# STEIN VARIATIONAL GRADIENT DESCENT ON INFINITE-DIMENSIONAL SPACE AND APPLICATIONS TO STATISTICAL INVERSE PROBLEMS*

JUNXIONG JIA†, PEIJUN LI‡, AND DEYU MENG§

**Abstract.** In this paper, we propose an infinite-dimensional version of the Stein variational gradient descent (iSVGD) method for solving Bayesian inverse problems. The method can generate approximate samples from posteriors efficiently. Based on the concepts of operator-valued kernels and vector-valued reproducing kernel Hilbert spaces, a rigorous definition is given for the infinite-dimensional objects, e.g., the Stein operator, which are proved to be the limit of finite-dimensional ones. Moreover, a more efficient iSVGD with preconditioning operators is constructed by generalizing the change of variables formula and introducing a regularity parameter. The proposed algorithms are applied to an inverse problem of the steady state Darcy flow equation. Numerical results confirm our theoretical findings and demonstrate the potential applications of the proposed approach in the posterior sampling of large-scale nonlinear statistical inverse problems.

**Key words.** statistical inverse problems, Bayes' method, variational inference method, Stein variational gradient descent, machine learning

**AMS subject classifications.** 65L09, 49N45, 62F15

**1. Introduction.** Driven by rapid algorithmic development and a steady increase of computer power, the Bayesian approach has enjoyed great popularity for solving inverse problems over the last decade. By transforming inverse problems into statistical inference problems, the approach provides a general framework to quantify uncertainties [1]. The posterior distribution automatically delivers an estimate of the statistical uncertainty in the reconstruction, and hence suggests "confidence" intervals that allow to reject or accept scientific hypotheses [44]. It has been widely used in many applications, e.g., artifact detecting in medical imaging [64].

The approach begins with establishing an appropriate Bayes model. When the parameters are in a finite-dimensional space, the finite-dimensional Bayesian method can be employed [56]. A comprehensive account of the finite-dimensional theory can be found in [32]. When the inferred parameters are in the infinite-dimensional space, the problems are more challenging since the Lebesgue measure cannot be defined rigorously in this case [15]. Recently, some attempts have been made to handle the issue. For example, a general framework was designed for the Bayesian formula and the general theory was applied to inverse problems of fluid mechanic equations [12]. A survey can be found in [53] on the basic framework of the infinite-dimensional Bayes'

†School of Mathematics and Statistics, Xi'an Jiaotong University, Xi'an 710049, China (jjx323@xjtu.edu.cn).

‡Department of Mathematics, Purdue University, West Lafayette, Indiana, 47907, USA (lipeijun@math.purdue.edu).

§Corresponding author. School of Mathematics and Statistics and Ministry of Education Key Lab of Intelligent Networks and Network Security, Xi'an Jiaotong University, Xi'an, 710049, China, and Macau Institute of Systems Engineering, Macau University of Science and Technology, Taipa, Macau (dymeng@mail.xjtu.edu.cn).

approach for solving inverse problems. Inverse problems of partial differential equations (PDEs) often involve infinite-dimensional spaces, and the infinite-dimensional Bayes' theory has recently attracted more attention [5, 13, 24, 45, 46].

As pointed out in [1], one of the challenges for the Bayesian approach is how to effectively extract information encoded in the posterior probability measure. To overcome the difficulty, the two main strategies are the point estimate method and the sampling method. The former is to find the maximum a posteriori (MAP) estimate which is equivalent to solve an optimization problem [5, 24]. In some situations, the MAP estimates are more desirable and computationally feasible than the entire posterior distribution [26, 55]. However, the point estimates cannot convey uncertainty information and are usually recognized as an incomplete Bayes' method. The sampling type methods, such as the well known Markov chain Monte Carlo (MCMC), are often used to extract posterior information. They are well studied in the finite-dimensional setting [35]. Although the MCMC methods are accurate and effective, they are usually not robust under mesh refinement [13]. Multiple dimension-independent MCMC-type algorithms have been proposed [13, 14, 20, 51]. However, these MCMC-type algorithms are computationally too expensive to be adopted in such an application as seismic exploration [21].

The finite-dimensional problems have been extensively studied and many efficient algorithms have been developed to quantify uncertainties effectively. In particular, the variational inference (VI) methods have been broadly investigated in machine learning [3, 43, 62, 63]. Under the mean-field assumption, the linear inverse problems were examined in [30, 29] by using a hierarchical formulation with Gaussian and centered-t noise distribution. The skewed-t noise distribution was considered for a similar setting in [23]. A new type of variational inference algorithm, called the Stein variational gradient descent (SVGD), was proposed in [39]. The method can achieve reliable uncertainty estimation by efficiently using an interacting repulsive mechanism. The SVGD has shown to be a fast and flexible method for solving challenging machine learning problems and inverse problems of PDEs [10, 11].

Compared with the finite-dimensional problems, the infinite-dimensional problems are much less studied for the variational inference (VI). When the approximate measures are restricted to be Gaussian, the novel Robbins–Monro algorithm was developed in [45, 46] from a calculus-of-variations viewpoint. It was shown in [54] that the Kullback–Leibler (KL) divergence between the stochastic processes is equal to the supremum of the KL divergence between the measures restricted to finite marginals. Meanwhile, they developed a VI method for functions parameterized by Bayesian neural networks. Under the classical mean-field assumption, a general VI framework defined on separable Hilbert spaces was proposed recently in [28]. A function space particle optimization method including the SVGD was developed in [61] to solve the particle optimization directly in the space of functions. The function space algorithm was also employed to solve computer vision problems, e.g., the context of semantic segmentation and depth estimation [9]. However, the function spaced SVGD assumes that the random functions can be parameterized by a finite number of parameters, e.g., parameterized by some neural networks [61]. Hence, the probability measures on functions are implicitly defined through the probability distributions of a finite number of parameters, instead of the expected infinite-dimensional function space.

This work concerns inverse problems of PDEs imposed on infinite dimensional function spaces. Motivated by the preconditioned Crank–Nicolson (pCN) algorithm [13], we aim to construct the SVGD on separable Hilbert spaces with random functions. Throughout, the iSVGD stands for SVGD defined on the infinite-dimensional

function space. The goal is to develop algorithms defined on Hilbert spaces and lay a foundation for appropriate discretizations. It contains three contributions:

(1) We investigate the Bayesian formula in infinite-dimensional spaces. The rigorous definition of the SVGD on separable Hilbert spaces is provided, the Stein operator is defined and the corresponding optimization problem on some Hilbert spaces is considered, and the finite-dimensional problem is proved to converge to the infinite-dimensional counterpart;

(2) By introducing vector-valued reproducing kernel Hilbert space (RKHS) and operator-valued kernel, we improve the iSVGD with precondition information (e.g., Hessian information operator), which can accelerate the iSVGD algorithm significantly. This is the first work on such an iSVGD algorithm with precondition information;

(3) Explicit numerical strategies are designed by using the finite-element approach. Through theoretical analysis and numerical examples, we demonstrate that the regularity parameter $s$ introduced in the abstract theory (see Assumptions 5 and 7 in Section 3.2) should belong to the interval $(0, 0.5)$ and be close to 0.5. The scalability of the algorithm depends only on the scalability of the forward and adjoint PDE solvers. Hence, the algorithm is applicable to solve large-scale inverse problems of PDEs.

The paper is organized as follows. The SVGD in finite-dimensional spaces is introduced in Section 2. Section 3 is devoted to the construction of the iSVGD. The basic concepts of operator-valued kernels and Hilbert scales are briefly reviewed; the Stein operator is defined on separable Hilbert spaces; it is shown that the infinite-dimensional version is indeed equivalent to the finite-dimensional version in some limit sense; Based on the Stein operator and the theory of reproducing kernel Hilbert space (RKHS), the update direction of the iSVGD is derived; In addition, the change of variables is studied and the iSVGD is constructed with preconditioning operators; a preliminary theoretical study is given for the corresponding continuous equations. In Section 4, the algorithm is applied to solve an inverse problem governed by the steady state Darcy flow equation. The paper is concluded with some general remarks and directions for future work in Section 5.

**2. A short review of SVGD.** Let $\mathcal{H}$ be a separable Hilbert space endowed with the Borel $\sigma$-algebra $\mathcal{B}(\mathcal{H})$. Denote by $\mathcal{G}$, $u$, and $\boldsymbol{d}$ the solution operator of some PDE, the model parameter, and the observation, respectively. We assume that $u \in \mathcal{H}$ and $\boldsymbol{d} \in \mathbb{R}^{N_d}$ with $N_d$ being a positive integer. The observation $\boldsymbol{d}$ is related to $\mathcal{G}(u)$ and the random noise $\boldsymbol{\epsilon}$ through some functions [32], e.g., the additive noise model or the multiplicative noise model. We refer to Section 4 for a specific example.

For statistical inverse problems, it is usually required to find a probability measure $\mu^{\boldsymbol{d}}$ on $\mathcal{H}$, which is known as the posterior probability measure and is specified by its density with respect to a prior probability measure $\mu_0$. The Bayesian formula on a Hilbert space is defined by

$$(1) \qquad \frac{d\mu^{\boldsymbol{d}}}{d\mu_0}(u) = \frac{1}{Z_{\boldsymbol{d}}} \exp\Big( - \Phi(u; \boldsymbol{d}) \Big),$$

where $\Phi \in C(\mathcal{H} \times \mathbb{R}^{N_d}; \mathbb{R})$ and $\exp(-\Phi(u; \boldsymbol{d}))$ is integrable with respect to $\mu_0$. The constant $Z_{\boldsymbol{d}}$ is chosen to ensure that $\mu^{\boldsymbol{d}}$ is indeed a probability measure. The prior measure $\mu_0 := \mathcal{N}(0, \mathcal{C}_0)$ is assumed to be a Gaussian measure defined on $\mathcal{H}$ with $\mathcal{C}_0$ being a self-adjoint, positive definite, and trace class operator. Let $(\lambda_k, \varepsilon_k)_{k=1}^{\infty}$ be the eigensystem of $\mathcal{C}_0$ satisfying $\mathcal{C}_0 \varepsilon_k = \lambda_k^2 \varepsilon_k$. Denote by $P^N$ and $Q^N$ the orthogonal

projections of $\mathcal{H}$ onto $X^N := \text{span}\{\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_N\}$ and $X^\perp := \text{span}\{\varepsilon_{N+1}, \varepsilon_{N+2}, \ldots\}$, respectively. Clearly, we have $Q^N = \text{Id} - P^N$. Let $u^N := P^N u \in X^N$ and $u^\perp := Q^N u \in X^\perp$. Define $\mathcal{C}_0^N = P^N \mathcal{C}_0 P^N$ and let $\mu_0^N = \mathcal{N}(0, \mathcal{C}_0^N)$ be a finite-dimensional Gaussian measure defined on $X^N$. Then an approximate measure $\mu^{\boldsymbol{d}N}$ on $X^N$ can be defined by

$$(2) \qquad \frac{d\mu^{\boldsymbol{d}N}}{d\mu_0^N}(u^N) = \frac{1}{Z_{\boldsymbol{d}}^N} \exp\Big(-\Phi(u^N; \boldsymbol{d})\Big),$$

where

$$Z_{\boldsymbol{d}}^N = \int_{X^N} \exp\Big(-\Phi(u^N; \boldsymbol{d})\Big)\mu_0^N(du^N).$$

Some more properties of the above approximate measure can be found in [16, Subsection 5.6]. The probability measure $\mu^{\boldsymbol{d}N}$ can be written as the pushforward of the posterior measure $\mu^{\boldsymbol{d}}$ on $\mathbb{R}^N$, i.e., $\mu^{\boldsymbol{d}N} = P_\#^N \mu^{\boldsymbol{d}} := \mu^{\boldsymbol{d}} \circ (P^N)^{-1}$. Hence the measure $\mu^{\boldsymbol{d}N}$ has a Lebesgue density denoted by $p^{\boldsymbol{d}N}$ with the following form:

$$(3) \qquad p^{\boldsymbol{d}N}(u^N) \propto \exp\Big(-\Phi(u^N; \boldsymbol{d}) - \frac{1}{2}\|u^N\|_{\mathcal{C}_0^N}^2\Big),$$

where $\|\cdot\|_{\mathcal{C}_0^N}$ represents $\|(\mathcal{C}_0^N)^{-1/2}\cdot\|_{\ell^2}$ with $\|\cdot\|_{\ell^2}$ standing for the usual $\ell^2$-norm. Obviously, the target distribution $\mu^{\boldsymbol{d}N}$ is the solution to the optimization problem defined on the set $\mathcal{P}_2(\mathbb{R}^N)$ of probability measures $\nu$ such that $\int \|u^N\|^2 d\nu^N(u^N) < \infty$ by:

$$(4) \qquad \min_{\nu^N \in \mathcal{P}_2(\mathbb{R}^N)} \text{KL}(\nu^N \| \mu^{\boldsymbol{d}N}),$$

where KL denotes the Kullback-Leibler (KL) divergence.

Now, we present the Stein variational gradient descent (SVGD) algorithm. Denote $\text{KL}(\cdot\|\mu^{\boldsymbol{d}N}) : \mathcal{P}_2(\mathbb{R}^N) \to [0, +\infty)$ as the functional $\nu^N \mapsto \text{KL}(\nu^N\|\mu^{\boldsymbol{d}N})$. In order to obtain samples from $\mu^{\boldsymbol{d}N}$, the SVGD applies a gradient descent-like algorithm to the functional $\text{KL}(\cdot\|\mu^{\boldsymbol{d}N})$. The standard gradient descent algorithm in the Wasserstein space applied to $\text{KL}(\cdot\|\mu^{\boldsymbol{d}N})$, at each iteration $\ell \geq 0$, is

$$(5) \qquad \nu_{\ell+1}^N = \Big(\text{Id} - \epsilon\nabla\log\Big(\frac{d\nu_\ell^N}{d\mu^{\boldsymbol{d}N}}\Big)\Big)_\# \nu_\ell^N,$$

where $\epsilon > 0$ is the step size. This corresponds to a forward Euler discretization of the gradient flow of $\text{KL}(\cdot\|\mu^{\boldsymbol{d}N})$ with respect to Stein geometry [18]. Instead of the Wasserstein gradient $\nabla\log\big(d\nu_\ell^N/d\mu^{\boldsymbol{d}N}\big)$ used in (5), the SVGD uses $P_{\nu_\ell^N}\nabla\log\big(d\nu_\ell^N/d\mu^{\boldsymbol{d}N}\big)$ to generate the following iteration:

$$(6) \qquad \nu_{\ell+1}^N = \Big(\text{Id} - \epsilon P_{\nu_\ell^N}\nabla\log\Big(\frac{d\nu_\ell^N}{d\mu^{\boldsymbol{d}N}}\Big)\Big)_\# \nu_\ell^N,$$

where $P_{\nu_\ell^N}$ is the same as that in Subsection 3.1 of [33]. Let $\mathcal{H}_K^N$ be an $N$-dimensional reproducing kernel Hilbert space (RKHS) [52] with the kernel function $K : \mathbb{R}^N \times \mathbb{R}^N \to \mathbb{R}$. To define $P_{\nu_\ell^N}$ rigorously, it is necessary to introduce the kernel integral operator based on the kernel function $K$, which will not be used in the rest of the paper. Hence,

we omit it and refer to [33] for the details. The reason for introducing the operator $P_{\nu_\ell^N}$ is that we have

(7) $\quad P_{\nu_\ell^N} \nabla \log \left( \dfrac{d\nu_\ell^N}{d\mu^{\boldsymbol{d}N}} \right)(\cdot) = -\mathbb{E}_{u^N \sim \nu_\ell^N} \left[ K(u^N, \cdot) \nabla_{u^N} \log p^{\boldsymbol{d}N}(u^N) + \nabla_{u^N} K(u^N, \cdot) \right]$

under some mild conditions. For every $\ell \geq 0$, let $u^{N,\ell}$ be distributed according to $\nu_\ell^N$. Using (6)–(7), we obtain a particle update scheme

(8) $\qquad\qquad\qquad u^{N,\ell+1} = u^{N,\ell} + \epsilon \phi_\ell^N(u^{N,\ell}),$

where

(9) $\qquad\qquad \phi_\ell^{N*}(\cdot) = \mathbb{E}_{u^N \sim q_\ell^N} \left[ K(u^N, \cdot) \nabla_{u^N} \log p^{\boldsymbol{d}N}(u^N) + \nabla_{u^N} K(u^N, \cdot) \right].$

The basic SVGD algorithm is given in Algorithm 1. Inspired by applications in machine learning, the SVGD type algorithms have been widely studied over the last few years [17, 18, 33, 38, 39, 42].

---

**Algorithm 1** Finite-dimensional Stein variational gradient descent

---

**Input:** A target probability measure with density function $p^{\boldsymbol{d}N}(u^N)$ and a set of particles $\{u_i^{N,0}\}_{i=1}^m$.
**Output:** A set of particles $\{u_i^N\}_{i=1}^m$ that approximates the target probability measure.
**for** iteration $\ell$ do

$$u_i^{N,\ell+1} \longleftarrow u_i^{N,\ell} + \epsilon_\ell \phi^*(u_i^{N,\ell}),$$

where

$$\phi^*(u^N) = \frac{1}{m} \sum_{j=1}^m \left[ K(u_j^{N,\ell}, u^N) \nabla_{u_j^{N,\ell}} \log p^{\boldsymbol{d}N}(u_j^{N,\ell}) + \nabla_{u_j^{N,\ell}} K(u_j^{N,\ell}, u^N) \right],$$

and $\epsilon_\ell$ is the step size at the $\ell$-th iteration.
**end for**

---

**3. SVGD on separable Hilbert spaces.** This section is devoted to the construction of iSVGD and the preconditioning operators. The corresponding continuity equations are provided for a preliminary theoretical study of the method.

**3.1. Hilbert scale and vector-valued RKHS.** For constructing iSVGD, we need to characterize the smoothness of functions that belong to some infinite dimensional spaces. The Sobolev spaces are usually employed to characterize the smoothness of functions. However, for presenting a general theory, we introduce the Hilbert scales defined by the prior covariance operator [19]. The reason is that different covariance operators employed in practical problems lead to the same form of Hilbert scales. However, they are related to different Sobolev spaces. Hence, the same form of the general theory can be flexibly adapted to different practical problems.

Let $\mathcal{C}_0 : \mathcal{H} \to \mathcal{H}$ be the covariance operator introduced in Section 2. Denote by $\mathcal{D}(\mathcal{C}_0)$ and $\mathcal{R}(\mathcal{C}_0)$ the domain and range of $\mathcal{C}_0$, respectively. Let $\mathcal{H} = \overline{\mathcal{R}(\mathcal{C}_0)} \oplus \mathcal{R}(\mathcal{C}_0)^\perp = \overline{\mathcal{R}(\mathcal{C}_0)}$ (the closure of $\mathcal{R}(\mathcal{C}_0)$). It is clear to note that $\mathcal{C}_0^{-1}$ is a densely

defined, unbounded, symmetric and positive-definite operator in $\mathcal{H}$. Let $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ and $\| \cdot \|_{\mathcal{H}}$ be the inner product and norm defined on the Hilbert space $\mathcal{H}$, respectively. Define the Hilbert scales $(\mathcal{H}^t)_{t \in \mathbb{R}}$ with $\mathcal{H}^t := \overline{\mathcal{S}_f}^{\| \cdot \|_{\mathcal{H}^t}}$, where

$$\mathcal{S}_f := \bigcap_{n=0}^{\infty} \mathcal{D}(\mathcal{C}_0^{-n}), \quad \langle u, v \rangle_{\mathcal{H}^t} := \langle \mathcal{C}_0^{-t/2} u, \mathcal{C}_0^{-t/2} v \rangle_{\mathcal{H}}, \quad \|u\|_{\mathcal{H}^t} := \left\| \mathcal{C}_0^{-t/2} u \right\|_{\mathcal{H}}.$$

The norms defined above possess the following properties (cf. [19, Proposition 8.19]).

LEMMA 1. *Let $(\mathcal{H}^t)_{t \in \mathbb{R}}$ be the Hilbert scale induced by the operator $\mathcal{C}_0$ given above. Then the following assertions hold:*

1. *Let $-\infty < s < t < \infty$. Then the space $\mathcal{H}^t$ is densely and continuously embedded into $\mathcal{H}^s$.*
2. *If $t \geq 0$, then $\mathcal{H}^t = \mathcal{D}(\mathcal{C}_0^{-t/2})$, and $\mathcal{H}^{-t}$ is the dual space of $\mathcal{H}^t$.*
3. *Let $-\infty < q < r < s < \infty$ then the interpolation inequality $\|u\|_{\mathcal{H}^r} \leq \|u\|_{\mathcal{H}^q}^{\frac{s-r}{s-q}} \|u\|_{\mathcal{H}^s}^{\frac{r-q}{s-q}}$ holds when $u \in \mathcal{H}^s$.*

Now, we introduce some basic notations of vector-valued reproducing kernel Hilbert space (RKHS). The following definition concerns the Hilbert space adjoint opertor [50].

DEFINITION 2. *Let $\mathcal{X}$ and $\mathcal{Y}$ be Banach spaces, and $T$ be a bounded linear operator from $\mathcal{X}$ to $\mathcal{Y}$. The **Banach space adjoint** of $T$, denoted by $T'$, is the bounded linear operator from $\mathcal{Y}^*$ to $\mathcal{X}^*$ and is defined by $(T'\ell)(u) = \ell(Tu)$ for all $\ell \in \mathcal{Y}^*$, $u \in \mathcal{X}$. Let $\mathcal{X}$ and $\mathcal{Y}$ be Hilbert spaces, and $C_1 : \mathcal{X} \to \mathcal{X}^*$ be the map that assigns to each $u \in \mathcal{X}$, the bounded linear functional $\langle u, \cdot \rangle_{\mathcal{X}}$ in $\mathcal{X}^*$. Let $C_2 : \mathcal{Y} \to \mathcal{Y}^*$ be defined similarly as $C_1$. Then the **Hilbert space adjoint** of $T$ is a map $T^* : \mathcal{Y} \to \mathcal{X}$ given by $T^* = C_1^{-1} T' C_2$.*

Next, we introduce operator-valued positive definite kernels, which constitute the framework for specifying vector-valued RKHS. Following Kadri et al. [31] to avoid topological and measurability issues, we focus on separable Hilbert spaces with reproducing operator-valued kernels whose elements are continuous functions. Denote by $\mathcal{X}$ and $\mathcal{Y}$ the separable Hilbert spaces and by $\mathcal{L}(\mathcal{X}, \mathcal{Y})$ the set of bounded linear operators from $\mathcal{X}$ to $\mathcal{Y}$. When $\mathcal{X} = \mathcal{Y}$, we write $\mathcal{L}(\mathcal{Y}, \mathcal{Y})$ briefly as $\mathcal{L}(\mathcal{Y})$.

DEFINITION 3. *(Operator-valued kernels) An $\mathcal{L}(\mathcal{Y})$-valued kernel $\boldsymbol{K}$ on $\mathcal{X} \times \mathcal{X}$ is an operator $\boldsymbol{K}(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \to \mathcal{L}(\mathcal{Y})$;*

1. *$\boldsymbol{K}$ is Hermitian if $\forall u, v \in \mathcal{X}$, $\boldsymbol{K}(u, v) = \boldsymbol{K}(v, u)^*$;*
2. *$\boldsymbol{K}$ is nonnegative on $\mathcal{X}$ if it is Hermitian and for every natural number $r$ and all $\{(u_i, v_i)_{i=1,\ldots,r}\} \in \mathcal{X} \times \mathcal{Y}$, the matrix with $ij$-th entry $\langle \boldsymbol{K}(u_i, u_j) v_i, v_j \rangle_{\mathcal{Y}}$ is nonnegative (positive-definite).*

DEFINITION 4. *(Vector-valued RKHS) Let $\mathcal{X}$ and $\mathcal{Y}$ be separable Hilbert spaces. A Hilbert space $\mathcal{F}$ of operators from $\mathcal{X}$ to $\mathcal{Y}$ is called a reproducing kernel Hilbert space if there is a nonnegative $\mathcal{L}(\mathcal{Y})$-valued kernel $\boldsymbol{K}$ on $\mathcal{X} \times \mathcal{X}$ such that*

1. *the operator $v \longmapsto \boldsymbol{K}(u, v) g$ belongs to $\mathcal{F}$ for all $v, u \in \mathcal{X}$ and $g \in \mathcal{Y}$;*
2. *for every $f \in \mathcal{F}$, $u \in \mathcal{X}$ and $g \in \mathcal{Y}$, we have $\langle f(u), g \rangle_{\mathcal{Y}} = \langle f(\cdot), \boldsymbol{K}(u, \cdot) g \rangle_{\mathcal{F}}$.*

Throughout the paper, we assume that the kernel $\boldsymbol{K}$ is *locally bounded and separately continuous*, which guarantee that $\mathcal{F}$ is a subspace of $C(\mathcal{X}, \mathcal{Y})$ (the vector space of continuous operators from $\mathcal{X}$ to $\mathcal{Y}$). If the kernel $\boldsymbol{K}$ is nice enough [7, 8], then it is the reproducing kernel of some Hilbert space $\mathcal{F}$.

