

# CONVERGENCE OF SVGD IN KL DIVERGENCE VIA APPROXIMATE GRADIENT FLOW

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This study investigates the convergence of Stein variational gradient descent (SVGD), which is used to approximate a target distribution based on a gradient flow on the space of probability distributions. The existing studies mainly focus on the convergence in the kernel Stein discrepancy, which doesn't imply weak convergence in many practical settings. To address this issue, we propose to introduce a novel analytical approach called  $(\epsilon, \delta)$ -approximate gradient flow, extending conventional concepts of approximation error for the Wasserstein gradient. With this approach, we show the sub-linear convergence of SVGD in Kullback–Leibler divergence under the discrete-time and infinite particle settings. Finally, we validate our theoretical findings through several numerical experiments.

## 1 INTRODUCTION

Sampling from an unnormalized target distribution, such as posterior distribution in Bayesian inference, is a fundamental problem in machine learning. The mainstream approaches for obtaining such samples is using Markov Chain Monte Carlo (MCMC) methods (Hastings, 1970; Welling & Teh, 2011) or approximating the target distribution by variational inference (VI) (Jordan et al., 1999; Blei et al., 2017). While MCMC provides guarantees of producing asymptotically unbiased samples from the target density, it tends to be computationally intensive (Robert & Casella, 2004). On the other hand, VI achieves a computationally efficient approximation of the target distribution through stochastic optimization under a simpler alternative distribution; however, it does not come with a guarantee of obtaining unbiased samples (Blei et al., 2017).

To alleviate such sample bias while maintaining computational efficiency of VI as much as possible, Liu & Wang (2016) introduced *Stein variational gradient descent* (SVGD), which allows the direct approximation of the target distribution without the need for alternative distributions. SVGD iteratively updates correlated samples, referred to as *particles*, by minimizing the Kullback–Leibler (KL) divergence between a distribution of particles and the target distribution through a gradient flow on the space of probability distributions. Since the Wasserstein gradient is intractable in practice, SVGD approximates it through a kernel method.

On the theoretical front, analysis has been actively conducted ever since Liu (2017) elucidated the asymptotic behavior of SVGD from the perspective of gradient flow within the reproducing kernel Hilbert space (RKHS). Korba et al. (2020) showed sub-linear convergence in kernel Stein discrepancy (KSD) under infinite particles assuming that KSD at each step is bounded. Salim et al. (2022) contributed a proof of sub-linear convergence in KSD without the necessity of bounded KSD assuming that the target distribution satisfies  $T_1$  inequality (Villani, 2008), and Sun et al. (2023) provided the proofs of this convergence property by relaxing the smoothness assumption of the target distribution. A common thread in these analyses is seeing SVGD's update rule as the approximation of the Wasserstein gradient in the RKHS and showing that the KL divergence to target distribution monotonically decreases like gradient descent. Beyond the infinite particle setting, Shi & Mackey (2023) has recently shown that the SVGD with  $n$  finite particles and an appropriate step size converges in KSD at the  $\mathcal{O}(1/\sqrt{\log \log n})$  order if the target distribution is sub-Gaussian with a Lipschitz score.

However, the convergence analysis in terms of KSD is insufficient to understand the weak convergence property of SVGD because the convergence in KSD holds under highly restrictive conditions for the kernel and the target distribution under practical settings as shown by Gorham & Mackey (2017). This fact underscores the importance of conducting convergence analysis using criteria

other than KSD to provide more realistic guarantees for the obtained particles. A natural candidate for the criterion is the KL divergence itself, which is the objective function of SVGD. Recently, Liu et al. (2023) showed that SVGD with finite particles achieves linear convergence in KL divergence under a very limited setting where the target distribution is Gaussian. However, the analytical approach presented in previous studies makes it difficult to conduct convergence analysis based on KL divergence in a more global setting. The reason for this lies in the fact that while the logarithmic Sobolev inequality (LSI) (Gross, 1975) is typically employed to show the linear convergence in KL divergence for a gradient flow in the space of probability distributions (Villani, 2008), it becomes apparent that the inequality similar to the LSI (see Eq. (7)) does not hold in practical settings (Duncan et al., 2023) when considering SVGD as a gradient flow in the RKHS.

In this study, we introduce a novel analytical approach that allows us to circumvent the aforementioned issue. A key idea in our analysis is to consider SVGD as an *approximation* of the gradient flow in the space of probability distributions, as opposed to the conventional analytical approach that views SVGD as a gradient flow in the RKHS. To express the degree of this approximation, we introduce a new concept called  $(\epsilon, \delta)$ -approximate gradient flow, which extends the concept of approximation error widely used in the gradient estimation context such as score gradient estimation (Lee et al., 2022; 2023) and particle-based VI (Liu et al., 2019; Dong et al., 2022).

With our concept, we offer new insights into the convergence of SVGD in the settings of discrete-time and an infinite number of particles. We first analyze the degree of the approximation error  $\{\epsilon, \delta\}$  between the Wasserstein gradient of the KL divergence and the update rule in SVGD by focusing on spectral decomposition specified via a kernel function. With this approximation error analysis, we show that SVGD exhibits sub-linear convergence in the KL divergence for the first time, to the best of our knowledge. At last, we conduct a numerical study to examine the convergence behavior of SVGD across various metrics and validate the soundness of our theoretical findings.

## 2 PRELIMINARIES

Random variables are denoted by capital letters like  $X$ , while deterministic values are denoted by lowercase letters like  $x$ . The Euclidean inner product and distance are expressed as  $\langle \cdot, \cdot \rangle$  and  $\| \cdot \|$ , respectively. Let  $\mathcal{X} = \mathbb{R}^d$  and let  $C^l(\mathcal{X}, \mathcal{Y})$  be the space of  $l$  continuously differentiable functions from  $\mathcal{X}$  to a Hilbert space  $\mathcal{Y}$ . We abbreviate  $C^l(\mathcal{X}, \mathbb{R})$  as  $C^l(\mathcal{X})$ . The set of smooth functions with compact support is expressed as  $C_c^\infty(\mathcal{X})$ . If  $\phi \in C^1(\mathcal{X})$ , its gradient is  $\nabla \phi$ . For  $\phi \in C^1(\mathcal{X}, \mathcal{X})$ , the Jacobian is represented as  $J\phi(x)$ , a  $d \times d$  matrix at each point  $x \in \mathcal{X}$ . We define  $\text{div}\phi(x) = \text{Tr}J\phi(x)$ . The Hilbert–Schmidt and operator norm of a matrix are denoted as  $\| \cdot \|_{\text{HS}}$  and  $\| \cdot \|_{\text{op}}$ .

### 2.1 WASSERSTEIN SPACE AND CONTINUITY EQUATION

Here we summarize some of the basics of optimal transport that underlie our analysis. We denote the set of probability measures on  $\mathcal{X}$  with finite second moments as  $\mathcal{P}_2(\mathcal{X})$ . For any  $\mu \in \mathcal{P}_2(\mathcal{X})$ , we express the set of measurable functions  $f : \mathcal{X} \rightarrow \mathcal{X}$  with  $\int \|f\|^2 d\mu < \infty$  as  $L^2(\mu)$ , with its norm and inner product as  $\| \cdot \|_{L^2(\mu)}$  and  $\langle \cdot, \cdot \rangle_{L^2(\mu)}$ . Given a measurable map  $T : \mathcal{X} \rightarrow \mathcal{X}$  and  $\mu$ , we denote the pushforward measure of  $\mu$  by  $T$  as  $T\#\mu \in \mathcal{P}_2(\mathcal{X})$ , which is characterized by  $\int \phi(T(x)) d\mu(x) = \int \phi(y) dT\#\mu(y)$  for any measurable and bounded function  $\phi$ . Given  $\mu, \nu \in \mathcal{P}_2(\mathcal{X})$ , the Wasserstein distance between  $\mu$  and  $\nu$  is defined as  $W_2^2(\mu, \nu) = \inf_{s \in \mathcal{S}(\mu, \nu)} \int \|x - y\|^2 ds(x, y)$ , where  $\mathcal{S}(\mu, \nu)$  is the set of couplings between  $\mu$  and  $\nu$ . This distance defines a metric on  $\mathcal{P}_2(\mathcal{X})$ , making  $(\mathcal{P}_2(\mathcal{X}), W_2)$  the Wasserstein space, which is complete and separable.

Now we introduce a *continuous equation*. Let  $T > 0$  and consider a weakly continuous map  $\mu : (0, T) \rightarrow \mathcal{P}_2(\mathcal{X})$ ,  $t \mapsto \mu_t$ . The family  $(\mu_t)_{t \in (0, T)}$  satisfies a continuity equation if there exists  $(v_t)_{t \in (0, T)}$  such that  $v_t \in L^2(\mu_t)$  and  $\frac{\partial \mu_t}{\partial t} + \text{div}(\mu_t v_t) = 0$  holds in the distribution sense (see Appendix B.1 for the formal meaning of *distribution sense*). A family  $(\mu_t)_{t \in (0, T)}$  that satisfies a continuity equation with integrable  $\|v_t\|_{L^2(\mu_t)}$  over  $(0, T)$  is referred to as *absolutely continuous*. Conversely, one can construct an absolutely continuous  $(\mu_t)_{t \in (0, T)}$  by selecting  $(v_t)_{t \in (0, T)}$  such that they meet the above condition.

While the Wasserstein space does not inherently possess the characteristics of a Riemannian manifold, it can be endowed with a Riemannian structure and interpretation (Otto, 2001). In this in-

terpretation, the tangent space of  $\mathcal{P}_2(\mathcal{X})$  at  $\mu_t$ , denoted as  $\mathcal{T}_{\mu_t}\mathcal{P}_2(\mathcal{X})$ , forms a subset of  $L^2(\mu_t)$ . When considering all possible  $(v_t)_{t \in (0, T)}$ , we call  $v_t$  that exhibits the minimal  $L^2(\mu_t)$  norm as the velocity field of  $(\mu_t)_{t \in (0, T)}$ . This minimality condition can be characterized by the requirement that  $v_t \in \mathcal{T}_{\mu_t}\mathcal{P}_2(\mathcal{X}) (\subset L^2(\mu_t))$ .

## 2.2 SAMPLING-BASED APPROXIMATION VIA GRADIENT FLOW OF KL DIVERGENCE

We aim to obtain samples from the density  $\pi(x) \propto e^{-V(x)}$  in  $\mathcal{P}_2(\mathcal{X})$  under the following assumption for the potential function  $V : \mathcal{X} \rightarrow \mathbb{R}$ .

**Assumption 1.** *The Hessian of  $V \in C^2(\mathcal{X})$ ,  $H_V$ , satisfies  $\|H_V\|_{\text{op}} \leq L$ .*

This task can be formulated as the optimization problem over a functional space, i.e., minimizing a functional, KL divergence of  $\mu$  from  $\pi$  defined on Wasserstein space, that is,

$$\min_{\mu \in \mathcal{P}_2(\mathcal{X})} \text{KL}(\mu|\pi), \quad \text{KL}(\mu|\pi) := \int \log \frac{d\mu}{d\pi}(x) d\mu(x), \quad (1)$$

where  $\text{KL}(\cdot|\pi) : \mathcal{P}_2(\mathcal{X}) \rightarrow [0, +\infty)$ ,  $\mu \mapsto \text{KL}(\mu|\pi)$  and  $\mu$  is absolutely continuous with respect to (w.r.t.)  $\pi$ . Thus, Radon–Nikodym<sup>1</sup> derivative  $d\mu/d\pi$  is available ( $\text{KL}(\mu|\pi) = +\infty$  otherwise).

As a method for solving Eq. (1), a gradient-descent-like algorithm utilizing the differential structure of the Wasserstein space and continuous equations (see Section 2.1) is often employed. Let the Wasserstein gradient of  $\text{KL}(\mu|\pi)$  at  $\mu$  be  $\nabla_{W_2}\text{KL}(\mu|\pi)$  (the formal definition is presented in Appendix B.1). We then consider how  $\text{KL}(\mu|\pi)$  evolves by the continuity equation, i.e.,

$$\frac{d}{dt} \text{KL}(\mu_t|\pi) = \langle \nabla_{W_2}\text{KL}(\mu_t|\pi), v_t \rangle_{L^2(\mu_t)}, \quad (2)$$

which shows that  $\text{KL}(\mu|\pi)$  is minimized by choosing  $v_t$  such that  $\langle \nabla_{W_2}\text{KL}(\mu_t|\pi), v_t \rangle_{L^2(\mu_t)} \leq 0$  and using the continuity equation. A natural choice is to use the Wasserstein gradient itself as  $v_t = -\nabla_{W_2}\text{KL}(\mu_t|\pi)$ , which results in  $\frac{d}{dt} \text{KL}(\mu_t|\pi) = -\|\nabla_{W_2}\text{KL}(\mu_t|\pi)\|_{L^2(\mu_t)}^2 \leq 0$ . According to the fact that the Wasserstein gradient of KL divergence is obtained as  $\nabla_{W_2}\text{KL}(\mu|\pi) = \nabla \log \frac{\mu}{\pi} \in L^2(\mu)$  (Ambrosio et al., 2005), we have

$$\frac{d}{dt} \text{KL}(\mu_t|\pi) = -\left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2. \quad (3)$$

Many existing studies analyzed Eq. (3) under the following assumption (Bakry et al., 2013).

**Assumption 2.** *We say that the target distribution  $\pi$  satisfies the LSI, if, for any  $\mu \in \mathcal{P}_2(\mathcal{X})$ , there exists a positive constant  $C_{\text{LS}}$  such that*

$$\text{KL}(\mu|\pi) \leq \frac{1}{C_{\text{LS}}} \left\| \nabla \log \frac{\mu}{\pi} \right\|_{L^2(\mu)}^2. \quad (4)$$

With the above inequality and Eq. (3), we have  $\text{KL}(\mu_t|\pi) \leq e^{-C_{\text{LS}}t} \text{KL}(\mu_0|\pi)$ , which implies linear convergence. However, it is difficult to deal with the continuous-time equation of Eq. (3), and thus discretization such as a forward Euler discretization (Ambrosio et al., 2005) is often used. This recursion is given by

$$\mu_{t+1} = \left( I - \gamma_t \nabla \log \frac{\mu_t}{\pi} \right) \# \mu_t, \quad (5)$$

at each iteration  $t$ <sup>2</sup>, where  $\gamma_t > 0$  is a stepsize and  $I$  is the identity map.

<sup>1</sup>Suppose that  $\mu$  is absolutely continuous w.r.t.  $\pi$ , i.e.,  $\mu \ll \pi$ . Then, there exists a function  $f$  such that, for any measurable set  $A$ ,  $\mu(A) = \int_A f(x) d\pi(x)$ . This function  $f$  is referred to as the Radon–Nikodym derivative of  $\mu$  w.r.t.  $\pi$ , denoted by  $f = d\mu/d\pi$  (Durrett, 2019).

<sup>2</sup>For the sake of readability, we adopt  $t$  to express both continuous and discrete time.

### 2.3 STEIN VARIATIONAL GRADIENT DESCENT

Performing optimization based on Eq. (5) is still difficult because  $\mu$  is often intractable and thus  $\nabla \log \frac{\mu}{\pi}$  is hard to compute. SVGD is one of the alternative gradient flow approaches to avoid this issue by projecting  $\nabla \log \frac{\mu}{\pi}$  into the reproducing kernel Hilbert space (RKHS) by a kernel function.

