := \{\mathbf{X} : \mathbb{R}^{d \times d} : \mathbf{X}^\top = \mathbf{X}, \mathbf{X} \succeq \mathbf{0}\}$, equipped with the affine-invariant Riemannian metric. The geodesic from \mathbf{X} to \mathbf{Y} is given by $\gamma(t) = \mathbf{X}^{1/2}(\mathbf{X}^{-1/2}\mathbf{Y}\mathbf{X}^{-1/2})^t\mathbf{X}^{1/2}$. At $\mathbf{X} \in \mathbb{S}_{++}^d$, the exponential map is derived as $\text{Exp}_{\mathbf{X}}(\mathbf{U}) = \mathbf{X} \exp(\mathbf{X}^{-1}\mathbf{U})$ for any $\mathbf{U} \in T_{\mathbf{X}}\mathbb{S}_{++}^d$. The logarithm map is $\text{Log}_{\mathbf{X}}(\mathbf{Y}) = \mathbf{X} \log(\mathbf{X}^{-1}\mathbf{Y})$. The Riemannian gradient of a function $f : \mathbb{S}_{++}^d \rightarrow \mathbb{R}$ is given by $\text{grad}f(\mathbf{X}) = \mathbf{X}\nabla f(\mathbf{X})\mathbf{X}$, where $\nabla f(\mathbf{X})$ is the Euclidean partial derivative of f at \mathbf{X} .

A.2. Sphere manifold. It is defined as $\mathcal{S}^{d-1} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$, which is an embedded submanifold of \mathbb{R}^d with the tangent space expression $T_{\mathbf{x}}\mathcal{S}^{d-1} = \{\mathbf{u} \in \mathbb{R}^d : \mathbf{x}^\top \mathbf{u} = 0\}$. It can be endowed with the standard inner product at the Riemannian metric, i.e., $\langle \mathbf{u}, \mathbf{v} \rangle_{\mathbf{x}} = \langle \mathbf{u}, \mathbf{v} \rangle_2$, for $\mathbf{u}, \mathbf{v} \in T_{\mathbf{x}}\mathcal{S}^{d-1}$. The orthogonal projection of any $\mathbf{v} \in \mathbb{R}^d$ to $T_{\mathbf{x}}\mathcal{S}^{d-1}$ is derived as $\text{Proj}_{\mathbf{x}}(\mathbf{v}) = \mathbf{v} - (\mathbf{x}^\top \mathbf{v})\mathbf{x}$. The exponential map along $u \in T_{\mathbf{x}}\mathcal{S}^{d-1}$ is $\text{Exp}_{\mathbf{x}}(\mathbf{u}) = \cos(\|\mathbf{v}\|_2)\mathbf{x} + \sin(\|\mathbf{v}\|_2) \frac{\mathbf{v}}{\|\mathbf{v}\|}$ and the logarithm map is $\text{Log}_{\mathbf{x}}(\mathbf{y}) = \arccos(\mathbf{x}^\top \mathbf{y}) \frac{\text{Proj}_{\mathbf{x}}(\mathbf{y}-\mathbf{x})}{\|\text{Proj}_{\mathbf{x}}(\mathbf{y}-\mathbf{x})\|_2}$. The Riemannian gradient of f is $\text{Proj}_{\mathbf{x}}(\nabla f(\mathbf{x}))$, where $\nabla f(\mathbf{x})$ is the Euclidean partial derivative of f at \mathbf{x} .

A.3. Stiefel manifold. It is the set $\text{St}(d, r) = \{\mathbf{X} \in \mathbb{R}^{d \times r} : \mathbf{X}^\top \mathbf{X} = \mathbf{I}\}$. It is a generalization of the sphere manifold to higher dimensions and can be similarly endowed with the standard inner product as metric $\langle \mathbf{U}, \mathbf{V} \rangle_{\mathbf{X}} = \langle \mathbf{U}, \mathbf{V} \rangle_2$. For the experiments, we consider the popular QR-based retraction for approximating the exponential map, i.e., $R_{\mathbf{X}}(\mathbf{U}) = \text{qf}(\mathbf{X} + \mathbf{U})$, where $\text{qf}(\cdot)$ returns the Q-factor from the QR decomposition for any tangent vector \mathbf{U} .

A.4. Doubly stochastic manifold. The doubly stochastic manifold (or coupling manifold) between two discrete probability measures $\mu = \sum_{i=1}^m a_i \delta_{\mathbf{x}_i}, \nu = \sum_{j=1}^n b_j \delta_{\mathbf{y}_j}$ is the set of couplings $\Pi(\mu, \nu) := \{\mathbf{\Gamma} \in \mathbb{R}^{m \times n} : \mathbf{\Gamma}_{i,j} > 0, \sum_i \mathbf{\Gamma}_{i,j} = b_j, \sum_j \mathbf{\Gamma}_{i,j} = a_i, \forall i, j\}$ endowed with the Fisher information Riemannian metric. The geometry has been developed in [20, 80, 54].

Algorithm B.1 Backtracking line-search

-
- 1: **Input:** Current iterate $p_t \in \mathcal{M}$, search direction $\xi_t \in T_{p_t}\mathcal{M}$, initial stepsize $\bar{\vartheta}$ and $r_1, \rho \in (0, 1)$.
 - 2: Initialize $\vartheta \leftarrow \bar{\vartheta}$.
 - 3: **while** $h(p_t) - h(\text{Exp}_{p_t}(\vartheta\xi_t)) < r_1\vartheta\langle -\text{grad}h(p_t), \xi_t \rangle_{p_t}$ **do**
 - 4: Set $\vartheta \leftarrow \rho\bar{\vartheta}$.
 - 5: **end while**
 - 6: **Output:** ϑ .
-

Without loss of any generality, we assume $\sum_i a_i = \sum_j b_j = 1$. The tangent space at $\mathbf{\Gamma} \in \Pi(\mu, \nu)$ is given by $T_{\mathbf{\Gamma}}\Pi(\mu, \nu) = \{\mathbf{U} \in \mathbb{R}^{m \times n} : \sum_i \mathbf{U}_{i,j} = \sum_j \mathbf{U}_{i,j} = 0, \forall i, j\}$. The Fisher information metric is defined as for $\mathbf{U}, \mathbf{V} \in T_{\mathbf{\Gamma}}\Pi(\mu, \nu)$, $\langle \mathbf{U}, \mathbf{V} \rangle_{\mathbf{\Gamma}} = \sum_{i,j} (\mathbf{U}_{i,j} \mathbf{V}_{i,j}) / \mathbf{\Gamma}_{i,j}$. For the experiments, we consider the Sinkhorn-based retraction. The Sinkhorn-Knopp algorithm [81] is a popular approach for balancing non-negative matrices to satisfy the row-sum and column sum constraint and later adapted to solve the optimal transport problem efficiently [68]. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{A}_{i,j} > 0$, and denote $\text{Sinkhorn}(\mathbf{A})$ as the output of applying the Sinkhorn-Knopp algorithm on \mathbf{A} with constraint defined by $\Pi(\mu, \nu)$, i.e., $\text{Sinkhorn}(\mathbf{A}) \in \Pi(\mu, \nu)$. Subsequently, the retraction is given by $R_{\mathbf{\Gamma}}(\mathbf{U}) = \text{Sinkhorn}(\mathbf{\Gamma} \odot \exp(\mathbf{U} \oslash \mathbf{\Gamma}))$, where \exp , \odot , and \oslash are elementwise exponential, product, and division operations, respectively.

Appendix B. Line-search methods and Wolfe conditions on Riemannian manifolds. In this section, we present the Riemannian versions of the Armijo, Wolfe, and strong Wolfe conditions [77].

DEFINITION B.1. Consider an iterative algorithm for minimizing $h : \mathcal{M} \rightarrow \mathbb{R}$, producing $p_{t+1} = \text{Exp}_{p_t}(\vartheta_t \xi_t)$ for some direction $\xi_t \in T_{p_t}\mathcal{M}$ and stepsize $\vartheta_t \in \mathbb{R}$. The Armijo condition is $h(p_t) - h(p_{t+1}) \geq r_1 \vartheta_t \langle -\text{grad}h(p_t), \xi_t \rangle$, for some $r_1 \in (0, 1)$. The (weak) Wolfe condition is the Armijo condition together with (B.1) and the strong Wolfe condition is the Armijo condition with (B.2), where

$$(B.1) \quad \langle \text{grad}h(p_{t+1}), \text{DExp}_{p_t}(\vartheta_t \xi_t)[\xi_t] \rangle_{p_{t+1}} \geq r_2 \langle \text{grad}h(p_t), \xi_t \rangle_{p_t}$$

$$(B.2) \quad |\langle \text{grad}h(p_{t+1}), \text{DExp}_{p_t}(\vartheta_t \xi_t)[\xi_t] \rangle_{p_{t+1}}| \leq r_2 |\langle \text{grad}h(p_t), \xi_t \rangle_{p_t}|$$

for some $r_2 \in (r_1, 1)$. Here, DExp is the differential of the exponential operation.

The backtracking line-search for satisfying the Armijo condition has been used in Riemannian steepest descent method [15].

One can generalize the analysis from the Euclidean space to show that there exists a stepsize that satisfy the three conditions for arbitrary direction ξ_t . The backtracking line-search for satisfying the Armijo condition is in Algorithm B.1. This has been used in Riemannian steepest descent method [15]. The procedures that return stepsizes satisfying the Wolfe conditions are in [75, 64].