Since the kernel is an important part of the SVGD, we provide some intuitive ideas about the operator-valued kernel. Let $u, v \in \mathcal{H}$ and $h > 0$ be a positive constant. To construct the infinite-dimensional SVGD, we may introduce a scalar-valued kernel $K(u, v) := \exp\left(-\frac{1}{h}\|u - v\|_{\mathcal{H}}^2\right)$ and consider the operator-valued kernel

$$\boldsymbol{K}(u, v) = K(u, v)\mathrm{Id}. \tag{10}$$

For example, we can take $\mathcal{H} = L^2(\Omega)$ with $\Omega$ being a bounded open domain and have

$$\|u - v\|_{\mathcal{H}}^2 = \int_{\Omega} |u(x) - v(x)|^2 dx. \tag{11}$$

However, for solving inverse problems of PDEs, it is useful to introduce some preconditioning operators which require to consider operator-valued kernels. Here, we illustrate this by a simple example. Let the prior measure $\mu_0 = \mathcal{N}(0, (\mathrm{Id} - \Delta)^{-2})$, where $\Delta$ is the Dirichlet Laplace operator and $\mathcal{H} = L^2(\Omega)$. Intuitively we have $\mathcal{H}^1 \approx H^2(\Omega)$, where $H^2(\Omega)$ is the usual Sobolev space. By the theory of Gaussian measures [48], we approximately have $\mu_0(H^2(\Omega)) = 0$ (not rigorously correct). Inspired by the pCN algorithm [13], we may choose the preconditioning operator $T = \mathrm{Id} - \Delta$. If we choose the Gaussian kernel as (10), then the transformed kernel function becomes

$$\boldsymbol{K}(u, v) = \exp\left(-\frac{1}{h}\|T(u - v)\|_{L^2}^2\right) T^{-1}(T^{-1})^*, \tag{12}$$

which is approximately equal to

$$\boldsymbol{K}(u, v) \approx \exp\left(-\frac{1}{h}\|u - v\|_{H^2}^2\right)(\mathrm{Id} - \Delta)^{-2}. \tag{13}$$

Obviously, the kernel function equals to zero when $u - v$ does not belong to $H^2(\Omega)$, i.e., $\|u - v\|_{H^2} < \infty$ when $u - v \in H^2(\Omega)$. Hence, the kernel function takes nonzero values and the algorithms can work only if the differences of any two particles reside in a measure zero set. In our opinion, this restriction seems too strong in the infinite-dimensional setting to make the particles over concentrated (see our numerical example in Section 4 to demonstrate this in details).

Based on the above discussion, we may introduce a parameter $s$ and have an approximate transformed kernel

$$\boldsymbol{K}(u, v) \approx \exp\left(-\frac{1}{h}\|u - v\|_{H^{2-2s}}^2\right)(\mathrm{Id} - \Delta)^{-2}. \tag{14}$$

However, to achieve this, we should not choose the original kernel (the kernel is not transformed by the operator $T$) to be the usual scalar-valued kernel. The original kernel may be chosen as $\boldsymbol{K}_0(u, v) = K_0(u, v)(\mathrm{Id} - \Delta)^{-2s}$, where $K_0(u, v) := e^{-\frac{1}{h}\|u - v\|_{L^2}}$ with $h > 0$ being a positive constant. In this setting, the preconditioning operator can be chosen as $T := (\mathrm{Id} - \Delta)^{1-s}$. These intuitive ideas indicate that it is necessary to construct the infinite-dimensional SVGD based on the more involved operator-valued kernel theory.

**3.2. iSVGD.** In this subsection, we present an infinite-dimensional version of the SVGD, i.e., iSVGD. For a function $u$, denote by $D_u$ and $D_{u_k}$ the Fréchet derivative and the directional derivative in the $k$th direction, respectively. For simplicity of

notation, we shall use $D$ and $D_k$ instead of $D_u$ and $D_{u_k}$, and write $\Phi(u; \boldsymbol{d})$ as $\Phi(u)$. Let

(15) $$V(u) = \Phi(u) + \frac{1}{2}\|u\|_{\mathcal{H}^1}^2,$$

where the potential functional $\Phi$ is required to satisfy the following assumptions.

ASSUMPTION 5. *Let $\mathcal{X}$ and $\mathcal{H}$ be two separable Hilbert spaces. For $s \in [0, 1]$, we assume $\mathcal{H}^{1-s} \subset \mathcal{X} \subset \mathcal{H}$. Let $M_1 \in \mathbb{R}^+$ be a positive constant. For each $u \in \mathcal{X} \subset \mathcal{H}$, we introduce $D\Phi : \mathcal{X} \to \mathcal{X}^*$ and $D^2\Phi : \mathcal{X} \to \mathcal{L}(\mathcal{X}, \mathcal{X}^*)$, then the functional $\Phi : \mathcal{X} \to \mathbb{R}$ satisfies*

$$-M_1 \leq \Phi(u) \leq M_2(\|u\|_{\mathcal{X}}),$$
$$\|D\Phi(u)\|_{\mathcal{X}^*} \leq M_3(\|u\|_{\mathcal{X}}),$$
$$\|D^2\Phi(u)\|_{\mathcal{L}(\mathcal{X}, \mathcal{X}^*)} \leq M_4(\|u\|_{\mathcal{X}}),$$

*where $M_2(\cdot)$, $M_3(\cdot)$, and $M_4(\cdot)$ are some monotonic non-decreasing functions.*

The above assumption is a local version of [16, Assumption 4], which can be verified for many problems, e.g., the Darcy flow model (Theorem 17 in Section 4). We now optimize $\phi$ in the unit ball of a general vector-valued RKHS $\mathcal{H}_{\boldsymbol{K}}$ with an operator valued kernel $\boldsymbol{K}(u, u') \in \mathcal{L}(\mathcal{Y})$:

(16) $$\phi_{\boldsymbol{K}}^* = \underset{\phi \in \mathcal{H}_{\boldsymbol{K}}}{\arg\max} \left\{ \mathbb{E}_{u \sim \mu}[\mathcal{S}\phi(u)], \text{ s.t. } \|\phi\|_{\mathcal{H}_K} \leq 1 \text{ and } D\phi : \mathcal{X} \to \mathcal{L}_1(\mathcal{X}, \mathcal{Y}) \right\},$$

where $\mathcal{S}$ is the generalized Stein operator defined formally as follows:

(17) $$\mathcal{S}\phi(u) = -\langle DV(u), \phi(u) \rangle_{\mathcal{Y}} + \sum_{k=1}^{\infty} D_k \langle \phi(u), e_k \rangle_{\mathcal{Y}},$$

and $\mathcal{L}_1(\mathcal{X}, \mathcal{Y})$ denotes the set of all trace class operators from $\mathcal{X}$ to $\mathcal{Y}$. For the convergence of the infinite sum, we illustrate it in Theorem 9. Here, $\{e_k\}_{k=1}^{\infty}$ stands for an orthonormal basis of space $\mathcal{Y}$ and $\mu$ is a probability measure defined on $\mathcal{H}$. Moreover, we assume that $\phi : \mathcal{X} \to \mathcal{Y}$ is Fréchet differentiable, and the derivative is continuous to ensure the validity of (16).

REMARK 6. *In the finite-dimensional case, the operator $D\phi(u)$ naturally belongs to $\mathcal{L}_1(\mathcal{X}, \mathcal{Y})$ (cf. [15, Appendix C]).*

The following assumption is also needed for the operator-valued kernels, which include many useful kernels, e.g., the radial basis function (RBF) kernel.

ASSUMPTION 7. *Let $\mathcal{X}$, $\mathcal{Y}$, and $\mathcal{H}$ be three separable Hilbert spaces. For $s \in [0, 1]$, we assume that $\mathcal{H}^{-s-1} \subset \mathcal{Y}$ and*

(18) $$\sup_{u \in \mathcal{X}} \|\boldsymbol{K}(u, u)\|_{\mathcal{L}(\mathcal{Y})} < \infty.$$

REMARK 8. *We mention that Condition (18) holds for the bounded scalar-valued kernel functionals since a scalar-valued kernel functional can be seen as a scalar-valued kernel functional composite with an identity operator as demonstrated in (10).*

To illustrate (16) and (17), we prove Theorem 9. For each particle $u$, we assume that $u \in \mathcal{H}^{1-s}$, which is based on the following two considerations:

- The SVGD with one particle is an optimization algorithm for finding maximum a posterior (MAP) estimate. The MAP estimate belongs to the separable Hilbert space $\mathcal{H}^1$.
- For the prior probability measure, the space $\mathcal{H}^1$ has zero measure [15]. Intuitively, if all particles belong to $\mathcal{H}^1$, the particles tend to concentrate around a small set that leads to unreliable estimates of statistical quantities. Hence, we may assume that the particles belong to a larger space containing $\mathcal{H}^1$.

THEOREM 9. *The generalized Stein operator (17) defined on $\mathcal{Y}$ can be obtained by taking $N \to \infty$ in the following finite-dimensional Stein operator:*

$$(19) \qquad \mathcal{S}^N \phi^N(u^N) = -\langle DV(u^N), \phi^N(u^N) \rangle_{\mathcal{Y}} + \sum_{k=1}^{N} D_k \langle \phi^N(u^N), e_k \rangle_{\mathcal{Y}},$$

*where $\phi^N := P^N \circ \phi$.*

*Proof.* By straightforward calculations, we have

$$
\begin{aligned}
\mathcal{S}\phi(u) - \mathcal{S}^N \phi^N(u^N) = & -\left( \langle DV(u), \phi(u) \rangle_{\mathcal{Y}} - \langle DV(u^N), \phi^N(u^N) \rangle_{\mathcal{Y}} \right) \\
(20) \qquad & + \left( \sum_{k=1}^{\infty} D_k \langle \phi(u), e_k \rangle_{\mathcal{Y}} - \sum_{k=1}^{N} D_k \langle \phi^N(u^N), e_k \rangle_{\mathcal{Y}} \right) \\
= & -\mathrm{I} + \mathrm{II}.
\end{aligned}
$$

For term I, we have

$$
\begin{aligned}
(21) \qquad \mathrm{I} = & \langle D(V(u) - V(u^N)), \phi^N(u^N) \rangle_{\mathcal{Y}} + \langle DV(u), \phi(u) - \phi^N(u^N) \rangle_{\mathcal{Y}} \\
= & \mathrm{I}_1(N) + \mathrm{I}_2(N).
\end{aligned}
$$

For term $\mathrm{I}_1(N)$, we find that

$$(22) \qquad \mathrm{I}_1(N) = \langle D(\Phi(u) - \Phi(u^N)), \phi^N(u^N) \rangle_{\mathcal{Y}} + \langle \mathcal{C}_0^{-1/2}(u - u^N), \mathcal{C}_0^{-1/2} \phi^N(u^N) \rangle_{\mathcal{Y}},$$

where the second term on the right-hand side is understood as the white noise mapping [48]. According to Assumptions 5 and 7, we know that

$$
\begin{aligned}
(23) \qquad \lim_{N \to \infty} \| D(\Phi(u) - \Phi(u^N)) \|_{\mathcal{Y}} \leq & \lim_{N \to \infty} C \| D(\Phi(u) - \Phi(u^N)) \|_{\mathcal{H}^{-1-s}} \\
\leq & \lim_{N \to \infty} C \| D(\Phi(u) - \Phi(u^N)) \|_{\mathcal{H}^{-1+s}} \\
\leq & \lim_{N \to \infty} C M_4 (2\|u\|_{\mathcal{X}}) \|u - u^N\|_{\mathcal{H}^{1-s}} = 0,
\end{aligned}
$$

where $C$ is a generic constant that can be different from line to line. Hence, we obtain

$$(24) \qquad \lim_{N \to \infty} \langle D(\Phi(u) - \Phi(u^N)), \phi^N(u^N) \rangle_{\mathcal{Y}} = 0.$$

Taking $u_m \in \mathcal{H}^2$ such that $u_m \to u$ in $\mathcal{H}^{1-s}$, we have

$$
\begin{aligned}
\langle \mathcal{C}_0^{-1/2}(u - u^N), \mathcal{C}_0^{-1/2} \phi^N(u^N) \rangle_{\mathcal{Y}} = & \lim_{m \to \infty} \langle \mathcal{C}_0^{-1/2}(u_m - u_m^N), \mathcal{C}_0^{-1/2} \phi^N(u^N) \rangle_{\mathcal{Y}} \\
= & \lim_{m \to \infty} \langle P^N \mathcal{C}_0^{-1}(u_m - u_m^N), \phi(u^N) \rangle_{\mathcal{Y}} \\
= & \lim_{m \to \infty} \langle \phi(\cdot), \boldsymbol{K}(u^N, \cdot) P^N \mathcal{C}_0^{-1}(u_m - u_m^N) \rangle_{\mathcal{H}_{\boldsymbol{K}}}.
\end{aligned}
$$

361    As for the last term in the above equality, we have the following estimates:

362    $\langle \phi(\cdot), \boldsymbol{K}(u^N, \cdot) P^N \mathcal{C}_0^{-1}(u_m - u_m^N) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \leq$

363    $\quad \langle \phi(\cdot), \phi(\cdot) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \langle \boldsymbol{K}(u^N, \cdot) P^N \mathcal{C}_0^{-1}(u_m - u_m^N), \boldsymbol{K}(u^N, \cdot) P^N \mathcal{C}_0^{-1}(u_m - u_m^N) \rangle_{\mathcal{H}_{\boldsymbol{K}}}$

364    $\quad \leq \langle \phi(\cdot), \phi(\cdot) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \langle P^N \boldsymbol{K}(u^N, u^N) P^N \mathcal{C}_0^{-1}(u_m - u_m^N), \mathcal{C}_0^{-1}(u_m - u_m^N) \rangle_{\mathcal{Y}}$

365    $\quad \leq C \langle \phi(\cdot), \phi(\cdot) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \| \mathcal{C}_0^{-1}(u_m - u_m^N) \|_{\mathcal{Y}}^2$

366
367    $\quad \leq C \langle \phi(\cdot), \phi(\cdot) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \| \mathcal{C}_0^{-\frac{1-s}{2}}(u_m - u_m^N) \|_{\mathcal{H}}^2.$

368    Replacing $u_m - u_m^N$ by $(u_m - u_m^N) - (u - u^N)$, we deduce

369    $\langle \mathcal{C}_0^{-1/2}(u - u^N), \mathcal{C}_0^{-1/2} \phi^N(u^N) \rangle_{\mathcal{Y}} = \lim_{m \to \infty} \langle \phi(\cdot), \boldsymbol{K}(u^N, \cdot) P^N \mathcal{C}_0^{-1}(u_m - u_m^N) \rangle_{\mathcal{H}_{\boldsymbol{K}}}$

370
371    (25)    $\qquad\qquad\qquad = \langle \phi(\cdot), \boldsymbol{K}(u^N, \cdot) P^N \mathcal{C}_0^{-1}(u - u^N) \rangle_{\mathcal{H}_{\boldsymbol{K}}}.$

372    Hence, we obtain

$$\lim_{N \to \infty} \langle \mathcal{C}_0^{-1/2}(u - u^N), \mathcal{C}_0^{-1/2} \phi^N(u^N) \rangle_{\mathcal{Y}}$$

373    (26)
$$= \lim_{N \to \infty} \langle \phi(\cdot), \boldsymbol{K}(u^N, \cdot) P^N \mathcal{C}_0^{-1}(u - u^N) \rangle_{\mathcal{H}_{\boldsymbol{K}}}$$
$$\leq \lim_{N \to \infty} \langle \phi(\cdot), \phi(\cdot) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \langle P^N \boldsymbol{K}(u^N, u^N) P^N \mathcal{C}_0^{-1}(u - u^N), \mathcal{C}_0^{-1}(u - u^N) \rangle_{\mathcal{Y}}$$

374
$$\leq C \langle \phi(\cdot), \phi(\cdot) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \lim_{N \to \infty} \| \mathcal{C}_0^{-\frac{1-s}{2}}(u - u^N) \|_{\mathcal{H}}^2 = 0.$$

375    Plugging (24) and (26) into (22), we arrive at $\lim_{N \to \infty} I_1(N) = 0$. For term $I_2(N)$, it
376    can be decomposed as follows:

377
378    (27)    $I_2(N) = \langle D\Phi(u), \phi(u) - \phi^N(u^N) \rangle_{\mathcal{Y}} + \langle \mathcal{C}_0^{-1/2} u, \mathcal{C}_0^{-1/2}(\phi(u) - \phi^N(u^N)) \rangle_{\mathcal{Y}}.$

379    It follows from the continuity of $\phi$ that we have $\lim_{N \to \infty} \langle D\Phi(u), \phi(u) - \phi^N(u^N) \rangle_{\mathcal{Y}} = 0$.
380    Using similar estimates as those for deriving (25), we obtain

381    (28)    $\langle \mathcal{C}_0^{-1/2} u, \mathcal{C}_0^{-1/2}(\phi(u) - \phi^N(u^N)) \rangle_{\mathcal{Y}}$
382    $\qquad\qquad = \langle \phi(\cdot), \boldsymbol{K}(u, \cdot) \mathcal{C}_0^{-1} u \rangle_{\mathcal{H}_{\boldsymbol{K}}} - \langle \phi(\cdot), \boldsymbol{K}(u^N, \cdot) P^N \mathcal{C}_0^{-1} u \rangle_{\mathcal{H}_{\boldsymbol{K}}}.$

383    By the continuity of $\boldsymbol{K}(\cdot, \cdot)$, we obtain

384    (29)    $\lim_{N \to \infty} \langle \mathcal{C}_0^{-1/2} u, \mathcal{C}_0^{-1/2}(\phi(u) - \phi^N(u^N)) \rangle_{\mathcal{Y}} = 0.$
385

386    Now, we conclude that $\lim_{N \to \infty} I_2(N) = 0$. For term II, we have

387    (30)    $\mathrm{II} = \sum_{k=1}^{N} D_k \langle \phi(u) - \phi(u^N), e_k \rangle_{\mathcal{Y}} + \sum_{k=N+1}^{\infty} D_k \langle \phi(u), e_k \rangle_{\mathcal{Y}}.$
388

389    Let $\{\varphi_k\}_{k=1}^{\infty}$ be an orthonormal basis in $\mathcal{X}$, and then we have

390    (31)    $\sum_{k=N+1}^{\infty} D_k \langle \phi(u), e_k \rangle_{\mathcal{Y}} = \sum_{k=N+1}^{\infty} \langle D\phi(u)\varphi_k, e_k \rangle_{\mathcal{Y}} \to 0 \quad \text{as } N \to \infty,$
391

where we use the condition $D\phi(u) \in \mathcal{L}_1(\mathcal{X}, \mathcal{Y})$. For the first term on the right-hand side of (30), we find that

$$\text{(32)} \qquad \sum_{k=1}^{N} D_k \langle \phi(u) - \phi(u^N), e_k \rangle_{\mathcal{Y}} = \sum_{k=1}^{N} \langle (D\phi(u) - D\phi(u^N))\varphi_k, e_k \rangle_{\mathcal{Y}}.$$

Due to the continuity of the Fréchet derivative of $\phi$, we know that the above summation goes to 0 as $N \to \infty$. Combining the estimates of I and II, we complete the proof. □

The following theorem gives explicitly the iSVGD update directions that are essential for the construction of iSVGD.

THEOREM 10. *Let $\boldsymbol{K}(\cdot, \cdot) : \mathcal{X}^2 \to \mathcal{L}(\mathcal{Y})$ be a positive definite kernel that is Fréchet differentiable on both variables. In addition, we assume that*

$$\text{(33)} \qquad \mathbb{E}_{u \sim \mu} \Big[ D_{u'} \boldsymbol{K}(u, u') \mathcal{C}_0^{-1/2} g + \sum_{k=1}^{\infty} D_{u_k} D_{u'} \boldsymbol{K}(u, u') e_k \Big]$$

*belongs to $\mathcal{L}_1(\mathcal{X}, \mathcal{Y})$ for each $u' \in \mathcal{X}$ and $g \in \mathcal{H}^{-s}$. Then, the optimal $\phi_{\boldsymbol{K}}^*$ in (16) is*

$$\text{(34)} \qquad \phi_{\boldsymbol{K}}^*(\cdot) \propto \mathbb{E}_{u \sim \mu} \Big[ \boldsymbol{K}(u, \cdot)(-D\Phi(u) - \mathcal{C}_0^{-1} u) + \sum_{k=1}^{\infty} D_{u_k} \boldsymbol{K}(u, \cdot) e_k \Big],$$

*where $\{e_k\}_{k=1}^{\infty}$ is an orthonormal basis of $\mathcal{Y}$ and the term $\boldsymbol{K}(u, \cdot) \mathcal{C}_0^{-1} u$ is understood in the following limiting sense:*

$$\text{(35)} \qquad \boldsymbol{K}(u, \cdot) \mathcal{C}_0^{-1} u := \lim_{m \to \infty} \boldsymbol{K}(u, \cdot) \mathcal{C}_0^{-1} u_m.$$

*Here the limit is taken in $\mathcal{H}_{\boldsymbol{K}}$ and $\{u_m\}_{m=1}^{\infty} \subset \mathcal{H}^2$ such that $\|\mathcal{C}_0^{-\frac{1-s}{2}} (u_m - u)\|_{\mathcal{H}} \to 0$ as $m \to \infty$.*

*Proof.* First, by taking $\phi(u)$ as an element in $\mathcal{H}_{\boldsymbol{K}}$, we have

$$\text{(36)} \qquad \langle DV(u), \phi(u) \rangle_{\mathcal{Y}} = \langle D\Phi(u), \phi(u) \rangle_{\mathcal{Y}} + \langle \mathcal{C}_0^{-1/2} u, \mathcal{C}_0^{-1/2} \phi(u) \rangle_{\mathcal{Y}} = \text{I} + \text{II},$$

where term II is understood as the white noise mapping. For term I, we have

$$\text{(37)} \qquad \text{I} = \langle \phi(\cdot), \boldsymbol{K}(u, \cdot) D\Phi(u) \rangle_{\mathcal{H}_{\boldsymbol{K}}},$$

where the proposition (2) in Definition 4 is employed. For term II, we take $u_m \in \mathcal{H}^2$ such that $\lim_{m \to \infty} \|\mathcal{C}_0^{-\frac{1-s}{2}} (u_m - u)\|_{\mathcal{H}} = 0$. It is clear to note that

$$\text{(38)} \qquad \langle \mathcal{C}_0^{-1/2} u_m, \mathcal{C}_0^{-1/2} \phi(u) \rangle_{\mathcal{Y}} = \langle \mathcal{C}_0^{-1} u_m, \phi(u) \rangle_{\mathcal{Y}} = \langle \phi(\cdot), \boldsymbol{K}(u, \cdot) \mathcal{C}_0^{-1} u_m \rangle_{\mathcal{H}_{\boldsymbol{K}}}.$$