Here, we briefly summarize the fundamental operations on the RKHS. Let  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a symmetric and positive semi-definite kernel and  $\mathcal{H}_0$  be its corresponding RKHS of real-valued functions  $\mathcal{X} \rightarrow \mathbb{R}$ . The inner product within  $\mathcal{H}_0$  is denoted as  $\langle \cdot, \cdot \rangle_{\mathcal{H}_0}$ , which satisfies  $f(x) = \langle f, k(\cdot, x) \rangle_{\mathcal{H}_0}$  ( $\forall f \in \mathcal{H}_0$ ) by the reproducing property of  $\mathcal{H}_0$ . We also define  $\mathcal{H}$  as the Cartesian product of  $\mathcal{H}_0$ , whose elements are expressed as  $f = (f_1, \dots, f_d)$  where  $f_i \in \mathcal{H}_0$  for  $i = 1, \dots, d$ . The inner product of  $f, g \in \mathcal{H}$  is given by  $\langle f, g \rangle_{\mathcal{H}} = \sum_{i=1}^d \langle f_i, g_i \rangle_{\mathcal{H}_0}$ . If  $\mu \in \mathcal{P}_2(\mathcal{X})$  and  $\int k(x, x) d\mu(x) < \infty$ , the integral operator associated to  $k$  and  $\mu$  can be defined as  $S_{\mu, k} f(x) := \int k(y, x) f(y) d\mu(y)$ , where  $S_{\mu, k} : L^2(\mu) \rightarrow \mathcal{H}$  and thus  $\mathcal{H} \subset L^2(\mu)$ <sup>3</sup>. We further define the inclusion map as  $\iota : \mathcal{H} \rightarrow L^2(\mu)$ , which is the adjoint of  $S_{\mu, k}$ . Under the map  $\iota$ , for  $f \in L^2(\mu)$  and  $g \in \mathcal{H}$ , we have  $\langle f, \iota g \rangle_{L^2(\mu)} = \langle \iota^* f, g \rangle_{\mathcal{H}} = \langle S_{\mu, k} f, g \rangle_{\mathcal{H}}$ , where  $\iota^*$  is the adjoint of  $\iota$ . We finally define the mapping function  $P_{\mu, k} : L^2(\mu) \rightarrow L^2(\mu)$ , where  $P_{\mu, k} = \iota S_{\mu, k}$ .

In SVGD, instead of using the Wasserstein gradient  $\nabla \log \frac{\mu}{\pi}$ , we employ  $-P_{\mu, k} \nabla \log \frac{\mu}{\pi}$  as  $v_t$  in Eq. (2), leading to the following discretized dynamics:

$$\mu_{t+1} = \left( I - \gamma_t P_{\mu, k} \nabla \log \frac{\mu_t}{\pi} \right) \# \mu_t. \quad (6)$$

The difference from Eq. (5) is that  $\nabla \log \frac{\mu}{\pi}$  is mapped by  $P_{\mu, k}$ . If a kernel function satisfies  $\lim_{\|x\| \rightarrow \infty} k(x, \cdot) \pi(x) = 0$ , by using an integration by parts (Liu, 2017), we can obtain  $P_{\mu, k} \nabla \log \frac{\mu}{\pi}(x) := -\int [\nabla \log \pi(y) k(y, x) + \nabla_y k(y, x)] d\mu(y)$ . By focusing on the continuous dynamics of the KL divergence, we have

$$\frac{d}{dt} \text{KL}(\mu_t | \pi) = - \left\langle \nabla \log \frac{\mu_t}{\pi}, P_{\mu_t, k} \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)} = - \left\| S_{\mu_t, k} \nabla \log \frac{\mu_t}{\pi} \right\|_{\mathcal{H}}^2 =: -I_{\text{stein}}(\mu_t | \pi),$$

where  $I_{\text{stein}}(\mu_t | \pi)$  is called as the Stein–Fisher (SF) information (Duncan et al., 2023). It is known that the square root of the SF information corresponds to the KSD. Now it is tempting to consider whether the inequality similar to LSI in Eq. (4) holds for the SF information presented below:

$$\text{KL}(\mu | \pi) \leq c I_{\text{stein}}(\mu | \pi), \quad (7)$$

where  $c$  is some positive constant. If this inequality holds, the linear convergence of SVGD holds. Unfortunately, the conditions for the validity of this inequality are not as evident as in the case of LSI and Duncan et al. (2023) has shown that Eq. (7) may not hold in many practical models with kernel functions like the RBF kernel, where the tail of  $\pi$  is exponential. Hence, showing the linear convergence of KL divergence in the geometry of  $\mathcal{H}$  is not straightforward. We refer to Liu (2017) and Duncan et al. (2023) for a detailed discussion of the geometry of SVGD.

Recently, Salim et al. (2022) showed the descent lemma,  $\text{KL}(\mu_{t+1} | \pi) \leq \text{KL}(\mu_t | \pi) - c \gamma I_{\text{stein}}(\mu_t | \pi)$  holds where  $c$  is some positive constant that depends on the problem. Although we can obtain the convergence in KSD from this inequality, the convergence KSD not necessarily means the weak convergence as discussed in Gorham & Mackey (2017).

## 3 APPROXIMATE GRADIENT FLOW

Here, we introduce a new concept of approximation for the Wasserstein gradient,  $(\epsilon, \delta)$ -approximate gradient flows (AGF). We then analyze the convergence of the KL divergence under our concept.

### 3.1 $(\epsilon, \delta)$ -APPROXIMATE GRADIENT FLOW

Let us assume that a gradient flow on the Wasserstein space exists, which is induced by some velocity  $v_t = g_{\mu_t}(x) \in L^2(\mu_t)$  for  $x \in \mathcal{X}$ . Here,  $g_{\mu_t}(x)$  represents a function of  $x$  only depending on  $\mu_t$ .

<sup>3</sup>We introduce  $S_{\mu, k}$  for vector inputs  $f = (f_1, \dots, f_d)$ . When  $f$  is a scalar ( $d = 1$ ), for simplicity, we consider  $S_{\mu, k}$  to be defined as applied to a single element, i.e.,  $S_{\mu, k} : L_0^2(\mu) \rightarrow \mathcal{H}_0$ , allowing us to abuse the notation, where  $L_0^2(\mu)$  is the set of a measurable function  $f_1 : \mathcal{X} \rightarrow \mathbb{R}$  with  $\int f_1^2 d\mu < \infty$ .

In the continuous-time setting, such a gradient flow is obtained via the continuity equation given as  $\frac{\partial \mu_t}{\partial t} + \text{div}(\mu_t g_{\mu_t}(x)) = 0$ . Under mild growth and regularity assumptions on  $g_{\mu_t}(x)$  (Ambrosio et al., 2005; Bonnet & Frankowska, 2021), the existence and uniqueness of a gradient flow by  $g_{\mu_t}$  is guaranteed. When considering discrete time, we assume that the recursion  $\mu_{t+1} = (I - \gamma_t g_{\mu_t}) \# \mu_t$  exists, which is similar to Eq. (6).

In the presence of these, we consider the time evolution of  $\text{KL}(\mu_t|\pi)$  under the velocity  $v_t = g_{\mu_t}(x)$  as in Section 2.2. In the continuous-time setting, we assume that  $\frac{d}{dt}\text{KL}(\mu_t|\pi) = \langle \nabla \log \frac{\mu_t}{\pi}, g_{\mu_t} \rangle_{L^2(\mu_t)}$ . As for the discrete-time setting, we assume the following inequality with a kind of descent property:

$$\text{KL}(\mu_{t+1}|\pi) \leq \text{KL}(\mu_t|\pi) - \eta_t \left\langle \nabla \log \frac{\mu_t}{\pi}, g_{\mu_t} \right\rangle_{L^2(\mu_t)}, \quad (8)$$

where  $\eta_t$  is some positive constant. Such a descent property holds both in the Wasserstein gradient flow (Ambrosio et al., 2005) and in SVGD as shown in Section 2.3.

From the above two (in)equalities, we can anticipate that when  $g_{\mu_t}(x)$  exhibits behavior close to that of  $\nabla \log \frac{\mu_t}{\pi}(x)$ , i.e.,  $\langle \nabla \log \frac{\mu_t}{\pi}, g_{\mu_t} \rangle_{L^2(\mu_t)} \geq 0$  is satisfied (recall the cosine similarity in the finite-dimensional case), the KL divergence does not increase with  $t$ . In SVGD, for example, we set  $g_{\mu_t} = P_{\mu_t, k} \nabla \log \frac{\mu_t}{\pi}$ , which satisfies  $\langle \nabla \log \frac{\mu_t}{\pi}, g_{\mu_t} \rangle_{L^2(\mu_t)} = I_{\text{stein}}(\mu_t|\pi) \geq 0$ .

However, the condition  $\langle \nabla \log \frac{\mu_t}{\pi}, g_{\mu_t} \rangle_{L^2(\mu_t)} \geq 0$  is insufficient for explicitly analyzing the convergence rate since it doesn't convey how accurate the approximation via  $g_{\mu_t}$  is. To overcome this situation, we introduce a new concept of the similarity between  $\nabla \log \frac{\mu_t}{\pi}(x)$  and  $g_{\mu_t}(x)$  as follows.

**Definition 1.** Suppose that  $\nabla \log \frac{\mu_t}{\pi}(x) < \infty$  (a.e.) and  $\|\nabla \log \frac{\mu_t}{\pi}\|_{L^2(\mu_t)} < \infty$  for all  $t$ . Then, we say a function  $g_{\mu_t}(x) \in L^2(\mu_t)$  is  $(\epsilon_t, \delta_t)$ -AGF if the following condition holds:

$$-\left\langle \nabla \log \frac{\mu_t}{\pi}, g_{\mu_t} \right\rangle_{L^2(\mu_t)} \leq -\epsilon_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 + \delta_t, \quad (9)$$

where  $\epsilon_t, \delta_t \geq 0$ .

Eq. (9) evaluates the approximation quality of  $g_{\mu_t}(x)$  for  $\nabla \log \frac{\mu_t}{\pi}$  via  $\{\epsilon_t, \delta_t\}$ , where  $\epsilon_t$  and  $\delta_t$  express the relative and absolute bias of approximating  $\nabla \log \frac{\mu_t}{\pi}$  by  $g_{\mu_t}(x)$ , respectively. This definition is motivated by the inexact gradient descent methods in finite-dimensional parameter space such as (Jaggi, 2013; Schmidt et al., 2011) and unifies some existing approximate flow methods (see Section 3.2).

Using the  $(\epsilon_t, \delta_t)$ -AGF, we can analyze the convergence in KL divergence qualitatively as follows.

**Lemma 1.** Suppose that Assumption 2 is satisfied. Then, under Eq. (8), for any  $T \in \mathbb{N}$ , we obtain  $\text{KL}(\mu_T|\pi) \leq \prod_{t=0}^{T-1} (1 - \eta_t \epsilon_t) \text{KL}(\mu_0|\pi) + \sum_{t=0}^{T-1} \delta_t \prod_{j=t+1}^{T-1} (1 - \eta_j \epsilon_j)$ .

*Proof.* By substituting Eq. (8) into Eq. (9) and applying the LSI, we obtain  $\text{KL}(\mu_{t+1}|\pi) \leq (1 - \eta_t \epsilon_t) \text{KL}(\mu_t|\pi) + \delta_t$ . By induction in the above, we obtain the claim.  $\square$

This lemma shows that  $\epsilon_t$  and  $\delta_t$  (as well as  $\eta_t$ ) significantly impact the convergence rate.

**Remark 1.** When  $\delta_t = 0$  and  $\eta_t \epsilon_t$  is independent of  $t$ , linear convergence is achieved, indicating that  $g_{\mu_t}(x)$  provides a precise approximation of  $\nabla \log \frac{\mu_t}{\pi}$ . When  $\delta = 0$  and  $\eta_t \epsilon_t = \mathcal{O}(1/t^\alpha)$  with a constant  $\alpha \in (0, 1]$ , it indicates sub-linear convergence, which implies that the approximation quality is not so significant but it is enough to ensure the convergence in KL divergence.

**Remark 2.** If  $\delta_t \neq 0$ , the convergence is biased in terms of KL divergence. However, by employing the technique in Lee et al. (2022), it remains feasible to mitigate the impact of bias on total variation.

### 3.2 RELATION TO EXISTING APPROXIMATE FUNCTIONAL GRADIENT FLOWS

In this section, we provide examples of AGF from existing studies. Dong et al. (2022) proposed the preconditioned functional gradient flow, where they considered approximating  $\nabla \log \frac{\mu_t}{\pi}$  by neural



networks (NNs). The authors also assumed that  $g_{\mu_t}(x)$ , which is the output of NNs, satisfies

$$\left\| g_{\mu_t} - \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 \leq \epsilon \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2, \quad (10)$$

where  $\epsilon < 1$ . This corresponds to a special case of our AGF with  $\epsilon_t := (1 - \epsilon)/2$  and  $\delta_t := 0$ , as confirmed by expanding the left-hand side of Eq. (10). According to Remark 1, the above inequality implies linear convergence in KL divergence. However, their method requires re-training NNs at each iteration, which yields difficulty in ensuring  $\delta = 0$  in practice. Conversely, it later becomes evident that SVGD achieves  $\delta_t = 0$  by using a kernel function that meets some conditions.

Lee et al. (2022; 2023) studied the score based diffusion models assuming that  $g_{\mu_t}(x) = s(x) + \log \mu_t(x)$ , where the  $\nabla \log \pi(x)$  in the Wasserstein gradient is approximated with some measurable function  $s(x)$  that satisfies

$$\left\| (s(\cdot) + \log \mu_t) - \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 \leq \delta. \quad (11)$$

The equation above corresponds to our AGF with  $\epsilon_t = 0$ , which signifies the presence of bias in the KL divergence (see Remark 2).

From the perspective of convergence analysis, the significant difference between these studies lies in the treatment of  $\{\epsilon, \delta\}$ . The convergence analysis in Dong et al. (2022), Lee et al. (2022), and Lee et al. (2023) assumes that  $g_{\mu_t}$  achieves sufficiently small  $\epsilon$  or  $\delta$  according to the criteria in Eq. (10) or (11). In our study, we take the opposite approach — identifying  $\{\epsilon, \delta\}$  that SVGD achieves under the AGF, and then evaluating its convergence properties.

## 4 APPLICATION TO STEIN VARIATIONAL GRADIENT DESCENT

In this section, we present the main result, the convergence of SVGD in KL divergence, obtained by applying the concept of  $(\epsilon, \delta)$ -AGF, and provide an overview of the proofs. Here,  $\mu_t$  represents the  $t$ -th output of the SVGD algorithm, where  $t \in \mathbb{N}$  is the number of iterations as shown in Eq. (6).

### 4.1 SUB-LINEAR CONVERGENCE OF SVGD IN KL DIVERGENCE

Our analyses are based on the following assumptions concerning the kernel function  $k$ .

**Assumption 3.** The feature map  $\nabla k(\cdot, x) : \mathcal{X} \rightarrow \mathcal{H}$  is continuous. Moreover, for all  $x \in \mathcal{X}$ , there exists  $B > 0$  such that  $\|k(\cdot, x)\|_{\mathcal{H}_0} \leq B$ ,  $\sum_{i=1}^d \|\partial_i k(\cdot, x)\|_{\mathcal{H}_0}^2 \leq B^2$ , and  $\sum_{i,j=1}^d \|\partial_i \partial_j k(\cdot, x)\|_{\mathcal{H}_0} \leq B^2$  hold.

**Assumption 4.** The kernel  $k$  is integrally strictly positive definite (ISPD), which means that  $\int \int k(x, y) d\rho(x) d\rho(y) > 0$  holds for all finite nonzero signed Borel measures  $\rho$ .

