

# STEREOGRAPHIC MARKOV CHAIN MONTE CARLO

BY JUN YANG<sup>1</sup>, KRZYSZTOF ŁATUSZYŃSKI<sup>2</sup>, AND GARETH O. ROBERTS<sup>2</sup>

<sup>1</sup>*Department of Mathematical Sciences, University of Copenhagen, Denmark.*

<sup>2</sup>*Department of Statistics, University of Warwick, United Kingdom.*

High-dimensional distributions, especially those with heavy tails, are notoriously difficult for off-the-shelf MCMC samplers: the combination of unbounded state spaces, diminishing gradient information, and local moves results in empirically observed “stickiness” and poor theoretical mixing properties – lack of geometric ergodicity. In this paper, we introduce a new class of MCMC samplers that map the original high-dimensional problem in Euclidean space onto a sphere and remedy these notorious mixing problems. In particular, we develop random-walk Metropolis type algorithms as well as versions of the Bouncy Particle Sampler that are uniformly ergodic for a large class of light and heavy-tailed distributions and also empirically exhibit rapid convergence in high dimensions. In the best scenario, the proposed samplers can enjoy the “blessings of dimensionality” that the convergence is faster in higher dimensions.

**1. Introduction.** Bayesian analysis relies heavily on Markov chain Monte Carlo (MCMC) methods to explore complex posterior distributions. In most typical settings, such distributions have support contained within a subset  $S$  of  $\mathbb{R}^d$  for some  $d > 0$ , and then it is natural to construct appropriate MCMC algorithms directly on  $S$ . In practice, this is how the vast majority of algorithms are constructed, although there are intrinsic problems with this approach. For instance, it is now well-established [Mengersen and Tweedie, 1996] that the popular vanilla MCMC workhorse, the random-walk Metropolis (RWM) algorithm fails to be uniformly ergodic for any target density  $\pi$  when  $S$  is unbounded. All existing generic MCMC methods are built upon local proposal mechanisms and are similarly afflicted. Lack of uniform ergodicity results in sensitivity of the algorithm’s convergence to its starting value, and potentially long burn-in periods.

Moreover, these convergence problems are exacerbated for target distributions with heavy (heavier than exponential) tails. For instance, in such cases RWM and the Metropolis-adjusted Langevin Algorithm (MALA) fail to be even geometrically ergodic (i.e., they converge at a rate slower than any geometric rate) [Roberts and Tweedie, 1996a, Jarner and Hansen, 2000, Roberts and Tweedie, 1996b]. In practice, this manifests itself on the algo-

---

*MSC2020 subject classifications:* Primary 60J20, 60J25, 65C05.

*Keywords and phrases:* random walk Metropolis, piecewise deterministic Markov processes, stereographic projection, uniform ergodicity, heavy tailed distributions, blessings of dimensionality.

rithm trajectories by the presence of infrequent excursions of heavy-tailed duration into the target distribution tails. This can lead to further theoretical and practical problems, e.g., the absence of a Central Limit Theorem (CLT) for all  $L^2$  functions, for example [Järner and Roberts, 2007], and instability of Monte Carlo estimators with large and difficult-to-quantify mean square errors which are highly sensitive to initial values.

Important modern innovations in MCMC algorithms have come from Piecewise Deterministic Markov Processes (PDMPs) [Bernard et al., 2009, Bouchard-Côté et al., 2018, Bierkens et al., 2019] which offer for the first time generic recipes for the construction of non-reversible MCMC and often yield substantial gains in algorithm efficiency as a result. Although not completely necessary, it turns out to be convenient and natural to construct PDMPs as continuous time algorithms. While the practical use of PDMPs for posterior exploration is still in its infancy, these methods offer substantial promise. However PDMPs are still localised algorithms and as such also suffer from the lack of uniform and/or geometric ergodicity on unfounded state spaces and/or heavy-tailed targets [Vasdekis and Roberts, 2022, 2023, Andrieu et al., 2021a,b].

In contrast, when the target density tails are exponential or lighter, RWM, MALA, and related algorithms are generally geometrically ergodic under weak regularity conditions [Roberts and Tweedie, 1996a, Järner and Hansen, 2000]. Some transformation strategies for achieving this are described in [Sherlock et al., 2010, Johnson and Geyer, 2012] with closely related strategies being proposed in [Kamatani, 2018], although these approaches are unlikely to regain uniform ergodicity. When  $S$  is bounded, then MCMC algorithms are generally easily shown to be uniformly ergodic (see for example [Mengersen and Tweedie, 1996] for RWM). This naturally suggests that transformations designed to compactify  $S$  might facilitate the construction of more robust families of MCMC algorithms. However, finding generic solutions to the construction of such transformations which lead to well-behaved densities on the transformed space is challenging.

It is, however, far easier to construct general transformations to transform  $\mathbb{R}^d$  to a *pre*-compact space. The most celebrated of such transformations is *stereographic projection* which was known to the ancient Egyptians, see [Coxeter, 1961] for a modern description. It maps  $\mathbb{R}^d \rightarrow \mathbb{S}^d \setminus N$  (where  $N$  denotes the *North Pole*), i.e., the  $d$ -sphere excluding its North Pole. This paper will explore the use of stereographically projected algorithms. One iteration of such an algorithm takes the current state  $\mathbf{x} \in \mathbb{R}^d$ , transforms to  $\mathbb{S}^d \setminus N$  through inverse stereographic projection, carrying out an appropriately constructed MCMC step on  $\mathbb{S}^d \setminus N$  before returning to  $\mathbb{R}^d$  by stereographic projection. The contributions of our work are as follows.

1. We present the Stereographic Projection Sampler (SPS), an efficient and practical implementation of the above programme for RWM. It employs a simple reprojection step to

ensure the Markov chain remains on  $\mathbb{S}^d \setminus N$ . We provide a dimension and scale dependent recipe for choosing the radius of  $S$ . See Section 2.1.

2. We prove that for continuous positive densities,  $\pi$  on  $\mathbb{R}^d$  with tails no heavier than those of a  $d$ -dimensional multivariate student's  $t$  distribution with  $d$  degrees of freedom, SPS is uniformly ergodic. The tail conditions on  $\pi$  are in fact necessary for geometric ergodicity. See Theorem 2.1 and Remark 2.1.
3. We give a high-dimensional analysis of SPS for the stylised family of product i.i.d. form targets, by maximizing the expected squared jumping distance (ESJD) as well as by showing for a variant of SPS that each component converges (as  $d \rightarrow \infty$ ) to a Langevin diffusion. This affords a direct and uniformly favorable comparison with the Euclidean RWM algorithm. See Section 5.
4. We also introduce the Stereographic Bouncy Particle Sampler (SBPS) which replaces the random walk move in SPS with a PDMP algorithm that follows great circle trajectories interspersed with abrupt direction changes. See Section 2.2.
5. We prove the uniform ergodicity of SBPS under weaker conditions than those for SPS, specifically requiring tails to be lighter than those of a  $d$ -dimensional multivariate student's  $t$  distribution with  $d - 1/2$  degrees of freedom. See Theorem 2.2 and Corollary 2.1.
6. For isotropic targets (i.e., spherically symmetric targets), both light- and heavy-tailed, we (informally) demonstrate that the proposed SPS and SBPS can enjoy a *blessings of dimensionality* effect which implies that SPS and SBPS converge arbitrarily faster in higher dimensions than their Euclidean analogues. See Section 3 and Section 6.
7. We introduce generalisations of stereographic projection which are suitable for elliptical targets. See Section 4. The framework we established in this paper opens opportunities for developing other mappings from  $\mathbb{R}^d$  to  $\mathbb{S}^d$  for other classes of targets. See Section 7.

In very recent work, [Lie et al., 2021] presents a related methodology for exploring distributions defined directly on a manifold, providing supporting theory showing that under suitable conditions their method's spectral gap is dimension-independent. However, the focus of their work is very different. Their work uses a different reprojection scheme based on [Cotter et al., 2013] and does not consider distributions defined directly on  $\mathbb{R}^d$ . Moreover, they consider densities that are absolutely continuous with respect to a Gaussian measure in the infinite-dimensional limit. This is arguably a very restrictive class. So their dimension-independent results are not really comparable with our results. An alternative reprojection scheme is introduced in [Zappa et al., 2018]. Compared to that method, the reprojection scheme used in this paper has the advantage that reprojection from a fixed point on the tangent space is always possible and is a much more natural approach for the hypersphere. The [Zappa et al., 2018] method does have advantages for use in more general manifolds, but this is not of any use to us here.

There are strong theoretical reasons for wanting to construct the algorithm dynamics directly on a manifold of positive curvature such as a hypersphere [Mangoubi and Smith, 2018, Ollivier, 2009, Mijatović et al., 2018]. However, quite naturally, the existing literature has concentrated primarily on the case of Brownian motion which clearly has the uniform invariant distribution on  $S$ , or on the geodesic walk designed specifically to target uniform distributions on manifolds. The uniform distribution on  $S$  maps via stereographic projection to the student's  $t$  distribution with  $d$  degrees of freedom on  $\mathbb{R}^d$ . Therefore, we can immediately lift theory from existing theory to (informally) show that SPS on such target densities has a dimension-free convergence time. In our paper, we go further. Rapid convergence extends to a large family of spherically symmetric target densities (essentially excluding only very heavy-tailed distributions).

However the proposed algorithms are not only applicable for stylized classes of target distributions. To showcase the practical value of stereographic samplers in more realistic statistical contexts, we shall demonstrate the utility of our methodology on a Bayesian analysis of a Cauchy regression model. Full specification of the model, its parameter values and other details will be given in Section 6.1.

EXAMPLE 1.1. Consider  $Y_i \sim \text{Cauchy}(\alpha + \beta^T X_i, \gamma)$ , where  $\{X_i, Y_i\}_{i=1}^n$  are the design matrices and responses, respectively,  $\alpha \in \mathbb{R}$ ,  $\beta \in \mathbb{R}^{d-2}$ ,  $\gamma \in \mathbb{R}^+$  are the parameters. Assume a flat prior for  $\alpha$  and  $\beta$ , a Gamma( $a, b$ ) prior for  $\gamma$  (i.e.,  $\pi_0(\gamma) \propto \gamma^{a-1} \exp(-b\gamma)$ ), the posterior can be written by

$$(1) \quad \pi(\alpha, \beta, \gamma) \propto \frac{\gamma^{a-1-n} \exp(-b\gamma)}{\prod_{i=1}^n \left\{ 1 + \left[ \frac{Y_i - (\alpha + \beta^T X_i)}{\gamma} \right]^2 \right\}}.$$

For sampling the posterior, we compare SPS and RWM in Fig. 1, where we plot the traceplots for  $\alpha$ ,  $\beta_1$  (the first coordinate of  $\beta$ ), and  $\log(\gamma)$ , as well as  $(\alpha, \beta_1)$  for both algorithms. From the figure, one can see clearly that: (1) RWM which is not geometric ergodic completely failed; (2) the proposed SPS which is uniformly ergodic converged extremely fast.

Our paper is structured as follows. Section 2 introduces both SPS and SBPS algorithms and presents formal ergodicity and uniform ergodicity results for both methods. Sections 3 to 5 provide additional results for SPS and its generalisations. In Section 3 a detailed analysis of SPS for isotropic densities is provided, while in Section 4, a generalised version of the algorithm is introduced. Section 4 also gives robustness results to departures from isotropy. In Section 5 we provide a high-dimensional analysis of SPS for the stylised family of product i.i.d. target densities and give positive comparisons to the Euclidean RWM algorithm. Numerical studies on both SPS and SBPS are provided in Section 6 to illustrate our theory and we conclude with the discussions on SPS and SBPS in Section 7. The supplementary materials contain all the technical proofs and some additional simulations.

Bayesian Cauchy Regression (d=11, n=15, accept rates SPS=0.57, RWM=0.69)

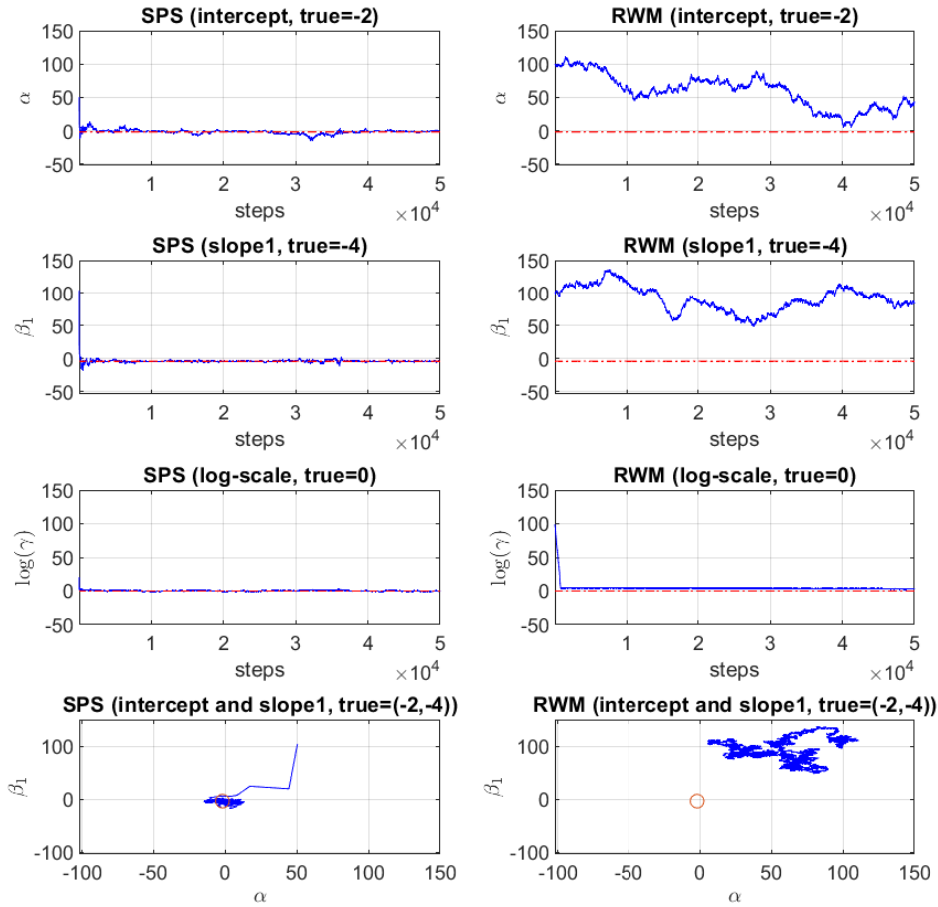


Figure 1: Bayesian Cauchy Regression. Red dotted lines and circles represent the parameters of the true sampling distribution.

## 2. Stereographic Markov Chain Monte Carlo.

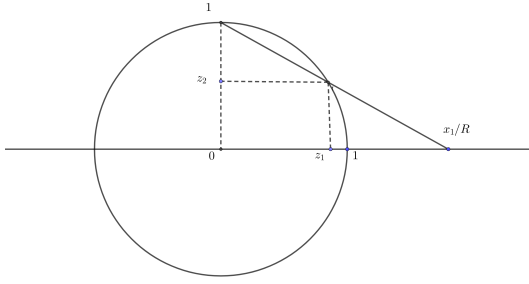
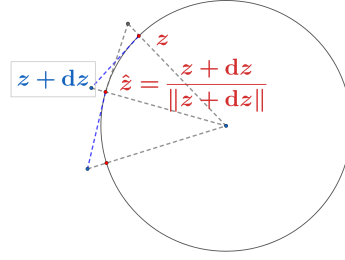
2.1. *Stereographic Projection Sampler (SPS).* We begin by giving a brief description of stereographic projection and then describe in detail two novel MCMC samplers that exploit the properties of stereographic projection.

Let  $\mathbb{S}^d$  denote the *unit sphere* in  $\mathbb{R}^{d+1}$  centered at the origin. A stereographic projection describes a bijection from  $\mathbb{S}^d \setminus \{(0, \dots, 0, 1)\}$  to  $\mathbb{R}^d$ . Within this paper, we shall restrict attention to projections indexed by a single parameter  $R \in \mathbb{R}^+$  and described by the mapping

$$x = \text{SP}(z) := \left( R \frac{z_1}{1 - z_{d+1}}, \dots, R \frac{z_d}{1 - z_{d+1}} \right)^T,$$

with Jacobian determinant at  $x \in \mathbb{R}^d$  satisfying

$$(2) \quad J_{\text{SP}}(x) \propto (R^2 + \|x\|^2)^d,$$

Figure 2: Illustration of Stereographic Projection  $SP : \mathbb{S} \rightarrow \mathbb{R}$ .Figure 3: Illustration of SPS proposing  $\hat{z} = SP^{-1}(\hat{X})$  from  $z = SP^{-1}(x)$ . By symmetry, the proposal distribution is the same as proposing  $z$  from  $\hat{z}$ .

and inverse  $SP^{-1} : \mathbb{R}^d \rightarrow \mathbb{S}^d \setminus \{(0, \dots, 0, 1)\}$  given by

$$(3) \quad z_i = \frac{2Rx_i}{\|x\|^2 + R^2}, \quad \forall 1 \leq i \leq d, \quad z_{d+1} = \frac{\|x\|^2 - R^2}{\|x\|^2 + R^2}.$$

See Fig. 2 for a geometric illustration of this stereographic projection (in the case  $d = 1$  for simplicity) and Section S8.13 for the proof of the Jacobian determinant Eq. (2).

Now suppose we wish to sample from a target density  $\pi(x)$  where  $x \in \mathbb{R}^d$ . Our aim is to take advantage of our stereographic bijection to construct an MCMC sampler directly on  $\mathbb{S}^d$ , projecting the output back onto  $\mathbb{R}^d$ . We denote the transformed target as  $\pi_S(z)$  then for  $x = SP(z)$  we have

$$(4) \quad \pi_S(z) \propto \pi(x)(R^2 + \|x\|^2)^d.$$

First, we just consider a random-walk Metropolis algorithm with a step size  $h$  on the unit sphere. Note that we describe our algorithms on the unit sphere. The radius  $R$  will only be taken into account when projecting to  $\mathbb{R}^d$ . Fig. 3 illustrates how the algorithm moves are constructed on the unit sphere. (Note that we could very easily have constructed more general Metropolis–Hastings algorithms.)

---

**Algorithm 1: Stereographic Projection Sampler (SPS)**


---

- Let the current state be  $X^d(t) = x$ ;
  - Compute the proposal  $\hat{X}$ :
    - Let  $z := SP^{-1}(x)$ ;
    - Sample independently  $d\tilde{z} \sim \mathcal{N}(0, h^2 I_{d+1})$ ;
    - Let  $dz := d\tilde{z} - \frac{(z^T \cdot d\tilde{z})z}{\|z\|^2}$  and  $\hat{z} := \frac{z+dz}{\|z+dz\|}$ ;
    - The proposal  $\hat{X} := SP(\hat{z})$ .
  - $X^d(t+1) = \hat{X}$  with probability  $1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(x)(R^2 + \|x\|^2)^d}$ ; otherwise  $X^d(t+1) = x$ .
- 

The symmetry of the re-projection mechanism in the proposal of the algorithm means that SPS is indeed a valid MCMC algorithm for  $\pi$ . More formally, we have the following result.

PROPOSITION 2.1. *If  $\pi(x)$  is positive and continuous in  $\mathbb{R}^d$ , then SPS gives rise to an ergodic Markov chain on  $\mathbb{R}^d$  with invariant distribution  $\pi$ .*

PROOF. See Section S8.1. □

We shall address the chain’s uniform ergodicity properties. Recall that a Markov chain  $X$  on state space  $E$  with transition kernel  $P$  is *uniformly ergodic* if  $\forall \epsilon > 0$ , there exists,  $N \in \mathbb{N}$  such that  $\|P^N(x, \cdot) - \pi\|_{\text{TV}} \leq \epsilon, \forall x \in E$ , where  $\pi$  denotes the chain’s (unique) invariant distribution and  $\|\cdot\|_{\text{TV}}$  represents total variation distance. Intuitively, for uniformly ergodic chains, we cannot have “arbitrarily bad” starting values.

It is long recognised that random-walk Metropolis (RWM) algorithms on bounded spaces are usually uniformly ergodic, while the same algorithms on unbounded spaces are never uniformly ergodic [Mengersen and Tweedie, 1996]. Therefore, given the compactness of  $\mathbb{S}^d$ , it is reasonable to hope that SPS might be uniformly ergodic under mild regularity conditions on  $\pi$ . Since the actual state space of SPS is in fact  $\mathbb{S}^d \setminus N$  which is not compact, this question is more complicated than in the Euclidean state space case as we need to consider the properties of the transformed density near  $N$ . However, our first main result confirms that we do get uniform ergodicity if and only if the transformed density on the sphere is bounded at  $N$ .

THEOREM 2.1. *If  $\pi(x)$  is positive and continuous in  $\mathbb{R}^d$ , then SPS is uniformly ergodic if and only if*

$$(5) \quad \sup_{x \in \mathbb{R}^d} \pi(x)(R^2 + \|x\|^2)^d < \infty.$$

PROOF. See Section S8.2. □

EXAMPLE 2.1. *If the target  $\pi(x)$  where  $x \in \mathbb{R}^d$  is multivariate student’s  $t$  distribution with degrees of freedom no smaller than  $d$ , then the SPS algorithm is uniformly ergodic.*

REMARK 2.1. *We make the following remarks:*

1. *The condition Eq. (5) is necessary: if  $\sup_{x \in \mathbb{R}^d} \pi(x)(R^2 + \|x\|^2)^d = \infty$ , then the chain is not even geometrically ergodic [Roberts and Tweedie, 1996a, Proposition 5.1];*
2. *The traditional RWM algorithm is not uniformly ergodic if the support of  $\pi$  is  $\mathbb{R}^d$  [Mengersen and Tweedie, 1996, Theorem 3.1] and not geometrically ergodic for any heavy-tailed target distribution [Järner and Hansen, 2000, Corollary 3.4];*
3. *The condition that  $\pi(x)$  is positive and continuous in  $\mathbb{R}^d$  in both Proposition 2.1 and Theorem 2.1 can be relaxed. We used it here just for the simplicity of the proof.*

2.2. *Stereographic Bouncy Particle Sampler (SBPS).* Many recent innovations in MCMC algorithm construction have focused on non-reversible methods, most particularly those de-

scribed by *piecewise deterministic Markov processes* (PDMPs) (see for example [Bouchard-Côté et al., 2018, Bierkens et al., 2019] and [Davis, 1984] for theoretical background). PDMPs are continuous-time processes that have stochastic jumps at event times of a point process, but where the state evolves deterministically between the event times. In this subsection, we shall demonstrate that we can readily incorporate these methods within our projective framework. We shall concentrate on a version of the *Bouncy Particle Sampler* (BPS) [Bouchard-Côté et al., 2018] as this adapts naturally to our context. The algorithm is described as follows.

PDMPs utilise an auxiliary random variable  $v$ , which in the case of the Stereographic Bouncy Particle Sampler (SBPS) has stationary distribution uniformly distributed on  $\mathbb{S}^d$  and independently of  $x$ . One feature of PDMP algorithms such as SBPS is the option to include *refresh* moves that contribute to the intensity of their constituting point process (according to some possibly  $x$ -dependent hazard rate) and which independently refresh  $v$  by sampling it from its invariant distribution (uniform on  $\mathbb{S}^d$ ). In the description below, we restrict ourselves to the case where this refresh rate is constant.

---

**Algorithm 2:** Stereographic Bouncy Particle Sampler (SBPS)

---

- Initialize  $z^{(0)} \in \mathbb{S}^d$  and  $v^{(0)}$  such that  $v^{(0)} \cdot z^{(0)} = 0$  and  $\|v^{(0)}\| = 1$ .
- Simulate BPS on unit sphere: for  $i = 1, 2, \dots$ 
  - Simulate bounce time  $\tau_{\text{bounce}}$  of a Poisson process of intensity

$$\chi(t) = \lambda(\sin(t)v^{(i-1)} + \cos(t)z^{(i-1)}, \cos(t)v^{(i-1)} - \sin(t)z^{(i-1)}),$$

where

$$\lambda(z, v) := \max\{0, [-v \cdot \nabla_z \log \pi_S(z)]\}.$$

- Simulate refreshment time  $\tau_{\text{refresh}} \sim \text{Exponential}(\lambda_{\text{refresh}})$ .
- Let  $\tau_i = \min\{\tau_{\text{bounce}}, \tau_{\text{refresh}}\}$  and

$$z^{(i)} = \sin(\tau_i)v^{(i-1)} + \cos(\tau_i)z^{(i-1)}.$$

- If  $\tau_i = \tau_{\text{refresh}}$ , sample new  $v^{(i)}$  independently

$$v^{(i)} \sim \text{Uniform}\{v : z^{(i)} \cdot v = 0, \|v\| = 1\}$$

- If  $\tau_i = \tau_{\text{bounce}}$ , compute

$$v^{(i)} = v_{\text{temp}} - 2 \left[ \frac{v_{\text{temp}} \cdot \tilde{\nabla}_z \log \pi_S(z^{(i)})}{\tilde{\nabla}_z \log \pi_S(z^{(i)}) \cdot \tilde{\nabla}_z \log \pi_S(z^{(i)})} \right] \tilde{\nabla}_z \log \pi_S(z^{(i)}),$$

where

$$v_{\text{temp}} = \cos(\tau_i)v^{(i-1)} - \sin(\tau_i)z^{(i-1)}$$

$$\tilde{\nabla}_z \log \pi_S(z^{(i)}) = \nabla_z \log \pi_S(z^{(i)}) - \left[ z^{(i)} \cdot \nabla_z \log \pi_S(z^{(i)}) \right] z^{(i)}.$$

- If  $\sum_{j=1}^i \tau_j \geq T$  (where  $T$  is some constant time), exit.
  - Return  $x = \text{SP}(z)$  where  $z$  denotes BPS on unit sphere.
-



Again, we demonstrate that SBPS is indeed a valid MCMC algorithm, at least under certain conditions on  $\pi$  and the refresh rate in the algorithm. Note that the conditions on  $\pi$  and  $\lambda_{\text{refresh}}$  in Proposition 2.2 are not required for invariance, only for ensuring  $\phi$ -irreducibility of the algorithm.

PROPOSITION 2.2. *Suppose that  $\lambda_{\text{refresh}} > 0$  and  $\pi > 0$  for all  $x \in \mathbb{R}^d$ . Then SBPS gives rise to an ergodic Markov chain on  $\mathbb{R}^d$  with invariant distribution for  $(x, v)$  with joint density on  $\mathbb{R}^d \times \mathbb{S}^d$ . The marginal density on  $\mathbb{R}^d$  is proportional to  $\pi(x)$ .*

PROOF. The proof of this result is routine (but tedious), following the lines of existing results for PDMPs in the literature. We give a sketch proof in Section S8.3.  $\square$

We shall prove a uniform ergodicity result analogous to that of Theorem 2.1. For a Markov process  $X$  with state space  $E$ , transition semi-group  $P$ , and invariant distribution  $\pi$ , it is uniformly ergodic if  $\forall \epsilon > 0$  there exists  $T$  such that  $\|P^T(x, \cdot) - \pi\|_{\text{TV}} \leq \epsilon, \forall x \in E$ .

THEOREM 2.2. *If  $\pi(x)$  is positive in  $\mathbb{R}^d$  with continuous first derivative in all components, then SBPS is uniformly ergodic if*

$$\limsup_{\{x: \|x\| \rightarrow \infty\}} \sum_{i=1}^d \left( \frac{\partial \log \pi(x)}{\partial x_i} x_i \right) + 2d < \frac{1}{2}.$$

PROOF. See Section S8.4.  $\square$

COROLLARY 2.1. *If the target  $\pi(x)$  where  $x \in \mathbb{R}^d$  is multivariate student's  $t$  distribution with degrees of freedom larger than  $d - \frac{1}{2}$ , then the SBPS algorithm is uniformly ergodic.*

PROOF. See Section S8.5.  $\square$

REMARK 2.2. *We make the following remarks:*

1. *This is the first known PDMP algorithm that is uniformly ergodic for a large family of target distributions including heavy-tailed targets. For comparison, the Euclidean BPS is only known to be geometrically ergodic under certain restrictive conditions on the target [Deligiannidis et al., 2019, Durmus et al., 2020] and is known not to be geometrically ergodic for any heavy-tailed target distribution [Vasdekis and Roberts, 2023].*
2. *We conjecture that the best possible condition for Theorem 2.2 is*

$$\limsup_{\{x: \|x\| \rightarrow \infty\}} \sum_{i=1}^d \left( \frac{\partial \log \pi(x)}{\partial x_i} x_i \right) + 2d < 1.$$

*We explain the reason for our conjecture in Section S8.6. Therefore, we conjecture that Corollary 2.1 is only loose by  $\frac{1}{2}$  degree. That is, if the target  $\pi(x)$  where  $x \in \mathbb{R}^d$  is multivariate student's  $t$  distribution with degree of freedom larger than  $d - 1$ , then the SBPS algorithm is conjectured to be uniformly ergodic.*

To finish this section, we present some typical sample paths obtained by implementing SBPS.

**EXAMPLE 2.2.** *In Fig. 4, we show the proposed SBPS without and with refreshment for standard Gaussian target in two dimensions. Note that the SBPS without refreshment is not irreducible in this case. The same issue exists for the Euclidean BPS for standard Gaussian targets. This issue can be fixed by adding the refreshment.*

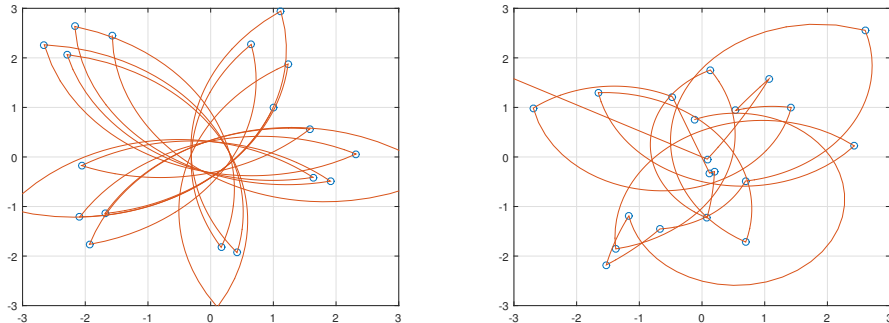


Figure 4: SBPS without (left) and with (right) refreshment for target distribution  $\mathcal{N}(0, I_2)$ . Note that the SBPS chain is not irreducible without refreshment.

**3. SPS: Isotropic Targets .** In this section, we consider *isotropic targets*, which is also called spherical symmetric targets<sup>1</sup>. That is,  $\pi(x)$  is only a function of  $\|x\|$ . Recall the mapping and Eq. (3) and the transformed target  $\pi_S(z)$  in Eq. (4). We can see that, for any isotropic target, the transformed target  $\pi_S(z)$  is only a function of the “latitude”  $z_{d+1}$ . In this sense, isotropic targets are the “best” targets for stereographically projected algorithms.

We shall assume that all second moments exist, and without loss of further generality, we suppose  $\pi(x)$  satisfies that

$$\mathbb{E}_{X \sim \pi}[\|X\|^2] = d.$$

Under this assumption, the “optimal” radius should be chosen as  $R = \sqrt{d}$ , since it maps the concentration region of  $\pi$  to the neighborhood of the “equator” of the sphere. Throughout this section, we study this best scenario. We will study the robustness to  $R$  and the optimal scaling of SPS for any given  $R$  in Section 5 for another family of targets.

Under the above assumptions, it suffices to study the path of the absolute value of the “latitude”  $z_{d+1}$  of SPS. Informally, the “stationary phase” of SPS is the period in which  $z_{d+1} = \mathcal{O}(d^{-1/2})$  and the “transient phase” is the period in which  $|z_{d+1}|$  is larger than

<sup>1</sup>Note that in some literature, isotropic distribution could be used to denote a distribution with zero mean and identity covariance matrix, which is different from our definition in this paper.

$\mathcal{O}(d^{-1/2})$ . Somehow surprisingly, by analyzing the proposed ‘‘latitude’’  $\hat{z}_{d+1}$ , we can (informally) show that the SPS enjoys the *blessings of dimensionality*: the number of iterations for the ‘‘latitude’’ of SPS to decrease to  $\mathcal{O}(d^{-1/2})$  decreases with the dimension  $d$ . This implies (informally) that the ‘‘transient phase’’ of SPS is  $\mathcal{O}(1)$  and it decreases with dimension.

3.1. *Analysis of the proposal distribution.* We assume the chain starts from either the North Pole or the South Pole. By the assumptions, the chain is in the ‘‘stationary phase’’ once it is around the ‘‘equator’’.

In the following, we give useful approximations for the proposed ‘‘latitude’’  $\hat{z}_{d+1}$  in both the transient phase and stationary phase. This can be used to analyze the behavior of isotropic targets. See Section S9.1 for some simulations using the result from Lemma 3.1.

LEMMA 3.1. *Let  $z_i$  be the current  $i$ -th coordinate and  $\hat{z}_i$  be the  $i$ -th coordinate of the proposal, where  $i = 1, \dots, d + 1$ . Then we have the following expression*

$$(6) \quad \hat{z}_i = \frac{1}{\sqrt{1 + h^2(U^2 + U_\perp^2)}} \left( z_i + \sqrt{1 - z_i^2} hU \right),$$

where  $U \sim \mathcal{N}(0, 1)$  and  $U_\perp^2 \sim \chi_{d-1}^2$  which is independent with  $U$ . Furthermore, if  $h = \mathcal{O}(d^{-1/2})$ , then we have the following coordinate-wise approximation

$$(7) \quad \hat{z}_i = \frac{1}{\sqrt{1 + h^2 U_\perp^2}} \left[ \left( 1 - \frac{1}{2} h^2 U^2 \right) z_i - \sqrt{1 - z_i^2} hU \right] + \mathcal{O}_{\mathbb{P}}(h^3 + dh^4 z_i).$$

As special cases of Eq. (7),  $z_{d+1}$  is the current ‘‘latitude’’ and  $\hat{z}_{d+1}$  be the proposed ‘‘latitude’’. Then, in the transient phase, if  $z_{d+1}^2 = 1 - o(h^2)$ , we have

$$(8) \quad \hat{z}_{d+1} = \frac{1}{\sqrt{1 + h^2(d-1)}} \left( 1 - \frac{1}{2} h^2 U^2 \right) z_{d+1} + \mathcal{O}_{\mathbb{P}}(d^{-1/2}),$$

and in the stationary phase, if  $z_{d+1} = \mathcal{O}(d^{-1/2})$ , we have

$$(9) \quad \hat{z}_{d+1} = \frac{1}{\sqrt{1 + h^2(d-1)}} (z_{d+1} - hU) + \mathcal{O}_{\mathbb{P}}(d^{-1}).$$

PROOF. See Section S8.7. □

The above lemma informally suggests that in the ‘‘transient phase’’, the proposed ‘‘latitude’’ is almost deterministic, whereas in the ‘‘stationary phase’’, the proposed ‘‘latitude’’ approximately follows an autoregressive process. For example, if we take  $h$  to be constant and  $i = d + 1$ , then Eq. (6) shows that the proposal ‘‘latitude’’  $\hat{z}_{d+1}$  decreases to  $\mathcal{O}(d^{-1/2})$  faster when dimension is larger since the concentration of  $U_\perp^2/d$ . If the proposals will be accepted by a probability bounded away from 0, the ‘‘transient phase’’ of SPS only involve  $\mathcal{O}(1)$  iterations and the number even decreases with the dimension. As an example, in the next subsection, we show a class of isotropic targets such that the SPS enjoys the *blessings of dimensionality*.

3.2. *Examples of Isotropic Targets.* We denote the multivariate student's  $t$  distributions by  $\pi_\nu(x)$  where  $\nu$  is the degree of freedom. We denote the standard multivariate Gaussian as the limit  $\pi_\infty(x)$ . That is,

$$\pi_\nu(x) \propto \left(1 + \frac{1}{\nu}\|x\|^2\right)^{-(\nu+d)/2}, \quad \pi_\infty(x) \propto \exp\left(-\frac{1}{2}\|x\|^2\right)$$

Then the logarithm of the likelihood ratio can be written as a function of  $z_{d+1}$  and  $\hat{z}_{d+1}$ :

$$\log \frac{\pi_\nu(\hat{X})}{\pi_\nu(x)} + d \log \frac{R^2 + \|\hat{X}\|^2}{R^2 + \|x\|^2} = d(g_{\nu/d}(z_{d+1}) - g_{\nu/d}(\hat{z}_{d+1})),$$

where  $g_k(z) := \frac{k+1}{2} \log(k + \frac{1+z}{1-z}) + \log(1-z)$ . If  $k = \nu/d \rightarrow \infty$ , we have  $g_k(z)$  converges to  $g_\infty(z) := \frac{1}{1-z} - \frac{1}{2} + \log(1-z)$  up to a constant. See Section S9.2 for the plots of the function of  $g_k(z)$  for different values of  $k$ .