Because

$$|\langle \phi(\cdot), \boldsymbol{K}(u, \cdot) \mathcal{C}_0^{-1} u_m \rangle_{\mathcal{H}_{\boldsymbol{K}}} - \langle \phi(\cdot), \boldsymbol{K}(u, \cdot) \mathcal{C}_0^{-1} u \rangle_{\mathcal{H}_{\boldsymbol{K}}}|^2$$

$$\leq \langle \phi(\cdot), \phi(\cdot) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \langle \boldsymbol{K}(u, \cdot) \mathcal{C}_0^{-1}(u_m - u), \boldsymbol{K}(u, \cdot) \mathcal{C}_0^{-1}(u_m - u) \rangle_{\mathcal{H}_{\boldsymbol{K}}}$$

$$= \langle \phi(\cdot), \phi(\cdot) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \langle \boldsymbol{K}(u, u) \mathcal{C}_0^{-1}(u_m - u), \mathcal{C}_0^{-1}(u_m - u) \rangle_{\mathcal{Y}}$$

$$\leq \langle \phi(\cdot), \phi(\cdot) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \langle \boldsymbol{K}(u, u) \mathcal{C}_0^{-1}(u_m - u), \mathcal{C}_0^{-1}(u_m - u) \rangle_{\mathcal{Y}}$$

$$\leq C \langle \phi(\cdot), \phi(\cdot) \rangle_{\mathcal{H}_{\boldsymbol{K}}} \|\mathcal{C}_0^{-\frac{1-s}{2}}(u_m - u)\|_{\mathcal{H}}^2,$$

431  we find that $\lim_{m\to\infty}\langle\phi(\cdot),\boldsymbol{K}(u,\cdot)\mathcal{C}_0^{-1}u_m\rangle_{\mathcal{H}_{\boldsymbol{K}}} = \langle\phi(\cdot),\boldsymbol{K}(u,\cdot)\mathcal{C}_0^{-1}u\rangle_{\mathcal{H}_{\boldsymbol{K}}}$. Hence, let
432  $m \to \infty$ in (38), we have

433
434  (39) $$\langle\mathcal{C}_0^{-1/2}u,\mathcal{C}_0^{-1/2}\phi(u)\rangle_{\mathcal{Y}} = \langle\phi(\cdot),\boldsymbol{K}(u,\cdot)\mathcal{C}_0^{-1}u\rangle_{\mathcal{H}_{\boldsymbol{K}}}.$$

435  Plugging (39) and (37) into (36), we obtain

436  (40) $$\begin{aligned}\langle DV(u),\phi(u)\rangle_{\mathcal{Y}} &=\langle\phi(\cdot),\boldsymbol{K}(u,\cdot)D\Phi(u) + \boldsymbol{K}(u,\cdot)\mathcal{C}_0^{-1}u\rangle_{\mathcal{H}_{\boldsymbol{K}}}\\ 437 &=\langle\phi(\cdot),\boldsymbol{K}(u,\cdot)DV(u)\rangle_{\mathcal{H}_{\boldsymbol{K}}}.\end{aligned}$$

438      Next, let us calculate the second term on the right-hand side of (17). A simple
439  calculation yields

440  (41) $$\sum_{k=1}^{\infty} D_k\langle\phi(u),e_k\rangle_{\mathcal{Y}} = \sum_{k=1}^{\infty} D_k\langle\phi(\cdot),\boldsymbol{K}(u,\cdot)e_k\rangle_{\mathcal{H}_{\boldsymbol{K}}}.$$

442  Since

443  (42) $$\begin{aligned}D_k\langle\phi(\cdot),\boldsymbol{K}(u,\cdot)e_k\rangle_{\mathcal{H}_{\boldsymbol{K}}} &= \lim_{\epsilon\to 0}\frac{1}{\epsilon}\langle\phi(\cdot),\boldsymbol{K}(u+\epsilon\varphi_k,\cdot)e_k - \boldsymbol{K}(u,\cdot)e_k\rangle_{\mathcal{H}_{\boldsymbol{K}}}\\ 444 &= \langle\phi(\cdot),D_k\boldsymbol{K}(u,\cdot)e_k\rangle_{\mathcal{H}_{\boldsymbol{K}}},\end{aligned}$$

445  we have

446  (43) $$\sum_{k=1}^{\infty} D_k\langle\phi(u),e_k\rangle_{\mathcal{Y}} = \left\langle\phi(\cdot),\sum_{k=1}^{\infty} D_k\boldsymbol{K}(u,\cdot)e_k\right\rangle_{\mathcal{H}_{\boldsymbol{K}}}.$$

448  Combining (40) and (43) with (17), we obtain

449  (44) $$\mathcal{S}\phi(u) = \left\langle\phi(\cdot), -\boldsymbol{K}(u,\cdot)DV(u) + \sum_{k=1}^{\infty} D_k\boldsymbol{K}(u,\cdot)e_k\right\rangle_{\mathcal{H}_{\boldsymbol{K}}}.$$

451  Thus, the optimization problem (16) possesses a solution $\phi_{\boldsymbol{K}}^*(\cdot)$ satisfying

452  (45) $$\phi_{\boldsymbol{K}}^*(\cdot) \propto \mathbb{E}_{u\sim\mu}\Big[ -\boldsymbol{K}(u,\cdot)DV(u) + \sum_{k=1}^{\infty} D_k\boldsymbol{K}(u,\cdot)e_k\Big].$$

454  Based on condition (33), we know that $D\phi_{\boldsymbol{K}}^*(u)$ belongs to $\mathcal{L}_1(\mathcal{X},\mathcal{Y})$ for each $u \in \mathcal{X}$,
455  which completes the proof.                                                          □

456      REMARK 11. *The optimal $\phi_{\boldsymbol{K}}^*$ is given in (34) which is consistent with the finite-*
457  *dimensional case. Since the first and second terms on the right-hand side of (34) are*
458  *similar, we may just focus on the second term which is usually named as the repul-*
459  *sive force term. For each $u,v \in \mathcal{X}$, consider $\boldsymbol{K}(u,v) := K(u,v)Id$ with $K(u,v) :=$*
460  *$\exp\left(-\frac{1}{h}\|u-v\|_{\mathcal{X}}^2\right)$. Then, we have*

461  (46) $$\begin{aligned}\sum_{k=1}^{\infty} D_{u_k}\boldsymbol{K}(u,v)e_k &=\sum_{k=1}^{\infty}\langle D_u K(u,v)e_k,\varphi_k\rangle_{\mathcal{X}}\\ 462 &=\sum_{k=1}^{\infty} -\frac{2}{h}\langle u-v,\varphi_k\rangle_{\mathcal{X}}K(u,v)e_k.\end{aligned}$$

463 *Projecting (46) on one particular coordinate $e_\ell$ with $\ell \in \mathbb{N}$, we obtain*

464 (47)
$$\left(\sum_{k=1}^{\infty} D_{u_k} \boldsymbol{K}(u,v)e_k\right)_\ell = \left\langle \sum_{k=1}^{\infty} -\frac{2}{h}\langle u-v, \varphi_k\rangle_{\mathcal{X}} K(u,v)e_k, e_\ell \right\rangle_{\mathcal{Y}}$$

465
$$= -\frac{2}{h}\langle u-v, \varphi_\ell\rangle_{\mathcal{X}} K(u,v),$$

466 *which is similar to the $\ell$th coordinate of $\nabla_{u^N} K(u^N, v^N)$ appearing in (9). Addition-*
467 *ally, we mention that the assumption (33) given in Theorem 10 can be verified for*
468 *many useful kernels. Detailed illustrations are provided in the supplementary material.*

469 By Theorem 10, we can construct a series of transformations as follows:

470 (48)
$$T_\ell(u) = u + \epsilon_\ell \mathbb{E}_{u'\sim\mu_\ell}\left[ -\boldsymbol{K}(u',u)DV(u') + \sum_{k=1}^{\infty} D_{(u')_k}\boldsymbol{K}(u',u)e_k \right]$$

472 with $\ell = 1, 2, \ldots$. In practice, we draw a set of particles $\{u_i^0\}_{i=1}^m$ from some initial
473 measure, and then iteratively update the particles with an empirical version of the
474 above transformation in which the expectation under $\mu_\ell$ is approximated by the em-
475 pirical mean of particles $\{u_i^\ell\}_{i=1}^m$ at the $\ell$-th iteration. The iSVGD is summarized in
476 Algorithm 2.

---

**Algorithm 2** Infinite-dimensional Stein variational gradient descent (iSVGD)

---

**Input:** A target probability measure $\mu^{\boldsymbol{d}}$ that is absolutely continuous w.r.t the
Gaussian measure $\mu_0 = \mathcal{N}(0, \mathcal{C}_0)$ with $\frac{d\mu^{\boldsymbol{d}}}{d\mu_0}(u) \propto \exp(-\Phi(u))$ and a set of particles
$\{u_i^0\}_{i=1}^m$.
**Output:** A set of particles $\{u_i\}_{i=1}^m$ that approximates the target probability mea-
sure.
**for** iteration $\ell$ **do**

$$u_i^{\ell+1} \longleftarrow u_i^\ell + \epsilon_\ell \phi^*(u_i^\ell),$$

where

$$\phi^*(u) = \frac{1}{m}\sum_{j=1}^m \left[ \boldsymbol{K}(u_j^\ell, u)(-D\Phi(u_j^\ell) - \mathcal{C}_0^{-1}u_j^\ell) + \sum_{k=1}^{\infty} D_{(u_j^\ell)_k}\boldsymbol{K}(u_j^\ell, u)e_k \right].$$

**end for**

---

477 **3.3. iSVGD with precondition information.** In the supplementary material,
478 the numerical experiments indicate that the SVGD without preconditioning operators
479 converges slowly for some inverse problems of PDEs. By the finite-dimensional SVGD
480 [58], it may accelerate the convergence and give reliable estimates efficiently by intro-
481 ducing preconditioning operators. For constructing the iSVGD with preconditioning
482 operators, let us begin with a theorem concerning the change of variables.

483 THEOREM 12. *Let $\mathcal{X}$ and $\mathcal{Y}$ be two separable Hilbert spaces, and let $\mathcal{F}_0$ be a RKHS*
484 *with a nonnegative $\mathcal{L}(\mathcal{Y})$-valued kernel $\boldsymbol{K}_0 : \mathcal{X} \times \mathcal{X} \to \mathcal{L}(\mathcal{Y})$. Let $\tilde{\mathcal{X}}$ and $\tilde{\mathcal{Y}}$ be two*
485 *separable Hilbert spaces, and $\mathcal{F}$ be the set of operators from $\tilde{\mathcal{X}}$ to $\tilde{\mathcal{Y}}$ given by*

486
487 (49)
$$\phi(u) = \boldsymbol{M}(u)\phi_0(t(u)) \quad \forall\, \phi_0 \in \mathcal{F}_0,$$

where $\boldsymbol{M} : \tilde{\mathcal{X}} \to \mathcal{L}(\mathcal{Y}, \tilde{\mathcal{Y}})$ is a fixed operator and is assumed to be an invertible operator for all $u \in \tilde{\mathcal{X}}$, and $t : \tilde{\mathcal{X}} \to \mathcal{X}$ is a fixed Fréchet differentiable one-to-one mapping. For all $\phi, \phi' \in \mathcal{F}$, we can identify a unique $\phi_0, \phi_0' \in \mathcal{F}_0$ such that $\phi(u) = \boldsymbol{M}(u)\phi_0(t(u))$ and $\phi'(u) = \boldsymbol{M}(u)\phi_0'(t(u))$. Define the inner product on $\mathcal{F}$ via $\langle \phi, \phi' \rangle_\mathcal{F} = \langle \phi_0, \phi_0' \rangle_{\mathcal{F}_0}$, and then $\mathcal{F}$ is also a vector-valued RKHS, whose operator-valued kernel is

$$(50) \qquad \boldsymbol{K}(u, u') = \boldsymbol{M}(u')\boldsymbol{K}_0(t(u), t(u'))\boldsymbol{M}(u)^*,$$

where $\boldsymbol{M}(u)^*$ denotes the Hilbert space adjoint.

*Proof.* Let $\{(u_i, g_i)_{i=1,\dots,N}\} \subset \tilde{\mathcal{X}} \times \tilde{\mathcal{Y}}$, and we have

$$(51) \qquad \begin{aligned} \langle \boldsymbol{K}(u_i, u_j)g_i, g_j \rangle_{\tilde{\mathcal{Y}}} &= \langle \boldsymbol{M}(u_j)\boldsymbol{K}_0(t(u_i), t(u_j))\boldsymbol{M}(u_i)^* g_i, g_j \rangle_{\tilde{\mathcal{Y}}} \\ &= \langle \boldsymbol{K}_0(t(u_i), t(u_j))\boldsymbol{M}(u_i)^* g_i, \boldsymbol{M}(u_j)^* g_j \rangle_{\mathcal{Y}}. \end{aligned}$$

Then, the nonnegativity of $\boldsymbol{K}(\cdot, \cdot)$ follows from the nonnegative property of $\boldsymbol{K}_0(\cdot, \cdot)$. To prove the theorem, it suffices to verify the two conditions shown in Definition 4. For every $u, v \in \tilde{\mathcal{X}}$ and $g \in \tilde{\mathcal{Y}}$, we consider the operator $f(v) = \boldsymbol{K}(u, v)g = \boldsymbol{M}(v)\boldsymbol{K}_0(t(u), t(v))\boldsymbol{M}(u)^* g$. Because of $\boldsymbol{M}(u)^* g \in \mathcal{Y}$, we easily obtain

$$\boldsymbol{K}_0(t(u), t(v))\boldsymbol{M}(u)^* g \in \mathcal{F}_0.$$

According to (49), we conclude that $f(\cdot) \in \mathcal{F}$.

Next, let us verify the reproducing property of $\boldsymbol{K}(\cdot, \cdot)$. For every $u \in \tilde{\mathcal{X}}, g \in \tilde{\mathcal{Y}}$, and $\phi \in \mathcal{F}$, we have

$$\begin{aligned} \langle \phi(u), g \rangle_{\tilde{\mathcal{Y}}} &= \langle \boldsymbol{M}(u)\phi_0(t(u)), g \rangle_{\tilde{\mathcal{Y}}} = \langle \phi_0(t(u)), \boldsymbol{M}(u)^* g \rangle_{\mathcal{Y}} \\ &= \langle \phi_0(\cdot), \boldsymbol{K}_0(t(u), \cdot)\boldsymbol{M}(u)^* g \rangle_{\mathcal{F}_0} \\ &= \langle \boldsymbol{M}(\cdot)\phi_0(t(\cdot)), \boldsymbol{M}(\cdot)\boldsymbol{K}_0(t(u), t(\cdot))\boldsymbol{M}(u)^* g \rangle_{\mathcal{F}} \\ &= \langle \phi(\cdot), \boldsymbol{K}(u, \cdot)g \rangle_{\mathcal{F}}, \end{aligned}$$

where the fourth equality follows from

$$\langle \phi, \phi' \rangle_\mathcal{F} = \langle \phi_0, \phi_0' \rangle_{\mathcal{F}_0}$$

with $\phi_0'(\cdot) = \boldsymbol{K}_0(t(u), \cdot)\boldsymbol{M}(u)^* g$. $\qquad \square$

Now we present a key result, which characterizes the change of kernels when applying invertible transformations on the iSVGD trajectory.

THEOREM 13. *Let $\mathcal{H}$, $\tilde{\mathcal{H}}$, $\mathcal{X}$, $\tilde{\mathcal{X}}$, $\mathcal{Y}$, and $\tilde{\mathcal{Y}}$ be separable Hilbert spaces satisfying $\mathcal{X} \subset \mathcal{Y}$, $\tilde{\mathcal{X}} \subset \tilde{\mathcal{Y}}$, $\mathcal{X} \subset \tilde{\mathcal{Y}}$, $\tilde{\mathcal{X}} \subset \mathcal{Y}$. Assume that Assumption 7 holds for the triples $(\mathcal{X}, \mathcal{Y}, \mathcal{H})$ and $(\tilde{\mathcal{X}}, \tilde{\mathcal{Y}}, \tilde{\mathcal{H}})$ with two fixed parameters $s \in [0, 1]$, respectively. Let $T \in \mathcal{L}(\mathcal{Y}, \tilde{\mathcal{Y}})$ and assume that $T$ is a bounded operator when restricted to be an operator from $\mathcal{X}$ to $\tilde{\mathcal{X}}$. Let $\mu$, $\mu^{\boldsymbol{d}}$ be two probability measures and $\tilde{\mu}$, $\tilde{\mu}^{\boldsymbol{d}}$ be the measures of $\tilde{u} = Tu$ when $u$ is drawn from $\mu$, $\mu^{\boldsymbol{d}}$, respectively. Introduce two Stein operators $\mathcal{S}$ and $\tilde{\mathcal{S}}$ as follows:*

$$\mathcal{S}\phi(u) = \langle -DV(u), \phi(u) \rangle_{\mathcal{Y}} + \sum_{k=1}^{\infty} D_k \langle \phi(u), e_k \rangle_{\mathcal{Y}}, \quad \forall u \in \mathcal{X},$$

$$\tilde{\mathcal{S}}\tilde{\phi}(\tilde{u}) = \langle -D_{\tilde{u}} V(T^{-1}\tilde{u}), \tilde{\phi}(\tilde{u}) \rangle_{\tilde{\mathcal{Y}}} + \sum_{k=1}^{\infty} D_{(\tilde{u})_k} \langle \tilde{\phi}(\tilde{u}), \tilde{e}_k \rangle_{\tilde{\mathcal{Y}}}, \quad \forall \tilde{u} \in \tilde{\mathcal{X}},$$

where $\{e_k\}_{k=1}^{\infty}$ and $\{\tilde{e}_k\}_{k=1}^{\infty}$ are orthonormal bases in $\mathcal{Y}$ and $\tilde{\mathcal{Y}}$, respectively. Then, we have

$$(52) \qquad \mathbb{E}_{u\sim\mu}[\mathcal{S}\phi(u)] = \mathbb{E}_{u\sim\tilde{\mu}}[\tilde{\mathcal{S}}\tilde{\phi}(u)] \quad \text{with } \phi(u) := T^{-1}\tilde{\phi}(Tu).$$

Therefore, in the asymptotics of infinitesimal step size ($\epsilon \to 0^+$), it is equivalent to running iSVGD with kernel $\boldsymbol{K}_0$ on $\tilde{\mu}$ and running iSVGD on $\mu$ with the kernel $\boldsymbol{K}(u,u') = T^{-1}\boldsymbol{K}_0(Tu,Tu')(T^{-1})^*$, in the sense that the trajectory of these two SVGD can be mapped to each other by the map $T$ (and its inverse).

*Proof.* Let us introduce a mapping defined by $u' = f(u) = u + \epsilon\phi(u)$. Denote $f_{\#}\mu$ as the probability measure $\mu \circ f^{-1}$. Let $\tilde{u}' \sim T_{\#}(f_{\#}\tilde{\mu})$ which is obtained by

$$
\begin{aligned}
\tilde{u}' = Tu' = T(u + \epsilon\phi(u)) &= T(T^{-1}\tilde{u} + \epsilon\phi(T^{-1}\tilde{u})) \\
(53) \qquad &= \tilde{u} + \epsilon T\phi(T^{-1}\tilde{u}) \\
&= \tilde{u} + \epsilon\tilde{\phi}(\tilde{u}),
\end{aligned}
$$

where we use the definition $\phi(u) = T^{-1}\tilde{\phi}(Tu)$ in (52). According to [39, Theorem 3.1 ] and [58, Theorem 3], we have $\mathbb{E}_{u^N\sim P^N_{\#}\mu}[\mathcal{S}^N\phi^N(u^N)] = \mathbb{E}_{u^N\sim P^N_{\#}\tilde{\mu}}[\tilde{\mathcal{S}}^N\tilde{\phi}^N(u^N)]$, where

$$\mathcal{S}^N\phi^N(u^N) = -\langle DV(u^N), \phi^N(u^N)\rangle_{\mathcal{Y}} + \sum_{k=1}^{N} D_k\langle\phi^N(u^N), e_k\rangle_{\mathcal{Y}},$$

$$\tilde{\mathcal{S}}^N\tilde{\phi}^N(\tilde{u}^N) = -\langle D_{\tilde{u}^N}V(T^{-1}\tilde{u}^N), \tilde{\phi}^N(\tilde{u}^N)\rangle_{\tilde{\mathcal{Y}}} + \sum_{k=1}^{N} D_{(\tilde{u}^N)_k}\langle\tilde{\phi}^N(\tilde{u}^N), \tilde{e}_k\rangle_{\tilde{\mathcal{Y}}}.$$

It is clear to note that there is no Jacobian matrix given by the transformation in $D_{\tilde{u}^N}V(T^{-1}\tilde{u}^N)$ since the Jacobian matrix does not depend on $\tilde{u}^N$ for linear mappings, i.e., the derivative is zero. Following the proof for Theorem 9, we take $N \to \infty$ and obtain $\mathbb{E}_{u\sim\mu}[\mathcal{S}\phi(u)] = \mathbb{E}_{u\sim\tilde{\mu}}[\tilde{\mathcal{S}}\tilde{\phi}(u)]$. From Theorem 12, when $\tilde{\phi}$ is in $\tilde{\mathcal{F}}$ with kernel $\boldsymbol{K}_0(u,u')$, $\phi$ is in $\mathcal{F}$ with kernel $\boldsymbol{K}(u,u')$. Therefore, maximizing $\mathbb{E}_{u\sim\mu}[\mathcal{S}\phi(u)]$ in $\mathcal{F}$ is equivalent to $\mathbb{E}_{u\sim\tilde{\mu}}[\tilde{\mathcal{S}}\tilde{\phi}(u)]$ in $\tilde{\mathcal{F}}$. This suggests that the trajectory of iSVGD on $\tilde{\mu}^{\boldsymbol{d}}$ with $\boldsymbol{K}_0$ and that on $\mu^{\boldsymbol{d}}$ with $\boldsymbol{K}$ are equivalent, which completes the proof. $\square$

REMARK 14. *Similar to the matrix-valued case [58], Theorem 13 suggests a conceptual procedure for constructing proper operator kernels to incorporate desirable preconditioning information. Different from the finite-dimensional case, the map $T$ is only allowed to be linear at this stage. For a nonlinear map, there is a Jacobian matrix in $\tilde{\mathcal{S}}^N\tilde{\phi}^N(\tilde{u}^N)$. It is difficult to analyze the limiting behavior of the Jacobian matrix related term. Practically, linear maps seem to be enough since even in the finite-dimensional case nonlinear maps will yield an unnatural algorithm [58].*

In the last part of this subsection, we provide some examples of preconditioning operators that are frequently used in statistical inverse problems.

