

---

# Grassmann Stein Variational Gradient Descent

---

**Xing Liu**

Imperial College London

**Harrison Zhu**

Imperial College London

**Jean-François Ton**

University of Oxford

**George Wynne**

Imperial College London

**Andrew Duncan**

Imperial College London

## Abstract

Stein variational gradient descent (SVGD) is a deterministic particle inference algorithm that provides an efficient alternative to Markov chain Monte Carlo. However, SVGD has been found to suffer from variance underestimation when the dimensionality of the target distribution is high. Recent developments have advocated projecting both the score function and the data onto real lines to sidestep this issue, although this can severely overestimate the epistemic (model) uncertainty. In this work, we propose *Grassmann Stein variational gradient descent* (GSVGD) as an alternative approach, which permits projections onto arbitrary dimensional subspaces. Compared with other variants of SVGD that rely on dimensionality reduction, GSVGD updates the projectors simultaneously for the score function and the data, and the optimal projectors are determined through a coupled Grassmann-valued diffusion process which explores favourable subspaces. Both our theoretical and experimental results suggest that GSVGD enjoys efficient state-space exploration in high-dimensional problems that have an intrinsic low-dimensional structure.

## 1 INTRODUCTION

Variational inference (VI) (Blei et al., 2017) is an optimisation-centric framework for approximating complex distributions that are intractable (only known

up to a scale factor): given any target distribution  $p(x)$ , VI searches over a user-defined class of distributions  $\mathcal{Q}$  for an optimal  $q(x) \in \mathcal{Q}$  that is closest (in terms of a discrepancy or divergence) to  $p(x)$ . Extensively applied in the field of Bayesian inference, VI has the attraction of being scalable to big datasets, although it leads to biased estimation unless  $\mathcal{Q}$  is broad enough to include the target distribution. This is in contrast to Markov chain Monte Carlo (MCMC) algorithms (Gilks et al., 1995), which allows asymptotically exact sampling from the true target distribution, but does not scale well to big datasets and high dimensionality due to long mixing times (Levin and Peres, 2017).

The choice of  $\mathcal{Q}$  in VI is crucial to guarantee a good approximation to  $p(x)$  while retaining computational tractability. The classical mean-field approximation, which assumes that the distributions in  $\mathcal{Q}$  have independent marginals, can be overly simplistic in many cases. To address this, a growing line of work in VI with Normalising Flows (NFs) (Rezende and Mohamed, 2015) seeks to construct an invertible map  $T$  so that the pushforward distribution of  $T(x)$ , with  $x \sim q$  and  $q \in \mathcal{Q}$ , will form a flexible approximation to  $p$ .

Stein variational gradient descent (SVGD) (Liu et al., 2016) was introduced as a particle-based variational inference method in which the map  $T$  seeks to move a set of particles along a vector field which is chosen within a reproducing kernel Hilbert space (RKHS) to optimally transport towards the target  $p(x)$ . Building on the kernel Stein discrepancy (Liu et al., 2016), the SVGD transport both drives the particles to the high-probability regions of  $p(x)$  and enforces repulsion between the particles to prevent mode-collapse. It hence has the advantage of being “particle-efficient” in that a small number of particles can achieve good approximation of  $p(x)$ . Although this has made SVGD a popular tool in a range of applications including meta-learning (Yoon et al., 2018) and learning diversified mixture models (Wang and Liu, 2019), it is found that the marginal variance of the resulting

---

Proceedings of the 25<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2022, Valencia, Spain. PMLR: Volume 151. Copyright 2022 by the author(s).

particles scales inversely with the dimension, resulting in under-estimation of the variance. This was studied analytically in Ba et al. (2022), in which the authors attributes this issue to (i) an uneven scale of variance for the two terms in the SVGD update, and (ii) a deterministic bias due to the dependence of the particle positions at each step.

Recently, Gong et al. (2021) proposed a *sliced* version of kernel Stein discrepancy (KSD) analogous to other sliced discrepancies (Kolouri et al., 2019), and an associated *sliced SVGD* (S-SVGD). This variant of SVGD decomposes the kernel-based dynamics along a sequence of 1-dimensional projections, called *slices*. Empirical results (Gong et al., 2021, Appendix J.5, Table 5) showed that the S-SVGD dynamic can mitigate the variance under-estimation issue of SVGD in high dimensions. However, S-SVGD is constrained to 1-dimensional projection, and only seeks the optimal slices for the data but keeps the slices for the score function deterministic. As we demonstrate in this paper, these can lead to inflated variance of the S-SVGD estimation.

In this work, we introduce a unified approach where a non-uniform probability distribution over the projectors for *both the score function and the data* are adaptively updated to emphasise directions in which there is largest discrepancy. This is achieved by introducing a stochastic diffusion process taking values in the Grassmann manifold (Bendokat et al., 2020) which evolves along with SVGD particles. By tuning the diffusion process parameters we can adjust the trade-off between exploration and exploitation of suitable projections. In addition, using projections onto higher dimensional spaces allows us to take into account the correlations/interactions between components and thus producing more accurate uncertainty estimates.

**Contributions:** Motivated by S-SVGD, we propose a novel algorithm, Grassmann SVGD (GSVGD), which employs  $m$ -dimensional projectors (compared with the 1-dimensional ones in S-SVGD) to evolve the particles towards the target distribution and to mitigate variance under-estimation. We show numerically that our method is competitive to SVGD and S-SVGD on high-dimensional synthetic and benchmark problems while more accurately estimating the epistemic uncertainty.

## 2 BACKGROUND

**Stein Variational Gradient Descent:** Let  $\mathcal{X} = \mathbb{R}^d$  and  $P$  be a probability measure over  $\mathcal{X}$  with smooth positive density  $p$  which we can evaluate up to a normalisation constant. We are interested in approximating  $P$  by transporting a known measure  $Q$ , defined over  $\mathcal{X}$  and with smooth density  $q$ , to  $P$  via



Figure 1: A summary of different SVGD algorithms using no projections (SVGD), 1-dimensional projections (S-SVGD), and  $m$ -dimensional projections (GSVGD) with an arbitrary  $1 \leq m \leq d$ , where  $d$  is the dimensionality of the problem.

a sequence of maps that minimise a given loss.

Choosing the loss to be the Kullback-Leibler (KL) divergence defined as  $\text{KL}(Q, P) = \int_{\mathcal{X}} q(x) \log(q(x)/p(x)) dx$  allows us to minimise the discrepancy  $\text{KL}(T_{\#}Q, P)$  over maps  $T: \mathcal{X} \rightarrow \mathcal{X}$  where  $T_{\#}Q$  is the pushforward of  $Q$  with respect to  $T$ . In Liu and Wang (2016), the authors choose a specific parametrisation  $T(x) = x + \varepsilon\phi(x)$  where  $\phi$  lies within the unit ball  $\mathcal{B}_k^d := \{\phi \in \mathcal{H}_k^d : \|\phi\|_{\mathcal{H}_k^d} \leq 1\}$  of the Cartesian product  $\mathcal{H}_k^d := \times_{i=1}^d \mathcal{H}_k$  of the RKHS  $\mathcal{H}_k$  associated with kernel  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ . The crucial observation of Liu and Wang (2016) is that the maximal rate of decay of this discrepancy with respect to  $\varepsilon$  is the kernel Stein discrepancy (KSD) (Chwialkowski et al., 2016; Liu et al., 2016) given by

$$\text{KSD}(Q, P) = \sup_{\phi \in \mathcal{B}_k^d} \mathbb{E}_Q[\mathcal{A}_p\phi(x)], \quad (1)$$

where  $x \sim Q$ ,  $s_p(x) = \nabla_x \log p(x)$  is the *score function* and  $\mathcal{A}_p\phi(x) = s_p(x)^\top \phi(x) + \nabla \cdot \phi(x)$  is the *Stein operator* associated to  $P$ . The key advantage of KSD is that it quantifies the difference between two measures  $P$  and  $Q$  where only the score function of  $P$  is available, which is the typical situation in Bayesian inference for complex models. The supremum in (1) is attained by the function  $\phi^*(\cdot) := \mathbb{E}_Q[\mathcal{A}_p k(\cdot, x)] = \mathbb{E}_Q[k(\cdot, x)s_p(x) + \nabla_x k(\cdot, x)]$ , which leads to the remarkable result (Liu and Wang, 2016, Theorem 3.3)  $\nabla_\phi \text{KL}(T_{\#}Q, P)|_{\phi=0} = -\phi^*(x)$  where  $\nabla_\phi$  is the functional derivative and  $T(x)$  is the map parameterised as above. This suggests that  $\phi^*$  is, according to KL divergence, the best choice of function to use in  $T$  to map  $Q$  to  $P$ .

The expression for the optimal  $\phi^*$  forms the basis of the SVGD algorithm. Particles  $X^t = (x_1^t, \dots, x_N^t)$  at step  $t$ , are evolved using the following update rule

$$x_i^{t+1} = x_i^t + \varepsilon \widehat{\phi}^*(x_i^t), \quad (2)$$

for  $i = 1, \dots, N$  and  $t = 0, 1, 2, \dots$ , where

$$\widehat{\phi}^*(\cdot) = \frac{1}{N} \sum_{j=1}^N [k(\cdot, x_j^t) \nabla_{x_j^t} \log p(x_j^t) + \nabla_{x_j^t} k(\cdot, x_j^t)] \quad (3)$$

is an estimator of  $\phi^*(\cdot)$  obtained by replacing  $Q$  with  $\frac{1}{N} \sum_{i=1}^N \delta_{x_i^t}$ . Intuitively, (2) is an Euler discretisation with step-size  $\epsilon$  of the following system of ODEs

$$\begin{aligned} \frac{dx_i}{dt}(t) = & \frac{1}{N} \sum_{j=1}^N [k(x_i(t), x_j(t)) \nabla_{x_j(t)} \log p(x_j(t)) \\ & + \nabla_{x_j(t)} k(x_i(t), x_j(t))], \quad i = 1, \dots, N, \end{aligned}$$

for particle positions  $x_1(t), \dots, x_N(t)$ .

**Tackling the Curse of Dimensionality:** The SVGD algorithm has shown to be asymptotically exact, in the sense that in the limit of infinite particles and vanishing step-size the particles will converge to the target density (Lu et al., 2019). However, in the finite particle regime, SVGD may fail to adequately explore the state-space, exhibiting mode collapse and variance under-estimation, particularly in high dimensions (Liu et al., 2016; Zhuo et al., 2018; Gong et al., 2021).

The deterioration in performance in high-dimensions can be understood from the update rule defined in (3). The RHS consists of two terms: a kernel-averaged score function  $\frac{1}{N} \sum_j k(\cdot, x_j^t) \nabla \log p(x_j^t)$  which attracts particles towards modes of the target density, and a repulsion term  $\frac{1}{N} \sum_j \nabla k(\cdot, x_j^t)$  which encourages particle diversity by pushing nearby particles away from each other. As has been observed in Liu and Wang (2016, Section 3.2) and studied analytically in (Zhuo et al., 2018, Section 3), the influence of the repulsion term drops dramatically with increasing dimension, effectively reducing SVGD to a gradient ascent method for  $\log p$ . As a result, the estimation of SVGD can suffer from a significant under-estimated variance for high-dimensional targets.

Various developments have been proposed to mitigate this issue; we now discuss the most relevant to our work, the sliced approach of Gong et al. (2021), and defer a review of the others to Section 5. Gong et al. (2021) addressed this problem by projecting both the particles and the score function along one-dimensional directions, known as *slices*, leading to the *max sliced kernel Stein discrepancy* (maxSKSD). This approach hinges on the fact that measuring KSD along finitely many slices is sufficient to capture all geometric information of two distributions, as long as optimal slices are used. Given a user-defined orthonormal basis  $O$  in  $\mathbb{R}^d$ , the *maxSKSD* takes the form

$$\begin{aligned} \text{maxSKSD}(Q, P) = & \sum_{r \in O} \sup_{g_r \in \mathbb{S}^{d-1}} \sup_{\phi \in \mathcal{B}_k} \mathbb{E}_Q [s_p^r(x) \phi(x^\top g_r) \\ & + r^\top g_r \nabla_{x^\top g_r} \phi(x^\top g_r)], \quad (4) \\ = & \sum_{r \in O} \sup_{g_r \in \mathbb{S}^{d-1}} \sup_{\phi \in \mathcal{B}_{k_r, g_r}} \mathbb{E}_Q [s_p(x)^\top \phi(x) + \nabla \cdot \phi(x)] \quad (5) \end{aligned}$$

where  $s_p^r(x) := s_p(x)^\top r$  is the score projected along a slicing direction  $r$  and  $\mathcal{B}_k$  is the unit ball of from

RKHS  $\mathcal{H}_k$  for some kernel  $k: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ . The subscript  $g_r$  is used to emphasise that for each  $r$  a new optimal  $g_r$  must be found. In (5),  $\mathcal{B}_{k_r, g_r}$  is the unit ball of the RKHS corresponding to the matrix valued kernel  $k_{r, g_r}(x, y) = r r^\top k(g_r^\top x, g_r^\top y)$ . This can be seen from standard results regarding matrix valued kernels (Paulsen and Raghupathi, 2016, Chapter 6) and this will be relevant to our later developments.