EXAMPLE 3.1. (*Multivariate student's  $t$  distribution with DoF  $\nu = d$* ) It can be easily verified that  $g_1(z) = \log(2)$ . Therefore, any proposal will be accepted since the acceptance rate is always 1 whatever  $h$  is.

EXAMPLE 3.2. (*Multivariate student's  $t$  distribution with DoF  $\nu > d$ , including Gaussian*) From Lemma 3.1, one can see informally: (i) in the “transient phase”, as the proposed “latitude” is almost deterministic which has higher target density, the acceptance probability is almost 1 starting from either the North Pole or the South Pole; (ii) In the “stationary phase”, as the proposed “latitude” approximately follows an autoregressive process, the acceptance rate is well-approximated by a positive constant. This suggests that as long as the target has a “lighter tail” than multivariate student's  $t$  with DoF  $d$ , the “transient phase” of SPS takes at most  $\mathcal{O}(1)$  steps. For comparison, for the standard multivariate Gaussian target, the “transient phase” of the Euclidean RWM takes  $\mathcal{O}(d)$  steps [Christensen et al., 2005].

REMARK 3.1. Our theoretical results don't cover the cases of SPS for targets with heavier tails, such as multivariate student's  $t$  with DoF  $\nu < d$ . In this case, the SPS cannot start from the North Pole since the first proposal of the SPS will be rejected with probability 1. One might consider to start from the South Pole. However, the SPS could get stuck at the South Pole if the proposal variance is large (even if the origin is the mode of the target density). See Section S8.12 for comments. In practice, we suggest choosing the initial state of the SPS as a random state uniformly sampled on the sphere.

## 4. SPS: Extension to Elliptical Targets.

4.1. *Extensions of Stereographic Projection.* Same as the previous section, we denote the state of the Markov chain by  $x = (x_1, \dots, x_d)^T \in \mathbb{R}^d$ . Suppose  $z = (z_1, \dots, z_{d+1})^T$  is the coordinates of a unit sphere in  $\mathbb{R}^{d+1}$  (that is,  $\|z\| = 1$ ).

Now, we map  $z \in \mathbb{S}^d$  to  $x \in \mathbb{R}^d$  by the generalized stereographic projection (GSP)

$$x = \text{GSP}(z) := Q \left( R\sqrt{\lambda_1} \frac{z_1}{1 - z_{d+1}}, \dots, R\sqrt{\lambda_d} \frac{z_d}{1 - z_{d+1}} \right)^T,$$

where  $R$  is the radius parameter,  $Q^T = Q^{-1} \in \mathbb{R}^{d \times d}$  is a rotation matrix,  $\{\lambda_1, \dots, \lambda_d\}$  are non-negative constants. That is, the GSP is obtained by “stretching” via  $\{\lambda_i\}$  and “rotating” via  $Q$  from the SP.

Defining the norm

$$\|x\|_{\Lambda, Q}^2 := x^T Q \Lambda^{-1} Q^T x,$$

where  $\Lambda = \text{Diag}\{\lambda_1, \dots, \lambda_d\}$ , we have the following Generalized Stereographic Projection Sampler (GSPS) with parameters  $R, h, \Lambda, Q$ .

---

**Algorithm 3:** Generalized Stereographic Projection Sampler (GSPS)

---

- Let the current state be  $X^d(t) = x$ ;
  - Compute the proposal  $\hat{X}$ :
    - Let  $z := \text{GSP}^{-1}(x)$ ;
    - Sample independently  $d\tilde{z} \sim \mathcal{N}(0, h^2 I_{d+1})$ ;
    - Let  $dz := d\tilde{z} - \frac{(z^T \cdot d\tilde{z})z}{\|z\|^2}$  and  $\hat{z} := \frac{z + dz}{\|z + dz\|}$ ;
    - The proposal  $\hat{X} := \text{GSP}(\hat{z})$ .
  - $X^d(t+1) = \hat{X}$  with probability  $1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|_{\Lambda, Q}^2)^d}{\pi(x)(R^2 + \|x\|_{\Lambda, Q}^2)^d}$ ; otherwise  $X^d(t+1) = x$ .
- 

Note that even though GSPS can be formulated alternatively as using “ellipsoid” instead of “sphere”, this is not our formulation. We still use the unit sphere to compute  $\hat{z}$  for GSPS and only replace the mapping SP by GSP. That is, we only “stretch” and “rotate” the proposal  $\hat{z}$  via GSP when projecting to  $\mathbb{R}^d$ .

The following example shows that we can extend isotropic targets to elliptical targets using the GSPS.

**EXAMPLE 4.1.** Suppose  $\pi(x)$  is multivariate student’s  $t$  with covariance matrix  $\Sigma = Q\Lambda Q^T$  where  $\Lambda = \text{Diag}\{\lambda_1, \dots, \lambda_d\}$ :

$$\pi(x) \propto \left( 1 + \frac{1}{\nu} x^T \Sigma^{-1} x \right)^{-(\nu+d)/2}$$

Then, for the GSPS with the corresponding  $Q$  and  $\Lambda$ , and  $R^2 = \sum_i \lambda_i$ , the acceptance rate is always 1 for any  $h$ .

The GSPS naturally suggests that one can estimate the covariance matrix  $\Sigma = Q\Lambda Q^T$  under the adaptive MCMC framework, which results in adaptive GSPS. Indeed, if both  $\Lambda$  and  $Q$  are known, then one can normalize GSPS and reduce it to SPS, which is GSPS with  $\Sigma = I_d$ .

As the robustness to estimations of the mean and the covariance matrix of the target is the key to the success of adaptive GSPS, we study the robustness of Gaussian targets in the next section.

4.2. *Robustness for Gaussian Targets.* We consider multivariate Gaussian targets with mean vector  $\mu$  and covariance matrix  $\Sigma$ :

$$\pi_{\mu, \Sigma}(x) \propto \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right),$$

where

$$\Sigma = \text{Diag}(\lambda_1, \dots, \lambda_d), \quad \mu = (\mu_1, \dots, \mu_d)^T.$$

We are interested in the robustness when  $\mu \neq 0$  and  $\Sigma \neq I_d$ .

For comparison with the traditional RWM in  $\mathbb{R}^d$  on the orders of stepsize, we recall that the optimal scaling theory gives an optimal stepsize of  $\mathcal{O}(d^{-1/2})$  for RWM [Roberts et al., 1997]. Our stepsize  $h$  is defined on the unit sphere. When projecting to  $\mathbb{R}^d$ , we multiply by  $R = \mathcal{O}(d^{1/2})$ . Therefore, the optimal stepsize of RWM roughly corresponds to  $h = \mathcal{O}(d^{-1})$  in our setting. In this subsection, we study how large the stepsize  $h$  can be so the expected acceptance probability of SPS goes to 1. This roughly corresponds to a stepsize of  $\mathcal{O}(hd^{1/2})$  in  $\mathbb{R}^d$ . Our results show that the order of the stepsize of SPS when projected back to  $\mathbb{R}^d$  is always no smaller than the order of the optimal stepsize of RWM.

**THEOREM 4.1.** *Assume there exist constants  $C < \infty$  and  $c > 0$  such that  $c \leq \lambda_i \leq C$  and  $|\mu_i| \leq C$  for all  $i = 1, \dots, d$ . Furthermore, assume  $R = d^{1/2}$  and*

$$(10) \quad \left| \sum_{i=1}^d \mu_i^2 - \sum_{i=1}^d (1 - \lambda_i) \right| = \mathcal{O}(d^\alpha),$$

where  $\alpha \leq 1$ . Then, under stationarity that  $X \sim \pi$ , the expected acceptance probability converges to 1 as  $d \rightarrow \infty$

$$\mathbb{E}_{X \sim \pi_{\mu, \Sigma}} \mathbb{E}_{\hat{X} | X} \left[ 1 \wedge \frac{\pi_{\mu, \Sigma}(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi_{\mu, \Sigma}(X)(R^2 + \|X\|^2)^d} \right] \rightarrow 1,$$

for all  $h$  such that

$$(11) \quad h = o\left(\frac{d^{-1}}{\sqrt{\max\left\{\frac{1}{d} \sum_i |1 - \lambda_i|, \frac{1}{d} \sum_i \mu_i^2\right\}}} \wedge d^{-(\frac{1}{2} \vee \alpha)}\right).$$

**PROOF.** See Section S8.8. □

Note that for the standard Gaussian target, we know that the acceptance rate goes to 1 as  $d \rightarrow \infty$  for any  $h$ . Recall that for traditional RWM in dimension  $k$ , the optimal stepsize is  $\mathcal{O}(k^{-1/2})$  in  $\mathbb{R}^k$  under the optimal scaling framework [Roberts et al., 1997], which corresponds to  $h = \mathcal{O}(k^{-1/2}d^{-1/2})$  in the SPS setting. Therefore, Theorem 4.1 suggests that

the “effective dimension” is determined by  $\{\mu_i\}$  and  $\{\lambda_i\}$ : as long as  $\{\lambda_i\}$  and  $\{\mu_i\}$  are uniformly bounded, the “effective dimension” of SPS is never larger than  $d$ , which is the “effective dimension” of the traditional RWM. Furthermore, we can compare the “effective dimension” for SPS with  $d$  using Theorem 4.1 under different settings.

EXAMPLE 4.2. *Suppose  $\mu_i = 0$  for all  $i = 1, \dots, d$ . Moreover,  $\lambda_1 = \lambda_2 = \dots = \lambda_k = 2$  and  $\lambda_{k+1} = \dots = \lambda_d = 1$ , then as long as  $k$  is a fixed number,  $\frac{1}{d} \sum_i |1 - \lambda_i| = \mathcal{O}(kd^{-1})$ . By Theorem 4.1, when  $h = o(k^{-1/2}d^{-1/2})$  then the acceptance rate goes to 1 as  $d \rightarrow \infty$ . This suggests that the “effective dimension” for SPS is no more than  $k$ .*

EXAMPLE 4.3. *Suppose  $\lambda_i = 0$  for all  $i = 1, \dots, d$ . Furthermore,  $\mu_1 = \mu_2 = \dots = 1$  and  $\mu_{k+1} = \dots = \mu_d = 0$  where  $k$  is a fixed number. Then we have  $\frac{1}{d} \sum_i \mu_i^2 = \mathcal{O}(kd^{-1})$ . By Theorem 4.1, the acceptance rate goes to 1 when  $h = o(k^{-1/2}d^{-1/2})$ , which implies the “effective dimension” for SPS is no more than  $k$ .*

One can consider Theorem 4.1 in the context of adaptive MCMC for Gaussian targets, in which  $\{\mu_i\}$  and  $\{1 - \lambda_i\}$  represent the “estimation errors” of the coordinate means and eigenvalues of the covariance matrix of the target. One can choose the “radius”  $R$  of the stereographic projection properly to satisfy Eq. (10) with  $\alpha = \frac{1}{2}$ . This is to properly scale  $R$  so that the “latitude” on the unit sphere of the SPS is  $\mathcal{O}_{\mathbb{P}}(d^{-1/2})$ . Then, according to Theorem 4.1, the “effective dimension” of SPS is smaller than  $d$  if  $\frac{1}{d} \sum_i |1 - \lambda_i| = o(1)$  and  $\frac{1}{d} \sum_i \mu_i^2 = o(1)$ . However, the above results on “effective dimension” are based on Theorem 4.1, which only hold for Gaussian targets. For more results on robustness and optimal acceptance rate for adaptive MCMC, we will study another family of targets in Section 5.

Finally, we show two examples that the result of Theorem 4.1 is tight. First, consider the special case that  $\mu_i = \mu > 0$  and  $\lambda_i = \sigma^2$  for all  $i$ . We assume  $\mu^2 = 1 - \sigma^2$  so that Eq. (10) holds for any  $\alpha \leq 1/2$ . Then Eq. (11) suggests that the acceptance probability goes to zero if  $h = o(d^{-1}/\mu)$ . On the other hand, the target in this special case is a product i.i.d. target with marginal distribution  $\mathcal{N}(\mu, 1 - \mu^2)$ . By our optimal scaling results in Section 5, we know that the acceptance probability does not go to zero if  $h = \mathcal{O}(d^{-1}/\mu)$  (see Lemma 5.1 and Corollary 5.2). Therefore, Eq. (11) is tight. Second, consider the special case that  $\mu_i = \mu = 0$  and  $\lambda_i = \sigma^2 \neq 1$  for all  $i$ . Then Eq. (10) holds for  $\alpha = 1$ . In this case, Eq. (11) suggests that the acceptance probability goes to zero if  $h = o(d^{-\alpha}) = o(d^{-1})$ . On the other hand, the target in this special case is a product i.i.d. target with marginal  $\mathcal{N}(0, \sigma^2)$  where  $\sigma^2 \neq 1$ . If we properly rescale  $R$  and  $\sigma^2$ , it reduces to the case of Lemma 5.1 where the target is the standard Gaussian but  $R \neq d^{1/2}$ . Therefore, by Lemma 5.1, the acceptance probability doesn't converge to zero if  $h = \mathcal{O}(d^{-1})$ . Therefore, Eq. (11) is again tight.

**5. SPS in High-dimensional Problems.** In this section, we shall give a case study of the behavior of SPS for high-dimensional target densities. We shall consider various limits of

TABLE 1  
Examples of  $C_\nu$  and  $C_\nu/(C_\nu - 1)$  for different  $\nu$  in Remark 5.1 and Remark 5.4

	$\nu = 3$	$\nu = 5$	$\nu = 10$	$\nu = 20$	$\nu = 50$	$\nu = 100$
$C_\nu$	7.1285	3.0187	1.7521	1.3336	1.1250	1.0612
$C_\nu/(C_\nu - 1)$	1.1632	1.4954	2.3297	3.9977	8.9990	17.3328

SPS as  $d \rightarrow \infty$  though to have tractability of this limit we need to consider a very specialised class of target densities. To this end, analogous to [Roberts et al., 1997] we assume the target  $\pi(x)$  has a product i.i.d. form.

5.1. *Assumptions on  $\pi$ .* We assume the target  $\pi(x)$  has a product i.i.d. form:

$$(12) \quad \pi(x) = \prod_{i=1}^d f(x_i).$$

Without loss of generality, we assume  $f$  is normalized such that

$$(13) \quad \mathbb{E}_f(X^2) = \int x^2 f(x) dx = 1, \quad \mathbb{E}_f(X^6) < \infty.$$

We further assume  $f'/f$  is Lipschitz continuous,  $\lim_{x \rightarrow \pm\infty} x f'(x) = 0$ , and

$$(14) \quad \mathbb{E}_f \left[ \left( \frac{f'(X)}{f(X)} \right)^8 \right] < \infty, \quad \mathbb{E}_f \left[ \left( \frac{f''(X)}{f(X)} \right)^4 \right] < \infty, \quad \mathbb{E}_f \left[ \left( \frac{X f'(X)}{f(X)} \right)^4 \right] < \infty.$$

REMARK 5.1. Under the assumption  $\mathbb{E}_f(X^2) = 1$ , by Cauchy–Schwarz inequality

$$\mathbb{E}_f \left[ ((\log f)')^2 \right] \geq 1$$

where the equality is achieved by standard Gaussian (or truncated standard Gaussian). For univariate student's  $t$  distribution with any DoF =  $\nu > 2$ , rescaling by a factor  $\sqrt{\frac{\nu-2}{\nu}}$ , we can obtain a target density  $f_\nu$  with  $\mathbb{E}_{f_\nu}(X^2) = 1$  and

$$C_\nu := \mathbb{E}_{f_\nu} \left[ ((\log f_\nu)')^2 \right] = \left( \frac{\nu}{\nu-2} \right) \left( \frac{\nu+1}{\nu} \right) \left( \frac{\nu+4}{\nu+3} \right) \sqrt{\frac{\nu+4}{\nu}} > 1.$$

where  $\nu \rightarrow \infty$  recovers the case for the standard Gaussian target. See Table 1 for values of  $C_\nu$  for different  $\nu$ . One can see that  $C_\nu$  is very close to 1 for even a medium size of  $\nu$ .

In this section, we consider  $f$  has full support in  $\mathbb{R}$ . Then the product i.i.d. target  $\pi$  is not isotropic unless it is the standard Gaussian.

5.2. *Acceptance Probability.* To make progress, we need a detailed understanding of the high-dimensional behavior of the acceptance probability of SPS. It turns out that to optimally apply SPS, we need to scale  $R$  to be  $\mathcal{O}(d^{1/2})$ . Not doing this will concentrate mass at either the North or South poles in an ultimately degenerate way. We shall thus assume  $R$  to be scaled in this way in what follows. Therefore, we consider  $R = \sqrt{\lambda d}$  for a fixed constant  $\lambda$ .



To simplify the final result, we replace the step size  $h$  by another parameter  $\ell$  via

$$(15) \quad h = \frac{1}{\sqrt{d-1}} \left[ \frac{1}{\left(1 - \frac{\ell^2}{2d} \frac{4\lambda}{(1+\lambda)^2}\right)^2} - 1 \right]^{1/2}$$

which implies  $\frac{1}{\sqrt{1+h^2(d-1)}} = 1 - \frac{\ell^2}{2d} \frac{4\lambda}{(1+\lambda)^2}$ . Note that  $\ell$  is simply a re-parameterization of  $h$ .

When  $\ell$  is a fixed constant,  $h$  is scaled as  $\mathcal{O}(d^{-1})$  since  $h \approx \frac{\ell}{d} \sqrt{4\lambda/(1+\lambda)^2}$ .

LEMMA 5.1. *Under the assumptions on  $\pi$  in Section 5.1, suppose the current state  $X \sim \pi$ , and the parameter of the algorithm  $R = \sqrt{\lambda d}$ , where  $\lambda > 0$  is a fixed constant. We re-parameterize  $h$  by  $\ell$  according to Eq. (15). Then, if either  $\lambda \neq 1$  or  $f$  is not the standard Gaussian density, there exists a sequence of sets  $\{F_d\}$  such that  $\pi(F_d) \rightarrow 1$  and*

$$\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(X)(R^2 + \|X\|^2)^d} - 1 \wedge \exp(W_{\hat{X}|X}) \right| \right] = o(d^{-1/4} \log(d)),$$

where  $W_{\hat{X}|X} \sim \mathcal{N}(\mu, \sigma^2)$  and

$$\mu = \frac{\ell^2}{2} \left\{ \frac{4\lambda}{(1+\lambda)^2} - \mathbb{E}_f \left[ ((\log f)')^2 \right] \right\}, \quad \sigma^2 = \ell^2 \left\{ \mathbb{E}_f \left[ ((\log f)')^2 \right] - \frac{4\lambda}{(1+\lambda)^2} \right\}.$$

PROOF. See Section S8.9. □

REMARK 5.2. *Lemma 5.1 requires either  $\lambda \neq 1$  or  $\pi$  is not the standard multivariate Gaussian. If  $\lambda = 1$  and  $\pi$  is the standard multivariate Gaussian, the case reduces to isotropic targets discussed in Section 3.2, so the Gaussian approximation in Lemma 5.1 doesn't hold.*

5.3. *Optimisation and robustness of SPS.* Note that Lemma 5.1 suggests that the expected acceptance probability in the stationary phase

$$\mathbb{E}_{X \sim \pi} \mathbb{E}_{\hat{X}|X} \left[ 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(X)(R^2 + \|X\|^2)^d} \right] \rightarrow 2\Phi\left(-\frac{\sigma}{2}\right).$$

Furthermore, as  $\mathbb{E}[\|\hat{X} - X\|^2] \rightarrow \ell^2$ , we obtain the commonly used approximation of Expected Squared Jumping Distance (ESJD):

$$(16) \quad \mathbb{E}[\|\hat{X} - X\|^2] \cdot \mathbb{E} \left[ 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(X)(R^2 + \|X\|^2)^d} \right] \rightarrow 2\ell^2 \cdot \Phi \left( -\frac{\ell}{2} \sqrt{\mathbb{E}_f \left[ ((\log f)')^2 \right] - \frac{4\lambda}{(1+\lambda)^2}} \right).$$

In this subsection we shall explicitly consider the optimisation of SPS, and take a close look at its relative performance in comparison to standard Euclidean RWM. All the results in this section to date have used the convenient restriction that  $\mathbb{E}_f(X^2) = 1$ . However, at this point we shall need to generalise this notion. Therefore consider density  $f$  to have mean  $m$  and variance  $s^2$ .

We consider SPS with  $R = \sqrt{\lambda(s^2 + m^2)d}$ . and stepsize  $h = \frac{1}{\sqrt{d-1}} \left[ \frac{1}{(1 - \frac{\ell^2}{2d} \frac{4\lambda}{(1+\lambda)^2})^2} - 1 \right]^{1/2}$ . We shall denote this algorithm's approximate limiting ESJD by  $E(\ell, m, s, \lambda)$  and its corresponding acceptance probability:  $A(\ell, m, s, \lambda)$ . The following result is a direct consequence of Eq. (16) applied to a scaled density (dividing by  $\sqrt{s^2 + m^2}$ ):

$$E(\ell, m, s, \lambda) = \sqrt{s^2 + m^2} \cdot 2\ell^2 \Phi \left( \frac{-\ell}{2} \sqrt{\tilde{I}(s, m, \lambda)} \right), \quad A(\ell, m, s, \lambda) = 2\Phi \left( \frac{-\ell}{2} \sqrt{\tilde{I}(s, m, \lambda)} \right),$$

where  $\tilde{I}(s, m, \lambda) = (s^2 + m^2)I - \frac{4\lambda}{(1+\lambda)^2}$  and  $I = \mathbb{E}_f((\log f)'(X)^2)$ .

We note that we can readily recover the the corresponding quantities for RWM:  $E(\ell, m, s, \infty)$  and  $A(\ell, m, s, \infty)$ . In particular, to consider robustness, we shall consider the relative performance ratio

$$\tilde{R}(m, s, \lambda) := \frac{\sup_{\ell} E(\ell, m, s, \lambda)}{\sup_{\ell} E(\ell, m, s, \infty)}$$

To get an expression for  $\tilde{R}$ , we shall follow the standard approach of [Roberts et al. \[1997\]](#) of expressing the efficiency in terms of acceptance rate. To that end, we define

$$\tilde{E}(A, m, s, \lambda) := E(\ell(A), m, s, \lambda)$$

where  $\ell$  is chosen to be the unique solution to  $A(\ell(A), m, s, \lambda) = A$ . A simple calculation yields the following.

**COROLLARY 5.1.** *We have  $\tilde{E}(A, m, s, \lambda) = A \cdot \Phi^{-1} \left( \frac{A}{2} \right) \cdot \frac{4(s^2+m^2)}{\tilde{I}(s, m, \lambda)}$  and  $\tilde{R}(m, s, \lambda) = \frac{(s^2+m^2)I}{\tilde{I}(s, m, \lambda)} = \frac{1}{1-\alpha\beta\gamma}$ , where  $\alpha = \frac{4\lambda}{(1+\lambda)^2}$ ,  $\beta = \frac{s^2}{s^2+m^2}$ ,  $\gamma = \frac{1}{s^2I}$ . Since all of  $\alpha, \beta, \gamma$  are in  $[0, 1]$ , In this situation, SPS is never less efficient than Euclidean RWM.*

**REMARK 5.3.** *Primarily, Corollary 5.1 is a robustness result. However it also highlights the (only) ways in which SPS can fail to outperform Euclidean RWM. The three constants  $\alpha, \beta, \gamma$  characterise different sensitivities of the algorithm.*

*$\alpha$  measures the penalty due to mis-specification of the sphere radius.  $4\lambda/(1+\lambda)^2$  is optimised at  $\lambda = 1$  when there is no penalty for misspecification of the target distribution dispersion.*

*$\beta$  describes the penalty due to mis-location of the hyper-sphere. It is seen that the optimal choice is to locate the sphere at the mean of the target distribution.*

*$\gamma$  is a distribution-specific penalty. It is straightforward by to check by functional calculus that  $\gamma \leq 1$  with equality achieved only in the case where the target density is Gaussian. Thus  $\gamma$  is characterising proximity to Gaussianity.*

*From this result, it can be seen that we can only achieve super-efficiency (convergence complexity more rapid than  $\mathcal{O}(d)$ ) when all three parameters,  $\alpha$ ,  $\beta$  and  $\gamma$  are close to (and converging to) unity.*

5.4. *Maximizing ESJD.* In this subsection, we consider  $\lambda = 1$  only (i.e.,  $R = \sqrt{d}$ ) and prove that the approximate limiting ESJD in Eq. (16) is indeed the limiting maximum ESJD. We will demonstrate (again analogously to [Roberts et al., 1997]) that its limit is optimised by targeting an acceptance probability of 0.234.

DEFINITION 5.1. *Expected Squared Jumping Distance (ESJD):*

$$\text{ESJD} := \mathbb{E}_{X \sim \pi} \mathbb{E}_{\hat{X} | X} \left[ \|\hat{X} - X\|^2 \left( 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(X)(R^2 + \|X\|^2)^d} \right) \right].$$

THEOREM 5.1. *Under the assumptions on the target in Section 5.1, suppose  $f$  is not the standard Gaussian density, the SPS chain is in the stationary phase, and the radius parameter is chosen as  $R = \sqrt{d}$  and re-parameterize  $h$  by  $\ell$  according to Eq. (15) with  $\lambda = 1$ . Then, as  $d \rightarrow \infty$ , we have  $\text{ESJD} \rightarrow 2\ell^2 \cdot \Phi \left( -\frac{\ell}{2} \sqrt{\mathbb{E}_f \left[ ((\log f)')^2 \right] - 1} \right)$ .*

PROOF. See Section S8.10. □

COROLLARY 5.2. *The maximum ESJD is approximately  $\frac{1.3}{\mathbb{E}_f \left[ ((\log f)')^2 \right] - 1}$ , which is achieved when the acceptance rate is about 0.234. The optimal  $\hat{\ell} \approx \frac{2.38}{\sqrt{\mathbb{E}_f \left[ ((\log f)')^2 \right] - 1}}$  and the optimal  $\hat{h} = \frac{1}{\sqrt{d-1}} \left[ \left( 1 - \frac{\hat{\ell}^2}{2d} \right)^{-2} - 1 \right]^{1/2} \approx \frac{\hat{\ell}}{\sqrt{d(d-1)}} \approx \frac{2.38}{\sqrt{d(d-1)}} \frac{1}{\sqrt{\mathbb{E}_f \left[ ((\log f)')^2 \right] - 1}}$ .*

REMARK 5.4. *Compared with the maximum ESJD of RWM, the maximum ESJD of SPS is  $\frac{\mathbb{E}_f \left[ ((\log f)')^2 \right]}{\mathbb{E}_f \left[ ((\log f)')^2 \right] - 1}$  times larger. For example, for the class of distributions defined in Remark 5.1 indexed by  $\nu$ , we can compute  $C_\nu / (C_\nu - 1)$ . See Table 1 for examples of  $C_\nu / (C_\nu - 1)$  for different  $\nu$ . One can see that the maximum ESJD of SPS can be much larger than the maximum ESJD of RWM even for a medium size of  $\nu$ .*

5.5. *Diffusion Limit.* Continuing our analogy to the Euclidean RWM case, we shall provide a diffusion limit result, giving a more explicit description of SPS for high-dimensional situations. However, it is difficult to obtain a diffusion limit directly for SPS. Instead, we shall slightly change the original SPS algorithm in a way that is asymptotically negligible (as  $d$  gets large) but which greatly facilitates our limiting diffusion approach. The revised algorithm is called RSPS.

Note that the difference between SPS and RSPS is that in RSPS two independent proposals  $\hat{X}'$  and  $\hat{X}''$  are first computed. Then the final proposal is composed by the first coordinate of  $\hat{X}'$  and other coordinates of  $\hat{X}''$ , i.e.,  $\hat{X} := (\hat{X}'_1, \hat{X}''_{2:d})$ . This design guarantees  $\hat{X}_1$  is independent with  $\hat{X}_{2:d}$  conditional on the current state, which is the technical condition needed to prove the following result on a diffusion limit.

---

**Algorithm 4: Revised Stereographic Projection Sampler (RSPS)**


---

- Let the current state be  $X^d(t) = x$ ;
  - Compute the proposal  $\hat{X}$ :
    - Let  $z := \text{SP}^{-1}(x)$ ;
    - Sample independently  $d\tilde{z}', d\tilde{z}'' \sim \mathcal{N}(0, h^2 I_{d+1})$ ;
    - Let  $dz' := d\tilde{z} - \frac{(z^T \cdot d\tilde{z}')z}{\|z\|^2}$  and  $dz'' := d\tilde{z} - \frac{(z^T \cdot d\tilde{z}'')z}{\|z\|^2}$ ;
    - Let  $\hat{z}' := \frac{z+dz'}{\|z+dz'\|}$  and  $\hat{z}'' := \frac{z+dz''}{\|z+dz''\|}$ ;
    - Two independent proposals  $\hat{X}' := \text{SP}(\hat{z}')$  and  $\hat{X}'' := \text{SP}(\hat{z}'')$ ;
    - The proposal  $\hat{X} := (\hat{X}'_1, \hat{X}''_{2:d})$ .
  - $X^d(t+1) = \hat{X}$  with probability  $1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(x)(R^2 + \|x\|^2)^d}$ ; otherwise  $X^d(t+1) = x$ .
- 

**THEOREM 5.2.** *Under the assumptions on  $\pi$  in Section 5.1, suppose  $f$  is not the standard Gaussian density and the RSPS chain  $\{X^d(t)\}$  starts from the stationarity, i.e.  $X^d(0) \sim \pi$ , and the radius parameter is chosen as  $R = \sqrt{d}$ . Writing  $X^d(t) = (X_1^d(t), \dots, X_d^d(t))$ , we let  $U^d(t) := X_1^d(\lfloor dt \rfloor)$  be the sequence of the first coordinates of  $\{X^d(t)\}$  sped-up by a factor of  $d$ . Then, as  $d \rightarrow \infty$ , we have  $U^d \Rightarrow U$ , where  $\Rightarrow$  denotes weak convergence in Skorokhod topology, and  $U$  satisfies the following Langevin SDE*

$$dU(t) = (s(\ell))^{1/2} dB(t) + s(\ell) \frac{f'(U(t))}{2f(U(t))} dt,$$

where  $s(\ell) := 2\ell^2 \Phi\left(-\ell \frac{\sqrt{\mathbb{E}_f[(\log f)']^2} - 1}{2}\right)$  is the speed measure for the diffusion process, and  $\Phi(\cdot)$  being the standard Gaussian cumulative density function.

**PROOF.** See Section S8.11. □

**COROLLARY 5.3.** *The optimal acceptance rate for RSPS is 0.234 and the maximum speed of the diffusion limit is  $s(\hat{\ell}) \approx \frac{1.3}{\mathbb{E}_f[(\log f)']^2 - 1}$  where  $\hat{\ell} \approx \frac{2.38}{\sqrt{\mathbb{E}_f[(\log f)']^2} - 1}$ . The optimal*

$$\hat{h} = \frac{1}{\sqrt{d-1}} \left[ \left(1 - \frac{\hat{\ell}^2}{2d}\right)^{-2} - 1 \right]^{1/2} \approx \frac{\hat{\ell}}{\sqrt{d(d-1)}} \approx \frac{2.38}{\sqrt{d(d-1)}} \frac{1}{\sqrt{\mathbb{E}_f[(\log f)']^2} - 1}.$$

**REMARK 5.5.** *Compared with the maximum speed of the diffusion limit of RWM [Roberts et al., 1997], the maximum speed of the diffusion limit of RSPS is  $\frac{\mathbb{E}_f[(\log f)']^2}{\mathbb{E}_f[(\log f)']^2 - 1}$  times larger.*

**REMARK 5.6.** *A reasonable conjecture is that the same diffusion limit holds for the original SPS algorithm. In order to establish the same diffusion limit, if we follow the same arguments as in the proof of Theorem 5.1, it is required to show the acceptance rate term becomes “asymptotically independent” with the first coordinate as a rate of  $\mathcal{O}(d^{-1/2})$ . However, our current technical arguments in Theorem 5.1 can achieve a rate of  $\mathcal{O}(d^{-1/8})$  which is not*

enough for establishing the weak convergence to a diffusion limit. Therefore, we only prove the diffusion limit for RSPS in this paper and leave the proof for SPS as an open problem.

**6. Simulations.** In this section, we study the proposed SPS and SBPS through numerical examples. In most of the examples, we consider  $d = 100$  dimensions and two choices of target distributions, the heavy-tailed multivariate student's  $t$  target with  $d$  degree of freedom and the standard Gaussian target. By default, we choose  $R = \sqrt{d}$  for SPS and SBPS. We refer to Section S9 for additional simulations such as different choices of  $R$ .

6.1. *SPS: Bayesian Cauchy regression.* Here we return to Example 1.1. In Fig. 1, we choose  $a = b = 0.1$ ,  $d = 11$ ,  $n = 15$ , and  $R = \sqrt{d}$ . Random design  $X_i \sim \mathcal{N}(0, I) \in \mathbb{R}^d$  and responses  $\{Y_i\}_{i=1}^n$  are generated by  $Y_i = \alpha_0 + \beta_0^T X_i + \epsilon_i$  where  $\alpha_0 = -2$ ,  $\beta_0 = (-4, -3, -2, -1, 0, 1, 2, 3, 4)^T$ ,  $\{\epsilon_i\}$  are i.i.d. zero-mean Cauchy distribution with scale  $\gamma_0 = 1$ . We take a logarithm transformation for  $\gamma$  and implement both SPS and RWM in  $\mathbb{R}^d$ . RWM starts from  $(100, 100, \dots, 100)$  and SPS starts from the North Pole. From the figure, one can see clearly that: (1) RWM which is not geometric ergodic completely failed; (2) the proposed SPS which is uniformly ergodic converged extremely fast.

6.2. *SPS: ESJD per dimension.* In this example, we study the ESJD for SPS and its robustness to the choice of the radius  $R$ . We first tune the proposal stepsize  $h$  of SPS to get different acceptance rates. Then we plot ESJD per dimension for varying acceptance rates as the efficiency curve of SPS. Fig. 5 shows eight efficiency curves for different choices of  $R$ . The target distribution is multivariate student's  $t$  distribution with  $d = 100$  degrees of freedom. In this setting,  $R = \sqrt{d}$  is optimal. The four subplots in the first column are for  $R < \sqrt{d}$  and the second column contains four cases of  $R > \sqrt{d}$ .

Although we do not plot the ESJD for RWM, the maximum ESJD per dimension for RWM is known to have an order of  $\mathcal{O}(d^{-1})$ . In all cases in Fig. 5, SPS has a much larger ESJD than RWM. For the two subplots in the first row of Fig. 5, since  $R$  is closer to  $\sqrt{d}$ , the acceptance rate cannot be lower than 0.5 whatever the proposal variance is. For all the other efficiency curves, an interesting observation is that the maximum ESJDs are always achieved when the acceptance rate is around 0.234. This suggests the optimal acceptance rate 0.234 is quite robust and not limited to product i.i.d. targets, which is similar to the case of optimal acceptance rate 0.234 for RWM [Yang et al., 2020].

6.3. *SBPS: ESS per Switch.* In this example, we study the efficiency curves of SBPS and BPS in terms of ESS per Switch versus the refresh rate. The first subplot of Fig. 6 contains the proportion of refreshments in all the  $N$  events for varying refresh rates. It is clear that as the refresh rate increases, the proportion of refreshments increases. In the other three subplots of Fig. 6, we plot the logarithm of ESS per Switch as a function of the refresh rate for three cases, the 1-st coordinate, the negative log-density, and the squared 1-st coordinate, respectively. For

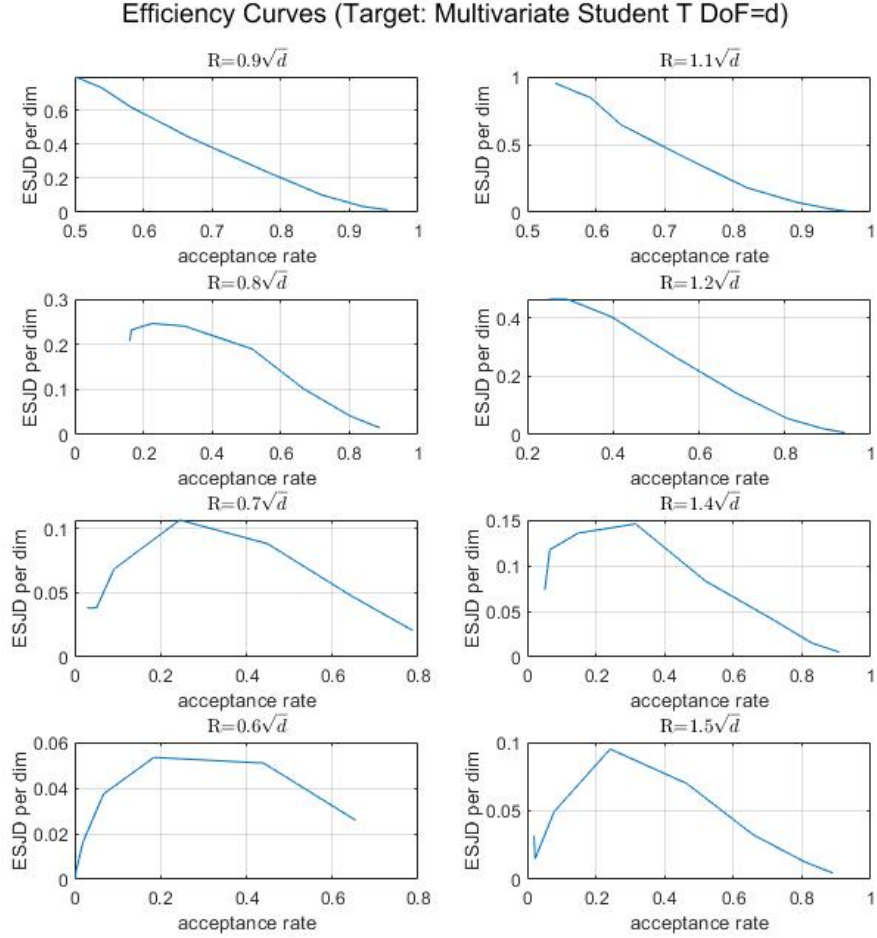


Figure 5: Efficiency Curves (ESJD per dimension) as functions of the acceptance rate of SPS for different choices of  $R$ . Target distribution is multivariate student's  $t$  distribution with DoF =  $d = 100$ . When  $R$  is close to  $\sqrt{d}$  (such as  $R = 0.9\sqrt{d}$  and  $R = 1.1\sqrt{d}$ ), the acceptance rate is always larger than 0.234.

each efficiency curve,  $N = 1000$  events are simulated. We use random initial states for our SBPS. For comparison, BPS starts from stationarity to avoid slow mixing. As SBPS and BPS are continuous-time processes, each unit time period is discretized into 5 samples. The target is standard Gaussian with  $d = 100$ . We refer to Section S9 for additional efficiency curves for heavy-tailed targets.