**Assumption 5.** The trace of a kernel is bounded for any  $\mu \in \mathcal{P}_2(\mathcal{X})$ , i.e.,  $\int k(x, x) d\mu(x) < \infty$ .

Under Assumption 5, the Hilbert–Schmidt operator  $P_{\mu, k}$  has positive eigenvalues  $\{\lambda_i\}$  (see Section 4.2 and Appendix C). We thus further pose the following assumption according to this fact.

**Assumption 6.** Eigenvalues  $\{\lambda_i\}$  are constant order w.r.t.  $t$  and strictly positive, i.e., there exist upper and lower bounds for  $\{\lambda_i\}$  that are independent of  $t$  and are greater than 0.

Assumptions 3-6 are satisfied in the RBF kernel commonly employed in SVGD. A detailed discussion of these assumptions can be found in Appendix A.

We now show the main contribution of this paper, which establishes the sub-linear convergence of SVGD in KL divergence.

**Theorem 1.** Suppose that Assumptions 1-6 are satisfied. Let  $\alpha > 1$  and the stepsize  $\gamma_t$  satisfies  $\gamma_t \leq \mathcal{O}(1/t^{2/3})$  and  $\gamma_t \leq (\alpha - 1)\alpha B^2(1 + \|\nabla V(0) + L\mathbb{E}_\pi\|x\| + L\sqrt{2C_{\text{LS}}^{-1}\text{KL}(\mu_0|\pi)}) (= C_\gamma)$  for all  $t$ . Then, SVGD is  $(c_0, 0)$ -AGF and for any  $T \in \mathbb{N}$ , we have

$$\text{KL}(\mu_T|\pi) \leq \prod_{t=0}^{T-1} (1 - c_0\gamma_t) \text{KL}(\mu_0|\pi), \quad (12)$$

where  $c_0(> 0)$  is a problem-dependent constant that is independent of  $t$ .

This theorem guarantees the sub-linear convergence of SVGD in KL divergence because  $\lim_{t \rightarrow \infty} \frac{\text{KL}(\mu_{t+1}|\pi)}{\text{KL}(\mu_t|\pi)} = \lim_{t \rightarrow \infty} 1 - c_0 \gamma_t = 1$ . Moreover, by setting  $\gamma_t = \frac{c_1}{t}$  for some positive constant  $c_1$  in the above, for example, we obtain  $\text{KL}(\mu_T|\pi) \leq \frac{\text{KL}(\mu_0|\pi)}{T^{c_1 c_0}}$ .

Before outlining the proof, we position our results in comparison to existing studies. As suggested by [Korba et al. \(2021\)](#) and [Duncan et al. \(2023\)](#), it is difficult for SVGD to achieve linear convergence in KL divergence and the difficulty also surfaces in our analysis. To show our results, the step size must be  $\gamma_t \leq \mathcal{O}(1/t^{2/3})$  to control  $\|\nabla \log \frac{\mu_t}{\pi}\|_{L^2(\mu_t)}$ , which highlights the difficulty of achieving convergence faster than sub-linear order. [While Huang et al. \(2023\) has shown the linear convergence in a continuous-time setting, the kernel function utilized in their study is specifically designed to guarantee linear convergence and thus it is not commonly employed in practice.](#) On the other hand, our result is established within the discrete time setting that corresponds to the SVGD algorithm, under realistic assumptions commonly met by the RBF kernel frequently adopted in SVGD.

Expanding our sight to other deterministic sampling methods based on [kernel functions](#), sub-linear convergence has been demonstrated in the kernel herding (e.g., [Chen et al. \(2010\)](#); [Bach et al. \(2012\)](#)) and Bayesian Quadrature context (e.g., [Briol et al. \(2015\)](#); [Futami et al. \(2019\)](#)) when employing infinite-dimensional kernel functions like the RBF kernel. Our results are consistent with these facts.

#### 4.2 SPECTRAL DECOMPOSITION AND $(\epsilon, \delta)$ -APPROXIMATION

The main objective here is to provide an overview of the proof focusing on how we detect  $\epsilon_t$  and  $\delta_t$  in the AGF. The complete proof is in [Appendix C](#).

To conduct analyses based on our AGF, we need to show that  $\|\nabla \log \frac{\mu_t}{\pi}\|_{L^2(\mu_t)}$  is bounded for all  $t$  in SVGD, which is guaranteed by the following lemma (see [Appendix C.2](#) for complete proof).

**Lemma 2.** *Suppose that Assumptions 1-3 and 5 are satisfied. Let  $\gamma_t$  satisfies  $\gamma_t \leq C_\gamma$  defined in Theorem 1. Then, there exists a positive problem-dependent constant  $c$  and is independent of  $t$  such that, for any  $t \in (0, T]$  we have  $\|\nabla \log \frac{\mu_t}{\pi}\|_{L^2(\mu_t)} \leq \|\nabla \log \frac{\mu_0}{\pi}\|_{L^2(\mu_0)} + c \sum_{t=0}^{t-1} \gamma_t$ .*

Now we are ready to begin the analysis of the convergence of SVGD based on AGF. Substituting  $g_{\mu_t}(x) = P_{\mu_t, k} \nabla \log \frac{\mu_t}{\pi}$  into [Eq. \(9\)](#) and multiplying both sides by  $\eta_t (> 0)$  yields

$$-\eta_t \left\langle \nabla \log \frac{\mu_t}{\pi}, P_{\mu_t, k} \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)} = -\eta_t I_{\text{stein}}(\mu_t|\pi) \leq -\epsilon_t \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 + \eta_t \delta_t. \quad (13)$$

According to the fact that  $\eta_t I_{\text{stein}}(\mu_t|\pi) \leq \text{KL}(\mu_0|\pi)$  (see [Appendix C](#)), we further obtain the following inequalities:

$$\text{KL}(\mu_0|\pi) \geq \eta_t I_{\text{stein}}(\mu_t|\pi) \geq \epsilon_t \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 - \eta_t \delta_t. \quad (14)$$

Therefore, our goal is to guarantee the existence of the above inequality. If [Eq. \(14\)](#) exists, we can qualitatively analyze the convergence in KL divergence by specifying  $\{\epsilon_t, \delta_t\}$  and utilizing the property of AGF shown in [Lemma 1](#) and [Remarks 1](#) and [2](#).

To focus on the discussion for detecting  $\{\epsilon, \delta\}$ , we first mention the necessary conditions for the existence of [Eq. \(14\)](#) w.r.t.  $\eta_t$  under our final results. As can be seen from [Theorem 1](#), we obtain  $\epsilon_t = c_0$  and  $\delta_t = 0$  through the proof that we explain later, where  $c_0$  is independent of  $t$ . In this case, from [Eq. \(13\)](#), it is necessary for  $\eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2$  to be uniformly upper bounded w.r.t.  $t$  to compensate for the convergence based on AGF. This condition can be satisfied by setting  $\eta_t$  such that it fulfills  $\gamma_t \leq \mathcal{O}(1/t^{2/3})$  from [Lemma 2](#) (see [Appendix C](#) for this derivation).

Our strategy is to show the boundedness of the following equality expressed as

$$\epsilon_t \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 - \eta_t I_{\text{stein}}(\mu_t|\pi) = \eta_t \left\langle \nabla \log \frac{\mu_t}{\pi}, (\epsilon_t I - P_{\mu_t, k}) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)}. \quad (15)$$

We adopt the spectral decomposition of the kernel operator to analyze the above. Since a Hilbert-Schmidt operator  $P_{\mu, k}$  is compact and self-adjoint, we have, for all  $i$ ,  $P_{\mu, k} \phi_i = \lambda_i \phi_i$ , where  $\phi_i \in$

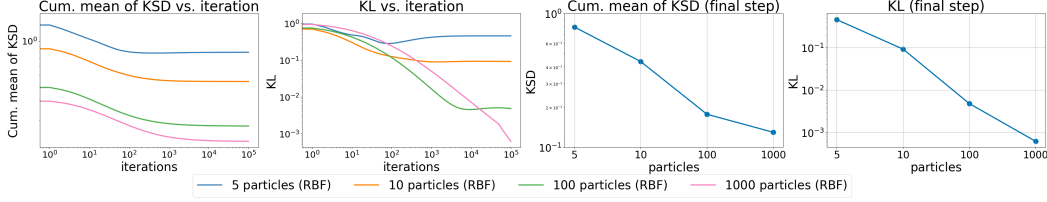


Figure 1: The convergence behavior in terms of  $\text{KL}(\mu_T|\pi)$  and  $\frac{1}{T} \sum_{t=1}^T I_{\text{stein}}(\mu_t|\pi)$  for all  $T$  under two-dimensional Gaussian distribution experiments ( $\beta = 0.67 \approx 2/3$ ).

$L^2(\mu)$  represents an eigenfunction that satisfies a complete orthonormal system (CONS), and  $\lambda_i$  is an eigenvalue corresponding to  $\phi_i$ . Even if these eigenvalues are ordered as  $\lambda_1 \geq \lambda_2 \geq \dots > 0$ , it does not compromise generality. Moreover, the kernel function can be decomposed into  $k(x, y) = \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(y)$ , where the convergence of this infinite series holds in the norm of  $\|\cdot\|_{L^2(\mu)}$ .

Defining  $v_t := \nabla \log \frac{\mu_t}{\pi}$  for simplicity in notation, we can obtain  $v_t = \sum_{i=1}^{\infty} \langle v_t, \phi_i \rangle_{L^2(\mu)} \phi_i$  and  $P_{\mu, k} v_t = \sum_{i=1}^{\infty} \lambda_i \langle v_t, \phi_i \rangle_{L^2(\mu)} \phi_i$  because the kernel function is dense in  $L^2(\mu)$  and thus its eigenvectors are complete. We provide the discussion for non-complete eigenvectors in Appendix A. Substituting these equalities into Eq. (15), we have

$$\epsilon_t \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 - \eta_t I_{\text{stein}}(\mu_t|\pi) = \eta_t \sum_{i=1}^{\infty} (\epsilon_t - \lambda_i) \langle v_t, \phi_i \rangle_{L^2(\mu_t)}^2. \quad (16)$$

In the right-hand side term of the above, there exists a index  $1 < j$  such that  $\lambda_j > \epsilon_t > \lambda_{j+1}$  by setting sufficiently small  $\epsilon_t$ . Hence, by regularizing  $\{\langle v_t, \phi_i \rangle_{L^2(\mu_t)}^2\}_{i=1}^{\infty}$ , we can render the left-hand side of Eq. (16) negative. For that purpose, we focus on the RKHS associated with  $k$  given as  $\mathcal{H} = \{f \in L^2(\mu) \mid f = \sum_{i=1}^{\infty} a_i \phi_i, \sum_{i=1}^{\infty} \lambda_i^{-1} \|a_i\|^2 < \infty, a_i \in \mathbb{R}\}$ , where  $\mathcal{H}$  is dense in  $L^2(\mu)$ . In this RKHS, there exists a function  $v_t^{(l)} \in \mathcal{H}$  such that the sequence of  $v_t^{(l)} \rightarrow v_t$  as  $l \rightarrow \infty$  in  $L^2(\mu)$  norm. Thus, by approximating the original  $v_t$  with  $v_t^{(l)}$  in  $\mathcal{H}$ , we can regularize  $\{\langle v_t, \phi_i \rangle_{L^2(\mu_t)}^2\}_{i=1}^{\infty}$ .

Under the regularized  $\{\langle v_t, \phi_i \rangle_{L^2(\mu_t)}^2\}_{i=1}^{\infty}$  in the above and sufficiently small  $\epsilon_t$ , we can obtain  $\epsilon_t \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 - \eta_t I_{\text{stein}}(\mu_t|\pi) < 0$ , which implies that  $\delta_t = 0$  in the AGF. From Assumption 6, we can show that  $\epsilon_t$  is the constant order w.r.t.  $t$  and express it as  $c_0$  (see Appendix C). This concludes the proof outline.

## 5 NUMERICAL EXPERIMENTS

In this section, we aim to confirm the validity of our theoretical results. We only show the results of the two-dimensional Gaussian experiments due to the page limitation. The details of the experimental settings and additional results including the Gaussian mixture can be seen in Appendix E.

We set the target distribution as the two-dimensional Gaussian distribution. We adopted the RBF kernel  $k(x, y) = \exp(-\frac{1}{h} \|x - y\|_2^2)$ , which is commonly used in practice and satisfies the assumptions in Section 4. The bandwidth  $h$  was selected by the median trick as in Liu & Wang (2016). To appropriately verify our theoretical analysis, we simply set the decaying step size  $\gamma_t = 1/(1+t^\beta)$  ( $= \mathcal{O}(1/t^\beta)$ ) suggested by Theorem 1 and did not use the Adagrad-based stepsize, which is adopted in related studies such as Korba et al. (2021) and others. We evaluated the KL divergence:  $\text{KL}(\mu_T|\pi)$  and the cumulative mean of KSD:  $\frac{1}{T} \sum_{t=1}^T I_{\text{stein}}(\mu_t|\pi)$ , which are theoretically guaranteed sub-linear convergence.

**Results:** From Figures 1 and 2, we can see that SVGD with the RBF kernel tends to achieve sub-linear convergence both in  $\text{KL}(\mu_T|\pi)$  and in  $\frac{1}{T} \sum_{t=1}^T I_{\text{stein}}(\mu_t|\pi)$ , which supports Theorem 1. As discussed in Appendix A, the bias remains in the KL divergence as we increase  $T$  since we used the finite particles and thus  $\delta_t \neq 0$  in AGF. Such a bias can be reduced by increasing the number of particles increases. Conversely, employing a substantial number of particles leads to



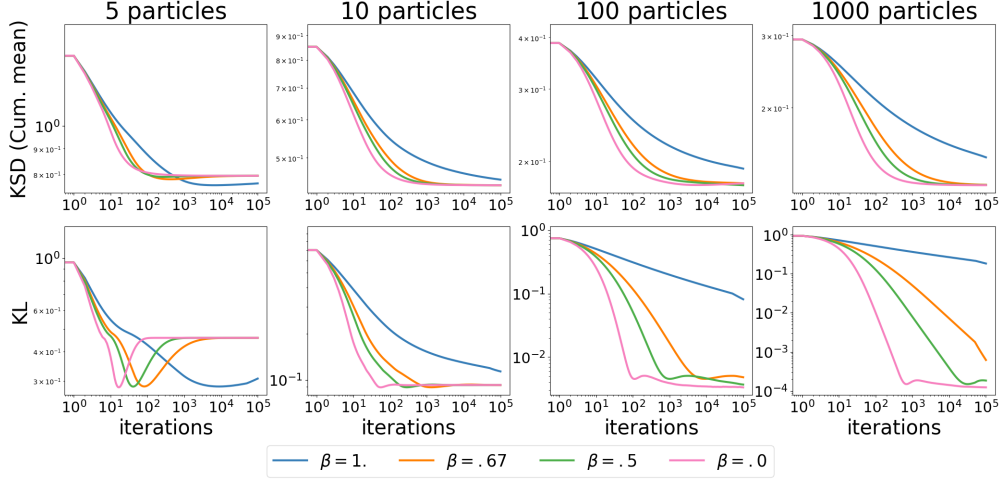


Figure 2: Convergence in  $\text{KL}(\mu_T|\pi)$  and  $\frac{1}{T} \sum_{t=1}^T I_{\text{stein}}(\mu_t|\pi)$  for all  $T$  under different particles and stepsize settings ( $\beta = \{0., 0.5, 0.67, 1.\}$ ).

slower convergence for both the KSD and KL divergence. This phenomenon may be attributed to the presence of exceedingly small eigenvalues of  $P_{\mu,k}$  when using a larger number of particles since the eigenvalues of the RBF kernel decay exponentially fast (Wainwright, 2019). Our numerical evaluation of eigenvalues can be seen in Appendix E. In other words, there exists a trade-off between the improvement in the approximation accuracy achieved by using a large number of particles and the convergence speed.