**3.3.1. Fixed preconditioning operator.** In Section 5 of [16], the Langevin equation was considered by using $\mathcal{C}_0$ as a preconditioner, and an analysis was carried out for the pCN algorithm. For the Newton based iterative method, we usually take the inverse of the second-order derivative of the objective functional as the preconditioning operator [41]. Here, we consider a linear operator $T$ that has similar properties as $\mathcal{C}_0^{-\frac{1-s}{2}}$. Specifically, we require

$$(54) \qquad T \in \mathcal{L}(\mathcal{H}^{1-s}, \mathcal{H}) \cap \mathcal{L}(\mathcal{H}^{-1-s}, \mathcal{H}^{-2}).$$

Then, we specify the Hilbert space appearing in Theorem 12 as $\mathcal{X} = \mathcal{H}^{1-s}$, $\mathcal{Y} = \mathcal{H}^{-1-s}$, $\tilde{\mathcal{X}} = \mathcal{H}$, $\tilde{\mathcal{Y}} = \mathcal{H}^{-2}$ with $s \in [0,1]$. For the kernel $\boldsymbol{K}_0(\cdot, \cdot) : \tilde{\mathcal{X}} \times \tilde{\mathcal{X}} \to \tilde{\mathcal{Y}}$, we assume that

$$
\sup_{\tilde{u} \in \mathcal{H}} \|\boldsymbol{K}_0(\tilde{u}, \tilde{u})\|_{\mathcal{L}(\mathcal{H}^{-2})} < \infty. \tag{55}
$$

It follows from Theorem 13 that we may use a kernel of the form

$$
\boldsymbol{K}(u, u') := T^{-1} \boldsymbol{K}_0(Tu, Tu')(T^{-1})^*, \tag{56}
$$

where $u, u' \in \mathcal{H}^{1-s}$. Obviously, the kernel $\boldsymbol{K}$ given above satisfies

$$
\sup_{u \in \mathcal{H}^{1-s}} \|T^{-1} \boldsymbol{K}_0(Tu, Tu)(T^{-1})^*\|_{\mathcal{L}(\mathcal{H}^{-1-s})} < \infty. \tag{57}
$$

As an example, we may take $\boldsymbol{K}_0$ to be the scalar-valued Gaussian RBF kernel composed with operator $\mathcal{C}_0^s$:

$$
\boldsymbol{K}_0(u, u') := \exp\left(-\frac{1}{h}\|u - u'\|_{\mathcal{H}}^2\right)\mathcal{C}_0^s, \tag{58}
$$

which yields

$$
\boldsymbol{K}(u, u') = \exp\left(-\frac{1}{h}\|T(u - u')\|_{\mathcal{H}}^2\right)T^{-1}\mathcal{C}_0^s(T^{-1})^*, \tag{59}
$$

where $h$ is a bandwidth parameter. Define $\boldsymbol{K}_0^T(u, u') := \boldsymbol{K}_0(Tu, Tu')$. Let $\mathcal{P} := T^{-1}\mathcal{C}_0^s(T^{-1})^*$. By simple calculations, we find that the iSVGD update direction of the kernel in (56) is

$$
\phi_{\boldsymbol{K}}^*(\cdot) = \mathcal{P}\mathbb{E}_{u \sim \mu}\left[\boldsymbol{K}_0^T(u, \cdot)(-D\Phi(u) - \mathcal{C}_0^{-1}u) + \sum_{k=1}^{\infty} D_k \boldsymbol{K}_0^T(u, \cdot)e_k\right] = \mathcal{P}\phi_{\boldsymbol{K}_0^T}^*, \tag{60}
$$

which is a linear transform of the iSVGD update direction of the kernel $\boldsymbol{K}_0^T$ with the operator $T^{-1}\mathcal{C}_0^s(T^{-1})^*$.

**3.3.2. The $\mathcal{C}_0$ operator.** Choosing $T := \mathcal{C}_0^{-\frac{1-s}{2}}$, we can see that the condition (54) holds. Given the Kernel $\boldsymbol{K}_0$ in (58), the kernel $\boldsymbol{K}$ defined in (59) can be written as

$$
\boldsymbol{K}(u, u') = \exp\left(-\frac{1}{h}\|\mathcal{C}_0^{-\frac{1-s}{2}}(u - u')\|_{\mathcal{H}}^2\right)\mathcal{C}_0.
$$

The operator $\mathcal{P}$ used in (60) is just $\mathcal{C}_0$. If there is only one particle, the iSVGD update direction is then reduced to $\phi_{\boldsymbol{K}}^*(\cdot) = \mathcal{C}_0(D\Phi(u) + \mathcal{C}_0^{-1}u)$.

**3.3.3. The Hessian operator.** For statistical inverse problems, the forward operator $\mathcal{G}$ is usually nonlinear, e.g., the inverse medium scattering problem [26, 27]. Around each particle $u_i$ with $i = 1, 2, \ldots, m$, the forward map can be approximated by the linearized map

$$
\mathcal{G}(u) \approx \mathcal{G}(u_i) + D\mathcal{G}(u_i)(u - u_i). \tag{61}
$$

Assume that the potential function $\Phi$ takes the form $\Phi(u) = \frac{1}{2}\|\Sigma^{-1/2}(\mathcal{G}(u) - d)\|_{\ell^2}^2$, where $\Sigma$ is a positive definite matrix. Using the approximate formula (61), we have

$$
V(u) \approx \tilde{V}(u) := \frac{1}{2}\|\Sigma^{-1/2}(D\mathcal{G}(u_i)u - D\mathcal{G}(u_i)u_i + \mathcal{G}(u_i) - d)\|_{\ell^2}^2 + \frac{1}{2}\|\mathcal{C}_0^{-1/2}u\|_{\mathcal{H}}^2.
$$

It follows from a simple calculation that $D^2\tilde{V}(u_i) = D\mathcal{G}(u_i)^*\Sigma^{-1}D\mathcal{G}(u_i) + \mathcal{C}_0^{-1}$. For the Newton-type iterative method, we can take the linear transformation $T = \mathcal{C}_0^{s/2}(\frac{1}{m}\sum_{i=1}^m (D\mathcal{G}(u_i)^*\Sigma^{-1}D\mathcal{G}(u_i) + \mathcal{C}_0^{-1}))^{1/2}$. If $\mathcal{G}$ is a linear operator (e.g., the examples in [25]), it is easy to verify the condition (54). For nonlinear problems, it is necessary to employ the regularity properties of the direct problems, which is beyond the scope of this work. Hence we will not verify this condition in this paper and leave it as a future work. With this choice of $T$, the kernel (59) and the iSVGD update direction (60) can be easily obtained. If there is only one particle, the iSVGD update direction is degenerated to the usual Newton update direction when evaluating MAP estimate.

**3.3.4. Mixture preconditioning.** Using a fixed preconditioning operator, we can not specify different preconditioning operators for different particles. Inspired by the mixture precondition [58], we propose an approach to achieve point-wise preconditioning. The idea is to use a weighted combination of several linear preconditioning operators. This involves leveraging a set of anchor points $\{v_\ell\}_{\ell=1}^m$, each of which is associated with a preconditioning operator $T_\ell$ (e.g., $T_\ell = \mathcal{C}_0^{s/2}(D\mathcal{G}(v_\ell)^*\Sigma^{-1}D\mathcal{G}(v_\ell) + \mathcal{C}_0^{-1})^{1/2}$). In practice, the anchor points $\{v_\ell\}_{\ell=1}^m$ can be set to be the same as the particles $\{u_i\}_{i=1}^m$. We then construct a kernel by $\boldsymbol{K}(u, u') = \sum_{\ell=1}^m \boldsymbol{K}_\ell(u, u')w_\ell(u)w_\ell(u')$, where

$$(62) \qquad \boldsymbol{K}_\ell(u, u') := T_\ell^{-1}\boldsymbol{K}_0(T_\ell u, T_\ell u')(T_\ell^{-1})^*,$$

and $w_\ell(u)$ is a positive scalar-valued function that determines the contribution of kernel $\boldsymbol{K}_\ell$ at point $u$. Here $w_\ell(u)$ should be viewed as a mixture probability, and hence should satisfy $\sum_{\ell=1}^m w_\ell(u) = 1$ for all $u$. In our empirical studies, we take

$$(63) \qquad w_\ell(u) = \frac{\exp\left(-\frac{1}{2}\|T_\ell(u - v_\ell)\|_{\mathcal{H}}^2\right)}{\sum_{\ell'=1}^m \exp\left(-\frac{1}{2}\|T_{\ell'}(u - v_{\ell'})\|_{\mathcal{H}}^2\right)}.$$

In this way, each point $u$ is mostly influenced by the anchor point closest to it, which allows to apply different preconditioning for different points. In addition, the iSVGD update direction has the form

$$(64) \qquad \phi_{\boldsymbol{K}}^*(\cdot) = \sum_{\ell=1}^m w_\ell(\cdot)\mathbb{E}_{u\sim\mu}\Big[-w_\ell(u)\boldsymbol{K}_\ell(u, \cdot)(D\Phi(u) + \mathcal{C}_0^{-1}u)$$
$$+ \sum_{k=1}^\infty D_k(w_\ell(u)\boldsymbol{K}_\ell(u, \cdot)e_k)\Big],$$

which is a weighted sum of a number of iSVGD update directions with linear preconditioning operators. The implementation details of (64) are given in the supplementary material.

REMARK 15. *For the kernel defined above, the particles should belong to the Hilbert space $\mathcal{H}^{1-s}$. Based on the studies the finite-dimensional problems [58], we may let the parameter $s$ be equal to $0$. However, when the parameter $s = 0$, each particle $u_i$ belongs to $\mathcal{H}^1$ which is the Cameron–Martin space of the prior measure. By the classical Gaussian measure theory [15], we know that $\mathcal{H}^1$ has zero measure. This fact implies that all of the particles belong to a set with zero measure, which may lead to too concentrated particles and deviates from our purpose. Hence we should choose $s > 0$ to ensure the effectiveness of the SVGD sampling algorithm. These observations are illustrated by our numerical experiments in Section 4.*

**3.4. Some insights about iSVGD.** We have constructed the well-defined iSVGD algorithms with or without preconditioning operators, which is the first step to extend the finite-dimensional SVGD to the infinite-dimensional space. Some mathematical studies have been carried out for the finite-dimensional SVGD, e.g., gradient flow on probability space [38] and mean field limit theory related to the macroscopic behavior [42]. These results provide in-depth understandings of the SVGD algorithm and motivate many new algorithms [37]. In this subsection, we intend to provide a preliminary mathematical study on the iSVGD under a simpler setting.

We consider the kernel operator $\boldsymbol{K}(u,v) := K(\|u-v\|_{\mathcal{H}})\mathrm{Id}$ with $u,v \in \mathcal{H}$ and $K(\cdot)$ being a scalar function. Let $m$ be the sample number and $V(u)$ be defined in (15). Similar to the finite-dimensional case, the iterative procedure in Algorithm 2 can be viewed as a particle system:

$$\frac{d}{dt}u_i(t) = -(\tilde{D}\boldsymbol{K} * \mu_m(t))(u_i(t)) - (\boldsymbol{K} * DV\mu_m(t))(u_i(t)),$$

(65)
$$\mu_m(t) = \frac{1}{m}\sum_{j=1}^{m}\delta_{u_j(t)},$$

$$u_i(0) = u_i^0, \quad i = 1, 2, \ldots, m,$$

where $\{u_i^0\}_{i=1}^m$ are the initial particles, $\delta_{u_i(t)}$ denotes the Dirac measure concentrated on $u_i(t)$ with $i = 1, 2, \ldots, m$, "$*$" denotes the usual convolution operator, and $\tilde{D}\boldsymbol{K}(u-v) = \sum_{k=1}^{\infty}D_{u_k}\boldsymbol{K}(u-v)e_k$. For convenience, we write the two convolution terms in the following forms:

$$(\tilde{D}\boldsymbol{K} * \mu_m(t))(u_i(t)) = \frac{1}{m}\sum_{j=1}^{m}\tilde{D}\boldsymbol{K}(u_i(t) - u_j(t)),$$

$$(\boldsymbol{K} * DV\mu_m(t))(u_i(t)) = \frac{1}{m}\sum_{j=1}^{m}\boldsymbol{K}(u_i(t) - u_j(t))DV(u_j(t)).$$

Similarly, we consider the weak form equation about the measure-valued function:

(66)
$$\frac{d}{dt}\langle\mu(t),\varphi\rangle = \langle\mu(t), L(\mu(t))\varphi\rangle,$$

$$\mu(0) = \nu,$$

where $\nu$ is the probability measure employed to generate initial particles, $\varphi$ is the test function, and

(67)
$$L(\mu(t))\varphi = \langle\tilde{D}\boldsymbol{K} * \mu(t), D\varphi\rangle_{\mathcal{H}} + \langle\boldsymbol{K} * DV\mu(t), D\varphi\rangle_{\mathcal{H}}.$$

Let $W^{1,2}(\mathcal{H},\mu)$ be the usual Sobolev space defined for a Gaussian measure $\mu$ [47].

THEOREM 16. *Let $\mu_0$ and $\Phi$ be the prior measure and potential function defined in (1), respectively. Assume $K(\cdot) \in W^{1,2}(\mathcal{H},\mu_0)$ and $e^{-\Phi(\cdot;\boldsymbol{d})} \in L^2(\mathcal{H},\mu_0)$. Then, the posterior measure $\mu^{\boldsymbol{d}}$ defined in (1) is an invariant solution to Eq. (66), i.e., when $\nu := \mu^{\boldsymbol{d}}$, the solution $\mu(t)$ of (66) is equal to $\mu^{\boldsymbol{d}}$.*

The proof is given in the supplementary material. Clearly, this theorem holds in the finite-dimensional setting. We point out that the integration by parts may not hold for the infinite-dimensional case. In the finite-dimensional setting, the analysis of

the corresponding particle system (65) and Eq. (66) have been given recently in [42]. It is sophisticated to define meaningful solutions for the above interacting particle system (65) and the measure-valued function equation (66), which are beyond the scope of this study and are left for future work. One of the major difficulties for the infinite-dimensional case is that $\mathcal{C}_0^{-1}$ (the precision operator of the prior measure) is usually an unbounded operator [16]. Nearly all of the estimates presented in [42] for the finite-dimensional case cannot be adopted for the infinite-dimensional setting.

Numerical experiments indicate that the SVGD without preconditioning operators can hardly provide accurate estimates for some inverse problems. The SVGD with preconditioning operators can accelerate the convergence and give reliable estimates efficiently. In addition, the unboundedness issue induced by the precision operator $\mathcal{C}_0^{-1}$ may be overcome by introducing preconditioning operators. A detailed analysis of the iSVGD with preconditioning operators may be a good starting point for future theoretical studies.

At the end of this subsection, we mention a critical difference between finite- and infinite-dimensional theories. It follows from Theorem 2.7 in [42] and Theorem 1.1 in [57] that the empirical measure constructed by particles in finite-dimensional SVGD can approximate the continuous counterpart with accuracy $\epsilon$ when the number of particles are of order $O(\epsilon^d)$, where $d$ is the discrete dimension. Obviously, an infinite number of particles is needed if the dimension $d$ goes to infinity, which indicates that the infinite-dimensional theory may be meaningless.

The above statement explains that not every finite-dimensional setting can be meaningfully generalized to the infinite-dimensional space. The assumption on prior measure is important for the infinite-dimensional theory (the current assumption may be slightly relaxed, e.g., the Besov type measure). According to the general analysis for the convergence and concentration of empirical measures given in [34], we believe that the prior measures used here can be approximated by the empirical measures under the Wasserstein distance on infinite-dimensional Hilbert space. Specifically, the estimate of the convergence speed is not relevent to the dimension when considering some finite-dimensional spaces as the projected infinite-dimensional space. If a theorem similar as Theorem 2.7 in [42] for the system (65)–(66) can be proved, we are able to confirm that the particles obtained by iSVGD can approximate the posterior measure for certain accuracy with particle numbers independent of the discrete dimension. However, it is higly non-trivial to carry out an in-depth study of the system (65)–(66) and is beyond the scope of the current work. In Subsection 6.3 of the supplementary material, we give a numerical illustration to address this issue.

**4. Applications.** The proposed framework is valid for Bayesian inverse problems governed by any systems of PDEs. Due to the page limitation, we present one example of an inverse problem governed by the steady state Darcy flow equation. The second example concerns an inverse problem of the Helmholtz equation and is given in the supplementary material.

Consider the following PDE model:

$$(68) \qquad \begin{aligned} -\nabla \cdot (e^u \nabla w) &= f \quad \text{in } \Omega, \\ w &= 0 \quad \text{on } \partial\Omega, \end{aligned}$$

where $\Omega \subset \mathbb{R}^2$ is a bounded Lipschitz domain, $f(x)$ denotes the sources, and $e^{u(x)}$ describes the permeability of the porous medium. This model is used as a benchmark problem in many works, e.g., the preconditioned Crank–Nicolson (pCN) algorithm [13] and the sequential Monte Carlo method [2]. We will compare the performance

726 of the proposed iSVGD approach with the pCN [13] and the randomized maximum a
727 posterior (rMAP) method [59].

728 **4.1. Basic settings and finite-element discretization.** For numerical im-
729 plementations, it is essential to compute all of the related gradients and Hessian
730 operators before discretization (i.e., pushing the discretization to the last step). A
731 direct calculation yields the gradient and Hessian operators of the operator-valued
732 kernel, but the adjoint method [41] needs to be employed for the potential $\Phi$ involv-
733 ing PDEs. More discussions on finite- and infinite-dimensional approaches can be
734 found in the supplementary material, which might be helpful for readers who are not
735 familiar with infinite-dimensional approach. Let $\mathcal{F}$ be the solution operator that maps
736 the parameter $u$ to the solution of (68), and $\mathcal{M}$ be the measurement operator defined
737 as $\boldsymbol{d} = \mathcal{M}(w) = (\ell_{x_1}(w), \ell_{x_2}(w), \dots, \ell_{x_{N_d}}(w))^T$, where

738 (69)
$$\ell_{x_j}(w) = \int_\Omega \frac{1}{2\pi\delta^2} e^{-\frac{1}{2\delta^2}\|x-x_j\|^2} w(x) dx$$
739

740 with $\delta > 0$ being a sufficiently small number and $x_i \in \Omega$ for $i = 1, \dots, N_d$. The forward
741 map can be defined as $\mathcal{G} := \mathcal{M} \circ \mathcal{F}$, and the problem can be written in the abstract
742 form $\boldsymbol{d} = \mathcal{G}(u) + \boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathrm{Id})$. Then we have $\Phi(u) = \frac{1}{2\sigma^2}\|\mathcal{M}(w) - \boldsymbol{d}\|^2$. The
743 gradient $D\Phi(u)$ acting in any direction $\tilde{u}$ is given by

744 (70)
$$\langle D\Phi(u), \tilde{u} \rangle = \int_\Omega \tilde{u} e^u \nabla w \cdot \nabla p \, dx,$$
745

746 where the adjoint state $p$ satisfies *the adjoint equation*

747 (71)
$$-\nabla \cdot (e^u \nabla p) = -\frac{1}{\sigma^2} \sum_{j=1}^{N_d} \frac{1}{2\pi\delta^2} e^{-\frac{1}{2\delta^2}\|x-x_j\|^2} (\ell_{x_j}(w) - d_j) \quad \text{in } \Omega,$$

748
$$p = 0 \quad \text{on } \partial\Omega.$$

749 The Hessian acting in direction $\tilde{u}$ and $\hat{u}$ reads

750 (72)
$$\langle\langle D^2\Phi(u), \hat{u}\rangle, \tilde{u}\rangle = \int_\Omega \hat{u}\tilde{u} e^u \nabla w \cdot \nabla p \, dx + \int_\Omega \tilde{u} e^u \nabla w \cdot \nabla \hat{p} \, dx$$
$$+ \int_\Omega \tilde{u} e^u \nabla p \cdot \nabla \hat{w} \, dx,$$
751

752 where the state $\hat{w}$ satisfies *the incremental forward equation*

753 (73)
$$-\nabla \cdot (e^u \nabla \hat{w}) = \nabla \cdot (\hat{u} e^u \nabla w) \quad \text{in } \Omega,$$
$$\hat{w} = 0 \quad \text{on } \partial\Omega,$$
754

755 and the state $\hat{p}$ satisfies *the incremental adjoint equation*

756 (74)
$$-\nabla \cdot (e^u \nabla \hat{p}) = \nabla \cdot (\hat{u} e^u \nabla p) - \frac{1}{2\pi\delta^2\sigma^2} \sum_{j=1}^{N_d} \hat{w} e^{-\frac{1}{2\delta^2}\|x-x_j\|^2} \quad \text{in } \Omega,$$

757
$$\hat{p} = 0 \quad \text{on } \partial\Omega.$$

758 In experiments, we choose $\Omega$ to be a rectangular domain $\Omega = [0,1]^2 \subset \mathbb{R}^2$, set
759 $\mathcal{H} = L^2(\Omega)$, and consider the prior measure $\mu_0 = \mathcal{N}(u_0, \mathcal{C}_0)$ with the mean function $u_0$

and the covariance operator $\mathcal{C}_0 := A^{-2}$, where $A = \alpha(I - \Delta)$ $(\alpha > 0)$ with the domain of $A$ given by $D(A) := \left\{ u \in H^2(\Omega) \ : \ \frac{\partial u}{\partial \boldsymbol{n}} = 0 \text{ on } \partial\Omega \right\}$. Here, $H^2(\Omega)$ is the usual Sobolev space. Assume that the mean function $u_0$ resides in the Cameron–Martin space of $\mu_0$.

Based on (70) and (72), we can prove the following results, which satisfy Assumptions 5. The proof is given in the supplementary material.

THEOREM 17. *Let $H^{-1}(\Omega)$ be the usual Sobolev space with the regularity index $-1$. Assume $\mathcal{X} = \mathcal{H}^{1-s}$ with the parameter $s < 0.5$, and then we have*

$$0 \leq \Phi(u) \leq C(1 + \|f\|_{H^{-1}})^2 e^{2\|u\|_{\mathcal{X}}},$$

$$\|D\Phi(u)\|_{\mathcal{X}^*} \leq C(1 + \|f\|_{H^{-1}})^2 e^{4\|u\|_{\mathcal{X}}},$$

$$\|D^2\Phi(u)\|_{\mathcal{L}(\mathcal{X},\mathcal{X}^*)} \leq C(1 + \|f\|_{H^{-1}})^2 e^{6\|u\|_{\mathcal{X}}}.$$

In the following, we use the Gaussian kernel, i.e., $\boldsymbol{K}(u, u') = \exp\left(-\frac{1}{h}\|u - u'\|_{\mathcal{H}}^2\right)$, for the iSVGD without preconditioning operators. For numerical examples with preconditioning operators, we employed the kernel given in Subsection 3.3.4.

For finite-dimensional approximations, we consider a finite-dimensional subspace $V_h$ of $L^2(\Omega)$ originating from the finite element discretization with the continuous Lagrange basis functions $\{\phi_j\}_{j=1}^n$, which correspond to the nodal points $\{x_j\}_{j=1}^n$, such that $\phi_j(x_i) = \delta_{ij}$ for $i, j \in \{1, \ldots, n\}$. Instead of statistically inferring parameter functions $u \in L^2(\Omega)$, we consider the approximation $u_h = \sum_{j=1}^n u_j\phi_j \in V_h$. Under this finite-dimensional approximation, we can employ the numerical method provided in [4] to discretize the prior, and construct finite-dimensional approximations of the Gaussian approximation of the posterior measure. Based on our analysis in Subsection 3.3, we need to calculate the fractional powers of the operator $\mathcal{C}_0$. Here, we employ the matrix transfer technique (MTT) [6]. The main idea of MTT is to indirectly discretize a fractional Laplacian using a discretization of the standard Laplacian. As discussed in [4], the operator $M$ is taken as

$$(75) \qquad M = (M_{ij})_{i,j=1}^n \quad \text{and} \quad M_{ij} = \int_\Omega \phi_i(x)\phi_j(x)dx, \quad i, j \in \{1, \ldots, n\}.$$

The matrix $M^{1/2}$ is approximated by the diagonal matrix $\mathrm{diag}(M_{11}^{1/2}, \ldots, M_{nn}^{1/2})$.

Finally, we mention that the finite element discretization is implemented by employing the open software FEniCS (Version 2019.1.0) [40]. All programs were run on a personal computer with Intel(R) Core(TM) i7-7700 at 3.60 GHz (CPU), 32 GB (memory), and Ubuntu 18.04.2 LTS (OS).

**4.2. Numerical results.** In the experiments, the noise level is fixed to be 1% since the goal is to test algorithms rather than demonstrating the Bayesian modeling. We compare the iSVGD with the mixture preconditioning operator (iSVGDMPO) with the preconditioned Crank–Nicolson (pCN) sampling algorithm [16] and the randomized maximum a posteriori (rMAP) algorithm [59]. Since the rMAP sampling algorithm is not accurate for nonlinear problems, we choose $\alpha = 0.5$ in the prior probability measure. It should be mentioned that we choose the anchor points in the iSVGDMPO just to be the same as the particles and the anchor points will be updated during the iterations. The initial particles of the iSVGD are generated from a probability measure by using the method proposed in [4].