According to Fig. 6, the ESS per Switch of SBPS is much larger than the ESS per Switch of BPS for all cases (actually the gap becomes larger in higher dimensions). For both the 1-st coordinate and the negative log-target density, the ESS per Switch of SBPS can be larger than 1 if the refresh rate is relatively low. For BPS, however, even starting from stationarity, the ESS per Switch is always much smaller than 1.

## Efficiency Curves: ESS/N vs refresh rate (Gaussian target)

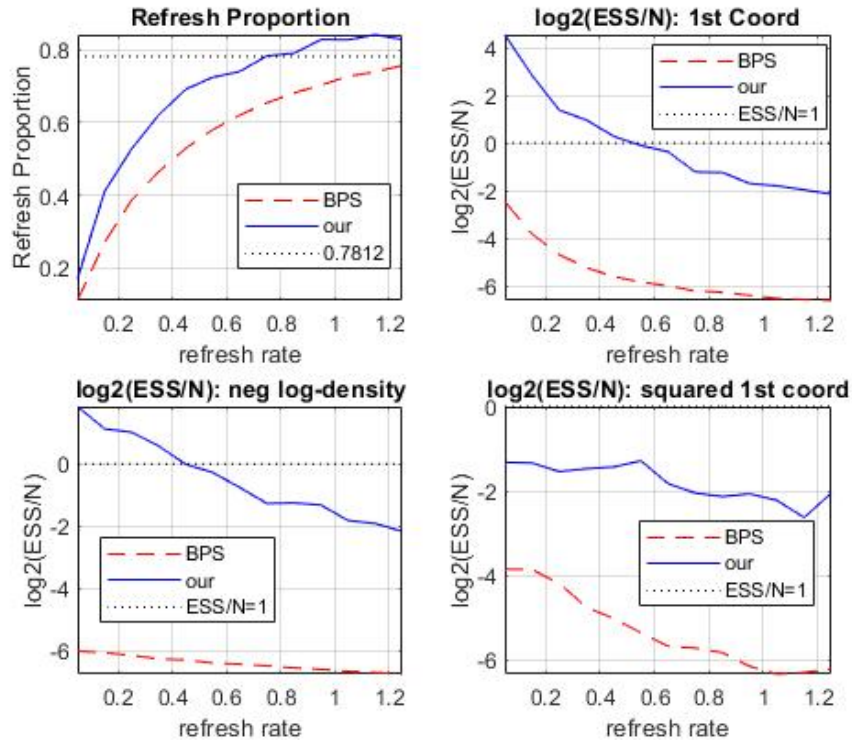


Figure 6: Efficiency curves (ESS per Switch) for SBPS and BPS for varying refresh rate.  $N = 1000$  events are simulated, Random initial values for SBPS and BPS starting from stationarity. Each unit time is discretized to 5 samples. Target distribution: standard Gaussian with  $d = 100$ . The first subplot is the proportion of refreshment events in all  $N$  events. The other three subplots are ESS for the 1-st coordinate, the negative log density, and the squared 1-st coordinate.

**7. Discussion.** We have explored the use of stereographically projected algorithms and developed SPS and SBPS that are uniformly ergodic for a large class of light and heavy-tailed targets and also exhibit fast convergence in high-dimensions. The framework established opens new opportunities for developing MCMC algorithms for sampling high-dimensional and/or heavy-tailed distributions. We finish the paper with some future directions.

- *Adaptive MCMC*: the proposed GSPS algorithm fits the adaptive MCMC framework very well [Roberts and Rosenthal, 2009]. The tuning parameters of GSPS such as the covariance matrix, the location as well as the radius of the sphere can be tuned adaptively. Empirical study of the sensitivity of GSPS to the tuning parameters and extending GSPS for sampling multimodal distributions [Pompe et al., 2020] are important future directions.
- *Quantitative bounds*: we know SPS is uniform ergodicity for a large class of targets and dimension-free for isotropic targets. However, we do not establish quantitative bounds for the mixing time. The quantitative bound and its dimension dependence (e.g., [Yang and Rosenthal, 2023]) would be an interesting direction for future work.

- *Scaling limit for SBPS*: we know SBPS is dimension-free for isotropic targets. For product i.i.d. targets, we establish optimal scaling results for SPS but not for SBPS. Recently, scaling limits for the traditional BPS have been developed in different settings [[Bierkens et al., 2022](#), [Deligiannidis et al., 2021](#)]. It would be interesting to obtain such scaling limits for SBPS for non-isotropic targets to compare with BPS directly.
- *Stereographic MALA, HMC, and others*: another future direction is to develop new MCMC samplers based on other popular MCMC algorithms, such as (Riemann) MALA and Hamiltonian Monte Carlo (HMC) [[Girolami and Calderhead, 2011](#)] and other PDMPs, or based on other mappings from  $\mathbb{R}^d$  to  $\mathbb{S}^d$  than the (generalized) stereographic projection.

**Acknowledgement.** JY was supported by Florence Nightingale Fellowship, Lockey Fund, and St Peter’s College Research Fund (O’Connor Fund) from University of Oxford. KŁ was supported by a Royal Society University Research Fellowship. GOR was supported by EPSRC grants Bayes for Health (R018561) CoSInES (R034710) PINCODE (EP/X028119/1), EP/V009478/1 and by the UKRI grant, OCEAN, EP/Y014650/1.



## Supplement to “Stereographic Markov Chain Monte Carlo”

Jun Yang, Krzysztof Łatuszyński, and Gareth O. Roberts

**S8. Proofs of Main Results.***S8.1. Proof of Proposition 2.1.*

PROOF. If  $\pi(x) > 0, \forall x \in \mathbb{R}^d$  and  $\pi$  being continuous, then  $\pi_S(z)$  is continuous and  $\pi_S(z) > 0$  for all  $z \in \mathbb{S}^d$  except the North Pole. According to [Meyn and Tweedie, 2012, Theorem 13.0.1], it suffices to prove the existence of a stationary distribution, aperiodicity,  $\pi$ -irreducibility, and positive Harris recurrence.

The existence of a stationary distribution comes from the detailed balance of Metropolis–Hastings chain. SPS is aperiodic by [Roberts and Smith, 1994, Theorem 3(i)] using the facts that the acceptance probability does not equal to zero and the proposal distribution as a Markov kernel is aperiodic. Furthermore,  $\pi$ -irreducibility comes from [Roberts and Smith, 1994, Theorem 3(ii)] using the fact that the proposal distribution as a Markov kernel is  $\pi$ -irreducible. Since  $\pi$  is finite, the chain is positive recurrent [Tierney, 1994, pp.1712]. Harris recurrence comes directly from [Tierney, 1994, Corollary 2].

□

*S8.2. Proof of Theorem 2.1.*

PROOF. If there exists  $x$  such that  $\sup_{x \in \mathbb{R}^d} \pi(x)(R^2 + \|x\|^2)^d = \infty$ , then the chain is not even geometric ergodic by [Roberts and Tweedie, 1996a, Proposition 5.1]. This implies that the condition of Eq. (5) is necessary. Therefore, it suffices to prove Eq. (5) is also a sufficient condition.

We prove uniform ergodicity by showing the whole  $\mathbb{S}^d$  is a small set. More precisely, we will show the 3-step transitional kernel satisfies the following minorization condition:

$$P^3(z, A) \geq \epsilon Q(A)$$

for all measurable set  $A$ , where  $Q(\cdot)$  is a fixed probability measure on  $\mathbb{S}^d$  and  $\epsilon > 0$ . The rationale behind establishing a 3-step minorization condition is that this facilitates the SPS chain being able to reach any point on the sphere, and there exists enough probability mass around “2-stop” paths to ensure that the chain can avoid stopping too close to the North Pole, the potentially problematic region.

We use the following notations:

1. We use  $\text{AC}(\epsilon) := \{z \in \mathbb{S}^d : z_{d+1} \geq 1 - \epsilon\}$  to denote the “Arctic Circle” with “size”  $\epsilon \geq 0$ .
2. We use  $\text{HS}(z, 0) := \{z' \in \mathbb{S}^d : z^T z' \geq 0\}$  to denote the “hemisphere” from  $z$ . and  $\text{HS}(z, \epsilon') := \{z' \in \mathbb{S}^d : z^T z' \geq \epsilon'\}$  to denote the area left after “cutting” the hemisphere’s boundary of “size”  $\epsilon' \geq 0$ .

Now consider the 3-step path  $z \rightarrow z_1 \rightarrow z_2 \rightarrow z' \in A$  where  $z \in \mathbb{S}^d$  is the starting point,  $z_1$  and  $z_2$  are two intermediate points, and  $z' \in A$  is final point. Denoting  $q(z, \cdot)$  as the proposal density, we have

$$p(z, z') \geq q(z, z') \left( 1 \wedge \frac{\pi_S(z')}{\pi_S(z)} \right).$$

Then the 3-step transition kernel can be bounded below by

$$\begin{aligned} P^3(z, A) &\geq \int_{z' \in A} \int_{z_2 \in \mathbb{S}^d} \int_{z_1 \in \mathbb{S}^d} q(z, z_1) q(z_1, z_2) q(z_2, z') \\ &\quad \cdot \left( 1 \wedge \frac{\pi_S(z_1)}{\pi_S(z)} \right) \left( 1 \wedge \frac{\pi_S(z_2)}{\pi_S(z_1)} \right) \left( 1 \wedge \frac{\pi_S(z')}{\pi_S(z_2)} \right) \nu(dz_1) \nu(dz_2) \nu(dz') \end{aligned}$$

where  $\nu(\cdot)$  denotes the Lebesgue measure on  $\mathbb{S}^d$ .

By the conditions of Theorem 2.1,  $\pi(x) > 0$  is continuous in  $\mathbb{R}^d$ , which implies  $\pi_S(z)$  is positive and continuous in any compact subset of  $\mathbb{S}^d \setminus N$  (recall that  $N$  denotes the North Pole). Then using the fact that a continuous function on a compact set is bounded, if we rule out  $\text{AC}(\epsilon)$ , then  $\pi_S(z)$  is bounded away from 0. That is,

$$\inf_{z \notin \text{AC}(\epsilon)} \pi_S(z) \geq \delta_\epsilon > 0, \forall \epsilon \in (0, 1),$$

where  $\delta_\epsilon \rightarrow 0$  as  $\epsilon \rightarrow 0$ .

Furthermore, since the surface area of  $\mathbb{S}^d$  is finite and  $\pi_S(z) \propto \pi(x)(R^2 + \|x\|^2)^d$  where  $x = \text{SP}(z)$ , we know  $\sup_{x \in \mathbb{R}^d} \pi(x)(R^2 + \|x\|^2)^d < \infty$  implies that  $M := \sup \pi_S(z) < \infty$ . Therefore, we have

$$\begin{aligned} P^3(z, A) &\geq \int_{z' \in A} \int_{z_2 \in \mathbb{S}^d \setminus \text{AC}(\epsilon)} \int_{z_1 \in \mathbb{S}^d \setminus \text{AC}(\epsilon)} q(z, z_1) q(z_1, z_2) q(z_2, z') \\ &\quad \cdot \left( 1 \wedge \frac{\delta_\epsilon}{M} \right)^2 \left( 1 \wedge \frac{\pi_S(z')}{M} \right) \nu(dz_1) \nu(dz_2) \nu(dz'). \end{aligned}$$

Next, we consider the term  $q(z, z_1) q(z_1, z_2) q(z_2, z')$ . For our algorithm, the proposal can cover the whole hemisphere except the boundary, by ‘‘cutting’’ the boundary a little bit, the proposal density will be bounded below. That is, we have

$$q(z, z') \geq \delta'_{\epsilon'} > 0, \quad \forall z' \in \text{HS}(z, \epsilon'), \epsilon' \in (0, 1),$$

where  $\delta'_{\epsilon'} \rightarrow 0$  as  $\epsilon' \rightarrow 0$ .

Therefore, if  $z_1 \in \text{HS}(z, \epsilon')$  and  $z_2 \in \text{HS}(z', \epsilon')$  then  $q(z, z_1) \geq \delta'_{\epsilon'}$  and  $q(z_2, z') \geq \delta'_{\epsilon'}$ . Then we have

$$\begin{aligned} &\int_{z_2 \in \text{HS}(z', \epsilon') \setminus \text{AC}(\epsilon)} \int_{z_1 \in \text{HS}(z, \epsilon') \setminus \text{AC}(\epsilon)} q(z, z_1) q(z_1, z_2) q(z_2, z') \nu(dz_1) \nu(dz_2) \\ &\geq (\delta'_{\epsilon'})^3 \int_{z_2 \in \text{HS}(z', \epsilon') \setminus \text{AC}(\epsilon)} \int_{z_1 \in \text{HS}(z, \epsilon') \setminus \text{AC}(\epsilon)} \mathbf{1}_{z_1 \in \text{HS}(z_2, \epsilon')} \nu(dz_1) \nu(dz_2). \end{aligned}$$

Note that 3 steps are enough to reach any point  $z'$  from any  $z$  on the sphere. It is clear that under the Lebesgue measure  $\nu$  on  $\mathbb{S}^d$ , we can define the following positive constant  $C$ :

$$C := \inf_{z, z' \in \mathbb{S}^d} \int_{z_2 \in \text{HS}(z', 0) \setminus \text{AC}(0)} \int_{z_1 \in \text{HS}(z, 0) \setminus \text{AC}(0)} \mathbf{1}_{z_1 \in \text{HS}(z_2, 0)} \nu(dz_1) \nu(dz_2) > 0.$$

Now consider the following function of  $\epsilon \geq 0$  and  $\epsilon' \geq 0$

$$C_{\epsilon, \epsilon'} := \inf_{z, z' \in \mathbb{S}^d} \int_{z_2 \in \text{HS}(z', \epsilon') \setminus \text{AC}(\epsilon)} \int_{z_1 \in \text{HS}(z, \epsilon') \setminus \text{AC}(\epsilon)} \mathbf{1}_{z_1 \in \text{HS}(z_2, \epsilon')} \nu(dz_1) \nu(dz_2).$$

It is clear that  $C \geq C_{\epsilon, \epsilon'}$ . The difference between them can be bounded by

$$\begin{aligned} C - C_{\epsilon, \epsilon'} &\leq \sup_{z \in \mathbb{S}^d} \nu(\{z_1 : z_1 \in (\text{HS}(z, 0) \setminus \text{HS}(z, \epsilon')) \cap (\text{AC}(\epsilon) \setminus \text{AC}(0))\}) \\ &\quad + \sup_{z' \in \mathbb{S}^d} \nu(\{z_2 : z_2 \in (\text{HS}(z', 0) \setminus \text{HS}(z', \epsilon')) \cap (\text{AC}(\epsilon) \setminus \text{AC}(0))\}) \\ &\quad + \sup_{z_2 \in \mathbb{S}^d} \nu(\{z_1 : z_1 \in \text{HS}(z_2, 0) \setminus \text{HS}(z_2, \epsilon')\}). \end{aligned}$$

Then clearly the upper bound is continuous w.r.t. both  $\epsilon$  and  $\epsilon'$  and goes to zero when  $\epsilon \rightarrow 0$  and  $\epsilon' \rightarrow 0$ . Therefore, there exists  $\epsilon > 0$  and  $\epsilon' > 0$  such that  $C - C_{\epsilon, \epsilon'} \leq \frac{1}{2}C$ , that is

$$C_{\epsilon, \epsilon'} \geq \frac{1}{2}C > 0.$$

Using such  $\epsilon$  and  $\epsilon'$  we have

$$\int_{z_2 \in \text{HS}(z', \epsilon') \setminus \text{AC}(\epsilon)} \int_{z_1 \in \text{HS}(z, \epsilon') \setminus \text{AC}(\epsilon)} q(z, z_1) q(z_1, z_2) q(z_2, z') \nu(dz_1) \nu(dz_2) \geq \frac{1}{2}C(\delta'_{\epsilon'})^3 > 0.$$

Then we have

$$P^3(z, A) \geq \left(1 \wedge \frac{\delta_\epsilon}{M}\right)^2 \frac{1}{2}C(\delta'_{\epsilon'})^3 \int_{z' \in A} \left(1 \wedge \frac{\pi_S(z')}{M}\right) \nu(dz')$$

Note that  $\pi_S(z') > 0$  almost surely w.r.t.  $\nu$ . Denoting

$$g(z') := \left(1 \wedge \frac{\delta_\epsilon}{M}\right)^2 \frac{1}{2}C(\delta'_{\epsilon'})^3 \left(1 \wedge \frac{\pi_S(z')}{M}\right),$$

we have established

$$P^3(z, A) \geq \int_{z' \in A} g(z') \nu(dz'),$$

where  $g(z) > 0$  almost surely w.r.t.  $\nu$ . Then we can define

$$\int_A g(z') \nu(dz') = \int_{\mathbb{S}^d} g(z') \nu(dz') \frac{\int_A g(z') \nu(dz')}{\int_{\mathbb{S}^d} g(z') \nu(dz')} =: \epsilon Q(A),$$

where  $\epsilon = \int_{\mathbb{S}^d} g(z') \nu(dz') > 0$  and  $Q(A) := \frac{\int_A g(z') \nu(dz')}{\int_{\mathbb{S}^d} g(z') \nu(dz')}$  is a probability measure. This established the desired minorization condition and completed the proof.  $\square$

S8.3. *Sketch proof of Proposition 2.2.*

PROOF. In this proof, for two vectors  $a$  and  $b$ , we use  $a \cdot b$  to denote the inner product. We denote the generator of our PDMP and its adjoint (see [Fearnhead et al., 2018, Section 2.2] or [Liggett, 2010]) as  $\mathcal{A}$  and  $\mathcal{A}^*$ , respectively. Then, we can divide the generator  $\mathcal{A}$  into two parts

$$\mathcal{A} = \mathcal{A}_D + \mathcal{A}_J,$$

where the first term  $\mathcal{A}_D$  is the generator of the deterministic dynamics (moving along the greatest circle) and the second term  $\mathcal{A}_J$  relates the jumping process (that is, the change in expectation at bounce or refreshment events).

Then it can be easily derived that, for suitable functions  $f(z, v)$  and  $g(z, v)$ , where  $v \perp z$  and  $\|v\| = \|z\| = 1$ , the generator  $\mathcal{A}_D$  and its adjoint  $\mathcal{A}_D^*$  of the deterministic process satisfy

$$\mathcal{A}_D f = \nabla_z f \cdot v - \nabla_v f \cdot z,$$

$$\mathcal{A}_D^* g = \nabla_v g \cdot z - \nabla_z g \cdot v.$$

Denoting  $|\mathbb{S}^{d-1}|$  as the volume of the  $\mathbb{S}^{d-1}$ , we substitute the following joint distribution on  $\mathbb{S}^d \times \mathbb{S}^d$ :

$$\pi(z, v) = \pi_S(z)p(v | z) = \pi_S(z) \frac{\mathbf{1}_{v \cdot z=0}}{|\mathbb{S}^{d-1}|}, \quad \forall (z, v) \in \mathbb{S}^d \times \mathbb{S}^d.$$

For any  $(z, v) \in \mathbb{S}^d \times \mathbb{S}^d$  such that  $z \perp v$ , we have

$$\begin{aligned} \mathcal{A}_D^*[\pi(z, v)] &= 0 \cdot z - \frac{1}{|\mathbb{S}^{d-1}|} \nabla_z \pi_S(z) \cdot v \\ &= -\pi_S(z) \frac{\mathbf{1}_{v \cdot z=0}}{|\mathbb{S}^{d-1}|} [v \cdot \nabla_z \log \pi_S(z)] \\ &= -\pi(z, v) [v \cdot \nabla_z \log \pi_S(z)]. \end{aligned}$$

For  $\mathcal{A}_J$ , since the refreshment clearly preserves  $\pi$  as the stationary distribution, we only need to consider the bounce events. When there is no refreshment, we can follow similar arguments as in [Fearnhead et al., 2018, Section 2.2 and 3.2] to get the adjoint of the second part of the generator

$$\mathcal{A}_J^*[\pi(z, v)] = -\pi(z, v)\lambda(z, v) + \int \lambda(z, v')q(v | z, v')\pi(z, v')dv',$$

where

$$q(\cdot | z, v) := \delta_{\tilde{v}}(\cdot), \quad \tilde{v} := v - 2 \left[ \frac{v \cdot \tilde{\nabla}_z \log \pi_S(z)}{\tilde{\nabla}_z \log \pi_S(z) \cdot \tilde{\nabla}_z \log \pi_S(z)} \right] \tilde{\nabla}_z \log \pi_S(z).$$

Therefore, it suffices to verify that  $\pi(z, v)$  satisfies  $\mathcal{A}^*[\pi(z, v)] = \mathcal{A}_D^*[\pi(z, v)] + \mathcal{A}_J^*[\pi(z, v)] = 0$ . That is,

$$\mathcal{A}^*[\pi(z, v)] = -\pi(z, v)[v \cdot \nabla_z \log \pi_S(z) + \lambda(z, v)] + \int \lambda(z, v')q(v | z, v')\pi(z, v')dv'.$$

Then it is enough to verify

$$p(v | z)[\lambda(z, v) + v \cdot \nabla_z \log \pi_S(z)] - \lambda(z, \tilde{v})p(\tilde{v} | z) = 0.$$

Note that  $\tilde{v} \perp z$  which implies that  $p(v | z) = p(\tilde{v} | z)$ . Therefore, it suffices to verify

$$\lambda(z, v) - \lambda(z, \tilde{v}) = -v \cdot \nabla_z \log \pi_S(z).$$

Recall that in our PDMP, we have  $\lambda(z, v) = \max\{0, [-v \cdot \nabla_z \log \pi_S(z)]\}$ . Therefore, we only need to verify

$$-v \cdot \nabla_z \log \pi_S(z) = \tilde{v} \cdot \nabla_z \log \pi_S(z).$$

Using the definition of  $\tilde{v}$  and  $\tilde{\nabla}_z \log \pi_S(z) = \nabla_z \log \pi_S(z) - [z \cdot \nabla_z \log \pi_S(z)]z$ , together with the fact that  $v \perp z$  and  $\tilde{v} \perp z$ , we have

$$-v \cdot \nabla_z \log \pi_S(z) = -v \cdot \tilde{\nabla}_z \log \pi_S(z) = \tilde{v} \cdot \tilde{\nabla}_z \log \pi_S(z) = \tilde{v} \cdot \nabla_z \log \pi_S(z),$$

which completes the proof.  $\square$

#### S8.4. Proof of Theorem 2.2.

PROOF. Recall the definition of the transformation between  $x \in \mathbb{R}^d$  and  $z \in \mathbb{S}^d$ :

$$x_i = R \frac{z_i}{1 - z_{d+1}}, \quad z_i = \frac{2Rx_i}{R^2 + \|x\|^2}, \quad \forall i = 1, \dots, d, \quad z_{d+1} = \frac{\|x\|^2 - R^2}{\|x\|^2 + R^2},$$

and the transformed target on  $\mathbb{S}^d$  is  $\pi_S(z) \propto \pi(x)(R^2 + \|x\|^2)^d$ . One version of  $\nabla_z \log \pi_S(z)$  (there is no unique representation since  $z \in \mathbb{R}^{d+1}$  but  $x \in \mathbb{R}^d$ ) can be derived as

$$\begin{aligned} \frac{\partial \log \pi_S(z)}{\partial z_i} &= \frac{\partial \log \pi(x)}{\partial x_i} \frac{R^2 + \|x\|^2}{2R}, \quad i = 1, \dots, d. \\ \frac{\partial \log \pi_S(z)}{\partial z_{d+1}} &= \left[ \frac{1}{d} \sum_{j=1}^d \frac{\partial \log \pi(x)}{\partial x_j} x_j + 1 \right] \frac{d}{2R^2} (R^2 + \|x\|^2). \end{aligned}$$

Recall that we define  $\tilde{\nabla}_z \log \pi_S(z)$  as the projection of (any version of)  $\nabla_z \log \pi_S(z)$  onto the (tangent plane of the) sphere:

$$\tilde{\nabla}_z \log \pi_S(z) = \nabla_z \log \pi_S(z) - [z \cdot \nabla_z \log \pi_S(z)]z.$$

Then the representation of  $\tilde{\nabla}_z \log \pi_S(z)$  is unique.

Denote  $(\cdot)_{d+1}$  as the  $(d+1)$ -th coordinate. One can verify that

$$\limsup_{\{z: z_{d+1} \rightarrow 1\}} \left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} = \limsup_{\{x: \|x\| \rightarrow \infty\}} \sum_{i=1}^d \left( \frac{\partial \log \pi(x)}{\partial x_i} x_i \right) + 2d.$$

Therefore, the condition is equivalent to

$$\limsup_{\{z: z_{d+1} \rightarrow 1\}} \left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} < \frac{1}{2}.$$

We have assumed  $\|v\| = 1$ . Define the Lyapunov function

$$f(z, v) = 2 + v_{d+1},$$

which is bounded since  $|v_{d+1}| \leq 1$ . Note that moving along the greatest circle on the sphere with constant speed  $\|v\| = 1$  implies the equations

$$\begin{aligned} z(t) &= z(0) \cos(t) + v(0) \sin(t), \\ v(t) &= v(0) \cos(t) - z(0) \sin(t). \end{aligned}$$

For simplicity, we will write  $z(0)$  and  $v(0)$  as  $z$  and  $v$  in the rest of the proof. For two vectors  $a$  and  $b$ , we will use  $a \cdot b$  to denote the inner product  $a^T b$ . The key property we will use is that  $\frac{df}{dt} = \frac{dv_{d+1}}{dt} = -z_{d+1}$ .

Now we can derive the (extended) generator of SBPS

$$\begin{aligned} \tilde{\mathcal{L}}f(z, v) &:= \lim_{t \rightarrow 0^+} \frac{1}{t} [f(z(t), v(t)) - f(z, v)] \\ &\quad + \int \lambda \psi(dw) [f(z, w) - f(z, v)] \\ &\quad + (-\nabla_z \log \pi_S(z) \cdot v)^+ [f(z, \tilde{v}) - f(z, v)] \\ &= -z_{d+1} - \lambda v_{d+1} + (-\nabla_z \log \pi_S(z) \cdot v)^+ (\tilde{v}_{d+1} - v_{d+1}), \end{aligned}$$

where we have used

$$\int \lambda \psi(dw) [f(z, w) - f(z, v)] = \int \lambda \psi(dw) [w_{d+1} - v_{d+1}] = -\lambda v_{d+1}.$$

Our goal is to establish a drift condition for the generator. Note that

$$\tilde{\mathcal{L}}f/f = \tilde{\mathcal{L}}f/(2 + v_{d+1}) \leq \max\{\tilde{\mathcal{L}}f, \tilde{\mathcal{L}}f/3\}$$

Therefore, we only need to focus on  $\tilde{\mathcal{L}}f$  and prove it is negative outside a small set. It suffices to prove  $\tilde{\mathcal{L}}f$  is negative in a small neighborhood of the North Pole.

Recall that

$$\tilde{v} := v - 2 \left[ \frac{v \cdot \tilde{\nabla}_z \log \pi_S(z)}{\tilde{\nabla}_z \log \pi_S(z) \cdot \tilde{\nabla}_z \log \pi_S(z)} \right] \tilde{\nabla}_z \log \pi_S(z).$$

Also, it can be verified that

$$v \cdot \tilde{\nabla}_z \log \pi_S(z) = v \cdot \nabla_z \log \pi_S(z).$$

Therefore, we have

$$\tilde{\mathcal{L}}f(z, v) = -z_{d+1} - \lambda v_{d+1} + \left( -\tilde{\nabla}_z \log \pi_S(z) \cdot v \right)^+ \left[ -2 \frac{v \cdot \tilde{\nabla}_z \log \pi_S(z)}{\|\tilde{\nabla}_z \log \pi_S(z)\|^2} \left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} \right].$$

Note that if  $\left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} \leq 0$  then  $\tilde{\mathcal{L}}f(z, v) \leq -z_{d+1} - \lambda v_{d+1}$ . Otherwise,  $\left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} > 0$ , then we maximize the last term over  $v$  to get

$$\left( -\tilde{\nabla}_z \log \pi_S(z) \cdot v \right)^+ \left[ -2 \frac{v \cdot \tilde{\nabla}_z \log \pi_S(z)}{\|\tilde{\nabla}_z \log \pi_S(z)\|^2} \left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} \right]$$

$$\begin{aligned} &\leq \|\tilde{\nabla}_z \log \pi_S(z)\| \left[ 2 \frac{\|\tilde{\nabla}_z \log \pi_S(z)\|}{\|\tilde{\nabla}_z \log \pi_S(z)\|^2} \left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} \right] \\ &= 2 \left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} \end{aligned}$$

where the upper bound is achieved if

$$v = \frac{-\tilde{\nabla}_z \log \pi_S(z)}{\|\tilde{\nabla}_z \log \pi_S(z)\|}.$$

Overall, using  $v_{d+1} \leq \sqrt{1 - z_{d+1}^2}$ , we have already shown that

$$\sup_{\{z: z_{d+1}=z_*\}} \tilde{\mathcal{L}}f(z, v) \leq -z_* + \lambda \sqrt{1 - z_*^2} + 2 \sup_{\{z: z_{d+1}=z_*\}} \left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1}^+$$

Finally, let  $z_* \rightarrow 1$ , the right hand side goes to

$$-1 + 2 \limsup_{\{z: z_{d+1} \rightarrow 1\}} \left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1}^+$$

which is negative because of the condition in Theorem 5.1 that

$$\limsup_{\{z: z_{d+1} \rightarrow 1\}} \left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} < \frac{1}{2}.$$

Therefore, we have shown the (extended) generator is negative in a neighborhood of the North Pole, which implies a drift condition for the (extended) generator. Finally, by [Down et al., 1995, Theorem 5.2(c)], as the Lyapunov function is bounded, the Markov process is uniformly ergodic, which completes the proof.  $\square$

### S8.5. Proof of Corollary 2.1.

PROOF. We present two versions of the proof. In this first version of the proof, we simply check the condition in Theorem 2.2. Let  $\nu$  to be the DoF of the multivariate student's  $t$  target  $\pi(x)$ . We have

$$\log \pi(x) = -\frac{\nu + d}{2} \log \left( 1 + \frac{1}{\nu} \|x\|^2 \right) + C$$

for some constant  $C$ . Then, one can get

$$\sum_{i=1}^d \frac{\partial \log \pi(x)}{\partial x_i} x_i = -\frac{\nu + d}{\nu} \frac{\|x\|^2}{1 + \|x\|^2/\nu}.$$

Taking  $\|x\| \rightarrow \infty$  yields

$$-\frac{\nu + d}{\nu} \nu < \frac{1}{2} - 2d$$

which completes the proof since

$$\nu > d - \frac{1}{2}.$$

The alternative proof is by following the proof of Theorem 2.2. Let the degree of freedom for the multivariate student's  $t$  target  $\pi(x)$  to be  $kd$  where  $k > 1 - \frac{1}{2d}$ . For simplicity of notations, we consider the case  $R = \sqrt{d}$  without loss of generality. Since  $\pi(x)$  is isotropic, we can verify that if  $R = \sqrt{d}$ , one version of  $\nabla \log \pi_S(z)$  (as the representation is not unique) is

$$\nabla \log \pi_S(z) = \left( 0, \dots, 0, \frac{-dz_{d+1}(k-1)}{(1-z_{d+1})[(k-1)(1-z_{d+1})+2]} \right)$$

Then it is easy to verify using the basic geometry of the sphere that if  $k \geq 1$  and  $z_{d+1} > 0$  then

$$\left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} \leq 0.$$

If  $k < 1$  and  $z_{d+1} > 0$  then

$$\frac{\left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1}}{\|\tilde{\nabla}_z \log \pi_S(z)\|} = \sqrt{1 - z_{d+1}^2}$$

and

$$\|\tilde{\nabla}_z \log \pi_S(z)\| = \sqrt{1 - z_{d+1}^2} \frac{-dz_{d+1}(k-1)}{(1-z_{d+1})[(k-1)(1-z_{d+1})+2]}$$

Denoting  $k = 1 - \epsilon$  and multiplying the two equations above yields

$$\left( \tilde{\nabla}_z \log \pi_S(z) \right)_{d+1} = \frac{d\epsilon z_{d+1}(1 - z_{d+1}^2)}{(1 - z_{d+1})[2 - \epsilon(1 - z_{d+1})]} = \frac{d\epsilon z_{d+1}(1 + z_{d+1})}{2 - \epsilon(1 - z_{d+1})}$$

Letting  $z_{d+1} \rightarrow 1$  yields

$$d\epsilon < d \frac{1}{2d} = \frac{1}{2}.$$

Therefore, the uniform ergodicity holds by Theorem 2.2.  $\square$

**S8.6. On the conjecture in Remark 2.2.** Our arguments include two parts. We first study a one-dimensional target distribution with density  $\pi(t) \propto |t|^{-\alpha}$ ,  $t \in [-1, 1]$  and  $\alpha > 0$ . We show that, for this target, the expected hitting time to the boundary (that is 1 or  $-1$ ) of the BPS starting from  $t = 0$  is finite (i.e., the BPS “escapes” the origin) if and only if  $\alpha < 1$ . Next, we show that, for target with the multivariate student's  $t$  distribution with DoF  $\nu$ , we can approximate the density on the greatest circle passing the North Pole by a form of  $\pi(t) \propto |t|^{-\alpha}$  if the chain is close to the North Pole. Then  $\alpha < 1$  corresponds to the condition in our conjecture in Remark 2.2, which is  $\nu > d - 1$  in this case.

**S8.6.1. One-dimensional target.** We consider the following one-dimensional target with density

$$f_0(t) \propto |t|^{-\alpha}, \quad \forall |t| < 1,$$



where  $\alpha > 0$ . Note that the density goes to infinity as  $t \rightarrow 0$ . To study the behavior of BPS, we consider the following “truncated” targets:

$$f_\epsilon(t) \propto (\epsilon + |t|)^{-\alpha}, \quad \forall t \in [-(1-\epsilon), 1-\epsilon]$$

and we will later let  $\epsilon \rightarrow 0$  to recover  $f_0(t)$ .