The idea of the maxSKSD approach is that the  $r \in O$  will slice the score function of the target, and for each of those slices the input is sliced according to some other directions  $g_r$ . It is known that, under some standard regularity conditions, maxSKSD discriminates measures (Gong et al., 2021, Corollary 3.1). Replacing KSD with maxSKSD in SVGD results in *sliced-SVGD* (S-SVGD; Gong et al. (2021)).

An issue with this procedure is that the basis  $O$  is a user-defined choice and may result in inefficient performance. Additionally, restricting  $g_r$  and  $r$  to one dimension might not capture the underlying structure as effectively as higher dimensional projections. Our approach, described in Section 4, will use projections onto subspaces that have dimension potentially greater than one. The framework of these subspace projections is described next. This framework also allows us to explore projections in the underlying geometry more efficiently by using Riemannian optimisation.

**Grassmann Manifold** Let  $m \in \{1, 2, \dots, d\}$ . A crucial ingredient of our proposed method is the *Grassmann manifold* of  $m$ -dimensional subspaces in  $\mathbb{R}^d$  (Bendokat et al., 2020), defined as

$$\text{Gr}(d, m) = \{E \subseteq \mathbb{R}^d : \dim(E) = m\}.$$

The intuition behind our proposed method is to define a KSD that seeks the worst possible discrepancy over all distinct subspaces of dimension  $m$ , where the worst-case subspace can be found with gradient-type optimisation on the Grassmann manifold. We briefly review the key ingredients of optimisation on the Grassmann manifold; a more detailed discussion is deferred to Appendix A.

To make sense of optimisation on the Grassmann manifold, we first need to represent each subspace  $E \in \text{Gr}(d, m)$  in memory. One way of doing so is by a projection operator that maps every  $x \in \mathbb{R}^d$  to an element in the subspace. To this end, we define a *projector* of rank  $m$  to be a  $d \times m$  matrix  $A$  with orthonormal columns, i.e.  $A^\top A = I_m$ , where  $I_m$  is the  $m \times m$  identity matrix. It then follows that  $AA^\top$  is a projection matrix since  $(AA^\top)^2 = AA^\top$  and  $(AA^\top)^\top = AA^\top$ . Denote by  $[A] = \{Ay : y \in \mathbb{R}^m\} \subset \mathbb{R}^d$  the image of  $A$ .

We represent each subspace  $E \in \text{Gr}(d, m)$  by *any* projector  $A$  for which  $[A] = E$ . This is always

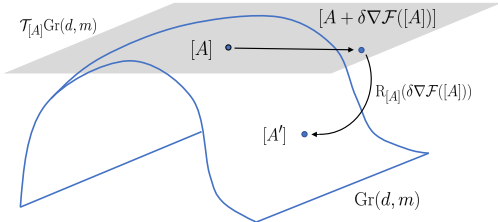


Figure 2: An illustration of a gradient ascent step on the Grassmann manifold.

possible since, for each projector  $A$  the subspace  $[A]$  is trivially an element of  $\text{Gr}(d, m)$ , and, conversely, for any  $E \in \text{Gr}(d, m)$ , we can construct a corresponding projector  $A$  by column-wise appending the elements of any orthonormal basis of  $E$ . Such  $A$  is unique up to an orthogonal transformation: for any projectors  $A, B$ , the subspaces  $[A] = [B]$  if and only if  $A = BC$  for some orthogonal matrix  $C$  in  $\mathbb{R}^{m \times m}$ . As we will show in Section 3, the proposed KSD does not depend on which projector is chosen so long as it corresponds to the same subspace  $E$ .

Given a representative  $A$  of  $[A] \in \text{Gr}(d, m)$ , the *tangent space* at  $[A]$  is defined as  $\mathcal{T}_{[A]}\text{Gr}(d, m) = \{\Delta \in \mathbb{R}^{d \times m} : A^\top \Delta = 0\}$ . We endow a Riemannian metric on  $\mathcal{T}_{[A]}\text{Gr}(d, m)$  by restricting the standard matrix inner product  $\langle \Delta, \tilde{\Delta} \rangle_0 = \text{Tr}(\Delta^\top \tilde{\Delta})$  to elements  $\Delta, \tilde{\Delta}$  in  $\mathcal{T}_{[A]}\text{Gr}(d, m)$ . It is then clear that  $\Pi_A = I_d - AA^\top$  is a projection operator of  $\mathbb{R}^{d \times m}$  onto  $\mathcal{T}_{[A]}\text{Gr}(d, m)$ :  $\Pi_A G = (I_d - AA^\top)G \in \mathcal{T}_{[A]}\text{Gr}(d, m)$ , for any  $G \in \mathbb{R}^{d \times m}$  (Boumal, 2020, Proposition 3.53).

Let  $\mathcal{F}$  be a function defined on  $\text{Gr}(d, m)$  that is smooth in a proper sense (Boumal, 2020, Section 3). Given a base point  $[A] \in \text{Gr}(d, m)$ , maximising  $\mathcal{F}$  requires (i) finding a vector field of steepest ascent of  $\mathcal{F}$  at  $[A]$ , and (ii) moving from  $[A]$  along that vector field without leaving the manifold.

On the Grassmann manifold, the direction of steepest ascent is no longer represented by the standard (Euclidean) gradient  $\nabla \mathcal{F}([A])$ . Instead, it is given by the *Riemannian gradient* of  $\mathcal{F}$  at  $[A]$ , given by

$$\text{grad} \mathcal{F} = \Pi_A \nabla \mathcal{F}([A]),$$

where  $[\nabla \mathcal{F}([A])]_{ij} = \partial \mathcal{F}([A]) / \partial A_{ij}$ . That is,  $\text{grad} \mathcal{F}$  is the projection of the Euclidean gradient  $\nabla \mathcal{F}([A])$  onto the tangent space at  $[A]$ . For (ii), one means of moving along a manifold in a given direction is by the *exponential map* (Boumal, 2020, Section 10.2). In our work, we use the *polar retraction* (Boumal, 2020, Definition 3.41), which is an accurate (second order) approximation of the exponential map that is amenable to computation, and is defined as

$$\text{R}_{[A]}(\Delta) = [UV^\top], \quad (6)$$

where  $[A] \in \text{Gr}(d, m)$ , and  $A + \Delta = USV^\top$  is a thin Singular Value Decomposition (see e.g. Boumal (2020, Section 9.6)).

To summarise, with a step-size  $\delta > 0$ , one step to maximise  $\mathcal{F}$  starting from a base point  $[A]$  is given by

$$[A'] = \text{R}_{[A]}(\delta \text{grad} \mathcal{F}).$$

See Fig.2 for an illustration. This is a generalisation of the standard gradient ascent to the Grassmann manifold (Boumal, 2020, Section 4.3).

### 3 GRASSMANN KERNEL STEIN DISCREPANCY

The starting point of our proposed approach is to introduce a novel Stein discrepancy which incorporates dimension reduction features through the use of a projector. To this end, for a projector  $A$  of rank  $m$  we introduce a matrix-valued kernel  $k_A$  arising from an embedding on  $\mathbb{R}^d$  into a  $m$ -dimensional subspace where  $m \leq d$

$$k_A(x, y) = AA^\top k(A^\top x, A^\top y), \quad x, y \in \mathbb{R}^d,$$

where  $k : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$  is a scalar-valued positive definite radial kernel on  $\mathbb{R}^m$ . That is,  $k(u, v) = \Psi(\|u - v\|_2)$ ,  $u, v \in \mathbb{R}^m$ , where  $\Psi : \mathbb{R} \rightarrow \mathbb{R}$  and  $\|\cdot\|_2$  is the Euclidean 2-norm. Before proceeding, since we shall be performing optimisation over  $\text{Gr}(d, m)$ , we must be sure that for any two projectors  $A, B$  that are equivalent in  $\text{Gr}(d, m)$  the kernels  $k_A, k_B$  are equal. The next lemma assures us of this (see Appendix B for the proof).

**Lemma 1.** *Let  $A, B$  be projectors of rank  $m$  with  $[A] = [B]$ , then  $k_A = k_B$ .*

Safe in the knowledge that  $k_A$  is well defined over  $\text{Gr}(d, m)$ , we may define an associated KSD, which reveals how this choice of  $k_A$  induces a projection of the data and score function.

$$\text{KSD}_A(Q, P) = \sup_{\phi \in \mathcal{B}_{k_A}} \mathbb{E}_Q[s_p(x) \cdot \phi(x) + \nabla \cdot \phi(x)] \quad (7)$$

$$= \sup_{\phi \in \mathcal{B}_k^m} \mathbb{E}_Q[(A^\top s_p(x)) \cdot \phi(A^\top x) + \nabla \cdot \phi(A^\top x)]. \quad (8)$$

where  $\mathcal{B}_{k_A}$  is the unit ball of the RKHS corresponding to  $k_A$  and  $\mathcal{B}_k^m$  is the unit ball of the  $m$ -times Cartesian product of  $\mathcal{H}_k$ , the RKHS corresponding to  $k$ . Contrasting this to (4), setting  $m = 1$  and  $r, g_r = A$  results in (8). The  $r^\top g_r$  term in (4) becomes  $I_m$  in (8) since  $A$  has orthonormal columns.

Analogous to the approach for KSD (see Liu et al.

(2016)), we exploit the kernel trick to express  $\text{KSD}_A$ :

$$\begin{aligned} \text{KSD}_A(Q, P) &= \mathbb{E}_{x, x' \sim Q} [(A^\top s_p(x)) \cdot A^\top s_p(x') k(A^\top x, A^\top x')] \\ &\quad + 2\mathbb{E}_{x, x' \sim Q} [(A^\top s_p(x)) \cdot \nabla_{x_2} k(A^\top x, A^\top x')] \\ &\quad + \mathbb{E}_{x, x' \sim Q} [\text{Tr}(\nabla_{x_1, x_2} k(A^\top x, A^\top x'))], \end{aligned} \quad (9)$$

where  $x, x'$  are i.i.d. random variables drawn from  $Q$ ,  $\nabla_{x_i} k(A^\top x, A^\top x')$  denotes the gradient of  $k$  with respect to the  $i$ -th argument,  $\nabla_{x_1, x_2} k$  is the matrix with  $ij$ -th entry  $\partial^2 k / \partial x_i \partial x_j$ , and  $\text{Tr}(\cdot)$  denotes the trace.

It is clear that  $\text{KSD}_A$  does not characterise probability measures for any fixed  $A$ . For example if  $P = \mathcal{N}(0, C)$  and  $Q = \mathcal{N}(0, C')$  are distinct probability measures such that  $A^\top C = A^\top C'$ , then  $\text{KSD}_A(Q, P) = 0$ . This motivates us to consider the worst-case discrepancy taken over all  $m$ -dimensional subspaces. We define the *Grassmann Kernel Stein Discrepancy* (GKSD) as

$$\text{GKSD}(Q, P) = \sup_{[A] \in \text{Gr}(d, m)} \text{KSD}_A(Q, P). \quad (10)$$

The following theorem guarantees that GKSD is able to discriminate distinct distributions. The assumptions on  $P, Q$  are standard (Gong et al., 2021).

**Theorem 2.** *Let  $P, Q$  be Borel measures with continuously differentiable densities  $p, q$  both supported on  $\mathbb{R}^d$  with  $\lim_{\|x\|_2 \rightarrow \infty} q(x) = 0$  and suppose the kernel  $k(u, v) = \Psi(\|u - v\|_2)$  is characteristic and bounded with continuous second order partial derivatives. Then the GKSD (10) equals 0 if and only if  $p = q$ .*

**Remark 1.** *A radial kernel  $k(x, y) = \Psi(\|x - y\|_2)$  is characteristic if and only if it is integrally strictly positive definite (Sriperumbudur et al., 2011). In particular, the assumptions of Theorem 2 hold for Gaussian and Inverse Multiquadric kernels.*

Intuitively, projecting on higher dimensions for  $s_p(x)$  will allow us to better account for the correlation between dimensions. With optimal projectors, correlations can be captured regardless of the projection dimension, since both GKSD and maxSKSD can discriminate distinct distributions. In practice, however, optimal projectors are not available, and hence projecting onto 1 dimension can lead to sub-optimal approximation due to weak correlation signals.

