

Riemannian Langevin Monte Carlo schemes for sampling PSD matrices with fixed rank ^{*}

Tianmin Yu[†], Shixin Zheng[‡], Jianfeng Lu[§], Govind Menon[†], AND Xiangxiong Zhang[‡]

Abstract. This paper introduces two explicit schemes to sample matrices from Gibbs distributions on $\mathcal{S}_+^{n,p}$, the manifold of real positive semi-definite (PSD) matrices of size $n \times n$ and rank p . Given an energy function $\mathcal{E} : \mathcal{S}_+^{n,p} \rightarrow \mathbb{R}$ and certain Riemannian metrics g on $\mathcal{S}_+^{n,p}$, these schemes rely on an Euler-Maruyama discretization of the Riemannian Langevin equation (RLE) with Brownian motion on the manifold. We present numerical schemes for RLE under two fundamental metrics on $\mathcal{S}_+^{n,p}$: (a) the metric obtained from the embedding of $\mathcal{S}_+^{n,p} \subset \mathbb{R}^{n \times n}$; and (b) the Bures-Wasserstein metric corresponding to quotient geometry. We also provide examples of energy functions with explicit Gibbs distributions that allow numerical validation of these schemes.

Key words. Langevin dynamics, sampling schemes, Bures-Wasserstein metric, Burer-Monteiro, embedded geometry, positive semi-definite matrices, Riemannian optimization

MSC codes.

1. Introduction.

1.1. Problem statement. Consider the space of real, symmetric positive semi-definite matrices with size $n \times n$ and rank p , denoted by

$$(1.1) \quad \mathcal{S}_+^{n,p} = \{X \in \mathbb{R}^{n \times n} \mid X = X^T, X \succeq 0, \text{rank}(X) = p\}.$$

Given an energy $\mathcal{E} : \mathcal{S}_+^{n,p} \rightarrow \mathbb{R}$ and a parameter $\beta > 0$ referred to as the inverse temperature, our goal is to sample efficiently from the Gibbs distribution

$$(1.2) \quad \rho_\beta(X) = \frac{1}{Z_\beta} e^{-\beta \mathcal{E}(X)} \rho_{\text{ref}}(X), \quad Z_\beta = \int_{\mathcal{S}_+^{n,p}} e^{-\beta \mathcal{E}(X')} \rho_{\text{ref}}(X') dX'.$$

Gibbs measures must be defined with respect to a base measure. In this work, we equip the space $\mathcal{S}_+^{n,p}$ with a Riemannian metric g and choose $\rho_{\text{ref}}(X) dX = \sqrt{\det g(X)} dX$ to be the canonical volume form associated to the metric g . This volume form is expressed in coordinates for the metrics studied in this paper in Section 4.

This sampling problem is related to the optimization problem $\min_{X \in \mathcal{S}_+^{n,p}} \mathcal{E}(X)$ since in the limit $\beta \rightarrow \infty$ the Gibbs distribution concentrates at the global minima of $\mathcal{E}(X)$. Minimization problems over the space $\mathcal{S}_+^{n,p}$ arise in many areas, especially semidefinite programming and machine learning, and have been studied extensively. Gibbs distributions originate in statistical physics, while the sampling problem may

^{*}

Funding: S.Z. and X.Z. are supported by NSF DMS-2208518. J.L. is supported in part by NSF DMS-2012286 and DMS-2309378. G.M. is supported in part by NSF DMS-2107205

[†]Division of Applied Mathematics, Brown University, Providence, RI (tianmin_yu@brown.edu, govind_menon@brown.edu).

[‡]Department of Mathematics, Purdue University, West Lafayette, IN (zheng513@purdue.edu, zhan1966@purdue.edu).

[§]Departments of Mathematics, Physics, and Chemistry, Duke University, Durham, NC (jianfeng@math.duke.edu).

also be seen as a stochastic variant of the optimization problem. For these reasons, the sampling problem has a broad range of applications; see Section 1.5 below.

The main contribution of this paper are efficient sampling schemes for ρ_β based on Langevin dynamics. Our approach builds on the geometric theory of optimization; in particular, we extend Riemannian optimization on $\mathcal{S}_+^{n,p}$ [34, 38] to Gibbs sampling as follows. In [34] it was recognized that two commonly used gradient descent schemes over $\mathcal{S}_+^{n,p}$ are time discretizations of *Riemannian* gradient flows, where $\mathcal{S}_+^{n,p}$ is equipped with the two natural Riemannian metrics listed below. We combine this observation with the theory of Brownian motion on Riemannian manifolds to obtain Riemannian Langevin equations and explicit sampling schemes.

The reader unfamiliar with these concepts should note that while the abstract theory serves to guide our work, the schemes presented in this paper may be implemented without requiring a complete understanding of the underlying theory. Further, while this paper is focused on the two numerical schemes below, the underlying framework can be used to extend other Riemannian gradient descent schemes to sampling schemes for the Gibbs measure. The new phenomenon that arises is the interplay between Brownian motion and curvature in the Riemannian Langevin equation. This interplay has been studied in depth by two of the authors (TY and GM) and their co-workers in recent papers for geometries used in optimization and physics [20, 28, 29].

1.2. Two Riemannian metrics on $\mathcal{S}_+^{n,p}$. Given $X \in \mathcal{S}_+^{n,p}$, let $X = YY^T$ be a low-rank decomposition where $Y \in \mathbb{R}^{n \times p}$. We use two fundamental metrics on $\mathcal{S}_+^{n,p}$ obtained from this parametrization, from the Euclidean metric for either the variable X or the variable Y through the use of Riemannian embedding and Riemannian submersion respectively. These are the two most natural ways of defining metrics on $\mathcal{S}_+^{n,p}$.

The flat metric for X corresponds to the embedded geometry of $\mathcal{S}_+^{n,p}$ in the Euclidean space $\mathbb{R}^{n \times n}$ [34]. Precisely, we consider the natural Riemannian embedding $\mathcal{S}_+^{n,p} \hookrightarrow \mathbb{R}^{n \times n}$ and use the Frobenius norm on $\mathbb{R}^{n \times n}$ to define a metric on $\mathcal{S}_+^{n,p}$. Denote it by g_E , then $g_E(A, B) = \text{Tr}(A^T B)$ for any two square matrices A, B in the tangent space of $\mathcal{S}_+^{n,p}$, where Tr denotes the trace of a matrix.

On the other hand, we may also use the flat geometry on Y to define a metric on $\mathcal{S}_+^{n,p}$. We observe that if $YY^T = X$, then it is also true that $\tilde{Y}\tilde{Y}^T = X$ where $\tilde{Y} = YO$ and $O \in \mathcal{O}_p$, the orthogonal group of dimension p . Thus, we may identify $\mathcal{S}_+^{n,p} \simeq \mathbb{R}_*^{n \times p} / \mathcal{O}_p$, as a quotient space, with a quotient map

$$\begin{aligned} \pi : \mathbb{R}_*^{n \times p} &\rightarrow \mathbb{R}_*^{n \times p} / \mathcal{O}_p \\ Y &\mapsto [Y] = \{YO \mid O \in \mathcal{O}_p\}. \end{aligned}$$

Here $\mathbb{R}_*^{n \times p}$ denotes full rank matrices.

The quotient space structure can be enhanced with a Riemannian metric through the use of Riemannian submersion. Roughly, the metric for X corresponds to the metric for Y in a manner that respects the splitting of the tangent space at Y into the space of the group action and its complement. If $\mathbb{R}_*^{n \times p}$ is equipped with Euclidean metric, then the metric induced by the submersion is often called the Bures-Wasserstein metric on $\mathcal{S}_+^{n,p} \simeq \mathbb{R}_*^{n \times p} / \mathcal{O}_p$, denoted by g_{BW} (see [2, 26, 27]).

1.3. Langevin dynamics and the Riemannian Langevin equation. We now explain how Langevin equations may be defined intrinsically on $(\mathcal{S}_+^{n,p}, g)$.

Let us first recall the Langevin equation on \mathbb{R}^n . Assume given a potential or energy function $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$ and let W_t denote the standard Wiener process on \mathbb{R}^n . The Langevin equation for the potential \mathcal{E} is the Itô differential equation

$$(1.3) \quad dx_t = -\nabla \mathcal{E}(x_t) dt + \sqrt{\frac{2}{\beta}} dW_t.$$

The Fokker-Planck equation describes the evolution of the probability density of x_t . With $\rho(x, t) dx = \mathbb{P}(x_t \in (x, x + dx))$, we have

$$(1.4) \quad \partial_t \rho = \frac{1}{\beta} \Delta \rho + \nabla \cdot (\rho \nabla \mathcal{E}).$$

The Gibbs density (with reference density being uniform with respect to Lebesgue measure) is the unique equilibrium of equation (1.4) under natural growth assumptions on the energy \mathcal{E} as $|x| \rightarrow \infty$.

The Langevin equation immediately yields a numerical scheme for (approximate) sampling from the Gibbs distribution. Fix a step size $\Delta t > 0$, let $t_k = k\Delta t$, $k = 0, 1, \dots$, and let x_k denote the numerical approximation to (1.3) at time t_k . The Euler-Maruyama scheme to approximate equation (1.3), also known as Langevin Monte Carlo in the statistics literature, is

$$(1.5) \quad x_{k+1} = x_k - \Delta t \nabla \mathcal{E}(x_k) + \sqrt{\frac{2\Delta t}{\beta}} \xi_k,$$

where $\xi_k = (\xi_k^1, \dots, \xi_k^n)$ is an i.i.d. sequence of standard Gaussian vectors in \mathbb{R}^n . This scheme is explicit. In order to extend it to sampling from (1.2) we must understand how to modify the Langevin equation on the Riemannian manifold $(\mathcal{S}_+^{n,p}, g)$.

First, the term $\nabla \mathcal{E}$ must be replaced by the Riemannian gradient, written as $\text{grad } \mathcal{E}$. The more subtle modification of equation (1.3) concerns the noise. The natural analogy is to replace the Wiener process W_t on \mathbb{R}^n with Brownian motion on the Riemannian manifold $(\mathcal{S}_+^{n,p}, g)$ at inverse temperature β , denoted $\mathbf{B}_t^{g, \beta}$. This yields the (formal) Riemannian Langevin equation on $(\mathcal{S}_+^{n,p}, g)$

$$(1.6) \quad d\mathbf{X}_t = -\text{grad } \mathcal{E}(\mathbf{X}_t) dt + d\mathbf{B}_t^{g, \beta}.$$

This equation is only formal because stochastic differential equations on manifolds must be defined using the Stratonovich formulation in order to ensure coordinate independence (Itô differentials do not satisfy the chain rule, while Stratonovich differentials do) [16, 19]. On the other hand, Itô differential equations are convenient for analysis as well as simulation. Thus, in formulating the Riemannian Langevin equation, it is necessary to first formulate the appropriate Stratonovich equation and then compute the deterministic Itô–Stratonovich correction. A central observation in our work is that this correction term is due to curvature and is explicitly computable for several Riemannian geometries relevant to optimization [17, 20, 28, 29].

1.4. Riemannian Langevin Monte Carlo sampling schemes. For the two metrics considered in this paper, the Itô–Stratonovich correction due to curvature may also be computed explicitly, yielding the SDEs in Section 2. The rigorous analysis of these SDEs is presented in the companion paper [36], and we focus on numerical algorithms in this paper. The Euler-Maruyama approximation to these SDEs yields the numerical sampling schemes listed below.

The SDEs also admit other numerical approximations. We have chosen the Euler-Maruyama schemes because these schemes are fully explicit, simple to state, implement and numerically validate. They are generalizations of the popular unadjusted Langevin Monte Carlo for sampling in Euclidean spaces. Further, these schemes reduce to deterministic Riemannian gradient descent methods in the limit $\beta \rightarrow \infty$.