Next, we consider  $T_\epsilon$  to be the first bouncing time when the starting state is from 0 with unit velocity for target  $f_\epsilon$ . Then we have

$$\mathbb{P}(T_\epsilon < 1 - \epsilon) = 1 - \exp\left(-\int_0^{1-\epsilon} (\log f_\epsilon)' dt\right)$$

Note that  $f_\epsilon(1-\epsilon) \propto 1$  and  $f_\epsilon(0) \propto \epsilon^{-\alpha}$ , both are up to the same constant. Therefore, we have

$$\int_0^{1-\epsilon} (\log f_\epsilon)' dt = \log(1) - \log(\epsilon^{-\alpha})$$

which implies

$$\mathbb{P}(T_\epsilon < 1 - \epsilon) = 1 - \epsilon^{-\alpha}.$$

Similarly, we can get the density of  $T_\epsilon$  as

$$f_{T_\epsilon}(t) = -\epsilon^\alpha \frac{d}{dt} (\epsilon + |t|)^{-\alpha} = \alpha \epsilon^\alpha (\epsilon + |t|)^{-\alpha-1} = \frac{\alpha}{\epsilon} \left(\frac{|t|}{\epsilon} + 1\right)^{-1-\alpha}, \quad \forall 0 < t < 1 - \epsilon.$$

Then, we define  $\tilde{T}_\epsilon$  as the conditional bouncing time which is  $T_\epsilon$  conditioned on  $T_\epsilon < 1 - \epsilon$ . Then we have the density of  $\tilde{T}_\epsilon$

$$f_{\tilde{T}_\epsilon}(t) = \frac{\alpha}{\epsilon(1-\epsilon^\alpha)} \left(\frac{t}{\epsilon} + 1\right)^{-1-\alpha}, \quad t \in [0, 1-\epsilon].$$

Now we consider the “hitting time” to the boundary (that is, to  $\pm(1-\epsilon)$ ). We can denote the “hitting time” as

$$S^\epsilon = \sum_{i=1}^{G^\epsilon} (2\tilde{T}_{\epsilon,i}) + (1-\epsilon),$$

where  $\tilde{T}_{\epsilon,i}$  are i.i.d. from density  $f_{\tilde{T}_\epsilon}$  and  $G^\epsilon$  is an independent geometric random variable with parameter  $\epsilon^\alpha$ . Note that

$$\mathbb{E}[G^\epsilon] = \epsilon^{-\alpha}.$$

Therefore, to study  $\lim_{\epsilon \rightarrow 0} \mathbb{E}[S^\epsilon] < \infty$ , it suffices to study

$$\lim_{\epsilon \rightarrow 0} \mathbb{E} \left[ \sum_{i=1}^{G^\epsilon} \tilde{T}_{\epsilon,i} \right] = \lim_{\epsilon \rightarrow 0} \epsilon^{-\alpha} \mathbb{E}[\tilde{T}_\epsilon].$$

Using the density of  $\tilde{T}_\epsilon$ , we have

$$\begin{aligned}
\mathbb{E}[\tilde{T}_\epsilon] &= \frac{\alpha}{\epsilon(1-\epsilon^\alpha)} \int_0^{1-\epsilon} t \left(\frac{t}{\epsilon} + 1\right)^{-1-\alpha} dt \\
&= \frac{\alpha}{\epsilon(1-\epsilon^\alpha)} \int_0^{1-\epsilon} \left[ \epsilon \left(\frac{t}{\epsilon} + 1\right) - \epsilon \right] \left(\frac{t}{\epsilon} + 1\right)^{-1-\alpha} dt \\
&= \frac{\alpha}{\epsilon(1-\epsilon^\alpha)} \left[ \int_0^{1-\epsilon} \epsilon \left(\frac{t}{\epsilon} + 1\right)^{-\alpha} dt - \int_0^{1-\epsilon} \epsilon \left(\frac{t}{\epsilon} + 1\right)^{-\alpha-1} dt \right] \\
&= \frac{\alpha}{\epsilon(1-\epsilon^\alpha)} \int_0^{1-\epsilon} \epsilon \left(\frac{t}{\epsilon} + 1\right)^{-\alpha} dt - \epsilon \\
&= \frac{\alpha\epsilon}{1-\epsilon^\alpha} \int_0^{\frac{1-\epsilon}{\epsilon}} (u+1)^{-\alpha} du - \epsilon \\
&= \begin{cases} \frac{\epsilon}{1-\epsilon} \log\left(\frac{1-\epsilon}{\epsilon} + 1\right) - \epsilon, & \text{if } \alpha = 1, \\ \frac{\alpha\epsilon}{1-\epsilon^\alpha} \frac{(1+\frac{1-\epsilon}{\epsilon})^{1-\alpha} - 1}{1-\alpha} - \epsilon, & \text{otherwise.} \end{cases} \\
&= \begin{cases} \frac{\epsilon}{1-\epsilon} \log\left(\frac{1}{\epsilon}\right) - \epsilon, & \text{if } \alpha = 1, \\ \frac{\alpha}{1-\alpha} \frac{\epsilon^\alpha - \epsilon}{1-\epsilon^\alpha} - \epsilon, & \text{otherwise.} \end{cases}
\end{aligned}$$

Therefore, letting  $\epsilon \rightarrow 0$ , if  $\alpha < 1$ , then

$$\epsilon^{-\alpha} \mathbb{E}[\tilde{T}_\epsilon] \rightarrow \epsilon^{-\alpha} \frac{\alpha}{1-\alpha} \epsilon^\alpha = \frac{\alpha}{1-\alpha} < \infty.$$

On the other hand, if  $\alpha = 1$ , then

$$\epsilon^{-\alpha} \mathbb{E}[\tilde{T}_\epsilon] = \epsilon^{-1} \left[ \frac{\epsilon}{1-\epsilon} \log\left(\frac{1}{\epsilon}\right) - \epsilon \right] = \frac{1}{1-\epsilon} \log\left(\frac{1}{\epsilon}\right) - 1 \rightarrow \infty.$$

Or if  $\alpha > 1$ , then

$$\epsilon^{-\alpha} \mathbb{E}[\tilde{T}_\epsilon] = \epsilon^{-\alpha} \left[ \frac{\alpha}{\alpha-1} \frac{\epsilon - \epsilon^\alpha}{1-\epsilon^\alpha} - \epsilon \right] \rightarrow \frac{1}{\alpha-1} \epsilon^{1-\alpha} \rightarrow \infty.$$

Therefore,  $\alpha < 1$  is required for the expected “escaping time” of the BPS to be finite.

**S8.6.2. Approximation round the North Pole.** We pick the greatest circle such that  $z_2 = z_3 = \dots = z_d = 0$  and  $z_1^2 + z_{d+1}^2 = 1$ . Suppose the target distribution

$$\pi_\nu(x) \propto \left(1 + \frac{\|x\|^2}{\nu}\right)^{-\frac{\nu+d}{2}}, \quad x \in \mathbb{R}^d.$$

Then, for any  $z \in \mathbb{S}^d$ , we have

$$\log \pi_S(z) = \log \pi_\nu(x) + d \log(d + \|x\|^2) = C - dg_{\nu/d}(z_{d+1}),$$

where  $C$  is a constant and

$$g_k(z_{d+1}) := \frac{k+1}{2} \log\left(k-1 + \frac{2}{1-z_{d+1}}\right) + \log(1-z_{d+1}).$$

Therefore, let  $k := \nu/d$ , we have

$$\pi_S(z) \propto \exp(-dg_k(z_{d+1})) = \left[ \left( k - 1 + \frac{2}{1 - z_{d+1}} \right)^{\frac{k+1}{2}} \right]^{-d} (1 - z_{d+1})^{-d}.$$

Now consider the picked greatest circle, for any  $z$  such that  $z_1^2 + z_{d+1}^2 = 1$ , we have

$$\begin{aligned} \pi_S(z_1, z_{d+1} \mid z_{2:d} = 0) &\propto \left[ \left( k - 1 + \frac{2}{1 - z_{d+1}} \right)^{\frac{k+1}{2}} \right]^{-d} (1 - z_{d+1})^{-d} \\ &= \left[ \left( k - 1 + \frac{2}{1 - \sqrt{1 - z_1^2}} \right)^{\frac{k+1}{2}} \right]^{-d} \left( 1 - \sqrt{1 - z_1^2} \right)^{-d}. \end{aligned}$$

If  $k$  is fixed, letting  $z_1 \rightarrow 0$  and using the approximation  $\sqrt{1 - z_1^2} \rightarrow 1 - z_1^2/2$ , we have

$$\left[ \left( k - 1 + \frac{4}{z_1^2} \right)^{\frac{k+1}{2}} \right]^d \left( \frac{z_1^2}{2} \right)^{-d} \rightarrow \left[ \left( \frac{4}{z_1^2} \right)^{\frac{k+1}{2}} \right]^{-d} \left( \frac{z_1^2}{2} \right)^{-d} \propto |z_1|^{(k+1)d-2d}$$

Therefore, if we write it as  $|z_1|^{-\alpha}$  then we have

$$\alpha = (1 - k)d = (1 - \nu/d)d = d - \nu.$$

Therefore,  $\alpha < 1$  corresponds to  $\nu > d - 1$ .

### S8.7. Proof of Lemma 3.1.

PROOF. All the randomness for generating the proposal comes from the standard multivariate Gaussian random variable with distribution  $\mathcal{N}(0, I_d)$  in the tangent space. We can denote it using the proposal variance  $h^2$  and  $d$  independent standard Gaussian random variables  $V_1, \dots, V_d \sim \mathcal{N}(0, 1)$ . Then we can use  $\sqrt{\sum_{j=1}^d V_j^2}$  to denote the Euclidean norm of the multivariate Gaussian random variable.

Suppose the current location of the chain is at  $z$ , and we are interested on the  $i$ -th coordinate of the proposal  $\hat{z}$ . If  $z_i \notin \{-1, 1\}$ , in the ‘‘tangent space’’ at  $z$ , there must be  $d - 1$  directions orthogonal to the direction of the  $i$ -th coordinate of the sphere. This leaves only one direction in the ‘‘tangent space’’ at  $z$  which is not orthogonal to the direction of the  $i$ -th coordinate of the sphere. We denote the marginal of the multivariate Gaussian to this direction as  $U_i$ . Clearly,  $U_i \sim \mathcal{N}(0, 1)$ .

Then, using some facts from basic geometry, the random walk in the tangent space without projection back to sphere changes the ‘‘latitude’’ from  $z_{d+1}$  to  $z_{d+1} + \sqrt{1 - z_{d+1}^2} h U_{d+1}$  with distance to the origin equals to  $\sqrt{1 + h^2 \sum_{j=1}^d V_j^2}$ . Therefore, basic geometry tells us after projection back to the sphere, the ‘‘latitude’’ of the proposal satisfies

$$\frac{\hat{z}_{d+1}}{1} = \frac{z_{d+1} + \sqrt{1 - z_{d+1}^2} h U_{d+1}}{\sqrt{1 + h^2 \sum_{j=1}^d V_j^2}}.$$

Similarly for other coordinates, we can get closed forms for any  $\hat{z}_i$  in terms of  $\sum_{j=1}^d V_j^2$  and  $U_i$  for  $i = 1, \dots, d+1$ . Writing  $\sum_{j=1}^d V_j^2 = U_i^2 + U_{i,\perp}^2$ , we have

$$\hat{z}_i = \frac{1}{\sqrt{1 + h^2 \sum_{j=1}^d V_j^2}} \left( z_i + \sqrt{1 - z_i^2} h U_i \right) = \frac{1}{\sqrt{1 + h^2 (U_i^2 + U_{i,\perp}^2)}} \left( z_i + \sqrt{1 - z_i^2} h U_i \right).$$

Furthermore, we have

$$\begin{aligned} \hat{z}_i &= \frac{1}{\sqrt{1 + h^2 \sum_{j=1}^d V_j^2}} \left( z_i + \sqrt{1 - z_i^2} h U_i \right) \\ &= \frac{\sqrt{1 + h^2 U_i^2}}{\sqrt{1 + h^2 \sum_{j=1}^d V_j^2}} \left( \frac{z_i}{\sqrt{1 + h^2 U_i^2}} + \frac{\sqrt{1 - z_i^2} h U_i}{\sqrt{1 + h^2 U_i^2}} \right) \end{aligned}$$

Intuitively, the term  $\frac{z_i}{\sqrt{1 + h^2 U_i^2}} + \frac{\sqrt{1 - z_i^2} h U_i}{\sqrt{1 + h^2 U_i^2}}$  comes from the change of  $i$ -th coordinate caused by  $U_i$ . The other term  $\frac{\sqrt{1 + h^2 U_i^2}}{\sqrt{1 + h^2 \sum_{j=1}^d V_j^2}}$  comes from the change of the  $i$ -th coordinate by the other  $d - 1$  orthogonal directions through the ‘‘curvature’’ of the sphere when projecting back to the sphere.

Note that  $\{U_i, i = 1, \dots, d+1\}$  are dependent, all of which can be written as linear combinations of  $\{V_i\}$ . If  $h = \mathcal{O}(d^{-1/2})$ , observing that  $\sum_{j=1}^d V_j^2 - U_i^2$  is chi-squared distributed with degree of freedom  $d - 1$ , we can write

$$\begin{aligned} \hat{z}_i &= \frac{\sqrt{1 + h^2 U_i^2}}{\sqrt{1 + h^2 \left( \sum_{j=1}^d V_j^2 - U_i^2 \right) + h^2 U_i^2}} \left( \frac{z_i}{\sqrt{1 + h^2 U_i^2}} + \frac{\sqrt{1 - z_i^2} h U_i}{\sqrt{1 + h^2 U_i^2}} \right) \\ &= \frac{\sqrt{1 + h^2 U_i^2} (1 + \mathcal{O}_{\mathbb{P}}(h^4 d))}{\sqrt{\left( 1 + h^2 \left( \sum_{j=1}^d V_j^2 - U_i^2 \right) \right) (1 + h^2 U_i^2)}} \left( \frac{z_i}{\sqrt{1 + h^2 U_i^2}} + \frac{\sqrt{1 - z_i^2} h U_i}{\sqrt{1 + h^2 U_i^2}} \right) \\ &= \left[ \frac{1}{\sqrt{1 + h^2 \left( \sum_{j=1}^d V_j^2 - U_i^2 \right)}} + \mathcal{O}_{\mathbb{P}}(h^4 d) \right] \left[ \left( 1 - \frac{1}{2} h^2 U_i^2 \right) \left( z_i + \sqrt{1 - z_i^2} h U_i \right) + \mathcal{O}_{\mathbb{P}}(h^3) \right] \\ &= \frac{1}{\sqrt{1 + h^2 U_{i,\perp}^2}} \left[ \left( 1 - \frac{1}{2} h^2 U_i^2 \right) z_i - \sqrt{1 - z_i^2} h U_i \right] + \mathcal{O}_{\mathbb{P}}(h^3 + h^4 d z_i), \end{aligned}$$

where  $U_{i,\perp}^2 \sim \chi_{d-1}^2$  which is independent with  $U_i$ . Therefore, in the stationary phase, if  $z_{d+1} = \mathcal{O}(d^{-1/2})$ , then we have

$$\begin{aligned} \hat{z}_{d+1} &= \frac{1}{\sqrt{1 + h^2 (d-1)}} \left[ 1 + \mathcal{O}_{\mathbb{P}}(d^{-1/2} + h^2 d^{1/2}) \right] (z_{d+1} - h U_{d+1}) + \mathcal{O}_{\mathbb{P}}(d^{-1}) \\ &= \left[ \frac{1}{\sqrt{1 + h^2 (d-1)}} + \mathcal{O}_{\mathbb{P}}(d^{-1/2}) \right] (z_{d+1} - h U_{d+1}) + \mathcal{O}_{\mathbb{P}}(d^{-1}) \end{aligned}$$

$$= \frac{1}{\sqrt{1+h^2(d-1)}} (z_{d+1} - hU_{d+1}) + \mathcal{O}_{\mathbb{P}}(d^{-1}).$$

Similarly, in the transient phase, if  $z_{d+1}^2 = 1 - o(h^2)$ , we have

$$\begin{aligned} \hat{z}_{d+1} &= \left[ \frac{1}{\sqrt{1+h^2(d-1)}} + \mathcal{O}_{\mathbb{P}}(d^{-1/2}) \right] \left( 1 - \frac{1}{2}h^2U^2 \right) z_{d+1} + \mathcal{O}_{\mathbb{P}}(h^4d) \\ &= \frac{1}{\sqrt{1+h^2(d-1)}} \left( 1 - \frac{1}{2}h^2U_{d+1}^2 \right) z_{d+1} + \mathcal{O}_{\mathbb{P}}(d^{-1/2}). \end{aligned}$$

□

### S8.8. Proof of Theorem 4.1.

PROOF. Throughout the proof, we assume  $h = o(d^{-1/2})$  so that  $dh^2 \rightarrow 0$ . For bounded random variables, convergence in probability implies convergence in  $L_1$ . Therefore, it suffices to show  $\left( 1 \wedge \frac{\pi_{\mu,\Sigma}(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi_{\mu,\Sigma}(X)(R^2 + \|X\|^2)^d} \right) \xrightarrow{\mathbb{P}} 1$ . Furthermore, for any  $\epsilon > 0$

$$\mathbb{P} \left( \left| 1 \wedge \frac{\pi_{\mu,\Sigma}(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi_{\mu,\Sigma}(X)(R^2 + \|X\|^2)^d} - 1 \right| > \epsilon \right) \leq \mathbb{P} \left( \left| \frac{\pi_{\mu,\Sigma}(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi_{\mu,\Sigma}(X)(R^2 + \|X\|^2)^d} - 1 \right| > \epsilon \right)$$

Therefore, it suffices to show  $\frac{\pi_{\mu,\Sigma}(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi_{\mu,\Sigma}(X)(R^2 + \|X\|^2)^d} \xrightarrow{\mathbb{P}} 1$ , or equivalently

$$\log \frac{\pi_{\mu,\Sigma}(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi_{\mu,\Sigma}(X)(R^2 + \|X\|^2)^d} \xrightarrow{\mathbb{P}} 0.$$

Defining  $\delta_i := \lambda_i^{-1} - 1$ , we can write  $\pi_{\mu,\Sigma}(x)$  as

$$\log \pi_{\mu,\Sigma}(x) = \log \pi_{0,I_d}(x) - \frac{1}{2} \sum_i \delta_i x_i^2 + \sum_i \mu_i (1 + \delta_i) x_i - \frac{1}{2} \sum_i \log(\lambda_i),$$

where  $\pi_{0,I_d}$  is just the standard multivariate Gaussian density. Furthermore, using  $R = d^{1/2}$  and

$$\log \frac{(d + \|\hat{X}\|^2)^d}{(d + \|X\|^2)^d} = d \log \frac{1 - z_{d+1}}{1 - \hat{z}_{d+1}},$$

we have

$$\begin{aligned} \log \frac{\pi_{\mu,\Sigma}(\hat{X})(d + \|\hat{X}\|^2)^d}{\pi_{\mu,\Sigma}(X)(d + \|X\|^2)^d} &= \log \frac{\pi_{\mu,\Sigma}(\hat{X})}{\pi_{\mu,\Sigma}(X)} + d \log \frac{1 - z_{d+1}}{1 - \hat{z}_{d+1}} \\ &= \log \frac{\pi_{0,I_d}(\hat{X})}{\pi_{0,I_d}(X)} + d \log \frac{1 - z_{d+1}}{1 - \hat{z}_{d+1}} - \frac{1}{2} \sum_i \delta_i (\hat{X}_i^2 - X_i^2) + \sum_i \mu_i (1 + \delta_i) (\hat{X}_i - X_i). \end{aligned}$$

Therefore, it suffices to show that if  $h$  satisfies the condition of Theorem 4.1 then we have all of the three results:

$$(17) \quad \log \frac{\pi_{0,I_d}(\hat{X})}{\pi_{0,I_d}(X)} + d \log \frac{1 - z_{d+1}}{1 - \hat{z}_{d+1}} \xrightarrow{\mathbb{P}} 0,$$

$$(18) \quad \sum_i \delta_i (\hat{X}_i^2 - X_i^2) \xrightarrow{\mathbb{P}} 0,$$

$$(19) \quad \sum_i \mu_i (1 + \delta_i) (\hat{X}_i - X_i) \xrightarrow{\mathbb{P}} 0.$$

The rest of the proof is just to show Eqs. (17) to (19) using Taylor expansions and Lemma 3.1.

S8.8.1. *Proof of Eq. (17).* Writing  $X = \mu + \Sigma^{1/2} \tilde{N}$  where  $\tilde{N}$  is a standard multivariate Gaussian vector, because of the conditions in Theorem 4.1 that there exist constants  $C < \infty$  and  $c > 0$  such that  $c \leq \lambda_i \leq C$  and  $|\mu_i| \leq C$ , we have  $\mu^T \Sigma^{1/2} \tilde{N} = \mathcal{O}_{\mathbb{P}}(d^{1/2})$ . Then using the fact  $\tilde{N}^T \Sigma \tilde{N} = \sum_i \lambda_i + \mathcal{O}_{\mathbb{P}}(d^{1/2})$ , we have

$$\begin{aligned} \|X\|^2 &= \|X - \mu\|^2 + \sum_i \mu_i^2 + 2\mu^T \Sigma^{1/2} \tilde{N} \\ &= \tilde{N}^T \Sigma \tilde{N} + \sum_i \mu_i^2 + 2\mu^T \Sigma^{1/2} \tilde{N} \\ &= \left( d - \sum_i (1 - \lambda_i) + \mathcal{O}_{\mathbb{P}}(d^{1/2}) \right) + \sum_i \mu_i^2 + \mathcal{O}_{\mathbb{P}}(d^{1/2}). \end{aligned}$$

Then we have  $\|X\|^2 = d + \sum_i \mu_i^2 - \sum_i (1 - \lambda_i) + \mathcal{O}_{\mathbb{P}}(d^{1/2})$  and

$$(20) \quad z_{d+1} = \frac{\|X\|^2 - d}{\|X\|^2 + d} = \frac{\sum_i \mu_i^2 - \sum_i (1 - \lambda_i)}{2d + \sum_i \mu_i^2 - \sum_i (1 - \lambda_i)} + \mathcal{O}_{\mathbb{P}}(d^{-1/2}).$$

By Eq. (10),  $|\sum_i \mu_i^2 - \sum_i (1 - \lambda_i)| = \mathcal{O}(d^\alpha)$ , we have  $\|X\|^2 = d + \mathcal{O}(d^\alpha) + \mathcal{O}_{\mathbb{P}}(d^{1/2})$  and  $z_{d+1} = \mathcal{O}(d^{\alpha-1}) + \mathcal{O}_{\mathbb{P}}(d^{-1/2})$ , where the  $\mathcal{O}(d^{\alpha-1})$  term represents  $\frac{\sum_i \mu_i^2 - \sum_i (1 - \lambda_i)}{2d + \sum_i \mu_i^2 - \sum_i (1 - \lambda_i)}$  which is bounded away from 1. From Eq. (6) we have

$$(21) \quad \hat{z}_i = (1 + \mathcal{O}_{\mathbb{P}}(h^2 d))(z_i + \mathcal{O}_{\mathbb{P}}(h)) = z_i + \mathcal{O}_{\mathbb{P}}(h) + \mathcal{O}_{\mathbb{P}}(h^2 d) z_i.$$

Using  $h = o(d^{-1/2})$ ,  $h = o(d^{-\alpha})$ , and  $z_{d+1} = \mathcal{O}(d^{\alpha-1}) + \mathcal{O}_{\mathbb{P}}(d^{-1/2})$  we have  $h^2 d z_{d+1} = o_{\mathbb{P}}(d^{-1})$  then from Eq. (21) we have

$$(22) \quad \hat{z}_{d+1} = z_{d+1} + \mathcal{O}_{\mathbb{P}}(h) + o_{\mathbb{P}}(d^{-1}).$$

Then we have

$$\begin{aligned} &\log \frac{\pi_{0,I_d}(\hat{X})}{\pi_{0,I_d}(X)} + d \log \frac{1 - z_{d+1}}{1 - \hat{z}_{d+1}} \\ &= d \left[ \left( \frac{1}{1 - z_{d+1}} - \frac{1}{1 - \hat{z}_{d+1}} \right) + \log \frac{1 - z_{d+1}}{1 - \hat{z}_{d+1}} \right] \\ &= d \left[ \frac{z_{d+1} - \hat{z}_{d+1}}{(1 - z_{d+1})(1 - \hat{z}_{d+1})} - \frac{z_{d+1} - \hat{z}_{d+1}}{1 - \hat{z}_{d+1}} \right] + \mathcal{O}_{\mathbb{P}}(dh^2) + \mathcal{O}_{\mathbb{P}}(d^3 h^4) z_{d+1}^2 \\ &= d \frac{(z_{d+1} - \hat{z}_{d+1}) z_{d+1}}{(1 - z_{d+1})(1 - \hat{z}_{d+1})} + \mathcal{O}_{\mathbb{P}}(dh^2) + \mathcal{O}_{\mathbb{P}}(d^3 h^4) z_{d+1}^2 \\ &= d \frac{[\mathcal{O}_{\mathbb{P}}(h) + \mathcal{O}_{\mathbb{P}}(h^2 d) z_{d+1}] [\mathcal{O}_{\mathbb{P}}(d^{-1/2}) + \mathcal{O}(d^{\alpha-1})]}{1 - \mathcal{O}_{\mathbb{P}}(d^{-1/2}) - \mathcal{O}(d^{\alpha-1}) - \mathcal{O}_{\mathbb{P}}(h)} + \mathcal{O}_{\mathbb{P}}(dh^2) + \mathcal{O}_{\mathbb{P}}(d^3 h^4) z_{d+1}^2 = o_{\mathbb{P}}(1), \end{aligned}$$

where the last equality holds because (i)  $dh^2 = o(1)$  as  $h = o(d^{-1/2})$ ; (ii) since  $h = o(d^{-\alpha})$ , we can get  $d^3 h^4 z_{d+1}^2$  equals to  $o_{\mathbb{P}}(d^{2\alpha-1})$  if  $\alpha \leq 1/2$  and equals to  $o_{\mathbb{P}}(d^{1-2\alpha})$  if  $1/2 < \alpha \leq 1$ . In both cases,  $\mathcal{O}_{\mathbb{P}}(d^4 h^4) z_{d+1}^2 = o_{\mathbb{P}}(1)$ ; (iii) similarly, we can verify in both cases  $\alpha \leq 1/2$  and  $1/2 < \alpha \leq 1$  we have  $d(h + h^2 dz_{d+1})(d^{-1/2} + d^{\alpha-1}) = o_{\mathbb{P}}(1)$ ; (iv) the term  $\mathcal{O}(d^{\alpha-1})$  is  $\frac{\sum_i \mu_i^2 - \sum_i (1-\lambda_i)}{2d + \sum_i \mu_i^2 - \sum_i (1-\lambda_i)}$  which is bounded away from 1.

S8.8.2. *Proof of Eq. (18).* First of all, we have

$$\hat{X}_i = \frac{d^{1/2}}{1 - \hat{z}_{d+1}} \hat{z}_i, \quad X_i = \frac{d^{1/2}}{1 - z_{d+1}} z_i.$$

Then we have

$$\begin{aligned} \sum_i \delta_i \left( \hat{X}_i^2 - X_i^2 \right) &= d \sum_i \delta_i \frac{\hat{z}_i^2 (1 - z_{d+1})^2 - z_i^2 (1 - \hat{z}_{d+1})^2}{(1 - z_{d+1})^2 (1 - \hat{z}_{d+1})^2} \\ &= d \sum_i \delta_i \left[ \frac{\hat{z}_i^2 - z_i^2}{(1 - \hat{z}_{d+1})^2} + z_i^2 \frac{2(\hat{z}_{d+1} - z_{d+1}) - (\hat{z}_{d+1}^2 - z_{d+1}^2)}{(1 - z_{d+1})^2 (1 - \hat{z}_{d+1})^2} \right]. \end{aligned}$$

Using Taylor expansion, we approximate  $z_i$  and  $\hat{z}_i$  using Eq. (7) in Lemma 3.1. First, we can replace  $\hat{z}_i$  by  $z_i$  :

$$(23) \quad \hat{z}_i = a_d^{(i)} z_i + a_d^{(i)} \sqrt{1 - z_i^2} h U_i - \frac{1}{2} a_d^{(i)} z_i h^2 U_i^2 + \mathcal{O}_{\mathbb{P}}(h^3 + dh^4 z_i),$$

where  $a_d^{(i)} = \frac{1}{\sqrt{1+h^2 U_{i,\perp}^2}} = \frac{1}{\sqrt{1+h^2(d-1)}} (1 + \mathcal{O}_{\mathbb{P}}(h^2 d^{1/2}) + \mathcal{O}_{\mathbb{P}}(d^{-1/2}))$  and  $U_i \sim \mathcal{N}(0, 1)$ . Then for given  $z$ , we can get

$$\begin{aligned} \hat{z}_i^2 - z_i^2 &= ((a_d^{(i)})^2 - 1) z_i^2 + 2(a_d^{(i)})^2 z_i \sqrt{1 - z_i^2} h U_i \\ &\quad + (a_d^{(i)})^2 (1 - 2z_i^2) h^2 U_i^2 + \mathcal{O}_{\mathbb{P}}(z_i h^3 + z_i^2 h^4 + h^6 + d^2 h^8 z_i^2). \end{aligned}$$

Then, we can replace  $z_i$  using  $z_{d+1}$  and  $X_i$  since

$$(24) \quad z_i = (1 - z_{d+1}) d^{-1/2} X_i.$$

By Eq. (20), for  $z_{d+1}$  we have

$$(25) \quad z_{d+1} = \frac{\sum_i \mu_i^2 - \sum_i (1 - \lambda_i)}{2d + \sum_i \mu_i^2 - \sum_i (1 - \lambda_i)} + \mathcal{O}_{\mathbb{P}}(d^{-1/2}).$$

Then, by Eq. (22) we have

$$(26) \quad \sum_i \delta_i \left( \hat{X}_i^2 - X_i^2 \right) = \mathcal{O}_{\mathbb{P}}(dh) \left( \frac{1}{d} \sum_i \delta_i X_i^2 \right) + d \sum_i \delta_i \frac{\hat{z}_i^2 - z_i^2}{(1 - \hat{z}_{d+1})^2}.$$

We first focus on the term  $d \sum_i \delta_i (\hat{z}_i^2 - z_i^2)$ . For given  $X$ , using Eq. (23), replacing  $z_i$  using Eq. (24), simplifying  $z_{d+1}$  with Eq. (25), we have

$$\begin{aligned} d \sum_i \delta_i (\hat{z}_i^2 - z_i^2) &= d \sum_i \delta_i [((a_d^{(i)})^2 - 1) X_i^2 / d + d^{-1/2} h X_i U_i \\ &\quad + (a_d^{(i)})^2 (1 - 2z_i^2) h^2 U_i^2 + \mathcal{O}_{\mathbb{P}}(z_i h^3 + h^4 + d^2 h^8)] \end{aligned}$$

Recall that  $U_i \sim \mathcal{N}(0, 1)$  but for  $i \neq j$ ,  $U_i$  and  $U_j$  are in general not independent since they are obtained by marginalizing the same multivariate Gaussian to two non-orthogonal directions. However, by basic geometry on the angle between the two marginalization directions corresponding to  $U_i$  and  $U_j$ , respectively, we know  $\mathbb{E}[U_i U_j | X] = \mathcal{O}(z_i z_j) = \mathcal{O}(d^{-1} X_i X_j) + \mathcal{O}(d^{-3/2})$ . That is, although  $U_i$  and  $U_j$  are independent if and only if either  $z_i = 0$  or  $z_j = 0$ , their correlation is quite small. Note that we only consider  $h = o(d^{-1/2})$  then  $a_d^{(i)} = a_d(1 + \mathcal{O}_{\mathbb{P}}(h^2 d^{1/2})) = a_d(1 + o_{\mathbb{P}}(d^{-1/2}))$  where  $a_d := \frac{1}{\sqrt{1+h^2(d-1)}}$ , we keep the leading terms of  $d \sum_i \delta_i(\hat{z}_i^2 - z_i^2)$ :

$$\begin{aligned}
d \sum_i \delta_i(\hat{z}_i^2 - z_i^2) &= d \sum_i \delta_i(a_d^2 - 1)X_i^2/d + d \sum_i \delta_i(d^{-1/2}h)X_i U_i + \mathcal{O}(dh^2) \sum_i \delta_i U_i^2 + o_{\mathbb{P}}(1) \\
&= d(a_d^2 - 1) \left( \frac{1}{d} \sum_i \delta_i X_i^2 \right) + \mathcal{O}(dh) \left( \frac{1}{d^{1/2}} \sum_i \delta_i X_i U_i \right) + \mathcal{O}(d^2 h^2) \left( \frac{1}{d} \sum_i \delta_i U_i^2 \right) + o_{\mathbb{P}}(1) \\
&= d(a_d^2 - 1) \left( \frac{1}{d} \sum_i \delta_i X_i^2 \right) + \mathcal{O}(dh) \mathcal{O}_{\mathbb{P}} \left( \sqrt{\frac{2}{d^2} \sum_{i \neq j} \delta_i \delta_j X_i^2 X_j^2 + \frac{1}{d} \sum_i \delta_i^2 X_i^2} \right) \\
&\quad + \mathcal{O}(d^2 h^2) \mathcal{O}_{\mathbb{P}} \left( \frac{1}{d} \sum_i \delta_i \right) + o_{\mathbb{P}}(1) \\
&= d(a_d^2 - 1) \left( \frac{1}{d} \sum_i \delta_i X_i^2 \right) + \mathcal{O}(dh) \left[ \mathcal{O}_{\mathbb{P}} \left( \sqrt{\frac{1}{d} \sum_i \delta_i^2 X_i^2} \right) + \mathcal{O}_{\mathbb{P}} \left( \frac{1}{d} \sum_i \delta_i X_i^2 \right) \right] \\
&\quad + \mathcal{O}(d^2 h^2) \mathcal{O}_{\mathbb{P}} \left( \frac{1}{d} \sum_i \delta_i \right) + o_{\mathbb{P}}(1).
\end{aligned}$$

Substituting to Eq. (26), using  $1 - a_d^2 = \mathcal{O}(h^2 d)$  and  $h = o(d^{-1/2})$ , we have

$$(27) \quad \sum_i \delta_i (\hat{X}_i^2 - X_i^2) = \mathcal{O}_{\mathbb{P}}(d^2 h^2 + dh) \left( \frac{1}{d} \sum_i \delta_i X_i^2 + \frac{1}{d} \sum_i \delta_i \right) + \mathcal{O}_{\mathbb{P}}(dh) \sqrt{\frac{1}{d} \sum_i \delta_i^2 X_i^2} + o_{\mathbb{P}}(1).$$

Now note that both  $\frac{1}{d} \sum_i \delta_i X_i^2$  and  $\frac{1}{d} \sum_i \delta_i^2 X_i^2$  concentrate to their means. Using  $X_i \sim \mathcal{N}(\mu_i, \lambda_i)$ , we have

$$\frac{1}{d} \sum_i \delta_i X_i^2 = \frac{1}{d} \sum_i \delta_i (\mu_i^2 + \lambda_i) + \mathcal{O}_{\mathbb{P}}(d^{-1/2}), \quad \frac{1}{d} \sum_i \delta_i^2 X_i^2 = \frac{1}{d} \sum_i \delta_i^2 (\mu_i^2 + \lambda_i) + \mathcal{O}_{\mathbb{P}}(d^{-1/2})$$

Therefore, substituting to Eq. (27), in order to guarantee  $\sum_i \delta_i (\hat{X}_i^2 - X_i^2) \xrightarrow{\mathbb{P}} 0$ , a sufficient condition for  $h$  in terms of  $\{\mu_i\}$  and  $\{\lambda_i\}$  is

$$(d^2 h^2 + dh) \cdot \left( \frac{1}{d} \sum_i \delta_i (\lambda_i + \mu_i^2) + \frac{1}{d} \sum_i \delta_i \right) + (dh) \cdot \sqrt{\frac{1}{d} \sum_i \delta_i^2 (\lambda_i + \mu_i^2)} = o(1).$$



Therefore, it is sufficient to assume

$$(28) \quad h = o \left( d^{-1} / \sqrt{\max \left\{ \left| \frac{1}{d} \sum_i \frac{\delta_i}{1 + \delta_i} \right|, \left| \frac{1}{d} \sum_i \delta_i \mu_i^2 \right|, \left| \frac{1}{d} \sum_i \frac{\delta_i^2}{1 + \delta_i} \right|, \left| \frac{1}{d} \sum_i \delta_i^2 \mu_i^2 \right|, \left| \frac{1}{d} \sum_i \delta_i \right| \right\} \wedge d^{-1/2}} \right).$$

S8.8.3. *Proof of Eq. (19).* We use a similar approach as the previous part. Using Taylor expansion and Eq. (23) to replace  $\hat{z}_i$ , we have

$$\begin{aligned} \hat{X}_i - X_i &= \frac{d^{1/2}}{1 - \hat{z}_{d+1}} \hat{z}_i - \frac{d^{1/2}}{1 - z_{d+1}} z_i \\ &= \frac{d^{1/2}}{1 - z_{d+1}} \left( a_d^{(i)} z_i + a_d^{(i)} \sqrt{1 - z_i^2} h U_i + \mathcal{O}_{\mathbb{P}}(h^2) \right) (1 + \mathcal{O}(\hat{z}_{d+1} - z_{d+1}) + \mathcal{O}_{\mathbb{P}}(h^2)) - \frac{d^{1/2}}{1 - z_{d+1}} z_i. \end{aligned}$$

Substituting to Eq. (19), and using Eq. (24) to replace  $z_i$  by  $X_i$ ,  $a_d$  to replace  $a_d^{(i)}$ , keeping the leading terms, we can get

$$(29) \quad \sum_i \mu_i (1 + \delta_i) (\hat{X}_i - X_i) = (a_d - 1) \sum_i [\mu_i (1 + \delta_i) X_i] + a_d h \sum_i \left[ \sqrt{d - X_i^2} \mu_i (1 + \delta_i) U_i \right] + o_{\mathbb{P}}(1).$$