For the current settings, the gradient descent based method seems hardly to find appropriate solutions in reasonable iterative steps. Hence, the optimization method
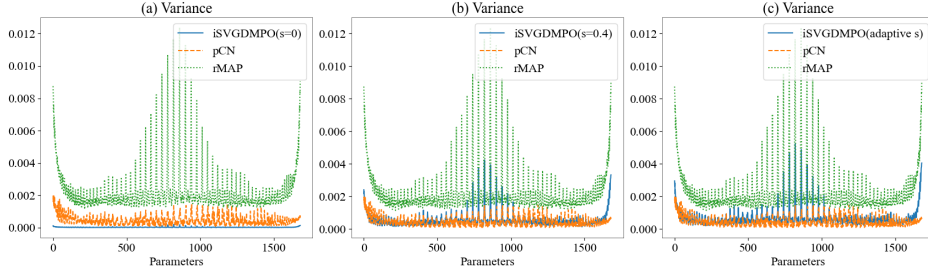
FIG. 1. *The comparison of the variances estimated by the pCN, rMAP, iSVGDMPO with different s. (a): $s = 0$; (b): $s = 0.4$; (c): adpatively chosen s.*

with preconditioning operators, e.g., the Newton-conjugate gradient method, is employed. The term $\mathbb{E}_{u' \sim \mu_\ell}[\boldsymbol{K}(u', u)DV(u')]$ in (48) is an averaged gradient descent component in the whole iterative term, which drives all of the particles to be concentrated. We anticipate that Algorithm 2 cannot work well due to the inefficiency of the gradient descent algorithm. Due to the page limitation, numerical results are given in the supplementary material, which show that Algorithm 2 does not perform well in some cases. This is one of the main motivations for us to study the iSVGD with preconditioning operators.

We compare the iSVGD with the mixture preconditioning operator (iSVGDMPO) with those obtained by the pCN and rMAP sampling algorithms. As illustrated in Remark 15, the parameter $s$ should not be zero. Intuitively, the particles should belong to a space with probability approximately equal to one under the prior measure $\mu_0$. By the Gaussian measure theory [15], we may take $s > 0.5$ since $\mu_0(\mathcal{H}^{1-s}) = 1$ for any $s > 0.5$. Since the posterior measure is usually concentrated on a small support set of the prior measure, the parameter $s$ should be slightly smaller than 0.5. Thus, we set $s = 0.3$ or 0.4 in our examples. Usually, the initial particles are scattered, and the variances of the initial particles are larger than the final particles obtained by the iSVGDMPO. We design the following adaptive empirical strategy for $s$:

$$(76) \qquad\qquad s = -0.5 \frac{\|\mathrm{var}\|_{\ell^2}}{\|\mathrm{var}_0\|_{\ell^2}} + 0.5,$$

where var is the current estimated variance, $\mathrm{var}_0$ is the estimated variance of the initial particles, and $\| \cdot \|_{\ell^2}$ is the usual $\ell^2$-norm. Obviously, for the initial particles, we have $s = 0$. The particles are forced to be concentrated. When the variance is reduced, the parameter $s$ approaches 0.5 to avoid that the particles are concentrated on a set with zero measure. Since the pCN is a dimension independent MCMC type sampling algorithm, we take the results obtained by the pCN as the baseline (accurate estimate). To make sure that the pCN algorithm yields an accurate estimate, we iterate $10^6$ steps and withdraw the first $10^5$ samples. Several different step-sizes are tried and the traces of some parameters are plotted, and then the most reliable one is picked as the baseline.

In Figure 1, we show the estimated variances obtained by the iSVGDMPO (blue solid line), rMAP (green dotted line), and the pCN (orange dashed line) sampling algorithms. The estimated variances of the iSVGDMPO are shown for $s = 0$ and $s = 0.4$ on the left and in the middle, respectively. On the right, we exhibit the estimated variances when the empirical adaptive strategy (76) is employed. As expected, the estimated variances are too small when $s = 0$, which indicates that the particles are
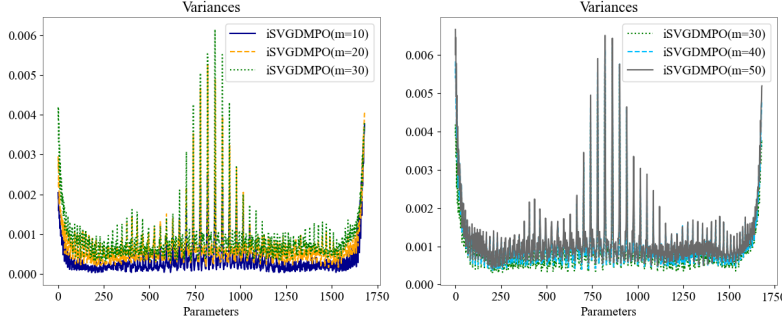
Fig. 2. *The comparison of the variances estimated by the iSVGDMPO with $s = 10, 20, 30, 40, 50$.*

842    concentrated on a small set. Choosing $s = 0.4$ or using the empirical strategy, we
843    obtain similar estimates, which is more similar to the baseline obtained by the pCN
844    compared with the estimates obtained by the rMAP.
845        One important question arises: how does $s$ influence the convergence of the
846    iSVGDMPO? The detailed numerical comparisons are given in the supplementary ma-
847    terial. Here we state the conclusions: The convergence speeds are similar for $s = 0.4$
848    and the adaptively chosen $s$. When specifying $s = 0.5$, the variances will gradually
849    approach the background truth, but the convergence speed seems much slower than
850    $s = 0.4$ or the adaptively chosen $s$. In the following numerical experiments, we use
851    the empirical adaptive strategy to specify the parameter $s$.
852        In addition, we provide three videos to exhibit the dynamic changing procedure of
853    the estimated variances in the supplementary material. The update perturbation with
854    and without repulsive force term are exhibited. These videos can further illustrate
855    our theoretical findings. We can see that the repulsive force terms indeed prevent the
856    particles from being over concentrated.
857        Apart from the parameter $s$, how many samples should be taken to guarantee
858    a stable statistical quantity estimate is important for using the iSVGDMPO. When
859    the particle number is too small, we cannot obtain reliable estimates. However, the
860    computational complexity increases when the particle number increases. In Figure 2,
861    we show the estimated variances when particle number equals to 10, 20, 30, 40, and
862    50. Denote by $m$ the number of samples. On the left in Figure 2, we show the results
863    obtained when $m = 10, 20, 30$. Obviously, when $m = 10$, the estimated variances are
864    significantly smaller than those obtained when $m = 20, 30$. On the right in Figure
865    2, we find that the estimated variances are similar when $m = 30, 40, 50$. Hence, it is
866    enough for our numerical examples to take $m = 20$ or 30, which attains a balance
867    between efficiency and accuracy. So far, we have only compared the variances with
868    different parameters in the iSVGDMPO. In the following, qualitative and quantitative
869    comparisons of other statistical quantities are provided to illustrate the effectiveness
870    of the iSVGDMPO.
871        Now, we specify the sampling number $m = 30$ and set the parameter $s$ by the
872    proposed empirical strategy (76). In Figure 3, we show the background truth and the
873    estimated mean and variance functions obtained by the pCN, rMAP, and iSVGDMPO,
874    respectively. The iterative number of the iSVGDMPO is set to be 30. From the first
875    line, we observe that the mean functions obtained by the rMAP and iSVGDMPO are
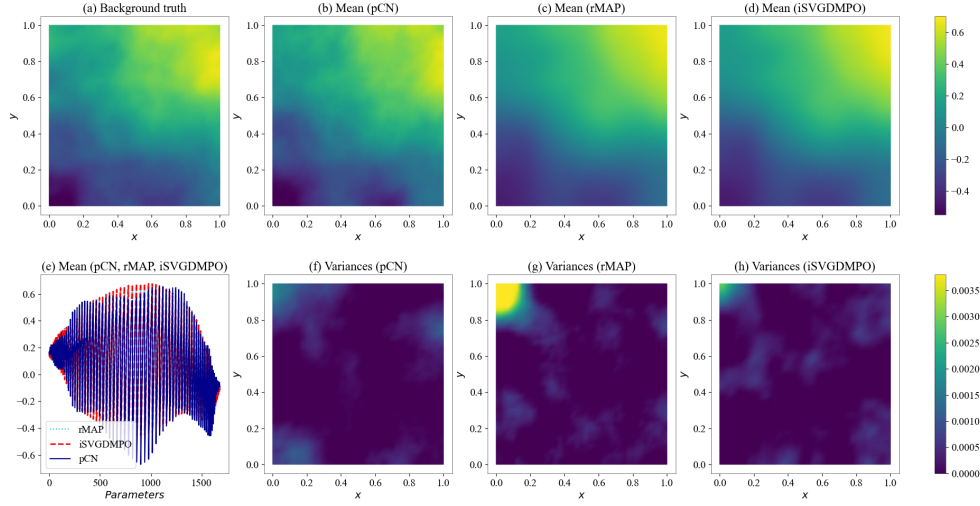876    similar, which are slightly smoother than the one obtained by the pCN algorithm.

FIG. 3. *The background truth and the estimated mean and variance functions by the pCN, rMAP, and iSVGDMPO. (a): The background truth; (b): The estimated mean function by the pCN; (c): The estimated mean function by the rMAP; (d): The estimated mean function by the iSVGDMPO; (e): The estimated mean function on mesh points by the pCN (blue solid line), rMAP (light blue dotted line), and iSVGDMPO (red dashed line); (f): The estimated variances by the pCN; (g): The estimated variances by the rMAP; (h): The estimated variances by the iSVGDMPO.*

This may be caused by the inexact matrix-free Newton-conjugate gradient algorithm [4]. As investigated in [59], many more powerful Newton-type algorithms can be employed to improve the performance both of the rMAP and iSVGDMPO. For the variances, the iSVGDMPO gives more reliable estimates compared with the rMAP, as can be seen from Figure 3 (f), (g), and (h).

Next, we provide some more comparisons of statistical quantities between the results obtained by the pCN, rMAP, and iSVGDMPO. The samples are discretization of functions. As introduced in [49], the mean, variance and covariance functions are the main statistics for functional data. The variance function denoted by $\text{var}_u(x)$ can be defined as $\text{var}_u(x) = \frac{1}{m}\sum_{i=1}^{m}(u_i(x) - \bar{u}(x))^2$, where $x \in \Omega$ is a point residing in the domain $\Omega$, $\bar{u}$ is the mean function, and $m$ is the sample number. The covariance function can be defined as $\text{cov}_u(x_1, x_2) = \frac{1}{m-1}\sum_{i=1}^{m}(u_i(x_1) - \bar{u}(x_1))(u_i(x_2) - \bar{u}(x_2))$, where $x_1, x_2 \in \Omega$ and $m, \bar{u}$ are defined as in $\text{var}_u(x)$. For simplicity, we compute these quantities on the mesh points and exhibit the results in Figure 4. In all of the subfigures in Figure 4, the estimates obtained by the pCN, rMAP, and iSVGDMPO are drawn in blue solid line, gray dotted line, and red dashed line, respectively. In Figure 4 (a), we show the variance function calculated on all of the mesh points, i.e., $\{\text{var}_u(x_i)\}_{i=1}^{N_g}$ ($N_g$ is the number of mesh points). In Figure 4 (c) and (e), we show the covariance function calculated on the pairs of points $\{(x_i, x_{i+50})\}_{i=1}^{N_g-50}$ and $\{(x_i, x_{i+100})\}_{i=1}^{N_g-100}$, respectively. Compared with the estimates given by the rMAP, we can find that the estimates obtained by the iSVGDMPO are visually more similar to the estimates provided by the pCN. In Figure 4 (b), (d), and (f), we provide the same estimates shown in (a), (c), and (e) with points indexing from 1000 to 1200, which give more detailed comparisons. The results also confirm that the iSVGDMPO provides more similar estimates to the pCN.
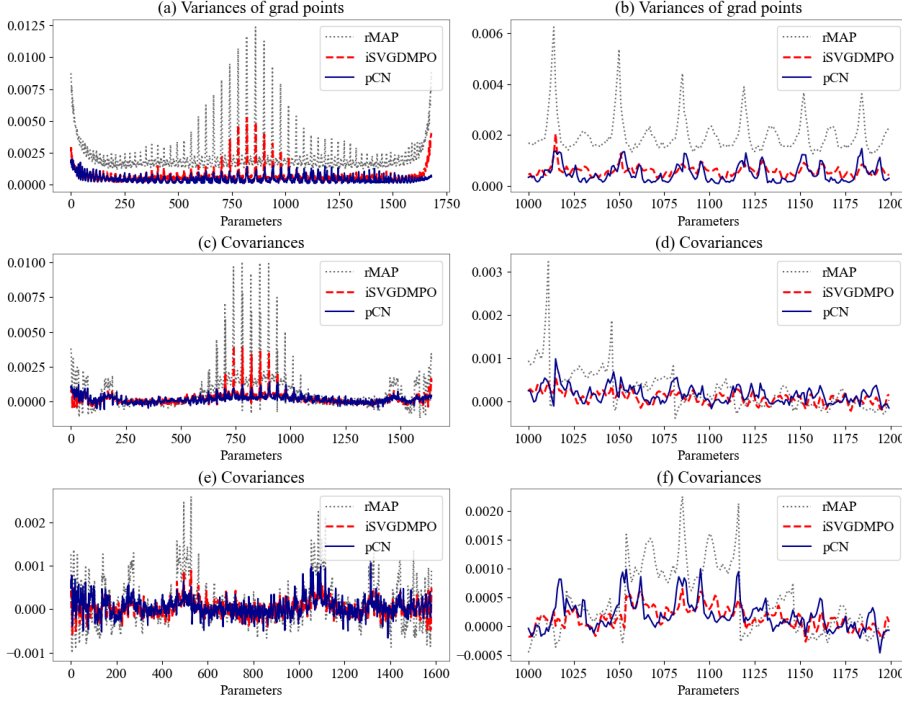
FIG. 4. *The estimated variances and covariances by the pCN (blue solid line), rMAP (gray dotted line), and iSVGDMPO (red dashed line). (a): The estimated variances $\{var_u(x_i)\}_{i=1}^{N_g}$ on all mesh points; (b): The estimated variances for mesh points with indexes from 1000 to 1200 (show details); (c): The estimated covariances $\{cov_u(x_i, x_{i+50})\}_{i=1}^{N_g-50}$ on mesh point pairs $\{(x_i, x_{i+50})\}_{i=1}^{N_g-50}$; (d): The estimated covariances shown in (c) with indexes from 1000 to 1200 (show details); (e): The estimated covariances $\{cov_u(x_i, x_{i+100})\}_{i=1}^{N_g-100}$ on mesh point pairs $\{(x_i, x_{i+100})\}_{i=1}^{N_g-100}$; (f): The estimated covariances shown in (e) with indexes from 1000 to 1200 (show details).*

902      In addition, a quantitative comparison among the pCN, rMAP, and iSVGDMPO
903 are given in Table 1. We compute the $\ell^2$-norm differences of the variance and covari-
904 ance functions on the mesh points obtained by the pCN, rMAP, and iSVGDMPO. In
905 the table, the notation $\text{cov}_u(x_i, x_{i+k})$ ($k = 10, 20, \ldots, 110$) means the covariance func-
906 tion values on the pair of mesh points $\{(x_i, x_{i+k})\}_{i=1}^{N_g}$. The numbers below this nota-
907 tion are the $\ell^2$ differences between the vectors obtained by the rMAP and iSVGDMPO
908 with the pCN, respectively. All of the $\ell^2$ differences of the iSVGDMPO with the pCN
909 are much smaller than the corresponding values of rMAP, which show the superiority
910 of the iSVGDMPO.

911     **5. Conclusion.** In this paper, the approximate sampling algorithm is proposed
912 for the infinite-dimensional Bayesian approach. We introduce the Stein operator on
913 Hilbert spaces and show that it is the limit of a particular finite-dimensional version.
914 Besides, we construct the update perturbation of the SVGD on infinite-dimensional
915 space (called iSVGD) by using the properties of operator-valued RKHS. To accelerate
916 the convergence speed of iSVGD, we investigate the change of variables formula and
917 introduced preconditioning operators. As examples, we present the fixed precondition-
918 ing operators and mixture preconditioning operators. Then, we calculate the explicit

TABLE 1
*The $\ell^2$-norm error of the variance and covariance functions on mesh points for the rMAP and iSVGDMPO (the estimates of the pCN are seen as the background truth).*

|  | $\text{var}_u(x_i)$ | $\text{cov}_u(x_i, x_{i+10})$ | $\text{cov}_u(x_i, x_{i+20})$ | $\text{cov}_u(x_i, x_{i+30})$ |
|---|---|---|---|---|
| rMAP | 0.00759 | 0.00100 | 0.00075 | 0.00092 |
| iSVGDMPO | 0.00038 | 0.00012 | 0.00009 | 0.00010 |
|  | $\text{cov}_u(x_i, x_{i+40})$ | $\text{cov}_u(x_i, x_{i+50})$ | $\text{cov}_u(x_i, x_{i+60})$ | $\text{cov}_u(x_i, x_{i+70})$ |
| rMAP | 0.00227 | 0.00038 | 0.00043 | 0.00056 |
| iSVGDMPO | 0.00015 | 0.00007 | 0.00006 | 0.00007 |
|  | $\text{cov}_u(x_i, x_{i+80})$ | $\text{cov}_u(x_i, x_{i+90})$ | $\text{cov}_u(x_i, x_{i+100})$ | $\text{cov}_u(x_i, x_{i+110})$ |
| rMAP | 0.00142 | 0.00029 | 0.00031 | 0.00047 |
| iSVGDMPO | 0.00012 | 0.00006 | 0.00006 | 0.00007 |

form of the update directions for the iSVGD with mixture preconditioning operators (iSVGDMPO). Finally, we apply the constructed algorithms to an inverse problem of the steady state Darcy flow equation. Comparing with the pCN and rMAP sampling algorithms, we demonstrate by numerical experiments that the proposed algorithms can generate accurate estimates efficiently.

The iSVGD is analyzed by studying the limiting behavior of the finite-dimensional objects. This work presents an infinite-dimensional version of the approach given in [58]. It is worth mentioning that our results not only provide an infinite-dimensional version but also indicate that an intuitive trivial generalization of algorithms given in [58] may not be suitable since particles will belong to a set with zero measure. Our results also show that it is necessary to introduce the parameter $s$, which has not been considered in the existing work.

The current work may be extended to combine the generalizations of the kernel using Hessian operators in the Wasserstein space [36]. The proposed approach may be combined with other algorithms, such as the accelerated information gradient flows [60] and the mean-field type MCMC algorithms [22], to generate new and more efficient algorithms. It is also interesting and important to do more theoretical studies, e.g., introduce infinite-dimensional Stein geometry [33] and develop systematic theories of the interacting particle system and the mean field limit equation [42]. We will report the progress on these aspects elsewhere in the future.

## REFERENCES

[1] S. ARRIDGE, P. MAASS, O. ÖKTEM, AND C.-B. SCHÖNLIEB, *Solving inverse problems using data-driven models*, Acta Numer., 28 (2019), pp. 1–174.

[2] A. BESKOS, A. JASRA, E. A. MUZAFFER, AND A. M. STUART, *Sequential Monte Carlo methods for Bayesian elliptic inverse problems*, Stat. Comput., 25 (2015), p. 727–737.

[3] C. M. BISHOP, *Pattern Recognition and Machine Learning*, Springer-Verlag, New York, NY, USA, 2006.

[4] T. BUI-THANH, O. GHATTAS, J. MARTIN, AND G. STADLER, *A computational framework for infinite-dimensional Bayesian inverse problems part I: The linearized case, with application to global seismic inversion*, SIAM J. Sci. Comput., 35 (2013), pp. A2494–A2523.

[5] M. BURGER AND F. LUCKA, *Maximum a posteriori estimates in linear inverse problems with log-concave priors are proper Bayes estimators*, Inverse Probl., 30 (2014), p. 114004.

[6] T. BUT-THANH AND Q. P. NGUYEN, *FEM-based discretization-invariant MCMC methods for PDE-constrained Bayesian inverse problems*, Inverse Probl. Imag., 10 (2016), pp. 943–975.

[7] C. CARMELI, E. D. VITO, AND A. TOIGO, *Vector valued reproducing kernel Hilbert spaces of*

*integrable functions and Mercer theorem*, Anal. Appl., 4 (2006), pp. 377–408.

[8] C. CARMELI, E. D. VITO, AND A. TOIGO, *Vector-valued reproducing kernel Hilbert spaces and universality*, Anal. Appl., 8 (2010), pp. 19–61.

[9] E. D. C. CARVALHO, R. CLARK, A. NICASTRO, AND P. H. J. KELLY, *Scalable uncertainty for computer vision with functional variational inference*, in CVPR, 2020, pp. 12003–12013.

[10] P. CHEN AND O. GHATTAS, *Stein variational reduced basis Bayesian inversion*, SIAM J. Sci. Comput., 43 (2021), pp. A1163–A1193.

[11] P. CHEN, K. WU, J. CHEN, T. O'LEARY-ROSEBERRY, AND O. GHATTAS, *Projected Stein variational Newton: a fast and scalable Bayesian inference method in high dimensions*, in NeurIPS, vol. 32, 2019.

[12] S. L. COTTER, M. DASHTI, J. C. ROBINSON, AND A. M. STUART, *Bayesian inverse problems for functions and applications to fluid mechanics*, Inverse Probl., 25 (2009), p. 115008.

[13] S. L. COTTER, G. O. ROBERTS, A. M. STUART, AND D. WHITE, *MCMC methods for functions: modifying old algorithms to make them faster*, Stat. Sci., 28 (2013), pp. 424–446.

[14] T. CUI, K. J. H. LAW, AND Y. M. MARZOUK, *Dimension-independent likelihood-informed MCMC*, J. Comput. Phys., 304 (2016), pp. 109–137.

[15] G. DAPRATO AND J. ZABCZYK, *Stochastic Equations in Infinite Dimensions*, Cambridge University Press, Cambridge, 1992.

[16] M. DASHTI AND A. M. STUART, *The Bayesian approach to inverse problems*, Handbook of Uncertainty Quantification, (2017), pp. 311–428.

[17] G. DETOMMASO, T. CUI, A. SPANTINI, AND Y. MARZOUK, *A Stein variational Newton method*, in NeurIPS, vol. 32, 2018.

[18] A. DUNCAN, N. NÜSKEN, AND L. SZPRUCH, *On the geometry of Stein variational gradient descent.* arXiv:1912.00894, 2019.

[19] H. W. ENGL, M. HANKE, AND A. NEUBAUER, *Regularization of Inverse Problems*, Springer, Netherlands, 1996.

[20] Z. FENG AND J. LI, *An adaptive independence sampler MCMC algorithm for Bayesian inferences of functions*, SIAM J. Sci. Comput., 40 (2018), pp. A1310–A1321.

[21] A. FICHTNER, *Full Seismic Waveform Modelling and Inversion*, Springer, New York, 2011.

[22] A. GARBUNO-INIGO, F. HOFFMANN, W. C. LI, AND A. M. STUART, *Interacting Langevin diffusions: gradient structure and ensemble Kalman sampler*, SIAM J. Appl. Dyn. Syst., 19 (2020), pp. 412–441.

[23] N. GUHA, X. WU, Y. EFENDIEV, B. JIN, AND B. K. MALICK, *A variational Bayesian approach for inverse problems with skew-t error distribution*, J. Comput. Phys., 301 (2015), pp. 377–393.

[24] T. HELIN AND M. BURGER, *Maximum a posteriori probability estimates in infinite-dimensional Bayesian inverse problems*, Inverse Probl., 31 (2015), p. 085009.

[25] J. JIA, J. PENG, AND J. GAO, *Posterior contraction for empirical Bayesian approach to inverse problems under non-diagonal assumption*, Inverse Probl. Imag., 15 (2020), pp. 201–228.

[26] J. JIA, B. WU, J. PENG, AND J. GAO, *Recursive linearization method for inverse medium scattering problems with complex mixture Gaussian error learning*, Inverse Probl., 35 (2019), p. 075003.

[27] J. JIA, S. YUE, J. PENG, AND J. GAO, *Infinite-dimensional Bayesian approach for inverse scattering problems of a fractional Helmholtz equation*, J. Funct. Anal., 275 (2018), pp. 2299–2332.