1.4.1. Scheme E for the embedded geometry. For the embedded manifold $(\mathcal{S}_+^{n,p}, g_E)$, the scheme is

$$(1.7) \quad X_{k+1} = P_{\mathcal{S}_+^{n,p}} \left[X_k - \Delta t \operatorname{grad} \mathcal{E}(X_k) + Q_k \left(\sqrt{\frac{2\Delta t}{\beta}} \begin{bmatrix} B_{11} & B_{12} \\ B_{12}^T & 0 \end{bmatrix} + \frac{\Delta t}{\beta} \sum_{i=1}^p \frac{1}{\lambda_i} \begin{bmatrix} 0 & 0 \\ 0 & I_{n-p} \end{bmatrix} \right) Q_k^T \right],$$

where $P_{\mathcal{S}_+^{n,p}}$ is the Euclidean projection to $\mathcal{S}_+^{n,p}$, and $X_k = Q_k \Lambda Q_k^T$ is the full SVD of $X_k = \mathcal{S}_+^{n,p}$ with eigenvalues $\lambda_1 \geq \dots \geq \lambda_p > 0$. The entries of B_{12} are i.i.d. drawn from $\sqrt{\frac{1}{2}}\mathcal{N}(0,1)$. The entries of the symmetric B_{11} are defined as follows: the diagonal entries are i.i.d. drawn from $\mathcal{N}(0,1)$, and off-diagonal entries are $b_{ij} = b_{ji} \sim \sqrt{\frac{1}{2}}\mathcal{N}(0,1)$. When $\beta = \infty$, equation (1.7) reduces to $X_{k+1} = P_{\mathcal{S}_+^{n,p}}(X_k - \Delta t \operatorname{grad} \mathcal{E}(X_k))$, which is the Riemannian gradient descent on $(\mathcal{S}_+^{n,p}, g_E)$, see [1, 38]. We refer to (1.7) as Scheme E.

In this scheme, the term $\frac{\Delta t}{\beta} \sum_{i=1}^p \frac{1}{\lambda_i} \begin{bmatrix} 0 & 0 \\ 0 & I_{n-p} \end{bmatrix}$ in equation (1.7) is the correction due to the mean curvature of the embedding of $\mathcal{S}_+^{n,p} \hookrightarrow \mathbb{R}^{n \times n}$.

1.4.2. Scheme BW for the Bures-Wasserstein metric. For the quotient manifold $(\mathbb{R}_*^{n \times p} / \mathcal{O}_p, g_{BW})$, the scheme is

$$(1.8) \quad Y_{k+1} = Y_k - \Delta t 2\nabla \mathcal{E}(Y_k Y_k^T) Y_k + \sqrt{\frac{2\Delta t}{\beta}} B_k + \frac{\Delta t}{\beta} U_k \left[\sum_{j:j \neq i} \frac{\sigma_i}{\sigma_i^2 + \sigma_j^2} \right]_{ii} V_k^T,$$

where B_k is n -by- p matrix with entries being i.i.d. standard Gaussian, $Y_k = U_k \Sigma_k V_k^T \in \mathbb{R}^{n \times p}$ is the compact SVD with singular values σ_i , and $\left[\sum_{j:j \neq i} \frac{\sigma_i}{\sigma_i^2 + \sigma_j^2} \right]_{ii}$ is the diagonal matrix whose i -th diagonal entry is $\sum_{j:j \neq i} \frac{\sigma_i}{\sigma_i^2 + \sigma_j^2}$. We refer to (1.8) as Scheme BW. The Riemannian Langevin Monte Carlo scheme (1.8) can be viewed as a natural extension of Burer-Monteiro gradient descent method

$$(1.9) \quad Y_{k+1} = Y_k - \Delta t 2\nabla \mathcal{E}(Y_k Y_k^T) Y_k,$$

which is the simplest low-rank gradient descent method for minimizing $\mathcal{E}(X)$ under the constraint $X \in \mathcal{S}_+^{n,p}$. It is clear that as $\beta \rightarrow \infty$, (1.8) reduces to (1.9). The Burer-Monteiro gradient descent method is equivalent to a Riemannian gradient descent method on the quotient manifold $\mathbb{R}_*^{n \times p} / \mathcal{O}_p$ with Bures-Wasserstein metric, see [38].

1.4.3. Gibbs distribution sampling and numerical validation. While the Gibbs distribution always has the same density function $e^{-\beta\mathcal{E}}$ with respect to ρ_{ref} , the reference density ρ_{ref} depends on the metric. Thus, the two schemes (1.7) and (1.8), generate samples for two different probability distributions. In order to validate our schemes, we choose energy functions that allow an explicit computation of these densities for both metric. These energy functions yield matrix integrals of independent analytic interest. They also allow side-to-side benchmarking for different Gibbs samplers on $\mathcal{S}_+^{n,p}$. We demonstrate the efficiency of sampling from these Gibbs distributions numerically. Further analysis on convergence to equilibrium as $t \rightarrow \infty$ using the Bakry-Emery criterion is considered in the companion paper [36].

Finally, while we do not discuss their convergence and efficiency approximating SDE as the step-size $\Delta t \rightarrow 0$; this is possible following existing approximation results [22, 8, 24].

1.5. Some applications and related work.

1.5.1. Applications of PSD matrices. Positive semi-definite (PSD) fixed rank matrices arise in many problems such as distance matrices [33] and covariance matrices in statistics, and have been used in applications including kernels in machine learning [30], semidefinite optimization [4], quantum information, etc. Riemannian optimization algorithms over $\mathcal{S}_+^{n,p}$ under different metrics have been well studied, e.g., see [34, 18, 27, 38] and references therein.

1.5.2. Langevin dynamics and Monte Carlo schemes on manifolds.

There is an extensive literature on Langevin dynamics in statistics and related areas, with interest in nonconvex optimization [6, 7], as well as machine learning such as generative models [12].

In recent years, there has been interest in studying Langevin diffusion and Monte Carlo Markov Chain (MCMC) schemes on manifolds [9, 10, 15, 3, 5, 37, 31, 13, 23, 24]. In this paper, we are interested in Riemannian Langevin Monte Carlo schemes on $\mathcal{S}_+^{n,p}$.

In the statistics literature, manifold Langevin schemes have been studied in [15, 5]. However, these schemes apply only to simpler embedded manifolds $\mathcal{M} \subset \mathbb{R}^n$ with explicit geodesics such as the sphere and Stiefel manifolds. The above schemes do not directly apply to the manifold $\mathcal{S}_+^{n,p}$, even for the embedded geometry. In [37], a sampling scheme using projection to surface is constructed; however, this is not a Langevin scheme.

In general, a Langevin scheme can be used for either optimization [35, 24], or Monte Carlo type numerical integration, which is common in Bayesian statistic. For optimization, stochastic optimization by Langevin dynamics with simulated annealing is an established approach [25]. In [6], underdamped Langevin schemes are shown to be much more efficient than the overdamped case (1.5). For sampling, Metropolis-adjusted Langevin algorithm [15] is often used. For simplicity, we focus on the simple schemes (1.7) and (1.8) without considering any of simulated annealing, underdamped Langevin, or Metropolis-adjustment, to which it is possible to extend our schemes. Though the Riemannian optimization on $\mathcal{S}_+^{n,p}$ can be easily extended to Hermitian PSD matrices of fixed rank [38], we remark that such an extension for Langevin dynamics would be significantly different.

1.6. Organization of the paper. In Section 2, we state the explicit formulae for the SDE (1.6) and Gibbs measure on the manifold $\mathcal{S}_+^{n,p}$ under two metrics g_E and g_{BW} . We then derive the schemes (1.7) and (1.8) in Section 3. The energy functions and Gibbs distributions used to benchmark the schemes are presented in Section 4. The numerical results are studied in Section 5.

2. Riemannian Langevin equations on $\mathcal{S}_+^{n,p}$. In this section, we state the Itô form of the Riemannian Langevin equation (1.6) for both Riemannian geometries studied in this paper. The theoretical basis for these SDEs is discussed at greater depth in [36]. The main ideas are as follows: (a) the abstract theory of Brownian motion on Riemannian manifolds is used to define the Riemannian Langevin equation in Stratonovich form for the metrics g_E and g_{BW} on $\mathcal{S}_+^{n,p}$; (b) the Itô-Stratonovich conversion rule is used to compute the associated Itô form of these SDEs and it is observed that the Itô-Stratonovich correction term corresponds to mean curvature. This approach yields the SDEs below. These SDEs are used to develop numerical schemes in Section 3.

2.1. The Riemannian Langevin equation for embedded geometry $(\mathcal{S}_+^{n,p}, g_E)$. ■

Let $X \in \mathcal{S}_+^{n,p}$ have the compact SVD $X = U\Lambda U^T$ with $U \in \mathbb{R}^{n \times p}$. Let $U_\perp \in \mathbb{R}^{n \times (n-p)}$ be a matrix with columns orthonormal to columns of U . The tangent space of $\mathcal{S}_+^{n,p}$ at $X = U\Lambda U^T \in \mathcal{S}_+^{n,p}$ is given by [34, 38]:

$$(2.1) \quad T_X \mathcal{S}_+^{n,p} = \left\{ [U \quad U_\perp] \begin{bmatrix} H & K^T \\ K & 0 \end{bmatrix} \begin{bmatrix} U^T \\ U_\perp^T \end{bmatrix} : \forall K \in \mathbb{R}^{(n-p) \times p}, \forall H \in \mathbb{R}^{p \times p}, H^T = H \right\}.$$

The induced metric g_E by the embedding $\mathcal{S}_+^{n,p} \hookrightarrow \mathbb{R}^{n \times n}$ is then defined as

$$g_E(A, B) = \text{Tr}(A^T B), \quad \forall A, B \in T_X \mathcal{S}_+^{n,p},$$

which is the Frobenius inner product for two matrices.

Equation (1.6) describes the evolution of a point $\mathbf{X}_t \in \mathcal{S}_+^{n,p}$ in abstract terms. We now rewrite it in a simpler equivalent form describing the evolution of the entries of the matrix entries $\{(X_t)_{ij}\}_{i,j=1}^n$ representing \mathbf{X}_t . Let us write $X = U\Lambda U^T$ for the compact singular value decomposition (SVD) of X . We further assume that the singular values $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ are written in decreasing order. We suppress the subscript t in the following equations, though the reader should note that U and Λ depend on X_t .

Then we find that the law of X_t is determined by the Itô differential equation

$$(2.2) \quad dX_t = -\text{grad } \mathcal{E}(X_t)dt + \sqrt{\frac{2}{\beta}} dW_t^{n,p,X_t} + \frac{1}{\beta} H(X_t)dt.$$

In this equation, the stochastic forcing W_t^{n,p,X_t} is the orthogonal projection of white noise in $\mathbb{R}^{n \times n}$ onto $T_{X_t} \mathcal{S}_+^{n,p}$. Precisely, given W_t^i for $1 \leq i \leq n$ and $W_t^{i,j}$ for $1 \leq i <$

$j \leq n$ independent standard one-dimensional Wiener process, we set

$$dW_t^{n,p,X_t} = [U \quad U_\perp] \begin{bmatrix} dW_t^1 & \cdots & \frac{1}{\sqrt{2}}dW_t^{1,p} & \frac{1}{\sqrt{2}}dW_t^{1,p+1} & \cdots & \frac{1}{\sqrt{2}}dW_t^{1,n} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\sqrt{2}}dW_t^{1,p} & \cdots & dW_t^p & \frac{1}{\sqrt{2}}dW_t^{p,p+1} & \cdots & \frac{1}{\sqrt{2}}dW_t^{p,n} \\ \frac{1}{\sqrt{2}}dW_t^{1,p+1} & \cdots & \frac{1}{\sqrt{2}}dW_t^{p,p+1} & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\sqrt{2}}dW_t^{1,n} & \cdots & \frac{1}{\sqrt{2}}dW_t^{p,n} & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} U^T \\ U_\perp^T \end{bmatrix},$$

The term $H(X_t)$ is the mean curvature of the embedding $\mathcal{S}_+^{n,p} \rightarrow \mathbb{R}^{n \times n}$. We adopt the convention in geometric analysis: the mean curvature is defined as the trace of the second fundamental form of the embedding. Explicitly, we have

$$(2.3) \quad H(X_t) = \left(\sum_{i=1}^p \frac{1}{\lambda_i} \right) [U \quad U_\perp] \begin{bmatrix} 0_{p \times p} & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & I_{n-p} \end{bmatrix} \begin{bmatrix} U^T \\ U_\perp^T \end{bmatrix}.$$

The following feature of equation (2.2) is fundamental. The stochastic forcing is the naive projection of white noise in the ambient space $\mathbb{R}^{n \times n}$ onto $T_{X_t} \mathcal{S}_+^{n,p}$. Intuitively, when one uses the Euler-Maruyama discretization, the role of this term is to update X_t by taking unbiased random steps in any direction in the tangent space. However, Itô calculus has a subtle interplay with the geometry of the embedding, and in order to keep X_t on the manifold $\mathcal{S}_+^{n,p}$, it is necessary to include the correction term given by the mean curvature.