## 6 LIMITATION & CONCLUSION

Ensuring the convergence of SVGD in KL divergence has proven challenging in finite and infinite particle settings. Furthermore, while many studies have provided convergence guarantees for SVGD in KSD, these do not necessarily ensure its weak convergence. As a *first strategy to address this issue*, we conducted the convergence analysis of SVGD under the ideal conditions of an infinite particle setting that guarantees an accurate gradient approximation. Then, we successfully elucidated the convergence of SVGD in KL divergence in this setting. Our finding suggests weak convergence of SVGD with infinite particles, affirming its capability to approximate the expectation by the target distribution without bias, akin to MCMC.

One of limitations in our paper is the challenge in furnishing a theoretical explanation for the convergence of SVGD when employing a finite number of particles. Extending our analysis to finite particle settings using AGF is being considered as our future study. The main challenge in this extension is expected to be in determining the values of  $\epsilon$  and  $\delta$ , primarily due to the unknown theoretical properties of gradient approximation on RKHS when dealing with correlated particles, as far as our current knowledge extends.

Another limitation is that Assumption 6 is rather strong. This assumption, introduced to ensure that  $c_0$  remains of constant order w.r.t.  $t$ , is difficult to justify in the infinite particle setting. The pursuit of convergence guarantees grounded in milder assumptions represents a crucial avenue for future research. Furthermore, we aspire for this study to catalyze further research endeavors that aim to furnish better convergence guarantee for SVGD in KL divergence.

**Ethics Statement:** Since our paper is fundamental research based on theory, we believe that it does not cause any potential harm, societal impact, or potentially harmful consequences.

## REFERENCES

- L. Ambrosio, N. Gigli, and G. Savaré. *Gradient flows: In metric spaces and in the space of probability measures*. Springer Science & Business Media, 2005.
- F. Bach, S. Lacoste-Julien, and G. Obozinski. On the equivalence between herding and conditional gradient algorithms. In *Proceedings of the 29th International Conference on Machine Learning*, pp. 1355–1362, 2012.
- D. Bakry, I. Gentil, and M. Ledoux. *Analysis and Geometry of Markov Diffusion Operators*, volume 348. Springer Science & Business Media, 2013.
- D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- B. Bonnet and H. Frankowska. Differential inclusions in Wasserstein spaces: The Cauchy-Lipschitz framework. *Journal of Differential Equations*, 271:594–637, 2021.
- F.-X. Briol, C. Oates, M. Girolami, and M. A. Osborne. Frank–Wolfe Bayesian quadrature: Probabilistic integration with theoretical guarantees. *Advances in Neural Information Processing Systems*, 28:1162–1170, 2015.
- C. Carmeli, E. De Vito, A. Toigo, and V. Umanitá. Vector valued reproducing kernel Hilbert spaces and universality. *Analysis and Applications*, 8(01):19–61, 2010.
- Y. Chen, M. Welling, and A. Smola. Super-samples from kernel herding. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pp. 109–116, 2010.
- H. Dong, X. Wang, L. Yong, and T. Zhang. Particle-based variational inference with preconditioned functional gradient flow. In *The Eleventh International Conference on Learning Representations*, 2022.
- A. Duncan, N. Nüsken, and L. Szpruch. On the geometry of Stein variational gradient descent. *Journal of Machine Learning Research*, 24(56):1–39, 2023.
- R. Durrett. *Probability: Theory and Examples*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 5 edition, 2019.
- F. Futami, Z. Cui, I. Sato, and M. Sugiyama. Bayesian posterior approximation via greedy particle optimization. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):3606–3613, 2019.
- J. Gorham and L. Mackey. Measuring sample quality with kernels. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 1292–1301, 2017.
- L. Gross. Logarithmic Sobolev inequalities. *American Journal of Mathematics*, 97(4):1061–1083, 1975.
- W. K. Hastings. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57(1):97–109, 1970.
- X. Huang, H. Dong, and C. Fang. Local KL convergence rate for Stein variational gradient descent with reweighted kernel, 2023. URL <https://openreview.net/forum?id=k2CRIF8tJ7Y>.
- M. Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. *Proceedings of the 30th International Conference on Machine Learning*, 28(1):427–435, 2013.
- M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. *Machine Learning*, 37(2):183–233, 1999.

- R. Jordan, D. Kinderlehrer, and F. Otto. The variational formulation of the Fokker–Planck equation. *SIAM journal on mathematical analysis*, 29(1):1–17, 1998.
- A. Korba, A. Salim, M. Arbel, G. Luise, and A. Gretton. A non-asymptotic analysis for Stein variational gradient descent. *Advances in Neural Information Processing Systems*, 33:4672–4682, 2020.
- A. Korba, P.-C. Aubin-Frankowski, S. Majewski, and P. Ablin. Kernel stein discrepancy descent. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pp. 5719–5730, 2021.
- H. Lee, J. Lu, and Y. Tan. Convergence for score-based generative modeling with polynomial complexity. *Advances in Neural Information Processing Systems*, 35:22870–22882, 2022.
- H. Lee, J. Lu, and Y. Tan. Convergence of score-based generative modeling for general data distributions. In *International Conference on Algorithmic Learning Theory*, pp. 946–985, 2023.
- C. Liu, J. Zhuo, P. Cheng, R. Zhang, and J. Zhu. Understanding and accelerating particle-based variational inference. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pp. 4082–4092, 2019.
- Q. Liu. Stein variational gradient descent as gradient flow. In *Advances in Neural Information Processing Systems*, volume 30, pp. 3118–3126, 2017.
- Q. Liu and D. Wang. Stein variational gradient descent: A general purpose Bayesian inference algorithm. In *Advances in Neural Information Processing Systems*, volume 29, pp. 2378–2386, 2016.
- T. Liu, P. Ghosal, K. Balasubramanian, and N. S. Pillai. Towards understanding the dynamics of Gaussian-Stein variational gradient descent. *arXiv preprint arXiv:2305.14076*, 2023.
- F. Otto. The geometry of dissipative evolution equations: The porous medium equation. *Communications in Partial Differential Equations*, 26:101–174, 2001.
- C. P. Robert and G. Casella. *Monte Carlo statistical methods*. Springer Texts in Statistics, 2004.
- L. Rosasco, M. Belkin, and E. De Vito. On learning with integral operators. *Journal of Machine Learning Research*, 11(30):905–934, 2010.
- A. Salim, L. Sun, and P. Richtarik. A convergence theory for SVGD in the population limit under Talagrand’s inequality T1. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pp. 19139–19152, 2022.
- M. Schmidt, N. Roux, and F. Bach. Convergence rates of inexact proximal-gradient methods for convex optimization. In *Advances in Neural Information Processing Systems*, volume 24, pp. 455–476, 2011.
- J. Shi and L. Mackey. A finite-particle convergence rate for Stein variational gradient descent. *arXiv preprint arXiv:2211.09721*, 2023.
- B. K. Sriperumbudur, A. Gretton, K. Fukumizu, B. Schölkopf, and G. RG Lanckriet. Hilbert space embeddings and metrics on probability measures. *The Journal of Machine Learning Research*, 11:1517–1561, 2010.
- B. K. Sriperumbudur, K. Fukumizu, and G. R.G. Lanckriet. Universality, characteristic kernels and RKHS embedding of measures. *Journal of Machine Learning Research*, 12(70):2389–2410, 2011.
- I. Steinwart and A. Christmann. *Support Vector Machines*. Springer Publishing Company, Incorporated, 1st edition, 2008.
- L. Sun, A. Karagulyan, and P. Richtarik. Convergence of Stein variational gradient descent under a weaker smoothness condition. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206, pp. 3693–3717, 2023.

- C. Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg, 2008.
- M. J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- M. Welling and Y. W. Teh. Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, pp. 681–688, 2011.

## A DISCUSSION REGARDING THE ASSUMPTIONS

It should be mentioned that Assumptions 3-5 is satisfied in the RBF kernel commonly used in SVGD. Assumption 3 requires the bounded twice differentiability of the kernel, which is stronger than existing work (Salim et al., 2022) but is required to control  $\|\nabla \log \frac{\mu_t}{\pi}\|_{L^2(\mu_t)}$  in Lemma 2. The existence of the constant  $B$  in this assumption depends on the choice of bandwidth for the RBF kernel and can be guaranteed by adopting standard selection methods such as the median trick (Liu & Wang, 2016). Assumption 5 holds true since the RBF kernel satisfies  $\int k(x, x) d\mu_t(x) = 1$ , while Assumption 6 is introduced as a technical assumption to ensure that  $c_0$  remains of constant order w.r.t.  $t$ . Confirming the validity of the latter assumption for infinite-dimensional particles is challenging. However, at the very least, it has been empirically observed that as the number of particles increases, there is a tendency for the eigenvalues to become independent of time (see Figures 7-10). In Section 5, we confirm that the sub-linear convergence in KL divergence guaranteed under these assumptions aligns with numerical experiments, while the need for better proof remains an important future task.

Assumption 4 assures that the projection  $P_{\mu,k}$  is injective and characteristic (Sriperumbudur et al., 2010), and also implies certain types of universality (Sriperumbudur et al., 2011), which hold in the RBF kernel. Thanks to these properties, we can expand the Wasserstein gradient via the eigenvectors as Eq. (16) since the eigenvectors of  $P_{\mu,k}$  can be CONS in  $L^2(\mu)$ . Furthermore, SVGD ensures  $\delta_t = 0$  under Assumption 4, a critical condition for convergence in KL divergence (see Section 3.2), which sets it apart from other approximate flow methods that face challenges in achieving  $\delta_t = 0$ . From these discussions, it is apparent that Assumption 4 plays a pivotal role to guarantee sub-linear convergence of SVGD in KL divergence, while being a stronger assumption than in existing studies (Korba et al., 2020; Salim et al., 2022; Sun et al., 2023).

If the eigenvectors of  $P_{\mu,k}$  is not CONS in  $L^2(\mu)$ , there is a bias to approximate  $v_t := \nabla \log \frac{\mu_t}{\pi}$  in RKHS. For instance, when we consider the null space of  $P_{\mu,k}$  as  $\text{Null}(P_{\mu,k}) := \{f \in L^2(\mu) | P_{\mu,k}f = 0\}$ , we have  $v_t = \sum_{i=1}^{\infty} \langle v_t, \phi_i \rangle_{L^2(\mu)} \phi_i + \Psi$ , where  $\Psi \in \text{Null}(P_{\mu,k})$  is some appropriately chosen function. In this case, we can express the right-hand side of Eq. (15) as  $\eta_t \sum_{i=1}^{\infty} (\epsilon_t - \lambda_i) \langle v_t, \phi_i \rangle_{L^2(\mu)}^2 + \eta_t \|\Psi\|_{L^2(\mu)}^2$ , which implies  $\delta_t = \|\Psi\|_{L^2(\mu)}^2 (\neq 0)$  in the AGF. Thus, from Remark 2, the KL divergence does not go to 0 even when we increase  $T$ . Such cases may arise when using linear or polynomial kernels in SVGD since these do not satisfy Assumption 4. Also, the eigenvectors are not guaranteed to be CONS when using a finite number of particles. If we draw  $m$  independent and identically distributed (i.i.d.) samples from  $\mu(x)$  and approximate the kernel operator, we can only access the first  $m$ -dimensional eigenvectors (Rosasco et al., 2010) and thus corresponding eigenvectors cannot be CONS in  $L^2(\mu)$ . Note that when considering the SVGD with finite particles, since the particles are not i.i.d. and correlated significantly, the approximation quality is much worse than the case of i.i.d. particles and results in larger  $\delta_t$ .

## B PRELIMINARIES

### B.1 GRADIENT FLOW

We are interested in sampling from the density  $\pi \in \mathcal{P}_2(\mathcal{X})$  and proportional as  $\pi \propto \exp^{-L(x)}$ , where  $V : \mathcal{X} \rightarrow \mathbb{R}$ . We assume that  $V \in C^2(\mathcal{X})$  and its Hessian  $H_V$  satisfies  $\|H_V\|_{\text{op}} \leq L$ .

We are interested in sampling from  $\pi$  and formalizing the task as the optimization problem. For that purpose, We define the Kullback-Leibler (KL) divergence. For any  $\mu, \pi \in \mathcal{P}_2(\mathcal{X})$ , KL divergence of  $\mu$  w.r.t  $\pi$  is defined by

$$\text{KL}(\mu|\pi) := \int \log \frac{d\mu}{d\pi}(x) d\mu(x) \quad (17)$$

if  $\mu$  is absolutely continuous w.r.t  $\pi$  and admits Radon-Nikodym derivative  $d\mu/d\pi$  and  $\text{KL}(\mu|\pi) = +\infty$  otherwise. We consider a functional  $\text{KL}(\cdot|\pi) : \mathcal{P}_2(\mathcal{X}) \rightarrow [0, +\infty)$ ,  $\mu \rightarrow \text{KL}(\mu|\pi)$  defined on over Wasserstein space:

$$\min_{\mu \in \mathcal{P}_2(\mathcal{X})} \mathcal{F}(\mu), \quad \mathcal{F}(\mu) := \text{KL}(\mu|\pi). \quad (18)$$



The advantage of using Wasserstein space is its differential structure to minimize this kind of functional.

Before introducing the differential structure, we introduce the continuous map. Let  $T > 0$  and consider a weakly continuous map  $\mu : (0, T) \rightarrow \mathcal{P}_2(\mathcal{X})$ ,  $t \rightarrow \mu_t$ . The family  $(\mu_t)_{t \in (0, T)}$  satisfies a continuity equation if there exists  $(v_t)_{t \in (0, T)}$  such that  $v_t \in L^2(\mu_t)$ :

$$\frac{\partial \mu_t}{\partial t} + \operatorname{div}(\mu_t v_t) = 0 \quad (19)$$

holds in the sense of distributions, i.e., for any  $\phi \in C_c^\infty(\mathcal{X})$ ,

$$\frac{d}{dt} \int \phi(x) d\mu_t(x) = \langle \nabla \phi, v_t \rangle_{L^2(\mu_t)} = \int \langle \nabla \phi(x), v_t(x) \rangle d\mu_t(x) \quad (20)$$

holds for any  $t \in (0, T)$ . And by integration parts, we have

$$\int \phi(x) \frac{\partial \mu_t(x)}{\partial t} + \int \phi(x) \operatorname{div}(\mu_t(x) v_t(x)) = 0. \quad (21)$$

Although the Wasserstein space is not a Riemannian manifold, it can be equipped with a Riemannian structure and interpretation (Otto, 2001) and the tangent space of  $\mathcal{P}_2(\mathcal{X})$  at  $\mu_t$  denoted  $\mathcal{T}_{\mu_t} \mathcal{P}_2(\mathcal{X})$  is a subset of  $L^2(\mu_t)$ . Under this setting, among all possible  $(v_t)_{t \in (0, T)}$ , we call  $v_t$  that has the minimal  $L^2(\mu_t)$  norm as the velocity field of  $(\mu_t)_{t \in (0, T)}$  and this minimality condition can be characterized by  $v_t \in \mathcal{T}_{\mu_t} \mathcal{P}_2(\mathcal{X}) \subset L^2(\mu_t)$ .

