

CONVERGENCE ASPECTS OF HYBRID KERNEL SVGD

Anonymous authors

Paper under double-blind review

ABSTRACT

Stein variational gradient descent (SVGD) is a particle based approximate inference algorithm. Many variants of SVGD have been proposed in recent years, including the hybrid kernel variant (h-SVGD), which has demonstrated promising results on image classification with deep neural network ensembles. In this paper, we demonstrate the ability of h-SVGD to alleviate variance collapse, a problem that SVGD is known to suffer from. Unlike other SVGD variants that alleviate variance collapse, h-SVGD does not incur additional computational cost, nor does it require the target density to factorise. We also develop the theory of h-SVGD by demonstrating the existence of a solution to the hybrid Stein partial differential equation. We highlight a special case in which h-SVGD is a kernelised Wasserstein gradient flow on a functional other than the Kullback-Leibler divergence, which is the functional describing the SVGD gradient flow. By characterising the fixed point in this special case, we show that h-SVGD does not converge to the target distribution in the mean field limit. Other theoretical results include a descent lemma and a large particle limit result. Despite the bias in the mean field limiting distribution, experiments demonstrate that h-SVGD remains competitive on high dimensional inference tasks whilst alleviating variance collapse.

1 INTRODUCTION

Stein variational gradient descent (SVGD) is a variational inference algorithm that generates samples from a target probability density (Liu & Wang, 2016). It has proven useful in many tasks in Bayesian inference and machine learning. SVGD evolves an interacting particle system until the particles resemble a sample from a target density. The dynamics of this system include a driving term that moves particles to regions of high probability, and a repulsive term that repels particles from one another. This repulsive term prevents particles from converging to the same mode. It has been shown that within a unit ball of a reproducing kernel Hilbert space (RKHS), the SVGD update direction optimally reduces the Kullback-Leibler (KL) divergence between the target density and the approximating density (Liu & Wang, 2016). The reproducing kernel of this RKHS appears in both the driving and repulsive terms, making the choice of kernel a key ingredient for SVGD.

SVGD is known to suffer from the curse of dimensionality through variance collapse (Wang et al., 2018; Ba et al., 2019) whereby the marginal variances of the particles underestimate the true marginal variances of the target in high dimensions. Capturing the variance of the posterior is critical for uncertainty quantification in Bayesian statistics since variance is often used as a measure of confidence in an estimate. Zhuo et al. (2018) explained that this phenomenon is due to the size of the repulsive term of the update direction scaling inversely with dimension. This enables the driving term to dominate in high dimensions, thereby forcing particles to converge to the mode(s) of the target. This insight suggests that strengthening the repulsive term in SVGD should lead to better variance estimation, an idea which we explore in Sections 3 and 4.

The theoretical properties of vanilla SVGD have been studied extensively. Liu (2017) showed that the empirical measure of the particles converges weakly to the target distribution. SVGD in the mean field regime has been described as a gradient flow on the KL divergence (Liu & Wang, 2016) and the chi-squared divergence (Chewi et al., 2020). Furthering this geometric point of view, Duncan et al. (2023) developed the Stein geometry along with its associated tangent spaces and geodesics, leading to guidelines for choosing kernels to improve convergence. The existence and uniqueness of the solution to the Stein partial differential equation (PDE) was established by Lu et al. (2019).

Various descent lemmas bounding the decrease in KL divergence at each iteration have also been established (Liu, 2017; Korba et al., 2020; Salim et al., 2022).

Many variants of SVGD have also been proposed in recent years, some offering improvements and others providing generalisations. Riemannian SVGD (Liu & Zhu, 2018) generalises SVGD by allowing for target densities on Riemannian manifolds, not just Euclidean spaces. Matrix SVGD (Wang et al., 2019) replaces the scalar valued kernel with a matrix valued kernel to incorporate preconditioning information and speed up particle exploration. Message passing SVGD (Zhuo et al., 2018) and graphical SVGD (Wang et al., 2018) focus on target densities that factorise according to a graph structure. This approach reduces variance collapse in high dimensions by converting the problem to a collection of low dimensional problems. Projected SVGD (Chen & Ghattas, 2020), sliced SVGD (Gong et al., 2021) and Grassman SVGD (Liu et al., 2022) also mitigate the issue of variance collapse by updating particles within lower dimensional subspaces, which comes at the expense of additional computation.