2.2. The Riemannian Langevin equation for Bures-Wasserstein geometry ($\mathcal{S}_+^{n,p}, g_{BW}$). The manifold $\mathcal{S}_+^{n,p}$ can also be viewed as a quotient manifold $\mathbb{R}_*^{n \times p} / \mathcal{O}_p$, for which the noncompact Stiefel manifold $\mathbb{R}_*^{n \times p}$ is called the *total space*. Denote the natural projection as

$$\pi : \mathbb{R}_*^{n \times p} \rightarrow \mathbb{R}_*^{n \times p} / \mathcal{O}_p.$$

For any $Y \in \mathbb{R}_*^{n \times p}$, the equivalence class containing Y is

$$[Y] = \pi^{-1}(\pi(Y)) = \{YO \mid O \in \mathcal{O}_p\},$$

which is an embedded submanifold of $\mathbb{R}_*^{n \times p}$ (see e.g., [1, Prop. 3.4.4]). The tangent space of $[Y]$ at Y is therefore a subspace of $T_Y \mathbb{R}_*^{n \times p}$ called the *vertical space* at Y , and is denoted by $\mathcal{V}_Y = \{Y\Omega \mid \Omega^T = -\Omega, \Omega \in \mathbb{R}^{p \times p}\}$, see [38].

Define

$$\begin{aligned} \theta : \mathbb{R}_*^{n \times p} &\rightarrow \mathcal{S}_+^{n,p} \\ Y &\mapsto YY^T. \end{aligned}$$

Then θ is invariant under the equivalence relation and induces a bijection $\tilde{\theta}$ on $\mathbb{R}_*^{n \times p} / \mathcal{O}_p$ such that $\theta = \tilde{\theta} \circ \pi$. For any function $\mathcal{E}(X)$ defined on $\mathcal{S}_+^{n,p}$, there is a function F defined on $\mathbb{R}_*^{n \times p}$ that induces \mathcal{E} : for any $X = YY^T \in \mathcal{S}_+^{n,p}$, $F(Y) := \mathcal{E} \circ \theta(Y) = \mathcal{E}(YY^T)$. This is summarized in the diagram below:

$$\begin{array}{ccccc} \mathbb{R}_*^{n \times p} & & & & \\ \downarrow \pi & \dashrightarrow^{\theta := \tilde{\theta} \circ \pi} & & & \\ \mathbb{R}_*^{n \times p} / \mathcal{O}_p & \xleftrightarrow{\tilde{\theta}} & \mathcal{S}_+^{n,p} & \xrightarrow{\mathcal{E}} & \mathbb{R} \end{array}$$

In particular, $\mathcal{S}_+^{n,p}$ is diffeomorphic to $\mathbb{R}_*^{n \times p} / \mathcal{O}_p$ under $\tilde{\theta}$, see [38]. For any $Y \in \mathbb{R}_*^{n \times p}$, the flat metric for the total space $\mathbb{R}_*^{n \times p}$, correction term corresponds to mean curvature.

$$g(a, b) = \text{Tr}(a^T b), \forall a, b \in T_Y \mathbb{R}_*^{n \times p} = \mathbb{R}^{n \times p}$$

induces a metric on the quotient manifold $\mathbb{R}_*^{n \times p} / \mathcal{O}_p$, which is called Bures-Wasserstein metric, see [27, 26, 38]. Another way to understand the Bures-Wasserstein metric at $X \in \mathcal{S}_+^{n,p} \simeq \mathbb{R}_*^{n \times p} / \mathcal{O}_p$ is via the map θ :

$$(2.4) \quad \begin{aligned} g_{BW}(A, B) &= \text{Tr}(ab^T) \quad \forall A, B \in T_X \mathcal{S}_+^{n,p}, a, b \in T_Y \mathbb{R}_*^{n \times p} \\ \text{s.t. } d\theta(Y)[a] &= A, d\theta(Y)[b] = B, a, b \in \text{Ker}(d\theta(Y))^\perp \end{aligned}$$

where X has decomposition $X = YY^T$, $d\theta(Y)[a] = Ya^T + aY^T$ is the differential of θ at Y , and $a \in \text{Ker}(d\theta(Y))^\perp \Leftrightarrow Y^T a = a^T Y$.

The Riemannian Langevin equation is now determined by the geometry of Riemannian submersion. We must obtain an Itô differential equation for Y_t , such that $X_t = Y_t Y_t^T$ is a matrix that has the same law as the solution to (1.6) in $(\mathcal{S}_+^{n,p}, g_{BW})$.

In comparison with equation (2.2), we see that the natural choice for white noise driving Y_t is white noise in $\mathbb{R}^{n \times p}$. This is the stochastic differential dW_t , where $W_t = \{W_t^{ij}\}_{1 \leq i \leq n, 1 \leq j \leq p}$ consists of np independent standard one-dimensional Wiener processes. However, as in equation (2.2) we must include a deterministic correction. This correction corresponds to mean curvature again, but in a more subtle way than (2.2). The equivalence class of Y such that $X = YY^T$ is a group orbit of \mathcal{O}_p embedded within $\mathbb{R}^{n \times p}$. The logarithm of the volume of this group orbit constitutes a natural Boltzmann entropy. It may be computed explicitly, and we find

$$(2.5) \quad S(Y) = \frac{1}{2} \sum_{i=1}^p \sum_{j=i+1}^p \log(\sigma_i^2 + \sigma_j^2)$$

where $\{\sigma_i\}_{i=1}^p$ are singular values of Y . It is known that $\nabla S(Y)$ is the mean curvature of the group orbit in $\mathbb{R}^{n \times p}$ [32, p.3505].

We then have the following Itô differential equation for Y_t such that $X_t = Y_t Y_t^T$ has the same law as the solution to (1.6).

$$(2.6) \quad dY_{ij} = -\frac{\partial \mathcal{E}(YY^T)}{\partial Y_{ij}} dt + \sqrt{\frac{2}{\beta}} dW_t^{ij} - \frac{1}{\beta} \frac{\partial S(Y)}{\partial Y_{ij}} dt, \quad 1 \leq i \leq n, 1 \leq j \leq p.$$

The correction term can be explicitly computed using the following

LEMMA 1. *If $Y \in \mathbb{R}_*^{n \times p}$ has SVD as $Y = Q\Sigma P^T$ with singular values σ_i , then the gradient of the correction term S is given by $\nabla S(Y) = Q\tilde{\Sigma}P^T$ where $\tilde{\Sigma}$ is a diagonal matrix with diagonal entries $\sum_{j \neq 1} \frac{\sigma_1}{\sigma_1^2 + \sigma_j^2}, \sum_{j \neq 2} \frac{\sigma_2}{\sigma_2^2 + \sigma_j^2}, \dots, \sum_{j \neq p} \frac{\sigma_p}{\sigma_p^2 + \sigma_j^2}$.*

3. Two Riemannian Langevin Monte Carlo schemes. To get a simple Riemannian Langevin Monte Carlo sampling scheme, we only consider convenient discretization and approximation methods, which can be easily and efficiently implemented. For the Brownian motion term, we consider the most straightforward and simplest discretization of the SDEs (2.2) and (2.6), i.e., the Euler-Maruyama type discretization.

One extra complication from the manifold constraint is how to approximate the exponential map. For optimization algorithms on Riemannian manifolds [1], retraction, which is at least a first order approximation to the exponential map, is often used. For instance, for approximating an ODE $\frac{d}{dt}\mathbf{X} = -\text{grad } \mathcal{E}(\mathbf{X})$ on a manifold \mathcal{M} , with any retraction operator $\mathcal{R}_{\mathcal{M}}$ mapping to \mathcal{M} , a simple forward Euler type approximation, or equivalently the Riemannian gradient descent method, is given by

$$\mathbf{X}_{k+1} = \mathcal{R}_{\mathcal{M}}[\mathbf{X}_{k+1} - \Delta t \text{grad } \mathcal{E}(\mathbf{X}_k)].$$

In particular, when combining the Euler-Maruyama type discretization for SDE and the simple Riemannian gradient descent by retraction, we get the two simple Riemannian Langevin Monte Carlo schemes as follows.

3.1. Scheme E for the embedded geometry.

3.1.1. The Riemannian gradient. For a given energy function $\mathcal{E}(X)$, its Riemannian gradient $\text{grad } \mathcal{E}(X)$ of at $X \in \mathcal{S}_+^{n,p}$, is the Euclidean projection of the Euclidean gradient $\nabla \mathcal{E}(X) \in \mathbb{R}^{n \times n}$ defined as $[\nabla \mathcal{E}(X)]_{ij} = \frac{\partial}{\partial X_{ij}} \mathcal{E}(X)$, onto the tangent space $T_X \mathcal{S}_+^{n,p}$, see [1, 34, 38]. It is straightforward to verify that $\nabla \mathcal{E}(X)$ is a symmetric matrix for any differentiable \mathcal{E} and any $X \in \mathcal{S}_+^{n,p}$. For any given $X \in \mathcal{S}_+^{n,p}$, let $X = U\Lambda U^T$ be its compact SVD. Let $P_U = UU^T$ and $P_{U_{\perp}} = U_{\perp}U_{\perp}^T = I - UU^T$. By derivations in [38], $\text{grad } \mathcal{E}(X)$ can be computed and represented as

$$\begin{aligned} \text{grad } \mathcal{E}(X) &= [U \quad U_{\perp}] \begin{bmatrix} U^T \nabla \mathcal{E}(X) U & U^T \nabla \mathcal{E}(X) U_{\perp} \\ U_{\perp}^T \nabla \mathcal{E}(X) U & 0 \end{bmatrix} \begin{bmatrix} U^T \\ U_{\perp}^T \end{bmatrix} \\ &= P_U \nabla \mathcal{E}(X) P_U + P_{U_{\perp}} \nabla \mathcal{E}(X) P_U + P_U \nabla \mathcal{E}(X) P_{U_{\perp}}. \end{aligned}$$

The compact implementation of computing $\text{grad } \mathcal{E}(X)$ is given in Algorithm 3.1.

Algorithm 3.1 Compact computation of the Riemannian gradient $\text{grad } \mathcal{E}(X)$

Require: The compact SVD of $X \in \mathcal{S}_+^{n,p}$: $X = U\Lambda U^T$
Ensure: $\text{grad } \mathcal{E}(X) = UHU^T + U_p U^T + UU_p^T \in T_X \mathcal{S}_+^{n,p}$
 $T \leftarrow \nabla \mathcal{E}(X) U$
 $H \leftarrow U^T T$
 $U_p \leftarrow T - UH$

3.1.2. The retraction by projection. Let $\mathcal{S}^{n \times n}$ denote symmetric matrices, then the Euclidean projection $P_{\mathcal{S}_+^{n,p}} : \mathcal{S}^{n \times n} \rightarrow \mathcal{S}_+^{n,p}$ is a convenient retraction operator, see [1, 34, 38]. A straightforward implementation is given in Algorithm 3.2.