## 4 GRASSMANN STEIN VARIATIONAL GRADIENT DESCENT

Equipped with this new form of kernel discrepancy we can perform a variant of SVGD which we call *Grassmann Stein Variational Gradient Descent* (GSVDG). More specifically, we seek to find

$$\phi_A^* \in \arg \sup_{\phi \in \mathcal{B}_{k_A}} - \frac{d}{d\epsilon} \text{KL}(T_\#^* Q, P) \Big|_{\epsilon=0}, \quad (11)$$

where  $T_\# = I + \epsilon\phi$ , and  $I$  is the identity map. The supremum on the RHS of (11) is attained by

$$\phi_A^*(\cdot) = \mathbb{E}_Q [AA^\top s_p(x) k(A^\top x, A^\top \cdot) + A \nabla_{x_1} k(A^\top x, A^\top \cdot)],$$

for which it is clear that

$$- \frac{d}{d\epsilon} \text{KL}(T_\#^* Q, P) \Big|_{\epsilon=0} = \text{KSD}_A(P, Q), \quad T_\#^* = I + \epsilon\phi_A^*.$$

To transport the particles efficiently, at every step, we seek to identify the subspace  $[A] \in \text{Gr}(d, m)$  with the largest discrepancy between  $P$  and  $Q$ , so that we select

$$[A^*] \in \arg \sup_{[A] \in \text{Gr}(d, m)} \text{KSD}_A(P, Q), \quad (12)$$

for which  $\text{KSD}_{A^*} = \text{GKSD}(P, Q)$ . In practice, the objective (12) can be (approximately) solved by searching for a critical point of  $\text{KSD}_A(P, Q)$  at which its Riemannian gradient vanishes. This is justified by the following proposition, which states that, when the probability measures  $P$  and  $Q$  differ in some low-dimensional subspace  $[A_0]$ , then  $[A_0]$  is indeed a solution sought by (12).

**Proposition 3.** *Suppose the conditions in Theorem 2 hold for the kernel  $k$  and measures  $P, Q$  with associated densities  $q, p$ . Assume further that the kernel satisfies a smoothness condition Assumption 1 in Appendix D.1, and that there exists  $[A_0] \in \text{Gr}(d, m)$  such that  $q, p$  admit the following decomposition*

$$q(x) \propto q^m(P_0 x) \xi(\Pi_{A_0} x) \quad (13)$$

$$p(x) \propto p^m(P_0 x) \xi(\Pi_{A_0} x), \quad (14)$$

where  $P_0 := A_0 A_0^\top$ ,  $\Pi_{A_0} = I_d - A_0 A_0^\top$  as before, and  $p^m, q^m$  and  $\xi$  are smooth, positive functions that are integrable on  $\mathbb{R}^d$  (i.e. they are unnormalised densities). Then  $\text{gradKSD}_A(P, Q)|_{A=A_0} = 0$ .

**Remark 2.** *The assumption on the decomposition (13) and (14) holds when the density  $q$  agrees with the target  $p$  except in the subspace  $[A_0]$ . An example is that  $q$  is a prior, and  $p$  is the posterior induced from  $q$  and a likelihood function that depends on  $x$  only through  $A_0 x$ . See Appendix D for more discussions and the proof.*

Sequentially solving the optimisation problems (11) and (12) every time-step yields a particle transport scheme which terminates if and only if the GKSD between the target distribution and the current particle distribution is zero, which implies  $P = Q$  under the conditions of Theorem 2.

Replacing  $Q$  by an empirical distribution of particles  $(x_1^0, \dots, x_N^0)$ , and denoting by  $(x_1^t, \dots, x_N^t)$  the set of particles obtained at step  $t$ , the scheme becomes

$$\begin{aligned} x_i^{t+1} &= x_i^t + \frac{\epsilon}{N} \sum_{j=1}^N [A_t A_t^\top s_p(x_j^t) k(A_t^\top x_j^t, A_t^\top x_i^t) \\ &\quad + A_t \nabla_{x_1} k(A_t^\top x_j^t, A_t^\top x_i^t)], \quad i = 1, \dots, N, \end{aligned} \quad (15)$$

$$A_t \in \arg \sup_{[A] \in \text{Gr}(d, m)} \alpha_t([A]), \quad (16)$$

---

**Algorithm 1** Grassmann Stein Variational Gradient Descent (GSVGD)
 

---

1: **Inputs:**  $M$  initialized projectors  $[A_{0,l}] \in \text{Gr}(d, m)$ ,  $l = 1, \dots, M$ ;  $N$  initialised particles  $\{x_i^0\}_{i=1}^N$ ; iteration number  $n_{\text{epochs}}$ ; step sizes  $\epsilon, \delta > 0$ .

2: **for**  $t = 1, 2, \dots, n_{\text{epochs}}$  **do**

3: Update particles  $x_i^{t+1} = x_i^t + \epsilon \sum_{l=1}^M \widehat{\phi}_{A_{t,l}}(x_i^t)$  for  $i = 1, \dots, N$ , where

$$\widehat{\phi}_{A_{t,l}}(x_i^t) = \frac{1}{N} \sum_{j=1}^N [A_{t,l} A_{t,l}^\top s_p(x_j^t) k(A_{t,l}^\top x_j^t, A_{t,l}^\top x_i^t) + A_{t,l} \nabla_{x_1} k(A_{t,l}^\top x_j^t, A_{t,l}^\top x_i^t)]. \quad (17)$$

4: Update projectors

$$A_{t+1,l} = \text{R}_{[A_{t,l}]} \left( \delta \Pi_{A_{t,l}} \nabla \alpha_t([A_{t,l}]) + \sqrt{2T\delta} \Pi_{A_{t,l}} \xi_{t,l} \right)$$

for  $l = 1, \dots, M$ , where  $\xi_{t,l} \in \mathbb{R}^{d \times m}$  has i.i.d.  $\mathcal{N}(0, 1)$  entries, and  $\text{R}_{[A_{t,l}]}$  is defined in (6).

5: **end for**

6: **Return:** Particles  $\{x_i^{n_{\text{epochs}}}\}_{i=1}^N$ .

---

where  $\alpha_t([A]) = \text{KSD}_A(Q_t, P)$ , for  $Q_t(dx) = \frac{1}{N} \sum_{i=1}^N \delta_{x_i^t}(dx)$ .

**Optimal Projections via SDEs:** Proposition 3 implies the “optimal” projection  $A_0$  is *one* critical point of  $\alpha_t$ , but it does not imply the critical point is unique. Indeed, the problem (12) is non-convex and there might be multiple local maximisers, from which finding global maxima is generally NP-hard. A well-known remedy is to add white noise to the gradient flow (Aluffi-Pentini et al., 1985; Geman and Hwang, 1986), allowing the trajectory to escape from local maxima. Mathematically, such methods can be formulated in terms of Stochastic Differential Equations (SDEs). In this particular instance, we explore the use of SDEs taking values on the Grassmann manifold. We consider stochastic dynamics defined by the solution of the following Stratonovich SDE

$$\begin{aligned} dA(t) &= \mu \text{grad} \alpha_t([A]) dt + \sqrt{2T\mu} \Pi_{A_t} \circ dW(t), \quad (18) \\ &= \mu \Pi_{A(t)} \nabla \alpha_t([A(t)]) dt + \sqrt{2T\mu} \Pi_{A(t)} \circ dW(t) \end{aligned}$$

where  $W(t)$  is a  $\mathbb{R}^{d \times m}$ -valued Brownian motion and  $\mu, T > 0$ . Combining the particle and projector evolution we can consider the full system as a discretisation of the following coupled ODE-SDE

$$\begin{aligned} \dot{x}_i(t) &= \frac{1}{N} A(t) \sum_{j=1}^N [ \\ &\quad A(t)^\top s_p(x_j(t)) k(A(t)^\top x_j(t), A(t)^\top x_i(t)) \\ &\quad + \nabla_{x_2} k(A(t)^\top x_j(t), A(t)^\top x_i(t))], \\ dA(t) &= \mu \Pi_{A(t)} \nabla \alpha_t([A(t)]) dt + \sqrt{2T\mu} \Pi_{A(t)} \circ dW(t). \end{aligned}$$

We see that the parameter  $T$  acts as a temperature, facilitating moving between local minima of the function  $\alpha$  by annealing the loss function. The parameter  $\mu$  controls the relative timescale between the evolution of  $A(t)$  and the particles  $(x_1^t, \dots, x_N^t)$ . In particular, when  $\mu \gg 1$ , the distribution of  $A(t)$  is approximated by the quasi-stationary distribution  $d\mu_t([A]) = e^{\alpha_t([A])/T} d[A]$ , where  $d[A]$  denotes the canonical measure on the Grassmann manifold. In the other extreme, when  $\mu \ll 1$  the projector  $A(t)$  will evolve on a slower timescale than the particles, so that the particles will reach their stationary configuration for a fixed projector before it evolves.

To simulate the Stratonovich SDE (18), we make use of the formulation of diffusions on a Riemannian manifold, making use of the associated retraction map, see Staneva and Younes (2017), Baxendale et al. (1976), and Belopolskaya and Dalecky (2012) where the  $\text{Gr}(d, m)$ -valued SDE (18) can be expressed in the form of a system of Ito SDEs

$$\begin{aligned} A(t) &= \text{R}_{[A(t)]}(B(t)), \quad (19) \\ dB(t) &= \mu \Pi_{A(t)} \nabla \alpha_t([A(t)]) dt + \sqrt{2T\mu} \Pi_{A(t)} dW_t. \quad (20) \end{aligned}$$

With step-size  $\delta > 0$ , an Euler-Maruyama discretisation for (20) results in

$$A_{t+1} = \text{R}_{[A_t]} \left( \Pi_{A_t} \nabla \alpha_t([A_t]) \delta + \sqrt{2T\delta} \Pi_{A_t} \xi_t \right), \quad (21)$$

where  $\xi_t$  is a  $d \times m$  matrix with independent Gaussian entries. Note that we “absorb” the constant  $\mu$  into the step-size  $\delta$  for simplicity.

**Batch Generalisation** An important observation is that at each step (16),  $x_i^{t+1}$  is different from  $x_i^t$  only in the image of  $A$ . This contrasts with both SVGD as well as S-SVGd, both of which update the particles along all dimensions. One extension is to use more than one projector  $A_{t,l} \in \mathbb{R}^{m_l \times d}$  for  $l = 1, \dots, M$ , where  $1 \leq M \leq d$  and  $\sum_{l=1}^M m_l \leq d$ . At each optimisation step, we can impose orthonormality using QR factorisation, so that the  $d \times d$  matrix  $A_t$  formed by column-wise concatenation of  $A_{t,1}, \dots, A_{t,M}$ , is orthogonal. However, this requires  $\mathcal{O}(d^3)$  operations, which can be prohibitive for large  $d$ . We therefore propose to impose orthonormality every 1000 steps, which we find sufficient in practice. The full algorithm is presented in Algorithm 1.

**Computational Complexity** Each iteration of the particles update in Algorithm 1 requires  $\mathcal{O}(MNdm(N+d))$  operations. Compared with the  $\mathcal{O}(Nd(N+d))$  cost for S-SVGd, the extra factor  $Mm$  arises from the fact that (17) involves matrix multiplication with  $A_{t,l}$ ,

which can be avoided in S-SVGD due to simplification by using the canonical basis. However, this extra computation can be largely reduced by using a fixed number of projectors  $M$  and projection dimension  $m$ . In our experiments, we find  $M = 20$  projectors are sufficient. A reasonable choice of  $m$ , however, will depend on the specific problem. Finally, each step of the projector update requires  $\mathcal{O}(Mdm^2)$  due to the retraction map, which reduces to the same cost as the slice update in S-SVGD when  $M = d$  and  $m = 1$ .

## 5 RELATED WORK

**Connections to S-SVGD** The major differences between GSVGD and S-SVGD are that (i) GSVGD allows projecting onto an arbitrary low-dimensional subspace, (ii) GSVGD updates the projectors via Riemannian optimisation on the Grassmann manifold, whereas S-SVGD optimises the slices in the unit ball of  $\mathbb{R}^d$ , and (iii) the projectors for the score function and the data are set to the same in GSVGD but not in S-SVGD. In particular, S-SVGD slices the score function with a fixed basis  $O$  (i.e. the canonical basis), and such arbitrary choice may result in inefficient exploration of the state-space if the basis vectors  $r \in O$  do not align with the latent dimensions of the target. On the other hand, optimising the projectors for the score function and the data separately would lead to a complex joint optimisation program that is practically challenging. GSVGD strikes a balance by using the same projector for the score function and the data.