[28] J. JIA, Q. ZHAO, D. MENG, AND Y. LEUNG, *Variational Bayes' method for functions with applications to some inverse problems*, SIAM J. Sci. Comput., 43 (2021), pp. A355–A383.

[29] B. JIN, *A variational Bayesian method to inverse problems with implusive noise*, J. Comput. Phys., 231 (2012), pp. 423–435.

[30] B. JIN AND J. ZOU, *Hierarchical Bayesian inference for ill-posed problems via variational method*, J. Comput. Phys., 229 (2010), pp. 7317–7343.

[31] H. KADRI, E. DUFLOS, P. PREUS, S. CANU, A. RAKOTOMAMONJY, AND J. AUDIFFREN, *Operator-valued kernels for learning from functional response data*, J. Mach. Learn. Res., 17 (2016), pp. 1–54.

[32] J. KAIPIO AND E. SOMERSALO, *Statistical and Computational Inverse Problems*, Springer-Verlag, New York, 2005.

[33] A. KORBA, A. SALIM, M. ARBEL, G. LUISE, AND A. GRETTON, *A non-asymptotic analysis for Stein variational gradient descent*, in NeurIPS, vol. 33, 2020.

[34] J. LEI, *Convergence and concentration of empirical measures under wasserstein distance in unbounded functional space*, Bernoulli, 26 (2020), pp. 767–798.

[35] D. A. LEVIN, Y. PERES, AND E. L. WILMER, *Markov Chains and Mixing Times*, American Mathematical Society, second ed., 2017.

[36] W. C. Li, *Hessian metric via transport information geometry*, J. Math. Phys, 62 (2021), p. 033301.

[37] C. Liu, J. Zhuo, P. Cheng, R. Zhang, and J. Zhu, *Understanding and accelerating particle-based variational inference*, in ICML, vol. 97, 2019, pp. 4082–4092.

[38] Q. Liu, *Stein variational gradient descent as gradient flow*, in NeurIPS, vol. 30.

[39] Q. Liu and D. Wang, *Stein variational gradient descent: A general purpose Bayesian inference algorithm*, in NeurIPS, vol. 29, 2016.

[40] A. Logg, K. A. Mardal, and G. N. Wells, *Automated Solution of Differential Equations by the Finite Element Method*, Springer, United Kingdom, 2012.

[41] J. C. D. los Reyes, *Numerical PDE-Constrained Optimization*, Springer, New York, 2015.

[42] J. Lu, Y. Lu, and J. Nolen, *Scaling limit of the Stein variational gradient descent: the mean field regime*, SIAM J. Math. Anal., 5 (2019), pp. 648–671.

[43] A. G. D. G. Matthews, *Scalable Gaussian process inference using variational methods*, PhD thesis, University of Cambridge, 9 2016.

[44] R. Nickl, *Betnstein-von Mises theorem for statistical inverse problems I: Schrödinger equation*, J. Eur. Math. Soc., 22 (2020), pp. 2697–2750.

[45] F. J. Pinski, G. Simpson, A. M. Stuart, and H. Weber, *Algorithms for Kullback-Leibler approximation of probability measures in infinite dimensions*, SIAM J. Sci. Comput., 37 (2015), pp. A2733–A2757.

[46] F. J. Pinski, G. Simpson, A. M. Stuart, and H. Weber, *Kullback-Leibler approximation for probability measures on infinite dimensional space*, SIAM J. Math. Anal., 47 (2015), pp. 4091–4122.

[47] G. D. Prato, *Kolmogorov Equations for Stochastic PDEs*, Birkhäuser Verlag, Basel, 2004.

[48] G. D. Prato, *An Introduction to Infinite-Dimensional Analysis*, Springer-Verlag, Berlin, 2006.

[49] J. O. Ramsay and B. W. Silverman, *Functional Data Analysis*, Springer, New York, second ed., 2005.

[50] M. Reed and B. Simon, *Functional Analysis I: Methods of Modern Mathematical Physics*, Elsevier (Singapore) Pte Ltd, revised and enlarged ed., 2003.

[51] A. Spantini, A. Solonen, T. Cui, J. Martin, L. Tenorio, and Y. Marzouk, *Optimal low-rank approximations of Bayesian linear inverse problems*, SIAM J. Sci. Comput., 37 (2015), pp. A2451–A2487.

[52] I. Steinwart and A. Christmann, *Support Vector Machines*, Springer, Germany, 2006.

[53] A. M. Stuart, *Inverse problems: A Bayesian perspective*, Acta Numer., 19 (2010), pp. 451–559.

[54] S. Sun, G. Zhang, J. Shi, and R. Grosse, *Functional variational Bayesian neural networks*, in ICLR, 2019.

[55] A. Tarantola, *Inverse Problem Theory and Methods for Model Parameter Estimation*, SIAM, United States, 2005.

[56] A. Tarantola and B. Valette, *Inverse problems = quset for information*, J. Geophys., 50 (1982), pp. 159–170.

[57] N. G. Trillos and D. Slepˇcev, *On the rate of convergence of empirical measures in ∞-transportation distance*, Canad. J. Math., 67 (2015), pp. 1358–1383.

[58] D. Wang, Z. Tang, C. Bajaj, and Q. Liu, *Stein variational gradient descent with matrix-valued kernels*, in NeurIPS, vol. 33, 2019.

[59] K. Wang, T. Bui-Thanh, and O. Ghattas, *A randomized maximum a posteriori method for posterior sampling of high dimensional nonlinear Bayesian inverse problems*, SIAM J. Sci. Comput., 40 (2018), pp. A142–A171.

[60] Y. Wang and W. C. Li, *Accelerated information gradient flows*. arXiv:1909.02102, 2020.

[61] Z. Wang, T. Ren, J. Zhu, and B. Zhang, *Function space particle optimization for Bayesian neural networks*, in ICLR, 2019.

[62] C. Zhang, J. Butepage, H. Kjellstrom, and S. Mandt, *Advances in variational inference*, IEEE T. Pattern Anal., 41 (2018), pp. 2008–2026.

[63] Q. Zhao, D. Meng, Z. Xu, W. Zuo, and Y. Yan, $l_1$-*norm low-rank matrix factorization by variational Bayesian method*, IEEE T. Neur. Net. Lear., 26 (2015), pp. 825–839.

[64] Q. Zhou, T. Yu, X. Zhang, and J. Li, *Bayesian inference and uncertainty quantification for medical image reconstruction with poisson data*, SIAM J. Imaging Sci., 13 (2020), pp. 29–52.

# SUPPLEMENTARY MATERIAL: STEIN VARIATIONAL GRADIENT DESCENT ON INFINITE-DIMENSIONAL SPACE AND APPLICATIONS TO STATISTICAL INVERSE PROBLEMS

JUNXIONG JIA, PEIJUN LI, AND DEYU MENG

ABSTRACT. In this supplementary material, we present the details for some of the results and examples given in the main text.

## CONTENTS

## 1. EXAMPLE OF KERNEL SATISFYING ASSUMPTION (33) IN THEOREM 10

In this section, we present an example of the kernel that satisfies the assumption (33) given in Theorem 10. Let us recall the assumption (33)

$$\mathbb{E}_{u \sim \mu}\Big[ D_{u'} \boldsymbol{K}(u, u') \mathcal{C}_0^{-1/2} g + \sum_{k=1}^{\infty} D_k D_{u'} \boldsymbol{K}(u, u') e_k \Big], \tag{1.1}$$

which belongs to $\mathcal{L}_1(\mathcal{X}, \mathcal{Y})$ for each $u' \in \mathcal{X}$ and $g \in \mathcal{H}^{-s}$.

Taking $\boldsymbol{K}(u, u') = K(u, u') \mathrm{Id}$ with

$$K(u, u') = \exp\left( -\frac{1}{h} \|u - u'\|_{\mathcal{X}}^2 \right)$$

being a scalar-valued kernel, we have

$$D_{u'} \boldsymbol{K}(u, u') = -\frac{1}{h} \langle u - u', \cdot \rangle_{\mathcal{X}} K(u, u'),$$

$$D_k D_{u'} \boldsymbol{K}(u, u') e_k = \frac{1}{h^2} (u_k - u'_k) \langle u - u', \cdot \rangle_{\mathcal{X}} K(u, u') e_k.$$

1

Let $\{\varphi_j\}_{j=1}^{\infty}$ be an orthonormal basis of $\mathcal{X}$, and recall that $\{e_j\}_{j=1}^{\infty}$ represents an orthonormal basis of $\mathcal{Y}$. Plugging the above formula into (1.1), we find that

$$\sum_{j=1}^{\infty}\langle\mathbb{E}_{u\sim\mu}\Big[D_{u'}\boldsymbol{K}(u,u')\mathcal{C}_0^{-1/2}g + \sum_{k=1}^{\infty}D_kD_{u'}\boldsymbol{K}(u,u')e_k\Big]\varphi_j, e_j\rangle = \mathbb{E}_{u\sim\mu}\Big\{\mathrm{I}+\mathrm{II}\Big\},$$

where

$$\mathrm{I} = -\frac{1}{h}\sum_{j=1}^{\infty}\langle\langle u-u', \varphi_j\rangle_{\mathcal{X}}K(u,u')\mathcal{C}_0^{-1/2}g, e_j\rangle_{\mathcal{Y}},$$

$$\mathrm{II} = \frac{1}{h^2}\sum_{j=1}^{\infty}\langle\sum_{k=1}^{\infty}(u_k-u_k')e_k\langle u-u', \varphi_j\rangle_{\mathcal{X}}, e_j\rangle_{\mathcal{Y}}K(u,u').$$

For term I, we have

$$
\begin{aligned}
\mathrm{I} \leq& \frac{1}{h}K(u,u')\Big(\sum_{j=1}^{\infty}\langle u-u', \varphi_j\rangle_{\mathcal{X}}^2\Big)^{1/2}\Big(\sum_{j=1}^{\infty}\langle\mathcal{C}_0^{-1/2}g, e_j\rangle\Big)^{1/2}\\
\leq& \frac{C}{h}K(u,u')\|u-u'\|_{\mathcal{X}}\|\mathcal{C}_0^{-1/2}g\|_{\mathcal{Y}}\\
\leq& \frac{C}{h}K(u,u')\|u-u'\|_{\mathcal{X}}\|g\|_{\mathcal{H}^{-s}} < \infty.
\end{aligned}
\tag{1.2}
$$

For term II, we have

$$
\begin{aligned}
\mathrm{II} \leq& \frac{1}{h^2}K(u,u')\sum_{j=1}^{\infty}\langle u-u', \varphi_j\rangle_{\mathcal{X}}\langle(u_j-u_j')e_j, e_j\rangle_{\mathcal{Y}}\\
\leq& \frac{1}{h^2}K(u,u')\Big(\sum_{j=1}^{\infty}\langle u-u', \varphi_j\rangle_{\mathcal{X}}^2\Big)^{1/2}\Big(\sum_{j=1}^{\infty}(u_j-u_j')^2\Big)^{1/2}\\
\leq& \frac{1}{h^2}K(u,u')\|u-u'\|_{\mathcal{X}}^2 < \infty.
\end{aligned}
\tag{1.3}
$$

Combining estimates (1.2) and (1.3) yields

$$\sum_{j=1}^{\infty}\langle\mathbb{E}_{u\sim\mu}\Big[D_{u'}\boldsymbol{K}(u,u')\mathcal{C}_0^{-1/2}g + \sum_{k=1}^{\infty}D_kD_{u'}\boldsymbol{K}(u,u')e_k\Big]\varphi_j, e_j\rangle < \infty, \tag{1.4}$$

which implies that (1.1) belongs to $\mathcal{L}_1(\mathcal{X}, \mathcal{Y})$. Taking $\mathcal{X} = \mathcal{H}^1, \mathcal{Y} = \mathcal{H}^{-1}$, $s = 0$, and projecting all of the quantities to $\mathcal{X}^N$, we then obtain the finite-dimensional SVGD as reviewed in Section 2 of the main text.

## 2. Implementation details for the mixture preconditioning

In Subsection 3.3, we present the mixture preconditioning operators, which can specify different preconditioning operators for different particles. Here, we provide some more implementation details.

In practice, we approximate the expectation $\mathbb{E}_{u \sim \mu}$ by empirical mean of particles $\{u_i\}_{i=1}^m$. Hence, the formula (56) reduces to

$$\phi_{\boldsymbol{K}}^*(\cdot) = \sum_{\ell=1}^m w_\ell(\cdot) \sum_{j=1}^m \left[ -w_\ell(u_j) \boldsymbol{K}_\ell(u_j, \cdot)(D_{u_j}\Phi(u_j) + \mathcal{C}_0^{-1}u_j) \right.$$
$$\left. + \sum_{k=1}^\infty D_k(w_\ell(u_j)\boldsymbol{K}_\ell(u_j, \cdot)e_k) \right]. \qquad (2.1)$$

Taking $\boldsymbol{K}_\ell$ in (54) with

$$T_\ell = \mathcal{C}_0^{s/2}(D\mathcal{G}(u_\ell)^*\Sigma^{-1}D\mathcal{G}(u_\ell) + \mathcal{C}_0^{-1})^{1/2},$$

we get

$$\boldsymbol{K}_\ell(u_j, \cdot)(D_{u_j}\Phi(u_j) + \mathcal{C}_0^{-1}u_j)$$
$$= (D\mathcal{G}(u_\ell)^*\Sigma^{-1}D\mathcal{G}(u_\ell) + \mathcal{C}_0^{-1})^{-1}(D_{u_j}\Phi(u_j) + \mathcal{C}_0^{-1}u_j)\exp\left(-\frac{1}{h}\|T_\ell(u_j - \cdot)\|_{\mathcal{H}}^2\right).$$

For the term $D_k(w_\ell(u_j)\boldsymbol{K}_\ell(u_j, \cdot)e_k)$, it is clear to note

$$D_k(w_\ell(u_j)\boldsymbol{K}_\ell(u_j, \cdot)e_k) = D_k w_\ell(u_j)\boldsymbol{K}_\ell(u_j, \cdot)e_k + w_\ell(u_j)D_k\boldsymbol{K}_\ell(u_j, \cdot)e_k. \qquad (2.2)$$

For the first term, we have

$$D_k w_\ell(u_j)\boldsymbol{K}_\ell(u_j, \cdot)e_k = -\langle T_\ell(u_j - u_\ell), T_\ell\varphi_k\rangle_{\mathcal{H}} w_\ell(u_j)\boldsymbol{K}_\ell(u_j, \cdot)e_k$$
$$- J_k w_\ell(u_j)\boldsymbol{K}_\ell(u_j, \cdot)e_k, \qquad (2.3)$$

where

$$J_k = \frac{\sum_{\ell'=1}^m \langle T_\ell(u_j - u_{\ell'}), T_\ell\varphi_k\rangle_{\mathcal{H}} \exp\left(-\frac{1}{2}\|T_{\ell'}(u_j - u_{\ell'})\|_{\mathcal{H}}^2\right)}{\sum_{\ell'=1}^m \exp\left(-\frac{1}{2}\|T_{\ell'}(u_j - u_{\ell'})\|_{\mathcal{H}}^2\right)}. \qquad (2.4)$$

For the second term, we have

$$w_\ell(u_j)D_k\boldsymbol{K}_\ell(u_j, \cdot)e_k = -\frac{2}{h}w_\ell(u_j)\langle T_\ell(u_j - \cdot), T_\ell\varphi_k\rangle_{\mathcal{H}}\boldsymbol{K}_\ell(u_j, \cdot)e_k. \qquad (2.5)$$

Combining (2.3), (2.4) and (2.5), we obtain

$$\sum_{k=1}^\infty D_k(w_\ell(u_j)\boldsymbol{K}_\ell(u_j, \cdot)e_k) = -\frac{2}{h}w_\ell(u_j)\sum_{k=1}^\infty \langle T_\ell(u_j - \cdot), T_\ell\varphi_k\rangle_{\mathcal{H}}\boldsymbol{K}_\ell(u_j, \cdot)e_k$$
$$- w_\ell(u_j)\sum_{k=1}^\infty \langle T_\ell(u_j - u_\ell), T_\ell\varphi_k\rangle_{\mathcal{H}}\boldsymbol{K}_\ell(u_j, \cdot)e_k \qquad (2.6)$$
$$- w_\ell(u_j)\sum_{\ell'=1}^m \sum_{k=1}^\infty \langle T_{\ell'}(u_j - \cdot), T_{\ell'}\varphi_k\rangle_{\mathcal{H}}\boldsymbol{K}_\ell(u_j, \cdot)e_k M_{\ell'},$$

where

$$M_{\ell'} = \frac{\exp\left(-\frac{1}{2}\|T_{\ell'}(u_j - u_{\ell'})\|_{\mathcal{H}}^2\right)}{\sum_{\ell''=1}^m \exp\left(-\frac{1}{2}\|T_{\ell''}(u_j - u_{\ell''})\|_{\mathcal{H}}^2\right)}. \qquad (2.7)$$

For specific examples, we have the explicit form

$$\sum_{k=1}^{\infty}\langle T_\ell(u_j-\cdot),T_\ell\varphi_k\rangle_{\mathcal{H}}\boldsymbol{K}_\ell(u_j,\cdot)e_k. \tag{2.8}$$

For example, we take $\mathcal{X}$, $\mathcal{Y}$, $\tilde{\mathcal{X}}$, and $\tilde{\mathcal{Y}}$ as in the fixed precondition case and specify $\boldsymbol{K}_\ell$ as in (51) with $T$ replaced by $T_\ell$. Then, we have

$$\sum_{k=1}^{\infty}\langle T_\ell(u_j-\cdot),T_\ell\varphi_k\rangle_{\mathcal{H}}\boldsymbol{K}_\ell(u_j,\cdot)e_k$$
$$= \exp\Big(-\frac{1}{h}\|T_\ell(u_j-\cdot)\|_{\mathcal{H}}^2\Big)T_\ell^{-1}\mathcal{C}_0^s(T_\ell^{-1})^*\mathcal{C}_0^{-s}T_\ell^*T_\ell(u_j-\cdot). \tag{2.9}$$

Hence, it is not required to calculate the orthonormal basis $\{e_k\}_{k=1}^{\infty}$ and $\{\varphi_i\}_{i=1}^{\infty}$ in spaces $\mathcal{Y}$ and $\mathcal{X}$ explicitly in the implementations.

## 3. PROOF OF THEOREM 16

Blow is the proof of Theorem 16.

*Proof.* Denote by $\mathcal{E}(\mathcal{H})$ the set of all the exponential functions and let

$$\varphi_h(x) := e^{i\langle x,h\rangle_{\mathcal{H}}}, \quad x,h \in \mathcal{H}. \tag{3.1}$$

By [18], the function space $\mathcal{E}(\mathcal{H})$ is dense in $L^2(\mathcal{H},\mu_0)$, where $\mu_0$ is the prior measure. Let $K_n, \varphi_n \in \mathcal{E}(\mathcal{H})$ satisfy

$$\lim_{n\to\infty}\|K_n-K\|_{W^{1,2}(\mathcal{H},\mu_0)}=0, \quad \lim_{n\to\infty}\|\psi_n-\exp(-\Phi)\|_{L^2(\mathcal{H},\mu_0)}=0.$$

For the prior probability measure, we have $\mathcal{C}_0\varepsilon_k = \lambda_k^2\varepsilon_k$ with $k=1,2,\ldots$, i.e., $\{\lambda_k^2,\varepsilon_k\}_{k=1}^{\infty}$ is the eigensystem of $\mathcal{C}_0$. It follows from [18, Lemma 1.5] that we have

$$\int_{\mathcal{H}}D_kK_n(u-\tilde{u})\psi_{n'}(\tilde{u})\mu_0(d\tilde{u})=-\int_{\mathcal{H}}K_n(u-\tilde{u})D_k\psi_{n'}(\tilde{u})\mu_0(d\tilde{u})$$
$$+\frac{1}{\lambda_k^2}\int_{\mathcal{H}}\tilde{u}_kK_n(u-\tilde{u})\psi_{n'}(\tilde{u})\mu_0(d\tilde{u}), \tag{3.2}$$

where $\tilde{u}_k = \langle u,\varepsilon_k\rangle_{\mathcal{H}}, k=1,2,\ldots$. By a simple calculation, we have

$$-\int_{\mathcal{H}}D_kK_n(u-\tilde{u})\psi_{n'}(\tilde{u})\mu_0(d\tilde{u})=\int_{\mathcal{H}}K_n(u-\tilde{u})\Big(D_k\psi_{n'}(\tilde{u})+\frac{\tilde{u}_k}{\lambda_k^2}\Big)\psi_{n'}(\tilde{u})\mu_0(d\tilde{u}).$$

Taking $n' \to \infty$ in the above equality leads to

$$\int_{\mathcal{H}}D_kK_n(u-\tilde{u})e^{-\Phi(\tilde{u};\boldsymbol{d})}-K_n(u-\tilde{u})D_kV(\tilde{u})e^{-\Phi(\tilde{u};\boldsymbol{d})}\mu_0(d\tilde{u})=0, \tag{3.3}$$

where $V(\cdot)$ is defined in (9). Taking $n \to \infty$, we arrive at

$$\int_{\mathcal{H}}D_kK(u-\tilde{u})e^{-\Phi(\tilde{u};\boldsymbol{d})}-K(u-\tilde{u})D_kV(\tilde{u})e^{-\Phi(\tilde{u};\boldsymbol{d})}\mu_0(d\tilde{u})=0, \tag{3.4}$$

which implies

$$\int_{\mathcal{H}}\langle DK(u-\tilde{u}),D\varphi(u)\rangle_{\mathcal{H}}+\langle K(u-\tilde{u})DV(\tilde{u}),D\varphi(u)\rangle_{\mathcal{H}}\mu^{\boldsymbol{d}}(d\tilde{u})=0, \tag{3.5}$$

where $\varphi \in \mathcal{E}(\mathcal{H})$ is a test function. Through simple calculations based on (3.5), we further obtain

$$\int_{\mathcal{H}} \langle DK * \mu^{\boldsymbol{d}}, D\varphi \rangle_{\mathcal{H}} + \langle K * DV\mu^{\boldsymbol{d}}, D\varphi \rangle_{\mathcal{H}} \mu^{\boldsymbol{d}}(du) = 0, \qquad (3.6)$$

which implies

$$\langle \mu^{\boldsymbol{d}}, L(\mu^d)\varphi \rangle = 0 \qquad (3.7)$$

with $L$ being defined in (59). Recalling the weak form of the equation (58), we complete the proof. $\qquad \square$

## 4. Proof of Theorem 17

Let $H_0^1(\Omega)$ and $H^{-1}(\Omega)$ be the usual Sobolev spaces. Consider the boundary value problem

$$\begin{aligned} -\nabla \cdot (e^u \nabla w) &= f \quad \text{in } \Omega, \\ w &= 0 \quad \text{on } \partial\Omega. \end{aligned} \qquad (4.1)$$

The following estimate is crucial to our proofs.