To select the appropriate  $(v_t)_{t \in (0, T)}$ , it is useful to use the differential structure of the Wasserstein space. Assume that given a proper lower semi-continuous functional  $\mathcal{F} : \mathcal{P}_2(\mathcal{X}) \rightarrow \mathbb{R}$  and  $\mu \in \mathcal{P}_2(\mathcal{X})$ ,  $\xi \in L^2(\mu)$  is a strong subdifferential of  $\mathcal{F}$  at  $\mu$  if for every  $\phi \in L^2(\mu)$  and for every  $\epsilon \in (0, 1]$ ,

$$\mathcal{F}(\mu) + \epsilon \langle \xi, \phi \rangle_{L^2(\mu)} + o(\epsilon) \leq \mathcal{F}((I + \epsilon \phi) \# \mu). \quad (22)$$

where  $I$  is the identity map.

Then the important consequence is that under the mild regularity conditions,  $W_2$  gradient of  $\mathcal{F}$  corresponds to the first variation of the functional. Assume that given a functional  $\mathcal{F} : \mathcal{P}_2(\mathcal{X}) \rightarrow \mathbb{R}$ , we call  $\frac{\delta \mathcal{F}}{\delta \mu}(\mu)$  as the first variation of  $\mathcal{F}$  at  $\mu$

$$\int \frac{\delta \mathcal{F}}{\delta \mu}(\mu) \phi d\xi = \frac{d}{d\epsilon} \mathcal{F}(\mu + \epsilon \xi) \Big|_{\epsilon=0} = \lim_{\epsilon \rightarrow 0} \frac{\mathcal{F}(\mu + \epsilon \xi) - \mathcal{F}(\mu)}{\epsilon} \quad (23)$$

for all  $\xi = \nu - \mu$  where  $\nu \in \mathcal{P}_2(\mathcal{X})$ .

From Lemma 10.4.1 in Ambrosio et al. (2005), given  $\mu \in \mathcal{P}_2(\mathcal{X})$ , which is absolutely continuous w.r.t. the Lebesgue measure and its density is  $C^1(\mathcal{X})$  and assume that  $\nabla_{W_2} \mathcal{F}(\mu)(x)$  belongs to the strong subdifferential of  $\mathcal{F}$  at  $\mu$ . Then it is given as

$$\nabla_{W_2} \mathcal{F}(\mu)(x) = \nabla \frac{\delta \mathcal{F}}{\delta \mu}(\mu)(x) \quad \text{for } \mu - \text{a.e. } x \in \mathbb{R}^d \quad (24)$$

and for every vector field  $\xi \in C_c^\infty(\mathbb{R}^d, \mathbb{R}^d)$ ,

$$\langle \nabla_{W_2} \mathcal{F}(\mu), \xi \rangle_{L^2(\mu)} = - \int_{\mathbb{R}^d} \nabla \frac{\delta \mathcal{F}}{\delta \mu}(\mu)(x) \operatorname{div}(\mu(x) \xi(x)) dx \quad (25)$$

and  $\nabla_{W_2} \mathcal{F}(\mu)$  belongs to  $\mathcal{T}_\mu \mathcal{P}(\mathbb{R}^d)$ , which is a subset of  $L^2(\mu)$ .

Now we are ready to leverage this differential structure of the Wasserstein space to minimize the functional  $\mathcal{F}$ . We consider how  $\mathcal{F}$  evolves by the continuity equation.

$$\dot{\mathcal{F}}(\mu_t) := \frac{d}{dt} \mathcal{F}(\mu_t) = \langle \nabla_{W_2} \mathcal{F}(\mu_t), v_t \rangle_{L^2(\mu_t)} \quad (26)$$

Thus, by choosing  $v_t$  such that  $\langle \nabla_{W_2} \mathcal{F}(\mu_t), v_t \rangle_{L^2(\mu_t)} \leq 0$ , we can minimize the functional by using the continuity equation. A natural choice is to use the Wasserstein gradient itself as  $v_t = -\nabla_{W_2} \mathcal{F}(\mu_t)$ , which results in

$$\dot{\mathcal{F}}(\mu_t) := \frac{d}{dt} \mathcal{F}(\mu_t) = -\|\nabla_{W_2} \mathcal{F}(\mu_t)\|_{L^2(\mu_t)}^2 \leq 0. \quad (27)$$

## B.2 KL DIVERGENCE

Now we focus on minimizing KL divergence. It is known that the Wasserstein gradient of the KL divergence is given as  $\nabla_{W_2} \text{KL}(\mu|\pi) = \nabla \log \frac{\mu}{\pi} \in L^2(\mu)$ . Then by setting  $v_t = -\nabla_{W_2} \text{KL}(\mu|\pi) = -\nabla \log \frac{\mu}{\pi}$  as the velocity field of the continuity equation, we have that

$$\frac{d}{dt} \text{KL}(\mu_t|\pi) = - \left\| \nabla \log \frac{\mu}{\pi} \right\|_{L^2(\mu_t)}^2. \quad (28)$$

When considering the forward discretization, we obtain the gradient descent algorithm in the Wasserstein space

$$\mu_{t+1} = \left( I - \gamma \nabla \log \frac{\mu_n}{\pi} \right) \# \mu_t \quad (29)$$

where  $\gamma > 0$  is a stepsize.

To analyze this time evolution in the continuous dynamics, many existing work assumes that  $\pi$  satisfies the logarithmic Sobolev inequality (LSI) (Bakry et al., 2013). The important consequence of this inequality is that there exists a positive constant  $C_{\text{LS}}$  such that

$$\text{KL}(\mu_t|\pi) \leq \frac{1}{C_{\text{LS}}} \left\| \nabla \log \frac{\mu}{\pi} \right\|_{L^2(\mu_t)}^2. \quad (30)$$

With this inequality, we have that

$$\text{KL}(\mu_t|\pi) \leq e^{-C_{\text{LS}}t} \text{KL}(\mu_0|\pi) \quad (31)$$

Thus, the Wasserstein gradient flow of the KL divergence achieves linear convergence.

How can we implement  $\nabla_{W_2} \text{KL}(\mu|\pi) = \nabla \log \frac{\mu}{\pi} \in L^2(\mu)$  in practice ? It is known that the Langevin dynamics are given as

$$dX_t = \nabla \log \pi(X_t) dt + \sqrt{2} dB_t, \quad (32)$$

where  $dB_t$  is the standard Brownian motion in  $\mathbb{R}^d$  can be seen as the implementation of the Wasserstein gradient flow in a probabilistic way. We can easily confirm that the probability density induced by SDE of Eq. 32 is given as the Fokker–Planck(FP) equation and that of the FP equation is equivalent to the continuity equation given by the Wasserstein gradient flow in a distributional sense Jordan et al. (1998).

## B.3 STEIN VARIATIONAL GRADIENT DESCENT (SVGD)

The LD method was successfully used as an MCMC method, however, it suffers from large bias when using obtained samples. To alleviate this, an alternative gradient flow approach has been developed. Among those methods, SVGD has been used extensively in practice. Since SVGD is a kind of projection of the Wasserstein gradient into the reproducing kernel Hilbert space (RKHS), here we first introduce the settings of RKHS.

Let a semi-positive definite kernel  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  and  $\mathcal{H}_0$  is its corresponding RKHS of real-valued functions  $\mathcal{X} \rightarrow \mathbb{R}$ . We express the inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}_0}$ . Due to the reproducing property  $\forall f \in \mathcal{H}_0, f = \langle f, k(x, \cdot) \rangle_{\mathcal{H}_0}$ . We define by  $\mathcal{H}$  as the Cartesian product of  $\mathcal{H}_0$ , its element  $f \in \mathcal{H}$ ,  $f = (f_1, \dots, f_d)$  and  $f_i \in \mathcal{H}_0$  for  $i = 1, \dots, d$ . Then the inner product of  $\mathcal{H}$  is given as  $\langle f, g \rangle_{\mathcal{H}} = \sum_{i=1}^d \langle f_i, g_i \rangle_{\mathcal{H}_0}$ . Let  $\mu \in \mathcal{P}_2(\mathcal{X})$  and if  $\int k(x, x) d\mu(x) < \infty$ , then the integral operator associated to  $k$  and  $\mu$  denoted by  $S_{\mu, k} : L^2(\mu) \rightarrow \mathcal{H}$  is defined as

$$S_{\mu, k} f := \int k(x, \cdot) f(x) d\mu(x). \quad (33)$$

By definition, note that  $\mathcal{H} \subset L^2(\mu)$ . We also define the inclusion map as  $\iota : \mathcal{H} \rightarrow L^2(\mu)$  and it is the adjoint of  $S_{\mu, k}$ . Thus for  $f \in L^2(\mu)$  and  $g \in \mathcal{H}$ , we have

$$\langle f, \iota g \rangle_{L^2(\mu)} = \langle \iota^* f, g \rangle_{\mathcal{H}} = \langle S_{\mu, k} f, g \rangle_{\mathcal{H}}. \quad (34)$$

We also define  $P_{\mu, k} : L^2(\mu) \rightarrow L^2(\mu)$ ,  $P_{\mu, k} = \iota S_{\mu, k}$ .

In SVGD, instead of using the Wasserstein gradient directly, we consider using  $P_{\mu,k}\nabla\log\frac{\mu}{\pi}$  as  $v_t$ . Then we obtain the discretized dynamics of SVGD

$$\mu_{t+1} = \left(I - \gamma P_{\mu,k}\nabla\log\frac{\mu_n}{\pi}\right) \# \mu_t \quad (35)$$

and if the kernel satisfies  $\lim_{\|x\|\rightarrow\infty} k(x, \cdot)\pi(x) = 0$ , we obtain

$$P_{\mu,k}\nabla\log\frac{\mu}{\pi}(\cdot) := - \int [\nabla\log\pi(x)k(x, \cdot) + \nabla_x k(x, \cdot)] d\mu(x) \quad (36)$$

by using an integration by parts (Liu, 2017). Then in the original work in Liu & Wang (2016), to approximate the expectation of  $\mu$ , we first prepare initial samples  $\{x_0^m\}_{m=1}^M$  and iteratively update them by a transformation

$$x_{t+1}^m = x_n^m - \frac{1}{M} \sum_{m'=1}^M \nabla\log\pi(x_n^{m'})k(x_n^{m'}, x_n^m) + \nabla_x k(x_n^{m'}, x_n^m) \quad (37)$$

where we approximate  $\mu_n$  by a finite set of particles  $\hat{\mu}_n = \frac{1}{M} \sum_{m=1}^M \delta_{x_n^m}$  where  $\delta_x$  is the Dirac measure with its mass at  $x$ . Originally, Liu & Wang (2016) derived the push forward as follows. Assume that the pushforward is given as  $T(x) = x - \gamma\phi(x)$ , where  $\gamma$  is a positive constant  $\phi(\cdot) \in \mathcal{H}$  is a perturbation direction. Then the update direction, which maximally decreases the Kullback–Leibler (KL) divergence between the particles and the target distribution,

$$\phi^*(x) = \arg\max_{\phi \in \mathcal{H}, \|\phi\|_{\mathcal{H}} \leq 1} \left\{ -\frac{d}{d\gamma} \text{KL}(T\#\mu|\pi)|_{\gamma=0} \right\}, \quad (38)$$

then the solution of this is  $\phi^*(\cdot) = S_{\mu,k}\nabla\log\frac{\mu}{\pi}$ .

Going back to the continuous dynamics,

$$\begin{aligned} \frac{d}{dt} \text{KL}(\mu_t|\pi) &= - \left\langle \nabla\log\frac{\mu_t}{\pi}, P_{\mu_t,k}\nabla\log\frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)} \\ &= - \left\langle \iota^* \nabla\log\frac{\mu_t}{\pi}, S_{\mu_t,k}\nabla\log\frac{\mu_t}{\pi} \right\rangle_{\mathcal{H}} \\ &= - \left\langle S_{\mu_t,k}\nabla\log\frac{\mu_t}{\pi}, S_{\mu_t,k}\nabla\log\frac{\mu_t}{\pi} \right\rangle_{\mathcal{H}} \\ &= - \left\| S_{\mu_t,k}\nabla\log\frac{\mu_t}{\pi} \right\|_{\mathcal{H}}^2 := -I_{\text{stein}}(\mu|\pi). \end{aligned} \quad (39)$$

where  $I_{\text{stein}}(\mu|\pi)$  is called as the Stein–Fisher (SF) information (Duncan et al., 2023). It is known that the square root of the SF information corresponds to the KSD. Then it is natural to examine the inequality like LSI as follows:

$$\text{KL}(\mu|\pi) \leq c I_{\text{stein}}(\mu|\pi). \quad (40)$$

where  $c$  is some positive constant and this is called Stein log-Sobolev inequality (Duncan et al., 2023). Unfortunately, the condition of this inequality is less clear compared to the LSI and Duncan et al. (2023) showed that it might fail to hold this inequality in many practical models like the RBF kernel function with exponential tail of  $\pi$ .

Thus it is difficult to consider the linear convergence of KL divergence under the geometry of  $\mathcal{H}$ . See Liu (2017) and Duncan et al. (2023) for a detailed discussion of the geometry of SVGD.

## C PROOFS OF THEORIES IN SECTION 4

### C.1 PROOF OF THEOREM 1

Under Assumptions 1 and 3-5, we have the following results from Theorem 3.2 in Salim et al. (2022):

$$\text{KL}(\mu_{t+1}|\pi) \leq \text{KL}(\mu_t|\pi) - \gamma_t \left( 1 - \frac{\gamma_t B^2(\alpha^2 + L)}{2} \right) I_{\text{stein}}(\mu_t|\pi), \quad (41)$$

where  $\eta_t := \gamma_t \left(1 - \frac{\gamma_t B^2(\alpha^2 + L)}{2}\right)$ . Substituting  $g_{\mu_t}(x) = P_{\mu_t, k} \nabla \log \frac{\mu_t}{\pi}$  into Eq. (9) yields

$$-\eta_t I_{\text{stein}}(\mu_t | \pi) \leq -\epsilon_t \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 + \eta_t \delta_t. \quad (42)$$

To show that the above inequality exists and to evaluate  $\epsilon_t$  and  $\delta_t$ , we study the following equality obtained via Eq. (42):

$$\epsilon_t \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 - \eta_t I_{\text{stein}}(\mu_t | \pi) = \eta_t \left\langle \nabla \log \frac{\mu_t}{\pi}, (\epsilon_t I - P_{\mu_t, k}) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)}. \quad (43)$$

As written in the main paper, we discuss the necessary conditions for the existence of Eq. (42) w.r.t.  $\eta_t$  under our final results. As can be seen from Theorem 1, we obtain  $\epsilon_t = c_0$  and  $\delta_t = 0$  through the proof that we explain later, where  $c_0$  is independent of  $t$ . In this case, as written in the main paper, it is necessary for  $\eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2$  to be uniformly upper bounded w.r.t.  $t$  to compensate for the convergence based on AGF. To satisfy this condition, from Lemma 2, when we set  $\gamma_t \leq \mathcal{O}(1/t^{2/3})$ , the second term is the order of  $\sum_t \gamma_t \leq \mathcal{O}(\eta^{1/3})$ . Then  $\left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 \leq \mathcal{O}(\eta^{2/3})$ . Then we can easily find that  $\eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 \leq \mathcal{O}(1)$  with respect to  $t$ .