**Other SVGD Variants** Liu and Zhu (2018) extends SVGD to tasks defined on Riemannian manifolds and Shi et al. (2021) proposes a variant applied to non-Euclidean spaces (such as a simplex) via mirror descent. These are fundamentally different from our work, which concerns problems on the Euclidean space, while evolving the projectors on the Grassmann manifold. Closer to our method is the pSVGD of Chen and Ghattas (2020), where the particles are projected onto a subspace defined by the top eigenvectors of a gradient information matrix. This method however suffers from practical issues, e.g. finding the eigen-decomposition would incur a computational cost that is cubic in the dimension. Other attempts to alleviate the curse-of-dimensionality problem of SVGD include Message Passing SVGD (Zhuo et al., 2018) and graphical SVGD (Wang et al., 2018), both of which are limited to problems where the target distribution is defined on a *probabilistic graphical model* (PGM) with a known Markov structure. In contrast, GSVGD can be applied to distributions of an arbitrary form.

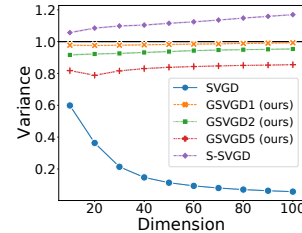


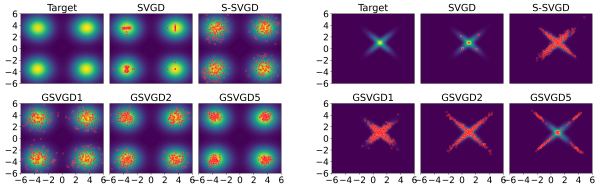
Figure 3: Estimating the dimension-averaged marginal variance of  $p(x) = \mathcal{N}(x; 0, I_d)$  across different dimensions  $d$ . GSVGD1 means GSVGD with 1-dimensional projections. Black solid line marks the true value. Results are averaged over 20 repetitions.

## 6 EXPERIMENTS

We study the uncertainty quantification property of GSVGD against existing methods such as SVGD and S-SVGD on a number of experiments. We conclude from our experiments that SVGD can underestimate the uncertainty, S-SVGD can overestimate it and GSVGD, with appropriate projection dimensions, yields the best estimates. For GSVGD, we use the batch approach of Algorithm 1, where at most  $M = 20$  projectors with the same projection dimension  $m$  are used. The temperature  $T$  is set to a small value and is gradually incremented to a larger one via an annealing scheme. The learning rates for all methods are tuned for optimal performance on the multimodal mixture example, and the same choices are used for the other experiments.

In all experiments, we use the Gaussian RBF kernel parameterised by  $k(x, x') = \exp(-\|x - x'\|_2^2 / (2\sigma^2))$  with bandwidth  $\sigma^2 = \text{med}^2 / (2 \log n)$ , where  $\text{med}$  is the median of the pairwise distance of the particles. This follows the heuristic in Liu and Wang (2016). To quantify sample quality, we use the *energy distance* (Székely and Rizzo, 2013) between the approximation and the ground truth. However, we find that the energy distance is not sensitive to differences in correlations between multivariate samples. Therefore, we also evaluate how well each method captures the dependence between multivariate samples by the covariance estimation error  $\|\hat{\Sigma} - \Sigma\|_2$ , where  $\|\cdot\|_2$  is the Frobenius norm,  $\Sigma$  is the ground truth sample covariance matrix given by a Hamiltonian Monte Carlo (HMC) (for non-synthetic data), and  $\hat{\Sigma}$  is the estimated covariance. For further details, see Appendix E. All the code is available at <https://github.com/ImperialCollegeLondon/GSVGD>.

**Multivariate Gaussian** The first toy example is a  $d$ -dimensional multivariate Gaussian  $p(x) = \mathcal{N}(x; 0, I_d)$ . For each method, 500 particles are initialized from  $\mathcal{N}(x; 2\mathbf{1}, 2I_d)$ , where  $\mathbf{1} \in \mathbb{R}^d$  denotes the vector of ones. This is a standard benchmark used to



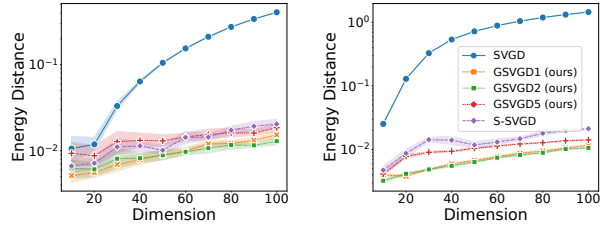
(a) Multimodal mixture (b) X-shaped mixture

Figure 4: Marginals of target and estimated densities in the first 2 dimensions. Red dots are the particles after 2000 iterations and the target density is shown by the contour. Both targets have dimension 50.

illustrate the diminishing variance issue of SVGD (Liu and Wang, 2016; Gong et al., 2021; Zhuo et al., 2018). We show empirically that GSVGD is able to mitigate this problem by plotting the dimension-averaged estimates for the marginal variance in Figure 3. We see that GSVGD estimates stay roughly at 1 for all dimensions as expected, whereas the SVGD estimates scale inversely with dimension. S-SVGD, on the other hand, over-estimates the variance. We remark, however, that GSVGD5 appears to be biased with a lower variance. This is due to the variance underestimation issue of SVGD persisting in 5 dimensional subspaces.

**Multimodal Mixture** In the second example, the target distribution is a mixture of 4  $d$ -dimensional Gaussian distributions  $p(x) = \sum_{k=1}^4 0.25\mathcal{N}(x; \mu_k, I_d)$  with uniform mixture ratios. The first two coordinates of the mean vectors are equally spaced on a circle, while the other coordinates are set to 0; see Figure 4(a) and Appendix E.2. Particles are initialized from  $\mathcal{N}(0, I_d)$  and only the first two dimensions need to be learned. The primary goal is to investigate whether GSVGD can efficiently recover this low-dimensional latent structure by adapting the projectors. As shown in 5(a), GSVGD1 and GSVGD2 outperforms S-SVGD in all dimensions in terms of the energy distance, while GSVGD5 is able to achieve a competitive performance. SVGD, on the other hand, fails to capture the uncertainty of the target, as shown in Figure 4(a).

**X-Shaped Mixture** The target in this experiment is a  $d$ -dimensional mixture of two correlated Gaussian distributions  $p(x) = 0.5\mathcal{N}(x; \mu_1, \Sigma_1) + 0.5\mathcal{N}(x; \mu_2, \Sigma_2)$ . The means  $\mu_1, \mu_2$  of each Gaussian have components equal to 1 in the first two coordinates and 0 otherwise, and the covariance matrices admit a correlated block diagonal structure (see Appendix E.3). The mixture hence manifests as an “X-shaped” density marginally in the first two dimensions (see Figure 4(b)). Figure 5(b) shows that GSVGD1 and GSVGD2 achieve a better approximation quality than both SVGD and S-SVGD across all dimensions. This is further verified by



(a) Multimodal mixture (b) X-shaped mixture

Figure 5: Energy distance between the target distribution and the particle estimation. Results are averaged over 20 repetitions. 95%-confidence intervals are shown by the shaded regions.

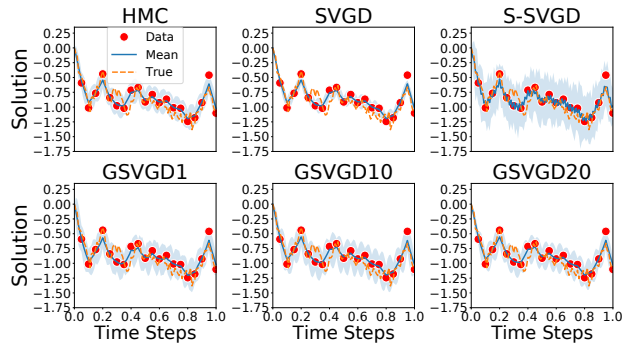


Figure 6: HMC, SVGD, S-SVGD and GSVGD solutions, as well as the posterior mean and 95% confidence interval after 1000 iterations.

plotting the first two coordinates of the final particles in 4 (b). On the other hand, GSVGD5 shows a slightly suboptimal performance for dimensions larger than 30. This is potentially because of the inefficiency due to projecting the particles to 5 dimensional subspace in the case we only care about the first 2.

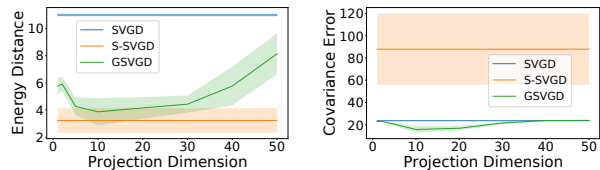
**Conditioned Diffusion Process** The next example is a benchmark that is often used to test inference methods in high dimensions (Cui et al., 2016; Chen and Ghattas, 2020; Detommaso et al., 2018). We consider a stochastic process  $u : [0, 1] \rightarrow \mathbb{R}$  governed by

$$du_t = \frac{10u(1-u^2)}{1+u^2}dt + dx_t, \quad u_t = 0, \quad (22)$$

where  $t \in (0, 1]$ , and the forcing term  $x = (x_t)_{t \geq 0}$  follows a Brownian motion so that  $x \sim q = \mathcal{N}(0, C)$  with  $C(t, t') = \min(t, t')$ . We observe noisy data  $y = (y_{t_1}, \dots, y_{t_{20}})^\top \in \mathbb{R}^{20}$  at 20 equi-spaced time points  $t_i = 0.05i$ , where  $y_{t_i} = u_{t_i} + \epsilon$  for  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  with  $\sigma = 0.1$ , and  $u_{t_i}$  is generated by solving (22) at a true Brownian path  $x_{t_{\text{true}}}$ . The goal is to use  $y$  to infer the forcing term  $x$  and thereby the solution state  $u$ .

The prior we use for  $x$  is the Brownian motion described above. The dynamic is discretized using an Euler-Maruyama scheme with step size  $\Delta t = 10^{-2}$ , leading





(a) Energy distance against projection dimensions (b)  $\|\hat{\Sigma} - \Sigma\|_2$  against projection dimensions

Figure 7: Metrics for Covertypes.

to a 100-dimensional inference problem. Due to the smoothness of the Brownian path, we expect the solution  $u$  to have an intrinsic low-dimensional structure. Figure 6 shows the mean and 95% credible intervals of the particle estimation of  $u$  after 1000 iterations. Compared with a reference HMC, all methods are able to estimate the mean accurately, but SVGD and S-SVGD either underestimate or overestimate the credible intervals. GSVGD projecting on 20 dimensions gives the most accurate uncertainty estimates. Choosing a good projection dimension is non-trivial even in this simple problem; however, by plotting the energy distance between the HMC reference and the results against different projection dimensions (Figure 11 of Appendix), we see that so long as we do not project onto very low or high dimensions, GSVGD achieves a greater degree of agreement with the reference HMC than SVGD and S-SVGD.

**Bayesian Logistic Regression** Following Liu and Wang (2016), we consider the Bayesian logistic regression model applied to the Covertypes dataset (Asuncion and Newman, 2007). We use HMC posterior samples as a gold standard for evaluating the posterior approximations. We ran all methods on a subset consisting 1,000 randomly selected data points 10 times. Similar to the synthetic experiments, we see in Figure 7 that there exists an optimal projection dimension (around  $m = 10$ ) such that GSVGD more accurately approximates the covariance matrix whilst yielding similar energy distances as S-SVGD. S-SVGD severely overestimates the uncertainty, also shown in Figure 8, whereas SVGD severely underestimates it.

## 7 CONCLUSION

We introduced Grassmann Stein Variational Gradient Descent (GSVGD), a novel extension to the SVGD algorithm, to tackle the curse-of-dimensionality problem. At each step, GSVGD seeks the subspace in which the proposal and the target has the largest discrepancy measured by the Grassmann Kernel Stein Discrepancy (GKSD). The evolution of the particles and the projector is determined through a coupled

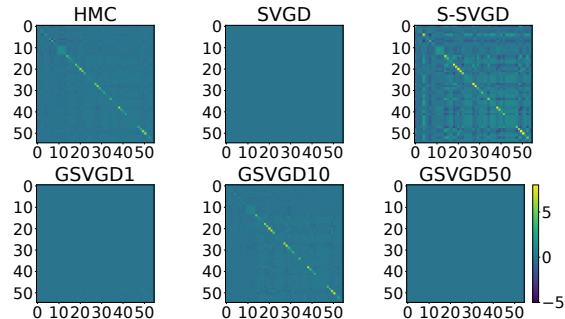


Figure 8: Covariance matrices for Covertypes.