Appendix C. Review of RGDA and RCEG. In this section, we provide the details of the Riemannian gradient descent ascent [28] and Riemannian corrected extra-gradient [92] algorithms for min-max optimization on manifolds.

RGDA simultaneously updates the variables in the direction of the min-max Riemannian gradient, i.e.,

$$x_{t+1} = \text{Exp}_{x_t}(-\eta_t \text{grad}_x f(x_t, y_t)), \quad y_{t+1} = \text{Exp}_{y_t}(\eta_t \text{grad}_y f(x_t, y_t)).$$

RCEG first updates the variables to the point (w_t, z_t) along the min-max Riemannian gradient. It then uses the obtained point to generate the final update, i.e.,

$$\begin{aligned} w_t &= \text{Exp}_{x_t}(-\eta_t \text{grad}_x f(x_t, y_t)), \\ z_t &= \text{Exp}_{y_t}(\eta_t \text{grad}_y f(x_t, y_t)), \\ x_{t+1} &= \text{Exp}_{w_t}(-\eta_t \text{grad}_x f(w_t, y_t) + \text{Log}_{w_t}(x_t)), \\ y_{t+1} &= \text{Exp}_{z_t}(\eta_t \text{grad}_y f(w_t, y_t) + \text{Log}_{z_t}(y_t)). \end{aligned}$$

In [92], only convergence for g-convex-concave functions is analyzed, where the authors show that RCEG converges sublinearly with averaged iterate under the fixed stepsize $\eta \leq \frac{1}{2L_1\tau_{\zeta,D}}$ where $\tau_{\zeta,D} > 1$ depends on the curvature and diameter of the domain. Thus, the analysis is only local with domain-dependent rate of convergence. The recent work [40] starts by showing average-iterate convergence of RCEG under g-convex-concave functions and last-iterate convergence under g-strongly-convex-concave functions. Nevertheless, similar assumptions on the bounded domain (and also the curvature) is required. The stepsize also requires to be carefully selected, which depends on the curvature and diameter bound. In addition, [40] proves convergence for RGDA under similar settings. For g-strongly-convex-concave functions, the last-iterate convergence of RGDA requires a diminishing stepsize, and for g-convex-concave functions, the average-iterate convergence of RGDA require a stepsize that again depends on the curvature and diameter bound.

Appendix D. Key propositions.

In this section, we derive the explicit expression for the Riemannian Hessian on the product manifold $\mathcal{M} = \mathcal{M}_x \times \mathcal{M}_y$ and show that the cross derivatives are adjoint with respect to the Riemannian metric.

PROPOSITION D.1 (Riemannian Hessian of product manifold). *Consider a product Riemannian manifold $\mathcal{M} = \mathcal{M}_x \times \mathcal{M}_y$ and $f : \mathcal{M} \rightarrow \mathbb{R}$. For any $p = (x, y) \in \mathcal{M}$ and $\xi = (u, v) \in T_x\mathcal{M}$, the Riemannian Hessian $\text{Hess}f(p)[\xi]$ is derived as*

$$\text{Hess}f(p)[\xi] = \begin{pmatrix} \text{Hess}_x f(x, y)[u] + \text{grad}_{y,x}^2 f(x, y)[v] \\ \text{grad}_{x,y}^2 f(x, y)[u] + \text{Hess}_y f(x, y)[v] \end{pmatrix}.$$

Proof. From standard analysis, the Levi-Civita connection on a product manifold $\mathcal{M} = \mathcal{M}_x \times \mathcal{M}_y$ (e.g., in [14, Exercise 5.4]) is given by

$$\nabla_{(U_x, U_y)}(V_x, V_y) = \left(\nabla_{U_x}^{(x)} V_x + D_y V_x[U_y], D_x V_y[U_x] + \nabla_{U_y}^{(y)} V_y \right),$$

where $V_x \in \mathfrak{X}(\mathcal{M}_x)$, $V_y \in \mathfrak{X}(\mathcal{M}_y)$ are vector fields on respective manifolds and D is the directional derivative. Further, $D_y V_x : \mathfrak{X}(\mathcal{M}_y) \rightarrow \mathfrak{X}(\mathcal{M}_x)$ and when evaluating at (x, y) , this is equivalently defined as $D_y V_x(x, \cdot)(y) : T_y\mathcal{M}_y \rightarrow T_x\mathcal{M}_x$, which is the directional derivative. $\nabla^{(x)}$, $\nabla^{(y)}$ are the Levi-Civita connections on \mathcal{M}_x , \mathcal{M}_y , respectively. Applying the definition of the Riemannian Hessian, $\text{Hess}f(p)[\xi] = \nabla_{\xi} \text{grad}f(p)$, we obtain the desired result. \square

PROPOSITION D.2. *For any $(x, y) \in \mathcal{M}_x \times \mathcal{M}_y$ and $(u, v) \in T_x\mathcal{M}_x \times T_y\mathcal{M}_y$, we have $\langle \text{grad}_{y,x}^2 f(x, y)[v], u \rangle_x = \langle \text{grad}_{x,y}^2 f(x, y)[u], v \rangle_y$. Equivalently, $\text{grad}_{y,x}^2 f(x, y)$ is the adjoint operator of $\text{grad}_{x,y}^2 f(x, y)$.*

Proof. Let $p = (x, y)$ and $\xi = (u, v)$, $\zeta = (w, z)$ for any $(u, v), (w, z) \in T_x\mathcal{M}_x \times T_y\mathcal{M}_y$. Then, from the self-adjoint property (symmetry) of the Riemannian Hessian,

we have

$$(D.1) \quad \langle \text{Hess}f(p)[\xi], \zeta \rangle_p = \langle \text{Hess}f(p)[\zeta], \xi \rangle_p,$$

for any ξ, ζ . Combining with Proposition D.1, the result (D.1) is equivalent to

$$\begin{aligned} & \langle \text{Hess}_x f(x, y)[u], w \rangle_x + \langle \text{grad}_{yx}^2 f(x, y)[v], w \rangle_x + \langle \text{grad}_{xy}^2 f(x, y)[u], z \rangle_y \\ & \quad + \langle \text{Hess}_y f(x, y)[v], z \rangle_y \\ = & \langle \text{Hess}_x f(x, y)[w], u \rangle_x + \langle \text{grad}_{yx}^2 f(x, y)[z], u \rangle_x + \langle \text{grad}_{xy}^2 f(x, y)[w], v \rangle_y \\ & \quad + \langle \text{Hess}_y f(x, y)[z], v \rangle_y. \end{aligned}$$

Given that Hess_x and Hess_y satisfy the self-adjoint property, we obtain

$$(D.2) \quad \begin{aligned} & \langle \text{grad}_{yx}^2 f(x, y)[v], w \rangle_x + \langle \text{grad}_{xy}^2 f(x, y)[u], z \rangle_y \\ = & \langle \text{grad}_{yx}^2 f(x, y)[z], u \rangle_x + \langle \text{grad}_{xy}^2 f(x, y)[w], v \rangle_y. \end{aligned}$$

We can see (D.2) holds for any choice of $(u, v), (w, z)$ and this only happens when $\langle \text{grad}_{yx}^2 f(x, y)[v], u \rangle_x = \langle \text{grad}_{xy}^2 f(x, y)[u], v \rangle_y$ holds for any (u, v) . To see this, consider the vectorization of the tangent vectors as $\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{z}$. We also denote $\mathbf{B}_{xy}, \mathbf{B}_{yx}$ as the matrix representation of the linear operators $\text{grad}_{xy}^2 f(x, y), \text{grad}_{yx}^2 f(x, y)$ at (x, y) respectively. Then (D.2) can be rewritten as

$$\mathbf{w}^\top \mathbf{G}_x \mathbf{B}_{yx} \mathbf{v} + \mathbf{z}^\top \mathbf{G}_y \mathbf{B}_{xy} \mathbf{u} = \mathbf{u}^\top \mathbf{G}_x \mathbf{B}_{yx} \mathbf{z} + \mathbf{v}^\top \mathbf{G}_y \mathbf{B}_{xy} \mathbf{w},$$

where $\mathbf{G}_x, \mathbf{G}_y$ are the (symmetric positive definite) metric tensors at x, y . This is equivalent to

$$\mathbf{z}^\top (\mathbf{G}_y \mathbf{B}_{xy} - \mathbf{B}_{yx}^\top \mathbf{G}_x) \mathbf{u} = \mathbf{v}^\top (\mathbf{G}_y \mathbf{B}_{xy} - \mathbf{B}_{yx}^\top \mathbf{G}_x) \mathbf{w},$$

which is satisfied for any $\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{z}$ and any $\mathbf{G}_x, \mathbf{G}_y$ as metric tensors. Hence, $\mathbf{G}_y \mathbf{B}_{xy} = \mathbf{B}_{yx}^\top \mathbf{G}_x$ and the proof is complete. \square

Remark D.3. Proposition D.2 shows that the Riemannian cross derivatives are symmetric with respect to Riemannian metric on respective manifolds. When $\mathcal{M}_x, \mathcal{M}_y$ are the Euclidean spaces, then Proposition D.2 is equivalent to the Schwarz's theorem of symmetric second-order derivatives.

Appendix E. Essential lemmas.