Algorithm 3.2 Computation of the retraction $P_{\mathcal{S}_+^{n,p}}(X + Z)$

Require: the compact SVD of X : $X = U\Lambda U^T \in \mathcal{S}_+^{n,p}$, $Z \in \mathcal{S}^{n \times n}$.
Ensure: $P_{\mathcal{S}_+^{n,p}}(X + Z) = Q_+ \Lambda_+ Q_+^T \in \mathcal{S}_+^{n,p}$.
 $(Q_+, \Lambda_+) = \text{svd}(X + Z)$
 $U_+ \leftarrow Q_+(\cdot, 1:p) \quad \Lambda_+ \leftarrow \Lambda_+(1:p, 1:p)$

3.1.3. A Riemannian Langevin Monte Carlo scheme. For approximating the SDE (2.2) on $(S_+^{n,p}, g_E)$, with the retraction operator and Euler-Maruyama method for SDE, we have the scheme (1.7), which can be more explicitly written as (3.1)

$$X_{k+1} = P_{S_+^{n,p}} \left(\begin{bmatrix} U & U_\perp \end{bmatrix} \begin{bmatrix} \Lambda - \Delta t U^T \nabla \mathcal{E}(X_k) U + \sqrt{\frac{2\Delta t}{\beta}} B_{11} & -\Delta t U^T \nabla \mathcal{E}(X_k) U_\perp + \sqrt{\frac{2\Delta t}{\beta}} B_{12} \\ -\Delta t U_\perp^T \nabla \mathcal{E}(X_k) U + \sqrt{\frac{2\Delta t}{\beta}} B_{12}^T & \frac{\Delta t}{\beta} \sum_{i=1}^p \frac{1}{\lambda_i} I_{n-p} \end{bmatrix} \begin{bmatrix} U^T \\ U_\perp^T \end{bmatrix} \right),$$

where $X_k = U \Lambda U^T$ is the compact SVD of $X_k \in S_+^{n,p}$ with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0$. The third term in the right hand side is the white noise term in the tangent space $T_{X_k} S_+^{n,p}$. Entries of $B_{12} \in \mathbb{R}^{p \times (n-p)}$ are i.i.d drawn from $\sqrt{\frac{1}{2}} \mathcal{N}(0, 1)$, and $B_{11} \in \mathbb{R}^{p \times p}$ are defined as follows.

$$(3.2) \quad B_{11} = \begin{bmatrix} \mathcal{N}(0, 1) & & & \\ & \ddots & b_{ij} & \\ & & b_{ji} & \ddots \\ & & & & \mathcal{N}(0, 1) \end{bmatrix}$$

with $b_{ij} = b_{ji} \sim \sqrt{\frac{1}{2}} \mathcal{N}(0, 1)$. The implementation details of the scheme (1.7) are given as follows in the Algorithm 3.3.

Algorithm 3.3 The Riemannian Langevin Monte Carlo scheme (1.7) for $(S_+^{n,p}, g_E)$

Require: initial iterate $X_1 \in S_+^{n,p}$; full SVD of X_1 : $X_1 = Q_1 \Lambda_1 Q_1^T$

1: **for** $k = 1, 2, \dots, N$ **do**

2: Compute Riemannian gradient

$$\xi_k := \text{grad } \mathcal{E}(X_k)$$

▷ See Algorithm 3.1

3: Compute noise term

$$B = \sqrt{\frac{2\Delta t}{\beta}} \begin{bmatrix} B_{11} & B_{12} \\ B_{12}^T & 0 \end{bmatrix} + \frac{\Delta t}{\beta} \sum_{i=1}^p \frac{1}{\lambda_i} \begin{bmatrix} 0 & 0 \\ 0 & I_{n-p} \end{bmatrix}$$

4: Obtain the new iterate by retraction

$$X_{k+1} = P_{S_+^{n,p}}(X_k - \Delta t \xi_k + Q_k B Q_k^T)$$

▷ See Algorithm 3.2

5: **end for**

REMARK 2. *The mean curvature correction term is necessary for avoiding rank deficient samples in the following sense. A sampling scheme on $S_+^{n,p}$ might generate a sample X with a rank numerically close to $p - 1$, and the mean curvature correction term in the scheme (1.7) would be huge if $\lambda_p \rightarrow 0$, thus it will force iterate X_k to stay away from the boundary of $S_+^{n,p}$.*

REMARK 3. *Notice that the complexity of computing SVD of $X + Z$ in Algorithm 3.2 would be $\mathcal{O}(n^3)$ in a naive implementation. For a Riemannian gradient method, if $Z \in T_{X_k} S_+^{n,p}$, a compact implementation of computing $P_{S_+^{n,p}}(X + Z)$ in [38] is only $\mathcal{O}(np^2) + \mathcal{O}(p^3)$, which is no longer possible for the Langevin Monte Carlo scheme (1.7) due to the mean curvature correction term in the normal space. On the other hand, if Lanczos type algorithm is used for computing to top p eigen-componenes of $X + Z$, it seems possible to explore the special structure in (3.1) to find a more efficient implementation, but we do not consider a more compact implementation in this paper.*

3.2. Scheme BW for the Bures-Wasserstein metric.

3.2.1. The Riemannian gradient and a simple retraction operator.

Given a smooth energy function $\mathcal{E}(X)$ defined on $\mathcal{S}_+^{n,p}$, the corresponding function h on $\mathbb{R}_*^{n \times p} / \mathcal{O}_p$ satisfies

$$(3.3) \quad \begin{aligned} h : \mathbb{R}_*^{n \times p} / \mathcal{O}_p &\rightarrow \mathbb{R} \\ \pi(Y) &\mapsto \mathcal{E}(\tilde{\beta}(\pi(Y))) = \mathcal{E}(\beta(Y)) = \mathcal{E}(YY^T). \end{aligned}$$

Observe that the function $F(Y) := \mathcal{E}(YY^T)$ satisfies $F(Y) = h \circ \pi(Y) = \mathcal{E} \circ \beta(Y)$. The Riemannian gradient of h at $\pi(Y)$ is a tangent vector in $T_{\pi(Y)} \mathbb{R}_*^{n \times p} / \mathcal{O}_p$. The next theorem is given in [1, Section 3.6.2], showing that the horizontal lift of $\text{grad } h(\pi(Y))$ can be obtained from the Riemannian gradient of F defined on $\mathbb{R}_*^{n \times p}$.

THEOREM 4. *The horizontal lift of the gradient of h at $\pi(Y)$ is the Riemannian gradient of F at Y . That is,*

$$\overline{\text{grad } h(\pi(Y))}_Y = \text{grad } F(Y).$$

For the Bures-Wasserstein metric, the following result is proven in [38]:

PROPOSITION 5. *Let \mathcal{E} be a smooth real-valued function defined on $\mathcal{S}_+^{n,p}$ and let $F : \mathbb{R}_*^{n \times p} \rightarrow \mathbb{R} : Y \mapsto \mathcal{E}(YY^T)$. Assume $YY^T = X$. Then the Riemannian gradient of F is given by*

$$\text{grad } F(Y) = 2\nabla \mathcal{E}(YY^T)Y$$

where $\nabla \mathcal{E}(\cdot)$ is the gradient of \mathcal{E} w.r.t. X .

In [26, Prop. A.8], the relationship between the horizontal lifts of the quotient tangent vector $\xi_{\pi(Y)}$ lifted at different representatives in $[Y]$ is given:

LEMMA 6. *Let η be a vector field on $\mathbb{R}_*^{n \times p} / \mathcal{O}_p$, and let $\bar{\eta}$ be the horizontal lift of η . Then for each $Y \in \mathbb{R}_*^{n \times p}$, we have*

$$\bar{\eta}_{YO} = \bar{\eta}_Y O$$

for all $O \in \mathcal{O}_p$.

The retraction on the quotient manifold $\mathbb{R}_*^{n \times p} / \mathcal{O}_p$ can be defined using the retraction on the total space $\mathbb{R}_*^{n \times p}$. For any $A \in T_Y \mathbb{R}_*^{n \times p}$ and a step size $\tau > 0$,

$$\bar{R}_Y(\tau A) := Y + \tau A,$$

is a retraction on $\mathbb{R}_*^{n \times p}$ if $Y + \tau A$ remains full rank, which is ensured for small enough τ . Then Lemma 6 indicates that \bar{R} satisfies the conditions of [1, Prop. 4.1.3], which implies that

$$(3.4) \quad R_{\pi(Y)}(\tau \eta_{\pi(Y)}) := \pi(\bar{R}_Y(\tau \bar{\eta}_Y)) = \pi(Y + \tau \bar{\eta}_Y)$$

defines a retraction on the quotient manifold $\mathbb{R}_*^{n \times p} / \mathcal{O}_p$ for a small enough step size $\tau > 0$.

Finally, we give an example of what these results imply by considering the Riemannian gradient descent method for minimizing $\mathcal{E}(X)$ over $(\mathcal{S}_+^{n,p}, g_{BW})$. With the simple

retraction (3.4), the Riemannian gradient descent method for minimizing the function $h[\pi(Y)]$ on $\mathbb{R}_*^{n \times p} / \mathcal{O}_p$ is given by

$$Y_{k+1} = Y_k - \Delta t 2 \nabla \mathcal{E}(Y_k Y_k^T) Y_k,$$

which is the simple Burer-Monteiro gradient descent method for minimizing $\mathcal{E}(X)$ over $\mathcal{S}_+^{n,p}$. See Section 5.1 in [38] for details.

3.2.2. A simple Riemannian Langevin Monte Carlo scheme. With the Euler-Maruyama discretization for SDE (2.6), and the simple retraction and Riemannian gradient as given previously, a simple Riemannian Langevin Monte Carlo scheme for approximating the Riemannian SDE (2.6) on the Riemannian manifold $(\mathcal{S}_+^{n,p}, g_{BW})$ can be given as

$$(3.5) \quad Y_{k+1} = Y_k - \Delta t 2 \nabla \mathcal{E}(Y_k Y_k^T) Y_k + \sqrt{\frac{2\Delta t}{\beta}} B_k + \frac{\Delta t}{\beta} U \left[\sum_{j:j \neq i} \frac{\sigma_i}{\sigma_i^2 + \sigma_j^2} \right]_{ii} V^T,$$

where B_k is n -by- p matrix with i.i.d. $\mathcal{N}(0, 1)$ entries and $Y_k = U \Sigma V^T$ is the compact SVD of Y with singular values $\sigma_i > 0$ for $i = 1, 2, \dots, p$.

Notice that all operations are performed in the space of size $n \times p$. For finding compact SVD of Y , one can first compute QR decomposition of Y , which costs $\mathcal{O}(np^2) + \mathcal{O}(p^3)$. Then compute SVD of size $p \times p$, which is $\mathcal{O}(p^3)$. So the complexity of this scheme is $\mathcal{O}(np^2) + \mathcal{O}(p^3)$ for each iteration. For large n and small p , Scheme BW should be cheaper than Scheme E in each iteration, but they generate different samples for different Gibbs distributions which depend on the metric, i.e., Scheme BW cannot replace Scheme E for generating Gibbs distribution defined by embedded geometry.

4. Examples with analytical formulae. In this section, we provide a few examples with analytical formulae so that they can be used in numerical experiments for testing the two schemes (3.1) and (3.5) on the Gibbs distribution.

For the rest of this section, $X = Q \Lambda Q^T \in \mathcal{S}_+^{n,p}$ denotes the full SVD with descending eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0$.

4.1. Scalar random variables. Let X be a random variable satisfying the Gibbs distribution on $\mathcal{S}_+^{n,p}$ with dimension $N = np - \frac{p(p-1)}{2}$ under either metric, then X is a matrix-valued random variable. For convenience, we consider a scalar random variable $D = D(X)$ as a function of $X \in \mathcal{S}_+^{n,p}$, e.g., $D = \|X\|_F$ where $\|\cdot\|_F$ is the matrix Frobenius norm.

We consider the distribution function for the scalar random variable D :

$$(4.1) \quad \Pr[D < d] = \frac{1}{Z_\beta} \int_{U_d} e^{-\beta \mathcal{E}} dV, \quad Z_\beta = \int_{\mathcal{M}} e^{-\beta \mathcal{E}} dV,$$

where $U_d := \{X \in \mathcal{S}_+^{n,p} | D(X) < d\}$ is the domain of integral. For simplicity we only consider symmetric functions such that the random variable D , the energy function \mathcal{E} , and the volume form are all invariant under the group action by the orthogonal group \mathcal{O}_n . We consider an energy function \mathcal{E} satisfying $\mathcal{E}(X) = \mathcal{E}(OXO^T)$, $\forall O \in \mathcal{O}_n$, so that Gibbs distribution function only depends on the spectrum of X when considering (4.1) with $D = \|X\|_F = \sqrt{\lambda_1^2 + \dots + \lambda_p^2}$. Since \mathcal{O}_n is an isometry group for both

metrics g_E and g_{BW} , the volume form dV in the two cases is also invariant under \mathcal{O}_n action.