From here, we analyze Eq. (43) by using the spectral decomposition of the kernel operator. Since a Hilbert–Schmidt operator  $P_{\mu, k}$  is compact and self-adjoint and we use the real-valued kernel function, we can decompose  $P_{\mu, k}$  by spectral theorem. Thus, we have, for all  $i$ ,

$$P_{\mu, k} \phi_i = \lambda_i \phi_i, \quad (44)$$

where  $\phi_i \in L^2(\mu)$  represents an eigenfunction that satisfies a complete orthonormal system (CONS), and  $\lambda_i$  is an eigenvalue corresponding to  $\phi_i$ . Even if these eigenvalues are ordered as  $\lambda_1 \geq \lambda_2 \geq \dots > 0$ , it does not compromise generality. Moreover, the kernel function can be decomposed into  $k(x, y) = \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(y)$ , where the convergence of this infinite series holds in the norm of  $\|\cdot\|_{L^2(\mu)}$ .

From the spectral theorem, all the eigenvalues of a positive-definite kernel function are positive real values and their multiplicity (the dimension of the eigenspace) is finite. Under Assumption 4 (ISPD assumption), a kernel function is dense in  $L^2(\mu)$  (Steinwart & Christmann, 2008; Sriperumbudur et al., 2011; Carmeli et al., 2010). Therefore, the eigenfunctions  $\{\phi_n\}$  are CONS in  $L^2(\mu)$ . Then, for any  $f \in L^2(\mu)$ , we have

$$f = \sum_{i=1}^{\infty} \langle f, \phi_i \rangle_{L^2(\mu)} \phi_i. \quad (45)$$

We should note that the above discussions overly simplify the eigenvalues and eigenvectors for vector functions because we treat the kernel function as the vector-valued one. As we mentioned in the main paper, since  $f = (f_1, \dots, f_d) \in L^2(\mu)$  is the vector valued function, each  $f_1, \dots, f_d \in L_0^2(\mu)$  are measurable function  $f_1 : \mathcal{X} \rightarrow \mathbb{R}$  with  $\int f_1^2 d\mu < \infty$ . Abusing the notation,  $P_{\mu, k} : L_0^2(\mu) \rightarrow \mathcal{H}_0$ , which projects the scalar function to  $\mathcal{H}_0$ . In this case, the eigenvalues and eigenvectors are given as

$$P_{\mu, k} \psi_i = \lambda_i \psi_i, \quad (46)$$

where  $\psi_i \in L_0^2(\mu)$  is an eigenfunction which is a scalar value function. When considering  $P_{\mu, k} : L^2(\mu) \rightarrow \mathcal{H}$ , the operator  $P_{\mu, k} \phi$  is regarded as the elementwise projection defined as  $P_{\mu, k} \phi_i = (\int k(x, y) \phi_1(y) d\mu(y), \dots, \int k(x, y) \phi_d(y) d\mu(y))$ . Since the eigenfunctions  $\{\psi_i\}$  are CONS in  $L_0^2(\mu)$  and the vector functions are in  $L^2(\mu)$ , we focus on each dimension and apply a spectral decomposition, that is,

$$f = (f_1, \dots, f_d) = \left( \sum_{i=1}^{\infty} \langle f_1, \psi_i \rangle_{L_0^2(\mu)} \psi_i, \dots, \sum_{i=1}^{\infty} \langle f_d, \psi_i \rangle_{L_0^2(\mu)} \psi_i \right). \quad (47)$$

By applying the kernel projection, we have

$$P_{\mu, k} f = (P_{\mu, k} f_1, \dots, P_{\mu, k} f_d) = \left( \sum_{i=1}^{\infty} \lambda_i \langle f_1, \psi_i \rangle_{L_0^2(\mu)} \psi_i, \dots, \sum_{i=1}^{\infty} \lambda_i \langle f_d, \psi_i \rangle_{L_0^2(\mu)} \psi_i \right). \quad (48)$$

The above equality corresponds to the setting when  $\phi_i = (\psi_i, \dots, \psi_i) \in L^2(\mu)$ . Then, for any  $f \in L^2(\mu)$ , we have

$$f = \sum_{i=1}^{\infty} \langle f, \phi_i \rangle_{L^2(\mu)} \phi_i, \quad (49)$$

and

$$P_{\mu,k} f = \sum_{i=1}^{\infty} \lambda_i \langle f, \phi_i \rangle_{L^2(\mu)} \phi_i. \quad (50)$$

From the above equalities, we can be seen as the same calculation for  $P_{\mu,t} : L_0^2(\mu) \rightarrow \mathcal{H}_0$  and  $P_{\mu,t} : L^2(\mu) \rightarrow \mathcal{H}$ .

In this paper, for simplicity, we do not work on  $P_{\mu,t} : L_0^2(\mu) \rightarrow \mathcal{H}_0$  with eigenvectors  $\{\psi_i\}$ , but work on  $P_{\mu,t} : L^2(\mu) \rightarrow \mathcal{H}$  with eigenvectors  $\{\phi_i\}$  as shown in Eq. (44). For completeness, we remark the norm calculation as follows:

$$\|f\|_{L^2(\mu)}^2 = \sum_{j=1}^d \|f_j\|_{L_0^2(\mu)}^2 = \sum_{j=1}^d \sum_{i=1}^{\infty} \langle f_j, \psi_i \rangle_{L_0^2(\mu)}^2 = \sum_{i=1}^{\infty} \sum_{j=1}^d \langle f_j, \psi_i \rangle_{L_0^2(\mu)}^2, \quad (51)$$

where  $\sum_{j=1}^d \langle f_j, \psi_i \rangle_{L_0^2(\mu)}^2$  is the multiply  $\psi_i$  in an elementwise way and summing it up. Based on the above, we have

$$\|P_{\mu,k} f\|_{L^2(\mu)}^2 = \sum_{j=1}^d \|P_{\mu,k} f_j\|_{L_0^2(\mu)}^2 = \sum_{j=1}^d \sum_{i=1}^{\infty} \lambda_i \langle f_j, \psi_i \rangle_{L_0^2(\mu)}^2 = \sum_{i=1}^{\infty} \lambda_i \sum_{j=1}^d \langle f_j, \psi_i \rangle_{L_0^2(\mu)}^2. \quad (52)$$

By using the spectral decomposition and the Hilbert–Schmidt theorem, we obtain

$$k(x, y) = \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(y), \quad (53)$$

where the convergence of this infinite series holds in the norm of  $\|\cdot\|_{L^2(\mu)}$ . If  $\mathcal{X}$  is a compact space, then from Mercer’s theorem, the convergence of the above infinite series holds absolutely and uniformly. However, we *do not* assume that  $\mathcal{X}$  is compact.

Let us show that, in our analysis, it is sufficient to study the convergence of  $k(x, y) = \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(y)$  in  $L^2(\mu)$  since our analysis is based on the norm of  $\|\cdot\|_{L^2(\mu)}$ . From Eq. (43), we have

$$\begin{aligned} & \eta_t \left\langle \nabla \log \frac{\mu_t}{\pi}, (\epsilon_t I - P_{\mu_t,k}) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)} \\ & \leq \eta_t \left\langle \nabla \log \frac{\mu_t}{\pi}, \left( \epsilon_t I - \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(\cdot) + \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(\cdot) - P_{\mu_t,k} \right) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)} \\ & \leq \eta_t \left\langle \nabla \log \frac{\mu_t}{\pi}, \left( \epsilon_t I - \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(\cdot) \right) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)} \\ & \quad + \eta_t \left\langle \nabla \log \frac{\mu_t}{\pi}, \left( \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(\cdot) - P_{\mu_t,k} \right) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)}, \quad (54) \end{aligned}$$

and the last term can be bounded by using the Cauchy–Schwartz inequality as follows:

$$\begin{aligned} & \left\langle \nabla \log \frac{\mu_t}{\pi}, \left( \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(\cdot) - P_{\mu_t,k} \right) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)} \\ & \leq \left\| \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(\cdot) - P_{\mu_t,k} \right\|_{\text{op}} \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2. \quad (55) \end{aligned}$$



Eq. (55) can be arbitrarily small since  $\|\nabla \log \frac{\mu_t}{\pi}\|_{L^2(\mu_t)}^2$  is bounded (Lemma 2), the operator norm is bounded by HS norm, and property of the convergence in  $L^2(\mu)$  norm. By setting  $\|\sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(\cdot) - P_{\mu_t,k}\|_{\text{op}} = \epsilon_0$ , we have

$$\begin{aligned} & \eta_t \left\langle \nabla \log \frac{\mu_t}{\pi}, (\epsilon_t I - P_{\mu_t,k}) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)} \\ & \leq \eta_t \left\langle \nabla \log \frac{\mu_t}{\pi}, \left( \epsilon_t I - \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(\cdot) \right) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)} + \epsilon_0 \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2, \end{aligned} \quad (56)$$

where  $\epsilon_0$  is arbitrarily small and negligible for the convergence. Thus, it is sufficient to focus on the convergence of the spectral decomposition of  $P_{\mu,k}$  in  $L^2(\mu)$  norm.

Defining  $v_t := \nabla \log \frac{\mu_t}{\pi}$  for simplicity in notation, we can obtain  $v_t = \sum_{i=1}^{\infty} \langle v_t, \phi_i \rangle_{L^2(\mu)} \phi_i$  and  $P_{\mu,k} v_t = \sum_{i=1}^{\infty} \lambda_i \langle v_t, \phi_i \rangle_{L^2(\mu)} \phi_i$  because the kernel function is dense in  $L^2(\mu)$  and thus its eigenvectors are complete. Using these, we obtain

$$\epsilon_t \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 - \eta_t I_{\text{stein}}(\mu_t | \pi) = \eta_t \sum_{i=1}^{\infty} (\epsilon_t - \lambda_i) \langle v_t, \phi_i \rangle_{L^2(\mu_t)}^2. \quad (57)$$

By setting  $\epsilon_t$  sufficiently small, there exists a index  $1 < j$  such that  $\lambda_j > \epsilon_t > \lambda_{j+1}$ . Hence, by regularizing  $\{\langle v_t, \phi_i \rangle_{L^2(\mu_t)}^2\}_{i=1}^{\infty}$ , we can render the left-hand side of Eq. (57) negative. For that purpose, we focus on the RKHS associated with  $k$  given as

$$\mathcal{H} = \left\{ f \in L^2(\mu) \mid f = \sum_{k=1}^{\infty} a_k \phi_k, \sum_{i=1}^{\infty} \frac{\|a_i\|^2}{\lambda_i}, a_i \in \mathbb{R} \right\}, \quad (58)$$

where  $\mathcal{H}$  is dense in  $L^2(\mu)$ . Thanks to this property of  $\mathcal{H}$ , there exists a function  $v^{(l)} \in \mathcal{H}$  such that the sequence of  $v^{(l)} \rightarrow v$  as  $l \rightarrow \infty$  in  $L^2(\mu)$  norm. We express such  $v^{(l)}$  as

$$v_t^{(l)} = \sum_{i=1}^{\infty} b_i^{(l)} \phi_i \in \mathcal{H} \quad (59)$$

Note that we have that  $\sum_{i=1}^{\infty} \frac{\|b_i^{(l)}\|^2}{\lambda_i} < \infty$  since  $v_t^{(l)} \in \mathcal{H}$ .

We consider replacing  $v_t$  by  $v_t^{(l)}$  to control the coefficient of  $\{\phi_n\}$ . By definition, for any  $\epsilon > 0$ , we can obtain  $\|v_t - v_t^{(l)}\|_{L^2(\mu_t)} < \epsilon$  by choosing sufficiently large  $l$ . According to this fact, we obtain

$$\begin{aligned} & \left\langle \nabla \log \frac{\mu_t}{\pi}, (\epsilon_t I - P_{\mu_t,k}) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)} \\ & = \left\langle v_t - v_t^{(l)} + v_t^{(l)}, (\epsilon_t I - P_{\mu_t,k}) v_t - v_t^{(l)} + v_t^{(l)} \right\rangle_{L^2(\mu_t)} \\ & = \left\langle v_t^{(l)}, (\epsilon_t I - P_{\mu_t,k}) v_t^{(l)} \right\rangle_{L^2(\mu_t)} + 2 \left\langle v_t^{(l)}, (\epsilon_t I - P_{\mu_t,k}) v_t - v_t^{(l)} \right\rangle_{L^2(\mu_t)} \\ & \quad + \left\langle v_t - v_t^{(l)}, (\epsilon_t I - P_{\mu_t,k}) v_t - v_t^{(l)} \right\rangle_{L^2(\mu_t)} \\ & = \left\langle v_t^{(l)}, (\epsilon_t I - P_{\mu_t,k}) v_t^{(l)} \right\rangle_{L^2(\mu_t)} + 2\epsilon_t \epsilon \|v_t\|_{L^2(\mu_t)}^2 + \epsilon_t \epsilon^2 \end{aligned} \quad (60)$$

According to Lemma 2,  $\|v_t\|_{L^2(\mu_t)}^2$  is bounded for  $t < \infty$ , and the second and the third term can be arbitrarily small by  $\epsilon$ . Thus, we now focus on the term  $\left\langle v_t^{(l)}, (\epsilon_t I - P_{\mu_t,k}) v_t^{(l)} \right\rangle_{L^2(\mu_t)}$  to analyze  $\left\langle \nabla \log \frac{\mu_t}{\pi}, (\epsilon_t I - P_{\mu_t,k}) \nabla \log \frac{\mu_t}{\pi} \right\rangle_{L^2(\mu_t)}$ .

From the above arguments and the fact in Eq. (59), we can be rewritten Eq. (57) as

$$\epsilon_t \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 - \eta_t I_{\text{stein}}(\mu_t | \pi) = \eta_t \sum_{i=1}^{\infty} (\epsilon_t - \lambda_i) (b_i^{(l)})^2 + o(\epsilon), \quad (61)$$

where the residual term  $o(\epsilon)$ , which comes from Eq. (60), can be arbitrarily small and thus it is trivial in our discussion.

Next, we consider to choose appropriate  $\epsilon_t$ . Recall that, by setting  $\epsilon_t$  sufficiently small, there exists a index  $1 < j$  such that  $\lambda_j > \epsilon_t > \lambda_{j+1}$ . Then, we can expand the right-hand side term of Eq. (61) as

$$\sum_{i=1}^{\infty} (\epsilon_t - \lambda_i) (b_i^{(l)})^2 = - \underbrace{\sum_{i=1}^j (\lambda_i - \epsilon_t) (b_i^{(l)})^2}_{=:A} + \underbrace{\sum_{i=j+1}^{\infty} (\epsilon_t - \lambda_i) (b_i^{(l)})^2}_{=:B}, \quad (62)$$

where  $A, B > 0$  by definition. We now show that there exists  $\epsilon_t > 0$  such that

$$\sum_{i=1}^{\infty} (\epsilon_t - \lambda_i) (b_i^{(l)})^2 = -A + B < 0. \quad (63)$$

This can be easily confirmed by the definition of  $b_i^{(l)}$ . Since  $v_t^{(l)} \in \mathcal{H}$ , we have  $\sum_{i=1}^{\infty} \frac{\|b_i^{(l)}\|^2}{\lambda_i} < \infty$ . This implies that  $\sup_{i \leq m} (b_m^{(l)})^2$  goes to 0 at least as the same order as  $\lambda_i$ . By setting  $\epsilon_t$  sufficiently small, the corresponding index  $j$  becomes large and then  $\{(b_i^{(l)})^2\}_{i \geq j+1}$  becomes small. With this procedure, we obtain

$$\epsilon_t \eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 - \eta_t I_{\text{stein}}(\mu_t | \pi) < 0. \quad (64)$$

This implies that  $\delta_t = 0$  in our AGF.