**Theorem 4.1.** *Let $u \in L^\infty(\Omega)$ and $f \in H^{-1}(\Omega)$, then Eq. (4.1) has a unique solution $w \in H_0^1(\Omega)$ satisfies*

$$\|w\|_{H_0^1(\Omega)} \leq Ce^{\|u\|_{L^\infty(\Omega)}} \|f\|_{H^{-1}(\Omega)}, \qquad (4.2)$$

*where $C$ is a positive constant independent of $u$.*

Using Theorem 4.1, we can derive the estimates for the adjoint, incremental forward, and incremental adjoint equations. For the adjoint equation, we have

$$\|p\|_{H_0^1(\Omega)} \leq Ce^{\|u\|_{L^\infty}} \left\| \sum_{j=1}^{N_d} e^{\frac{1}{2\delta^2}\|x-x_j\|^2} (\ell_{x_j}(w) - d_j) \right\|_{L^2}. \qquad (4.3)$$

Let $|\Omega|$ be the volume of domain $\Omega$. Since

$$\left\| e^{-\frac{1}{2\delta^2}\|x-x_j\|^2} \right\|_{L^2} \leq |\Omega|^{1/2} \qquad (4.4)$$

and

$$\ell_{x_j}(w) \leq |\Omega|^{1/2} \|w\|_{L^2} \quad \text{for } j = 1, \ldots, N_d, \qquad (4.5)$$

we deduce

$$\begin{aligned} \|p\|_{H_0^1(\Omega)} &\leq \frac{|\Omega|^{1/2}}{2\pi\delta^2} \sum_{j=1}^{N_d} (\|\boldsymbol{d}\| + |\Omega|^{1/2}\|w\|_{L^2}) \\ &\leq \frac{N_d|\Omega|^{1/2}}{2\pi\delta^2} \left( \|\boldsymbol{d}\| + C|\Omega|^{1/2}e^{\|u\|_{L^\infty}} \|f\|_{H^{-1}} \right) \\ &\leq C \left( 1 + e^{\|u\|_{L^\infty}} \|f\|_{H^{-1}} \right), \end{aligned} \qquad (4.6)$$

which implies

$$\|p\|_{H_0^1(\Omega)} \leq C(1 + \|f\|_{H^{-1}})e^{2\|u\|_{L^\infty}}. \qquad (4.7)$$

For the incremental forward equation, we have

$$
\begin{aligned}
\|\hat{w}\|_{H_0^1} &\le Ce^{\|u\|_{L^\infty}} \|\nabla \cdot (\hat{u}e^u \nabla w)\|_{H^{-1}} \\
&\le Ce^{\|u\|_{L^\infty}} \|\hat{u}e^u \nabla w\|_{L^2} \\
&\le Ce^{2\|u\|_{L^\infty}} \|\hat{u}\|_{L^\infty} \|\nabla w\|_{L^2} \\
&\le Ce^{3\|u\|_{L^\infty}} \|f\|_{H^{-1}} \|\hat{u}\|_{L^\infty}.
\end{aligned}
\tag{4.8}
$$

Similarly, based on Theorem 4.1, we have

$$
\|\hat{p}\|_{H_0^1} \le Ce^{\|u\|_{L^\infty}} [I_1 + I_2],
\tag{4.9}
$$

where

$$
I_1 = \|\nabla \cdot (\hat{u}e^u \nabla p)\|_{H^{-1}},
\tag{4.10}
$$

$$
I_2 = \frac{1}{2\pi\delta^2\sigma^2} \sum_{j=1}^{N_d} \|\ell_{x_j}(\hat{w})e^{-\frac{1}{2\delta^2}\|x-x_j\|^2}\|_{L^2}.
\tag{4.11}
$$

For $I_1$, we have

$$
\begin{aligned}
I_1 &\le \|\hat{u}e^u \nabla p\|_{L^2} \le e^{\|u\|_{L^\infty}} \|\hat{u}\|_{L^\infty} \|\nabla p\|_{L^2} \\
&\le C(1 + \|f\|_{H^{-1}})e^{3\|u\|_{L^\infty}} \|\hat{u}\|_{L^\infty}.
\end{aligned}
\tag{4.12}
$$

For $I_2$, we get

$$
I_2 \le C\|\hat{w}\|_{L^2} \le Ce^{3\|u\|_{L^\infty}} \|f\|_{H^{-1}} \|\hat{u}\|_{L^\infty}.
\tag{4.13}
$$

Combining (4.9) with estimates of $I_1$ and $I_2$, we obtain the estimate of the adjoint equation

$$
\|\hat{p}\|_{H_0^1} \le C(1 + \|f\|_{H^{-1}})e^{4\|u\|_{L^\infty}} \|\hat{u}\|_{L^\infty}.
\tag{4.14}
$$

It is clear to note that

$$
\|u\|_{L^\infty} \le C\|u\|_{\mathcal{H}^{1-s}} = C\|u\|_{\mathcal{X}}
\tag{4.15}
$$

holds for $s < 0.5$, which can be deduced based on similar arguments given in [10, Lemma 16 or Theorem 28]. Since the Hilbert scale is based on the covariance operator $\mathcal{C}_0$ [1, 13], the space $\mathcal{H}^{1-s}$ is different from the one introduced in [10]. The space $\mathcal{H}^{1-s}$ in our paper is approximately equal to the space $\mathcal{H}^{2(1-s)}$ defined in [10]. Next, we give the three estimates shown in Theorem 17.

First is to estimate $\Phi(u)$. A simple calculation gives

$$
\begin{aligned}
\Phi(u) &= \frac{1}{2\sigma^2} \|\mathcal{M}(w) - \boldsymbol{d}\|^2 \le C(1 + \|w\|_{L^2})^2 \\
&\le C(1 + \|f\|_{H^{-1}})^2 e^{2\|u\|_{\mathcal{X}}},
\end{aligned}
\tag{4.16}
$$

where the last inequality used estimates (4.2) and (4.15).

Next is to estimate $D\Phi(u)$. For any $\tilde{u} \in \mathcal{X}$, we get

$$
\begin{aligned}
\langle D\Phi(u), \tilde{u} \rangle = \int_\Omega \tilde{u}e^u \nabla w \cdot \nabla p \, dx &\le \|\tilde{u}\|_{L^\infty} e^{\|u\|_{L^\infty}} \|\nabla w\|_{L^2} \|\nabla p\|_{L^2} \\
&\le C(1 + \|f\|_{H^{-1}})^2 e^{4\|u\|_{L^\infty}} \|\tilde{u}\|_{L^\infty} \\
&\le C(1 + \|f\|_{H^{-1}})^2 e^{4\|u\|_{\mathcal{X}}} \|\tilde{u}\|_{\mathcal{X}},
\end{aligned}
\tag{4.17}
$$

where estimates (4.2) and (4.7) are used to derive the second inequality and estimate (4.15) is used for obtaining the third inequality. Clearly, it follows from (4.17) that

$$\|D\Phi(u)\|_{\mathcal{X}^*} \leq C(1 + \|f\|_{H^{-1}})^2 e^{4\|u\|_{\mathcal{X}}}. \tag{4.18}$$

It is also required to estimate $D^2\Phi(u)$. For any $\tilde{u}, \hat{u} \in \mathcal{X}$, we obtain

$$\langle\langle D^2\Phi(u), \hat{u}\rangle, \tilde{u}\rangle = I_1 + I_2 + I_3, \tag{4.19}$$

where

$$I_1 = \int_\Omega \hat{u}\tilde{u}e^u \nabla w \cdot \nabla p dx, \tag{4.20}$$

$$I_2 = \int_\Omega \tilde{u}e^u \nabla w \cdot \nabla \hat{p} dx, \tag{4.21}$$

$$I_3 = \int_\Omega \tilde{u}e^u \nabla p \cdot \nabla \hat{w} dx. \tag{4.22}$$

For $I_1$, we have

$$
\begin{aligned}
I_1 &\leq \|\hat{u}\|_{L^\infty}\|\tilde{u}\|_{L^\infty} e^{\|u\|_{L^\infty}}\|\nabla w\|_{L^2}\|\nabla p\|_{L^2} \\
&\leq C\|\hat{u}\|_{L^\infty}\|\tilde{u}\|_{L^\infty} e^{\|u\|_{L^\infty}}(1 + \|f\|_{H^{-1}})^2 e^{3\|u\|_{L^\infty}} \\
&\leq C(1 + \|f\|_{H^{-1}})^2 e^{4\|u\|_{\mathcal{X}}}\|\hat{u}\|_{\mathcal{X}}\|\tilde{u}\|_{\mathcal{X}},
\end{aligned} \tag{4.23}
$$

where (4.2) and (4.7) are used for deriving the second inequality and (4.15) is employed to derive the third inequality. By similar calculations, we obtain from (4.3), (4.8), (4.14), and (4.15) that

$$I_2 \leq C(1 + \|f\|_{H^{-1}})^2 e^{6\|u\|_{\mathcal{X}}}\|\hat{u}\|_{\mathcal{X}}\|\tilde{u}\|_{\mathcal{X}} \tag{4.24}$$

and

$$I_3 \leq C(1 + \|f\|_{H^{-1}})^2 e^{6\|u\|_{\mathcal{X}}}\|\hat{u}\|_{\mathcal{X}}\|\tilde{u}\|_{\mathcal{X}}. \tag{4.25}$$

substituting (4.23), (4.24), and (4.25) into (4.19), we obtain

$$\langle\langle D^2\Phi(u), \hat{u}\rangle, \tilde{u}\rangle \leq C(1 + \|f\|_{H^{-1}})^2 e^{6\|u\|_{\mathcal{X}}}\|\hat{u}\|_{\mathcal{X}}\|\tilde{u}\|_{\mathcal{X}}. \tag{4.26}$$

Hence

$$\|D^2\Phi(u)\|_{\mathcal{L}(\mathcal{X}, \mathcal{X}^*)} \leq C(1 + \|f\|_{H^{-1}})^2 e^{6\|u\|_{\mathcal{X}}}, \tag{4.27}$$

which completes the proof.

## 5. MORE NUMERICAL RESULTS FOR THE DARCY FLOW MODEL

In this section, we provide more numerical results for the Darcy flow model given in Section 4 of the main text. We intend to answer the following two questions: how do different optimization methods affect the estimates; how does $s$ influence the convergence speed of the algorithm.
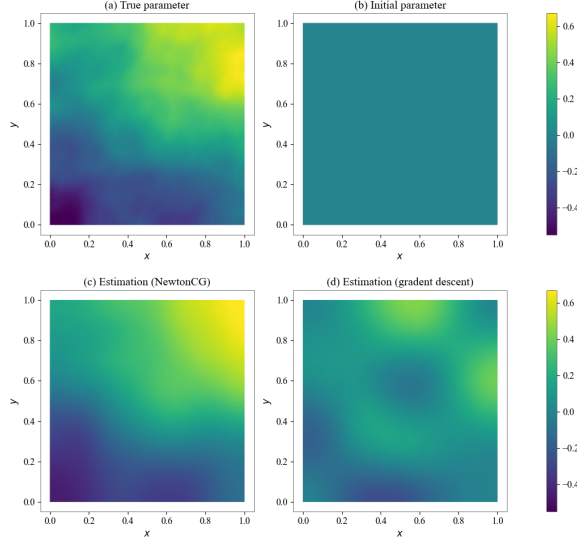
FIGURE 1. Comparison between the results obtained by the gradient descent (GD) and IMFNCG algorithms for the Darcy flow model. (a): The background truth; (b): The initial guess of the parameter; (c): The MAP estimate obtained by the IMFNCG algorithm with 10 steps; (d): The MAP estimate obtained by the GD algorithm with 1000 steps.

5.1. **Comparison of optimization methods.** We compare different optimization methods for solving the inverse problem of the Darcy flow equation. Specifically, we present the maximum a posteriori (MAP) estimate obtained by the gradient descent (GD) algorithm and an inexact matrix-free Newton-conjugate gradient (IMFNCG) algorithm. The latter is suitable for computing large-scale inverse problems. For more details about the IMFNCG algorithm, we refer to [5, 21] and references therein. The step length of GD and IMFNCG are determined by the Armijo line search, and the initial guess is set to be a zero function.

Figure 1 shows the estimates obtained by the GD and IMFNCG. On the top left, we show the background truth function $u$. On the top right, we show the initial zero function. In the second row, we show the MAP estimates obtained by the IMFNCG and GD algorithms, respectively. It can be seen that the IMFNCG algorithm with only 10 steps of iteration gives a reasonable estimate. However, the GD algorithm with Armijo line search cannot provide an accurate estimate even after 1000 iterative steps. The iSVGD sampling algorithm with no precondition is reduced to the GD algorithm when only one particle is considered. Hence, it is expected that the iSVGD sampling algorithm cannot work well since particles can hardly concentrate due to the inefficiency of the optimization procedure. Figure 2 exhibits the estimates of the variance and covariance functions calculated on mesh points by iSVGD and iSVGDMPO when the initial particles are generated by Gaussian approximation of the posterior measure [5]. The results shown in Figure 2 confirm our intuition.

In addition, these numerical results verify that it is necessary to introduce the iSVGD with preconditioning operators to enhance the optimization procedure.
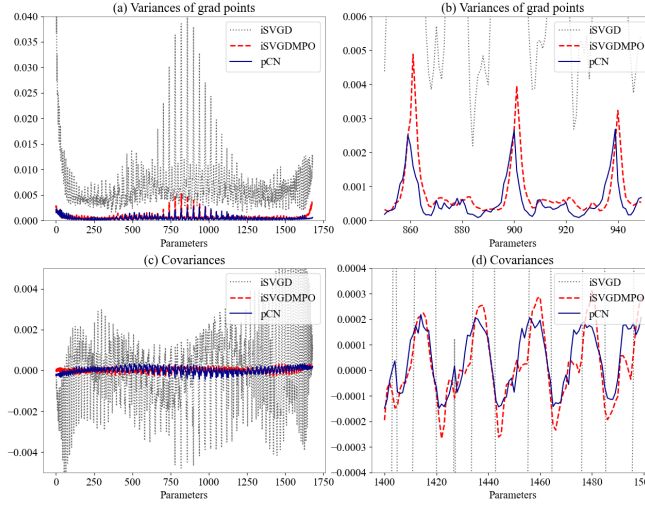
FIGURE 2. Comparison of the variances and covariance estimated by the pCN, iSVGD (1000 iterative steps) and iSVGDMPO (25 iterative steps) for the Darcy flow model. (a): Variances of grad points estimated by pCN, iSVGD and iSVGDMPO (adpative s); (b): Local enlarged draw of variances in (a); (c): Covariances of point with coordinate $(0.465, 0.035)$ with all other points on the grid estimated by pCN, iSVGD and iSVGDMPO (adpative s); (d): Local enlarged draw of covariances in (c).

Only with an efficient optimization procedure, the concentrate force (i.e., the first term in the bracket of (40)) and the repulsive force (i.e., the second term in the bracket of (40)) can sufficiently play their roles to provide accurate samplings.

5.2. **Convergence speed comparison for different values of $s$.** When choosing a kernel and the prior measure as in Section 4 of the main text, the particles should belong to the Hilbert space $\mathcal{H}^{1-s}$. From the analysis, we know that $\mu_0(\mathcal{H}^{1-s}) = 0$ or 1, when $s = 0$ or $s > 0.5$, respectively. The intuitive idea for specifying the parameter $s$ can be explained as follows:

(1) The particles should not belong to a set with zero measure, which may lead to inaccurate estimates;
(2) The particles should reside in a small support region of the prior probability measure.

Based on the above two criteria, we may choose $s$ around 0.5. Here, we provide some numerical results to answer the important question: how does $s$ influence the convergence speed of the iSVGDMPO algorithm.

Figure 3 show the detailed comparisons for the Darcy flow model. We present the estimated variances when the iterative numbers equal to 10, 20 and 30 in (a), (b) and (c) of Figure 3, respectively. In (d), (e) and (f) of Figure 3, we depict the estimated variances only for some parameters, which provide more detailed illustrations. In these figures, estimated variances for $s = 0, 0.4, 0.5$, and the adaptively chosen one are shown, which indicate the following results:
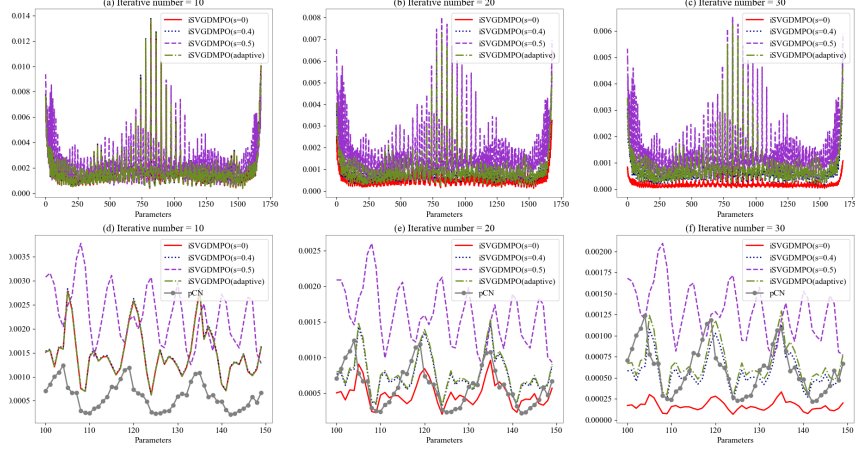
FIGURE 3. For $s = 0, 0.4, 0.5$ or adaptively chosen $s$, comparison for the estimated variances of iSVGDMPO when iterative numbers are $10, 20, 30$, respectively. (a): The estimated variances when iterative number equal to 10; (b): The estimated variances when iterative number equal to 20; (c): The estimated variances when iterative number equal to 30; (e): The estimated variances (part of the parameters) for the pCN and iSVGDMPO (iterative number equal to 10); (f): The estimated variances (part of the parameters) for the pCN and iSVGDMPO (iterative number equal to 20); (g): The estimated variances (part of the parameters) for the pCN and iSVGDMPO (iterative number equal to 30).

(1) When the iterative number is smaller than 10, the convergence speeds for $s = 0, 0.4$, and that adaptively chosen are almost the same. The convergence speed for $s = 0.5$ is obviously slower than other cases;

(2) When the iterative number approximates 30, the estimated variances for $s = 0$ is much smaller than the estimations given by the pCN and iSVGDMPO algorithm with $s = 0.4, 0.5$, and the adaptively chosen $s$.

In summary, the convergence speeds are similar for $s = 0.4$ or that adaptively chosen. The obtained estimates, at least for the variance function, are more accurate when the results of pCN are chosen as the background truth. In the main text, the comparisons for other statistical quantities are given when the parameter $s$ is specified adaptively. When specifying $s = 0.5$, the variances will gradually approach the background truth, but the convergence speed seems much slower than $s = 0.4$ or the adaptively chosen $s$.

## 6. DISCUSSIONS ON THE FINITE- AND INFINITE-DIMENSIONAL APPROACHES

Since SVGD is constructed usually for the finite-dimensional problems in the field of machine learning, it would be better for us to provide some detailed explanations about finite- and infinite-dimensional approaches, which should be useful for readers who are not familiar with the infinite-dimensional approach.

6.1. **General illustration.** The SVGD algorithm is related to optimization problems since it reduces to an optimization problem for computing maximum a posterior estimate when only one particle is considered. In the following, we firstly recall some discussions from the perspective of PDE-constrained optimization problems. For PDE-constrained optimization problems, there are two typical approaches:

- *Discretize-then-optimize*: Discretize the PDEs to formulate a finite dimensional optimization problem, then all of the optimization techniques developed on finite-dimensional space can be applied.
- *Optimize-then-discretize*: Formulate infinite-dimensional optimization problems and construct the optimization schemes on some appropriate infinite-dimensional spaces. The discretizations are pushed to the last step to generate practical numerical schemes.

*Discretize-then-optimize* and *optimize-then-discretize* are the finite- and infinite-dimensional approaches mentioned in the main context, respectively. For the advantages of the approach of *optimize-then-discretize*, we refer to page 43–44 of [15] and Chapters 2 and 3 of [12]. More specifically, the advantages of infinite-dimensional approach are mainly two-folds:

- It is important to have a better understanding of the function space structure of the numerical algorithms in order to design optimal numerical schemes for related PDEs (e.g., when forward PDEs are not self-adjoint, we may need to design certain numerical schemes to calculate forward PDEs and adjoint PDEs then to calculate the gradient).
- The approach is *mesh independent*. The mesh independence implies that the convergence behavior (e.g., convergence rate and number of iterations) of an infinite-dimensional method reflects the behavior of properly discretized problems, when the mesh size is sufficiently small.

Another method for solving inverse problems of PDEs is the Bayesian inverse methods studied in the current work. Similar to the PDE-constrained optimization methods, the Bayesian inverse methods also contain two typical approachs:

- *Discretize-then-Bayesianize*: The PDEs are initially discretized to approximate the original problem in some finite-dimensional space, and the reduced approximate problem is then solved by using the Bayes' method.
- *Bayesianize-then-discretize*: The Bayes' formula and algorithms are initially constructed on infinite-dimensional space, and after the infinite dimensional algorithm is built, some finite-dimensional approximation is carried out.

*Discretize-then-Bayesianize* and *Bayesianize-then-discretize* are the finite- and infinite-dimensional approaches mentioned in the main contexts, respectively. Similar as the optimization case, these two approaches both have their own advantages and disadvantages, and also either could be suggested to be used dependent on the specific properties of the investigated inverse problems of PDEs. By our understanding, the advantages of *Bayesianize-then-discretize* are similar as the case of *Optimize-then-discretize*:

- It is important to have a better understanding of the function space structures in order to design optimal numerical schemes of PDEs, especially when the gradient information is employed. To design sampling algorithms,

infinite-dimensional theories will be helpful to design appropriate discretization of probability measures.

- *Bayesianize-then-discretize* approach is *mesh independent*. The sampling efficiency will not highly depend on the dimension of the discretization, which is an important expected property for solving inverse problems of PDEs.

The book [14] provides a comprehensive discussions on the finite-dimensional approach, i.e., *discretize-then-Bayesianize* approach. For the infinite-dimensional approach, we refer to [20, 9, 5, 4, 8] and the references there in.

6.2. **A simple example.** In Subsection 6.1, general discussions on finite- and infinite-dimensional approaches are given, which can hardly provide some intuitions on the differences of the numerical schemes. In the following, we consider a simple example that illustrates the implementation differences between "optimize (Bayesianize)-then-discretize" and "discretize-then-optimize (Bayesianize)" approach. Let us consider the following equation:

$$\begin{cases} -0.1\Delta w + w = u, & \text{in } \Omega, \\ u = 0, & \text{on } \partial\Omega, \end{cases} \tag{6.1}$$

where $\Omega = [0,1]^2$. Denote the forward operator $\mathcal{F}(u) := w$ and the measurement operator $\mathcal{M}(w) := (w(x_1), \ldots, w(x_{N_d}))^T$ where $\{x_i\}_{i=1}^{N_d}$ reside in $\Omega$ and $N_d$ is a positive integer. Define $\mathcal{G} := \mathcal{M} \circ \mathcal{F}$. We then have the following formulation:

$$\boldsymbol{d} = \mathcal{G}(u) + \boldsymbol{\epsilon}, \tag{6.2}$$

where $\boldsymbol{d}$ is the noisy data and $\boldsymbol{\epsilon}$ is the random noise. The simplest way for estimating $u$ from $\boldsymbol{d}$ is to solve the following minimization problem:

$$\min_u F(u) \tag{6.3}$$

with $F(u) := \frac{1}{2}\|\mathcal{G}(u) - \boldsymbol{d}\|_{\ell^2}^2$. Now, we employ the finite-element method to discretize the above problem. Denote the finite element mass matrix by $\boldsymbol{M}$, the stiffness matrix of equation (6.1) by $\boldsymbol{K}$, and the measurement matrix by $\boldsymbol{S}$. The forward operator $\mathcal{G}$ then has the following discretized form:

$$\boldsymbol{d} = \boldsymbol{S}\boldsymbol{K}^{-1}\boldsymbol{M}\boldsymbol{u} + \boldsymbol{\epsilon}, \tag{6.4}$$

where $\boldsymbol{u}$ is the discretized vector of the function $u$.

**Discretize-then-optimize (Bayesianize)**: For using *discretize-then-optimize (Bayesianize)* approach, we need to formulate the following discrete problem:

$$\min_{\boldsymbol{u}} \frac{1}{2}\|\boldsymbol{S}\boldsymbol{K}^{-1}\boldsymbol{M}\boldsymbol{u} - \boldsymbol{d}\|_{\ell^2}^2. \tag{6.5}$$

Using the gradient descent method, we obtain the following iterative scheme:

$$\boldsymbol{u}_{k+1} = \boldsymbol{u}_k - \gamma(\boldsymbol{S}\boldsymbol{K}^{-1}\boldsymbol{M})^T(\boldsymbol{S}\boldsymbol{K}^{-1}\boldsymbol{M}\boldsymbol{u} - \boldsymbol{d}), \tag{6.6}$$

where $\gamma$ is the step size.