Now we first focus on the term  $(a_d - 1) \sum_i [\mu_i (1 + \delta_i) X_i]$  in Eq. (29). Using the fact that  $a_d - 1 = \mathcal{O}(h^2 d)$  and  $X_i \sim \mathcal{N}(\mu_i, \lambda_i)$  are independent random variables, we have

$$(a_d - 1) \sum_i [\mu_i (1 + \delta_i) X_i] = \mathcal{O}(h^2 d^2) \left[ \frac{1}{d} \sum_i \mu_i^2 (1 + \delta_i) + \frac{1}{d} \mathcal{O}_{\mathbb{P}} \left( \sqrt{\sum_i \mu_i^2 (1 + \delta_i)^2 \lambda_i} \right) \right].$$

Therefore,  $(a_d - 1) \sum_i [\mu_i (1 + \delta_i) X_i] \xrightarrow{\mathbb{P}} 0$  if  $(h^2 d^2) \cdot \frac{1}{d} \sum_i \mu_i^2 (1 + \delta_i) = o(1)$ . Next, we focus on the term  $a_d h \sum_i \left[ \sqrt{d - X_i^2} \mu_i (1 + \delta_i) U_i \right]$  in Eq. (29). Using  $U_i \sim \mathcal{N}(0, 1)$  conditional on  $X$  and variance decomposition formula, we have

$$\begin{aligned} &\text{var} \left[ a_d h \sum_i \left( \sqrt{d - X_i^2} \mu_i (1 + \delta_i) U_i \right) \right] \\ &= \mathbb{E} \left[ a_d^2 h^2 \sum_i (d - X_i^2) \mu_i^2 (1 + \delta_i)^2 \right] + \mathcal{O}(d^{-1}) \left( \sum_i a_d h \sqrt{d} \mu_i (1 + \delta_i) \right)^2 \\ &= \mathcal{O}(d^2 h^2) \left[ \frac{1}{d} \sum_i \left( 1 - \frac{\mu_i^2 + \lambda_i}{d} \right) \mu_i^2 (1 + \delta_i)^2 + \left( \frac{1}{d} \sum_i \mu_i (1 + \delta_i) \right)^2 \right]. \end{aligned}$$

Therefore,  $a_d h \sum_i \left( \sqrt{d - X_i^2} \mu_i (1 + \delta_i) U_i \right) \xrightarrow{\mathbb{P}} 0$  if both  $(d^2 h^2) \cdot \frac{1}{d} \sum_i \mu_i^2 (1 + \delta_i)^2 = o(1)$  and  $(dh) \left| \frac{1}{d} \sum_i \mu_i (1 + \delta_i) \right| = o(1)$ . Overall, in order to make  $\sum_i \mu_i (1 + \delta_i) (\hat{X}_i - X_i)$  to converge to 0 in probability, we need Eq. (29) to be  $o_{\mathbb{P}}(1)$ . Then, a sufficient condition for  $h$  in terms of  $\{\mu_i\}$  and  $\{\lambda_i\}$  is

$$d^2 h^2 \max \left\{ \frac{1}{d} \sum_i \mu_i^2 (1 + \delta_i), \frac{1}{d} \sum_i \mu_i^2 (1 + \delta_i)^2 \right\} \rightarrow 0, \quad dh \left| \frac{1}{d} \sum_i \mu_i (1 + \delta_i) \right| \rightarrow 0.$$

Therefore, it is sufficient if

$$(30) \quad h = o\left(d^{-1}/\sqrt{\max\left\{\left|\frac{1}{d}\sum_i \frac{\delta_i}{1+\delta_i}\right|, \left|\frac{1}{d}\sum_i \delta_i \mu_i^2\right|, \left|\frac{1}{d}\sum_i \delta_i^2 \mu_i^2\right|, \left|\frac{1}{d}\sum_i \mu_i^2\right|\right\} \wedge d^{-1/2}}\right).$$

Overall, combing Eq. (28) and Eq. (30), since we only need the order of  $h$ , Eq. (28) can be relaxed to

$$h = o\left(d^{-1}/\sqrt{\frac{1}{d}\sum_i |1 - \lambda_i| \wedge d^{-1/2}}\right)$$

and Eq. (30) can be relaxed to

$$h = o\left(d^{-1}/\sqrt{\max\left\{\frac{1}{d}\sum_i |1 - \lambda_i|, \frac{1}{d}\sum_i \mu_i^2\right\} \wedge d^{-1/2}}\right).$$

Combing the above two with  $h = o(d^{-\alpha})$  completes the proof.  $\square$

### S8.9. Proof of Lemma 5.1.

PROOF. Note that we assumed  $h = \mathcal{O}(d^{-1})$ . We use  $z$  and  $\hat{z}$  to denote the corresponding coordinates of  $X$  and  $\hat{X}$  on the unit sphere. We first “upgrade” Eq. (7) in Lemma 3.1 to a convergence in  $L_1$  statement. Define

$$\bar{z}_i := \frac{1}{\sqrt{1+h^2U_\perp^2}} \left[ \left(1 - \frac{1}{2}h^2U^2\right) z_i - \sqrt{1-z_i^2}hU \right].$$

and the facts that  $\hat{z}_i$  is bounded and  $\bar{z}_i$  has a finite exponential moment, applying Cauchy-Schwarz inequality and Markov inequality, we have

$$\begin{aligned} \mathbb{E}[|\hat{z}_i - \bar{z}_i|] &\leq \mathbb{E}[|\hat{z}_i - \bar{z}_i \mathbf{1}_{|\bar{z}_i| \leq C \log(d)}|] + \mathbb{E}[|\bar{z}_i \mathbf{1}_{|\bar{z}_i| > C \log(d)}|] \\ &= C \log(d) \mathcal{O}(h^3 + dh^4 z_i) + \sqrt{\mathbb{E}[\bar{z}_i^2]} \sqrt{\mathbb{P}(\exp(|\bar{z}_i|) > d^C)} \\ &= \mathcal{O}(\log(d)h^3) + \mathcal{O}(d^{-C/2}) = \mathcal{O}(\log(d)h^3), \end{aligned}$$

where  $C$  is chosen as a large enough constant so that  $d^{-C/2} = o(\log(d)h^3)$ . Therefore, Eq. (7) in Lemma 3.1 can be rewritten as a convergence in  $L_1$  statement:

$$(31) \quad \mathbb{E}_{\hat{z}_i|z_i} \left[ \left| \hat{z}_i - \left( \frac{1}{\sqrt{1+h^2U_\perp^2}} \left[ \left(1 - \frac{1}{2}h^2U^2\right) z_i - \sqrt{1-z_i^2}hU \right] \right) \right| \right] = \mathcal{O}(\log(d)h^3).$$

Next, we want to rule out some subsets of “bad states”, including those states with the “latitude”  $z_{d+1}$  too close to 1. That is, we define a sequence of “typical sets”  $\{F_d\}$  with  $\pi(F_d) \rightarrow 1$  such that when  $x \in F_d$ , the corresponding  $z_{d+1}$  is bounded away from 1. We

define

$$\begin{aligned}
F_d := & \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d} \sum_{i=1}^d [(\log f(x_i))']^2 - \mathbb{E}_f [((\log f)')^2] \right| < d^{-1/2} \log(d) \right\} \\
& \cap \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d} \sum_{i=1}^d [(\log f(x_i))''] - \mathbb{E}_f [((\log f)'')] \right| < d^{-1/2} \log(d) \right\} \\
& \cap \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d} \sum_{i=1}^d x_i (\log f(x_i))' - \mathbb{E}_f [X(\log f)'] \right| < d^{-1/2} \log(d) \right\} \\
& \cap \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d} \sum_{i=1}^d x_i^2 - \mathbb{E}_f (X^2) \right| < d^{-1/2} \log(d) \right\}
\end{aligned}$$

Note that by the assumption  $\lim_{x \rightarrow \pm\infty} x f'(x) = 0$ , we have

$$\mathbb{E}_f [X(\log f)'] = \int \left[ \frac{x}{f(x)} f'(x) \right] f(x) dx = \int x f'(x) dx = 0 - \int f(x) dx = -1.$$

Therefore our assumptions on  $\pi$  imply

$$\pi(F_d^c) = \mathbb{P}_\pi(X \notin F_d) = \mathcal{O} \left( \left( \frac{d^{1/2}}{\log(d)d^{1/2}} \right)^4 + \left( \frac{d^{1/2}}{\log(d)d^{1/2}} \right)^3 \right) = o(1).$$

Using the definition of  $\{F_d\}$ , we can assume  $X \in F_d$  and only consider the expectation over  $\hat{X}$  given  $X$ . By the definition of  $\{F_d\}$ , we know  $z_{d+1}$  is bounded away from 1.

Next, we want to rule out the subset of “bad proposals” where the proposal “latitude”  $\hat{z}_{d+1}$  is too close to 1. Under stationarity, we know  $R^2 \frac{1+z_{d+1}}{1-z_{d+1}} = \lambda d \frac{1+z_{d+1}}{1-z_{d+1}} = d + \mathcal{O}_{\mathbb{P}}(d^{1/2})$ , which implies  $z_{d+1} = \frac{1-\lambda}{1+\lambda} + \mathcal{O}_{\mathbb{P}}(d^{-1/2})$  and  $\frac{1}{1-z_{d+1}} = \frac{1+\lambda}{2\lambda} + \mathcal{O}_{\mathbb{P}}(d^{-1/2})$ . Furthermore, for given  $X \in F_d$ , since we know

$$\hat{z}_{d+1} = z_{d+1} + \mathcal{O}_{\mathbb{P}}(h),$$

uniformly on  $X \in F_d$ , for all large enough  $d$ , we have

$$z_{d+1} + h \log(d) = \frac{1-\lambda}{1+\lambda} + \mathcal{O}(\log(d)d^{-1/2}) + h \log(d) < 1 - \epsilon,$$

where  $\frac{\lambda}{1+\lambda} < \epsilon < 1$  is a fixed constant. Furthermore, we have

$$\sup_{X \in F_d} \mathbb{P}(\hat{z}_{d+1} > 1 - \epsilon) \leq \sup_{X \in F_d} \mathbb{P}(\hat{z}_{d+1} > z_{d+1} + h \log(d)) = o(d^{-1/2}).$$

Therefore, we only need consider the cases such that  $\hat{z}_{d+1}$  is no smaller than  $1 - \epsilon$  since

$$\begin{aligned}
& \sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(X)(R^2 + \|X\|^2)^d} - 1 \wedge \exp(W_{\hat{X}|X}) \right| \right] \\
& \leq \sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| 1 \wedge \left( \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(X)(R^2 + \|X\|^2)^d} \mathbf{1}_{\{\hat{z}_{d+1} \leq 1-\epsilon\}} \right) - 1 \wedge \exp(W_{\hat{X}|X}) \right| \right] \\
& \quad + \sup_{X \in F_d} \mathbb{P}(\hat{z}_{d+1} > 1 - \epsilon)
\end{aligned}$$

$$\leq \sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| \log \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(X)(R^2 + \|X\|^2)^d} \mathbf{1}_{\{\hat{z}_{d+1} \leq 1-\epsilon\}} - W_{\hat{X}|X} \right| \right] + o(d^{-1/2}),$$

where the last step is because  $1 \wedge e^x$  is a 1-Lipschitz function.

In the rest of the proof, all we will use for approximating  $\hat{z}_{d+1}$  is Eq. (31). Since  $\frac{\lambda}{1+\lambda} < \epsilon < 1$ , Eq. (31) is also true if we replace  $\hat{z}_{d+1}$  by  $\hat{z}_{d+1} \mathbf{1}_{\{\hat{z}_{d+1} \leq 1-\epsilon\}}$ . That is, for truncated  $\hat{z}_{d+1}$ :

$$(32) \quad \sup_{X \in F_d} \mathbb{E}_{\hat{z}_{d+1}|z_{d+1}} \left[ \left| \hat{z}_{d+1} \mathbf{1}_{\{\hat{z}_{d+1} \leq 1-\epsilon\}} - \bar{z}_{d+1} \right| \right] = \mathcal{O}(\log(d)h^3).$$

Therefore, in the rest of the proof, we simply work on the cases such that  $\hat{z}_{d+1} \leq 1-\epsilon$  and sometimes omit the term  $\mathbf{1}_{\{\hat{z}_{d+1} \leq 1-\epsilon\}}$  for simplicity. The goal is then to consider  $X \in F_d$  and establish a Gaussian approximation result for

$$\left( \sum_{i=1}^d [\log f(\hat{X}_i) - \log f(X_i)] + d[\log(1-z_{d+1}) - \log(1-\hat{z}_{d+1})] \right) \mathbf{1}_{\{\hat{z}_{d+1} < 1-\epsilon\}}.$$

Note that by Eq. (31), we know  $\mathbb{E}[\|\hat{z}_i - z_i\|^3] = \mathcal{O}(h^3)$ ,  $\mathbb{E}[\|\hat{z}_i - z_i\|(\hat{z}_{d+1} - z_{d+1})^2] = \mathcal{O}(h^3)$ ,  $\mathbb{E}[\|\hat{z}_{d+1} - z_{d+1}\|(\hat{z}_i - z_i)^2] = \mathcal{O}(h^3)$ . Furthermore, we only need to consider  $\frac{1}{1-\hat{z}_{d+1}}$  is bounded as  $\hat{z}_{d+1}$  is bounded away from 1. Then, using

$$(33) \quad \begin{aligned} \hat{X}_i - X_i &= \frac{R}{1-\hat{z}_{d+1}} \hat{z}_i - \frac{R}{1-z_{d+1}} z_i \\ &= \frac{R}{(1-z_{d+1})(1-\hat{z}_{d+1})} (\hat{z}_{d+1} - z_{d+1}) z_i + \left( \frac{R(\hat{z}_{d+1} - z_{d+1})}{(1-z_{d+1})(1-\hat{z}_{d+1})} + \frac{R}{1-z_{d+1}} \right) (\hat{z}_i - z_i), \end{aligned}$$

together with Taylor expansion and mean value theorem, uniformly on  $\forall X \in F_d$ , we have

$$\mathbb{E}_{\hat{X}|X} \left[ \left| \sum_{i=1}^d \log \frac{f(\hat{X}_i)}{f(X_i)} - d \log \frac{1-\hat{z}_{d+1}}{1-z_{d+1}} - (\text{Term I}) - (\text{Term II}) \right| \cdot \mathbf{1}_{\{\hat{z}_{d+1} \leq 1-\epsilon\}} \right] = \mathcal{O}(dh^3),$$

where

$$\begin{aligned} \text{Term I} &:= \sum_i \left[ (\log f(X_i))'(\hat{X}_i - X_i) \right] + \frac{d}{1-z_{d+1}} (\hat{z}_{d+1} - z_{d+1}), \\ \text{Term II} &:= \frac{1}{2} \left\{ \sum_i \left[ (\log f(X_i))''(\hat{X}_i - X_i)^2 \right] + \frac{d}{(1-z_{d+1})^2} (\hat{z}_{d+1} - z_{d+1})^2 \right\}. \end{aligned}$$

We will see that, similar to the CLT arguments used for analyzing RWM [Roberts et al., 1997, Yang et al., 2020] for general targets, the mean of the Gaussian approximation result is determined by the sum of the mean of Term I and the mean of Term II whereas the variance of the Gaussian approximation result is determined by the variance of Term I only. In the rest of the proof, we first approximate Terms I and II separately, then combine them in the end.

**S8.9.1. Term I.** Using Eq. (33) and Eq. (31), with the definition of  $F_d$ , we can get the following approximation of  $\hat{X}_i - X_i$ :

$$\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| \hat{X}_i - X_i - \frac{1+\lambda}{2\lambda} (\hat{z}_{d+1} - z_{d+1}) X_i + \frac{R}{1-z_{d+1}} (\hat{z}_i - z_i) \right| \cdot \mathbf{1}_{\{\hat{z}_{d+1} \leq 1-\epsilon\}} \right] = \mathcal{O}(d^{1/2}h^2).$$

Next, let  $\hat{z} - z = (\hat{z}_1 - z_1, \dots, \hat{z}_{d+1} - z_{d+1})^T$ , we can approximate the ‘‘Term I’’ by

$$\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| (\text{Term I}) - \frac{R}{1 - z_{d+1}} \cdot (\text{Inner Product Term}) \right| \cdot \mathbf{1}_{\{\hat{z}_{d+1} \leq 1 - \epsilon\}} \right] = o(d^{-1/2}),$$

where the ‘‘Inner Product Term’’ represents

$$\left( (\log f(X_1))', \dots, (\log f(X_d))', \sum_{i=1}^d \left( \frac{1 + \lambda}{2\lambda} (\log f(X_i))' z_i + \frac{1}{R} \right) \right) \cdot (\hat{z} - z).$$

The key observation is that: we can see the inner product term as the ‘‘projection’’ of  $\hat{z} - z$  to the particular ‘‘direction’’ defined by

$$\tilde{v} := \left( (\log f(X_1))', \dots, (\log f(X_d))', \sum_{i=1}^d \left( \frac{1 + \lambda}{2\lambda} (\log f(X_i))' z_i + \frac{1}{R} \right) \right).$$

This observation allows us to just study the distribution of Term I through projections of the current locations and the proposal for this particular direction. Note that the approximation results of Eq. (31) and its proof hold not only for all coordinates but also for projections to any direction due to the symmetry of the sphere.

The ‘‘projection’’ of  $z$  onto  $\tilde{v}$  is

$$\begin{aligned} \tilde{z} &:= \frac{\tilde{v} \cdot z}{\|\tilde{v}\|} = \frac{\left[ \sum_{i=1}^d ((\log f(X_i))' z_i) \right] + z_{d+1} \left[ \sum_{i=1}^d \left( \frac{1 + \lambda}{2\lambda} (\log f(X_i))' z_i + \frac{1}{R} \right) \right]}{\sqrt{\sum_{i=1}^d ((\log f(X_i))')^2 + \left[ \sum_{i=1}^d \left( \frac{1 + \lambda}{2\lambda} (\log f(X_i))' z_i + \frac{1}{R} \right) \right]^2}} \\ &= \frac{\left( \frac{2\lambda}{1 + \lambda} + z_{d+1} \right) \left[ \sum_{i=1}^d \left( \frac{1 + \lambda}{2\lambda} (\log f(X_i))' z_i + \frac{1}{R} \right) \right] - \frac{2\lambda}{1 + \lambda} \frac{d}{R}}{\sqrt{\sum_{i=1}^d ((\log f(X_i))')^2 + \left[ \sum_{i=1}^d \left( \frac{1 + \lambda}{2\lambda} (\log f(X_i))' z_i + \frac{1}{R} \right) \right]^2}} \in [-1, 1] \end{aligned}$$

Then we can write

$$\begin{aligned} \|\tilde{v}\|^2 - |\tilde{v} \cdot z|^2 &= \sum_{i=1}^d ((\log f(X_i))')^2 + \left[ \sum_{i=1}^d \left( \frac{1 + \lambda}{2\lambda} (\log f(X_i))' z_i + \frac{1}{R} \right) \right]^2 \\ &\quad - \left( \frac{2\lambda}{1 + \lambda} + z_{d+1} \right)^2 \left[ \sum_{i=1}^d \left( \frac{1 + \lambda}{2\lambda} (\log f(X_i))' z_i + \frac{1}{R} \right) \right]^2 - \left( \frac{2\lambda}{1 + \lambda} \right)^2 \frac{d^2}{R^2} \\ &\quad + \frac{2\lambda}{1 + \lambda} \frac{2d \left( \frac{2\lambda}{1 + \lambda} + z_{d+1} \right)}{R} \left[ \sum_{i=1}^d \left( \frac{1 + \lambda}{2\lambda} (\log f(X_i))' z_i + \frac{1}{R} \right) \right]. \end{aligned}$$

Note that by the definition of  $F_d$ , we know  $\sum_{i=1}^d (\log f(X_i))' X_i + d = o(d)$  and  $\frac{2\lambda}{1 + \lambda} + z_{d+1} = 1 + \mathcal{O}(d^{-1/3})$ . Then, we have

$$\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| (\|\tilde{v}\|^2 - |\tilde{v} \cdot z|^2) - (\text{Term III}) \right| \cdot \mathbf{1}_{\{\hat{z}_{d+1} \leq 1 - \epsilon\}} \right] = o(d^{1/2}),$$

where the ‘‘Term III’’ is defined by

$$\text{Term III} := \sum_{i=1}^d ((\log f(X_i))')^2 + \frac{2\lambda}{1 + \lambda} \frac{2d}{R} \sum_{i=1}^d \left( \frac{1 + \lambda}{2\lambda} (\log f(X_i))' z_i \right) + \left[ \frac{4\lambda}{1 + \lambda} - \left( \frac{2\lambda}{1 + \lambda} \right)^2 \right] \frac{d^2}{R^2}.$$

Let  $\tilde{z}'$  be the ‘‘projection’’ of  $\hat{z}$  onto  $\tilde{v}$ . Then, the inner product term can be written as  $\tilde{v} \cdot (\hat{z} - z) = \|\tilde{v}\|(\tilde{z}' - \tilde{z})$ . Next, we can approximate the ‘‘projection’’ of  $\hat{z}$  onto  $\tilde{v}$  using the equivalent result of Eq. (31) for  $\tilde{z}$

$$\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| \tilde{z}' - \tilde{a}_d \left( \tilde{z} - \sqrt{1 - \tilde{z}^2} h \tilde{U} \right) \right| \right] = \mathcal{O}(h^2),$$

where  $\tilde{U}$  is a standard Gaussian and  $\tilde{a}_d := \frac{1}{\sqrt{1+h^2\tilde{U}_1^2}} = a_d(1 + \mathcal{O}_{\mathbb{P}}(h^2d^{1/2})) = a_d(1 + \mathcal{O}_{\mathbb{P}}(d^{-3/2}))$  where  $a_d := \frac{1}{\sqrt{1+(d-1)h^2}}$ . Note that we can also ‘‘upgrade’’ the relation between  $\tilde{a}_d$  and  $a_d$  to  $L_1$  statement:  $\mathbb{E}[|\tilde{a}_d - a_d|] = \mathcal{O}(d^{-3/2} \log(d))$ . Then by triangle inequality

$$\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| \tilde{z}' - a_d \left( \tilde{z} - \sqrt{1 - \tilde{z}^2} h \tilde{U} \right) \right| \right] = \mathcal{O}(h^2 + d^{-3/2} \log(d)).$$

Therefore, we have

$$\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| \|\tilde{v}\|(\tilde{z}' - \tilde{z}) - W_1 \right| \cdot \mathbf{1}_{\{\hat{z}_{d+1} \leq 1-\epsilon\}} \right] = \mathcal{O}(d^{-1/2} \log(d)),$$

where  $W_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  with  $\mu_1 := (a_d - 1) \|\tilde{v}\| \tilde{z}$  and  $\sigma_1^2 := (1 - \tilde{z}^2) \|\tilde{v}\|^2 a_d^2 h^2$ . Substituting  $1 - \tilde{z}^2 = \frac{\|\tilde{v}\|^2 - \|\tilde{v} \cdot z\|^2}{\|\tilde{v}\|^2}$ , using the definition of  $F_d$ , we can compute the mean and variance of  $W_1$  as

$$\begin{aligned} \mu_1 &= (a_d - 1) \frac{1 + \lambda}{2\lambda} \left[ \sum_i (\log f(X_i))' X_i + \frac{1 - \lambda}{1 + \lambda} d \right] + \mathcal{O}(d^{-1/2} \log(d)) \\ &= (a_d - 1) \frac{1 + \lambda}{2\lambda} \left[ d \mathbb{E}[(\log f(X_1))' X_1] + \frac{1 - \lambda}{1 + \lambda} d \right] + \mathcal{O}(d^{-1/2} \log(d)), \\ \sigma_1^2 &= \frac{(1 + \lambda)^2}{4\lambda} da_d^2 h^2 \left[ \sum_i ((\log f(X_i))')^2 + \frac{4}{1 + \lambda} \sum_i (\log f(X_i))' X_i + \frac{4d}{(1 + \lambda)^2} \right] \\ &\quad + \mathcal{O}(d^{-1/2} \log(d)) \\ &= \frac{(1 + \lambda)^2}{4\lambda} da_d^2 h^2 \left[ d \mathbb{E} \left[ ((\log f(X_1))')^2 \right] + \frac{4d}{1 + \lambda} \mathbb{E} [(\log f(X_1))' X_1] + \frac{4d}{(1 + \lambda)^2} \right] \\ &\quad + \mathcal{O}(d^{-1/2} \log(d)). \end{aligned}$$

**S8.9.2. Term II.** Next, we approximate Term II in the Taylor expansion. Clearly, the variance of Term II can be ignored since it is of the order  $o(d^{-1})$ . We only need to focus on the expectation. We start with computing

$$\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ (\log f(X_i))'' (\hat{X}_i - X_i)^2 \cdot \mathbf{1}_{\{\hat{z}_{d+1} \leq 1-\epsilon\}} \right].$$

Note that

$$\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| (\hat{X}_i - X_i)^2 - (\text{Term IV}) \right| \cdot \mathbf{1}_{\{\hat{z}_{d+1} \leq 1-\epsilon\}} \right] = o(d^{1/2} h^2),$$

where the ‘‘Term IV’’ is defined as

$$\left( \frac{1 + \lambda}{2\lambda} \right)^2 (\hat{z}_{d+1} - z_{d+1})^2 X_i^2 + \frac{R^2}{(1 - z_{d+1})^2} (\hat{z}_i - z_i)^2 + \frac{1 + \lambda}{2\lambda} \frac{2RX_i}{1 - z_{d+1}} (\hat{z}_{d+1} - z_{d+1})(\hat{z}_i - z_i).$$

Taking expectation over  $\hat{X}$  given  $X$  and using

$$\begin{aligned}\mathbb{E}_{\hat{X}|X}[\hat{z}_i - z_i] &= (a_d - 1)z_i + o(d^{-1}) \\ \mathbb{E}_{\hat{X}|X}[(\hat{z}_i - z_i)^2] &= a_d^2 h^2 + [(a_d - 1)^2 - a_d^2 h^2]z_i^2 + o(d^{-1}),\end{aligned}$$

substituting to the ‘‘Term IV’’, then uniformly over  $X \in F_d$ , we have

$$\begin{aligned}\mathbb{E}_{\hat{X}|X}[(\hat{X}_i - X_i)^2] &= \left(\frac{1+\lambda}{2\lambda}\right)^2 X_i^2 [a_d^2 h^2(1 - z_{d+1}^2) + (1 - a_d)^2 z_{d+1}^2] \\ &\quad + \frac{R^2}{(1 - z_{d+1})^2} [a_d^2 h^2(1 - z_i^2) + (1 - a_d)^2 z_i^2] \\ &\quad + \frac{1+\lambda}{2\lambda} \frac{2RX_i}{1 - z_{d+1}} (1 - a_d)^2 z_i z_{d+1} + o(d^{-1}) + o(d^{1/2}h^2).\end{aligned}$$

This implies that, uniformly over  $X \in F_d$ , we have

$$\begin{aligned}\sum_i \mathbb{E}_{\hat{X}|X}[(\hat{X}_i - X_i)^2] &= \left(\frac{1+\lambda}{2\lambda}\right)^2 \frac{d}{\lambda} [a_d^2 h^2(1 - z_{d+1}^2) + (1 - a_d)^2 z_{d+1}^2] \\ &\quad + \frac{R^2}{(1 - z_{d+1})^2} da_d^2 h^2 - da_d^2 h^2 + (1 - a_d)^2 d + \frac{1+\lambda}{2\lambda} \frac{d}{\lambda} (1 - a_d)^2 z_{d+1} + \mathcal{O}(d^{-1/2}) \\ &= \frac{R^2}{(1 - z_{d+1})^2} da_d^2 h^2 + \mathcal{O}(d^{-1/2}) = (1 - a_d^2) d \frac{(1+\lambda)^2}{4\lambda} + \mathcal{O}(d^{-1/2}).\end{aligned}$$

Therefore, uniformly on  $F_d$ , we have

$$\begin{aligned}&\sum_{i=1}^d \mathbb{E}_{\hat{X}|X} \left[ (\log f(X_i))'' (\hat{X}_i - X_i)^2 \right] + \frac{d}{(1 - z_{d+1})^2} \mathbb{E}_{\hat{X}|X} [(\hat{z}_{d+1} - z_{d+1})^2] \\ &= \frac{R^2}{(1 - z_{d+1})^2} a_d^2 h^2 \left[ \sum_{i=1}^d (\log f(X_i))'' - \left[ 1 - \left(\frac{1+\lambda}{2\lambda}\right)^2 (1 - z_{d+1}^2) \right] \sum_{i=1}^d ((\log f(X_i))'' z_i^2) \right] \\ &\quad + \frac{R^2}{(1 - z_{d+1})^2} (1 - a_d)^2 \left[ 1 + \left(\frac{1+\lambda}{2\lambda}\right)^2 z_{d+1}^2 + \frac{1+\lambda}{\lambda} z_{d+1} \right] \sum_{i=1}^d ((\log f(X_i))'' z_i^2) \\ &\quad + \frac{1 + z_{d+1}}{1 - z_{d+1}} da_d^2(h)h^2 + \frac{d(1 - a_d)^2 z_{d+1}^2}{(1 - z_{d+1})^2} + \mathcal{O}(d^{-1/2}).\end{aligned}$$

Finally, using the definition of  $F_d$ , we can have

$$\begin{aligned}&2 \cdot \text{Term II} + \mathcal{O}(d^{-1/2} \log(d)) \\ &= \frac{(1+\lambda)^2}{4\lambda} da_d^2 h^2 \sum_i \mathbb{E} [(\log f(X_i))''] - a_d^2 h^2 \left(1 - \frac{1}{\lambda}\right) \sum_i \mathbb{E} [(\log f(X_i))'' X_i^2] \\ &\quad + \frac{(1+\lambda)^2}{4\lambda} (1 - a_d)^2 \sum_i \mathbb{E} [(\log f(X_i))'' X_i^2] + \frac{1}{\lambda} da_d^2 h^2 + \frac{(1-\lambda)^2}{4\lambda^2} (1 - a_d)^2 d.\end{aligned}$$

**S8.9.3. Combining Term I and Term II.** Finally, we combine Term I and Term II. Using the assumption  $\lim_{x \rightarrow \pm\infty} x f'(x) = 0$  which implies  $\lim_{x \rightarrow \pm\infty} f'(x) = 0$ , one can show

$$\mathbb{E}_f[(\log f)'' + ((\log f)')^2] = \int \left[ \frac{f(x)f''(x) - (f'(x))^2}{f(x)^2} + \left( \frac{1}{f(x)} f'(x) \right)^2 \right] f(x) dx = \int f''(x) dx = 0,$$

$$\mathbb{E}_f[X(\log f)'] = \int \left[ \frac{x}{f(x)} f'(x) \right] f(x) dx = \int x f'(x) dx = 0 - \int f(x) dx = -1.$$

Using the above results and  $(d-1)a_d^2 h^2 = 1 - a_d^2$ , we have

$$\mu_1 = -(1 - a_d) \frac{1 + \lambda}{2\lambda} d \mathbb{E}_f[X(\log f)'] - (1 - a_d) \frac{1 - \lambda}{2\lambda} d + \mathcal{O}(d^{-1/2} \log(d)),$$

$$\mathbb{E}[\text{Term II}] = (1 - a_d) \frac{(1 + \lambda)^2}{4\lambda} d \mathbb{E}_f[(\log f)''] + \mathcal{O}(d^{-1/2} \log(d)),$$

$$\begin{aligned} \frac{\sigma_1^2}{2} &= (1 - a_d) \frac{(1 + \lambda)^2}{4\lambda} d \mathbb{E}_f[((\log f)')^2] + (1 - a_d) \frac{1 + \lambda}{\lambda} d \mathbb{E}_f[X(\log f)'] \\ &\quad + (1 - a_d) \frac{d}{\lambda} + \mathcal{O}(d^{-1/2} \log(d)). \end{aligned}$$

Therefore, the mean and variance of the Gaussian distribution satisfy

$$\mu = \mu_1 + \mathbb{E}[\text{Term II}] + \mathcal{O}(d^{-1/2} \log(d)) = d(1 - a_d) \frac{(1 + \lambda)^2}{4\lambda} \left\{ \frac{4\lambda}{(1 + \lambda)^2} + \mathbb{E} \left[ \frac{\partial^2 \log \pi}{\partial x^2} \right] \right\},$$

$$\frac{\sigma^2}{2} = \frac{\sigma_1^2}{2} + \mathcal{O}(d^{-1/2} \log(d)) = d(1 - a_d) \frac{(1 + \lambda)^2}{4\lambda} \left\{ \mathbb{E} \left[ \left( \frac{\partial \log \pi}{\partial x} \right)^2 \right] - \frac{4\lambda}{(1 + \lambda)^2} \right\}.$$

Using  $1 - a_d = \frac{\ell^2}{2d} \frac{4\lambda}{(1 + \lambda)^2}$ , we have

$$\mu = \frac{\ell^2}{2} \left\{ \frac{4\lambda}{(1 + \lambda)^2} - \mathbb{E}_f \left[ ((\log f)')^2 \right] \right\}, \quad \sigma^2 = \ell^2 \left\{ \mathbb{E}_f \left[ ((\log f)')^2 \right] - \frac{4\lambda}{(1 + \lambda)^2} \right\}.$$

Therefore, we have shown a Gaussian random variable  $W_{\hat{X}|X} \sim \mathcal{N}(\mu, \sigma^2)$  which satisfies

$$\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left( (\text{Term I} + \text{Term II}) \cdot \mathbf{1}_{\{\hat{z}_{d+1} \leq 1 - \epsilon\}} - W_{\hat{X}|X} \right)^2 \right] = \mathcal{O}(d^{-1/2} \log(d)).$$

Combining everything together, we have shown

$$\begin{aligned} &\sup_{X \in F_d} \mathbb{E}_{\hat{X}|X} \left[ \left| \log \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(X)(R^2 + \|X\|^2)^d} \mathbf{1}_{\{\hat{z}_{d+1} \leq 1 - \epsilon\}} - W_{\hat{X}|X} \right| \right] \\ &= \mathcal{O}(\sqrt{d^{-1/2} \log(d)}) = o(d^{-1/4} \log(d)), \end{aligned}$$

which completes the proof.  $\square$

### S8.10. Proof of Theorem 5.1.

**PROOF.** By the i.i.d. assumption, it suffices to consider the first coordinate as

$$\text{ESJD} = d \cdot \mathbb{E}_{X \sim \pi} \mathbb{E}_{\hat{X}|X} \left[ (\hat{X}_1 - X_1)^2 \left( 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(X)(R^2 + \|X\|^2)^d} \right) \right].$$



We first follow a similar approach as the proof of Lemma 5.1. We define a sequence of “typical sets”  $\{F_d\}$  such that  $\pi(F_d) \rightarrow 1$ . Define the sequence of sets  $\{F_d\}$  by

$$(34) \quad \begin{aligned} F_d := & \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d-1} \sum_{i=2}^d [(\log f(x_i))']^2 - \mathbb{E}_f [((\log f)')^2] \right| < d^{-1/8} \right\} \\ & \cap \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d-1} \sum_{i=2}^d [(\log f(x_i))''] - \mathbb{E}_f [((\log f)'')] \right| < d^{-1/8} \right\} \\ & \cap \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d-1} \sum_{i=2}^d x_i (\log f(x_i))' - \mathbb{E}_f [X (\log f)'] \right| < d^{-1/8} \right\} \\ & \cap \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d} \sum_{i=1}^d x_i^2 - \mathbb{E}_f (X^2) \right| < d^{-1/6} \log(d) \right\} \\ & \cap \left\{ x \in \mathbb{R}^d : |x_1| < d^{1/5} \right\}. \end{aligned}$$

Following the arguments in Section S8.9, there exists a constant  $0 < \epsilon < 1$  such that  $\sup_{X \in F_d} \mathbb{P}(\hat{z}_{d+1} > 1 - \epsilon) = o(1)$ . We define

$$\bar{F}_d := \{x \in \mathbb{R}^d : z_{d+1} \leq 1 - \epsilon\}.$$

We will rule out the “bad” proposals such that  $\hat{X} \in \bar{F}_d^c$ . By Cauchy–Schwarz inequality

$$\begin{aligned} & d \cdot \mathbb{E}_{X \sim \pi} \mathbb{E}_{\hat{X} | X} \left[ (\hat{X}_1 - X_1)^2 \mathbf{1}_{\{X \notin F_d\} \cup \{\hat{X} \notin \bar{F}_d\}} \right] \\ & \leq d \sqrt{\mathbb{E}[(\hat{X}_1 - X_1)^4]} \sqrt{\mathbb{P}(\{X \notin F_d\} \cup \{\hat{X} \notin \bar{F}_d\})} = o(d^2 h^2) = o(1). \end{aligned}$$