The following lemmas generalize [1, Lemmas 17, 28] to linear operators, specifically in terms of the Riemannian Hessian operator. We first highlight that for two operators T, T^* that are adjoint, we have $\lambda(T \circ T^*) = \lambda(T^* \circ T) = \sigma^2(T) = \sigma^2(T^*)$.

LEMMA E.1. *Consider the Riemannian Hessian $\text{Hess}f(p)$ where $p = (x, y) \in \mathcal{M}_x \times \mathcal{M}_y$. Suppose $\text{Hess}_y f(x, y) = 0$. Then, $\lambda_{|\min|}(\text{Hess}f(p)) \geq \frac{\sigma_{\min}^2(B_{xy})}{\sqrt{2\sigma_{\min}^2(B_{xy}) + \|H_x\|_x^2}}$.*

Proof. We consider the operator $\text{Hess}f(p) \circ \text{Hess}f(p)$ and study its eigenvalue. First, we see that for any $p = (x, y) \in \mathcal{M}_x \times \mathcal{M}_y$ and $\xi = (u, v) \in T_x \mathcal{M}_x \times T_y \mathcal{M}_y$, we have

$$\text{Hess}f(p)[\xi] = \begin{pmatrix} \text{Hess}_x f(x, y)[u] + \text{grad}_{yx}^2 f(x, y)[v] \\ \text{grad}_{xy}^2 f(x, y)[u] \end{pmatrix},$$

and therefore,

$$\begin{aligned} & \text{Hess}f(p)[\text{Hess}f(p)[\xi]] \\ &= \begin{pmatrix} \text{Hess}_x f(x, y)[\text{Hess}_x f(x, y)[u]] + \text{Hess}_x f(x, y)[\text{grad}_{yx}^2 f(x, y)[v]] \\ \quad + \text{grad}_{yx}^2 f(x, y)[\text{grad}_{xy}^2 f(x, y)[u]] \\ \text{grad}_{xy}^2 f(x, y)[\text{Hess}_x f(x, y)[u]] + \text{grad}_{xy}^2 f(x, y)[\text{grad}_{yx}^2 f(x, y)[v]] \end{pmatrix}. \end{aligned}$$

Suppose (δ, ξ) is an eigenpair of the operator $\text{Hess}f(p) \circ \text{Hess}f(p)$, which gives

$$(E.1) \quad \begin{aligned} & \text{Hess}_x f(x, y)[\text{Hess}_x f(x, y)[u]] + \text{Hess}_x f(x, y)[\text{grad}_{yx}^2 f(x, y)[v]] \\ & \quad + \text{grad}_{yx}^2 f(x, y)[\text{grad}_{xy}^2 f(x, y)[u]] = \delta u, \end{aligned}$$

$$(E.2) \quad \text{grad}_{xy}^2 f(x, y)[\text{Hess}_x f(x, y)[u]] + \text{grad}_{xy}^2 f(x, y)[\text{grad}_{yx}^2 f(x, y)[v]] = \delta v.$$

Let $B_{xy} = \text{grad}_{xy}^2 f(x, y)$, $B_{yx} = \text{grad}_{yx}^2 f(x, y)$, and $H_x = \text{Hess}_x f(x, y)$. Suppose $\delta < \frac{\sigma_{\min}^4(B_{xy})}{2\sigma_{\min}^2(B_{xy}) + \|H_x\|_x^2} < \sigma_{\min}^2(B_{xy})$. Then, we have $B_{xy} \circ B_{yx} - \delta \text{id}$ is invertible where we use the fact that B_{xy} and B_{yx} are adjoint. Hence, from (E.2) we have $v = -(B_{xy} \circ B_{yx} - \delta \text{id})^{-1} \circ (B_{xy} \circ H_x)[u]$. Substituting the expression of v into (E.1) yields

$$(E.3) \quad \left(H_x \circ (\text{id} - B_{yx} \circ (B_{xy} \circ B_{yx} - \delta \text{id})^{-1} \circ B_{xy}) \circ H_x + B_{yx} \circ B_{xy} - \delta \text{id} \right) [u] = 0.$$

We next show that when

$$(E.4) \quad \delta < \frac{\sigma_{\min}^4(B_{xy})}{2\sigma_{\min}^2(B_{xy}) + \|H_x\|_x^2} < \sigma_{\min}^2(B_{xy}),$$

then (E.3) does not have a nontrivial solution in u (i.e., $u \neq 0$), which leads to a contradiction that ξ is an eigenvector. It suffices to show that for any δ satisfying the condition (E.4), the following inequality

$$(E.5) \quad \frac{-\delta \|H_x\|_x^2}{\sigma_{\min}^2(B_{xy}) - \delta} + \sigma_{\min}^2(B_{xy}) - \delta > 0,$$

holds, which violates (E.3). Here, we highlight B_{xy} is the adjoint of B_{yx} , and therefore, the eigenvalues $\lambda_i(\text{id} - B_{yx} \circ (B_{xy} \circ B_{yx} - \delta \text{id})^{-1} \circ B_{xy}) = \frac{-\delta}{\sigma_i^2(B_{xy}) - \delta} < 0$ from the singular value decomposition of B_{xy} . The roots of (E.5) are

$$\begin{aligned} r_1 &= \sigma_{\min}^2(B_{xy}) + \frac{1}{2} \|H_x\|_x^2 - \sqrt{(\sigma_{\min}^2(B_{xy}) + \frac{1}{2} \|H_x\|_x^2)^2 - \sigma_{\min}^4(B_{xy})} \\ r_2 &= \sigma_{\min}^2(B_{xy}) + \frac{1}{2} \|H_x\|_x^2 + \sqrt{(\sigma_{\min}^2(B_{xy}) + \frac{1}{2} \|H_x\|_x^2)^2 - \sigma_{\min}^4(B_{xy})}. \end{aligned}$$

One can show for any $c_1 > 0$, $4c_2 < c_1^2$, then $\frac{2c_2}{c_1} < c_1 - \sqrt{c_1^2 - 4c_2}$. Let $c_1 = \sigma_{\min}^2(B_{xy}) + \frac{1}{2} \|H_x\|_x^2$, $c_2 = \frac{1}{4} \sigma_{\min}^4(B_{xy})$, we have the smaller root satisfies $r_1 > \frac{\sigma_{\min}^4(B_{xy})}{2\sigma_{\min}^2(B_{xy}) + \|H_x\|_x^2} > \delta$. Hence, there does not exist $u \neq 0$ that satisfies (E.3), which implies $\delta \geq \frac{\sigma_{\min}^4(B_{xy})}{2\sigma_{\min}^2(B_{xy}) + \|H_x\|_x^2}$. This completes the proof. \square

LEMMA E.2. Consider the Riemannian Hessian $\text{Hess}f(p)$, where $p = (x, y) \in \mathcal{M}_x \times \mathcal{M}_y$. Let $H_x := \text{Hess}_x f(x, y)$, $H_y := \text{Hess}_y f(x, y)$, $B_{xy} := \text{grad}_{xy}^2 f(x, y)$, and

$$\begin{aligned} a &= 2\sigma_{\min}^2(B_{xy}) + \lambda_{|\min|}^2(H_x) + \lambda_{|\min|}^2(H_y), \\ b &= \left(\sigma_{\min}^2(B_{xy}) + \lambda_{|\min|}^2(H_x)\right) \left(\sigma_{\min}^2(B_{xy}) + \lambda_{|\min|}^2(H_y)\right) \\ &\quad - \sigma_{\max}^2(B_{xy})(\|H_x\|_x + \|H_y\|_y)^2. \end{aligned}$$

Suppose that $b > 0$. Then, $\lambda_{|\min|}(\text{Hess}f(p)) \geq \sqrt{\frac{b}{a}}$.

Proof. Similarly to Lemma E.1, we consider the operator $\text{Hess}f(p) \circ \text{Hess}f(p)$, i.e.,

$$\text{Hess}f(p)[\xi] = \begin{pmatrix} \text{Hess}_x f(x, y)[u] + \text{grad}_{yx}^2 f(x, y)[v] \\ \text{Hess}_y f(x, y)[v] + \text{grad}_{xy}^2 f(x, y)[u] \end{pmatrix},$$

and

$$\begin{aligned} &\text{Hess}f(p)[\text{Hess}f(p)[\xi]] \\ &= \begin{pmatrix} \text{Hess}_x f(x, y)[\text{Hess}_x f(x, y)[u] + \text{Hess}_x f(x, y)[\text{grad}_{yx}^2 f(x, y)[v]] \\ + \text{grad}_{yx}^2 f(x, y)[\text{Hess}_y f(x, y)[v] + \text{grad}_{yx}^2 f(x, y)[\text{grad}_{xy}^2 f(x, y)[u]] \\ \text{Hess}_y f(x, y)[\text{Hess}_y f(x, y)[v] + \text{Hess}_y f(x, y)[\text{grad}_{xy}^2 f(x, y)[u]] \\ + \text{grad}_{xy}^2 f(x, y)[\text{Hess}_x f(x, y)[u] + \text{grad}_{xy}^2 f(x, y)[\text{grad}_{yx}^2 f(x, y)[v]] \end{pmatrix}. \end{aligned}$$