In this work, we study a variant called hybrid kernel Stein variational gradient descent (h-SVGD). The name comes from its use of two distinct kernels for the driving and repulsive terms with the aim of mitigating variance collapse. This variant was originally proposed by D’Angelo et al. (2021) in the context of training deep neural network ensembles by sampling from the distribution of network parameters. In that setting, two particles may parameterise networks with very similar outputs despite being far apart in the weight space. Their insight was to encourage functional diversity between networks in the ensemble by using a standard kernel in the driving term, but a functional kernel in the repulsive term. In this neural network ensemble setting, h-SVGD demonstrated better performance than other variants on image classification. Annealed SVGD (D’Angelo & Fortuin, 2021) may also be considered an example of h-SVGD. In this variant, the driving kernel is a scalar multiple $\gamma(\ell) \in [0, 1]$ of the repulsive kernel, and this factor $\gamma(\ell)$ gradually increases to 1 as the iteration ℓ increases. Numerical experiments show that annealed SVGD improves the ability of particles to escape local modes. Scaling one of the update terms has also been used as a computational technique to aid other SVGD variants when training Bayesian neural networks (Gong et al., 2021).

Although preliminary numerical experiments have shown the benefits of h-SVGD (D’Angelo et al., 2021), the theoretical results for SVGD do not directly apply in the hybrid kernel setting due to the presence of a second kernel. In this paper, we address this theoretical gap and reinforce the practical benefits of h-SVGD through the following contributions.

- We establish existence of a solution to the hybrid Stein PDE and a kernelised Wasserstein gradient flow interpretation. Through the study of this gradient flow, we demonstrate that h-SVGD does *not* converge to the target distribution in the mean field limit.
- Other theoretical contributions include quantifying the rate of dissipation of the gradient flow functional, along with a discrete time version, otherwise known as a descent lemma.
- Despite not converging to the target distribution, numerical experiments show that h-SVGD can mitigate variance collapse in the finite particle regime at negligible additional cost, whilst remaining competitive at high dimensional inference tasks.

In Section 2 we clarify notation, recall necessary theory, and outline the vanilla and hybrid SVGD algorithms. Section 3 contains the theoretical contributions, with proofs relegated to Appendix A. Numerical experiments are in Section 4 with additional experiments and details in Appendix B.

2 BACKGROUND

2.1 NOTATION

Let $\mathcal{X} \subseteq \mathbb{R}^d$. Let π denote the target probability density on \mathcal{X} and let $\mathbf{s}_\pi(\mathbf{x}) = \nabla_{\mathbf{x}} \log \pi(\mathbf{x})$. We will often write π for the corresponding measure as well. Assume that $\pi(\mathbf{x}) = e^{-V(\mathbf{x})}$ for some potential V . Let $\mathcal{P}(\mathcal{X})$ be the set of probability measures on \mathcal{X} and let $\mathcal{P}_V(\mathcal{X})$ denote the subset where $\|\mu\|_{\mathcal{P}_V} := \int_{\mathcal{X}} (1 + V(\mathbf{x})) d\mu(\mathbf{x}) < \infty$. For $p \geq 1$, let $\mathcal{P}_p(\mathcal{X})$ be the subset where $\|\mu\|_{\mathcal{P}_p} := \int_{\mathcal{X}} \|\mathbf{x}\|^p d\mu(\mathbf{x}) < \infty$ and define the Wasserstein p -distance between the two measures $\mu, \nu \in \mathcal{P}_p(\mathcal{X})$ as $W_p(\mu, \nu) := (\inf_{\gamma \in \Gamma(\mu, \nu)} \int \|\mathbf{x} - \mathbf{y}\|^p d\gamma(\mu, \nu))^{1/p}$, where $\Gamma(\mu, \nu)$ is the set of couplings between μ and ν . Let $L_c^\infty(\mathcal{X})$ denote the set of probability densities bounded almost

everywhere with compact support. Let $L^2(\mu)$ denote the set of functions that are square integrable with respect to the measure μ . Given $\mu \in \mathcal{P}(\mathcal{X})$ and a smooth, invertible transform $T : \mathcal{X} \rightarrow \mathcal{X}$, let $T_{\#}\mu$ denote the pushforward measure of μ through T . The KL divergence between two measures $\mu, \nu \in \mathcal{P}(\mathcal{X})$ is denoted by $\text{KL}(\mu \parallel \nu)$. Fourier transforms are denoted by \mathcal{F} .