Notice that Q and Λ can be used as coordinates of the manifold $\mathcal{S}_+^{n,p}$. The volume form expressed by coordinates Q and Λ is given by

$$(4.2) \quad dV = \sqrt{\det g} \left(\prod_{i=1}^p d\lambda_i \right) d\mu_{\mathcal{O}_n},$$

where $\mu_{\mathcal{O}_n}$ is the Haar measure on \mathcal{O}_n , and g is the matrix of metric g_E or g_{BW} expressed under coordinate Q and λ . For g_E its determinant $\det g$ is

$$(4.3) \quad \det g = \left(\prod_{1 \leq i < j \leq p} |\lambda_i - \lambda_j|^2 \right) \left(\prod_{1 \leq i \leq p} \lambda_i^{2(n-p)} \right),$$

and for g_{BW} it is

$$(4.4) \quad \det g = \left(\prod_{1 \leq i < j \leq p} \frac{|\lambda_i - \lambda_j|^2}{\lambda_i + \lambda_j} \right) \left(\prod_{1 \leq i \leq p} \lambda_i^{(n-p)} \right).$$

So for g_E the distribution $\Pr[D < d]$ is expressed as

$$(4.5) \quad \begin{aligned} \Pr[D < d] &= \frac{1}{Z_\beta} \int_{\|X\|_F < d} e^{-\beta \mathcal{E}} dV \\ &\propto \int_{\substack{\sum_{i=1}^p \lambda_i^2 < d^2 \\ \lambda_i > 0, i=1, \dots, p}} e^{-\beta \mathcal{E}(\lambda_1, \dots, \lambda_p)} \left(\prod_{1 \leq i < j \leq p} |\lambda_i - \lambda_j| \right) \left(\prod_{1 \leq i \leq p} \lambda_i^{n-p} \right) d\lambda_1 \cdots d\lambda_p, \end{aligned}$$

where we have used the fact that the integrand does not depend on the coordinate $Q \in \mathcal{O}_n$, so the integral of $\mu_{\mathcal{O}_n}$ only provides a constant coefficient. As we could always renormalize $\Pr[D < d]$ by considering the quotient $\frac{\Pr[D < d]}{\Pr[D < \infty]}$, we only need the dependence of the integral on parameter d .

Similarly, for the Bures-Wasserstein metric g_{BW} we have

$$(4.6) \quad \Pr[D < d] \propto \int_{\substack{\sum_{i=1}^p \lambda_i^2 < d^2 \\ \lambda_i > 0, i=1, \dots, p}} e^{-\beta \mathcal{E}(\lambda_1, \dots, \lambda_p)} \left(\prod_{1 \leq i < j \leq p} \frac{|\lambda_i - \lambda_j|}{\sqrt{\lambda_i + \lambda_j}} \right) \left(\prod_{1 \leq i \leq p} \lambda_i^{\frac{n-p}{2}} \right) d\lambda_1 \cdots d\lambda_p$$

Next we give a few energy functions.

4.2. Example I: $\mathcal{E}(X) = \frac{1}{2}\|X\|_F^2$. This is the simplest example. Using the general expression above, for embedded geometry g_E we have

$$\begin{aligned}
\Pr[D < d] &\propto \int_{\substack{\sum_{i=1}^p \lambda_i^2 < d^2 \\ \lambda_i > 0, i=1, \dots, p}} e^{-\frac{\beta}{2} \sum_{i=1}^p \lambda_i^2} \left(\prod_{1 \leq i < j \leq p} |\lambda_i - \lambda_j| \right) \left(\prod_{1 \leq i \leq p} \lambda_i^{n-p} \right) d\lambda_1 \cdots d\lambda_p \\
&= \int_0^d e^{-\frac{\beta}{2} \rho^2} \rho^{N-1} \left(\int_{S_+^{p-1}} \prod_{1 \leq i < j \leq p} |\omega_i - \omega_j| \prod_{i=1}^p |\omega_i|^{n-p} \prod_{i=1}^p d\omega \right) d\rho \\
&= \left(\int_{S_+^{p-1}} \prod_{1 \leq i < j \leq p} |\omega_i - \omega_j| \prod_{i=1}^p |\omega_i|^{n-p} \prod_{i=1}^p d\omega \right) \int_0^d e^{-\frac{\beta}{2} \rho^2} \rho^{N-1} d\rho \\
(4.7) \quad &\propto \int_0^d e^{-\frac{\beta}{2} \rho^2} \rho^{N-1} d\rho,
\end{aligned}$$

where we have used the spherical coordinate for $(\lambda_1, \dots, \lambda_p) = \rho\omega$, with $\rho = \sqrt{\sum_{i=1}^p \lambda_i^2}$ being the radius and $\omega \in S_+^{p-1} = S^{p-1} \cap \mathbb{R}_+^p$ being the coordinate on the positive orthant of unit sphere.

For g_{BW} , similarly we have

$$(4.8) \quad \Pr[D < d] \propto \int_0^d e^{-\beta \rho^2} \rho^{\frac{N}{2}-1} d\rho.$$

Now we can see that $\beta D^2 = \beta \|X\|_F^2$ is subject to $\chi^2(N)$ distribution for the embedded metric g_E , and $\chi^2(\frac{N}{2})$ distribution for the Bures-Wasserstein metric.

4.3. Example II: $\mathcal{E}(X) = \text{Tr}(X \log X)$. We consider the von Neumann entropy

$$\mathcal{E}(X) = \text{Tr}(X \log X) = \sum_{i=1}^p \lambda_i \log \lambda_i$$

and construct a more interesting example. The minimizers of $\mathcal{E}(X) = \text{Tr}(X \log X)$ on $\mathcal{S}_+^{n,p}$ are matrices $X \in \mathcal{S}_+^{n,p}$ with spectrum $\lambda_1 = \dots = \lambda_p = e^{-1}$.

The random variable we consider is still $D = \|X\|_F$. Since $\mathcal{E}(X) = \text{Tr}(X \log X) = \sum_{i=1}^p \lambda_i \log \lambda_i$ only depends on spectrum, the argument in the previous section about

integral on \mathcal{O}_n still applies. Similar to (4.7), for g_E we have

$$\begin{aligned} \Pr(D < d) &= \int_{\substack{\sum_{i=1}^p \lambda_i^2 < d^2 \\ \lambda_i > 0, i=1, \dots, p}} e^{-\beta \sum_{i=1}^p \lambda_i \log \lambda_i} \prod_{1 \leq i < j \leq p} |\lambda_i - \lambda_j| \prod_{i=1}^p |\lambda_i|^{n-p} \prod_{i=1}^p d\lambda_i \\ &= \int_{\substack{\sum_{i=1}^p \lambda_i^2 < d^2 \\ \lambda_i > 0, i=1, \dots, p}} \prod_{1 \leq i < j \leq p} |\lambda_i - \lambda_j| \prod_{i=1}^p |\lambda_i|^{n-p-\beta\lambda_i} \prod_{i=1}^p d\lambda_i, \end{aligned}$$

and for g_{BW} we have

$$\Pr(D < d) = \int_{\substack{\sum_{i=1}^p \lambda_i^2 < d^2 \\ \lambda_i > 0, i=1, \dots, p}} \prod_{1 \leq i < j \leq p} \frac{|\lambda_i - \lambda_j|}{\sqrt{\lambda_i + \lambda_j}} \prod_{i=1}^p |\lambda_i|^{\frac{n-p}{2} - \beta\lambda_i} \prod_{i=1}^p d\lambda_i.$$

Although we do not have a closed expression for both cases, such integrals can be easily approximated by an accurate quadrature.

4.4. Example III: $\mathcal{E}(X) = \frac{1}{2} \|X - A\|_F^2$. We consider a quadratic function $\mathcal{E}(X) = \frac{1}{2} \|X - A\|_F^2$ where $A \in \mathcal{S}_+^{n,p}$ with $D = \|X - A\|_F$. In this example, \mathcal{O}_n symmetry does not hold, and we can only make an estimate of the distribution function.

The random variable D we are considering now is $D = \|X - A\|_F$, its distribution function is evaluated as

$$(4.9) \quad \Pr(D < d) \propto \int_{U_d} e^{-\frac{\beta}{2} D^2} dV$$

where $U_d = \{X \in \mathcal{S}_+^{n,p} | D(X) < d\}$. Using delta function, formally we can simplify the integral to

$$\begin{aligned} (4.10) \quad \Pr(D < d) &\propto \int_{\mathcal{M}} \mathbf{1}_{\{D < d\}} e^{-\frac{\beta}{2} D^2} dV \\ &= \int_{\mathcal{M}} \left(\int_0^\infty \mathbf{1}_{\{\rho < d\}} e^{-\frac{\beta}{2} \rho^2} \delta(D - \rho) d\rho \right) dV \\ &= \int_0^\infty \mathbf{1}_{\{\rho < d\}} e^{-\frac{\beta}{2} \rho^2} \left(\int_{\mathcal{M}} \delta(D - \rho) dV \right) d\rho \\ &= \int_0^d e^{-\frac{\beta}{2} \rho^2} \left(\int_{\mathcal{M}} \frac{d}{d\rho} \mathbf{1}_{\{D = \rho\}} dV \right) d\rho \\ &= \int_0^d e^{-\frac{\beta}{2} \rho^2} \frac{d}{d\rho} \left(\int_{\mathcal{M}} \mathbf{1}_{\{D = \rho\}} dV \right) d\rho \\ &= \int_0^d e^{-\frac{\beta}{2} \rho^2} \frac{d}{d\rho} V_D(\rho) d\rho \end{aligned}$$

(4.11)

where $V_D(\rho) = \int_{\mathcal{M}} \mathbf{1}_{\{D < \rho\}} dV = \int_{D < \rho} dV$.

In general it is difficult to calculate $\int_{D < \rho} dV$, but we consider the following approximation. Consider the volume of the ball $B_A^{n,p}(r) = B_A(r) \cap \mathcal{S}_+^{n,p}$, where

$$B_A(r) = \{X \in \mathcal{S}^{n \times n} : \|X - A\|_F < r\}.$$

It is difficult to compute $\text{Vol}(B_A^{n,p}(r))$, but we propose the following estimate, for fixed $A \in \mathcal{S}_+^{n,p}$:

$$(4.12) \quad \text{Vol}(B_{cA}^{n,p}(r)) \approx \alpha r^N, \quad c \gg 1,$$

where α is a constant that does not depend on r , N is the dimension of $\mathcal{S}_+^{n,p}$. For g_E , α is exactly the volume of unit ball in \mathbb{R}^N , while for g_{BW} , α depends on dimension N and $A \in \mathcal{S}_+^{n,p}$.

For the embedded geometry, the approximation (4.12) can be justified by the following arguments:

1. The second fundamental form \mathbf{II}_{cA} of the manifold is vanishing for fixed A and $c \rightarrow \infty$. See [36].
2. The Riemannian curvature tensor of ambient space $\mathcal{S}^{n \times n}$ is 0. Applying the Gauss equation [11, Prop 3.1] we can express the Riemannian curvature tensor R of $(\mathcal{S}_+^{n,p}, g_E)$ in terms of its second fundamental form \mathbf{II} :

$$(4.13) \quad \langle R(\mathbf{x}, \mathbf{y})\mathbf{z}, \mathbf{w} \rangle = -\langle \mathbf{II}(\mathbf{x}, \mathbf{z}), \mathbf{II}(\mathbf{y}, \mathbf{w}) \rangle + \langle \mathbf{II}(\mathbf{x}, \mathbf{w}), \mathbf{II}(\mathbf{y}, \mathbf{z}) \rangle,$$

$$\mathbf{x}, \mathbf{y}, \mathbf{z}, \mathbf{w} \in T\mathcal{S}_+^{n,p}, \quad \langle U, V \rangle = \text{Tr}(UV^T) \text{ is the metric in } \mathcal{S}^{n \times n}.$$

Thus, with vanishing \mathbf{II} we have vanishing Riemannian curvature tensor, and zero sectional curvature.