Finally, we show that we can choose  $\epsilon_t$  so as to be independent of  $t$ . We would like to find the **smallest** index  $m \in \mathbb{N}$  such that

$$\sum_{i=1}^m (b_i^{(l)})^2 > \sum_{i=m+1}^{\infty} (b_i^{(l)})^2, \quad (65)$$

holds. **Thus, by definition**

$$\sum_{i=1}^{m-1} (b_i^{(l)})^2 \leq \sum_{i=m}^{\infty} (b_i^{(l)})^2, \quad (66)$$

**holds.** Actually, we can find such index if we pick up sufficiently large  $m$  because  $\sup_{i \leq m} (b_m^{(l)})^2$  goes to 0 at least as the same order as  $\lambda_i$  from  $\sum_{i=1}^{\infty} \frac{\|b_i^{(l)}\|^2}{\lambda_i} < \infty$ . Under such  $m$ , we obtain

$$\begin{aligned} \sum_{i=1}^{\infty} (\epsilon_t - \lambda_i) (b_i^{(l)})^2 &= - \sum_{i=1}^m (\lambda_i - \epsilon_t) (b_i^{(l)})^2 + \sum_{i=m+1}^{\infty} (\epsilon_t - \lambda_i) (b_i^{(l)})^2 \\ &\leq - \sum_{i=1}^m \lambda_i (b_i^{(l)})^2 + 2\epsilon_t \sum_{i=1}^m (b_i^{(l)})^2 \\ &\leq (-\lambda_m + 2\epsilon_t) \sum_{i=1}^m (b_i^{(l)})^2. \end{aligned} \quad (67)$$

By setting  $\epsilon_t \leq \lambda_m/2$  and taking  $o(\epsilon)(>0)$  into account, we have the relationship as in Eq. (65).

If  $m$  **monotonically increases w.r.t. some sequence of  $t$** , the  $m$ -th eigenvalue  $\lambda_m$  would become smaller as  $t$  increases, leading to a decreasing upper bound of  $\epsilon_t \leq \lambda_m/2$  as  $t$  increases. To ensure convergence with the LSI, we must demonstrate that  $m$  satisfying Eq. (65) and (66) does not monotonically increase w.r.t. the iteration  $t$ . This can be shown through a proof by contradiction as follows.

We assume that there exists a subsequence of  $t$ ,  $\{t_k\}_{k=1}^{\infty}$ , such that  $m$  satisfying Eq. (65) and (66) **monotonically increases**. Also, for any  $k$ , let  $m$  corresponding to  $t_k$  be denoted as  $m_{t_k}$ . Furthermore,

we define the right-hand side of Eq. (66) as  $S_m$ , where  $S_m := \sum_{i=m}^{\infty} (b_i^{(l)})^2$ . From the definition of the index  $m$ , we have  $m_{t_1} < m_{t_2} < \dots < m_{t_k} < \dots$ , for all  $k$ . In this case, there exists a sequence in  $S_{m_{t_1}}, \dots, S_{m_{t_k}}, \dots$  that goes to 0 because  $\sup_{i \leq m} (b_m^{(l)})^2$  goes to 0 at least as the same order as  $\lambda_i$  from  $\sum_{i=1}^{\infty} \frac{\|b_i^{(l)}\|^2}{\lambda_i} < \infty$ . However, a contradiction arises since  $S_{m_{t_k}}$  approaching 0 does not satisfy Eq. (66).

This contradiction suggests that  $m$  might increase for some range of  $t$ , but it is upper bounded w.r.t.  $t$ . This allows us to identify a largest  $m$  as  $m'$  that does not depend on  $t$ . Recalling Assumption 6, which provides a strictly positive lower bound for  $\lambda_i$  denoted as  $\hat{\lambda}_{m'}$ , we see that this lower bound is also independent of  $t$ . Combining the above discussions, we can establish that  $\epsilon_t$  can be upper-bounded by  $\hat{\lambda}_{m'}$ , meaning that  $\epsilon_t \leq \hat{\lambda}_{m'}$ . This implies that we can set a positive constant  $c_0$  as an upper bound for  $\epsilon_t$ , and this constant is independent of  $t$ .

In conclusion, we have

$$\text{KL}(\mu_{t+1}|\pi) \leq \text{KL}(\mu_t|\pi) - c_0\eta_t \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)}^2 \leq (1 - \tilde{c}_0\eta_t) \text{KL}(\mu_t|\pi), \quad (68)$$

where we used the LSI and summarized the LSI constants in  $\tilde{c}_0$ . By recursively applying the above inequality, we obtain

$$\text{KL}(\mu_T|\pi) \leq \prod_{t=1}^{T-1} (1 - c_0\eta_t) \text{KL}(\mu_0|\pi). \quad (69)$$

We final note that if we have  $\eta_t = c_1/T$ ,

$$\prod_{t=1}^T \left(1 - \frac{c_0 c_1}{t}\right) \leq \left(1 - \frac{1}{T} \sum_{t=1}^T \frac{c_0 c_1}{t}\right)^T \leq e^{-\sum_{t=1}^T \frac{c_0 c_1}{t}} \leq e^{-c_0 c_1 \log T} \leq \frac{1}{T^{c_0 c_1}}. \quad (70)$$

This concludes the proof.

## C.2 PROOF OF LEMMA 2

For simplicity, we express  $\eta_t$  as  $\gamma$  is this proof. Let us define the mapping  $\rho_s := \phi_s \# \mu_t$  for  $s \in [0, \gamma]$  with  $\phi_s := I - s v_t$  and  $v_t = \nabla \log \frac{\mu_t}{\pi}$ . Then, by the change of variable formula, we have

$$\rho_s = |J\phi_s(\phi_s^{-1}(x))|^{-1} \mu_t(\phi_s^{-1}(x)), \quad (71)$$

where  $\rho_\gamma = \mu_{t+1}$ .

Our goal is to bound the following equality:

$$\begin{aligned} \left\| \nabla \log \frac{\rho_s}{\pi} \right\|_{L^2(\rho_s)}^2 &= \int \left\langle \nabla \log \frac{\rho_s}{\pi}(x), \nabla \log \frac{\rho_s}{\pi}(x) \right\rangle d\rho_s(x) \\ &= \int \left\langle \nabla \log \frac{\mu_t(x) |J\phi_s(x)|^{-1}}{\pi(\phi_s(x))}, \nabla \log \frac{\mu_t(x) |J\phi_s(x)|^{-1}}{\pi(\phi_s(x))} \right\rangle d\mu_t(x). \end{aligned} \quad (72)$$

First, we apply the mean value theorem (as known as Taylor expansion of order 1) to  $\psi(s) := \nabla \log \frac{\mu_t(x) |J\phi_s(x)|^{-1}}{\pi(\phi_s(x))}$  as a function of  $s$ . According to this theorem, there exists a constant  $c \in [0, \gamma]$  such that

$$\psi(s) = \psi(0) + \gamma \frac{d}{ds} \psi(s) \Big|_{s=c}. \quad (73)$$

This implies

$$\begin{aligned} \left\| \nabla \log \frac{\mu_{t+1}}{\pi} \right\|_{L^2(\rho_{t+1})} &= \left\| \nabla \log \frac{\rho_\gamma}{\pi} \right\|_{L^2(\rho_\gamma)} = \left\| \nabla \log \frac{\mu_t(x) |J\phi_\gamma(x)|^{-1}}{\pi(\phi_\gamma(x))} \right\|_{L^2(\mu_t)} \\ &\leq \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)} + \gamma \left\| \frac{d}{ds} \psi(s) \Big|_{s=c} \right\|_{L^2(\mu_t)}, \end{aligned} \quad (74)$$

where we used the triangle inequality. The second term in the right-hand side of Eq. (74) can be expressed as

$$\begin{aligned}
& \left\| \frac{d}{ds} \psi(s) \right\|_{L^2(\mu_t)} \\
&= \left\| -\frac{d}{ds} \nabla \log |I - sJv_t| + \frac{d}{ds} \nabla V(\phi_s(x)) \right\|_{L^2(\mu_t)} \\
&= \left\| \nabla \text{Tr} [(J\phi_s(x))^{-1} Jv_t(x)] + \nabla \langle \nabla V(\phi_s(x)), v_t(x) \rangle \right\|_{L^2(\mu_t)} \\
&= \left\| \nabla \sum_{ij} ((J\phi_s(x))^{-1})_{ij} (Jv_t(x))_{ji} + \nabla^2 V(\phi_s(x)) v_t(x) + Jv_t(x) \nabla V(\phi_s(x)) \right\|_{L^2(\mu_t)}, \quad (75)
\end{aligned}$$

where  $\nabla = (\partial_1, \dots, \partial_d)^\top$  and  $v_t = (v_1, \dots, v_d)^\top$ . For derivation, we first swap the time derivative and the gradient, and using the time derivative of  $\Psi(s)$  that is shown in Appendix B of Salim et al. (2022) and Sun et al. (2023).

To bound Eq. (75), we use the following existing bounds. From Lemma C.1. in Salim et al. (2022), we have

$$\|v_t(x)\| \leq B^2 \left( 1 + 2L \sqrt{\frac{2\text{KL}(\mu_0|\pi)}{C_{\text{LS}}}} + LW_2(\mu_0, \delta_x^*) \right) =: C_1, \quad (76)$$

where  $x^* = \arg \min_{x \in \mathcal{X}} V(x)$ . In Salim et al. (2022), they assumed that  $T_1$  inequality holds for  $\pi$ , whereas we assumed that the LSI holds. Since  $T_1$  inequality is satisfied if the LSI is available, their bound also holds in our setting. Also, from Appendix B in Salim et al. (2022), we have

$$\|Jv_t(x)\|_{\text{HS}}^2 \leq C_1^2. \quad (77)$$

Moreover, from the proof of Lemma C.1. in Salim et al. (2022), for any  $\mu \in \mathcal{P}_2(\mathcal{X})$ , we have

$$\|\nabla V\|_{L^2(\mu)} \leq 2L \sqrt{\frac{2\text{KL}(\mu_0|\pi)}{C_{\text{LS}}}} + LW_2(\mu_0, \delta_x^*) =: C_2. \quad (78)$$

From Appendix B in Salim et al. (2022), we have

$$\|(J\phi_s(x))^{-1}\|_{\text{HS}} \leq \alpha. \quad (79)$$

While the upper bound of  $\|(J\phi_s(x))^{-1}\|_{\text{op}}$  is presented above, the bound of  $\|(J\phi_s(x))^{-2}\|_{\text{HS}}$  can be derived almost similar way as shown in Appendix B in Salim et al. (2022) as follows.

$$\|(J\phi_s(x))^{-2}\|_{\text{HS}} \leq \alpha^2. \quad (80)$$

In addition to these upper bounds (Eqs. (76)-(80)), we use the following fact:

$$\begin{aligned}
\sum_{i,j,k=1}^d (\partial_i \partial_j v_k(x))^2 &= \sum_{i,j,k=1}^d \langle \partial_i \partial_j k(x, \cdot), v_k \rangle_{\mathcal{H}_0}^2 = \sum_{i,j,k=1}^d \|\partial_i \partial_j k(x, \cdot)\|_{\mathcal{H}_0}^2 \|v_k\|_{\mathcal{H}_0}^2 \\
&= \sum_{i,j,k=1}^d \|\partial_i \partial_j k(x, \cdot)\|_{\mathcal{H}_0}^2 \|v_t\|_{\mathcal{H}}^2 \\
&= B^2 \|v_t\|_{\mathcal{H}}^2 \leq C_1, \quad (81)
\end{aligned}$$

where we used Assumption 3 in the last line.

To upper bound Eq. (75), we apply triangle inequality. Then we focus on the square of the first term in Eq. (75),

$$\begin{aligned}
& \sum_k (\partial_k (\sum_{ij} ((J\phi_s(x))^{-1})_{ij} (Jv_t(x))_{ji}))^2 \\
&= \sum_k ((\sum_{ij} ((J\phi_s(x))^{-2})_{ij} \partial_k (-tJv_t(x))_{ij} (Jv_t(x))_{ji}) + (\sum_{ij} ((J\phi_s(x))^{-1})_{ij} \partial_k (Jv_t(x))_{ji}))^2 \\
&\leq \sum_k ((\sum_{ij} \alpha^2 |\partial_k (-tJv_t(x))_{ij}| C_1^2 + (\alpha |\partial_k (Jv_t(x))_{ji}|))^2 \\
&\leq \sum_k ((\sum_{ij} (\alpha + t\alpha^2 C_1^2) |\partial_k (Jv_t(x))_{ij}|))^2 \\
&\leq (\alpha + t\alpha^2 C_1^2)^2 d^2 \sum_{ijk} (\partial_k (Jv_t(x))_{ij})^2 \\
&\leq (\alpha + t\alpha^2 C_1^2)^2 d^2 B^2 \leq (\alpha + \gamma\alpha^2 C_1^2)^2 d^2 B^2 := D_1^2
\end{aligned} \tag{82}$$

where  $D_1 > 0$  and it is not depend on  $t$ .

Now we are ready for bounding Eq. (75). First, the second and third term in Eq. (75) can be upper bounded by

$$\begin{aligned}
& \|\nabla^2 V(\phi_s(x))v_t(x) + Jv_t(x)\nabla V(\phi_s(x))\|_{L^2(\mu_t)} \\
&\leq \|\nabla^2 V(\phi_s(x))v_t(x)\|_{L^2(\mu_t)} + \|Jv_t(x)\nabla V(\phi_s(x))\|_{L^2(\mu_t)} \\
&\leq \|\nabla^2 V(\phi_s(x))\|_{\text{op}} \|v_t(x)\|_{L^2(\mu_t)} + \|Jv_t(x)\|_{\text{op}} \|\nabla V(\phi_s(x))\|_{L^2(\mu_t)} \\
&\leq LC_1 + C_1 C_2.
\end{aligned} \tag{83}$$

Substituting the upper bounds in Eqs. (82) and (83), we have

$$\begin{aligned}
\left\| \frac{d}{ds} \psi(s) \right\|_{L^2(\mu_t)} &= \left\| -\frac{d}{ds} \nabla \log |I - sJv_t| + \frac{d}{ds} \nabla V(\phi_s(x)) \right\|_{L^2(\mu_t)} \\
&\leq D_1 + LC_1 + C_1 C_2 =: C_3.
\end{aligned} \tag{84}$$

Thus,

$$\left\| \nabla \log \frac{\mu_{t+1}}{\pi} \right\|_{L^2(\mu_{t+1})} = \left\| \nabla \log \frac{\rho_\gamma}{\pi} \right\|_{L^2(\rho_\gamma)} \leq \left\| \nabla \log \frac{\mu_t}{\pi} \right\|_{L^2(\mu_t)} + \gamma C_3 \tag{85}$$

By induction, we have the following results:

$$\left\| \nabla \log \frac{\mu_T}{\pi} \right\|_{L^2(\mu_T)} \leq \left\| \nabla \log \frac{\mu_0}{\pi} \right\|_{L^2(\mu_0)} + \sum_{t=0}^{T-1} \gamma_t C_3. \tag{86}$$

This completes the proof.