Suppose (δ, ξ) is an eigenpair of the operator $\text{Hess}f(p) \circ \text{Hess}f(p)$, which gives

$$\begin{aligned} &\text{Hess}_x f(x, y)[\text{Hess}_x f(x, y)[u] + \text{Hess}_x f(x, y)[\text{grad}_{yx}^2 f(x, y)[v]] \\ &+ \text{grad}_{yx}^2 f(x, y)[\text{Hess}_y f(x, y)[v] + \text{grad}_{yx}^2 f(x, y)[\text{grad}_{xy}^2 f(x, y)[u]] = \delta u, \end{aligned} \tag{E.6}$$

$$\begin{aligned} &\text{Hess}_y f(x, y)[\text{Hess}_y f(x, y)[v] + \text{Hess}_y f(x, y)[\text{grad}_{xy}^2 f(x, y)[u]] \\ &+ \text{grad}_{xy}^2 f(x, y)[\text{Hess}_x f(x, y)[u] + \text{grad}_{xy}^2 f(x, y)[\text{grad}_{yx}^2 f(x, y)[v]] = \delta v. \end{aligned} \tag{E.7}$$

Denote $T_x := H_x \circ H_x + B_{yx} \circ B_{xy} - \delta \text{id}$ and similarly for $T_y := H_y \circ H_y + B_{xy} \circ B_{yx} - \delta \text{id}$, where $H_x = \text{Hess}_x f(x, y)$, $H_y = \text{Hess}_y f(x, y)$ and $B_{xy} = \text{grad}_{xy}^2 f(x, y)$, $B_{yx} = \text{grad}_{yx}^2 f(x, y)$. Then, we can simplify (E.6) and (E.7) as

$$\begin{aligned} T_x[u] &= -(H_x \circ B_{yx} + B_{yx} \circ H_y)[v] \\ T_y[v] &= -(H_y \circ B_{xy} + B_{xy} \circ H_x)[u] \end{aligned} \tag{E.8}$$

Suppose $\delta < \frac{b}{a}$. Then, we can show T_y is invertible. This is because, for any $c_1 > 0$, $4c_2 < c_1^2$, we have $\frac{2c_2}{c_1} < c_1 - \sqrt{c_1^2 - 4c_2}$. From the definition of a and b and setting $c_1 = a$, $c_2 = b$, we have

$$\begin{aligned} \frac{2b}{a} &< 2\sigma_{\min}^2(B_{xy}) + \lambda_{\min}(H_x \circ H_x) + \lambda_{\min}(H_y \circ H_y) \\ &\quad - \sqrt{(\lambda_{\min}(H_x \circ H_x) - \lambda_{\min}(H_y \circ H_y))^2 + 4\sigma_{\max}^2(B_{xy})(\|H_x\|_x + \|H_y\|_y)^2} \\ &< 2\sigma_{\min}^2(B_{xy}) + \lambda_{\min}(H_x \circ H_x) + \lambda_{\min}(H_y \circ H_y) \\ &\quad - |\lambda_{\min}(H_x \circ H_x) - \lambda_{\min}(H_y \circ H_y)| \\ &\leq 2\sigma_{\min}^2(B_{xy}) + 2\lambda_{\min}(H_y \circ H_y), \end{aligned}$$

where we emphasize that B_{yx} is the adjoint to B_{xy} and hence $\lambda(B_{yx} \circ B_{xy}) = \lambda(B_{xy} \circ B_{yx}) = \sigma^2(B_{xy}) = \sigma^2(B_{yx})$.

Hence, $\delta < \frac{b}{a} < \sigma_{\min}^2(B_{xy}) + \lambda_{\min}(H_y \circ H_y)$ and $T_y = H_y \circ H_y + B_{xy} \circ B_{yx} - \delta \text{id}$ is invertible, because $\lambda_{\min}(T_y) \geq \sigma_{\min}^2(B_{xy}) + \lambda_{\min}(H_y \circ H_y) - \delta > 0$ by Weyl's inequality. Thus, (E.8) gives $v = -T_y^{-1} \circ (H_y \circ B_{xy} + B_{xy} \circ H_x)[u]$. Substituting this expression for v into the first equation of (E.8) yields

$$(E.9) \quad \left(T_x - (H_x \circ B_{yx} + B_{yx} \circ H_y) \circ T_y^{-1} \circ (H_y \circ B_{xy} + B_{xy} \circ H_x) \right) [u] = 0.$$

Nevertheless, we can verify when $\delta < \frac{b}{a}$, (E.9) does not have any nontrivial solution for u , which gives a contradiction. Specifically, we show the following inequality is always satisfied under the condition on δ ,

$$(E.10) \quad \begin{aligned} & (\lambda_{\min}(H_y \circ H_y) + \sigma_{\min}^2(B_{xy}) - \delta)^{-1} \sigma_{\max}^2(B_{xy}) (\|H_x\|_x + \|H_y\|_y)^2 \\ & < \lambda_{\min}(H_x \circ H_x) + \sigma_{\min}^2(B_{xy}) - \delta, \end{aligned}$$

which violates (E.9) for any $u \neq 0$, because (E.10) would imply that

$$\lambda_{\min} \left(\left(T_x - (H_x \circ B_{yx} + B_{yx} \circ H_y) \circ T_y^{-1} \circ (H_y \circ B_{xy} + B_{xy} \circ H_x) \right) \right) > 0,$$

subsequently (E.9) implies $u = 0$, hence, $\xi = 0$, a contradiction. It remains to show that under $\delta < \frac{b}{a}$, (E.10) is satisfied. That is, the roots of (E.10) are given by $\frac{1}{2}(a \pm \sqrt{a^2 - 4b})$. We have shown that $\delta < \frac{b}{a} < \frac{1}{2}(a - \sqrt{a^2 - 4b})$. This implies (E.10) is always satisfied and results in a contradiction. Hence, $\delta \geq \frac{b}{a}$, which completes the proof. \square

Appendix F. Analysis of RHM with conjugate gradient and trust-region update steps. We provide the details on convergence analysis of minimizing the Riemannian Hamiltonian with the Riemannian conjugate gradient and trust-region methods, i.e., we consider Algorithm 3.1 with the update step $\xi(p_t)$ computed as conjugate gradient direction and trust-region step.

F.1. RHM with conjugate gradient (RHM-CG).

THEOREM F.1 (Linear convergence of RHM-CG). *Under the same settings as in Theorem 3.7, consider Algorithm 3.1 with conjugate gradient direction $\xi(p_t)$ where β_t (used in update) and η_t are chosen such that $\langle \xi(p_t), -\text{grad}\mathcal{H}(p_t) \rangle \geq c \|\text{grad}\mathcal{H}(p_t)\|_{p_t}^2$ for some $c > 0$ and the Armijo condition (Definition B.1) is satisfied. Let $\tilde{\eta} = \min_{i=0, \dots, t} \eta_i$. Then, iterates p_t satisfy $\|\text{grad}f(p_t)\|_{p_t}^2 \leq (1 - 2r_1\tilde{\eta}c\delta)^t \|\text{grad}f(p_0)\|_{p_0}^2$.*

Proof. From the Armijo condition, we have for the stepsize η_t ,

$$\begin{aligned} \mathcal{H}(p_{t+1}) - \mathcal{H}(p_t) & \leq r_1 \eta_t \langle \text{grad}\mathcal{H}(p_t), \zeta(p_t) \rangle \\ & \leq -r_1 \eta_t c \|\text{grad}\mathcal{H}(p_t)\|_{p_t}^2 \leq -2r_1 \eta_t c \delta \mathcal{H}(p_t) \leq -2r_1 \tilde{\eta} c \delta \mathcal{H}(p_t), \end{aligned}$$

where the last inequality follows from the definition of $\tilde{\eta}$ and $\mathcal{H}(p_t) \geq 0$ for all p_t . Applying the result recursively completes the proof. \square

We notice that the bound only requires a descent direction and a sufficient function decrease. Hence, we suspect a tighter bound exists when analyzing specific types of conjugate gradient (with different β_t types).

We also highlight that most, if not all, types of conjugate gradient methods satisfy the conditions in Theorem F.1. See more discussions in [76]. As an example, consider

the *Fletcher-Reeves-type CG* [22] with $\beta_t = \frac{\|\text{grad}\mathcal{H}(p_t)\|_{p_t}^2}{\|\text{grad}\mathcal{H}(p_{t-1})\|_{p_{t-1}}^2}$. If the stepsize η_t is chosen to satisfy the strong Wolfe conditions (Definition B.1) with $0 < r_1 < r_2 < 1/2$, then from [77, Lemma 4.1], the conditions in Theorem F.1 are satisfied with $\langle \xi(p_t), -\text{grad}\mathcal{H}(p_t) \rangle \geq \frac{1-2r_2}{1-r_2} \|\text{grad}\mathcal{H}(p_t)\|^2$.

F.2. RHM with trust-region (RHM-TR). For the Riemannian trust-region (TR) method, the update step $\xi(p_t)$ is computed by (approximately) solving the trust-region subproblem on the tangent space [3], i.e.,

$$(F.1) \quad \xi(p_t) = \arg \min_{\xi \in T_{p_t}\mathcal{M}: \|\xi\|_{p_t} \leq \Delta_t} \widehat{m}_{p_t}(\xi) = \mathcal{H}(p_t) + \langle \text{grad}\mathcal{H}(p_t), \xi \rangle_{p_t} + \frac{1}{2} \langle H_t[\xi], \xi \rangle_{p_t},$$

where $H_t : T_{p_t}\mathcal{M} \rightarrow T_{p_t}\mathcal{M}$ is a self-adjoint linear operator that approximates the Hessian $\text{Hess}\mathcal{H}(p_t)$. Depending on how much decrease is provided by the obtained direction, we either accept or reject the trust-region step and modify the radius Δ_t .