ODE-SDE system, which allows trade-offs between exploration and exploitation. The GSVGD algorithm can be easily extended to a batched version that uses more than one projector. Our experiments demonstrate that GSVGD is able to achieve improved performance over existing methods when the target distribution has an intrinsic low-dimensional structure, especially more accurately quantifying the epistemic uncertainty.

**Limitations of GSVGD:** One limitation of GSVGD is the choice of projection dimension and number of projectors to use in order to attain superior performance. One possible way is to perform a grid search over a range of feasible values and then assess the performance using metrics that measure particle diversity. We leave this as an open question for future work. In addition, our Grassmann-valued SDE introduces additional hyperparameters, which potentially require tuning on a case-by-case basis for each problem. In Section 6 we present some heuristics for tuning GSVGD hyperparameters that we found yielded consistent results, but we would recommend future work to study the sensitivity of the hyperparameters.

## Acknowledgements

XL was supported by the President’s PhD Scholarships of Imperial College London and the EPSRC StatML CDT programme EP/S023151/1. HZ was supported by the EPSRC Centre for Doctoral Training in Modern Statistics and Statistical Machine Learning EP/S023151/1, the Department of Mathematics of Imperial College London and Cervest Limited. JFT was supported by the EPSRC and MRC through the OxWaSP CDT programme EP/L016710/1. GW was supported by an EPSRC Industrial CASE award EP/S513635/1 in partnership with Shell UK Ltd. AD work was supported by Wave 1 of The UKRI Strategic Priorities Fund under the EPSRC Grant EP/W006022/1, particularly the “Foundations of Ecosystems of Digital Twins” theme within that grant & The Alan Turing Institute

## References

- Absil, P.-A., Mahony, R., and Sepulchre, R. (2008). *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, Princeton, NJ.
- Aluffi-Pentini, F., Parisi, V., and Zirilli, F. (1985). Global optimization and stochastic differential equations. *Journal of optimization theory and applications*, 47(1):1–16.
- Asuncion, A. and Newman, D. (2007). UCI machine learning repository.
- Ba, J., Erdogdu, M. A., Ghassemi, M., Sun, S., Suzuki, T., Wu, D., and Zhang, T. (2022). Understanding the Variance Collapse of SVGD in High Dimensions. In *International Conference on Learning Representations*.
- Baxendale, P. et al. (1976). Measures and Markov processes on function spaces. *Mémoires de la Société Mathématique de France*, 46(131-141):3.
- Belopolskaya, Y. I. and Dalecky, Y. L. (2012). *Stochastic equations and differential geometry*, volume 30. Springer Science & Business Media.
- Bendokat, T., Zimmermann, R., and Absil, P.-A. (2020). A Grassmann Manifold Handbook: Basic Geometry and Computational Aspects. *arXiv preprint arXiv:2011.13699*.
- Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. (2017). Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877.
- Boumal, N. (2020). An introduction to optimization on smooth manifolds. *Available online, May*.
- Chen, P. and Ghattas, O. (2020). Projected Stein Variational Gradient Descent. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H., editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1947–1958. Curran Associates, Inc.
- Chwialkowski, K., Strathmann, H., and Gretton, A. (2016). A kernel test of goodness of fit. In Balcan, M. F. and Weinberger, K. Q., editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2606–2615, New York, New York, USA. PMLR.
- Cui, T., Law, K. J., and Marzouk, Y. M. (2016). Dimension-independent likelihood-informed MCMC. *Journal of Computational Physics*, 304:109–137.
- Dang, K.-D., Quiroz, M., Kohn, R., Minh-Ngoc, T., and Villani, M. (2019). Hamiltonian Monte Carlo with energy conserving subsampling. *Journal of machine learning research*, 20.
- Detommaso, G., Cui, T., Marzouk, Y., Spantini, A., and Scheichl, R. (2018). A Stein variational Newton method. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Geman, S. and Hwang, C.-R. (1986). Diffusions for global optimization. *SIAM Journal on Control and Optimization*, 24(5):1031–1043.
- Gilks, W., Richardson, S., and Spiegelhalter, D. (1995). *Markov Chain Monte Carlo in Practice*. Chapman & Hall/CRC Interdisciplinary Statistics. Taylor & Francis.
- Gong, W., Li, Y., and Hernández-Lobato, J. M. (2021). Sliced Kernelized Stein Discrepancy. In *International Conference on Learning Representations*.
- Hoffman, M. D., Gelman, A., et al. (2014). The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *J. Mach. Learn. Res.*, 15(1):1593–1623.
- Kolouri, S., Nadjahi, K., Simsekli, U., Badeau, R., and Rohde, G. (2019). Generalized Sliced Wasserstein Distances. *Advances in Neural Information Processing Systems*, 32:261–272.
- Levin, D. A. and Peres, Y. (2017). *Markov chains and mixing times*, volume 107. American Mathematical Soc.
- Liu, C. and Zhu, J. (2018). Riemannian Stein variational gradient descent for Bayesian inference. In *The 32nd AAAI Conference on Artificial Intelligence*, pages 3627–3634, New Orleans, Louisiana USA. AAAI press.
- Liu, Q., Lee, J., and Jordan, M. (2016). A Kernelized Stein Discrepancy for Goodness-of-fit Tests. In Balcan, M. F. and Weinberger, K. Q., editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 276–284, New York, New York, USA. PMLR.
- Liu, Q. and Wang, D. (2016). Stein Variational Gradient Descent: A General Purpose Bayesian Inference Algorithm. In Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Lu, J., Lu, Y., and Nolen, J. (2019). Scaling limit of the Stein variational gradient descent: The mean field regime. *SIAM Journal on Mathematical Analysis*, 51(2):648–671.
- Paulsen, V. I. and Raghupathi, M. (2016). *An Introduction to the Theory of Reproducing Kernel Hilbert Spaces*. Cambridge University Press.

- Rezende, D. and Mohamed, S. (2015). Variational inference with normalizing flows. In *International Conference on Machine Learning*, pages 1530–1538. PMLR.
- Scovel, C., Hush, D., Steinwart, I., and Theiler, J. (2010). Radial kernels and their reproducing kernel Hilbert spaces. *Journal of Complexity*, 26(6):641–660.
- Shi, J., Liu, C., and Mackey, L. (2021). Sampling with Mirrored Stein Operators. *arXiv preprint arXiv:2106.12506*.
- Sriperumbudur, B. K., Fukumizu, K., and Lanckriet, G. R. (2011). Universality, Characteristic Kernels and RKHS Embedding of Measures. *Journal of Machine Learning Research*, 12(7).
- Staneva, V. and Younes, L. (2017). Learning shape trends: parameter estimation in diffusions on shape manifolds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 38–46.
- Székely, G. J. and Rizzo, M. L. (2013). Energy statistics: A class of statistics based on distances. *Journal of statistical planning and inference*, 143(8):1249–1272.
- Wang, D. and Liu, Q. (2019). Nonlinear Stein Variational Gradient Descent for Learning Diversified Mixture Models. In Chaudhuri, K. and Salakhutdinov, R., editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 6576–6585. PMLR.
- Wang, D., Zeng, Z., and Liu, Q. (2018). Stein Variational Message Passing for Continuous Graphical Models. In Dy, J. G. and Krause, A., editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10–15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 5206–5214. PMLR.
- Welling, M. and Teh, Y. W. (2011). Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688. Citeseer.
- Yoon, J., Kim, T., Dia, O., Kim, S., Bengio, Y., and Ahn, S. (2018). Bayesian Model-Agnostic Meta-Learning. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Zhuo, J., Liu, C., Shi, J., Zhu, J., Chen, N., and Zhang, B. (2018). Message passing Stein variational gradient descent. In *International Conference on Machine Learning*, pages 6018–6027. PMLR.

## Supplementary Material: Grassmann Stein Variational Gradient Descent

### A FURTHER BACKGROUND MATERIALS ON THE GRASSMANN MANIFOLD

We provide a brief overview of the Grassmann manifold and optimisation on the Grassmann manifold. We focus only on the essential components required to understand the proposed method; a thorough treatment of these topics is beyond the scope of the paper. For more details, we refer the interested readers to [Bendokat et al. \(2020\)](#) regarding the theoretical perspectives of the Grassmann manifold, and [Boumal \(2020\)](#); [Absil et al. \(2008\)](#) regarding optimisation on Grassmann and other manifolds.

**The Stiefel and Grassmann Manifolds** The *Grassmann manifold*  $\text{Gr}(d, m)$  ([Boumal, 2020](#), Section 9) is the space of all  $m$ -dimensional linear subspaces in  $\mathbb{R}^d$ . Each element of  $\text{Gr}(d, m)$  is an equivalence class of  $d \times m$  matrices whose columns form an orthonormal basis for the same subspace. To this end, define the *Stiefel manifold*  $\text{St}(d, m) := \{A \in \mathbb{R}^{d \times m} : A^\top A = I_m\}$ , where  $I_m$  is the  $m \times m$  identity matrix. That is,  $\text{St}(d, m)$  is the set of projectors of rank  $m$  as defined in Section 2. The Grassmann manifold is built from the Stiefel manifold by identifying matrices whose columns span the same subspace. Formally, define an equivalence relationship  $\sim$  as  $A \sim B \iff A = BC$  for any  $C \in O(m)$ , the group of orthogonal matrices in  $\mathbb{R}^{m \times m}$ .  $\text{Gr}(d, m)$  is the quotient space

$$\text{Gr}(d, m) = \text{St}(d, m) / \sim = \{[A] : A \in \text{St}(d, m)\}, \quad \text{where } [A] = \{B \in \text{St}(d, m) : A \sim B\}.$$

The Grassmann manifold is an instance of *quotient manifolds* on which we can define smooth functions; see [Boumal \(2020, Chapter 9\)](#) for a full treatment of quotient manifolds and [Boumal \(2020, Chapter 3\)](#) for the precise definition of smoothness of a manifold-valued function. This allows the use of gradient-type optimisation to find the critical points of a smooth function. A function  $\mathcal{F}$  defined on  $\text{Gr}(d, m)$  differs from one defined on  $\text{St}(d, m)$  in that the former must preserve invariance of the chosen projector; that is,  $\mathcal{F}([A])$  depends only on the subspace  $[A] \in \text{Gr}(d, m)$  but not on  $A$ . An example is the matrix-valued kernel defined in Lemma 1:  $A \mapsto k_A(x, y) = AA^\top k(A^\top x, A^\top y)$ ,  $x, y \in \mathbb{R}^d$ .

**Tangent Spaces and the Riemannian Gradient** To define the Riemannian gradient  $\text{grad}\mathcal{F}([A])$  of a smooth function  $\mathcal{F}$  on  $\text{Gr}(d, m)$ , we endow  $\text{Gr}(d, m)$  with the standard inner matrix product  $g_{[A]}(\Delta, \tilde{\Delta}) := \langle \Delta, \tilde{\Delta} \rangle_0 = \text{Tr}(\Delta^\top \tilde{\Delta})$  restricted to  $\Delta, \tilde{\Delta} \in \text{Gr}(d, m)$ . The Riemannian gradient is then defined by the relationship

$$g_{[A]}(\text{grad}\mathcal{F}([A]), \Delta) = \left. \frac{d}{d\delta} \mathcal{F}(R_{[A]}(\delta\Delta)) \right|_{\delta=0},$$

for  $\Delta \in \mathcal{T}_{[A]}\text{Gr}(d, m) := \{\Delta \in \mathbb{R}^{d \times m} : A^\top \Delta = 0\}$ , and  $R_{[A]}(\cdot)$  is a *retraction* ([Boumal, 2020](#), Definition 3.41) that send  $[A]$  along a chosen direction while remaining on the Grassmann manifold. In particular,  $\mathcal{T}_{[A]}\text{Gr}(d, m)$  is called the *tangent space* at  $[A]$  ([Boumal, 2020](#), Definition 3.10). For the Grassmann manifold, the Riemannian gradient is simply the orthogonal projection of the standard gradient  $\nabla\mathcal{F}([A])$  to the tangent space:  $\text{grad}\mathcal{F}([A]) = \Pi_A \nabla\mathcal{F}([A])$ , where  $\Pi_A = I_d - AA^\top$ , and  $[\nabla\mathcal{F}([A])]_{ij} = \frac{\partial}{\partial A_{ij}} \mathcal{F}([A])$ .