Therefore, we have

$$(35) \quad \begin{aligned} \frac{\text{ESJD}}{d} &= \mathbb{E} \left\{ (\hat{X}_1 - X_1)^2 \mathbf{1}_{\{X \in F_d, \hat{X} \in \bar{F}_d\}} \left( 1 \wedge \exp \left[ \log \frac{\pi(\hat{X})}{\pi(X)} + \log \frac{(R^2 + \|\hat{X}\|^2)^d}{(R^2 + \|X\|^2)^d} \right] \right) \right\} + o(d^{-1}) \\ &\leq \sup_{x \in \bar{F}_d} \mathbb{E}_{\hat{X} | x} \left\{ (\hat{X}_1 - x_1)^2 \left( 1 \wedge \exp \left[ \log \frac{\pi(\hat{X})}{\pi(x)} + \log \frac{(R^2 + \|\hat{X}\|^2)^d}{(R^2 + \|x\|^2)^d} \right] \right) \mathbf{1}_{\{\hat{X} \in \bar{F}_d\}} \right\} + o(d^{-1}) \end{aligned}$$

In order to remove the dependence on  $\hat{X}_1$  from  $\|\hat{X}\|^2$ , we construct a coupling  $\tilde{X} \stackrel{d}{=} \hat{X}$  such that only  $\hat{X}_{2:d} = \tilde{X}_{2:d}$  and  $\tilde{X}_1$  is identical distributed as  $\hat{X}_1$ . Furthermore,  $\tilde{X}_1$  is independent with  $\hat{X}_1$  conditional on  $\hat{X}_{2:d}$ :

$$\tilde{X} := (\tilde{X}_1, \dots, \tilde{X}_d), \quad \tilde{X}_{2:d} = \hat{X}_{2:d}, \quad (\tilde{X}_1 \perp\!\!\!\perp \hat{X}_1 \mid \hat{X}_{2:d}), \quad \tilde{X}_1 \stackrel{d}{=} \hat{X}_1.$$

Using a similar argument, we can replace  $\mathbf{1}_{\{\hat{X} \in \bar{F}_d\}}$  by  $\mathbf{1}_{\{\hat{X} \in \bar{F}_d, \tilde{X} \in \bar{F}_d\}}$  in the first term of the r.h.s. of Eq. (35), which gives

$$(36) \quad \sup_{x \in F_d} \mathbb{E}_{(\hat{X}, \tilde{X})|x} \left\{ (\hat{X}_1 - x_1)^2 \mathbf{1}_{\{\hat{X} \in \bar{F}_d, \tilde{X} \in \bar{F}_d\}} \cdot \left[ 1 \wedge \exp \left( \log \frac{f(\hat{X}_1)}{f(x_1)} + \log \frac{(R^2 + \|\hat{X}\|^2)^d}{(R^2 + \|\tilde{X}\|^2)^d} + \sum_{i=2}^d \log \frac{f(\hat{X}_i)}{f(x_i)} + \log \frac{(R^2 + \|\hat{X}\|^2)^d}{(R^2 + \|x\|^2)^d} \right) \right] \right\}.$$

Using the fact that  $1 \wedge \exp(\cdot)$  is 1-Lipschitz, for two random variables  $A$  and  $B$ , if  $\sup_{x \in F_d} \mathbb{E}[(\hat{X}_1 - x_1)^2 | A - B | \mathbf{1}_{\{\hat{X} \in \bar{F}_d, \tilde{X} \in \bar{F}_d\}}] = o(d^{-1})$ , then

$$\begin{aligned} & \sup_{x \in F_d} \mathbb{E} \left\{ (\hat{X}_1 - x_1)^2 \mathbf{1}_{\{\hat{X} \in \bar{F}_d, \tilde{X} \in \bar{F}_d\}} \left| [1 \wedge \exp(A)] - [1 \wedge \exp(B)] \right| \right\} \\ & \leq \sup_{x \in F_d} \mathbb{E} \left\{ (\hat{X}_1 - x_1)^2 | A - B | \mathbf{1}_{\{\hat{X} \in \bar{F}_d, \tilde{X} \in \bar{F}_d\}} \right\} = o(d^{-1}). \end{aligned}$$

Alternatively, if  $\sup_{x \in F_d} \mathbb{E}[(A - B)^2 \mathbf{1}_{\{\hat{X} \in \bar{F}_d, \tilde{X} \in \bar{F}_d\}}] = o(1)$ , then by Cauchy–Schwarz inequality

$$\begin{aligned} & \sup_{x \in F_d} \mathbb{E} \left\{ (\hat{X}_1 - x_1)^2 | A - B | \mathbf{1}_{\{\hat{X} \in \bar{F}_d, \tilde{X} \in \bar{F}_d\}} \right\} \\ & \leq \sup_{x \in F_d} \sqrt{\mathbb{E} \left[ (\hat{X}_1 - x_1)^4 \right]} \sqrt{\mathbb{E} \left[ (A - B)^2 \mathbf{1}_{\{\hat{X} \in \bar{F}_d, \tilde{X} \in \bar{F}_d\}} \right]} = o(dh^2) = o(d^{-1}). \end{aligned}$$

We next use the above approximations several times to approximate Eq. (36). We first approximate  $\log \frac{(R^2 + \|\hat{X}\|^2)^d}{(R^2 + \|\tilde{X}\|^2)^d}$  by  $\frac{\hat{X}_1^2 - \tilde{X}_1^2}{2}$ , since by Taylor expansion and mean value theorem, using the fact  $R^2 = d$ , one can show

$$\sup_{x \in F_d} \mathbb{E} \left[ (\hat{X}_1 - x_1)^2 \left| \log \frac{(R^2 + \|\hat{X}\|^2)^d}{(R^2 + \|\tilde{X}\|^2)^d} - \frac{\hat{X}_1^2 - \tilde{X}_1^2}{2} \right| \mathbf{1}_{\{\hat{X} \in \bar{F}_d, \tilde{X} \in \bar{F}_d\}} \right] = \mathcal{O}(d^{-1/2} dh^2),$$

and

$$\sup_{x \in F_d} \mathbb{E} \left[ (\hat{X}_1 - x_1)^2 \left| \frac{\tilde{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2} \right| \mathbf{1}_{\{\hat{X} \in \bar{F}_d, \tilde{X} \in \bar{F}_d\}} \right] = \mathcal{O}(d^{-1/2} dh^2).$$

Next, we consider the conditional expectation over  $\tilde{X}$  of the term  $\sum_{i=2}^d \log \left( \frac{f(\tilde{X}_i)}{f(x_i)} \right)$ . Since  $\sup_{x \in F_d} \mathbb{E} \left[ \left| \sum_{i=2}^d (\tilde{X}_i - x_i)^3 \right|^2 \mathbf{1}_{\{\tilde{X} \in \bar{F}_d\}} \right] = o(d^{-1})$ , we can approximate  $\sum_{i=2}^d \log \left( \frac{f(\tilde{X}_i)}{f(x_i)} \right)$  by the first two terms of its Taylor expansion. Using Taylor expansion and mean value theorem, we have

$$(37) \quad \begin{aligned} & \sup_{x \in F_d} \mathbb{E}_{\tilde{X}|x} \left[ \left( \sum_{i=2}^d \log \left( \frac{f(\tilde{X}_i)}{f(x_i)} \right) - (\text{Term I} + \text{Term II}) \right)^2 \mathbf{1}_{\{\tilde{X} \in \bar{F}_d\}} \right] \\ & = \mathcal{O} \left( \sup_{x \in F_d} \mathbb{E}_{\tilde{X}|x} \left[ \left| \sum_{i=2}^d (\tilde{X}_i - x_i)^3 \right|^2 \mathbf{1}_{\{\tilde{X} \in \bar{F}_d\}} \right] \right) = o(d^{-1}), \end{aligned}$$

where

$$\text{Term I} + \text{Term II} := \sum_{i=2}^d \left[ (\log f(x_i))'(\tilde{X}_i - x_i) + \frac{1}{2}(\log f(x_i))''(\tilde{X}_i - x_i)^2 \right].$$

Furthermore, following the Gaussian approximation arguments in the proof of Lemma 5.1, we have the following result:

$$\sup_{x \in F_d} \mathbb{E}_{\tilde{X}|x} \left[ \left( (\text{Term I} + \text{Term II}) + \log \frac{(R^2 + \|\tilde{X}\|^2)^d}{(R^2 + \|x\|^2)^d} - W_{\tilde{X}|x} \right)^2 \mathbf{1}_{\{\tilde{X} \in \bar{F}_d\}} \right] = o(1)$$

where  $W_{\tilde{X}|x} \sim \mathcal{N}(\mu, \sigma^2)$  with

$$\mu = \frac{\ell^2}{2} \left\{ 1 - \mathbb{E}_f \left[ ((\log f)')^2 \right] \right\}, \quad \sigma^2 = \ell^2 \left\{ \mathbb{E}_f \left[ ((\log f)')^2 \right] - 1 \right\}.$$

Therefore, after the above approximations, we have simplified Eq. (36) to

(38)

$$\sup_{x \in F_d} \mathbb{E}_{(\hat{X}_1, \tilde{X})|x} \left\{ (\hat{X}_1 - x_1)^2 \mathbf{1}_{\{\hat{X} \in \bar{F}_d\}} \left[ 1 \wedge \exp \left( \log \frac{f(\hat{X}_1)}{f(x_1)} + \frac{\hat{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2} + W_{\tilde{X}|x} \right) \right] \right\}.$$

Note that the randomness of  $\log \frac{f(\hat{X}_1)}{f(x_1)} + \frac{\hat{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2}$  only comes from  $\hat{X}_1$ , and the randomness of  $W_{\tilde{X}|x}$  comes from  $\tilde{X}$ . However,  $W_{\tilde{X}|x}$  is not independent with  $\hat{X}_1$ , because of the weak dependence of  $\hat{X}_1$  with  $\tilde{X}_{2:d} = \hat{X}_{2:d}$ . Indeed, it can be argued that  $W_{\tilde{X}|x}$  is asymptotically independent with  $\hat{X}_1$ . However, we wish  $W_{\tilde{X}|x}$  to be independent with  $\hat{X}_1$  in order to make further progress.

In Section S8.10.1, we will construct another random variable  $\tilde{W}$ , which is identically distributed with  $W_{\tilde{X}|x}$ , but independent with  $\hat{X}_1$ . Furthermore,  $\tilde{W}$  is “physically close” to  $W_{\tilde{X}|x}$  such that

$$(39) \quad \sup_{x \in F_d} \mathbb{E} \left[ \left( \tilde{W} - W_{\tilde{X}|x} \right)^2 \right] = o(1).$$

Note that the existence of  $\tilde{W}$  is very intuitive: if  $W_{\tilde{X}|x}$  is asymptotically independent with  $\hat{X}_1$ , then there should be an identically distributed random variable  $\tilde{W}$ , which is independent with  $\hat{X}_1$  and dependent with  $W_{\tilde{X}|x}$ . Furthermore,  $W_{\tilde{X}|x}$  becomes “physically closer and closer” to  $\tilde{W}$ .

Using Eq. (39) and the Lipschitz property of  $1 \wedge \exp(\cdot)$ , we can simplify Eq. (38) to

$$\begin{aligned}
(40) \quad & \sup_{x \in F_d} \mathbb{E}_{(\hat{X}_1, \tilde{X})|x} \left\{ (\hat{X}_1 - x_1)^2 \mathbf{1}_{\{\hat{X} \in \bar{F}_d\}} \left[ 1 \wedge \exp \left( \log \frac{f(\hat{X}_1)}{f(x_1)} + \frac{\hat{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2} + W_{\tilde{X}|x} \right) \right] \right\} \\
& \rightarrow \sup_{x \in F_d} \mathbb{E}_{(\hat{X}, \tilde{X})|x} \left\{ (\hat{X}_1 - x_1)^2 \mathbf{1}_{\{\hat{X} \in \bar{F}_d\}} \left[ 1 \wedge \exp \left( \log \frac{f(\hat{X}_1)}{f(x_1)} + \frac{\hat{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2} + \tilde{W} \right) \right] \right\} \\
& \rightarrow \sup_{x \in F_d} \mathbb{E}_{(\tilde{W}, \hat{X}_1)|x} \left\{ (\hat{X}_1 - x_1)^2 \left[ 1 \wedge \exp \left( \log \frac{f(\hat{X}_1)}{f(x_1)} + \frac{\hat{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2} + \tilde{W} \right) \right] \right\} \\
& = \sup_{x \in F_d} \mathbb{E}_{\hat{X}_1|x} \left\{ (\hat{X}_1 - x_1)^2 \mathbb{E}_{\tilde{W}|\hat{X}_1, x} \left[ 1 \wedge \exp \left( \log \frac{f(\hat{X}_1)}{f(x_1)} + \frac{\hat{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2} + \tilde{W} \right) \right] \right\}.
\end{aligned}$$

LEMMA S1. [*Roberts et al., 1997, Proposition 2.4*] If  $W \sim \mathcal{N}(\mu, \sigma^2)$  then

$$\mathbb{E}_W[1 \wedge \exp(W)] = \Phi(\mu/\sigma) + \exp(\mu + \sigma^2/2)\Phi(-\sigma - \mu/\sigma)$$

where  $\Phi(\cdot)$  is the standard normal cumulative distribution function.

By Lemma S1, the inside expectation of Eq. (40) is

$$\begin{aligned}
M(\hat{X}_1) &:= \mathbb{E}_{\tilde{W}|\hat{X}_1, x} \left[ 1 \wedge \exp \left( \log \frac{f(\hat{X}_1)}{f(x_1)} + \frac{\hat{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2} + \tilde{W} \right) \right] \\
&= \Phi \left( -\frac{\sigma}{2} + \frac{\log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + \frac{\hat{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2}}{\sigma} \right) \\
&\quad + \exp \left( \log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + \frac{\hat{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2} \right) \cdot \Phi \left( -\frac{\sigma}{2} - \frac{\log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + \frac{\hat{X}_1^2 - \mathbb{E}[\tilde{X}_1^2]}{2}}{\sigma} \right).
\end{aligned}$$

Note that  $\mathbb{E}[\tilde{X}_1^2] = x_1^2 + \mathcal{O}(d^{-1})$ . Then it can be easily verified that

$$M(x_1) \rightarrow 2\Phi \left( -\frac{\sigma}{2} \right).$$

Therefore, substituting  $M(\hat{X}_1)$  to Eq. (40) and by Taylor expansion, as  $d \rightarrow \infty$ , we have

$$\begin{aligned}
\text{ESJD} &\rightarrow d \cdot \left\{ \mathbb{E}[(\hat{X}_1 - x_1)^2] \cdot 2\Phi \left( -\frac{\sigma}{2} \right) \right\} \\
&\rightarrow 2\ell^2 \cdot \Phi \left( -\frac{\ell}{2} \sqrt{\mathbb{E}_f \left[ ((\log f)')^2 \right] - 1} \right),
\end{aligned}$$

which completes the proof. Therefore, it suffices to show the construction of  $\tilde{W}$ .

S8.10.1. *Construction of  $\tilde{W}$ .* For simplicity of notations, we will write  $W_{\tilde{X}|x}$  as  $W$ . Recall the following definitions in Section S8.7: all the randomness comes from a standard Gaussian  $\mathcal{N}(0, I_d)$ . We can write it as  $d$  independent standard  $\mathcal{N}(0, 1)$  as  $V_1, \dots, V_d$ . We have

also defined  $U_1, \dots, U_{d+1}$  which are marginal Gaussian random variables with  $U_i \sim \mathcal{N}(0, 1)$ . Recall that  $U_1, \dots, U_{d+1}$  are usually not independent. However, in the stationary phase when  $R = d^{1/2}$ ,  $\{U_i\}$  are “almost independent” since their correlations are very small. Particularly, if  $z_i \rightarrow 0$  then  $U_i$  is asymptotically independent with all  $\{U_j : j \neq i\}$ .

Recall that in Section S8.7, we have derived that the proposed  $\hat{z}_i$  can be written as

$$\hat{z}_i = \frac{\sqrt{1 + h^2 U_i^2}}{\sqrt{1 + h^2 \left( \sum_{j=1}^d V_j^2 \right)}} \left( \frac{z_i}{\sqrt{1 + h^2 U_i^2}} + \frac{\sqrt{1 - z_i^2} h U_i}{\sqrt{1 + h^2 U_i^2}} \right), \quad i = 1, \dots, d + 1.$$

Also,  $\hat{X}_i = d^{1/2} \frac{\hat{z}_i}{1 - \hat{z}_{d+1}}$  as  $R = d^{1/2}$ . To emphasize the sources of randomness, we write

$$\hat{z}_i = \mathcal{F}_{\hat{z}_i} \left( \sum_{j=1}^d V_j^2, U_i \right), \quad i = 1, \dots, d + 1,$$

to denote that the randomness of  $\hat{z}_i$  comes from  $\sum_j V_j^2$  and  $U_i$ . Similarly, we can write  $\hat{X}_i$  as

$$\hat{X}_i = \mathcal{F}_{\hat{X}_i} \left( \sum_{j=1}^d V_j^2, U_i, U_{d+1} \right), \quad i = 1, \dots, d$$

to denote the randomness of  $\hat{X}_i$  comes from  $\sum_j V_j^2$ ,  $U_i$ , and  $U_{d+1}$ .

Now recall that  $W$  is a function of  $\tilde{X}$ , where  $\tilde{X}_1$  is conditional independent with  $\hat{X}_1$ , and  $\tilde{X}_i = \hat{X}_i$  for all  $i = 2, \dots, d$ . We can therefore write  $W$  as a random function of  $\hat{X}_{2:d}$ . Therefore, we can write the following:

$$W = \mathcal{F}_W \left( \sum_{j=1}^d V_j^2, U_{2:d+1} \right), \quad \hat{X}_1 = \mathcal{F}_{\hat{X}_1} \left( \sum_{j=1}^d V_j^2, U_1, U_{d+1} \right).$$

Now, it is clear that the dependence of  $W$  on  $\hat{X}_1$  has three sources:

1. Shared source of randomness,  $\sum_{j=1}^d V_j^2$ , which a chi-squared random variable with degree of freedom  $d$ ;
2. The weak dependence between  $U_{2:d}$  and  $U_1$ . Note that  $U_1, U_{2:d}$  are correlated and joint Gaussian distributed. When  $z_1 \rightarrow 0$ , then  $U_1$  becomes (asymptotically) independent with  $U_{2:d}$ .
3. Shared source of randomness  $U_{d+1}$ , which becomes (asymptotically) independent with  $U_{1:d}$  when  $z_{d+1} \rightarrow 0$ .

In order to replace  $W$  by an identical distributed  $\tilde{W}$ , that is also independent with  $\hat{X}_1$ , we do three coupling arguments as follows:

1. First coupling argument: we first replace  $\sum_{j=1}^d V_j^2$  by a constant  $d$  to define

$$W' := \mathcal{F}_W(d, U_{2:d+1})$$

Following the proof of Lemma 5.1, it's very easy to check that  $W'$  has the same distribution as  $W$ . The key is to show  $W'$  and  $W$  are “physically close” so that

$$\sup_{x \in F_d} \mathbb{E}[(W' - W)^2] = o(1).$$

Actually, we can show a better rate that

$$\sup_{x \in F_d} \mathbb{E}[(W' - W)^2] = \mathcal{O}(d^{-1}).$$

The proof is delayed to Section S8.10.2.

2. Second coupling argument: we remove the dependence of  $U_{2:d+1}$  on  $U_1$ . Note that  $U_{2:d+1}$  are “almost independent” with  $U_1$ . We “orthogonalize” these Gaussian random variables  $U_{2:d+1}$  by the following decomposition:

$$U_i = c_i^{\parallel} U_1 + c_i^{\perp} U_i^{\perp}, \quad i = 2, \dots, d+1$$

where  $c_i^{\parallel}$  is the component that  $U_i$  is “parallel” with  $U_1$ , and  $c_i^{\perp}$  is the “orthogonal” (that is, independent) component. Now we replace the  $U_1$  in the decomposition by an independent copy  $\tilde{U}_1$ . That is, we define

$$U'_i := c_i^{\parallel} \tilde{U}_1 + c_i^{\perp} U_i^{\perp}, \quad i = 2, \dots, d+1,$$

where  $\tilde{U}_1$  is independent standard Gaussian random variable. Note that  $U'_i$  is highly correlated with  $U_i$  since the component  $c_i^{\perp} U_i^{\perp}$  remains. If  $z_1 = 0$ , for example, we recover  $U'_i = U_i$ . Next, we define

$$W'' := \mathcal{F}_W(d, U'_{2:d+1}).$$

It is clear that  $W''$  is identically distributed as  $W'$ . When  $z_1 = 0$ , we recover  $|W'' - W'| = 0$ . Intuitively,  $W''$  is “physically close” to  $W'$ , when  $z_1$  is small. Actually, our definition of  $F_d$  guarantees that for all  $x \in F_d$  we have  $x_1 = \mathcal{O}(d^{1/5})$  which implies  $z_1 = \mathcal{O}(d^{1/5-1/2})$  is indeed small. Then, our goal is to show

$$\sup_{x \in F_d} \mathbb{E}[(W'' - W')^2] = o(1).$$

Actually, we can prove a slightly better rate, which is

$$\sup_{x \in F_d} \mathbb{E}[(W' - W'')^2] = \mathcal{O}(z_1^2) = \mathcal{O}(d^{-3/5}).$$

The proof is delayed to Section S8.10.3.

3. Third coupling argument: finally, we construct  $\tilde{W}$  based on  $W''$ . Recall that  $W'' = \mathcal{F}_W(d, U'_{2:d+1})$ . In order to construct  $\tilde{W}$  which is independent with  $\hat{X}_1$ , it suffices to replace the source of  $U'_{d+1}$  by an independent copy  $\tilde{U}_{d+1}$ . We first “orthogonalize”  $U'_{2:d}$  using the same way as the previous coupling argument:

$$U'_i = \tilde{c}_i^{\parallel} U'_{d+1} + \tilde{c}_i^{\perp} U_i^{\perp}, \quad i = 2, \dots, d.$$

and then replace the component  $U'_{d+1}$  by an independent copy  $\tilde{U}_{d+1}$  to define  $\tilde{U}_{2:d}$  by

$$\tilde{U}_i := \tilde{c}_i^{\parallel} \tilde{U}_{d+1} + \tilde{c}_i^{\perp} U_i^{\perp}, \quad i = 2, \dots, d.$$

Finally, we define

$$\tilde{W} := \mathcal{F}_W(d, \tilde{U}_{2:d+1}).$$

It is clear that  $\tilde{W}$  is identically distributed as  $W''$  and it is now independent with  $\hat{X}_1$ . Then the goal is to show

$$\sup_{x \in F_d} \mathbb{E}[(\tilde{W} - W'')^2] = o(1).$$

It can be seen that the ‘‘orthogonalization’’ part of the coupling argument is the same as the second coupling. We then expect this part cause a distance of order  $\mathcal{O}_{\mathbb{P}}(z_{d+1})$ , which is indeed the case. By the definition of  $F_d$ , for all  $x \in F_d$  we have  $\|x_i\|^2 - d = o(d^{6/7})$  so  $z_{d+1} = o(d^{6/7-1}) = o(d^{-1/7})$ . The other part of the coupling argument is to replace the direct dependence on  $U'_{d+1}$  by  $\tilde{U}_{d+1}$ , which we can prove a distance of order  $\mathcal{O}(d^{-1/8})$ . Overall, we can show

$$\sup_{x \in F_d} \mathbb{E}[(\tilde{W} - W'')^2] = \mathcal{O}(d^{-1/4}) + \mathcal{O}(z_{d+1}^2) = \mathcal{O}(d^{-1/4}) + o(d^{-2/7}).$$

The proof is delayed to Section [S8.10.4](#).

Overall, we have shown the construction of  $\tilde{W}$ , which has the same distribution as  $W$  and is independent with  $\hat{X}_1$ . Most importantly, we have shown

$$\begin{aligned} \sup_{x \in F_d} \mathbb{E}[(\tilde{W} - W)^2] &\leq \sup_{x \in F_d} \mathbb{E}[ (|\tilde{W} - W''| + |W'' - W'| + |W' - W|)^2 ] \\ &= \mathcal{O}(d^{-1/4}) + o(d^{-2/7}) + \mathcal{O}(d^{-3/5}) + \mathcal{O}(d^{-1}) = o(1). \end{aligned}$$

**S8.10.2. The first coupling argument.** We first define the replacement of  $\hat{z}_i$  by

$$z'_i := \frac{\sqrt{1 + h^2 U_i^2}}{\sqrt{1 + h^2 d}} \left( \frac{z_i}{\sqrt{1 + h^2 U_i^2}} + \frac{\sqrt{1 - z_i^2} h U_i}{\sqrt{1 + h^2 U_i^2}} \right), \quad i = 1, \dots, d + 1,$$

in which we only replace  $\sum_j V_j^2$  in  $\hat{z}_i$  by a constant  $d$ . It is clear that  $\hat{z}_i$  and  $z'_i$  are highly correlated, we can write

$$(41) \quad \hat{z}_i = (1 + \mathcal{O}_{\mathbb{P}}(h^2 d^{1/2})) \frac{\sqrt{1 + h^2 U_i^2}}{\sqrt{1 + h^2 d}} \left( \frac{z_i}{\sqrt{1 + h^2 U_i^2}} + \frac{\sqrt{1 - z_i^2} h U_i}{\sqrt{1 + h^2 U_i^2}} \right).$$

Following the proof of Lemma [5.1](#), it is obvious that  $W'$  has the same distribution as  $W$ , which means the mean of  $W - W'$  is zero. Therefore, we only need to focus how fast the variance of  $W - W'$  goes to zero.

Recall that in the proof of  $W$  using Lemma 5.1, the variance of  $W$  is dominated by the variance of the following inner product term

$$\frac{R}{1 - z_{d+1}} \underbrace{\left( 0, (\log f(x_2))', \dots, (\log f(x_d))', \sum_{i=2}^d \left( (\log f(x_i))' z_i + \frac{1}{R} \right) \right)}_{\text{inner product}} \cdot (\hat{z} - z)$$

When we prove the variance of  $W - W'$ , by replacing  $\hat{z} - z$  to the difference of  $\hat{z} - z$  and  $z' - z$ , which is  $\hat{z} - z'$ , following the same arguments as the proof of Lemma 5.1, the variance of  $W - W'$  is then determined by the variance of

$$\frac{R}{1 - z_{d+1}} \left( 0, (\log f(x_2))', \dots, (\log f(x_d))', \sum_{i=2}^d \left( (\log f(x_i))' z_i + \frac{1}{R} \right) \right) \cdot (\hat{z} - z')$$

where  $\hat{z} - z' = (\hat{z}_1 - z'_1, \dots, \hat{z}_{d+1} - z'_{d+1})^T$ . Substituting the definition of  $\hat{z}$  and  $z'$ , using  $\mathcal{O}(h^2 d^{1/2}) = \mathcal{O}(d^{-3/2})$ , we get the order of variance as  $\mathcal{O} \left( \left( \sum_{i=1}^d (\log f(x_i))' x_i \right)^2 \right) \mathcal{O}(d^{-3})$ .

By the definition of  $F_d$ , we know  $\sum_{i=1}^d (\log f(x_i))' x_i = \mathcal{O}(d)$ . Therefore, we have

$$\sup_{x \in F_d} \mathbb{E}[(W - W')^2] = \mathcal{O}(d^2) \mathcal{O}(d^{-3}) = \mathcal{O}(d^{-1}).$$

**S8.10.3. The second coupling argument.** By the construction of  $W''$ , it is clear that  $W''$  has the same distribution as  $W'$ . Define the replacement of  $z'_i$  by

$$z''_i := \frac{\sqrt{1 + h^2 (U'_i)^2}}{\sqrt{1 + h^2 d}} \left( \frac{z_i}{\sqrt{1 + h^2 (U'_i)^2}} + \frac{\sqrt{1 - z_i^2} h U'_i}{\sqrt{1 + h^2 (U'_i)^2}} \right), \quad i = 1, \dots, d + 1,$$

Then we can again follow the proof of Lemma 5.1 and focus on the variance of  $W' - W''$  (since the mean of  $W' - W''$  is clearly zero). Then the variance of  $W' - W''$  is determined by the variance of

$$\frac{R}{1 - z_{d+1}} \left( 0, (\log f(x_2))', \dots, (\log f(x_d))', \sum_{i=2}^d \left( (\log f(x_i))' z_i + \frac{1}{R} \right) \right) \cdot (z'' - z')$$

where  $z'' - z' = (z''_1 - z'_1, \dots, z''_{d+1} - z'_{d+1})^T$ . Using the fact that

$$U_i - U'_i = c_i^\parallel (U_1 - \tilde{U}_1), \quad i = 2, \dots, d + 1,$$

we can get that the variance of  $W' - W''$  is bounded by

$$\left\{ d \sum_{i=2}^d ((\log f(x_i))')^2 + \left[ \sum_{i=2}^d ((\log f(x_i))' x_i + 1) \right]^2 \right\} \left\{ 2h^2 \sum_{i=2}^d (c_i^\parallel)^2 + o(h^2) \right\}$$

Note that, by the definition of  $F_d$ , we have

$$d \sum_{i=2}^d ((\log f(x_i))')^2 = \mathcal{O}(d^2), \quad \left[ \sum_{i=2}^d ((\log f(x_i))' x_i + 1) \right]^2 = o(d^2).$$



Using some basic geometry to analyze the ‘‘angle’’ between  $U_i$  and  $U_1$ , we know  $c_i^\parallel = \mathcal{O}(z_i z_1)$ . Therefore, using  $h^2 = \mathcal{O}(d^{-2})$ , we have

$$\sup_{x \in F_d} \mathbb{E}[(W' - W'')^2] = \mathcal{O} \left( \sum_{i=2}^d (c_i^\parallel)^2 \right) = \mathcal{O} \left( \sum_{i=2}^d z_i^2 z_1 \right) = \mathcal{O}(z_1^2).$$

Finally, for all  $x \in F_d$ , we have  $x_1 = \mathcal{O}(d^{1/5})$  which implies that  $z_1 = \mathcal{O}(d^{1/5-1/2}) = \mathcal{O}(d^{-3/10})$ .

**S8.10.4. The third coupling argument.** By the construction of  $\tilde{W}$ , it is clear that the distribution of  $\tilde{W}$  is the same as the distribution of  $W''$ . Therefore, the mean of  $\tilde{W} - W''$  is zero.

We define the replacement of  $z_i''$  by

$$\tilde{z}_i := \frac{\sqrt{1 + h^2 \tilde{U}_i^2}}{\sqrt{1 + h^2 \bar{d}}} \left( \frac{z_i}{\sqrt{1 + h^2 \tilde{U}_i^2}} + \frac{\sqrt{1 - z_i^2 h \tilde{U}_i}}{\sqrt{1 + h^2 \tilde{U}_i^2}} \right), \quad i = 1, \dots, d+1.$$

Again, following the proof of Lemma 5.1, the variance of  $\tilde{W} - W''$  is determined by the variance of

$$\begin{aligned} & \frac{R}{1 - z_{d+1}} \left( 0, (\log f(x_2))', \dots, (\log f(x_d))', \sum_{i=2}^d \left( (\log f(x_i))' z_i + \frac{1}{R} \right) \right) \cdot (\tilde{z} - z'') \\ &= \underbrace{\frac{R}{1 - z_{d+1}} \left( 0, (\log f(x_2))', \dots, (\log f(x_d))', 0 \right) \cdot (\tilde{z} - z'')}_{\text{the first term}} \\ & \quad + \underbrace{\frac{R}{1 - z_{d+1}} \left( 0, \dots, 0, \sum_{i=2}^d \left( (\log f(x_i))' z_i + \frac{1}{R} \right) \right) \cdot (\tilde{z} - z'')}_{\text{the second term}} \end{aligned}$$

where  $\tilde{z} - z'' = (\tilde{z}_1 - z_1'', \dots, \tilde{z}_{d+1} - z_{d+1}'')$ . Note that we have decomposed the inner product to the sum of two terms. We compute the variances of the two terms separately.

Similar to the previous coupling argument, the variance of the first term is bounded by

$$\left[ d \sum_{i=2}^d ((\log f(x_i))')^2 \right] \left[ 2h^2 \sum_{i=2}^d (c_i^\parallel)^2 + o(h^2) \right] = \mathcal{O}(z_{d+1}^2).$$

The second term can be written as

$$\sum_{i=2}^d [(\log f(x_i))' x_i + 1] (1 + \mathcal{O}_{\mathbb{P}}(z_{d+1})) (\tilde{z}_{d+1} - z_{d+1}'').$$

By the definition of  $F_d$ , we know

$$\sum_{i=2}^d [(\log f(x_i))' x_i + 1] = d\mathcal{O}(d^{-1/8}), \quad z_{d+1} = o(d^{-1/7})$$

Furthermore, by the definition of  $\tilde{z}_{d+1}$  and  $z''_{d+1}$ , both have finite exponential moments. Moreover, we have

$$\mathbb{E} \left[ \left| \tilde{z}_{d+1} - z''_{d+1} - \sqrt{1 - z_{d+1}^2} h(\tilde{U}_{d+1} - U'_{d+1}) \right| \right] = \mathcal{O}(d^{-2}).$$

Therefore, we have the second term

$$\begin{aligned} & \sum_{i=2}^d [(\log f(x_i))' x_i + 1] (1 + \mathcal{O}_{\mathbb{P}}(z_{d+1})) (\tilde{z}_{d+1} - z''_{d+1}) \\ & = d\mathcal{O}(d^{-1/8})(1 + o(d^{-1/7}))h\mathcal{O}_{\mathbb{P}}(1) = \mathcal{O}_{\mathbb{P}}(d^{-1/8}), \end{aligned}$$

and its variance is of order  $\mathcal{O}(d^{-1/4})$ . Overall, adding the variance of the second term to the variance of the first term yields

$$\sup_{x \in F_d} \mathbb{E}[(\tilde{W} - W'')^2] = \mathcal{O}(z_{d+1}^2) + \mathcal{O}(d^{-1/4}) = o(d^{-2/7}) + \mathcal{O}(d^{-1/4}).$$

This completes the proof. □

### S8.11. Proof of Theorem 5.2.

PROOF. We follow the framework of [Roberts et al., 1997] using the generator approach [Ethier and Kurtz, 1986]. Define the (discrete-time) generator of  $x$  by

$$(G_d V)(x) := d\mathbb{E}_{\hat{X}} \left\{ [V(\hat{X}) - V(x)] \left( 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(x)(R^2 + \|x\|^2)^d} \right) \right\}$$

for any function  $V$  for which this definition makes sense. In the Skorokhod topology, it doesn't cause any problem to treat  $G_d$  as a continuous time generator. We shall restrict attention to test functions such that  $V(x) = V(x_1)$ . We show uniform convergence of  $G_d$  to  $G$ , the generator of the limiting (one-dimensional) Langevin diffusion, for a suitable large class of real-valued functions  $V$ , where, for some fixed function  $h(\ell)$ ,

$$(GV)(x_1) := h(\ell) \left\{ \frac{1}{2} V''(x_1) + \frac{1}{2} [(\log f)'(x_1)] V'(x_1) \right\}.$$

Since  $f'/f$  is Lipschitz, by [Ethier and Kurtz, 1986, Thm 2.1 in Ch.8], a core for the generator has domain  $C_c^\infty$ , which is the class of continuous functions with compact support such that all orders of derivatives exist. This enable us to restrict attentions to functions  $V \in C_c^\infty$  such that  $V(x) = V(x_1)$ .

Define the sequence of sets  $\{F_d\}$  by

$$\begin{aligned}
(42) \quad F_d &:= \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d-1} \sum_{i=2}^d [(\log f(x_i))']^2 - \mathbb{E}_f [((\log f)')^2] \right| < d^{-1/8} \right\} \\
&\cap \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d-1} \sum_{i=2}^d [(\log f(x_i))''] - \mathbb{E}_f [((\log f)'')] \right| < d^{-1/8} \right\} \\
&\cap \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d-1} \sum_{i=2}^d x_i (\log f(x_i))' - \mathbb{E}_f [X (\log f)'] \right| < d^{-1/8} \right\} \\
&\cap \left\{ x \in \mathbb{R}^d : \left| \frac{1}{d} \sum_{i=1}^d x_i^2 - \mathbb{E}_f (X^2) \right| < d^{-1/6} \log(d) \right\} \\
&\cap \left\{ x \in \mathbb{R}^d : |x_1| < d^{1/5} \right\}
\end{aligned}$$

Then, for fixed  $t$ , by the union bound, Markov inequality, and the assumptions Eq. (14), we have

$$\begin{aligned}
&\mathbb{P} \left( X^d(\lfloor ds \rfloor) \notin F_d, \exists 0 \leq s \leq t \right) \leq td \mathbb{P}_\pi (X \notin F_d) \\
&= \mathcal{O} \left( td \left( \frac{d^{1/8}}{(d-1)^{1/2}} \right)^4 + td \left( \frac{d^{1/6}}{\log(d)d^{1/2}} \right)^3 + td \left( \frac{1}{d^{1/5}} \right)^6 \right) = \mathcal{O}(t(\log(d))^{-3}).
\end{aligned}$$

Therefore, for any fixed  $t$ , if  $d \rightarrow \infty$ , the probability of all  $\{X^d(\lfloor ds \rfloor), 0 \leq s \leq t\}$  are in  $F_d$  goes to 1. It then suffices to consider only  $x \in F_d$ .