2.2 REPRODUCING KERNEL HILBERT SPACES

A function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is positive definite if $\sum_{i,j} a_i k(\mathbf{x}_i, \mathbf{x}_j) a_j > 0$ for any $a_1, \dots, a_d \in \mathbb{R}$ and $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathcal{X}$. Given a Hilbert space \mathcal{H} of functions $\phi : \mathcal{X} \rightarrow \mathbb{R}$, a function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is said to be a reproducing kernel for \mathcal{H} if it satisfies the reproducing property, $\phi(\mathbf{x}) = \langle \phi, k(\mathbf{x}, \cdot) \rangle_{\mathcal{H}}$ for all $\phi \in \mathcal{H}$. A positive definite $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ admits a unique Hilbert space \mathcal{H} of functions $\phi : \mathcal{X} \rightarrow \mathbb{R}$ for which the Dirac functionals $\delta_{\mathbf{x}} : \mathcal{H} \rightarrow \mathbb{R}, \delta_{\mathbf{x}}\phi = \phi(\mathbf{x})$ are all continuous and k is a reproducing kernel. This Hilbert space is called the reproducing kernel Hilbert space (RKHS) of k and it is equal to the closure of the span of $\{k(\mathbf{x}, \cdot) : \mathbf{x} \in \mathcal{X}\}$. Let $\mathcal{H}^d = \mathcal{H} \times \dots \times \mathcal{H}$ denote the Hilbert space of functions $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ whose components are all in \mathcal{H} , and equip it with the usual inner product $\langle \phi, \psi \rangle_{\mathcal{H}^d} = \sum_{i=1}^d \langle \phi_i, \psi_i \rangle_{\mathcal{H}}$. Given two kernels $k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, let $\mathcal{H}_1, \mathcal{H}_2$ denote their respective RKHS. An important kernel used throughout this paper is the radial basis function (RBF) kernel $k_{\text{RBF}}(\mathbf{x}, \mathbf{y}; h) := \exp(-\|\mathbf{x} - \mathbf{y}\|_2^2 / (2h))$ with bandwidth $h > 0$. For a thorough treatment of RKHS we refer the reader to (Aronszajn, 1950; Steinwart & Christmann, 2008; Berlinet & Thomas-Agnan, 2011).

2.3 STEIN VARIATIONAL GRADIENT DESCENT

The key result from (Liu & Wang, 2016) identifies a transform $T : \mathcal{X} \rightarrow \mathcal{X}$ that optimally decreases the KL divergence from an arbitrary probability measure to π . More precisely, let \mathcal{H} be an RKHS with kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and consider transforms of the form $T(\mathbf{x}) = \mathbf{x} + \epsilon \phi(\mathbf{x})$ where $\epsilon > 0$ and ϕ is in the unit ball $\{\phi \in \mathcal{H}^d : \|\phi\|_{\mathcal{H}^d} \leq 1\}$. The maximum value of

$$-\nabla_{\epsilon} \text{KL}(T_{\#}\mu \parallel \pi)|_{\epsilon=0} \quad (1)$$

occurs at $\phi_{\mu, \pi}^k / \|\phi_{\mu, \pi}^k\|_{\mathcal{H}^d}$, where

$$\phi_{\mu, \pi}^k(\cdot) := \mathbb{E}_{\mathbf{x} \sim \mu} [k(\mathbf{x}, \cdot) \mathbf{s}_{\pi}(\mathbf{x}) + \nabla_{\mathbf{x}} k(\mathbf{x}, \cdot)]. \quad (2)$$

When μ is an empirical distribution (i.e. a sum of Dirac measures), the expectation in (2) can be computed exactly by summing over the particles of each Dirac measure. Using this observation, the SVGD algorithm starts with an initial set of N particles $(\mathbf{x}_0^i)_{i=1}^N$ and iteratively applies the transform T with (2) as the update direction. At each iteration ℓ , this yields a set of particles $(\mathbf{x}_{\ell}^i)_{i=1}^N$ and a corresponding empirical distribution $\mu_{\ell} = \frac{1}{N} \sum_i \delta_{\mathbf{x}_{\ell}^i}$. This is captured in Algorithm 1. The intention is that after sufficiently many iterations, the set of particles will resemble samples from π and expectations of the form $\mathbb{E}_{\mathbf{x} \sim \pi} h(\mathbf{x})$ can be approximated by $\mathbb{E}_{\mathbf{x} \sim \mu_{\ell}} h(\mathbf{x}) = \frac{1}{N} \sum_i h(\mathbf{x}_{\ell}^i)$. We also recall the definition of the kernelised Stein discrepancy (KSD) from (Liu et al., 2016),

$$\mathbb{S}_k(\mu, \pi) := \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \mu} [(s_{\pi}(\mathbf{x}) - s_{\mu}(\mathbf{x}))^{\top} k(\mathbf{x}, \mathbf{y}) (s_{\pi}(\mathbf{y}) - s_{\mu}(\mathbf{y}))]. \quad (3)$$

Algorithm 1 Stein Variational Gradient Descent (Liu & Wang, 2016)

Input: A target probability distribution π , a kernel k , an initial set of particles $(\mathbf{x}_0^i)_{i=1}^N$ in \mathcal{X} , and a sequence of step sizes (ϵ_{ℓ}) .