## D DETAILS OF EXPERIMENTAL SETTINGS

We set the target distribution as  $p(x) = \mathcal{N}(x|\mu^*, \Sigma^*)$  with  $\mu^* = [1, 1]$  and  $\Sigma^* = \text{diag}(1, 1)$ , where  $\text{diag}$  is a diagonal matrix. For the Gaussian mixture distribution, we set  $p(x) = \frac{2}{3}\mathcal{N}(x|\mu_1^*, \Sigma_1^*) + \frac{1}{3}\mathcal{N}(x|\mu_2^*, \Sigma_2^*)$  with  $\mu_1^* = [2, -2]$ ,  $\mu_2^* = [-2, 2]$ ,  $\Sigma_1^* = \Sigma_2^* = \text{diag}(1, 1)$ , which is the extension of the experimental settings in Liu & Wang (2016) for the two-dimensional setting. We generated the initial particles from  $\mathcal{N}(x|\mu_0, \Sigma_0)$  with  $\mu_0 = [0, 0]$  and  $\Sigma_0 = \text{diag}(1, 1)$  or  $\mu_0 = [-5, -5]$  and  $\Sigma_0 = \text{diag}(1, 1)$  for the Gaussian and the Gaussian mixture experiments, respectively.

We adopted the RBF kernel  $k(x, y) = \exp(-\frac{1}{h}\|x - x'\|_2^2)$ , which is commonly used in practice and satisfies the assumptions in Section 4. The bandwidth  $h$  was selected by the median trick, i.e.,  $\text{med}^2 / \log n$  as in Liu & Wang (2016), where  $\text{med}$  is the median of the pairwise distance between the current particles.



As we stated in Section 5, we simply set the decaying step size  $\gamma_t = 1/(1 + t^\beta)$  ( $= \mathcal{O}(1/t^\beta)$ ) suggested by Theorem 1 and did not use the Adagrad-based stepsize, which is adopted in related studies such as Korba et al. (2021) and others. We set the initial stepsize as  $\gamma_0 = 0.01$  for all experiments. We evaluated the KL divergence:  $\text{KL}(\mu_T|\pi)$  and the cumulative mean of KSD:  $\frac{1}{T} \sum_{t=1}^T I_{\text{stein}}(\mu_t|\pi)$ , which are theoretically guaranteed sub-linear convergence.

We conducted our experiments based on the above settings using  $\{5, 10, 100, 1000\}$  particles.

## E ADDITIONAL EXPERIMENTAL RESULTS

In this section, we provide the additional experimental results.

We confirmed in Section 5 that SVGD with the RBF kernel tends to achieve sub-linear convergence both in  $\text{KL}(\mu_T|\pi)$  and in  $\frac{1}{T} \sum_{t=1}^T I_{\text{stein}}(\mu_t|\pi)$  (see Figures 1 and 2). As for the Gaussian mixture settings, we can observe the same behavior (see Figures 3 and 4), which also supports Theorem 1.

As discussed in Appendix A, the bias in the KL divergence persists as we increase the value of  $T$  because we utilized a finite number of particles, leading to  $\delta_t \neq 0$  in AGF. Such a bias can be reduced by increasing the number of particles increases (see Figures 1-4, 5, and 6). On the flip side, even in the Gaussian mixture experiments, using a large number of particles results in slower convergence for both the KSD and KL divergence.

This phenomenon can be attributed to the existence of extremely small eigenvalues of  $P_{\mu,k}$  when a larger number of particles is used, as the eigenvalues of the RBF kernel decay exponentially fast (Wainwright, 2019). To confirm this fact, we measured the eigenvalues of the Gram matrix obtained from the RBF kernel function at three points: the initial stage of learning ( $t = 1$ ), the midpoint ( $t = 5 * 10^4$ ), and the final stage ( $t = 10^5$ ). We summarize these results in Figures 7 and 8. We can see that the exponential decay of the eigenvalues tends to be occurring as the number of particles increases.

Assumption 6, which states that the eigenvalues have a strictly positive lower bound and upper bound that are independent of  $t$ , is crucial for showing the sub-linear convergence of SVGD in KL divergence under the setting of an infinite number of particles. Since it is difficult to theoretically show this fact, we instead conducted numerical experiments to confirm that the dependence of the upper bound (maximum value) and lower bound (minimum value) of the eigenvalues on the variable  $t$  diminishes as the number of particles increases. As a metric for measuring time-dependence, we employed the following growth rate for the time interval  $[t_1, t_2]$  ( $t_2 > t_1$ ):

$$\frac{|\tilde{\lambda}_{t_2} - \tilde{\lambda}_{t_1}|}{t_2 - t_1}, \quad \frac{|\bar{\lambda}_{t_2} - \bar{\lambda}_{t_1}|}{t_2 - t_1},$$

where  $\{\tilde{\lambda}_{t_j}, \bar{\lambda}_{t_j}\}$  ( $j = 1, 2$ ) is the maximum and minimum value of the eigenvalues at iteration  $t_j$ . In the above metric, if the dependence of the eigenvalues for  $t$  is small, meaning that the changes in  $\{\tilde{\lambda}_{t_j}, \bar{\lambda}_{t_j}\}$  with the progression of  $t$  are small, then the value will be small. Furthermore, when the overall behavior indicates minimal fluctuation in this value throughout the training process, it signifies that the changes in the eigenvalues with the progression of  $t$  are small. This, in turn, reflects the small dependence of  $\{\tilde{\lambda}_{t_j}, \bar{\lambda}_{t_j}\}$  w.r.t.  $t$ .

We summarized the experimental results under the following three case:  $(t_1, t_2) = (0, 5 * 10^4)$  (refer to case (i)),  $(5 * 10^4, 10^5)$  (refer to case (ii)), and  $(0, 10^5)$  (refer to case (iii)) in Figures 9 and 10. To begin with, we can see that the change between the midpoint and endpoint of  $t$ , i.e., the case (ii), yields the smallest value for all particle settings. This suggests that during this period, SVGD is gradually approaching convergence, which seems to align with the results of  $\frac{1}{T} \sum_{t=1}^T I_{\text{stein}}(\mu_t|\pi)$  in Figures 1-4. On the other hand, in the low number of particles setting, there is a significant difference in the amount of change for each configuration, whereas when a large number of particles are used, this difference becomes small. This fact illustrates that with an increase in the number of particles, the fluctuations of  $\{\tilde{\lambda}_{t_j}, \bar{\lambda}_{t_j}\}$  w.r.t.  $t$  decrease. From these observations, it is expected that in the setting of an infinite number of particles, the dependence of the upper and lower bounds of the eigenvalues on the iteration  $t$  becomes significantly small, which implies the constant order of eigenvalues.

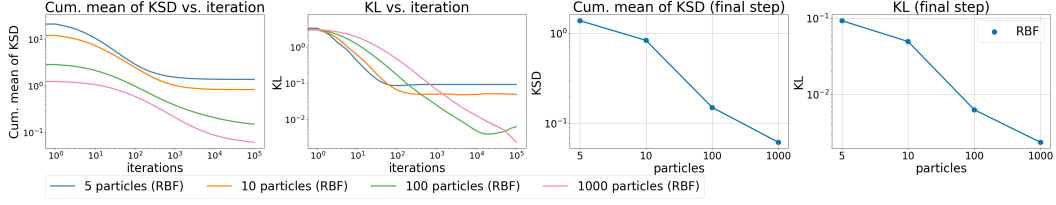


Figure 3: The convergence behavior in terms of  $KL(\mu_T|\pi)$  and  $\frac{1}{T} \sum_{t=1}^T I_{\text{stein}}(\mu_t|\pi)$  for all  $T$  under two-dimensional Gaussian mixture experiments ( $\beta = 0.67 \approx 2/3$ ).

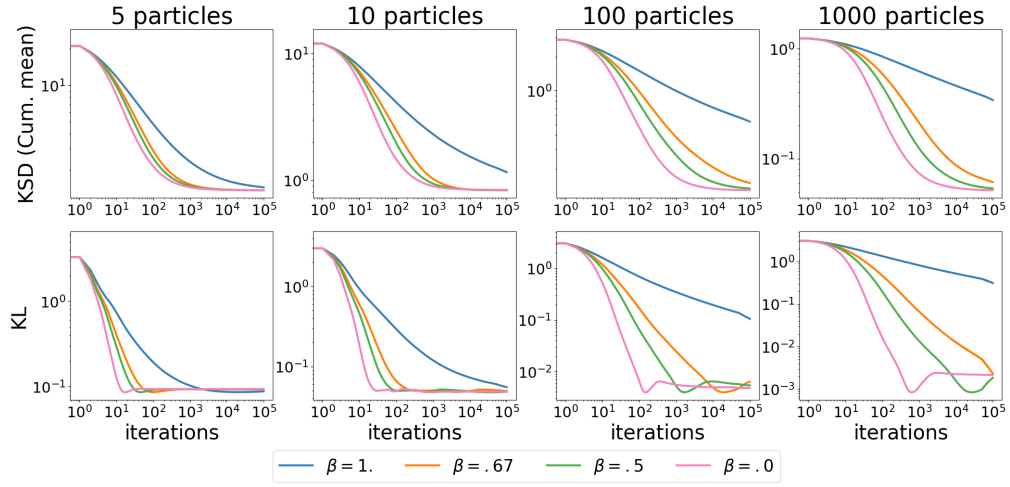


Figure 4: Convergence in  $KL(\mu_T|\pi)$  and  $\frac{1}{T} \sum_{t=1}^T I_{\text{stein}}(\mu_t|\pi)$  for all  $T$  under different particles and stepsize settings ( $\beta = \{0., 0.5, 0.67, 1.\}$ ).

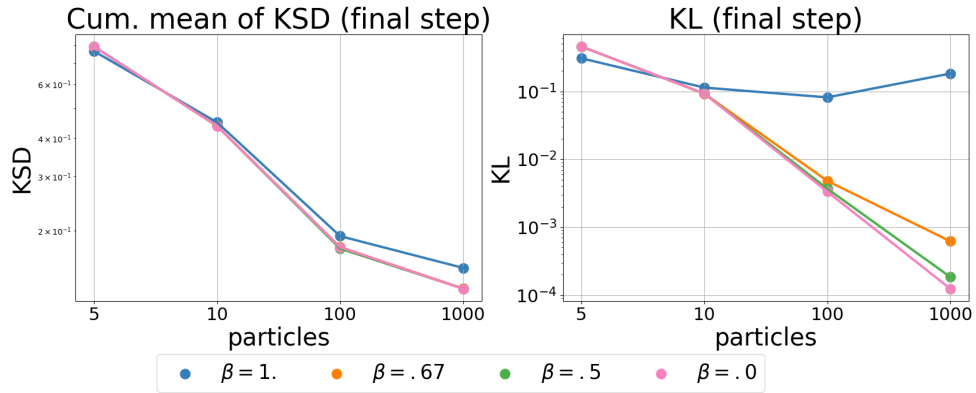


Figure 5: Change in convergence for variations in the order of stepsize with  $\beta = \{0., 0.5, 0.67, 1.\}$  under Gaussian distribution estimation experiments.

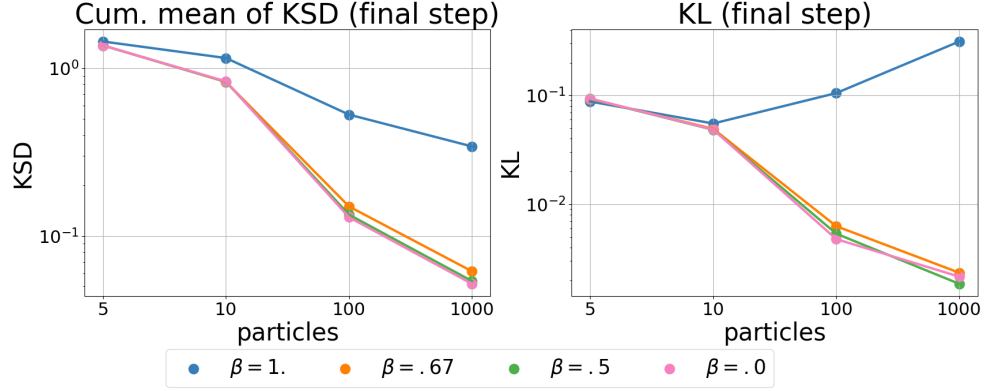


Figure 6: Change in convergence for variations in the order of stepsize with  $\beta = \{0., 0.5, 0.67, 1.\}$  under Gaussian distribution estimation experiments.

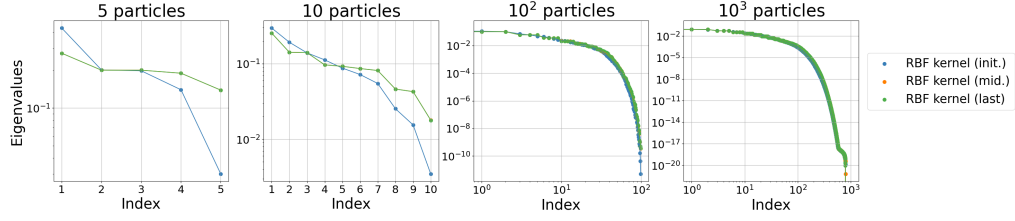


Figure 7: Eigenvalues of the Gram matrix in the two-dimensional Gaussian distribution experiments.

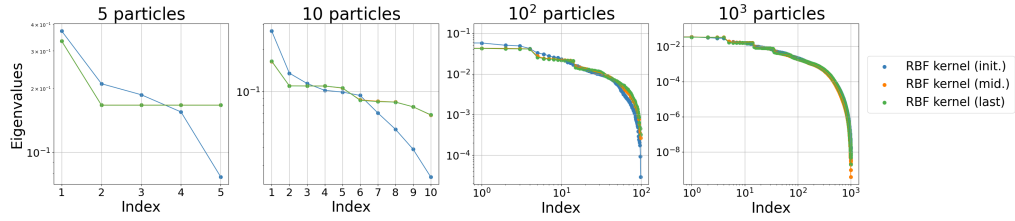


Figure 8: Eigenvalues of the Gram matrix in the two-dimensional Gaussian mixture experiments.

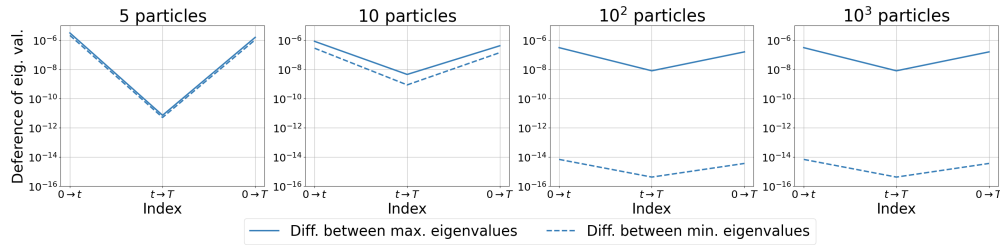


Figure 9: Difference between {maximum, minimum} eigenvalues of the Gram matrix in the two-dimensional Gaussian distribution experiments. In this figure,  $(t, T)$  represents  $(5 * 10^4, 10^5)$ .

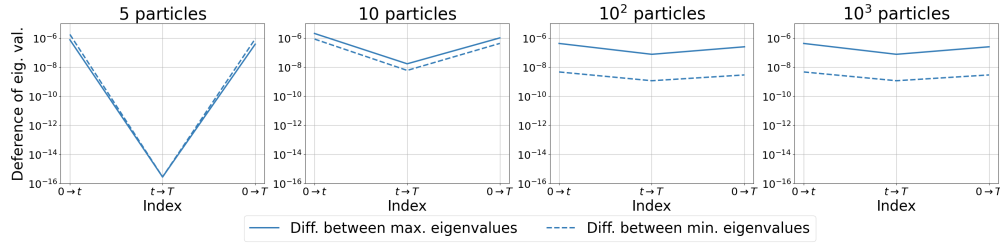


Figure 10: Difference between {maximum, minimum} eigenvalues of the Gram matrix in the two-dimensional Gaussian mixture experiments. In this figure,  $(t, T)$  represents  $(5 * 10^4, 10^5)$ .