Next, we decompose the expectation over  $\hat{X} = (\hat{X}_1, \dots, \hat{X}_d)$  into expectations over  $\hat{X}_1$  and  $\hat{X}_{2:d} := (\hat{X}_2, \dots, \hat{X}_d)$ . Recall that  $\hat{X}_1 = \hat{X}'_1$ ,  $\hat{X}_{2:d} = \hat{X}''_{2:d}$ , and  $\hat{X}'$  is independent with  $\hat{X}''$ . Therefore,  $\hat{X}_1$  is independent with  $\hat{X}_{2:d}$ . Then, we have

$$\begin{aligned}
(GdV)(x) &= d \mathbb{E}_{\hat{X}_1} \left\{ [V(\hat{X}_1) - V(x_1)] \mathbb{E}_{\hat{X}_{2:d}} \left[ 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(x)(R^2 + \|x\|^2)^d} \right] \right\} \\
&= d \mathbb{E}_{\hat{X}_1} \left\{ [V(\hat{X}_1) - V(x_1)] \mathbb{E}_{\hat{X}'_1 | \hat{X}_{2:d}} \mathbb{E}_{\hat{X}_{2:d}} \left[ 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(x)(R^2 + \|x\|^2)^d} \right] \right\} \\
&= d \mathbb{E}_{\hat{X}_1} \left\{ [V(\hat{X}_1) - V(x_1)] \mathbb{E}_{\hat{X}''} \left[ 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(x)(R^2 + \|x\|^2)^d} \right] \right\}.
\end{aligned}$$

Now we focus on the inner expectation

$$\begin{aligned}
&\mathbb{E}_{\hat{X}''} \left[ 1 \wedge \frac{\pi(\hat{X})(R^2 + \|\hat{X}\|^2)^d}{\pi(x)(R^2 + \|x\|^2)^d} \right] \\
&= \mathbb{E}_{\hat{X}''} \left[ 1 \wedge \exp \left\{ d \log \left( \frac{R^2 + \|\hat{X}\|^2}{R^2 + \|\hat{X}''\|^2} \right) + \sum_{i=1}^d \log \left( \frac{f(\hat{X}_i)}{f(x_i)} \right) + d \log \left( \frac{R^2 + \|\hat{X}''\|^2}{R^2 + \|x\|^2} \right) \right\} \right].
\end{aligned}$$

Following the proof of Theorem 5.1, for any  $x \in F_d$ , we can replace  $\sum_{i=2}^d \log \left( \frac{f(\hat{X}_i)}{f(x_i)} \right) + d \log \left( \frac{R^2 + \|\hat{X}''\|^2}{R^2 + \|x\|^2} \right)$  by  $W \sim \mathcal{N}(\mu, \sigma^2)$  is a Gaussian random variable with

$$\mu = \frac{\ell^2}{2} \left\{ 1 - \mathbb{E}_f \left[ ((\log f)')^2 \right] \right\}, \quad \sigma^2 = \ell^2 \left\{ \mathbb{E}_f \left[ ((\log f)')^2 \right] - 1 \right\}.$$

To keep the dependence on  $\hat{X}_1$ , we denote  $r(\hat{X}_1) := d \mathbb{E}_{\hat{X}''} \left[ \log \left( \frac{R^2 + \|\hat{X}\|^2}{R^2 + \|\hat{X}''\|^2} \right) \right]$ . Using Taylor expansion we can verify that  $r(x_1) = \mathcal{O}(d^{-1})$  and  $r'(x_1) = x_1 + o(d^{-1/2})$ . Therefore, it suggests that we can approximate

$$d \log \left( \frac{R^2 + \|\hat{X}\|^2}{R^2 + \|\hat{X}''\|^2} \right) + \sum_{i=1}^d \log \left( \frac{f(\hat{X}_i)}{f(x_i)} \right) + d \log \left( \frac{R^2 + \|\hat{X}''\|^2}{R^2 + \|x\|^2} \right)$$

by  $\log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + r(\hat{X}_1) + W$ . Hence, we define

(43)

$$(\tilde{G}_d V)(x) := d \mathbb{E}_{\hat{X}_1} \left\{ [V(\hat{X}_1) - V(x_1)] \cdot \mathbb{E}_{\hat{X}''} \left[ 1 \wedge \exp \left\{ \log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + r(\hat{X}_1) + W \right\} \right] \right\}.$$

Since  $V \in C_c^\infty$ , we have for some  $Z_1 \in (x_1, \hat{X}_1)$  or  $(\hat{X}_1, x_1)$  that

$$V(\hat{X}_1) - V(x_1) = V'(x_1)(\hat{X}_1 - x_1) + \frac{1}{2} V''(Z_1)(\hat{X}_1 - x_1)^2.$$

This implies that  $\mathbb{E}_{\hat{X}_1} [V(\hat{X}_1) - V(x_1)] = \mathcal{O}(x_1 d^{-1}) = o(d^{-1/2})$  and  $\mathbb{E} \left[ \left| V(\hat{X}_1) - V(x_1) \right| \right] = \mathcal{O}(d^{-1/2})$ , since  $x_1 = \mathcal{O}(d^{1/5})$  for  $x \in F_d$ . Using the fact that the function  $1 \wedge \exp(x)$  is Lipschitz [Roberts et al., 1997, Proposition 2.2], we can have

$$\begin{aligned} \sup_{x \in F_d} \left| (G_d V)(x) - (\tilde{G}_d V)(x) \right| &= d \sup_{x \in F_d} \mathbb{E}_{\hat{X}_1} [V(\hat{X}_1) - V(x_1)] \mathcal{O}(d^{-1/2}) \\ &\quad + d \sup_{x \in F_d} \mathbb{E}_{\hat{X}_1} \left[ \left| V(\hat{X}_1) - V(x_1) \right| \left( o(d^{-1/2}) \right) \right] = o(1). \end{aligned}$$

Therefore, we can now concentrate on  $(\tilde{G}_d V)(x)$  defined in Eq. (43) for  $x \in F_d$ .

Note that conditional on  $\hat{X}_1$ , the term inside the inner expectation of Eq. (43) is Gaussian distributed, since  $\log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + r(\hat{X}_1) + W \sim \mathcal{N}(\mu', \sigma^2)$  where  $\mu' := \mu + \log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + r(\hat{X}_1)$ . Therefore, by Lemma S1, we have

$$\begin{aligned} M(\hat{X}_1) &:= \mathbb{E}_{\hat{X}''} \left[ 1 \wedge \exp \left\{ \log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + r(\hat{X}_1) + W \right\} \right] \\ &= \Phi \left( -\frac{\sigma}{2} + \frac{\log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + r(\hat{X}_1)}{\sigma} \right) \\ &\quad + \exp \left( \log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + r(\hat{X}_1) \right) \cdot \Phi \left( -\frac{\sigma}{2} - \frac{\log \left( \frac{f(\hat{X}_1)}{f(x_1)} \right) + r(\hat{X}_1)}{\sigma} \right). \end{aligned}$$

Since  $r(x_1) = \mathcal{O}(d^{-1})$  and  $r'(x_1) = x_1 + o(d^{-1/2})$ , it can be verified that

$$(44) \quad M(x_1) \rightarrow 2\Phi\left(-\frac{\sigma}{2}\right), \quad M'(x_1) \rightarrow \Phi\left(-\frac{\sigma}{2}\right) \left[ (\log f(x_1))' + x_1 + o(d^{-1/2}) \right].$$

Furthermore, one can verify  $M''(x_1)$  is bounded since  $\Phi(\cdot)$ ,  $\Phi'(\cdot)$ , and  $\Phi''(\cdot)$  are all bounded functions.

Therefore, by mean value theorem, there exist constants  $K_1$  and  $K_2$  such that

$$\begin{aligned} & d[V(\hat{X}_1) - V(x_1)]M(\hat{X}_1) \\ &= d \left[ V'(x_1)(\hat{X}_1 - x_1) + \frac{1}{2}V''(x_1)(\hat{X}_1 - x_1)^2 + K_1(\hat{X}_1 - x_1)^3 \right] \\ & \quad \cdot \left[ M(x_1) + M'(x_1)(\hat{X}_1 - x_1) + K_2(\hat{X}_1 - x_1)^2 \right] \\ &= dV'(x_1)M(x_1)(\hat{X}_1 - x_1) \\ & \quad + d \left[ \frac{1}{2}V''(x_1)M(x_1) + V'(x_1)M'(x_1) \right] (\hat{X}_1 - x_1)^2 + \mathcal{O}(|d(\hat{X}_1 - x_1)^3|). \end{aligned}$$

Taking expectations over  $\hat{X}_1$ , using Eq. (44), as  $d \rightarrow \infty$ , we have

$$\begin{aligned} (\tilde{G}_d V)(x) &= \mathbb{E}_{\hat{X}_1} \left[ d[V(\hat{X}_1) - V(x_1)]M(\hat{X}_1) \right] \\ &\rightarrow dV'(x_1)M(x_1)\mathbb{E}[\hat{X}_1 - x_1] + d \left[ \frac{1}{2}V''(x_1)M(x_1) + V'(x_1)M'(x_1) \right] \mathbb{E}[(\hat{X}_1 - x_1)^2] \\ &\rightarrow -\frac{\ell^2}{2}x_1V'(x_1)M(x_1) + \ell^2 \left[ \frac{1}{2}V''(x_1)M(x_1) + V'(x_1)M'(x_1) \right] \\ &\rightarrow \Phi\left(-\frac{\sigma}{2}\right) \left[ -\ell^2x_1V'(x_1) + \ell^2V''(x_1) + \ell^2V'(x_1)[(\log f(x_1))' + x_1 + o(d^{-1/2})] \right] \\ &\rightarrow 2\ell^2\Phi\left(-\frac{\sigma}{2}\right) \left[ \frac{1}{2}V''(x_1) + \frac{1}{2}[(\log f)'(x_1)]V'(x_1) \right] \\ &\rightarrow 2\ell^2\Phi\left(-\frac{\ell\sqrt{\mathbb{E}_f\left[\left((\log f)'\right)^2\right]} - 1}{2}\right) \left[ \frac{1}{2}V''(x_1) + \frac{1}{2}[(\log f)'(x_1)]V'(x_1) \right]. \end{aligned}$$

This completes the proof.  $\square$

S8.12. *Comments on Remark 3.1.* To illustrate the issue that for heavy tail targets, the SPS might stuck if starting from the South Pole, we consider multivariate student's  $t$  targets. For these targets, one can easily derive that the ‘‘equator’’ is always a stationary point of the density (which is maximum if  $\nu < d$  and minimum if  $\nu > d$ ).

We first consider the case that  $k = \nu/d < 1$  is a small constant, from Lemma 3.1 and Fig. 8, one can see that in the *transient phase*, the acceptance rate goes to 0 exponentially fast as  $d \rightarrow \infty$ . For example, when  $k = 0.01$  if  $z = -1$  and  $\hat{z} = 0$ , then the acceptance rate is roughly  $\exp(-1.7d) \approx 1/5^d$ .

Next, we consider an approximation for the cases  $\nu = \mathcal{O}(1)$  and  $d \rightarrow \infty$ . Consider current location  $x = 0$  which corresponding to the south pole and the ‘‘typical’’ proposal  $\tilde{x}$  satisfying  $\|\tilde{x}\|^d \approx d$ . For multivariate student’s  $t$  target with DoF  $\nu = \mathcal{O}(1)$ , we have

$$\log \frac{\pi(\tilde{x})}{\pi(x)} \approx -\frac{\nu + d}{2} \log\left(1 + \frac{d}{\nu - 2}\right) = \mathcal{O}(d \log(d))$$

and  $\log \frac{(R^2 + \|\tilde{x}\|^2)^d}{(R^2 + \|x\|^2)^d} = \mathcal{O}(d)$ . Therefore, the acceptance probability for a ‘‘typical’’ proposal starting from the South Pole is

$$\min \left\{ 1, \frac{\pi(\tilde{x}) (R^2 + \|\tilde{x}\|^2)^d}{\pi(x) (R^2 + \|x\|^2)^d} \right\} \approx 1/(d^d),$$

which implies for multivariate student’s  $t$  targets with DoF  $\nu = \mathcal{O}(1)$ , the SPS chain initialized at the South Pole can stuck when the proposal variance is large.

**S8.13. Proof of the Jacobian determinant Eq. (2).** The Jacobian determinant of the stereographic projection is well-known. For example, see [Gehring et al., 2017, Lemma 3.2.1]. One way to derive it is to calculate the Jacobian by comparing the ratio of volumes. Let  $x = \text{SP}(z)$  and  $y = \text{SP}(w)$ , one can get

$$\begin{aligned} \|z - w\|^2 &= \sum_{i=1}^d \|z_i - w_i\|^2 + \|z_{d+1} - w_{d+1}\|^2 \\ &= \sum_{i=1}^d 4 \left( \frac{Rx_i}{R^2 + \|x\|^2} - \frac{Ry_i}{R^2 + \|y\|^2} \right)^2 + 4 \left( \frac{R^2}{R^2 + \|x\|^2} - \frac{R^2}{R^2 + \|y\|^2} \right)^2 \\ &= \frac{4R^2}{(R^2 + \|x\|^2)(R^2 + \|y\|^2)} \|x - y\|^2, \end{aligned}$$

where the last equality is obtained using

$$\sum_i [x_i(R^2 + \|y\|^2) - y_i(R^2 + \|x\|^2)]^2 = -R^2(\|y\|^2 - \|x\|^2)^2 + \|x - y\|^2(R^2 + \|x\|^2)(R^2 + \|y\|^2)$$

Then we know the Euclidean distance from  $\mathbb{R}^d$  to  $\mathbb{S}^d$  is scaled by  $\frac{2R}{\sqrt{(R^2 + \|x\|^2)(R^2 + \|y\|^2)}}$ . Therefore, comparing ratio of volume elements on  $\mathbb{R}^d$  and  $\mathbb{S}^d$ , the Jacobian determinant is proportional to  $(R^2 + \|x\|^2)^d$ .

**S9. Additional Simulations.** We have included some numerical examples in Section 6. In this section, we give additional simulations for SPS and SBPS. We use ACF to denote the *sample autocorrelation function*, which is an estimate of the autocorrelation function  $\text{acf}(k) := \frac{\mathbb{E}_\pi[(X_i - \mathbb{E}[X_i])(X_{i+k} - \mathbb{E}[X_{i+k}])]}{\sqrt{\text{var}_\pi(X_i)}\sqrt{\text{var}_\pi(X_{i+k})}}$ .

**S9.1. Some simulations using the result from Lemma 3.1.** We plot in Fig. 7 for some simulations using the result from Lemma 3.1.

**S9.2. The function  $g_k(\cdot)$  used in Section 3.2.** We plot the function  $g_k$  in Fig. 8 for different values of  $k$ .

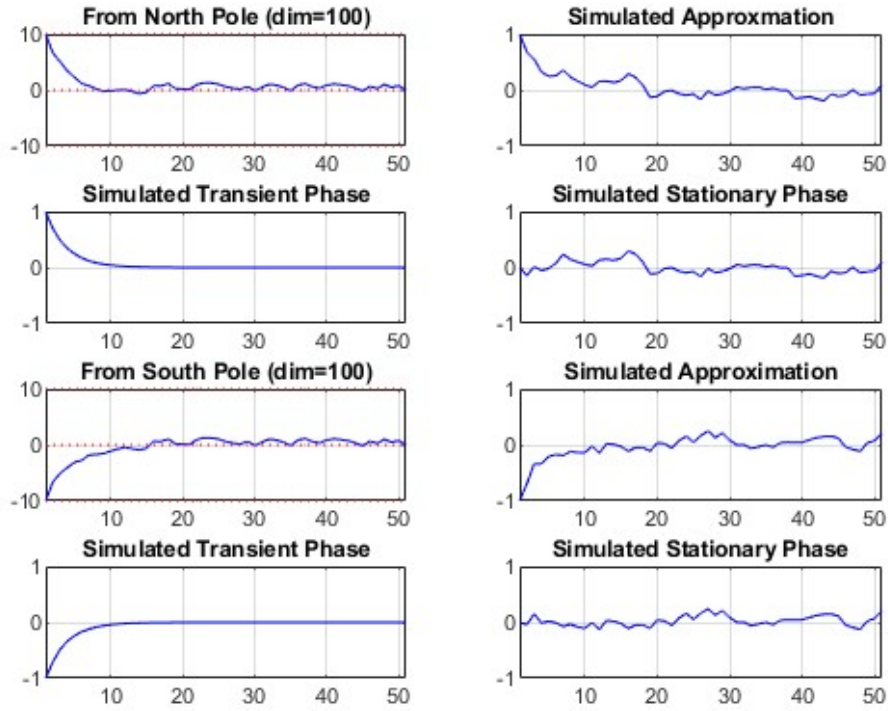


Figure 7: We consider  $h = 0.1$  and  $d = 100$  and two cases for the initial state: from the North Pole (up) and the South Pole (down). For each case, the first subplot is the traceplot of the proposal latitudes. The other three are approximations using Eq. (7), Eq. (8), and Eq. (9) from Lemma 3.1.

**S9.3. SPS: traceplot and ACF.** In this example, Fig. 9 shows traceplots and ACFs of SPS (first column) and RWM (the last three columns). SPS starts from the North Pole and RWM starts from different initial states  $(c, c, \dots, c)$  where  $c = 10, 20, 50$ , respectively (see the last three columns of Fig. 9). For SPS starting from the South Pole compared with RWM, we refer to Section S9.4. We plot the traceplots and ACFs of the 1-st coordinate and negative log-target density, as well as traceplots of the first two coordinates. The target is standard Gaussian in  $d = 100$  dimensions. The proposal variance for RWM is tuned such that the acceptance rate is about 0.234, which is known to be the optimal acceptance rate [Roberts et al., 1997]. For SPS, the acceptance rate is roughly 0.78, which is the lowest acceptance rate possible.

According to Fig. 9, the SPS mixes almost immediately. Indeed, the transient phase of SPS in  $d = 100$  dimensions is less than 10 iterations. Compared with SPS, the mixing time of RWM relies on the initial value. For example, in the last column of Fig. 9, the RWM hasn't mixed after 5000 iterations. In the stationary phase, with an acceptance rate of 0.78, SPS generates almost uncorrelated samples according to the ACFs. Compared with SPS, the samples from RWM with an optimal acceptance rate of 0.234 are highly correlated in high dimensions.

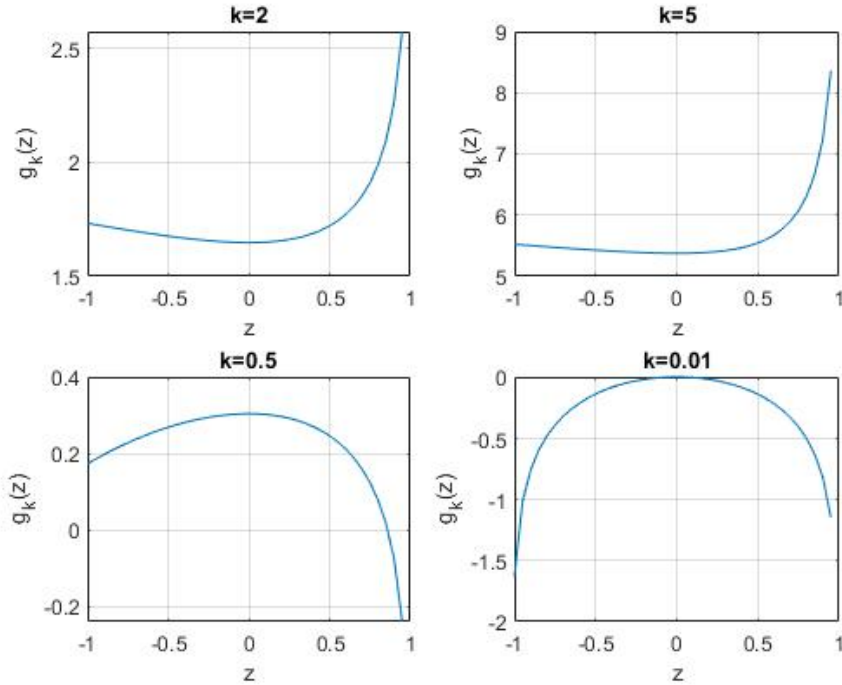


Figure 8: The function  $g_k(z)$  used in Section 3.2 for different values of  $k$ .

S9.4. *SPS: burn-in starting from the South Pole.* We have shown the traceplots and ACFs for starting from the North pole in Fig. 9. In this example, Fig. 10 shows traceplots and ACFs of SPS (the first column) starting from the South Pole and RWM (the second column) starting from the target mode. We plot the traceplots and ACFs of the 1-st coordinate and negative log-target density, as well as traceplots of the first two coordinates. The target is standard Gaussian in  $d = 100$  dimensions. The proposal variance for RWM is tuned such that the acceptance rate is about 0.234, which is known to be the optimal acceptance rate [Roberts et al., 1997]. For SPS, the acceptance rate is roughly 0.78, which is the lowest acceptance rate possible.

According to Fig. 10, the SPS mixes almost immediately. Indeed, the transient phase of SPS in  $d = 100$  dimensions is less than 10 iterations. Compared with SPS, even starting from the target mode, the first 100 iterations are the transient phase of RWM. In the stationary phase, with an acceptance rate of 0.78, SPS generates almost uncorrelated samples according to the ACFs. Compared with SPS, the samples from RWM with an optimal acceptance rate of 0.234 are highly correlated in high dimensions.

S9.5. *SPS versus GSPS.* In this example, we study SPS when the covariance matrix of the target distribution does not equal the identity matrix. We choose  $R = \sqrt{d}$  and the target is multivariate student's  $t$  with covariance matrix  $\Sigma$  that satisfies the sum of eigenvalues of  $\Sigma$  is  $d$ , i.e.,  $\lambda_1 + \dots + \lambda_d = d$ . We compare the decay of the performance of SPS with the



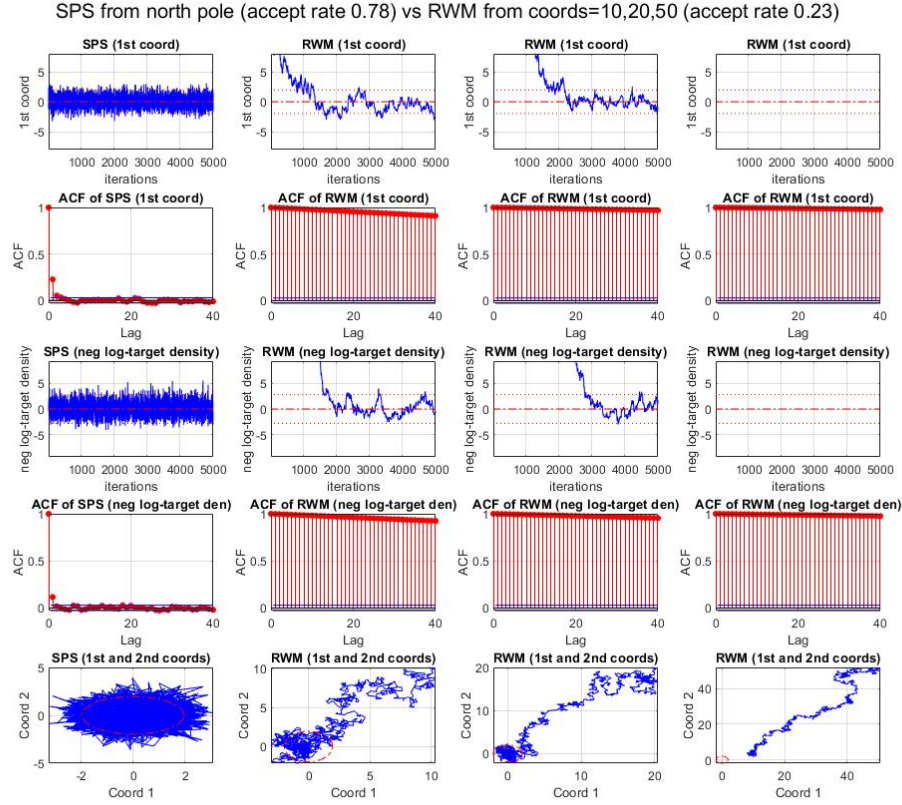


Figure 9: Traceplots and ACFs of the 1-st coordinate, negative log-target density, and the first two coordinates, for standard Gaussian target in 100 dimensions: SPS starts from the north pole (first column) vs RWM starts from different initial states  $(c, c, \dots, c)$  where  $c = 10, 20, 50$  (the last three columns).

performance of GSPS which uses the covariance matrix  $\Sigma$ . In the simulations, we consider the cases that a sequence of covariance matrices  $\Sigma_1, \Sigma_2, \dots, \Sigma_{50}$  are all block-diagonal and we add different numbers of diagonal sub-blocks which equal to a particular  $2 \times 2$  matrix. For example, for the case of one sub-block, denoted by  $\Sigma_1$ :

$$\Sigma_1 := \begin{pmatrix} 1 & 0.8 & 0 & \cdots & 0 \\ 0.8 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \cdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

One can easily verify the eigenvalues are

$$\Lambda_1 = (0.2, 1.8, 1, 1, \dots, 1)^T$$

SPS from south pole (acceptance rate 0.78) versus RWM from origin (acceptance rate 0.23)

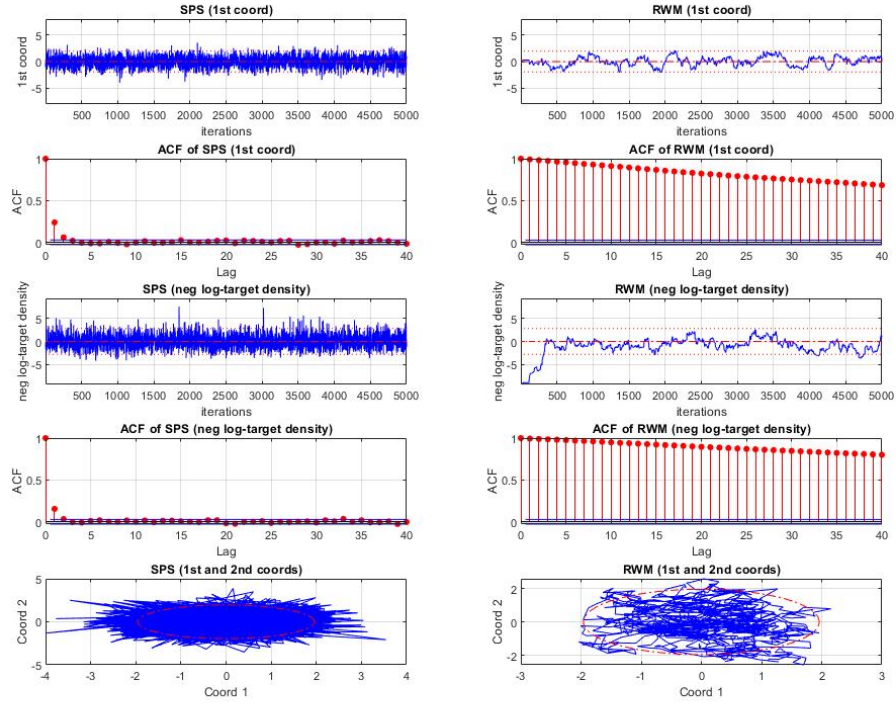


Figure 10: Traceplots and ACFs of the 1-st coordinate, negative log-target density, and the first two coordinates, for standard Gaussian target in dimension 100: SPS from south pole (the first column) vs RWM from the origin (the second column).

and the first two coordinates have rotation with angle  $3\pi/4$ , since

$$\begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix} = \begin{pmatrix} \cos \frac{3\pi}{4} & -\sin \frac{3\pi}{4} \\ \sin \frac{3\pi}{4} & \cos \frac{3\pi}{4} \end{pmatrix} \begin{pmatrix} 0.2 & 0 \\ 0 & 1.8 \end{pmatrix} \begin{pmatrix} \cos \frac{3\pi}{4} & -\sin \frac{3\pi}{4} \\ \sin \frac{3\pi}{4} & \cos \frac{3\pi}{4} \end{pmatrix}^T$$

Thus, we have the corresponding rotation matrix used by GSPS:

$$Q_1 := \begin{pmatrix} \cos \frac{3\pi}{4} & -\sin \frac{3\pi}{4} & 0 & \cdots & 0 \\ \sin \frac{3\pi}{4} & \cos \frac{3\pi}{4} & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \cdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1. \end{pmatrix}$$

Similarly, we can replace the diagonal terms by 2 sub-blocks to define  $\Sigma_2$  and the corresponding rotation matrix:

$$\Sigma_2 := \begin{pmatrix} 1 & 0.8 & 0 & 0 & 0 \cdots 0 \\ 0.8 & 1 & 0 & 0 & 0 \cdots 0 \\ 0 & 0 & 1 & 0.8 & 0 \cdots 0 \\ 0 & 0 & 0.8 & 1 & 0 \cdots 0 \\ 0 & 0 & 0 & 0 & 1 \cdots 0 \\ \vdots & \vdots & \cdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 \cdots 1 \end{pmatrix}, \quad Q_2 := \begin{pmatrix} \cos \frac{3\pi}{4} & -\sin \frac{3\pi}{4} & 0 & 0 & 0 \cdots 0 \\ \sin \frac{3\pi}{4} & \cos \frac{3\pi}{4} & 0 & 0 & 0 \cdots 0 \\ 0 & 0 & \cos \frac{3\pi}{4} & -\sin \frac{3\pi}{4} & 0 \cdots 0 \\ 0 & 0 & \sin \frac{3\pi}{4} & \cos \frac{3\pi}{4} & 0 \cdots 0 \\ 0 & 0 & 0 & 0 & 1 \cdots 0 \\ \vdots & \vdots & \cdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 \cdots 1 \end{pmatrix}.$$

We can define  $\Sigma_3, \dots, \Sigma_{50}$  and the corresponding rotation matrices used for GSPS in the same way.

We run simulations for SPS for certain targets without using the true covariance matrices and we tune the acceptance rate of SPS to be 0.234 (or the closest possible). For comparison with the GSPS, we assume GSPS uses the correct covariance matrices of the targets, but the proposal variances  $h^2$  are chosen to be *the same* as those used for SPS. We then consider the performance of GSPS as the *benchmark*, since GSPS provides the optimal performance for SPS by using the information of the true target covariance matrix with the same proposal variance. We consider ACFs of two cases: the “latitude” and the 1-st coordinate. In Fig. 11, the first and third columns show the ACFs of the latitude and first coordinate of SPS for  $d = 100$  dimensional multivariate student’s  $t$  targets with different covariance matrices indexed by the number of sub-blocks (the proposal variances are tuned to target the 0.234 acceptance rate). For the same proposal variances, the second and fourth columns show the performance of the GSPS (the “benchmark”) using the true target covariance matrices. From Fig. 11, one can see that the SPS performs well when the number of sub-blocks is small. For the first and second rows, the acceptance rates of SPS cannot be as low as 0.234 and they equal 0.48 and 0.30, respectively. The ACF decreases slower as the lag when the number of sub-blocks increases. From the “benchmark” in the second and fourth columns, we know the proposal variance  $h^2$  becomes smaller. However, even in the case of  $\Sigma_{50}$  (the last row), the ACF still decreases much faster than RWM Fig. 9 for Gaussian targets. Note that in this example, all the targets are heavy-tailed and we know RWM is not even geometrically ergodic so the performance of RWM would be much worse than in Fig. 9. Overall, we can conclude that SPS is much better than RWM even when the covariance matrix is  $\Sigma_{50}$ . The performance can be significantly improved by GSPS even without changing the proposal variance. This suggests that adaptively tuning the parameters of GSPS (in the adaptive MCMC framework) is very promising.

**S9.6. The definition of Effective Sample Size (ESS) for SBPS.** We study the PDMP algorithms using the notion of Effective Sample Size (ESS) per Switch. We first recall the defini-

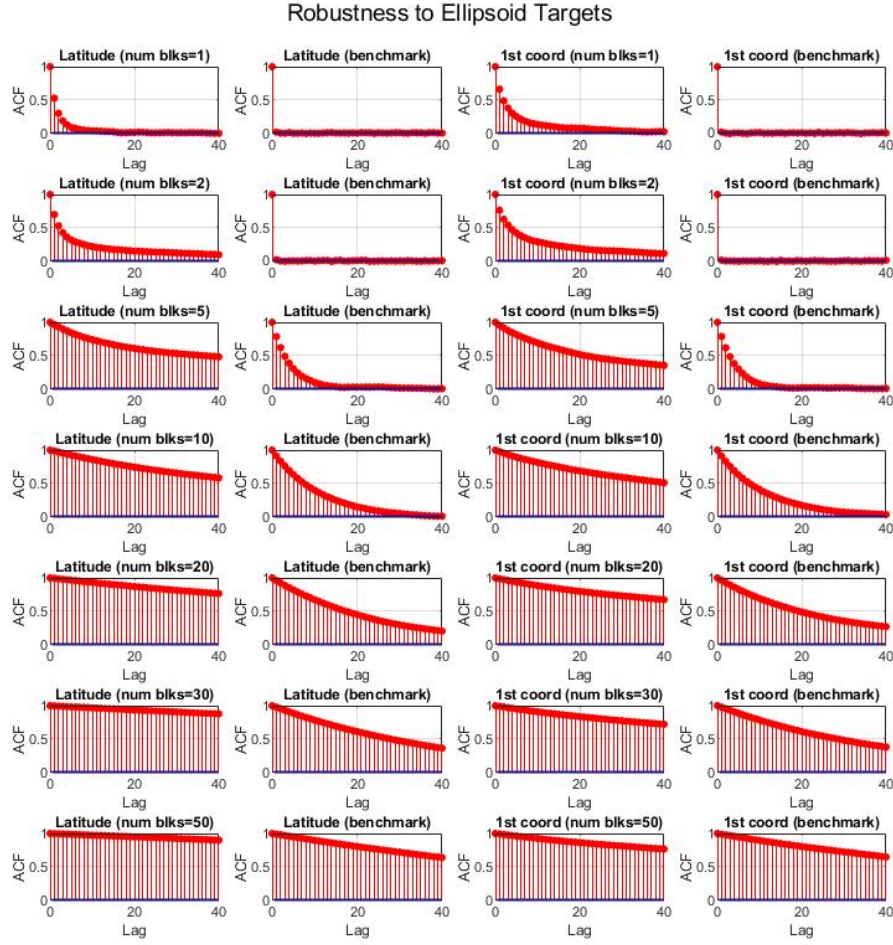


Figure 11: SPS versus GSPS: ACFs of latitude (first two columns) and 1-st coordinate (last two columns) for SPS (first and third columns) and GSPS (second and fourth columns). The proposal variances of SPS are chosen to target the 0.234 acceptance rate (for the first and second rows, the acceptance rates are 0.48 and 0.30, respectively). The proposal variances of GSPS are chosen to be the same as those for SPS to provide the “benchmark” performance. Target is 100-dim ellipsoid targets using  $2 \times 2$  blocks; GSPS (benchmark) uses the true covariance matrices.

tion of ESS in [Bierkens et al., 2019, supplemental material]. Consider a function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ , for a continuous process  $Y(t)$ , by CLT

$$\frac{1}{\sqrt{t}} \int_0^t [g(Y(s)) - \pi(h)] ds \rightarrow \sigma_g^2$$

where  $\sigma_h^2$  is called *asymptotic variance*.

The asymptotic variance can be estimated using batches of length  $T/B$

$$\widehat{\sigma}_g^2 = \frac{1}{B-1} \sum_{i=1}^B (X_i - \bar{X})^2,$$

where  $\bar{X} = \frac{1}{B} \sum_i X_i$  and

$$X_i := \sqrt{\frac{B}{T}} \int_{(i-1)T/B}^{iT/B} g(Y(s)) ds.$$

We also estimate the mean and variance under  $\pi$  by

$$\widehat{\pi}(g) := \frac{1}{T} \int_0^T g(Y(s)) ds, \quad \widehat{\text{var}}_{\pi}(g) := \frac{1}{T} \int_0^T g(Y(s))^2 ds - \left(\widehat{\pi}(g)\right)^2.$$

Then the ESS is estimated by

$$\widehat{\text{ESS}} := \frac{T \widehat{\text{var}}_{\pi}(g)}{\widehat{\sigma}_g^2}$$

The ESS per Switch is estimated by

$$\text{ESS per Switch} = \frac{\widehat{\text{ESS}}}{\text{Number of Events Simulated}}$$

We consider three cases for  $g$ , one is the first coordinate, the second one is the *negative log-target density*:

$$g(t) = \sqrt{d} \left( \frac{\|Y(t)\|^2}{d} - 1 \right) \sim \mathcal{N}(0, 2),$$

where the distribution holds for standard Gaussian targets. The third one is the square of the first coordinate.