THEOREM F.2 (Convergence of RHM-TR). *Under the same settings as in Theorem 3.7 with $L = L_0L_1 + L_2^2$, consider Algorithm 3.1 with $\xi(p_t)$ given by solving (F.1) with truncated conjugate gradient. Assume further that $\|H_t - \text{Hess}\mathcal{H}(p_t)\|_{p_t} \leq L_H \|\text{grad}\mathcal{H}(p_t)\|_{p_t}$. Let $c = \min_{i=0, \dots, t} \frac{\Delta_i}{L_0L_1}$ and $\widetilde{L} = L_H L_0 L_1 + L$. Then, the iterates p_t satisfy $\|\text{grad}f(p_t)\|_{p_t}^2 \leq (1 - \frac{1}{2} \min\{c, 1/\widetilde{L}\} \rho' \delta)^t \|\text{grad}f(p_0)\|_{p_0}^2$.*

Under an additional Lipschitzness condition on $\nabla^2 \widehat{\mathcal{H}}_p$, we can show around the global minima p^ , there exists $\theta > 0, T > 0$ such that for all $t > T$, the convergence is superlinear with $d(p_{t+1}, p^*) \leq \theta d^2(p_t, p^*)$.*

Proof. First from Assumption 2, $\|\text{grad}\mathcal{H}(p_t)\|_{p_t} = \|\text{Hess}f(p_t)[\text{grad}f(p_t)]\|_{p_t} \leq L_1 L_0$ and the operator norm of H_t is bounded as

$$\|H_t\|_{p_t} \leq \|H_t - \text{Hess}\mathcal{H}(p_t)\|_{p_t} + \|\text{Hess}\mathcal{H}(p_t)\|_{p_t} \leq L_H L_0 L_1 + L.$$

Also, the trust-region direction $\xi(p_t)$ returned by the truncated conjugate gradient method satisfies a so-called *Cauchy decrease inequality* [3, eq. (7.14)], which gives

$$\begin{aligned} \widehat{m}_{p_t}(0) - \widehat{m}_{p_t}(\xi(p_t)) &\geq \frac{1}{2} \|\text{grad}\mathcal{H}(p_t)\|_{p_t} \min \left\{ \Delta_t, \frac{\|\text{grad}\mathcal{H}(p_t)\|_{p_t}}{\|H_t\|_{p_t}} \right\} \\ &\geq \frac{1}{2} \|\text{grad}\mathcal{H}(p_t)\|_{p_t} \min \left\{ c \|\text{grad}\mathcal{H}(p_t)\|_{p_t}, \frac{\|\text{grad}\mathcal{H}(p_t)\|_{p_t}}{\|H_t\|_{p_t}} \right\} \\ &\geq \frac{1}{2} \min \left\{ c, \frac{1}{L_H L_0 L_1 + L} \right\} \|\text{grad}\mathcal{H}(p_t)\|_{p_t}^2 \\ &\geq \frac{1}{2} \min \left\{ c, \frac{1}{L_H L_0 L_1 + L} \right\} \delta \mathcal{H}(p_t). \end{aligned}$$

where the second inequality follows from the definition of c and Assumption 2 where Furthermore, from the acceptance rule,

$$\mathcal{H}(p_{t+1}) - \mathcal{H}(p_t) \leq \rho' (\widehat{m}_{p_t}(\xi(p_t)) - \widehat{m}_{p_t}(0)) \leq -\frac{1}{2} \min \left\{ c, \frac{1}{L_H L_0 L_1 + L} \right\} \rho' \delta \mathcal{H}(p_t).$$

Hence, the linear convergence is proved by recursively applying the result. The super-linear convergence simply follows from [3, Theorem 7.4.11] around any local minima. \square

Appendix G. On geodesic quadratic bilinear optimization.

We first show an important result on the orthogonality of the min-max Riemannian gradient and Riemannian gradient of the Riemannian Hamiltonian for any g-bilinear function on arbitrary manifolds.

PROPOSITION G.1. *Let $f(x, y)$ be a g-bilinear function on $\mathcal{M} = \mathcal{M}_x \times \mathcal{M}_y$. Denote $G(p) = (\text{grad}_x f(x, y), -\text{grad}_y f(x, y)) \in T_p \mathcal{M}$ for $p = (x, y) \in \mathcal{M}$ as the min-max Riemannian gradient. Then for any $p \in \mathcal{M}$, we have $\langle G(p), \text{grad} \mathcal{H}(p) \rangle_p = 0$ where $\mathcal{H}(p) = \frac{1}{2} \|\text{grad} f(p)\|_p^2$ is the Riemannian Hamiltonian of f .*

Proof. First it is aware that for any g-bilinear function, we have $\text{Hess}_x f(x, y) = \text{Hess}_y f(x, y) = 0$. Hence, from Proposition 3.2 and D.1, we show

$$\text{grad} \mathcal{H}(p) = \text{Hess} f(p)[\text{grad} f(p)] = \begin{pmatrix} \text{grad}_{yx}^2 f(x, y)[\text{grad}_y f(x, y)] \\ \text{grad}_{xy}^2 f(x, y)[\text{grad}_x f(x, y)] \end{pmatrix}.$$

Finally, we have

$$\begin{aligned} \langle G(p), \text{grad} \mathcal{H}(p) \rangle_p &= \langle \text{grad}_x f(x, y), \text{grad}_{yx}^2 f(x, y)[\text{grad}_y f(x, y)] \rangle_x \\ &\quad + \langle -\text{grad}_y f(x, y), \text{grad}_{xy}^2 f(x, y)[\text{grad}_x f(x, y)] \rangle_y \\ &= \langle \text{grad}_y f(x, y), \text{grad}_{xy}^2 f(x, y)[\text{grad}_x f(x, y)] \rangle_y \\ &\quad + \langle -\text{grad}_y f(x, y), \text{grad}_{xy}^2 f(x, y)[\text{grad}_x f(x, y)] \rangle_y = 0 \end{aligned}$$

where we apply Proposition D.2. \square

Proof of Proposition 7.1. First, the expression of geodesic curve connecting any $\mathbf{X}_0, \mathbf{X}_1 \in \mathcal{M}_{\text{SPD}}$ is given by $\gamma(t) = \mathbf{X}_0^{1/2} (\mathbf{X}_0^{-1/2} \mathbf{X}_1 \mathbf{X}_0^{-1/2})^t \mathbf{X}_0^{1/2}$. From [89, Proposition 5.7], we see $\log \det(\mathbf{X})$ is geodesic linear. That is, for the geodesic $\gamma(t)$ joining $\mathbf{X}_0, \mathbf{X}_1$ with $\gamma(0) = \mathbf{X}_0, \gamma(1) = \mathbf{X}_1$, it can be shown that $\log \det(\gamma(t)) = (1-t) \log \det(\mathbf{X}_0) + t \log \det(\mathbf{X}_1)$. It remains to show $(\log \det(\mathbf{X}))^2$ is geodesic convex, which is equivalent to show $\frac{d^2(\log \det(\gamma(t)))^2}{dt^2} \geq 0$ for all $t \in [0, 1]$ (second order characterization of geodesic convexity [89]). Specifically, we show

$$(G.1) \quad \frac{d^2(\log \det(\gamma(t)))^2}{dt^2} = 2(\log \det(\mathbf{X}_1) - \log \det(\mathbf{X}_0))^2 \geq 0.$$

The equality in (G.1) holds when $\mathbf{X}_0 \neq \mathbf{X}_1$ while $\det(\mathbf{X}_0) = \det(\mathbf{X}_1)$ and hence $\frac{d^2(\log \det(\gamma(t)))^2}{dt^2} > 0$ is not always satisfied. Similar arguments hold for g-concavity with respect to \mathbf{Y} . \square

Proof of Proposition 7.2. The Riemannian gradient of f is derived as

$$\begin{aligned} \text{grad}_{\mathbf{X}} f(\mathbf{X}, \mathbf{Y}) &= (c_l \log \det(\mathbf{Y}) + 2c_q \log \det(\mathbf{X})) \mathbf{X} \\ \text{grad}_{\mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) &= (c_l \log \det(\mathbf{X}) - 2c_q \log \det(\mathbf{Y})) \mathbf{Y}. \end{aligned}$$

Under the affine-invariant metric, the Hamiltonian is given by

$$\mathcal{H}(\mathbf{X}, \mathbf{Y}) = \frac{(4c_q^2 + c_l^2)d}{2} \left((\log \det(\mathbf{X}))^2 + (\log \det(\mathbf{Y}))^2 \right).$$