3. Vanishing extrinsic curvature and intrinsic curvature means that the neighborhood is approximately an Euclidean space, so the ball $B_{cA}^{n,p}(r)$ is approximately just a ball in \mathbb{R}^N and has volume αr^N , with α being the volume of a unit ball.

For the g_{BW} metric, following similar arguments, we can get the same approximation (4.12). We emphasize that the approximation (4.12) is accurate only if c is large enough. Putting all this together, when A has eigenvalues $\lambda_1 \geq \dots \geq \lambda_p \gg 1$, we have the following

$$(4.14) \quad \Pr(D < d) \propto \int_{D < d} e^{-\frac{\beta}{2} D^2} dV = \int_0^d e^{-\frac{\beta}{2} \rho^2} \frac{d}{d\rho} \left(\int_{D < \rho} dV \right) d\rho \propto \int_0^d e^{-\frac{\beta}{2} \rho^2} \rho^{N-1} d\rho,$$

where \propto stands for *being approximately proportional to*.

4.5. MCMC numerical integration. It is well known that MCMC can be used for integrating a function numerically, and that one of the main advantages is that the convergence rate is independent of the dimension. Both schemes in this paper are MCMC type sampling schemes on the manifold. Suppose we have generated samples X_i satisfying the Gibbs distribution on the manifold, e.g.,

$$X_i \sim \frac{1}{Z_\beta} e^{-\beta \mathcal{E}(X)} dV_g,$$

where $Z_\beta = \int_{S_+^{n,p}} e^{-\beta\mathcal{E}(X)} dV$ is an unknown normalization factor and dV is the volume form depending on the metric. Then for approximating the integral of a nice function $f(X)$ on the same manifold $\int_{S_+^{n,p}} f(X)dV$, we can use

$$(4.15) \quad \frac{1}{m} \sum_{i=1}^m f(X_i) e^{\beta\mathcal{E}(X_i)} \approx \frac{\int_{S_+^{n,p}} f(X) dV}{\int_{S_+^{n,p}} e^{-\beta\mathcal{E}(X)} dV} = \frac{1}{Z_\beta} \int_{S_+^{n,p}} f(X) dV,$$

because each $f(X_i) e^{\beta\mathcal{E}(X_i)}$ is a random variable with expectation

$$\mathbb{E} \left[f(X_i) e^{\beta\mathcal{E}(X_i)} \right] = \frac{1}{Z_\beta} \int_{S_+^{n,p}} f(X) e^{\beta\mathcal{E}(X)} e^{-\beta\mathcal{E}(X)} dV,$$

and the left hand side is a random variable with expectation

$$\mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m f(X_i) e^{\beta\mathcal{E}(X_i)} \right] = \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[f(X_i) e^{\beta\mathcal{E}(X_i)} \right] = \frac{1}{Z_\beta} \int_{S_+^{n,p}} f(X) dV,$$

where the expectation $\mathbb{E}[\cdot]$ is taken w.r.t. Gibbs distribution under corresponding metric.

So using the generated samples X_i , we can approximate the integral $\int_{S_+^{n,p}} f(X) dV$ up to a constant Z_β that does not depend on $f(X)$. Notice that the additional advantage of Monte Carlo type quadrature on a manifold is that we do not need to know what dV is. On the other hand, Z_β cannot be approximated by the same approach. Though we do not consider any specific application for numerical integration, equation (4.15) can be used as one way to validate the Riemannian Langevin Monte Carlo schemes.

For the following special functions, it is possible to calculate exact integrals. For the energy function $\mathcal{E}(X) = \frac{1}{2} \|X\|_F^2$, and a special integrand $f(X) = \|X\|_F^k e^{-\frac{\alpha}{m} \|X\|_F^m}$ with $k > -N$, $m > 2$, $\alpha > 0$, using the results in 4.2, the distribution of $D = \|X\|_F$ is

$$(4.16) \quad \text{for metric } g_E : \Pr[D < d] \propto \int_0^d e^{-\frac{\beta}{2} \rho^2} \rho^{N-1} d\rho,$$

$$(4.17) \quad \text{for metric } g_{BW} : \Pr[D < d] \propto \int_0^d e^{-\frac{\beta}{2} \rho^2} \rho^{\frac{N}{2}-1} d\rho,$$

so the integral on the manifold could be expressed by expectation of a random variable, which leads to

$$\begin{aligned}
\text{for } g_E : \frac{1}{Z_\beta} \int_{\mathcal{S}_+^{n,p}} f(X) dV &= \mathbb{E}[f(X) e^{\frac{\beta}{2} \|X\|_F^2}] = \mathbb{E}[D^k e^{-\frac{\alpha}{m} D^m} e^{\frac{\beta}{2} D^2}] \\
(4.18) \quad &= \frac{\int_0^\infty \rho^k e^{-\frac{\alpha}{m} \rho^m + \frac{\beta}{2} \rho^2} \rho^{N-1} e^{-\frac{\beta}{2} \rho^2} d\rho}{\int_0^\infty \rho^{N-1} e^{-\frac{\beta}{2} \rho^2} d\rho} = \frac{\frac{1}{m} (\alpha/m)^{-\frac{k+N}{m}} \Gamma((k+N)/m)}{\frac{1}{2} (\beta/2)^{-N/2} \Gamma(N/2)}
\end{aligned}$$

$$\begin{aligned}
\text{for } g_{BW} : \frac{1}{Z_\beta} \int_{\mathcal{S}_+^{n,p}} f(X) dV &= \mathbb{E}[f(X) e^{\frac{\beta}{2} \|X\|_F^2}] = \mathbb{E}[D^k e^{-\frac{\alpha}{m} D^m} e^{\frac{\beta}{2} D^2}] \\
(4.19) \quad &= \frac{\int_0^\infty \rho^k e^{-\frac{\alpha}{m} \rho^m + \frac{\beta}{2} \rho^2} \rho^{\frac{N}{2}-1} e^{-\frac{\beta}{2} \rho^2} d\rho}{\int_0^\infty \rho^{\frac{N}{2}-1} e^{-\frac{\beta}{2} \rho^2} d\rho} = \frac{\frac{1}{m} (\alpha/m)^{-\frac{k+N/2}{m}} \Gamma((k+N/2)/m)}{\frac{1}{2} (\beta/2)^{-N/4} \Gamma(N/4)}.
\end{aligned}$$

5. Numerical tests. In this section we test the samples generated by the two Riemannian Langevin Monte Carlo schemes (3.1) and (3.5) on the examples constructed in the previous section. The samples are generated by the following procedure: we run the iterative schemes (3.1) or (3.5) for sufficiently many \tilde{m} iterations then take the last m iterates as the samples for the Gibbs distribution. Both \tilde{m} and m should be chosen such that the $(\tilde{m} - m)$ -th iterate has already reached equilibrium e.g., \tilde{m} is 6,000,000 and m is 5,000,000 for specially chosen energy functions and parameters β .

Now suppose we have generated samples $X_i \in \mathcal{S}_+^{n,p}$ ($i = 1, \dots, m$) for either metric. In order to test or show the numerical convergence to the Gibbs distribution, we will consider two kinds of numerical tests.

The first kind of tests is to test on the scalar random variable $D(X) = \|X\|_F$ or $D(X) = \|X - A\|_F$ as described in Section 4. Then we compare the cumulative distribution function (CDF) of the random variable D with its empirical CDF calculated from the MCMC samples.

Denote the true CDF of D by $F_D(t) := \Pr(D \leq t)$. The empirical CDF of samples is

$$\hat{F}_D(t) := \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{D(X_i) \leq t},$$

where $\mathbb{1}_{D(X_i) \leq t}$ takes value 1 if $D(X_i) \leq t$, and value 0 if otherwise. The Kolmogorov–Smirnov test statistic (K-S statistic) is defined by

$$(5.1) \quad KS_D := \sup_t |F_D(t) - \hat{F}_D(t)|.$$

In our numerical tests, we compute the KS statistic by taking the maximum difference of F_D and \hat{F}_D at 100 equally spaced points in the interval $[0, t_{max}]$ where $F_D(t_{max}) \approx 1$.

The second kind of tests is on the integral examples in Section 4.5, let X be a random variable satisfying Gibbs distribution on the manifold $\mathcal{S}_+^{n,p}$ under either metric. Define

$$\mu := \mathbb{E} \left(f(X) e^{\beta \mathcal{E}(X)} \right) = \frac{1}{Z_\beta} \int_{\mathcal{S}_+^{n,p}} f(X) dV.$$

Given m samples $X_i \in \mathcal{S}_+^{n,p}$, we define

$$(5.2) \quad \hat{\mu}_m := \frac{1}{m} \sum_{i=1}^m f(X_i) e^{\beta \mathcal{E}(X_i)}.$$

Notice that samples generated by MCMC are not independent. If we assume

$$\sigma^2 := \text{var} \left(f(X_1) e^{\beta \mathcal{E}(X_1)} \right) + 2 \sum_{k=1}^{\infty} \text{cov} \left(f(X_1) e^{\beta \mathcal{E}(X_1)}, f(X_{1+k}) e^{\beta \mathcal{E}(X_{1+k})} \right) < \infty,$$

then by the Markov Chain Central Limit Theorem[21, 14], as $m \rightarrow \infty$, we have

$$(5.3) \quad \sqrt{m}(\hat{\mu}_m - \mu) \rightarrow \mathcal{N}(0, \sigma^2)$$

where the convergence is in the sense of distribution. Thus if $m \gg 1$, $\frac{\hat{\mu}_m - \mu}{\mu}$ roughly follows the distribution $\mathcal{N}(0, \mathcal{O}(\frac{1}{m}))$ and the relative error term $\left| \frac{\hat{\mu}_m - \mu}{\mu} \right|$ roughly follows the folded normal distribution with mean $\mathcal{O}(\frac{1}{\sqrt{m}})$ and variance $\mathcal{O}(\frac{1}{m})$. Hence we can use $\hat{\mu}_m$ defined in (5.2) to estimate $\mu = \frac{1}{Z_\beta} \int_{\mathcal{S}_+^{n,p}} f(X) dV$, and the relative error is $\mathcal{O}(\frac{1}{\sqrt{m}})$.

5.1. Numerical validation of the scalar variable $D(X)$. The manifold $\mathcal{S}_+^{n,p}$ has dimension $N = np - p(p-1)/2$. For both metrics, we consider three examples in Section 4 with special energy functions \mathcal{E} in the Gibbs distribution $e^{-\beta \mathcal{E}}$ and the CDF for the scalar variable $D(X)$:

1. Example I: $\mathcal{E}(X) = \frac{1}{2} \|X\|_F^2$ with the CDF for $D(X) = \|X\|_F$:

$$\text{For } g_E : \quad F_D(t) = \Pr(\|X\|_F \leq t) \propto \int_0^t e^{-\frac{\beta}{2} \rho^2} \rho^{N-1} d\rho,$$

$$\text{For } g_{BW} : \quad F_D(t) = \Pr(\|X\|_F \leq t) \propto \int_0^t e^{-\frac{\beta}{2} \rho^2} \rho^{N/2-1} d\rho.$$

2. Example II: $\mathcal{E}(X) = \text{Tr}(X \log X)$ with the CDF $F_D(t) = \Pr(\|X\|_F \leq t)$ for $D(X) = \|X\|_F$:

$$\text{For } g_E : \quad F_D(t) \propto \int_{\substack{\sum_{i=1}^p \lambda_i^2 < t^2 \\ \lambda_i > 0, i=1, \dots, p}} \prod_{1 \leq i < j \leq p} |\lambda_i - \lambda_j| \prod_{i=1}^p |\lambda_i|^{n-p-\beta \lambda_i} \prod_{i=1}^p d\lambda_i,$$

$$\text{For } g_{BW} : \quad F_D(t) \propto \int_{\substack{\sum_{i=1}^p \lambda_i^2 < t^2 \\ \lambda_i > 0, i=1, \dots, p}} \prod_{1 \leq i < j \leq p} \frac{|\lambda_i - \lambda_j|}{\sqrt{\lambda_i + \lambda_j}} \prod_{i=1}^p |\lambda_i|^{\frac{n-p-1}{2} - \beta \lambda_i} \prod_{i=1}^p d\lambda_i.$$

which is a p -fold integral and can be approximated accurately by quadrature such as Simpson's rule for relatively small values of p , e.g., $p = 2, 3$.