Output: A set of particles $(\mathbf{x}_{\ell}^i)_{i=1}^N$ in \mathcal{X} whose empirical distribution approximates π .
for iteration ℓ **do**

$$\mathbf{x}_{\ell+1}^i \leftarrow \mathbf{x}_{\ell}^i + \epsilon_{\ell} \hat{\phi}_{\mu_{\ell}, \pi}^*(\mathbf{x}_{\ell}^i), \quad \forall i = 1, \dots, N$$

$$\hat{\phi}_{\mu_{\ell}, \pi}^*(\mathbf{x}) = \frac{1}{N} \sum_{j=1}^N \left(k(\mathbf{x}_{\ell}^j, \mathbf{x}) \nabla_{\mathbf{x}_{\ell}^j} \log \pi(\mathbf{x}_{\ell}^j) + \nabla_{\mathbf{x}_{\ell}^j} k(\mathbf{x}_{\ell}^j, \mathbf{x}) \right) \quad (4)$$

2.4 HYBRID KERNEL STEIN VARIATIONAL GRADIENT DESCENT

The SVGD update in (4) contains two terms, each using the same kernel. The first term, often referred to as the driving term, uses the score function to move particles towards regions of high probability density, and the repulsive term prevents particles from collapsing at the modes. The h-SVGd variant proposed by D’Angelo et al. (2021) uses a different kernel in each term. Let k_1 denote the kernel that appears alongside the score function, and let k_2 denote the repulsive kernel. For the remainder of this paper, k_1 and k_2 will both be positive definite. We present h-SVGd in Algorithm 2.

Algorithm 2 Hybrid Kernel Stein Variational Gradient Descent

Input: A target probability distribution π , two kernels k_1, k_2 , an initial set of particles $(\mathbf{x}_0^i)_{i=1}^N$ in \mathcal{X} , and a sequence of step sizes (ϵ_ℓ) .

Output: A set of particles $(\mathbf{x}^i)_{i=1}^N$ in \mathcal{X} whose empirical distribution approximates π .
for iteration ℓ **do**

$$\mathbf{x}_{\ell+1}^i \leftarrow \mathbf{x}_\ell^i + \epsilon_\ell \hat{\phi}_{\mu_\ell, \pi}^*(\mathbf{x}_\ell^i), \quad \forall i = 1, \dots, N$$

$$\hat{\phi}_{\mu_\ell, \pi}^*(\mathbf{x}) = \frac{1}{N} \sum_{j=1}^N \left(k_1(\mathbf{x}_\ell^j, \mathbf{x}) \nabla_{\mathbf{x}_\ell^j} \log \pi(\mathbf{x}_\ell^j) + \nabla_{\mathbf{x}_\ell^j} k_2(\mathbf{x}_\ell^j, \mathbf{x}) \right) \quad (5)$$

3 THEORETICAL RESULTS

3.1 DEFINITIONS AND ASSUMPTIONS

A function $f : \mathcal{X} \rightarrow \mathbb{R}$ is in the Stein class of π if it is smooth and satisfies the identity $\int_{\mathbf{x} \in \mathcal{X}} \nabla_{\mathbf{x}} (f(\mathbf{x}) \pi(\mathbf{x})) d\mathbf{x} = 0$. A function $\mathbf{f} = (f_1, \dots, f_d) : \mathcal{X} \rightarrow \mathbb{R}^d$ is in the Stein class of π if each f_i is in the Stein class of π . A kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is in the Stein class of π if it has continuous second order partial derivatives and both $k(\mathbf{x}, \cdot)$ and $k(\cdot, \mathbf{y})$ are in the Stein class of π for all $\mathbf{x}, \mathbf{y} \in \mathcal{X}$. The hybrid Stein operator acts on a pair of kernels $k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ by

$$\mathcal{S}_\pi \otimes (k_1, k_2)(\mathbf{x}, \cdot) := k_1(\mathbf{x}, \cdot) \mathbf{s}_\pi(\mathbf{x}) + \nabla_{\mathbf{x}} k_2(\mathbf{x}, \cdot),$$