The gradient of Hamiltonian is given by $\text{grad}_{\mathbf{X}} \mathcal{H}(\mathbf{X}, \mathbf{Y}) = (4c_q^2 + c_l^2)d \log \det(\mathbf{X}) \mathbf{X}$

and $\text{grad}_{\mathbf{Y}}\mathcal{H}(\mathbf{X}, \mathbf{Y}) = (4c_q^2 + c_l^2)d \log \det(\mathbf{Y})\mathbf{Y}$. Next, we verify

$$\begin{aligned} & \frac{1}{2} \left(\|\text{grad}_{\mathbf{X}}\mathcal{H}(\mathbf{X}, \mathbf{Y})\|_{\mathbf{X}}^2 + \|\text{grad}_{\mathbf{Y}}\mathcal{H}(\mathbf{X}, \mathbf{Y})\|_{\mathbf{Y}}^2 \right) \\ &= \frac{(4c_q^2 + c_l^2)^2 d^3}{2} \left((\log \det(\mathbf{X}))^2 + (\log \det(\mathbf{Y}))^2 \right) \\ &= (4c_q^2 + c_l^2)d^2 \mathcal{H}(\mathbf{X}, \mathbf{Y}). \end{aligned}$$

In addition, from the definition of global saddle point in (1.2), the pair $(\mathbf{X}^*, \mathbf{Y}^*)$ where $\det(\mathbf{X}^*) = \det(\mathbf{Y}^*) = 1$, satisfies $f(\mathbf{X}^*, \mathbf{Y}^*) = 0$. Thus, we have

$$f(\mathbf{X}^*, \mathbf{Y}) = -c_q(\log \det(\mathbf{Y}))^2 \leq f(\mathbf{X}^*, \mathbf{Y}^*) \leq c_q(\log \det(\mathbf{X}))^2 = f(\mathbf{X}, \mathbf{Y}^*)$$

for all $\mathbf{X}, \mathbf{Y} \in \mathbb{S}_{++}^d$. Hence, the proof is complete. \square

Finally, we show that the geodesic-bilinear problem does not satisfy the min-max Riemannian PL condition on the function f . To this end, we first need to define the Riemannian min-max PL condition below.

DEFINITION G.2 (Riemannian min-max PL condition). *For a min-max problem $\min_{x \in \mathcal{M}_x} \max_{y \in \mathcal{M}_y} f(x, y)$, the objective satisfies the Riemannian min-max PL condition if for a global saddle point (x^*, y^*) , there exists a constant $\delta > 0$ such that*

$$\begin{aligned} \frac{1}{2} \|\text{grad}_x f(x', y)\|_{x'}^2 &\geq \delta (f(x', y) - f(x^*, y)), \quad \forall y \in \mathcal{M}, \\ \frac{1}{2} \|\text{grad}_y f(x, y')\|_{y'}^2 &\geq \delta (f(x, y') - f(x, y^*)), \quad \forall x \in \mathcal{M}. \end{aligned}$$

Definition G.2 is equivalent to stating that the objective $f(x, y)$ satisfies the Riemannian PL in x and $-f(x, y)$ satisfies the Riemannian PL in y . Such definition is natural as it includes geodesic strongly convex strongly concave functions as special cases.

LEMMA G.3. *The g -bilinear function $f(\mathbf{X}, \mathbf{Y}) = \log \det(\mathbf{X}) \log \det(\mathbf{Y})$ does not satisfy Definition G.2.*

Proof. We show the case for \mathbf{X} . A similar statement also holds for \mathbf{Y} . As the global saddle point $(\mathbf{X}^*, \mathbf{Y}^*)$ satisfies $\det(\mathbf{X}^*) = \det(\mathbf{Y}^*) = 1$, we have $f(\mathbf{X}^*, \mathbf{Y}) = 0$. In addition, the Riemannian gradient is $\text{grad}_{\mathbf{X}} f(\mathbf{X}', \mathbf{Y}) = \mathbf{X}' \log \det(\mathbf{Y})$ with $\|\text{grad}_{\mathbf{X}} f(\mathbf{X}', \mathbf{Y})\|_{\mathbf{X}'}^2 = (\log \det(\mathbf{Y}))^2$. On the other hand, the right-hand-side in Definition G.2 is $f(\mathbf{X}', \mathbf{Y}) - f(\mathbf{X}^*, \mathbf{Y}) = \log \det(\mathbf{X}') \log \det(\mathbf{Y})$. It is clear that $\frac{1}{2}(\log \det(\mathbf{Y}))^2$ is not necessarily larger than $\delta \log \det(\mathbf{X}') \log \det(\mathbf{Y})$ for $\delta > 0$ and for all $\mathbf{Y} \in \mathbb{S}_{++}^d$. Hence, the claim follows. \square

Appendix H. Additional experiment results.

H.1. Optimality gap for geodesic quadratic bilinear optimization. We include additional convergence results in Fig. 5 on the optimality gap for the geodesic quadratic bilinear optimization problem in Section 7.1.

H.2. Results of RHM-SGD for orthonormal GAN. We show the sample collapse of RHM-SGD in Fig. 6.

H.3. Trace-logarithm bilinear optimization. We consider the ‘bilinear’ example of [92] on the symmetric positive definite (SPD) manifold (endowed with the affine-invariant metric), i.e.,

$$f(\mathbf{X}, \mathbf{Y}) = \text{tr}(\text{Log}_{\mathbf{X}}(\mathbf{X}_0) \text{Log}_{\mathbf{Y}}(\mathbf{Y}_0))$$

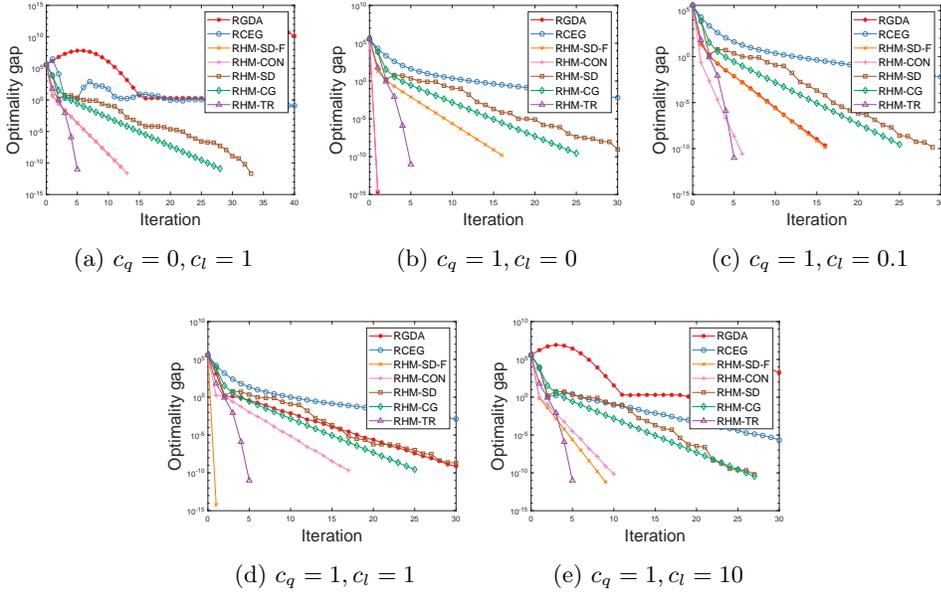


Fig. 5: Experiments comparing optimality gap on the geodesic quadratic bilinear problem (7.1) with $d = 30$, under different weights c_q, c_l . We observe that the RHM algorithms show a good rate of convergence in all the settings. In particular, RHM-SD-F and RHM-CON significantly outperforms RCEG in all the settings indicating better theoretical rates.



Fig. 6: Generated samples from RHM-SGD at $1, 2, 3 \times 10^4$ iterations from left to right. We see although RHM-SGD converges in Hamiltonian, the generated samples collapse to a single point (zoom the figures to see the single point).

for $\mathbf{X}_0, \mathbf{Y}_0 \in \mathbb{S}_{++}^d$, where $\text{Log}_{\mathbf{M}}(\mathbf{M}') = \{\mathbf{M} \log(\mathbf{M}^{-1} \mathbf{M}')\}_{\mathbf{S}}$ is the logarithm map on the SPD manifold with $\log(\cdot)$ representing the matrix principal logarithm. When the manifold is simply the Euclidean space, the logarithm map reduces to $\text{Log}_{\mathbf{M}}(\mathbf{M}') = \mathbf{M}' - \mathbf{M}$. Hence, this resembles a bilinear problem on the manifold.

For the experiment setting, we consider $\gamma = 0.2$ for RHM-CON and $\mathbf{X}_0 = \mathbf{Y}_0 = \mathbf{I}$. The convergence results are shown in Fig. 7, where we notice that both RGDA and RCEG oscillate while all the RHM algorithms are convergent. RHM-CON and RHM-SD-F converge rapidly initially but subsequently have a slow rate of convergence due to the hardness of the problem. RHM-CG, on the other hand, has a faster rate of convergence.

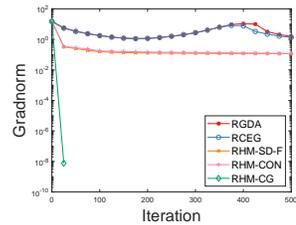


Fig. 7: Trace-logarithm bilinear problem on the SPD manifold. RGDA and RCEG diverge while RHM algorithms are convergent (though RHM with steepest descent has a slower rate of convergence).

Acknowledgments. Pawan Kumar acknowledges the support of Microsoft Academic Partnership Grant (MAPG) 2021.