3. Example III: $\mathcal{E}(X) = \frac{1}{2} \|X - A\|_F^2$ where $A \in \mathcal{S}_+^{n,p}$ has eigenvalues $\lambda_1 \geq \dots \geq \lambda_p \gg 1$, with the CDF for $D(X) = \|X - A\|_F$:

$$\text{For both } g_E \text{ and } g_{BW} : F_D(t) = \Pr(\|X - A\|_F \leq t) \propto \int_0^t e^{-\frac{\beta}{2}\rho^2} \rho^{N-1} d\rho.$$

In implementation of the scheme, the step size Δt and β in schemes (3.1) and (3.5) are two parameters that need to be tuned to reach equilibrium with reasonable computing time. We first use a numerically stable Δt then adjust β so that the noise term has reasonable variance. And of course one needs a sufficient large number of iterations for schemes (3.1) and (3.5) to reach their equilibrium state, and a sufficient large number m of samples to observe numerical convergence toward the Gibbs distribution through the scalar random variable D , e.g., the KS statistic (5.1) should be small. See Figure 1, Figure 2, Figure 3, and Figure 4 for the numerical results.

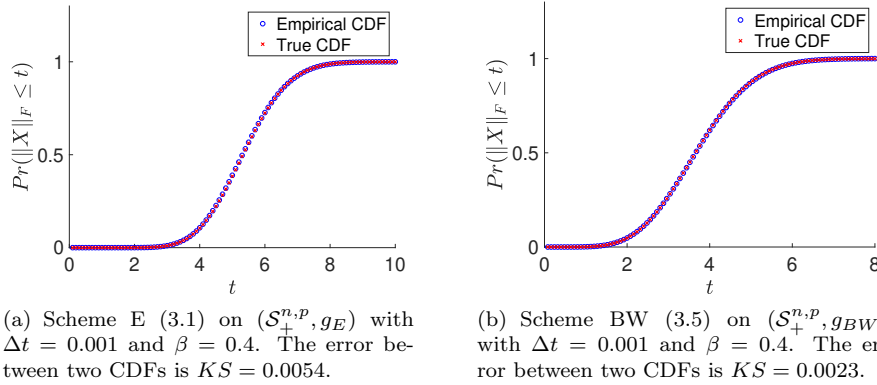


FIG. 1. Example I: $\mathcal{E}(X) = \frac{1}{2} \|X\|_F^2$, $n = 5, p = 3$ and manifold dimension is $N = 12$. The empirical CDF is computed by 5E6 MCMC samples generated after 6E6 iterations of the Riemannian Langevin Monte Carlo schemes. Both CDFs of scheme E and scheme BW are evaluated at 100 equally spaced points on $[0, 10]$ and $[0, 8]$, respectively, and the difference can be measured by the KS statistic (5.1).

5.2. MCMC numerical integration. We consider special cases $k = 0, m = 2$ in the examples (4.18) and (4.19), then (4.18) reduces to $(\frac{\beta}{\alpha})^{N/2}$ and (4.19) reduces to $(\frac{\beta}{\alpha})^{N/4}$. In other words, we may verify the numerical convergence of samples X_i to Gibbs distribution by verifying

$$(5.4) \quad \text{For } g_E : \frac{1}{m} \sum_{i=1}^m e^{-\frac{\alpha-\beta}{2} \|X_i\|_F^2} \rightarrow \left(\frac{\beta}{\alpha}\right)^{N/2},$$

$$(5.5) \quad \text{For } g_{BW} : \frac{1}{m} \sum_{i=1}^m e^{-\frac{\alpha-\beta}{2} \|X_i\|_F^2} \rightarrow \left(\frac{\beta}{\alpha}\right)^{N/4}.$$

In Figure 5 we indeed observe the $\mathcal{O}(1/\sqrt{m})$ for the relative error of numerical integration.

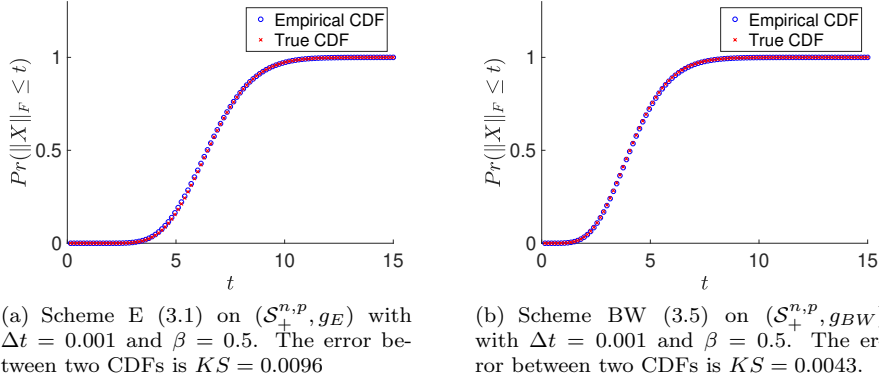


FIG. 2. *Example II: $\mathcal{E}(X) = \text{Tr}(X \log X)$, $n = 5, p = 3$ and manifold dimension is $N = 12$. The empirical CDF is computed by 5E6 MCMC samples generated after 6E6 iterations of the Riemannian Langevin Monte Carlo schemes. Both CDFs are evaluated at 100 equally spaced points on $[0, 15]$, and the difference can be measured by the KS statistic (5.1).*

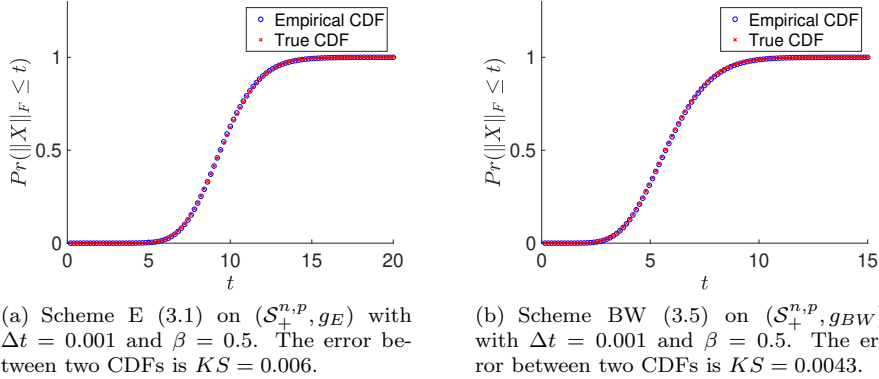
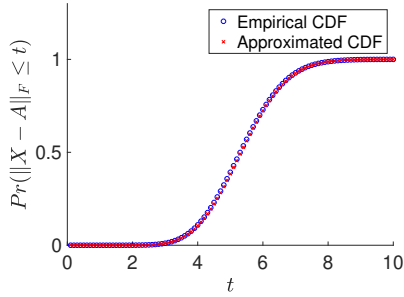


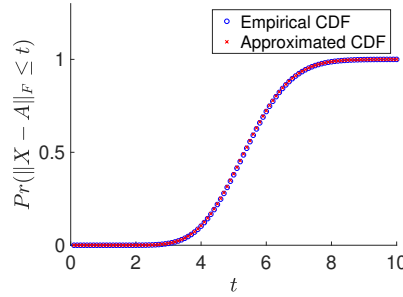
FIG. 3. *Example II: $\mathcal{E}(X) = \text{Tr}(X \log X)$, $n = 10, p = 2$ and manifold dimension is $N = 19$. The empirical CDF is computed by 5E6 MCMC samples generated after 6E6 iterations of the Riemannian Langevin Monte Carlo schemes. Both CDFs of scheme E and scheme BW are evaluated at 100 equally spaced points on $[0, 20]$ and $[0, 15]$, respectively, and the difference can be measured by the KS statistic (5.1).*

5.3. A numerical study of the convergence to equilibrium. The general mathematical theory of convergence of a Langevin equation to its equilibrium measure has been well studied; we consider the specific case of the RLE studied here in the companion paper [36]. One particular application of the two Riemannian Langevin Monte Carlo schemes is to use them to numerically study the SDE solutions, e.g., by taking very small time steps, a Riemannian Langevin Monte Carlo scheme approximates the Riemannian Langevin equation on the manifold. We have shown comparison of the Langevin equation on $(\mathcal{S}_+^{n,p}, g_E)$, $(\mathcal{S}_+^{n,p}, g_{BW})$, $\mathbb{R}^{n \times n}$ in Figure 6, in which we can see interesting differences between two metrics. With all three figures in Figure 6, we can see that the SDE on $(\mathcal{S}_+^{n,p}, g_{BW})$ has a much faster convergence to its Gibbs measure than the SDE on $(\mathcal{S}_+^{n,p}, g_E)$.

6. Conclusion. We have constructed two efficient Riemannian Langevin Monte Carlo schemes for sampling PSD matrices of fixed rank from the Gibbs distribution

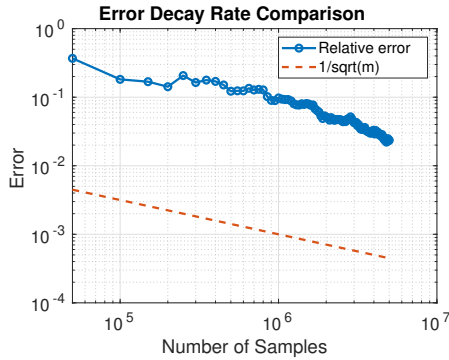


(a) Scheme E (3.1) on $(\mathcal{S}_+^{n,p}, g_E)$ with $\Delta t = 0.001$ and $\beta = 0.4$. The error between two CDFs is $KS = 0.0084$.

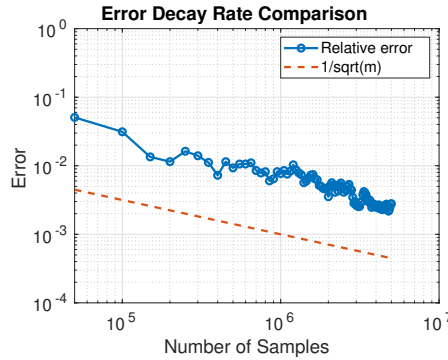


(b) Scheme BW (3.5) on $(\mathcal{S}_+^{n,p}, g_{BW})$ with $\Delta t = 2E-7$ and $\beta = 0.4$. The error between two CDFs is $KS = 0.0052$.

FIG. 4. *Example III: $\mathcal{E}(X) = \frac{1}{2} \|X - A\|_F^2$, $n = 5, p = 3$ and manifold dimension is $N = 12$. The nonzero eigenvalues of A are equally spaced between 10000 and 20000. The empirical CDF is computed by 5E6 MCMC samples generated after 6E6 iterations of the Riemannian Langevin Monte Carlo schemes. Both CDFs are evaluated at 100 equally spaced points on $[0, 10]$, and the difference can be measured by the KS statistic (5.1).*



(a) Integration on $(\mathcal{S}_+^{n,p}, g_E)$ via samples generated by Scheme E (3.1) with $\Delta t = 0.001$ and $\beta = 0.4$.