**Riemannian Gradient Descent** To find a critical point of  $\mathcal{F}$  where its Riemannian gradient is zero, *Riemannian gradient descent* (RGD), which generalised the standard gradient descent to a Riemannian manifold, can be used. In RGD, we start from  $[A_0] \in \text{Gr}(d, m)$  and iterate  $[A_{t+1}] = R_{[A_t]}(-\delta \text{grad}\mathcal{F}([A_t]))$  for  $t = 0, 1, \dots$  and step size  $\delta > 0$ . Intuitively,  $\text{grad}\mathcal{F}([A_t])$  finds the direction of steepest descent at  $[A_t]$ , and the retraction  $R_{[A_t]}(\cdot)$  is the action of moving along  $\text{Gr}(d, m)$  in that direction.

## B PROOF of LEMMA 1

*Proof.* One can show (e.g. see [Bendokat et al. \(2020\)](#)) that  $B \in [A]$  if and only if there exists an orthogonal matrix  $C \in \mathbb{R}^{m \times m}$  such that  $B = AC$ . Then, using the orthogonality of  $C$  we obtain

$$\begin{aligned} k_B(x, y) &= (AC)(AC)^\top k((AC)^\top x, (AC)^\top y) = ACC^\top A^\top k(C^\top A^\top x, C^\top A^\top y) \\ &= AA^\top \Psi (\|C^\top (A^\top x - A^\top y)\|_2) \\ &= AA^\top \Psi (\|A^\top x - A^\top y\|_2) = k_A(x, y), \end{aligned}$$

for all  $x, y \in \mathbb{R}^d$ . □

## C PROOF of THEOREM 2

The idea of the proof is to rewrite GKSD as a double integral type discrepancy, similar to [Liu et al. \(2016, Theorem 3.6\)](#), and the proof shall largely follow. We will then conclude that the discrepancy distinguishes  $p, q$  given the assumptions on the kernel.

The main difference to [Liu et al. \(2016, Theorem 3.6\)](#) is that we shall be using a matrix-valued kernel formed from the projection  $A$ , as written in (8), rather than a standard scalar-valued kernel. This means we must introduce some extra notation. First, given a density  $q$  on  $\mathbb{R}^d$ , we will be using the operator  $\mathcal{A}_q f(x) = s_q(x)^\top f(x) + \nabla \cdot f(x)$  for functions  $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$  and abuse notation so that for functions  $F: \mathbb{R}^d \rightarrow \mathbb{R}^{d \times m}$  we have  $\mathcal{A}_q F(x) = (\mathcal{A}_q F_1(x), \dots, \mathcal{A}_q F_m(x)) \in \mathbb{R}^{1 \times m}$  where  $F_i: \mathbb{R}^d \rightarrow \mathbb{R}^d$  is the  $i$ -th column of  $F$ .

The next lemma is a vectorised version of [Liu et al. \(2016, Lemma 2.2, Lemma 2.3\)](#) and the proof follows by applying those results elementwise to  $F$ .

**Lemma 4.** *Under the assumptions on  $q$  in Theorem 2, if  $F: \mathbb{R}^d \rightarrow \mathbb{R}^{d \times m}$  has bounded entries then*

$$\mathbb{E}_Q [(s_p(x) - s_q(x))^\top F(x)] = \mathbb{E}_Q [\mathcal{A}_p F(x)].$$

**Remark 3.** *The boundedness condition on  $F$  is to ensure it lies in the Stein class of  $q$  ([Liu et al., 2016, Definition 2.1](#)), since  $q$  is assumed to be supported on all of  $\mathbb{R}^d$ , so that [Liu et al. \(2016, Lemma 2.2, Lemma 2.3\)](#) can be applied. A similar condition on  $F$  can be imposed when the support of  $q$  is a compact subset  $\mathcal{X}$  of  $\mathbb{R}^d$ .*

Now we derive the analogy to [Liu et al. \(2016, Theorem 3.6\)](#) using Lemma 4, the proof is largely the same as the proof of [Liu et al. \(2016, Theorem 3.6\)](#), which we include for clarity.

**Lemma 5.** *Under the assumptions in Theorem 2 the projected KSD can be expressed as follows*

$$\text{GKSD}(Q, P) = \sup_{[A] \in \text{Gr}(d, m)} \mathbb{E}_{x, x' \sim Q} [(A^\top \delta_{p, q}(x))^\top A^\top \delta_{p, q}(x') k(A^\top x, A^\top x')], \quad (23)$$

where  $\delta_{p, q}(x) = s_p(x) - s_q(x)$ .

*Proof.* Define  $v(x, x') = \mathcal{A}_p(Ak(A^\top \cdot, A^\top x'))(x) \in \mathbb{R}^{1 \times m}$ . Since  $k$  is assumed to be bounded, we can apply Lemma 4 with  $F(x) = Ak(A^\top x, A^\top x')$  to yield

$$\begin{aligned} \mathbb{E}_{x, x' \sim Q} [(A^\top \delta_{p, q}(x))^\top A^\top \delta_{p, q}(x') k(A^\top x, A^\top x')] &= \mathbb{E}_{x, x' \sim Q} [v(x, x') A^\top \delta_{p, q}(x')] \\ &= \mathbb{E}_{x, x' \sim Q} [\delta_{p, q}(x')^\top A v(x, x')^\top]. \end{aligned}$$

The assumption on the second order partial derivatives of the kernel and the proof of [Liu et al. \(2016, Theorem 3.6\)](#) assure us that we can apply Lemma 4 again with  $F(x') = Av(x, x')^\top$  to get

$$\mathbb{E}_{x, x' \sim Q} [(A^\top \delta_{p, q}(x))^\top A^\top \delta_{p, q}(x') k(A^\top x, A^\top x')] = \mathbb{E}_{x, x' \sim Q} [\mathcal{A}_p(Av(x, \cdot)^\top)(x')].$$

Straightforward calculations along with  $A^\top A = I_m$  show that

$$\begin{aligned} \mathcal{A}_p(Av(x, \cdot)^\top)(x') &= (A^\top s_p(x))^\top (A^\top s_p(x')) k(A^\top x, A^\top x') + (A^\top s_p(x))^\top \nabla_{x_2} k(A^\top x, A^\top x') \\ &\quad + \nabla_{x_1} k(A^\top x, A^\top x')^\top A^\top s_p(x') + \text{Tr}(\nabla_{x_1, x_2} k(A^\top x, A^\top x')), \end{aligned}$$

where  $\nabla_{x_i}$  denotes the gradient of  $k$  with respect to the  $i^{\text{th}}$  argument, and  $\nabla_{x_1, x_2} k$  is the matrix with  $ij$ -th entry  $\partial^2 k / \partial x_i \partial x_j$ . This completes the proof since it is the integrand of (7). □

*Proof of Theorem 2.* First, if  $p = q$  then Lemma 5 immediately tells us that the discrepancy is zero since  $\delta_{p,q} = 0$ . Now suppose the discrepancy is zero. Take any projector  $A$ , then  $\|A^\top x - A^\top x'\|_2 \leq \|x - x'\|_2$  since  $AA^\top$  is a projection matrix. As  $k$  is a radial kernel we know that  $\Psi$  is monotonically decreasing (Scovel et al., 2010, Theorem 1.1). Therefore  $k(Ax, Ax') \geq k_d(x, x')$ , where  $k_d = \Psi(\|x - x'\|_2)$ . By Lemma 5,

$$0 = \text{GKSD}(Q, P) \geq \sup_{[A] \in \text{Gr}(d,m)} \mathbb{E}_{x, x' \sim Q} [(A^\top \delta_{p,q}(x))^\top A^\top \delta_{p,q}(x') k_d(x, x')]. \quad (24)$$

Now choose  $A$  as the matrix formed by stacking the first  $m$  canonical basis elements. Since  $k$  is characteristic and radial, we know that  $\Psi(\|x - x'\|_2) = L_\nu(\|x - x'\|_2)$  where  $L_\nu$  is the Laplace transform of a measure  $\nu$  on  $[0, \infty)$  with full support (Sriperumbudur et al., 2011, Proposition 5). Therefore  $k_d(x, x') = \Psi(\|x - x'\|_2) = L_\nu(\|x - x'\|_2)$  and  $k_d$  is also characteristic and hence integrally strictly positive definite.

Using this we can conclude from (24) that  $A^\top \delta_{p,q} = 0$ , which implies  $A^\top s_q = A^\top s_p$  so the first  $m$  entries of the two score functions are the same. Repeat this argument by setting  $A$  to be the next  $m$  canonical basis elements until we have checked all canonical basis elements. We then deduce that the score functions are equal in all coordinates, which implies  $p = q$  as required.  $\square$

## D PROOF of PROPOSITION 3

We prove that, provided the densities of the distributions  $Q, P$  differ only in a subspace identified by a projector  $A_0$  of rank  $m$ , the Riemannian gradient  $\text{grad}\alpha([A])$  is zero at  $A = A_0$ . That is, the optimal subspace  $[A_0]$  is indeed a solution sought by the objective

$$\sup_{[A] \in \text{Gr}(d,m)} \text{KSD}_A(Q, P),$$

where  $\text{KSD}_A(Q, P)$  is given by (9).

In the rest of this section, we first state an assumption on the form of the kernel in D.1 that is mild but can greatly simplify the proof. In D.2, we then provide intuition on the assumption of the decomposibility of the densities of  $Q$  and  $P$  and give a concrete example. The proof is presented in D.3.

### D.1 Assumption on the Kernel

In Theorem 2, we have assumed the kernel takes the form  $k(u, v) = \Psi(\|u - v\|_2)$ , i.e. it is a *radial kernel* defined in Sriperumbudur et al. (2011). To simplify the form of  $\text{grad}\alpha([A])$  in our proof, we introduce the following reformulation of  $\Psi$  and impose a smoothness condition.

**Assumption 1** (Smoothness of kernel). *There exists a continuously differentiable function  $\Phi : \mathbb{R} \rightarrow \mathbb{R}$  for which  $\Psi(s) = \Phi(s^2)$  for all  $s \geq 0$ . That is, the kernel has the form  $k(u, v) = \Phi(\|u - v\|_2^2)$ .*

**Remark 4.** *The continuous differentiability condition on  $\Phi$  is for convenience only and is mild. In particular, both Gaussian RBF and Inverse Multiquadric kernels satisfy Assumption 1.*

### D.2 Assumption on the Densities

Intuitively, the decomposition (13) and (14) states that the the candidate density  $q$  and the target  $p$  are distinct only in the subspace  $[A_0]$ . It can be viewed as a type of sparsity assumption. Both the multimodal mixture and the X-shaped mixture examples in the experiments satisfy this assumption, where it is easy to check that such functions  $q^m, p^m, \xi$  can be found with  $[A_0] = \{ae_1 + be_2 : a, b \in \mathbb{R}\}$  being the subspace spanned by the first two canonical basis vectors  $e_1, e_2$ . See Appendix E.2 and E.3 for the precise definition of the target densities in these two examples.

More generally, this assumption holds when  $q(x)$  is a prior, and  $p(x) \propto q(x)f(y|P_0x)$  is the posterior induced by observational data  $y$  and a likelihood function  $f(y|A_0A_0^\top x)$  that depends on  $x$  only through its projection onto the subspace  $[A_0]$ . When the columns of  $A_0$  are canonical basis vectors, this means the data-generating process only depends on a subset of the parameters  $x$ . A similar setting is also considered by Chen and Ghattas (2020), where the authors argue that such a likelihood function is a reasonable approximation when the variation of the likelihood is negligible outside of some eigenspace of a gradient information matrix.

### D.3 Proof

We begin by deriving an expression for the Euclidean and Riemannian gradients of  $\alpha([A])$  with respect to a projector  $A$ .