**REMARK S9.1.** *If in a simulation of an event, there's only one evaluation of the full likelihood, then ESS per Switch is the same as ESS per epoch.*

**S9.7. SBPS: traceplot and ACF.** In this example, we study SBPS via the traceplots and ACFs for the 1-st coordinate, the negative log density, and the squared 1-st coordinate, respectively. In Fig. 12, the target is multivariate student's  $t$  distribution with  $d = 100$  degrees of freedom. Since SBPS is a continuous-time process, for the traceplot and ACF, we further discretize every unit period into 5 samples.  $N = 1000$  events are simulated with a low refresh rate of 0.2. For other settings such as high refresh rate, light-tailed targets, and different covariance matrices, we will study in later subsections for additional simulations for SBPS.

There are several interesting observations from Fig. 12. The ACFs for all three cases have certain periodic behaviors and can take negative values for the first two cases. We also include the ESS per Switch for the three cases. As a result of negative ACFs, the ESS per Switch is larger than 1 in the first two cases. This suggests asymptotic variance for estimating the 1-st coordinate or the negative log density using SBPS is even smaller than the variance using  $N$  independent samples from the target.

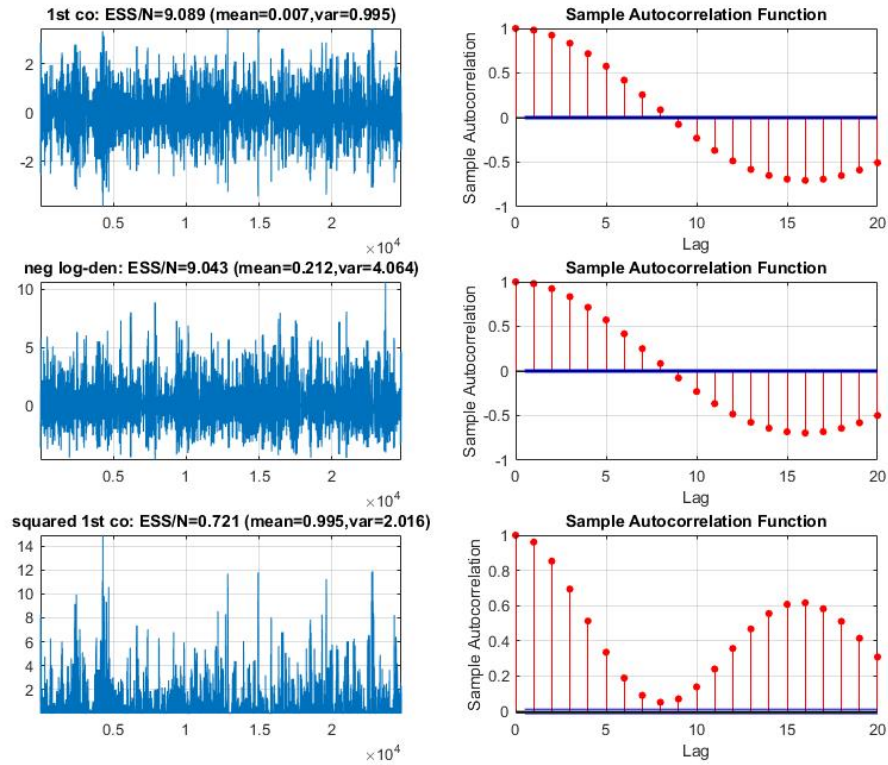
Trace Plots and ACFs: Refresh Rate=0.2 (random initial,  $d=100$ ,  $N=1000$ ,  $\text{DoF}/d=1.0$ ,  $\text{discretize}=0.2$ )

Figure 12: Trace plots and ACFs for the 1-st coordinate, the negative log-density, and the squared 1-st coordinate. Target distribution is multivariate student's  $t$  distribution with  $\text{DoF} = d = 100$ . Every unit period is discretized into 5 samples.  $N = 1000$  events are simulated. Low refresh rate = 0.2. The ESS per switch is larger than 1 and the ACF can be negative in the first two cases.

S9.8. *SBPS: traceplots and ACFs (cont.)*. In Fig. 12, we have studied SBPS by the traceplots and ACFs for the multivariate student's  $t$  target with  $\text{DoF} = d = 100$  for low refresh rate. In this example, we study other settings. Fig. 13 shows the traceplots and ACFs for multivariate student's  $t$  target with large refresh rate), Fig. 14 and Fig. 15 are those for standard Gaussian target (with  $d = 100$ ) using small and large refresh rates, respectively. We consider three cases: the 1-st coordinate, the negative log density, and the squared 1-st coordinate, respectively. Since SBPS is a continuous time process, for the traceplot and ACF, we further discretize every unit of time into 5 samples.  $N = 1000$  events are simulated. A low refresh rate corresponds to 0.2 and a high refresh rate corresponds to 2.

From Fig. 14, there are several similar interesting observations as in Fig. 12. The ACFs for all three cases have certain periodic behaviors and can take negative values for the first two cases. As a result of negative ACFs, the ESS per Switch is larger than 1 in the first two cases. For high refresh rate cases such as in Fig. 13 and Fig. 15, the periodic behavior and negativity of ACFs disappear and the corresponding ESS per Switch is smaller than 1, (note that the performance for SBPS is always significantly better than the traditional BPS in the

simulations). This suggests that a larger proportion of bounce events are helpful since bounce events use the information of the gradients of log-target density. As the refreshment events for SBPS are very efficient in the transient phase, the refresh rate can be chosen to be very low. This is not the case for BPS. The refreshments in the transient phase for BPS are inefficient in high dimensions.

Trace Plots and ACFs: Refresh Rate=2.0 (random initial,  $d=100$ ,  $N=1000$ ,  $\text{DoF}/d=1.0$ ,  $\text{discretize}=0.2$ )

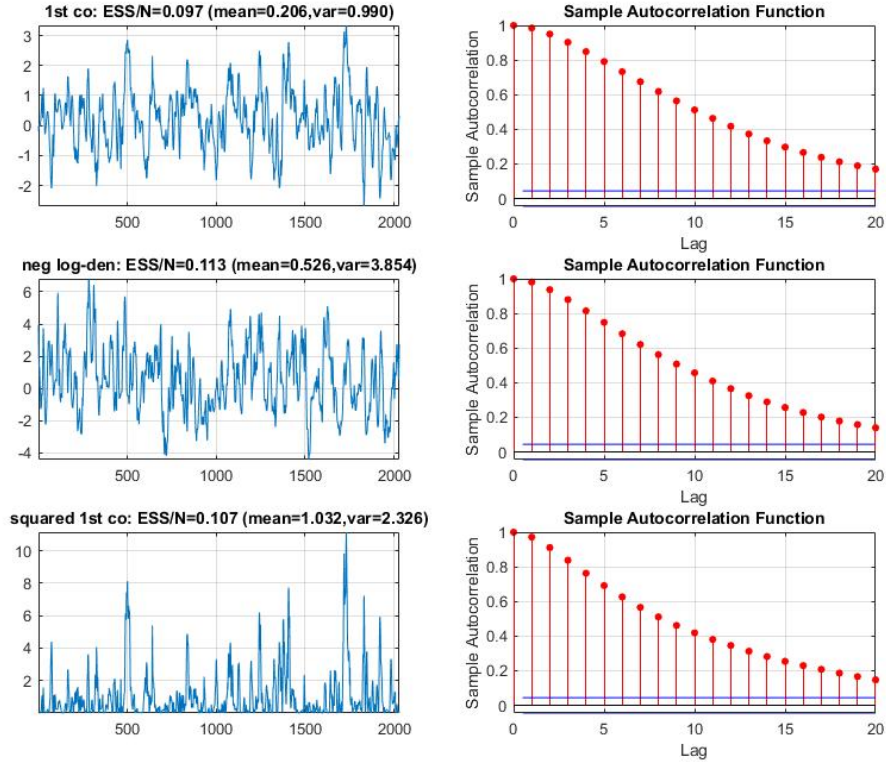


Figure 13: Trace plots and ACFs for the 1-st coordinate, the negative log-density, and the squared 1-st coordinate. Target distribution is multivariate student's  $t$  with  $\text{DoF} = d = 100$ . Every unit period is discretized into 5 samples.  $N = 1000$  events are simulated. High refresh rate = 2.

**S9.9. SBPS: ESS per Switch for multivariate student's  $t$  target.** We have studied the efficiency curves of SBPS and BPS in terms of ESS per Switch versus the refresh rate for the Gaussian target in Fig. 6. In this example, we consider the case of the multivariate student's  $t$  target with  $\text{DoF}$  equals  $d$ . The first subplot of Fig. 16 contains the proportion of refreshments in all the  $N$  events for varying refresh rates. Note that there are no bounce events for the SBPS so all the events for SBPS are refreshments. In the other three subplots of Fig. 16, we plot the logarithm of ESS per Switch as a function of the refresh rate for three cases, the 1-st coordinate, the negative log-density, and the squared 1-st coordinate. For each efficiency curve,  $N = 1000$  events are simulated, random initial value for SBPS and BPS starts from

Trace Plots and ACFs: Refresh Rate=0.2 (random initial,  $d=100$ ,  $N=1000$ , DoF/ $d=\text{Inf}$ , discretize=0.2)

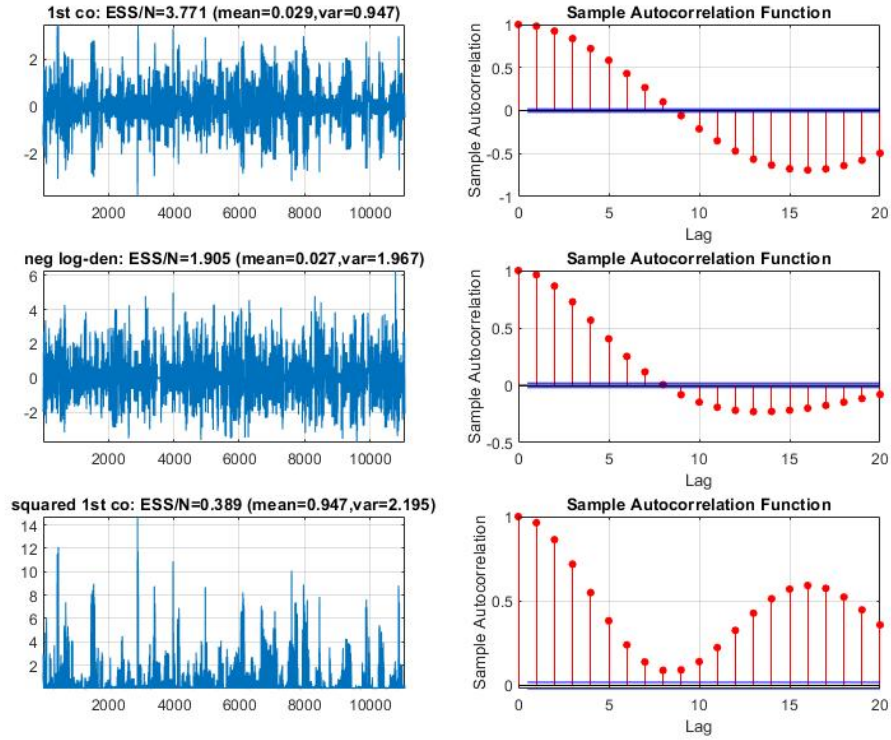


Figure 14: Trace plots and ACFs for the 1-st coordinate, the negative log-density, and the squared 1-st coordinate. The target distribution is standard Gaussian with  $d = 100$ . Every unit period is discretized into 5 samples.  $N = 1000$  events are simulated. Low refresh rate = 0.2.

stationarity. As SBPS and BPS are continuous-time processes, each unit time is discretized into 5 samples.

According to Fig. 16, the ESS per Switch of SBPS is much larger than the ESS per Switch of BPS for all cases (actually the gap is larger than the cases for Gaussian target, and the gap becomes larger in higher dimensions). For all three cases, the ESS per Switch of SBPS can be larger than 1 if the refresh rate is relatively low, this is also better than the cases for the Gaussian target. For BPS, however, even starting from stationarity, the ESS per Switch is always much smaller than 1.

S9.10. *SBPS: robustness to target variance.* Previously, in Fig. 14 and Fig. 15, we have studied the traceplots and ACFs of SBPS for standard Gaussian target with  $d = 100$  using small and large refresh rates, respectively. In this example, we repeat the same numerical experiments except that, instead of considering standard Gaussian targets, we consider Gaussian targets with larger or smaller variances. Since we still choose  $R = \sqrt{d}$ , this is equivalent to studying the robustness of SBPS to the choice of the radius of the sphere. Other than the target variances, numerical experiment settings are the same as in Fig. 14 and Fig. 15.



Trace Plots and ACFs: Refresh Rate=2.0 (random initial,  $d=100$ ,  $N=1000$ ,  $\text{DoF}/d=\text{Inf}$ ,  $\text{discretize}=0.2$ )

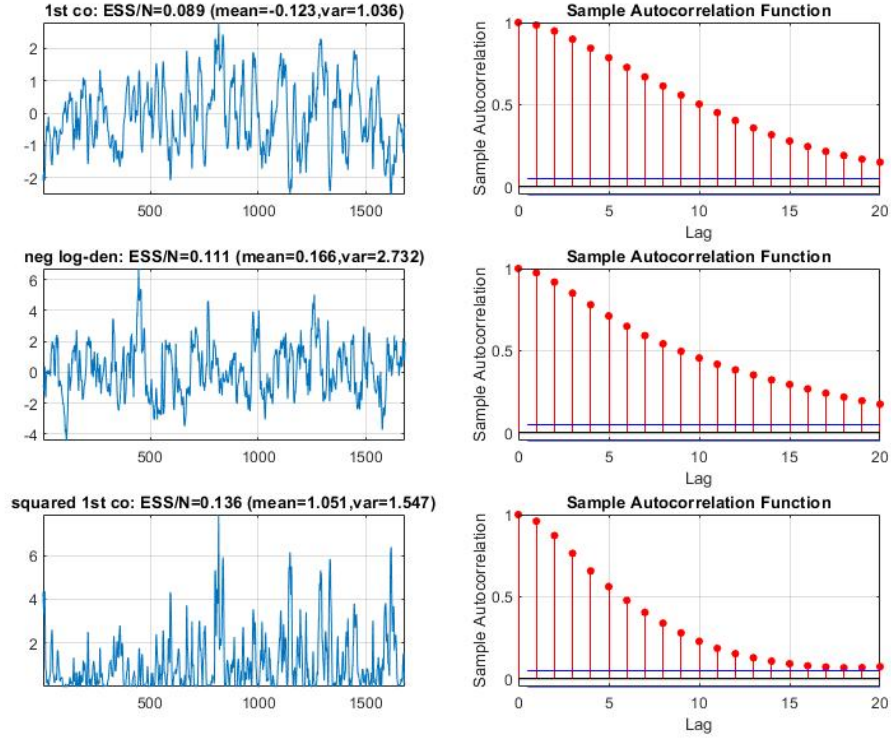


Figure 15: Trace plots and ACFs for the 1-st coordinate, the negative log-density, and the squared 1-st coordinate. The target distribution is standard Gaussian with  $d = 100$ . Every unit period is discretized into 5 samples.  $N = 1000$  events are simulated. High refresh rate = 2.

We consider two cases, Gaussian target  $\mathcal{N}(0, 1.5I_d)$  (see Fig. 17 for small refresh rate and Fig. 18 for large refresh rate) and Gaussian target  $\mathcal{N}(0, 0.7I_d)$  (see Fig. 19 for small refresh rate and Fig. 20 for large refresh rate). Comparing with the cases when the variance is 1 in Fig. 14 and Fig. 15, we can see that the traceplots and ACFs for the first coordinate and squared first coordinate are not affected too much. For the negative log-target density, the traceplots and ACFs behave quite differently. When the target covariance matrix is  $1.5I_d$ , more bounce events occur when the SBPS is moving to the equator. As a result, the SBPS stays almost all the time in the “Northern Hemisphere”. When the target covariance matrix is  $0.7I_d$ , the situation is exactly the opposite: the SBPS stays almost all the time in the “Southern Hemisphere”. The traceplots and ACFs in all four figures show that the performance of SBPS is quite robust to the target variance.

### Efficiency Curves: ESS/N vs refresh rate (Student-T( $d$ ) target)

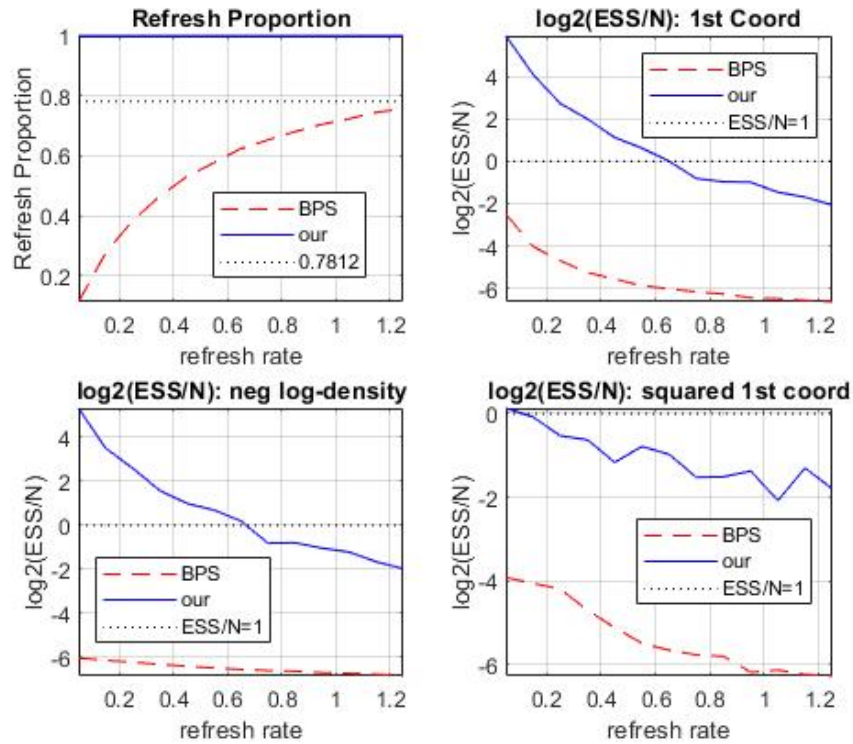


Figure 16: Efficiency: ESS per Switch for SBPS and BPS for varying refresh rate.  $N = 1000$  events are simulated, Random initial values for SBPS and BPS start from stationarity. Each unit time is discretized to 5 samples. Target distribution: Multivariate student's  $t$  with DoF =  $d = 100$ . The first subplot is the proportion of refreshment events in all  $N$  events. The other three subplots are ESS for the 1-st coordinate, the negative log density, and the squared 1-st coordinate, respectively.

Trace Plots and ACFs: Gaussian(0.0,1.5) (refresh=0.2 , d=100, N=1000, discretize=0.2)

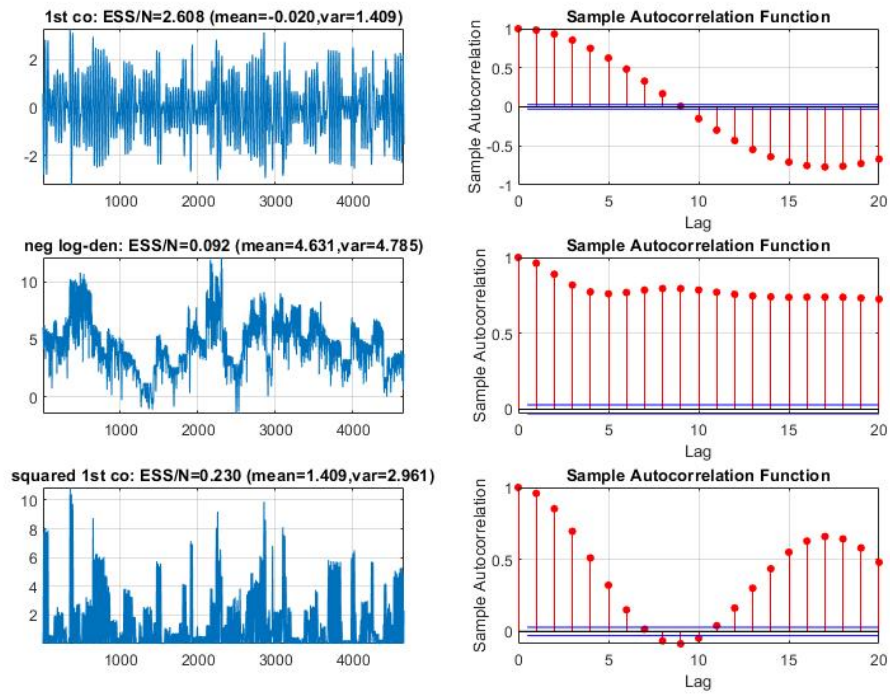


Figure 17: Trace plots and ACFs for the 1-st coordinate, the negative log-density, and the squared 1-st coordinate. Target distribution is  $\mathcal{N}(0, 1.5I_d)$  with  $d = 100$  and we choose  $R = d^{1/2}$ . Every unit period is discretized into 5 samples.  $N = 1000$  events are simulated. Low refresh rate = 0.2.

Trace Plots and ACFs: Gaussian(0.0,1.5) (refresh=2.0 , d=100, N=1000, discretize=0.2)

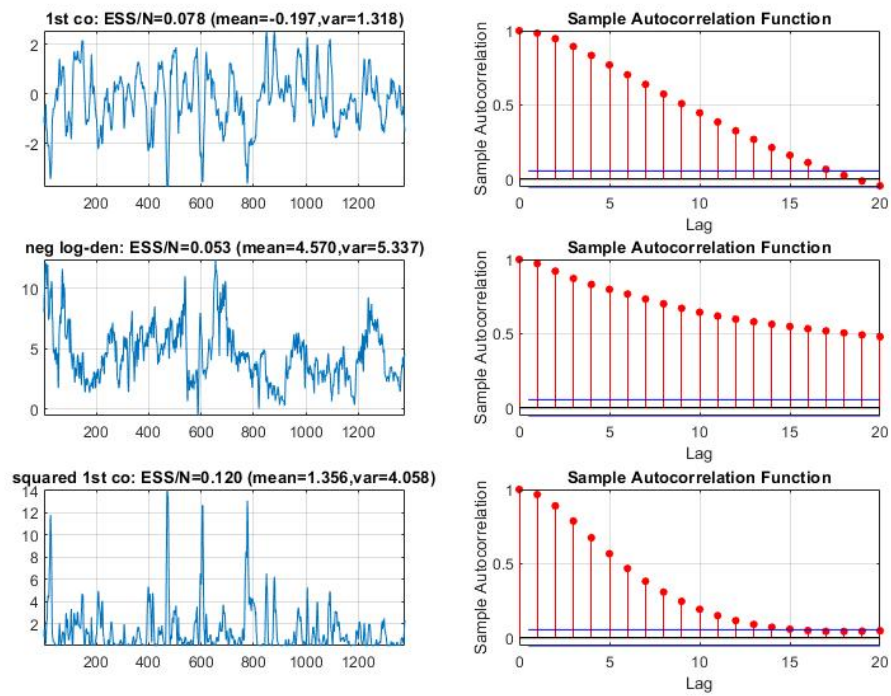


Figure 18: Trace plots and ACFs for the 1-st coordinate, the negative log-density, and the squared 1-st coordinate. Target distribution is  $\mathcal{N}(0, 1.5I_d)$  with  $d = 100$  and we choose  $R = d^{1/2}$ . Every unit time is discretized into 5 samples.  $N = 1000$  events are simulated. High refresh rate = 2.

Trace Plots and ACFs: Gaussian(0.0,0.7) (refresh=0.2 , d=100, N=1000, discretize=0.2)

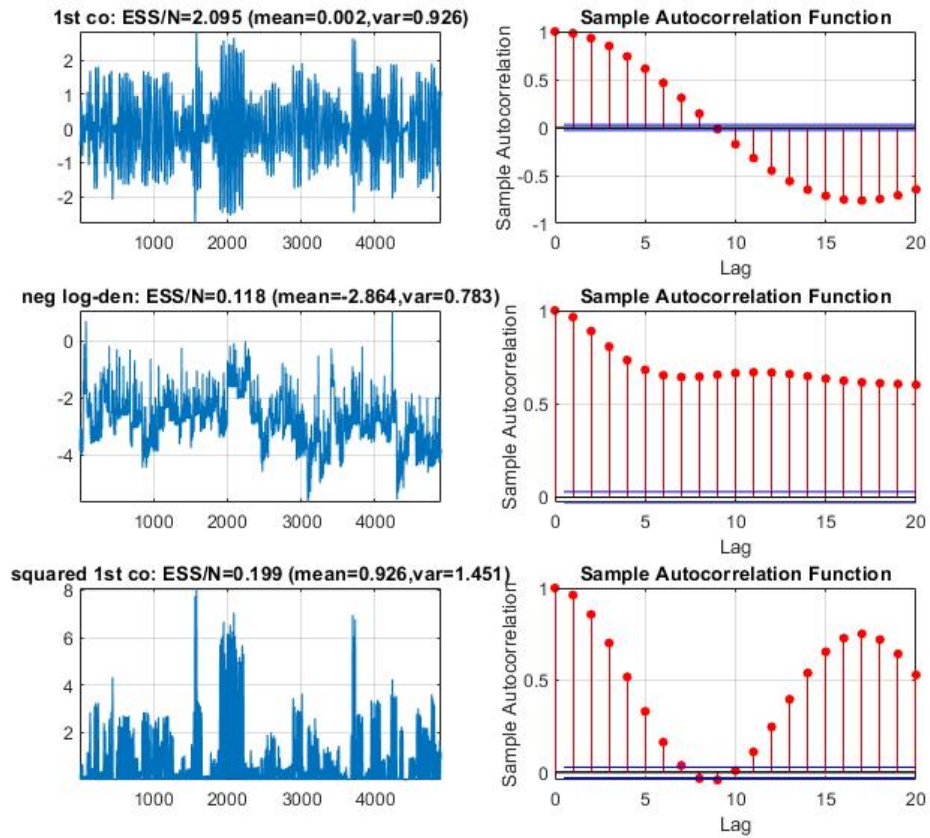


Figure 19: Trace plots and ACFs for the 1-st coordinate, the negative log-density, and the squared 1-st coordinate. Target distribution is  $\mathcal{N}(0, 0.7I_d)$  with  $d = 100$  and we choose  $R = d^{1/2}$ . Every unit of time is discretized into 5 samples.  $N = 1000$  events are simulated. Low refresh rate = 0.2.

Trace Plots and ACFs: Gaussian(0.0,0.7) (refresh=2.0 , d=100, N=1000, discretize=0.2)

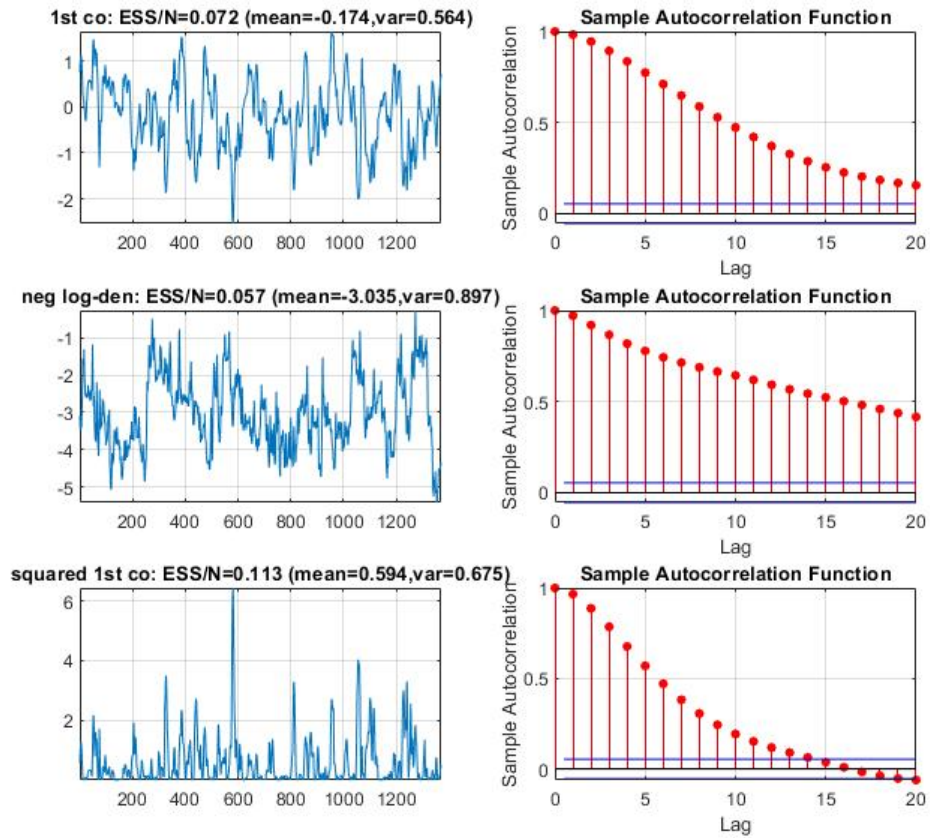


Figure 20: Trace plots and ACFs for the 1-st coordinate, the negative log-density, and the squared 1-st coordinate. Target distribution is  $\mathcal{N}(0, 0.7I_d)$  with  $d = 100$  and we choose  $R = d^{1/2}$ . Every unit period is discretized into 5 samples.  $N = 1000$  events are simulated. High refresh rate = 2.

## REFERENCES

- Christophe Andrieu, Paul Dobson, and Andi Q Wang. Subgeometric hypocoercivity for piecewise-deterministic markov process monte carlo methods. *Electronic Journal of Probability*, 26:1–26, 2021a.
- Christophe Andrieu, Alain Durmus, Nikolas Nüsken, and Julien Roussel. Hypocoercivity of piecewise deterministic markov process-monte carlo. *The Annals of Applied Probability*, 31(5):2478–2517, 2021b.
- Etienne P Bernard, Werner Krauth, and David B Wilson. Event-chain monte carlo algorithms for hard-sphere systems. *Physical Review E*, 80(5):056704, 2009.
- Joris Bierkens, Paul Fearnhead, and Gareth Roberts. The zig-zag process and super-efficient sampling for bayesian analysis of big data. *The Annals of Statistics*, 47(3):1288–1320, 2019.
- Joris Bierkens, Kengo Kamatani, and Gareth O Roberts. High-dimensional scaling limits of piecewise deterministic sampling algorithms. *The Annals of Applied Probability*, 32(5):3361–3407, 2022.
- Alexandre Bouchard-Côté, Sebastian J Vollmer, and Arnaud Doucet. The bouncy particle sampler: A nonreversible rejection-free markov chain monte carlo method. *Journal of the American Statistical Association*, 113(522):855–867, 2018.
- Ole F Christensen, Gareth O Roberts, and Jeffrey S Rosenthal. Scaling limits for the transient phase of local Metropolis–Hastings algorithms. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):253–268, 2005.
- Simon L Cotter, Gareth O Roberts, Andrew M Stuart, and David White. Mcmc methods for functions: modifying old algorithms to make them faster. *Statistical Science*, 28(3):424–446, 2013.
- Harold Scott Macdonald Coxeter. *Introduction to geometry*. 1961.
- M. H. A. Davis. Piecewise-deterministic Markov processes: a general class of nondiffusion stochastic models. *J. Roy. Statist. Soc. Ser. B*, 46(3):353–388, 1984. ISSN 0035-9246. With discussion.
- George Deligiannidis, Alexandre Bouchard-Côté, and Arnaud Doucet. Exponential ergodicity of the bouncy particle sampler. *The Annals of Statistics*, 47(3):1268–1287, 2019.
- George Deligiannidis, Daniel Paulin, Alexandre Bouchard-Côté, and Arnaud Doucet. Randomized hamiltonian monte carlo as scaling limit of the bouncy particle sampler and dimension-free convergence rates. *The Annals of Applied Probability*, 31(6):2612–2662, 2021.
- Douglas Down, Sean P Meyn, and Richard L Tweedie. Exponential and uniform ergodicity of markov processes. *The Annals of Probability*, 23(4):1671–1691, 1995.
- Alain Durmus, Arnaud Guillin, and Pierre Monmarché. Geometric ergodicity of the bouncy particle sampler. *The Annals of Applied Probability*, 30(5):2069–2098, 2020.
- Stewart N. Ethier and Thomas G. Kurtz. *Markov processes: Characterization and Convergence*. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics. Wiley, New York, 1986. ISBN 0-471-08186-8. Characterization and convergence.
- Paul Fearnhead, Joris Bierkens, Murray Pollock, Gareth O Roberts, et al. Piecewise deterministic markov processes for continuous-time monte carlo. *Statistical Science*, 33(3):386–412, 2018.
- Frederick W Gehring, Gaven J Martin, and Bruce P Palka. *An introduction to the theory of higher-dimensional quasiconformal mappings*, volume 216. American Mathematical Soc., 2017.
- Mark Girolami and Ben Calderhead. Riemann manifold langevin and hamiltonian monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(2):123–214, 2011.
- Søren F Jarner and Ernst Hansen. Geometric ergodicity of metropolis algorithms. *Stochastic processes and their applications*, 85(2):341–361, 2000.
- Søren F Jarner and Gareth O Roberts. Convergence of heavy-tailed monte carlo markov chain algorithms. *Scandinavian Journal of Statistics*, 34(4):781–815, 2007.

- Leif T Johnson and Charles J Geyer. Variable transformation to obtain geometric ergodicity in the random-walk metropolis algorithm. *The Annals of Statistics*, pages 3050–3076, 2012.
- Kengo Kamatani. Efficient strategy for the markov chain monte carlo in high-dimension with heavy-tailed target probability distribution. *Bernoulli*, 24(4B):3711–3750, 2018.
- Han Cheng Lie, Daniel Rudolf, Björn Sprungk, and Timothy John Sullivan. Dimension-independent markov chain monte carlo on the sphere. *Scandinavian Journal of Statistics*, 2021.
- Thomas Milton Liggett. *Continuous time Markov processes: an introduction*, volume 113. American Mathematical Soc., 2010.
- Oren Mangoubi and Aaron Smith. Rapid mixing of geodesic walks on manifolds with positive curvature. *The Annals of Applied Probability*, 28(4):2501–2543, 2018.
- Kerrie L Mengersen and Richard L Tweedie. Rates of convergence of the Hastings and Metropolis algorithms. *The Annals of Statistics*, 24(1):101–121, 1996.
- Sean P Meyn and Richard L Tweedie. *Markov chains and stochastic stability*. Springer Science & Business Media, 2012.
- Aleksandar Mijatović, Veno Mramor, and Gerónimo Uribe Bravo. Projections of spherical brownian motion. *Electronic Communications in Probability*, 23:1–12, 2018.
- Yann Ollivier. Ricci curvature of markov chains on metric spaces. *Journal of Functional Analysis*, 256(3):810–864, 2009.
- Emilia Pompe, Chris Holmes, and Krzysztof Łatuszyński. A framework for adaptive mcmc targeting multimodal distributions. *The Annals of Statistics*, 48(5):2930–2952, 2020.
- Gareth O Roberts and Jeffrey S Rosenthal. Examples of adaptive mcmc. *Journal of computational and graphical statistics*, 18(2):349–367, 2009.
- Gareth O Roberts and Adrian FM Smith. Simple conditions for the convergence of the gibbs sampler and metropolis-hastings algorithms. *Stochastic processes and their applications*, 49(2):207–216, 1994.
- Gareth O Roberts and Richard L Tweedie. Geometric convergence and central limit theorems for multidimensional hastings and Metropolis algorithms. *Biometrika*, 83(1):95–110, 1996a.
- Gareth O Roberts and Richard L Tweedie. Exponential convergence of langevin distributions and their discrete approximations. *Bernoulli*, pages 341–363, 1996b.
- Gareth O Roberts, Andrew Gelman, and Walter R Gilks. Weak convergence and optimal scaling of random walk Metropolis algorithms. *The Annals of Applied Probability*, 7(1):110–120, 1997.
- Chris Sherlock, Paul Fearnhead, and Gareth O Roberts. The random walk metropolis: linking theory and practice through a case study. *Statistical Science*, 25(2):172–190, 2010.
- Luke Tierney. Markov chains for exploring posterior distributions. *The Annals of Statistics*, 22(4):1701–1728, 1994.
- Giorgos Vasdekis and Gareth O Roberts. A note on the polynomial ergodicity of the one-dimensional zig-zag process. *Journal of Applied Probability*, 59(3):895–903, 2022.
- Giorgos Vasdekis and Gareth O Roberts. Speed up zig-zag. *The Annals of Applied Probability*, 33(6A):4693–4746, 2023.
- Jun Yang and Jeffrey S Rosenthal. Complexity results for mcmc derived from quantitative bounds. *The Annals of Applied Probability*, 33(2):1459–1500, 2023.
- Jun Yang, Gareth O Roberts, and Jeffrey S Rosenthal. Optimal scaling of random-walk Metropolis algorithms on general target distributions. *Stochastic Processes and their Applications*, 130(10):6094–6132, 2020.
- Emilio Zappa, Miranda Holmes-Cerfon, and Jonathan Goodman. Monte carlo on manifolds: sampling densities and integrating functions. *Communications on Pure and Applied Mathematics*, 71(12):2609–2647, 2018.