REFERENCES

- [1] J. ABERNETHY, K. A. LAI, AND A. WIBISONO, *Last-iterate convergence rates for min-max optimization: Convergence of hamiltonian gradient descent and consensus optimization*, in International Conference on Algorithmic Learning Theory, vol. 132, PMLR, 2021, pp. 3–47.
- [2] P.-A. ABSIL, C. G. BAKER, AND K. A. GALLIVAN, *Trust-region methods on Riemannian manifolds*, *Found. Comput. Math.*, 7 (2007), pp. 303–330.
- [3] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization algorithms on matrix manifolds*, Princeton University Press, 2009.
- [4] L. ADOLPHS, H. DANESHMAND, A. LUCCHI, AND T. HOFMANN, *Local saddle point optimization: A curvature exploitation approach*, in International Conference on Artificial Intelligence and Statistics, PMLR, 2019, pp. 486–495.
- [5] J. AFLALO, A. BEN-TAL, C. BHATTACHARYYA, J. S. NATH, AND S. RAMAN, *Variable sparsity kernel learning*, *J. Mach. Learn. Res.*, 12 (2011), pp. 565–592.
- [6] N. AGARWAL, N. BOUMAL, B. BULLINS, AND C. CARTIS, *Adaptive regularization with cubics on manifolds*, *Math. Program.*, 188 (2021), pp. 85–134.
- [7] M. ARJOVSKY, S. CHINTALA, AND L. BOTTOU, *Wasserstein generative adversarial networks*, in International Conference on Machine Learning, PMLR, 2017, pp. 214–223.
- [8] D. BALDUZZI, S. RACANIÈRE, J. MARTENS, J. FOERSTER, K. TUYLS, AND T. GRAEPEL, *The mechanics of n -player differentiable games*, in International Conference on Machine Learning, PMLR, 2018, pp. 354–363.
- [9] N. BANSAL, X. CHEN, AND Z. WANG, *Can we gain more from orthogonality regularizations in training deep networks?*, in Advances in Neural Information Processing Systems, vol. 31, 2018.
- [10] R. BERGMANN, *Manopt.jl: Optimization on manifolds in julia*, *Journal of Open Source Software*, 7 (2022), p. 3866.
- [11] D. P. BERTSEKAS, *Constrained optimization and Lagrange multiplier methods*, Academic Press, 2014.
- [12] R. BHATIA, *Positive definite matrices*, Princeton University Press, 2009.
- [13] N. BOUMAL, *Riemannian trust regions with finite-difference Hessian approximations are globally convergent*, in International Conference on Geometric Science of Information, Springer, 2015, pp. 467–475.
- [14] N. BOUMAL, *An introduction to optimization on smooth manifolds*, Available online, May, 3 (2020).
- [15] N. BOUMAL, P.-A. ABSIL, AND C. CARTIS, *Global rates of convergence for nonconvex optimization on manifolds*, *IMA J. Numer. Anal.*, 39 (2019), pp. 1–33.
- [16] N. BOUMAL, B. MISHRA, P.-A. ABSIL, AND R. SEPULCHRE, *Manopt, a Matlab toolbox for optimization on manifolds*, *J. Mach. Learn. Res.*, 15 (2014), pp. 1455–1459.
- [17] A. BROCK, J. DONAHUE, AND K. SIMONYAN, *Large scale GAN training for high fidelity natural image synthesis*, in International Conference on Learning Representations, 2018.

- [18] B. CHASNOV, L. RATLIFF, E. MAZUMDAR, AND S. BURDEN, *Convergence analysis of gradient-based learning in continuous games*, in Uncertainty in Artificial Intelligence, PMLR, 2020, pp. 935–944.
- [19] M. COGSWELL, F. AHMED, R. GIRSHICK, L. ZITNICK, AND D. BATRA, *Reducing overfitting in deep networks by decorrelating representations*, in International Conference on Learning Representations, 2016.
- [20] A. DOUIK AND B. HASSIBI, *Manifold optimization over the set of doubly stochastic matrices: A second-order geometry*, IEEE Trans. Signal Process., 67 (2019), pp. 5761–5774.
- [21] L. EL GHAOU AND H. LEBRET, *Robust solutions to least-squares problems with uncertain data*, SIAM J. Matrix Anal. Appl., 18 (1997), pp. 1035–1064.
- [22] R. FLETCHER AND C. M. REEVES, *Function minimization by conjugate gradients*, The computer journal, 7 (1964), pp. 149–154.
- [23] K. G., *The extragradient method for finding saddle points and other problems*, Ekonomika i Matematicheskie Metody, 12 (1976), pp. 747–756.
- [24] I. GOODFELLOW, J. POUGET-ABADIE, M. MIRZA, B. XU, D. WARDE-FARLEY, S. OZAIR, A. COURVILLE, AND Y. BENGIO, *Generative adversarial nets*, in Advances in Neural Information Processing Systems, vol. 27, 2014.
- [25] A. HAN AND J. GAO, *Improved variance reduction methods for Riemannian non-convex optimization*, IEEE Trans. Pattern Anal. Mach. Intell., (2021).
- [26] A. HAN, B. MISHRA, P. JAWANPURIA, AND J. GAO, *On Riemannian optimization over positive definite matrices with the Bures-Wasserstein geometry*, in Advances in Neural Information Processing Systems, vol. 34, 2021.
- [27] I. HOREV, F. YGER, AND M. SUGIYAMA, *Geometry-aware principal component analysis for symmetric positive definite matrices*, in Asian Conference on Machine Learning, PMLR, 2016, pp. 1–16.
- [28] F. HUANG, S. GAO, AND H. HUANG, *Gradient descent ascent for min-max problems on Riemannian manifolds*, arXiv:2010.06097, (2020).
- [29] L. HUANG, X. LIU, B. LANG, A. W. YU, Y. WANG, AND B. LI, *Orthogonal weight normalization: Solution to optimization over multiple dependent Stiefel manifolds in deep neural networks*, in Thirty-Second AAAI Conference on Artificial Intelligence, 2018.
- [30] M. HUANG, S. MA, AND L. LAI, *A Riemannian block coordinate descent method for computing the projection robust Wasserstein distance*, in International Conference on Machine Learning, PMLR, 2021, pp. 4446–4455.
- [31] W. HUANG, P.-A. ABSIL, AND K. A. GALLIVAN, *A Riemannian symmetric rank-one trust-region method*, Math. Program., 150 (2015), pp. 179–216.
- [32] W. HUANG, P.-A. ABSIL, K. A. GALLIVAN, AND P. HAND, *ROPTLIB: an object-oriented C++ library for optimization on Riemannian manifolds*, ACM Trans. Math. Software, 44 (2018), pp. 1–21.
- [33] P. JAWANPURIA, M. LAPIN, M. HEIN, AND B. SCHIELE, *Efficient output kernel learning for multiple tasks*, in Advances in Neural Information Processing Systems, 2015.
- [34] P. JAWANPURIA AND B. MISHRA, *A unified framework for structured low-rank matrix learning*, in International Conference on Machine Learning, 2018.
- [35] P. JAWANPURIA AND J. S. NATH, *A convex feature learning formulation for latent task structure discovery*, in International Conference on Machine Learning, 2012.
- [36] P. JAWANPURIA, J. S. NATH, AND G. RAMAKRISHNAN, *Efficient rule ensemble learning using hierarchical kernels*, in International Conference on Machine Learning, 2011.
- [37] P. JAWANPURIA, J. S. NATH, AND G. RAMAKRISHNAN, *Generalized hierarchical kernel learning*, J. Mach. Learn. Res., 16 (2015), pp. 617–652.
- [38] P. JAWANPURIA, N. T. V. SATYA DEV, AND B. MISHRA, *Efficient robust optimal transport: formulations and algorithms*, in IEEE Conference on Decision and Control, 2021.
- [39] C. JIN, P. NETRAPALLI, AND M. JORDAN, *What is local optimality in nonconvex-nonconcave minimax optimization?*, in International Conference on Machine Learning, PMLR, 2020, pp. 4880–4889.
- [40] M. I. JORDAN, T. LIN, AND E.-V. VLATAKIS-GKARAGKOUNIS, *First-order algorithms for min-max optimization in geodesic metric spaces*, arXiv:2206.02041, (2022).
- [41] H. KARIMI, J. NUTINI, AND M. SCHMIDT, *Linear convergence of gradient and proximal-gradient methods under the Polyak-Lojasiewicz condition*, in Joint European Conference on Machine Learning and Knowledge Discovery in Databases, Springer, 2016, pp. 795–811.
- [42] H. KASAI, H. SATO, AND B. MISHRA, *Riemannian stochastic recursive gradient algorithm*, in International Conference on Machine Learning, PMLR, 2018, pp. 2516–2524.
- [43] M. KOCHUROV, R. KARIMOV, AND S. KOZLUKOV, *Geoopt: Riemannian optimization in pytorch*, in ICML 2020 Workshop on Graph Representation Learning and Beyond, 2020.