**Lemma 6.** *Suppose that Assumption 1 holds, and that  $p, q$  satisfy the same assumptions in Theorem 2. Then, given a projector  $A$ , the Euclidean gradient of  $\alpha([A])$  at  $A$  is*

$$\begin{aligned} \nabla\alpha([A]) = & 2\mathbb{E}_{x, x' \sim Q} [k(A^\top x, A^\top x') \delta_{p,q}(x) \delta_{p,q}(x')^\top A \\ & + \Phi'(\|A^\top x - A^\top x'\|_2^2) \delta_{p,q}(x)^\top A A^\top \delta_{p,q}(x') (x - x')(x - x')^\top A], \end{aligned} \quad (25)$$

where the expectation is taken over i.i.d. copies  $x, x' \sim Q$ ,  $\Phi'$  is the first derivative of  $\Phi$ , and  $[\nabla\alpha([A])]_{ij} = \frac{\partial}{\partial A_{ij}}\alpha([A])$ . Moreover, letting  $\Pi_A := I_d - A A^\top$ , the Riemannian gradient of  $\alpha([A])$  at  $A$  is

$$\begin{aligned} \text{grad}\alpha([A]) = & \Pi_A \nabla\alpha([A]) \\ = & 2\Pi_A \mathbb{E}_{x, x' \sim Q} [k(A^\top x, A^\top x') \delta_{p,q}(x') \delta_{p,q}(x)^\top A \\ & + \Phi'(\|A^\top x - A^\top x'\|_2^2) \delta_{p,q}(x')^\top A A^\top \delta_{p,q}(x) (x - x')(x - x')^\top A]. \end{aligned} \quad (26)$$

*Proof of Lemma 6.* Letting  $G$  be the right hand side of (25). We will show that, for any  $t \in \mathbb{R}$  and projectors  $A, B$  of rank  $m$ ,

$$\alpha([A + tB]) - \alpha([A]) = t \cdot \text{Tr}(G^\top B) + \mathcal{O}(t^2).$$

Let  $w := x - x'$  so that  $k(A^\top x, A^\top x') = \Phi(\|A^\top w\|_2^2)$  for a projector  $A$ . Starting with  $\alpha([A])$  in the form of (23),

$$\begin{aligned} \alpha([A + tB]) - \alpha([A]) = & \mathbb{E}_{x, x' \sim Q} \left[ \delta_{p,q}(x)^\top (A + tB) (A + tB)^\top \delta_{p,q}(x') \Phi \left( \|(A + tB)^\top w\|_2^2 \right) \right] \\ & - \mathbb{E}_{x, x' \sim Q} \left[ \delta_{p,q}(x)^\top A A^\top \delta_{p,q}(x') \Phi \left( \|A^\top w\|_2^2 \right) \right] \\ = & \mathbb{E}_{x, x' \sim Q} \left[ \delta_{p,q}(x)^\top A A^\top \delta_{p,q}(x') \left( \Phi \left( \|(A + tB)^\top w\|_2^2 \right) - \Phi \left( \|A^\top w\|_2^2 \right) \right) \right] \\ & + 2t \mathbb{E}_{x, x' \sim Q} \left[ \delta_{p,q}(x)^\top A B^\top \delta_{p,q}(x') \Phi \left( \|(A + tB)^\top w\|_2^2 \right) \right] \\ & + \mathcal{O}(t^2). \end{aligned} \quad (27)$$

We now seek a linear approximation of the function  $A \mapsto \Phi(\|A^\top w\|_2^2)$ , where  $A$  is a projector of rank  $m$ . Since  $\Phi$  is continuously differentiable by assumption, it admits a first-order Taylor expansion at  $\|A^\top w\|_2^2$ , which gives

$$\Phi \left( \|(A + tB)^\top w\|_2^2 \right) - \Phi \left( \|A^\top w\|_2^2 \right) \quad (28)$$

$$\begin{aligned} = & \Phi \left( \|A^\top w\|_2^2 + t \left( 2w^\top A B^\top w + t \|B^\top w\|_2^2 \right) \right) - \Phi \left( \|A^\top w\|_2^2 \right) \\ = & \Phi \left( \|A^\top w\|_2^2 \right) + t \Phi' \left( \|A^\top w\|_2^2 \right) \left( 2w^\top A B^\top w + t \|B^\top w\|_2^2 \right) - \Phi \left( \|A^\top w\|_2^2 \right) + \mathcal{O}(t^2) \\ = & t \Phi' \left( \|A^\top w\|_2^2 \right) (2w^\top A B^\top w) + \mathcal{O}(t^2). \end{aligned} \quad (29)$$

The rest of the proof for the form of the Euclidean gradient follows by substituting (29) into (27) and using

properties of the trace operator:

$$\begin{aligned}
 \alpha([A + tB]) - \alpha([A]) &= 2t\mathbb{E}_{x,x'\sim Q} \left[ \Phi' \left( \|A^\top w\|_2^2 \right) \delta_{p,q}(x)^\top AA^\top \delta_{p,q}(x') w^\top AB^\top w \right] \\
 &\quad + 2t\mathbb{E}_{x,x'\sim Q} \left[ \delta_{p,q}(x)^\top AB^\top \delta_{p,q}(x') \Phi \left( \|A^\top w\|_2^2 \right) \right] + \mathcal{O}(t^2) \\
 &= 2t\mathbb{E}_{x,x'\sim Q} \left[ \Phi' \left( \|A^\top w\|_2^2 \right) \delta_{p,q}(x)^\top AA^\top \delta_{p,q}(x') \text{Tr}(ww^\top AB^\top) \right] \\
 &\quad + 2t\mathbb{E}_{x,x'\sim Q} \left[ \Phi \left( \|A^\top w\|_2^2 \right) \text{Tr}(\delta_{p,q}(x')\delta_{p,q}(x)^\top AB^\top) \right] + \mathcal{O}(t^2) \\
 &= t \cdot \text{Tr} \left( 2\mathbb{E}_{x,x'\sim Q} \left[ \Phi' \left( \|A^\top w\|_2^2 \right) \delta_{p,q}(x)^\top AA^\top \delta_{p,q}(x') ww^\top AB^\top \right] \right) \\
 &\quad + t \cdot \text{Tr} \left( 2\mathbb{E}_{x,x'\sim Q} \left[ \Phi \left( \|A^\top w\|_2^2 \right) \delta_{p,q}(x')\delta_{p,q}(x)^\top AB^\top \right] \right) + \mathcal{O}(t^2) \\
 &= t \cdot \text{Tr}(G^\top B) + \mathcal{O}(t^2),
 \end{aligned}$$

where in the last equality we have substituted the definition of  $G$  and used the fact  $\text{Tr}(GB^\top) = \text{Tr}(G^\top B)$ . Finally, it can be shown that  $\text{Gr}(d, m)$  is a Riemannian submanifold of the Euclidean manifold, so the Riemannian gradient is the orthogonal projection of the Euclidean gradient to the tangent spaces (see e.g. Boumal (2020, Proposition 3.53)). That is,  $\text{grad}\alpha([A]) = (I_d - AA^\top)\nabla\alpha([A])$ , so (26) follows.  $\square$

*Proof of Proposition 3.* Decomposing (26) into two terms, we have

$$\begin{aligned}
 \text{grad}\alpha([A]) &= 2\Pi_A\mathbb{E}_{x,x'\sim Q} [k(A^\top x, A^\top x')\delta_{p,q}(x')\delta_{p,q}(x)^\top A] \\
 &\quad + 2\Pi_A\mathbb{E}_{x,x'\sim Q} [\Phi'(\|A^\top x - A^\top x'\|_2^2) \delta_{p,q}(x')^\top AA^\top \delta_{p,q}(x)(x - x')(x - x')^\top A], \quad (30)
 \end{aligned}$$

where  $\Pi_A = I_d - AA^\top$ . We will show that both terms equal to zero at  $A = A_0$ . Firstly, define  $\tilde{\xi}_{\Pi_{A_0}}(x) = \xi(\Pi_{A_0}x)$  and  $P_0 = A_0A_0^\top$ . The score function of the candidate density is

$$\begin{aligned}
 \nabla_x \log q(x) &= \nabla_x \log(q^m \circ P_0)(x) + \nabla_x \log \tilde{\xi}_{\Pi_{A_0}}(x) \\
 &= P_0 \nabla_x \log q^m(P_0x) + \nabla_x \log \tilde{\xi}_{\Pi_{A_0}}(x),
 \end{aligned}$$

where  $\nabla_x \log q^m(P_0x)$  denotes the gradient of  $\log q^m$  evaluated at  $P_0x$ , and where the second line follows from the chain rule and the symmetry of  $P_0$ . Similarly, the score function of the target density takes the form

$$\nabla_x \log p(x) = P_0 \nabla_x \log q^m(P_0x) + \nabla_x \log \xi(x).$$

Taking the difference,

$$\delta_{p,q}(x) = P_0(\nabla_x \log p^m(P_0x) - \nabla_x \log q^m(P_0x)).$$

Since  $\Pi_{A_0}P_0 = 0$ , we conclude that  $\Pi_{A_0}\delta_{p,q}(x) = 0$ , and the first term of (30) is zero.

For the second term, we define  $a := \delta_{p,q}^\top(x')A_0A_0^\top\delta_{p,q}(x)$  and  $b := \Phi'(\|A^\top x - A^\top x'\|_2^2)$ . Since both the terms  $\delta_{p,q}(x)$  and  $A_0^\top x - A_0^\top x' = A_0^\top(P_0x - P_0x')$  are deterministic given  $(P_0x, P_0x')$ , so are  $a$  and  $b$ . We further introduce the notations  $w := P_0x, w' := P_0x'$  and  $v := \Pi_{A_0}x, v' := \Pi_{A_0}x'$ . Since  $q^m, p^m$  and  $\xi$  are unnormalised densities over  $\mathbb{R}^d$ , we can invoke the tower rule to simplify the second term as

$$\begin{aligned}
 \Pi_{A_0}\mathbb{E}_{x,x'\sim Q} [ab(x - x')(x - x')^\top A_0] &= \mathbb{E}_{x,x'\sim Q} [ab(\Pi_{A_0}x - \Pi_{A_0}x')(P_0x - P_0x')^\top A_0] \\
 &= \mathbb{E}_{w,w'\sim q^m} [\mathbb{E}_{v,v'\sim q^\perp} [ab(v - v')(w - w')^\top A_0 \mid w, w']] \\
 &= \mathbb{E}_{w,w'\sim q^m} [ab\mathbb{E}_{v,v'\sim q^\perp} [v - v' \mid w, w'] (w - w')^\top A_0],
 \end{aligned}$$

where the first line follows from the fact that  $A_0 = A_0A_0^\top A_0 = P_0A_0$ . Since  $x, x'$  are i.i.d. copies, we have  $\mathbb{E}_{v,v'\sim q^\perp} [v - v' \mid w, w'] = \mathbb{E}_{q^\perp} [v \mid w] - \mathbb{E}_{q^\perp} [v' \mid w'] = 0$ , so the second term of (30) is also zero. This shows that  $\text{grad}\alpha([A]) = 0$  at  $A = A_0$ , as required.  $\square$

## E EXPERIMENTAL DETAILS

In this section, we provide implementation details for the experiments in Section 6.



## E.1 Setups

**Learning Rates** The best step sizes of the gradient method for the particle update in all methods, and for the projector or slice update in GSVG and S-SVG, are selected over a grid of values from  $10^{-4}$  to 1 such that they give the minimal energy distance between the particle estimation and the ground truth. This is done for the 50-dimensional multimodal example only, and the same set of learning rates is then used for all experiments. This is done for all methods to ensure a fair comparison.

**Temperature  $T$**  For the temperature parameter  $T$  in (21), we adopt an annealing scheme where we start with  $T = T_0$  and gradually increment it to a large value,  $T_{\text{large}}$ . Intuitively, setting  $T$  small would lead to faster convergence, but hinders exploration of the maxima of  $\alpha_t([A]) = \text{KSD}_A(Q_t, P)$ . On the other hand, a large  $T$  means the system (20) is effectively sampling projectors at random from the Grassmann manifold, thus allowing better exploration of the modes at the expense of a slower convergence. The annealing scheme allows trade-offs between such exploration and exploitation. Starting from  $T_0$ , at each iteration  $t$ , we multiply  $T$  by 10 if the change  $|\gamma_t - \gamma_{t-1}|$  in the *particle-averaged magnitude* (Zhuo et al., 2018) of the update function, defined as  $\gamma_t := \frac{1}{N} \sum_{i=1}^N \left\| \sum_{l=1}^M \widehat{\phi}_{A_t, l}(x_i^t) \right\|_{\infty}$ , is less than a pre-specified threshold, where  $N$  is the sample size,  $M$  is the number of projectors, and the update function  $\widehat{\phi}_{A_t, l}$  is defined in (17). In the experiments, we used  $T_0 = 10^{-4}$  and  $T_{\text{large}} = 10^6$ , and set the threshold to be  $10^{-4}M$ . This choice of the threshold is motivated by the dependence of the GSVG particle update on the number of projectors.