**Optimize (Bayesianize)-then-discretize**: For using *optimize (Bayesianize)-then-discretize* approach, we need to firstly formulate the infinite-dimensional problem:

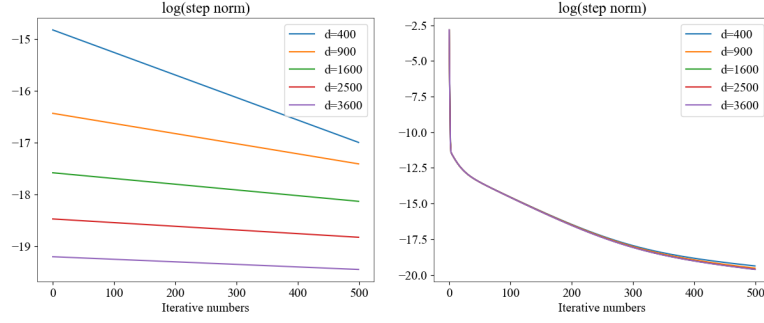$$\min_u \frac{1}{2}\|\mathcal{G}(u) - \boldsymbol{d}\|_{\ell^2}^2. \tag{6.7}$$

FIGURE 4. Left: Logarithm of the step norms computed by "discretize-then-optimize (Bayesianize)" approach with different discretized dimensions $d = 400, 900, 1600, 2500, 3600$. Right: Logarithm of the step norms computed by "optimize (Bayesianize)-then-discretize" approach with different discretized dimensions $d = 400, 900, 1600, 2500, 3600$.

Then we derive the gradient descent iteration on infinite-dimensional space to obtain:

$$u_{k+1} = u_k - \gamma \mathcal{G}^*(\mathcal{G}(u) - \boldsymbol{d}), \tag{6.8}$$

where $\mathcal{G}^*$ is the adjoint-operator of $\mathcal{G}$. According to Subsection 3.3 of [5], we may obtain the following iterative scheme on finite-dimensional space:

$$\boldsymbol{u}_{k+1} = \boldsymbol{u}_k - \gamma \boldsymbol{M}^{-1}(\boldsymbol{S}\boldsymbol{K}^{-1}\boldsymbol{M})^T(\boldsymbol{S}\boldsymbol{K}^{-1}\boldsymbol{M}\boldsymbol{u} - \boldsymbol{d}). \tag{6.9}$$

Comparing iterative schemes (6.6) and (6.9), we can see the difference. At a glance, this is a small difference. However, such a small difference leads to different behaviors of the two iterative schemes. We implement the two iterative schemes with different discretized dimensions $d = 20 \times 20, 30 \times 30, 40 \times 40, 50 \times 50, 60 \times 60$ to visually see such different behavior. The step size $\gamma$ is set to be 0.01 for all of the iterative schemes. We define the step norm as follows:

$$\text{The } k\text{-th step norm} = \|u_{k+1} - u_k\|_{L^2}. \tag{6.10}$$

In the left of Figure 4, we draw the step norms of the iterative scheme (6.6). We can see that the step norms decay rapidly when the dimension grows. This indicate that the convergence speed of iterative scheme (6.6) depends highly on the discretized dimension. In contrast, the step norms of the iterative scheme (6.9) are almost the same for different discretizations. From this simple toy example, we can see that the "discretize-then-optimize (Bayesianize)" approach can hardly keep the infinite-dimensional natural. Hence, it usually lacks *mesh independence* property. However, the "optimize (Bayesianize)-then-discretize" approach pushes the discretization implementation to the final step which makes it easier to catch the infinite-dimensional natural of the inverse problems of PDEs. The finally obtained algorithm usually has *mesh independence* property, which is important for solving inverse problems of PDEs.

6.3. **Mesh independence of iSVGD.** In Subsection 6.2, we just provide a simple example. For the proposed iSVGD, we need more standard techniques that can be found in some typical literatures [5, 11, 15, 19, 21, 6]. The lecture notes provided in "https://uvilla.github.io/inverse17/" are also beneficial for taking implementations.
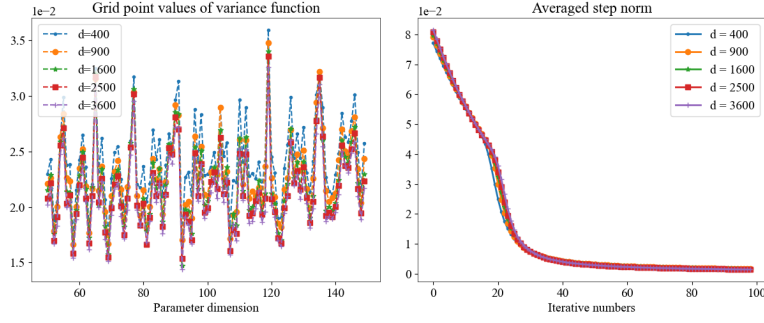
FIGURE 5. Left: Pointwise sample variances computed by different discretized dimensions $d = 400, 900, 1600, 2500, 3600$ (All of the variance functions are projected on a grid with $d = 400$ for comparison). Right: Decay of the averaged step norm $\frac{1}{m} \sum_{i=1}^{m} \|u_i^{\ell+1} - u_i^{\ell}\|_{L^2}$ w.r.t. the number of iterations for different discretized dimensions $d = 400, 900, 1600, 2500, 3600$.

Now, let us illustrate that the proposed iSVGD algorithm possesses the *mesh independence* property. That is to say, if the finite element mesh is refined, we indeed need more computational resources since the computations of each partial differential equation are more expensive. However, it might not need more iterations and particles when the finite element mesh is refined since discrete problems derived by refined mesh also approximate the infinite-dimensional formulation. For clearly illustrating this, we choose different discretized grids such that the dimensions of the function parameter are $d = 20 \times 20, 30 \times 30, 40 \times 40, 50 \times 50, 60 \times 60$. Using the same settings as in Section 4 of the main text, we only change the discretized dimension to see how discretized dimensions affect the behavior of the algorithm. In Figure 5, we show the numerical results which demonstrate the *mesh independence* as expected for *Bayesianize-then-discretize* approach.

Specifically speaking, we draw the variance functions with different discretized dimensions in the left of Figure 5. The variance functions are calculated by the iSVGDMPO with discretized dimensions $d = 400, 900, 1600, 2500, 3600$. When the algorithm generates the final particles, we calculate the variance functions and project the estimated variance functions on a mesh with dimension $d = 20 \times 20$. Then, we draw part of the grid point values of the variance functions calculated by different meshes. From the figure, it can be seen that the grid point values are similar. This validates that the estimated variance function obtained by iSVGDMPO is not sensitive to one particular discretization. Similar to other *mesh independence* methods such as rMAP used for comparison in our numerical experiments, it may be difficult to obtain exactly the same values due to the quantities being evaluated approximately, especially the gradients and Hessian operators, are not evaluated accurately. For discretize-first type methods, we can calculate the gradients of the discretized system exactly. However, the gradients and Hessian operators defined on infinite-dimensional space could only be calculated approximately.

In the right of Figure 5, we draw the averaged step norm defined as follows:

$$\frac{1}{m} \sum_{i=1}^{m} \|u_i^{\ell+1} - u_i^{\ell}\|_{L^2},\tag{6.11}$$

where $u_i^{\ell}$ stands for the $i$th particle at the $\ell$th iteration and $m$ is the number of particles. Obviously, the averaged step norms are similar for different discretized dimensions. It is evident that the curves under different discretized dimensions can hardly be distinguished, indicating that the algorithm has *mesh independence* property. The convergence speed is not affected by discretized dimensions, which is not true for many algorithms developed under the finite-dimensional setting.

At last, we should admit that more theoretical works are needed to ensure the *mesh independence* property. Specifically speaking, we may need to do further research based on Subsection 3.4 in the main text on infinite-dimensional particle interacting system and the measure-valued evolution equation. The well-posedness of these complicated equations should be proved and a theorem like Theorem 2.7 in [16] needs to be established. Along this direction, we may consult to the studies on the semilinear Mckean–Vlasov stochastic evolution equation in Hilbert space [2] and the theoretical analysis of the pCN algorithm [17].

## 7. Numerical results for the Helmholtz equation

In this section, we present numerical experiments for the Helmholtz equation

$$-\Delta w - e^{2u} w = 0 \text{ in } \Omega,$$
$$\frac{\partial w}{\partial \boldsymbol{n}} = g \text{ in } \partial \Omega,\tag{7.1}$$

where $w$ is the acoustic field, $u$ is the logarithm of the distributed wave number field on $\Omega$ ($\Omega$ is a bounded domain), $\boldsymbol{n}$ is the unit outward normal on $\partial \Omega$, and $g$ is the prescribed Neumann source on the boundary. The boundary value problem (7.1) may not have a unique solution due to possible resonances [7]. Hence, we can hardly verify Assumption 6 for this example. However, this model was studied for the randomized maximum a posteriori (rMAP) method [21], which is an approximate method used for our comparison in the main text. From the proof in Subsection 4, we may verify Assumption 6 under more suitable settings for the inverse medium scattering problem, e.g., Lemma 2.3 in [3] gives a similar estimate to the Darcy flow model.

Basic settings and the finite-element discretization are similar to the Darcy flow model considered in the main text. The only difference is that the measurement data are collected on the boundary of the domain, i.e., $x_i \in \partial \Omega$ for $i = 1, \ldots, N_d$. Figure 6 shows the estimates of the variances obtained by the pCN, rMAP, and iSVGDMPO with parameter $s = 0, 0.4$, or choosing adaptively according to formula (68) in the main text. Similarly, the estimated variances are too small when $s = 0$, which implies that the particles are concentrated on a small set. When $s$ is taken as 0.4 or chosen adaptively, we obtain similar estimates, which is more similar to the baseline provided by pCN compared with the estimates obtained by the rMAP. As in the main text, we use the empirical adaptive strategy to specify the parameter $s$ in the following.

As for the sample numbers, we also compare the estimated variances when the particle number $m$ equals to $10, 20, 30, 40$, and $50$. On the left and right in Figure
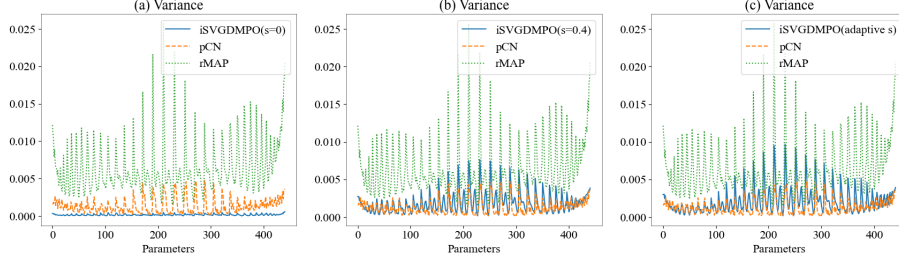
FIGURE 6. Comparison of the variances estimated by the pCN, rMAP, iSVGDMPO with different $s$ for the Helmholtz equation model. (a): Variances estimated by pCN, rMAP, and iSVGDMPO ($s = 0$); (b): Variances estimated by pCN, rMAP, and iSVGDMPO ($s = 0.4$); (c): Variances estimated by pCN, rMAP, and iSVGDMPO (adpative s).
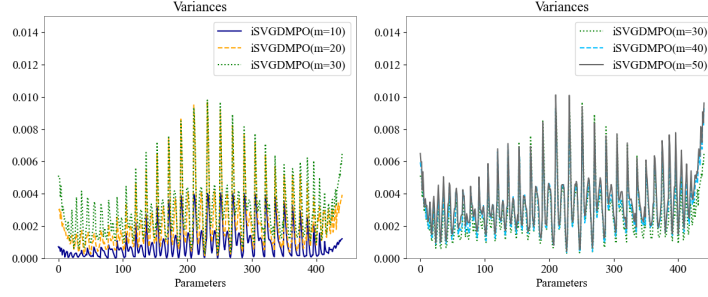


FIGURE 7. Comparison of the variances estimated by the iSVGDMPO when $s = 10, 20, 30, 40, 50$ for the Helmholtz equation model.

7, we show the results obtained when $m = 10, 20, 30$ and $m = 30, 40, 50$, respectively. Obviously, we find that $m = 10$ is not enough to give reliable estimates and the estimated variance functions are similar when $m = 30, 40, 50$. Hence, for the Helmholtz problem, it is enough to take $m = 20$ or $30$ for our numerical examples, which attains a fine balance between efficiency and accuracy.

For the following numerical experiments, we take $m = 30$ and set the parameter $s$ by the empirical strategy (68) as presented in the main text. In Figure 8, we demonstrate the background truth and the estimated mean and variance functions obtained by the pCN, rMAP, and iSVGDMPO, respectively. The iterative number of iSVGDMPO is set to be 30. The same observation can be made from the results. The mean functions obtained by the rMAP and iSVGDMPO are similar, which are slightly smoother than the one obtained by the pCN algorithm. Regarding the variance function, it can be seen from (f), (g), and (h) of the figure that the iSVGDMPO gives more reliable estimates than the rMAP does.

Now, we provide some more comparisons of statistical quantities among the results obtained by the pCN, rMAP, and iSVGDMPO. Similarly, we compute variance and covariance functions on the mesh points and exhibit the results in Figure 9. In all of the subfigures in Figure 9, the estimates obtained by the pCN, rMAP, and iSVGDMPO are drawn in blue solid line, gray dotted line, and red dashed line, respectively. All the notations here are the same as those used in the main text.
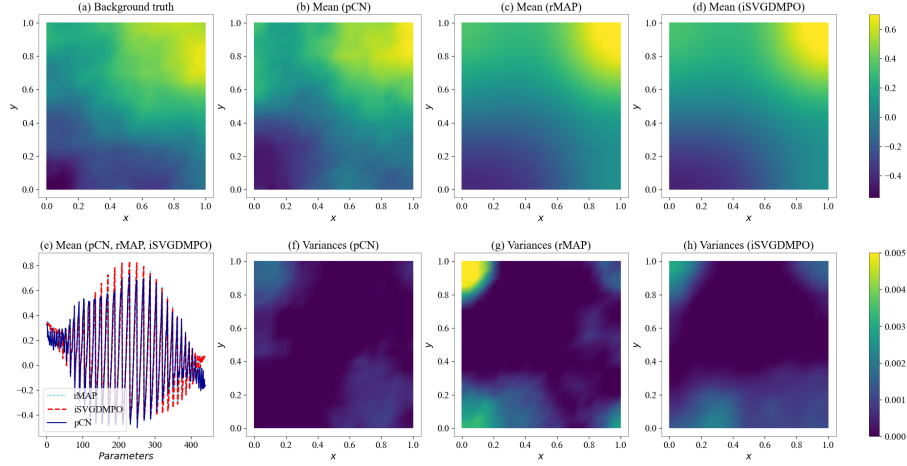
FIGURE 8. The background truth and estimated mean and variance functions by pCN, rMAP, and iSVGDMPO for the Helmholtz equation model. (a): Background truth; (b): Estimated mean function by pCN; (c): Estimated mean function by rMAP; (d): estimated mean function by iSVGDMPO; (e): Estimated mean function on mesh points by pCN (blue solid line), rMAP (light blue dotted line), and iSVGDMPO (red dashed line); (f): Estimated variances by pCN; (g): Estimated variances by rMAP; (h): Estimated variances by iSVGDMPO.

We can also obtain the same conclusions from the results: the estimates obtained by the iSVGDMPO are visually more similar to the estimates provided by the pCN compared with the results obtained by the rMAP.

TABLE 1. The $\ell^2$-norm error of variance and covariance functions on mesh points for the rMAP and iSVGDMPO (estimates of the pCN are seen as the background truth)

|  | $\text{var}_u(x_i)$ | $\text{cov}_u(x_i, x_{i+5})$ | $\text{cov}_u(x_i, x_{i+10})$ |
|---|---|---|---|
| rMAP | 0.01525 | 0.00155 | 0.00237 |
| iSVGDMPO | 0.00092 | 0.00026 | 0.00063 |
|  | $\text{cov}_u(x_i, x_{i+15})$ | $\text{cov}_u(x_i, x_{i+20})$ | $\text{cov}_u(x_i, x_{i+25})$ |
| rMAP | 0.00295 | 0.00353 | 0.00153 |
| iSVGDMPO | 0.00036 | 0.00059 | 0.00035 |
|  | $\text{cov}_u(x_i, x_{i+30})$ | $\text{cov}_u(x_i, x_{i+35})$ | $\text{cov}_u(x_i, x_{i+40})$ |
| rMAP | 0.00148 | 0.00154 | 0.00219 |
| iSVGDMPO | 0.00055 | 0.00037 | 0.00048 |

Besides these visual comparisons, a quantitative comparison of the differences among the pCN, rMAP, and iSVGDMPO are also given in Table 1. Again, all the notations have the same meaning as those used in the main text. It can be seen
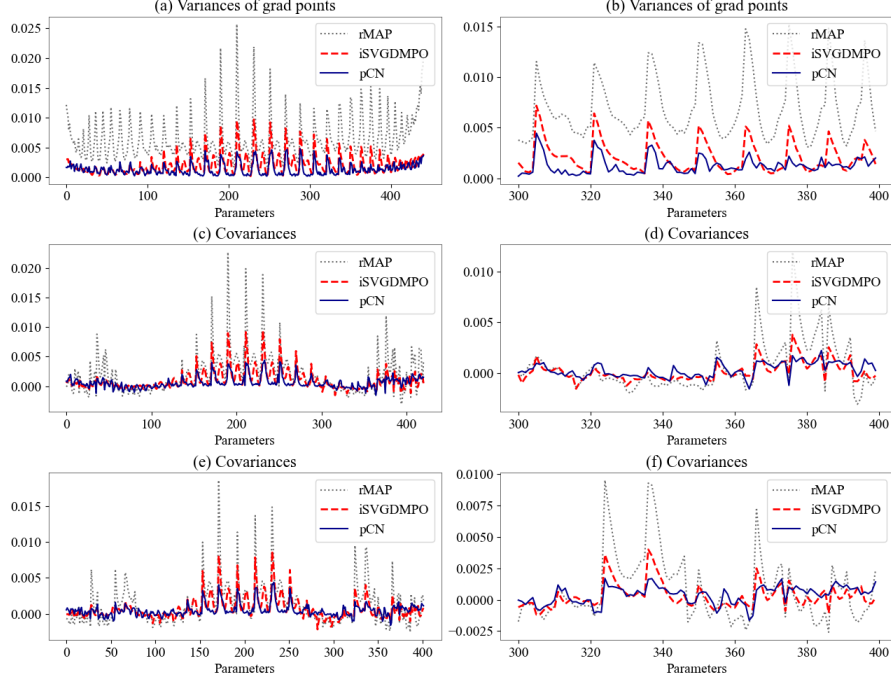
FIGURE 9. The estimated variances and covariances by the pCN (blue solid line), rMAP (gray dotted line), and iSVGDMPO (red dashed line). (a): Estimated variances $\{\text{var}_u(x_i)\}_{i=1}^{N_g}$ on all mesh points; (b): Estimated variances for mesh points with indexes from 300 to 400 (show details); (c): Estimated covariances $\{\text{cov}_u(x_i, x_{i+20})\}_{i=1}^{N_g-20}$ on mesh point pairs $\{(x_i, x_{i+20})\}_{i=1}^{N_g-20}$; (d): Estimated covariances shown in (c) with indexes from 300 to 400 (show details); (e): Estimated covariances $\{\text{cov}_u(x_i, x_{i+40})\}_{i=1}^{N_g-40}$ on mesh point pairs $\{(x_i, x_{i+40})\}_{i=1}^{N_g-40}$; (f): Estimated covariances shown in (e) with indexes from 300 to 400 (show details).

from Table 1 that all the $\ell^2$-norm differences of the iSVGDMPO with the pCN are evidently smaller than the corresponding values of the rMAP.

## REFERENCES

[1] S. Agapiou, S. Larsson, and A. M. Stuart. Posterior contraction rates for the Bayesian approach to linear ill-posed inverse problems. *Stoch. Proc. Appl.*, 123(10):3828–3860, 2013.

[2] N. U. Ahmed and X. Ding. A semilinear Mckean-Vlasov stochastic evolution equation in Hilbert space. *Stoch. Proc. Appl.*, 60:65–85, 1995.

[3] G. Bao and P. Li. Inverse medium scattering for the Helmholtz equation at fixed frequency. *Inverse Probl.*, 21(5):1621–1641, 2005.

[4] A. Beskos, A. Jasra, E. A. Muzaffer, and A. M. Stuart. Sequential Monte Carlo methods for Bayesian elliptic inverse problems. *Stat. Comput.*, 25:727–737, 2015.

[5] T. Bui-Thanh, O. Ghattas, J. Martin, and G. Stadler. A computational framework for infinite-dimensional Bayesian inverse problems part I: The linearized case, with application to global seismic inversion. *SIAM J. Sci. Comput.*, 35(6):A2494–A2523, 2013.

[6] T. But-Thanh and Q. P. Nguyen. FEM-based discretization-invariant MCMC methods for PDE-constrained Bayesian inverse problems. *Inverse Probl. Imag.*, 10(4):943–975, 2016.

[7] D. Colton and R. Kress. *Inverse Acoustic and Electromagnetic Scattering Theory*. Springer, Cham, fourth edition, 2019.

[8] S. L. Cotter, G. O. Roberts, A. M. Stuart, and D. White. MCMC methods for functions: modifying old algorithms to make them faster. *Stat. Sci.*, 28(3):424–446, 2013.

[9] Simon L Cotter, Massoumeh Dashti, James Cooper Robinson, and Andrew M Stuart. Bayesian inverse problems for functions and applications to fluid mechanics. *Inverse Probl.*, 25(11):115008, 2009.

[10] M. Dashti and A. M. Stuart. The Bayesian approach to inverse problems. *Handbook of Uncertainty Quantification*, pages 311–428, 2017.

[11] O. Ghattas and K. Willcox. Learning physics-based models from data: perspectives from inverse problems and model reduction. *Acta Numer.*, 30:445–554, 2021.

[12] M. Hinze, R. Pinnau, M. Ulbrich, and S. Ulbrich. *Optimization with PDE Constraints*. Springer Netherlands, 2009.

[13] J. Jia, J. Peng, and J. Gao. Posterior contraction for empirical Bayesian approach to inverse problems under non-diagonal assumption. *Inverse Probl. Imag.*, 15(2):201–228, 2020.

[14] J. Kaipio and E. Somersalo. *Statistical and Computational Inverse Problems*. Springer-Verlag, New York, 2005.

[15] Juan Carlos De los Reyes. *Numerical PDE-Constrained Optimization*. Springer, New York, 2015.

[16] J. Lu, Y. Lu, and J. Nolen. Scaling limit of the Stein variational gradient descent: the mean field regime. *SIAM J. Math. Anal.*, 5(2):648–671, 2019.

[17] N. S. Pillai, A. M. Stuart, and A. H. Thiery. Noisy gradient flow from a random walk in Hilbert space. *Stoch. Partial. Differ.*, 2(2):196–232, 2014.

[18] G. D. Prato. *Kolmogorov Equations for Stochastic PDEs*. Birkhäuser Verlag, Basel, 2004.

[19] A. Spantini, A. Solonen, T. Cui, J. Martin, L. Tenorio, and Y. Marzouk. Optimal low-rank approximations of Bayesian linear inverse problems. *SIAM J. Sci. Comput.*, 37(6):A2451–A2487, 2015.

[20] A. M. Stuart. Inverse problems: A Bayesian perspective. *Acta Numer.*, 19:451–559, 2010.

[21] K. Wang, T. Bui-Thanh, and O. Ghattas. A randomized maximum a posteriori method for posterior sampling of high dimensional nonlinear Bayesian inverse problems. *SIAM J. Sci. Comput.*, 40(1):A142–A171, 2018.

School of Mathematics and Statistics, Xi'an Jiaotong University, Xi'an, 710049, China

*Email address*: jjx323@xjtu.edu.cn

Department of Mathematics, Purdue University, West Lafayette, Indiana, 47907, USA

*Email address*: lipeijun@math.purdue.edu

School of Mathematics and Statistics, Xi'an Jiaotong University, Xi'an, 710049, China

*Email address*: dymeng@mail.xjtu.edu.cn