- [44] Y. LECUN, L. BOTTOU, Y. BENGIO, AND P. HAFFNER, *Gradient-based learning applied to document recognition*, Proceedings of the IEEE, 86 (1998), pp. 2278–2324.
- [45] J. LI, K. BALASUBRAMANIAN, AND S. MA, *Stochastic zeroth-order Riemannian derivative estimation and optimization*, arXiv:2003.11238, (2020).
- [46] T. LIN, C. FAN, N. HO, M. CUTURI, AND M. JORDAN, *Projection robust wasserstein distance and Riemannian optimization*, in Advances in Neural Information Processing Systems, vol. 33, 2020, pp. 9383–9397.
- [47] N. LOIZOU, H. BERARD, A. JOLICOEUR-MARTINEAU, P. VINCENT, S. LACOSTE-JULIEN, AND I. MITLIAGKAS, *Stochastic hamiltonian gradient methods for smooth games*, in International Conference on Machine Learning, PMLR, 2020, pp. 6370–6381.
- [48] D. MADRAS, E. CREAGER, T. PITASSI, AND R. ZEMEL, *Learning adversarially fair and transferable representations*, in International Conference on Machine Learning, PMLR, 2018, pp. 3384–3393.
- [49] A. MADRY, A. MAKELOV, L. SCHMIDT, D. TSIPRAS, AND A. VLADU, *Towards deep learning models resistant to adversarial attacks*, in International Conference on Learning Representations, 2018.
- [50] E. V. MAZUMDAR, M. I. JORDAN, AND S. S. SASTRY, *On finding local nash equilibria (and only local nash equilibria) in zero-sum games*, arXiv:1901.00838, (2019).
- [51] M. MEGHWANSHI, P. JAWANPURIA, A. KUNCHUKUTTAN, H. KASAI, AND B. MISHRA, *McTorch, a manifold optimization library for deep learning*, arXiv:1810.01811, (2018).
- [52] P. MERTIKOPOULOS, C. PAPADIMITRIOU, AND G. PILIOURAS, *Cycles in adversarial regularized learning*, in Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms, SIAM, 2018, pp. 2703–2717.
- [53] L. MESCHEDER, S. NOWOZIN, AND A. GEIGER, *The numerics of GANs*, in Advances in Neural Information Processing Systems, vol. 30, 2017.
- [54] B. MISHRA, N. SATYADEV, H. KASAI, AND P. JAWANPURIA, *Manifold optimization for non-linear optimal transport problems*, arXiv:2103.00902, (2021).
- [55] A. MOKHTARI, A. OZDAGLAR, AND S. PATTATHIL, *A unified analysis of extra-gradient and optimistic gradient methods for saddle point problems: Proximal point approach*, in International Conference on Artificial Intelligence and Statistics, PMLR, 2020, pp. 1497–1507.
- [56] A. MOKHTARI, A. E. OZDAGLAR, AND S. PATTATHIL, *Convergence rate of $O(1/k)$ for optimistic gradient and extragradient methods in smooth convex-concave saddle point problems*, SIAM J. Optim., 30 (2020), pp. 3230–3251.
- [57] R. D. MONTEIRO AND B. F. SVAITER, *On the complexity of the hybrid proximal extragradient method for the iterates and the ergodic mean*, SIAM J. Optim., 20 (2010), pp. 2755–2787.
- [58] R. D. MONTEIRO AND B. F. SVAITER, *Complexity of variants of tseng’s modified fb splitting and korpelevich’s methods for hemivariational inequalities with applications to saddle-point and convex optimization problems*, SIAM J. Optim., 21 (2011), pp. 1688–1720.
- [59] S.-M. MOOSAVI-DEZFOOLI, A. FAWZI, O. FAWZI, AND P. FROSSARD, *Universal adversarial perturbations*, in Proceedings of the Conference on Computer Vision and Pattern Recognition, 2017, pp. 1765–1773.
- [60] J. MÜLLER, R. KLEIN, AND M. WEINMANN, *Orthogonal Wasserstein GANs*, arXiv:1911.13060, (2019).
- [61] A. NEMIROVSKI, *Prox-method with rate of convergence $O(1/t)$ for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems*, SIAM J. Optim., 15 (2004), pp. 229–251.
- [62] J. V. NEUMANN, *Zur theorie der gesellschaftsspiele*, Math. Ann., 100 (1928), pp. 295–320.
- [63] M. NIMISHAKAVI, P. JAWANPURIA, AND B. MISHRA, *A dual framework for low-rank tensor completion*, in Advances in Neural Information Processing Systems, 2018.
- [64] J. NOCEDAL AND S. J. WRIGHT, *Numerical optimization*, Springer, 1999.
- [65] F.-P. PATY AND M. CUTURI, *Subspace robust wasserstein distances*, in International Conference on Machine Learning, 2019.
- [66] F.-P. PATY AND M. CUTURI, *Subspace robust wasserstein distances*, in International Conference on Machine Learning, PMLR, 2019, pp. 5072–5081.
- [67] X. PENNEC, *Manifold-valued image processing with SPD matrices*, in Riemannian Geometric Statistics in Medical Image Analysis, Elsevier, 2020, pp. 75–134.
- [68] G. PEYRÉ, M. CUTURI, ET AL., *Computational optimal transport: With applications to data science*, Foundations and Trends® in Machine Learning, 11 (2019), pp. 355–607.
- [69] B. T. POLYAK, *Gradient methods for minimizing functionals*, Zhurnal Vychislitel’noi Matematiki i Matematicheskoi Fiziki, 3 (1963), pp. 643–653.
- [70] L. D. POPOV, *A modification of the Arrow-Hurwicz method for search of saddle points*, Mathematical Notes of the Academy of Sciences of the USSR, 28 (1980), pp. 845–848.

- [71] A. RAKOTOMAMONJY, F. BACH, S. CANU, AND Y. GRANDVALET, *Simplekl*, J. Mach. Learn. Res., 9 (2008), pp. 2491–2521.
- [72] W. RING AND B. WIRTH, *Optimization methods on Riemannian manifolds and their application to shape space*, SIAM J. Optim., 22 (2012), pp. 596–627.
- [73] R. T. ROCKAFELLAR, *Monotone operators and the proximal point algorithm*, SIAM J. Control Optim., 14 (1976), pp. 877–898.
- [74] H. L. ROYDEN AND P. FITZPATRICK, *Real analysis*, vol. 32, Macmillan New York, 1988.
- [75] H. SATO, *A Dai–Yuan-type Riemannian conjugate gradient method with the weak Wolfe conditions*, Computational Optimization and Applications, 64 (2016), pp. 101–118.
- [76] H. SATO, *Riemannian conjugate gradient methods: General framework and specific algorithms with convergence analyses*, arXiv:2112.02572, (2021).
- [77] H. SATO, *Riemannian Optimization and Its Applications*, Springer, 2021.
- [78] H. SATO, H. KASAI, AND B. MISHRA, *Riemannian stochastic variance reduced gradient algorithm with retraction and vector transport*, SIAM J. Optim., 29 (2019), pp. 1444–1472.
- [79] F. SCHÄFER AND A. ANANDKUMAR, *Competitive gradient descent*, Advances in Neural Information Processing Systems, 32 (2019).
- [80] D. SHI, J. GAO, X. HONG, S. BORIS CHOY, AND Z. WANG, *Coupling matrix manifolds assisted optimization for optimal transport problems*, Mach. Learn., 110 (2021), pp. 533–558.
- [81] R. SINKHORN, *A relationship between arbitrary positive matrices and doubly stochastic matrices*, The Annals of Mathematical Statistics, 35 (1964), pp. 876–879.
- [82] M. SION, *On general minimax theorems.*, Pacific J. Math., 8 (1958), pp. 171–176.
- [83] O. SMIRNOV, *TensorFlow RiemOpt: a library for optimization on Riemannian manifolds*, arXiv:2105.13921, (2021).
- [84] S. SRA AND R. HOSSEINI, *Conic geometric optimization on the manifold of positive definite matrices*, SIAM J. Optim., 25 (2015), pp. 713–739.
- [85] L. STACHÓ, *Minimax theorems beyond topological vector spaces*, Acta Sci. Math.(Szeged), 42 (1980), pp. 157–164.
- [86] J. TOWNSEND, N. KOEP, AND S. WEICHWALD, *Pymanopt: A python toolbox for optimization on manifolds using automatic differentiation*, J. Mach. Learn. Res., 17 (2016), pp. 1–5, <http://jmlr.org/papers/v17/16-177.html>.
- [87] P. TSENG, *On linear convergence of iterative methods for the variational inequality problem*, J. Comput. Appl. Math., 60 (1995), pp. 237–252.
- [88] C. UDRISTE, *Convex functions and optimization methods on Riemannian manifolds*, vol. 297, Springer Science & Business Media, 2013.
- [89] N. K. VISHNOI, *Geodesic convex optimization: Differentiation on manifolds, geodesics, and convexity*, arXiv:1806.06373, (2018).
- [90] J. WANG, Y. CHEN, R. CHAKRABORTY, AND S. X. YU, *Orthogonal convolutional neural networks*, in Proceedings of the Conference on Computer Vision and Pattern Recognition, 2020, pp. 11505–11515.
- [91] H. ZHANG, S. J. REDDI, AND S. SRA, *Riemannian SVRG: Fast stochastic optimization on Riemannian manifolds*, in Advances in Neural Information Processing Systems, vol. 29, 2016.
- [92] P. ZHANG, J. ZHANG, AND S. SRA, *Minimax in geodesic metric spaces: Sion’s theorem and algorithms*, arXiv:2202.06950, (2022).
- [93] P. ZHOU, X.-T. YUAN, AND J. FENG, *Faster first-order methods for stochastic non-convex optimization on Riemannian manifolds*, in International Conference on Artificial Intelligence and Statistics, PMLR, 2019, pp. 138–147.