**Initialisation of Projectors** The projectors in GSVG are initialized with matrices formed by one-hot vectors. This is similar to the setup in S-SVG, where the slices are initialized to be the canonical basis elements. When more than one projector is used in GSVG, the projection dimensions,  $m$ , are kept fixed for all projectors for simplicity. The number of projectors,  $M$ , is chosen such that the projection dimensions sum up to  $d$  and is capped at 20, where  $d$  is the dimension of the problem. That is,  $M = \min(20, \lfloor d/m \rfloor)$ , where  $\lfloor a \rfloor$  denotes the largest integer smaller than or equal to  $a$ .

**Other Setups** For SVG, we follow exactly the same setups as in Liu and Wang (2016). For S-SVG, we follow mostly the same setups in Algorithm 2 and Appendix G of Gong et al. (2021), except we updated the slices at every step instead of only when the particles have moved a sufficiently far distance from the previous step, as suggested in their paper. As noted in their paper, this trick is to prevent over-fitting of the slices to small samples, which is unlikely an issue in our experiments due to the large sample size we chose. Also, we remark that, for GSVG, the coupled ODE-SDE system introduced in Section 4 suggests a principled way to balance the evolution of the particles and projectors in GSVG by varying the temperature  $T$ .

We ran each experiment with 20 random seeds, except Bayesian logistic regression which was repeated 10 times. The mean and 95% confidence intervals are reported in all figures. A No-U-Turn Sampler (NUTS) (Hoffman et al., 2014) is used as a gold standard in the conditioned diffusion and Bayesian logistic regression experiments. Each method ran for 2000 iterations, which we found to be sufficient for convergence (see Figure 9).

## E.2 Multimodal Mixture

The mean vectors of the multimodal mixture are defined as  $\mu_k = (\sqrt{5} \cos(2k\pi/K + \pi/4), \sqrt{5} \sin(2k\pi/K + \pi/4), 0, \dots, 0)^{\top} \in \mathbb{R}^d$ , for  $k = 1, \dots, 4$ ; see Figure 4(a). That is, the target in the first two dimensions is a mixture of 4 Gaussian distributions with modes equally spaced on a circle of radius  $\sqrt{5}$ , whilst in the other coordinates it is a standard Gaussian.

## E.3 X-Shaped Mixture

The covariance matrices in the X-shaped mixture example have the following correlated block diagonal structure:

$$\Sigma_1 = \begin{pmatrix} 1 & \delta & \mathbf{0} \\ \delta & 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & I_{d-2} \end{pmatrix}, \quad \Sigma_2 = \begin{pmatrix} 1 & -\delta & \mathbf{0} \\ -\delta & 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & I_{d-2} \end{pmatrix},$$

where  $\delta = 0.95$  controls the correlation between the first two dimensions.

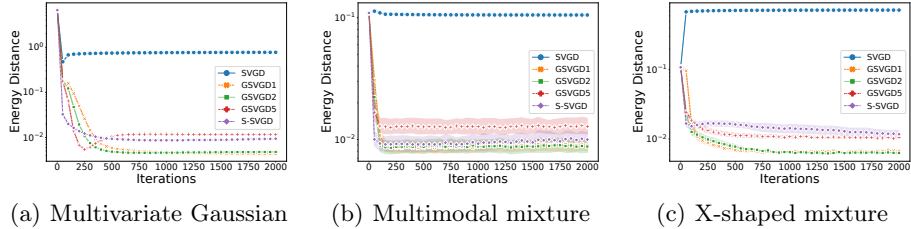


Figure 9: Convergence of particles in the (a) multivariate Gaussian, (b) multimodal mixture and (c) X-shaped experiments. The dimensionality is 50 in all cases.

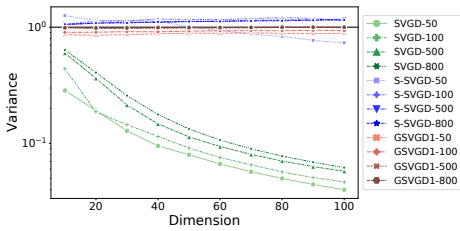


Figure 10: Variance estimates for different sample sizes, where the true value is shown by the black solid line. SVGD-50 means 50 particles are used to estimate the variance.

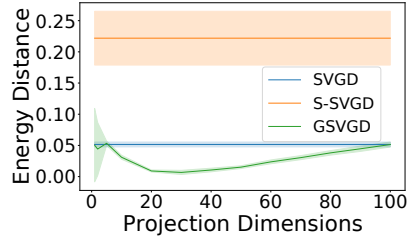


Figure 11: Conditioned Diffusion Process: Energy distance between HMC and GSVGD solutions against different projection dimensions. SVGD and S-SVGD are also included for comparison.

### E.4 Bayesian Logistic Regression

We follow the same setting as Liu and Wang (2016) by placing a Gaussian prior  $p_0(w|\alpha) = \mathcal{N}(w; 0, \alpha^{-1})$  on the regression weights  $w$  with  $p_0(\alpha) = \text{Gamma}(\alpha; 1, 0.01)$ . We are interested in the posterior  $p(x|D)$  of  $x = [w, \alpha]$ . The model was tested on a subset of the Covtype dataset. The original dataset consists of 581,012 data points and 54 features with binary labels, which is too large for NUTS. As a result, a subset of 1,000 randomly selected data points was used, although other methods that are more scalable to data size, such as the stochastic gradient Langevin dynamics (SGLD) of Welling and Teh (2011) or the HMC-ECS of Dang et al. (2019), could have been used to allow inference on the entire dataset.

To evaluate the results, we used both the energy distance and the covariance estimation error  $\|\hat{\Sigma} - \Sigma\|_2$ , where  $\|A\|_2 = \sqrt{\sum_{i,j=1}^d A_{ij}^2}$  is the Frobenius norm of a matrix  $A \in \mathbb{R}^{d \times d}$ , and  $\hat{\Sigma}, \Sigma$  are respectively the sample covariance matrices of the particle estimation and of a HMC run treated as the ground truth. The covariance estimation error was used because we found that the energy distance is not sensitive enough to differences in the second moments. From Figure 7, we observe an overestimated covariance between distinct variables compared with the HMC reference, as shown by the large values (in absolute term) of the off-diagonal entries. On the other hand, GSVGD projecting onto  $m = 10$  dimensional subspaces showed greater alignment, despite having a similar energy distance as S-SVGD (Figure 6). Again, we emphasise that the performance of GSVGD deteriorates when the projection dimension is extreme, e.g. when  $m$  equals 1 or the full dimension of the problem ( $m = 55$ ).

## F SUPPLEMENTARY FIGURES

In this section, we include supplementary figures for the experiments in Section 6 as well as some ablation studies.

### F.1 Synthetic Experiments

Figure 9 supplements Figure 5, and shows the convergence of SGVD, S-SVGD and GSVGD to the target distribution in the three synthetic experiments in Section 6. We can see that the smallest energy distance between the particle estimation and the ground truth is achieved by GSVGD1 and GSVGD2, followed by S-SVGD and

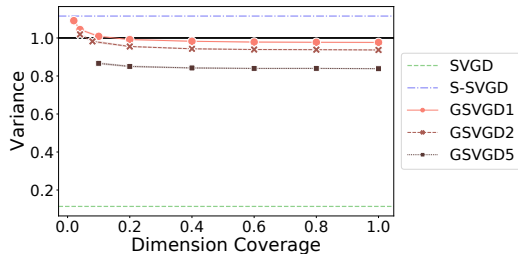


Figure 12: Ablation study on the number of projectors  $M$ . The target is a 50-dimensional multivariate Gaussian. Black solid line marks the true value.

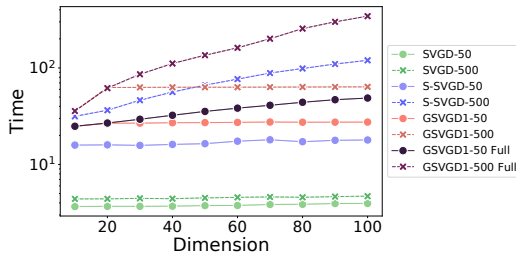


Figure 13: Run time (in seconds) against the dimension of the multivariate Gaussian target.

GSVGD5. GSVGD with any projection dimension has comparable per-iteration convergence rate as S-SVGD except for the multivariate Gaussian target. We believe this is because of the extra noise introduced via the coupled ODE-SDE system described in Section 4, which added unnecessary complexities for this simple, single-mode target.

## F.2 Conditioned Diffusion Process

Figure 11 plots the energy distance between the particle estimation of each method and a HMC sampler against different projection dimensions,  $m$ . We see that there exists a feasible region of projection dimensions where GSVGD achieves a greater degree of agreement with the reference HMC, compared with SVGD and S-SVGD. This is expected given the problem setup: the forcing term is a Brownian motion with covariance  $C(t, t') = \min(t, t')$ , and hence is correlated between time points. This correlation means the solution state  $u$  will admit a low-dimensional structure.

Although in this case an appropriate choice of projection dimension could be inferred from the problem setup, selecting the best projection dimension for a general problem is an open question. Nonetheless, the shape of the curves in Figure 11 suggests that a feasible region should not cover the extreme ends.

## F.3 Ablation Studies

**Sample Size** To evaluate the influence of sample size on the performance of GSVGD, we show the dimension-averaged variance estimates for the multivariate Gaussian target. We follow the same experimental setup in Section 6: for each of  $m = 1, 2, 5$ , we run each method for 2000 iterations and compute the dimension-averaged marginal variance of the particle estimation of the last iteration. Figure 10 shows the estimates for GSVGD1, SVGD and S-SVGD when the target is 50-dimensional. GSVGD with other projection dimensions is not included in the figure for clarity. We see that GSVGD1 remains robust with the sample size and achieves accurate estimation. In contrast, SVGD severely underestimates the variance in high dimensions regardless of the sample size, and S-SVGD shows slight overestimation.

**Number of Projectors** In Figure 12, we perform an ablation study of the performance of GSVGD with respect to the number of projectors  $M$  used in Algorithm 1. The target is again the 50-dimensional multivariate Gaussian distribution. The estimate of the dimension-averaged marginal variance is plotted against the *dimension coverage*, which is defined as  $(M * m)/d$ ; that is, it is the total ranks of the  $M$  projectors as a proportion of dimension. For a given  $m$ , we choose  $M = m, 10, 20, \dots, 50$ . We also include the result for  $M = 2$  and 5 when  $m = 1$ , and  $M = 2$  when  $m = 2$ .

We observe that the variance estimate of GSVGD improves as  $M$  increases until the dimension coverage reaches roughly 0.2, beyond which the performance stabilises. This suggests that using more projectors does not necessarily improve the estimation quality of GSVGD. It also justifies our choice of  $M = \min(20, \lfloor d/m \rfloor)$  (corresponding to a dimension coverage of 0.4 for  $m = 1$ , 0.8 for  $m = 2$  and 1.0 for  $m = 5$ ). We remark that the estimated variance of GSVGD1 and GSVGD2 are consistently better than S-SVGD and SVGD for all choices of  $M$ .

**Time Complexities** We compare the run time of GSVGD against its competitors in Figure 13. In GSVGD, we used  $M = \min(20, \lfloor d/m \rfloor)$  projectors of rank  $m = 1$  (GSVGD1-50, GSVGD1-500, where GSVGD1-50 means 50 particles were used). We have capped  $M$  to be at most 20 instead of setting it to the dimensionality  $d$  because this was the default setting in all experiments; however, since S-SVGD uses all  $d$  slices, we also include the run time of GSVGD1 when using  $M = d$  projectors (GSVGD1-50 Full, GSVGD1-500 Full) for a fair comparison.

We see that the run time of GSVGD1-50 Full and GSVGD1-500 Full grows an order-of-magnitude faster than S-SVGD even though they use the same projection dimension and number of projectors/slices. This is due to the extra matrix computation in the particle update of GSVGD, which in turn arises from the fact that the projectors for the score functions are allowed to vary in GSVGD but not in S-SVGD. In particular, the time complexity for a fixed sample size is  $\mathcal{O}(d^3)$  for GSVGD1-50 Full and GSVGD1-500 Full, in contrast to  $\mathcal{O}(d^2)$  for S-SVGD. However, capping  $M$  to be at most 20 can greatly reduce the run time of GSVGD in high dimensions. In fact, with such a choice of  $M$  the computational complexity reduces to  $\mathcal{O}(d^2)$ , as can be seen from the curves GSVGD1-50 and GSVGD1-500.

Finally, the computational burden of the projector update in GSVGD1-50 Full and GSVGD1-500 Full ( $\mathcal{O}(d^2)$ ) is less than that of the particle update. Hence, the main computational bottleneck of GSVGD arises from the particle update (17).