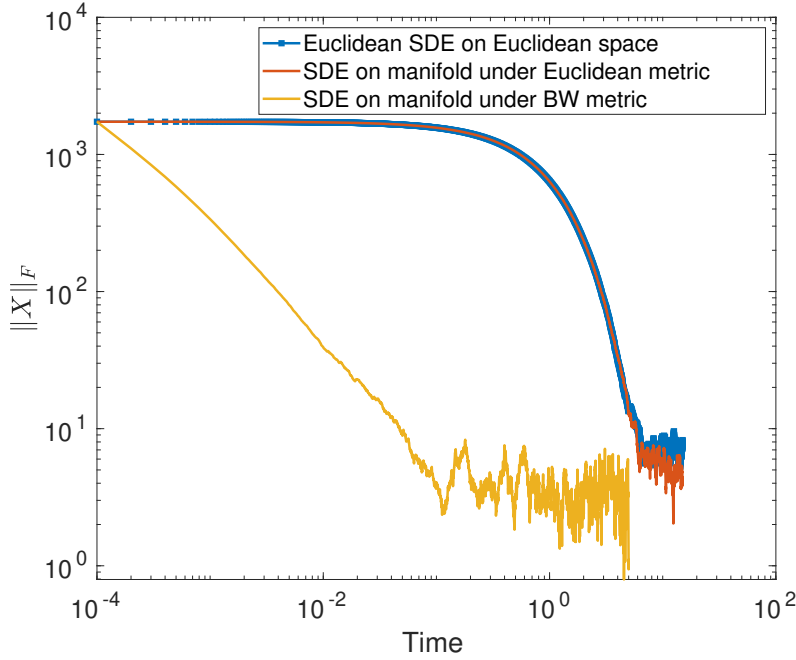
(b) Integration on $(\mathcal{S}_+^{n,p}, g_{BW})$ via samples generated by Scheme BW (3.5) with $\Delta t = 0.001$ and $\beta = 0.4$.

FIG. 5. *Convergence rate of the relative error of $\left| \frac{\hat{\mu}_m - \mu}{\mu} \right|$ MCMC integration on the manifold with $n = 10, p = 2$ and dimension $N = 19$. Parameters are $\alpha = 0.75, \beta = 0.4$, for which it is a numerical integration of the function $f(X) = \frac{1}{2} \|X\|_F^2$ on the manifold $\mathcal{S}_+^{n,p}$. The error shown is the averaged one of 12 independent runs.*

on the manifold $\mathcal{S}_+^{n,p}$ equipped with two fundamental metrics. We have also provided several examples for which these sampling schemes can be numerically validated.

REFERENCES

- [1] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization algorithms on matrix manifolds*, Princeton University Press, 2008.
- [2] R. BHATIA, T. JAIN, AND Y. LIM, *On the Bures-Wasserstein distance between positive definite matrices*, 2017, <https://arxiv.org/abs/1712.01504>.
- [3] M. BRUBAKER, M. SALZMANN, AND R. URTASUN, *A Family of MCMC Methods on Implicitly Defined Manifolds*, in Proceedings of the Fifteenth International Conference on Artificial



(a) Convergence of SDEs to the equilibrium.

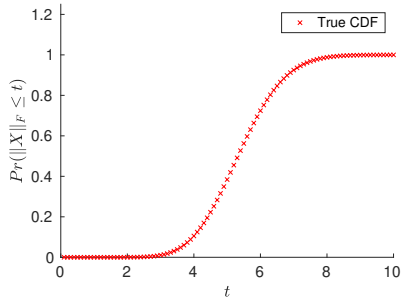
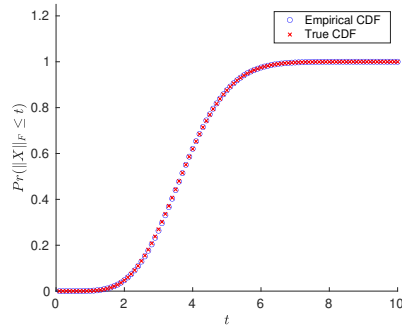

 (b) The true CDF of $D = \|X\|_F$ on $(\mathcal{S}_+^{n,p}, g_E)$. This implies the blue curve in Figure (a) has not reached equilibrium at Time=1.

 (c) The true CDF of $D = \|X\|_F$ on $(\mathcal{S}_+^{n,p}, g_{BW})$, and empirical CDF of 5E6 samples from the Scheme BW after Time=5.

FIG. 6. For $\mathcal{S}_+^{n,p}$ with $n = 5, p = 3$, the dimension is $N = 12$, and $\mathbb{R}^{n \times n}$ has dimension 25. The Gibbs measure is $e^{-\beta \mathcal{E}(x)}$ with $\beta = 0.4$ and $\mathcal{E} = \|X\|_F^2$. The time step size is $1e-4$. The initial guess is a random PSD matrix of rank-3 with eigenvalues = $[1000, 1000, 1000]$.

- Intelligence and Statistics, N. D. Lawrence and M. Girolami, eds., vol. 22 of Proceedings of Machine Learning Research, La Palma, Canary Islands, 21–23 Apr 2012, PMLR, pp. 161–172.
- [4] S. BURER AND R. D. MONTEIRO, *Local minima and convergence in low-rank semidefinite programming*, *Mathematical programming*, 103 (2005), pp. 427–444.
 - [5] S. BYRNE AND M. GIROLAMI, *Geodesic Monte Carlo on embedded manifolds*, *Scandinavian Journal of Statistics*, 40 (2013), pp. 825–845.
 - [6] X. CHENG, N. S. CHATTERJI, P. L. BARTLETT, AND M. I. JORDAN, *Underdamped Langevin MCMC: A non-asymptotic analysis*, in *Conference on learning theory*, PMLR, 2018, pp. 300–323.

- [7] X. CHENG, D. YIN, P. BARTLETT, AND M. JORDAN, *Stochastic gradient and Langevin processes*, in International Conference on Machine Learning, PMLR, 2020, pp. 1810–1819.
- [8] X. CHENG, J. ZHANG, AND S. SRA, *Efficient Sampling on Riemannian Manifolds via Langevin MCMC*, in Advances in Neural Information Processing Systems, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, eds., vol. 35, Curran Associates, Inc., 2022, pp. 5995–6006.
- [9] G. CICCOTTI, R. KAPRAL, AND E. VANDEN-EIJNDEN, *Blue moon sampling, vectorial reaction coordinates, and unbiased constrained dynamics*, ChemPhysChem, 6 (2005), pp. 1809–1814.
- [10] G. CICCOTTI, T. LELIEVRE, AND E. VANDEN-EIJNDEN, *Projection of diffusions on submanifolds: Application to mean force computation*, Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences, 61 (2008), pp. 371–408.
- [11] M. P. DO CARMO AND J. FLAHERTY FRANCIS, *Riemannian geometry*, vol. 6, Springer, 1992.
- [12] Y. DU AND I. MORDATCH, *Implicit generation and modeling with energy based models*, Advances in Neural Information Processing Systems, 32 (2019).
- [13] R. GE, H. LEE, J. LU, AND A. RISTESKI, *Efficient sampling from the bingham distribution*, in Algorithmic Learning Theory, PMLR, 2021, pp. 673–685.
- [14] C. J. GEYER, *Markov chain monte carlo lecture notes*, Course notes, Spring Quarter, 80 (1998).
- [15] M. GIROLAMI AND B. CALDERHEAD, *Riemann manifold langevin and hamiltonian monte carlo methods*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 73 (2011), pp. 123–214.
- [16] E. P. HSU, *Stochastic analysis on manifolds*, vol. 38 of Graduate Studies in Mathematics, American Mathematical Society, Providence, RI, 2002.
- [17] C.-P. HUANG, D. INAUEN, AND G. MENON, *Motion by mean curvature and Dyson Brownian motion*, Electron. Commun. Probab., 28 (2023), pp. 1–10, <https://doi.org/10.1214/23-ECP540>.
- [18] W. HUANG AND X. ZHANG, *Solving PhaseLift by low-rank Riemannian optimization methods for complex semidefinite constraints*, SIAM Journal on Scientific Computing, 39 (2017), pp. B840–B859.
- [19] N. IKEDA AND S. WATANABE, *Stochastic differential equations and diffusion processes*, vol. 24 of North-Holland Mathematical Library, North-Holland Publishing Co., Amsterdam; Kodansha, Ltd., Tokyo, second ed., 1989.
- [20] D. INAUEN AND G. MENON, *Stochastic Nash evolution*, 2023, <https://arxiv.org/abs/TBD>.
- [21] G. L. JONES, *On the Markov chain central limit theorem*, Probability Surveys, 1 (2004), pp. 299 – 320, <https://doi.org/10.1214/154957804100000051>, <https://doi.org/10.1214/154957804100000051>.
- [22] P. E. KLOEDEN AND E. PLATEN, *Numerical solution of stochastic differential equations*, vol. 23 of Applications of Mathematics (New York), Springer-Verlag, Berlin, 1992.
- [23] J. LEAKE, C. MCSWIGGEN, AND N. K. VISHNOI, *Sampling matrices from harish-chandra-itzkson-zuber densities with applications to quantum inference and differential privacy*, in Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing, 2021, pp. 1384–1397.
- [24] M. B. LI AND M. A. ERDOGDU, *Riemannian langevin algorithm for solving semidefinite programs*, 2023, <https://arxiv.org/abs/2010.11176>.
- [25] J. S. LIU AND J. S. LIU, *Monte Carlo strategies in scientific computing*, vol. 75, Springer, 2001.
- [26] E. MASSART AND P.-A. ABSIL, *Quotient geometry with simple geodesics for the manifold of fixed-rank positive-semidefinite matrices*, SIAM Journal on Matrix Analysis and Applications, 41 (2020), pp. 171–198.
- [27] E. MASSART, J. M. HENDRICKX, AND P.-A. ABSIL, *Curvature of the manifold of fixed-rank positive-semidefinite matrices endowed with the büres-wasserstein metric*, in Geometric Science of Information: 4th International Conference, GSI 2019, Toulouse, France, August 27–29, 2019, Proceedings, Springer, 2019, pp. 739–748.
- [28] G. MENON AND T. YU, *The Riemannian Langevin equation and conic programs*, 2023, <https://arxiv.org/abs/2302.11653>.
- [29] G. MENON AND T. YU, *Siegel Brownian motion*, 2023, <https://arxiv.org/abs/TBD>.
- [30] G. MEYER, S. BONNABEL, AND R. SEPULCHRE, *Regression on fixed-rank positive semidefinite matrices: a Riemannian approach*, The Journal of Machine Learning Research, 12 (2011), pp. 593–625.
- [31] A. MOITRA AND A. RISTESKI, *Fast Convergence for Langevin with Matrix Manifold Structure*, in ICLR 2020 Workshop on Integration of Deep Neural Models and Differential Equations, 2020.

- [32] T. PACINI, *Mean curvature flow, orbits, moment maps*, Trans. Amer. Math. Soc., 355 (2003), pp. 3343–3357, <https://doi.org/10.1090/S0002-9947-03-03307-5>, <https://doi.org/10.1090/S0002-9947-03-03307-5>.
- [33] A. TASISSA AND R. LAI, *Exact reconstruction of euclidean distance geometry problem using low-rank matrix completion*, IEEE Transactions on Information Theory, 65 (2018), pp. 3124–3144.
- [34] B. VANDEREYCKEN, P.-A. ABSIL, AND S. VANDEWALLE, *Embedded geometry of the set of symmetric positive semidefinite matrices of fixed rank*, in 2009 IEEE/SP 15th Workshop on Statistical Signal Processing, IEEE, 2009, pp. 389–392.
- [35] P. XU, J. CHEN, D. ZOU, AND Q. GU, *Global convergence of langevin dynamics based algorithms for nonconvex optimization*, 2020, <https://arxiv.org/abs/1707.06618>.
- [36] T. YU, S. ZHENG, J. LU, G. MENON, AND X. ZHANG, *Riemannian Langevin equations for PSD matrices of fixed rank*, 2023, <https://arxiv.org/abs/TBD>.
- [37] E. ZAPPA, M. HOLMES-CERFON, AND J. GOODMAN, *Monte Carlo on manifolds: sampling densities and integrating functions*, Communications on Pure and Applied Mathematics, 71 (2018), pp. 2609–2647.
- [38] S. ZHENG, W. HUANG, B. VANDEREYCKEN, AND X. ZHANG, *Riemannian optimization using three different metrics for Hermitian PSD fixed-rank constraints: an extended version*, 2022, <https://arxiv.org/abs/2204.07830>.