provided k_1 and k_2 both belong to the Stein class of π . This reduces to the Stein operator defined in (Liu et al., 2016) when $k_1 = k_2$. Motivated by the h-SVGd update in (5), define the update direction

$$\phi_{\mu, \pi}^{k_1, k_2}(\cdot) := \mathbb{E}_{\mathbf{x} \sim \mu} [\mathcal{S}_\pi \otimes (k_1, k_2)(\mathbf{x}, \cdot)], \quad (6)$$

and write $\phi^* = \phi_{\mu, \pi}^{k_1, k_2} / \|\phi_{\mu, \pi}^{k_1, k_2}(\cdot)\|_{\mathcal{H}_1}$. Let $\mathbf{G}(\cdot; k_1, \mu, \pi) := \mathbb{E}_{\mathbf{x} \sim \mu} [k_1(\mathbf{x}, \cdot) \mathbf{s}_\pi(\mathbf{x})]$ denote the gradient term and $\mathbf{R}(\cdot; k_2, \mu) := \mathbb{E}_{\mathbf{x} \sim \mu} [\nabla_{\mathbf{x}} k_2(\mathbf{x}, \cdot)]$ the repulsive term. We can then write $\phi_{\mu, \pi}^{k_1, k_2}(\cdot) = \mathbf{G}(\cdot; k_1, \mu, \pi) + \mathbf{R}(\cdot; k_2, \mu)$. The update transform

$$\mathbf{T}_{\mu, \pi}^{k_1, k_2}(\mathbf{x}) = \mathbf{x} + \epsilon \phi_{\mu, \pi}^{k_1, k_2}(\mathbf{x}) \quad (7)$$

and the map $\Phi_\pi^{k_1, k_2} : \mu \mapsto (\mathbf{T}_{\mu, \pi}^{k_1, k_2})_{\#} \mu$ characterise the h-SVGd dynamics. For each ℓ , define

$$\mu_{\ell+1}^N := \Phi_\pi^{k_1, k_2}(\mu_\ell^N), \quad \mu_{\ell+1}^\infty := \Phi_\pi^{k_1, k_2}(\mu_\ell^\infty), \quad (8)$$

where μ_0^N is the empirical measure of the initial particles $(\mathbf{x}_0^i)_{i=1}^N$ drawn i.i.d. from some μ_0^∞ .

All technical assumptions required in the theorems throughout this section are detailed here for completeness. The first set of assumptions relate to the potential of the target distribution.

(A1) $V \in C^\infty(\mathcal{X})$, $V \geq 0$, and $\lim_{|\mathbf{x}| \rightarrow \infty} V(\mathbf{x}) = +\infty$.

(A2) There exist constants $C_V > 0$ and $q > 1$ such that for all $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,

$$|\nabla V(\mathbf{x})|^q \leq C_V(1 + V(\mathbf{x}))$$

and

$$\sup_{\theta \in [0,1]} |\nabla^2 V(\theta \mathbf{x} + (1 - \theta) \mathbf{y})|^q \leq C_V(1 + V(\mathbf{x}) + V(\mathbf{y})).$$

(A3) For any $\alpha, \beta > 0$, there exists a constant $C_{\alpha, \beta} > 0$ such that

$$|\mathbf{y}| \leq \alpha |\mathbf{x}| + \beta \implies (1 + |\mathbf{x}|)(|\nabla V(\mathbf{y})| + |\nabla^2 V(\mathbf{y})|) \leq C_{\alpha, \beta}(1 + V(\mathbf{x})).$$

(A4) The Hessian H_V of V is well-defined and $\|H_V\|_{\text{op}} \leq M$ for some $M > 0$.

Assumptions (A1), (A2) and (A3) are identical to those in (Lu et al., 2019). Assumption (A4) is identical to Assumption (A2) in (Korba et al., 2020), and Assumption 2.1 in (Salim et al., 2022). We remark that Assumptions (A2) and (A3) simply control the decay of the tails of π . The above assumptions are all satisfied by normal distributions, mixtures thereof, as well as many continuous distributions from the exponential family. The following assumptions are also required to ensure sufficient regularity of the kernel functions.

(B1) There exist symmetric functions $K_1, K_2 : \mathcal{X} \rightarrow \mathbb{R}$ such that $k_1(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x} - \mathbf{y})$, $k_2(\mathbf{x}, \mathbf{y}) = K_2(\mathbf{x} - \mathbf{y})$, K_1 is C^2 with bounded derivatives, and K_2 is C^4 with bounded derivatives. By symmetric, we mean $K_i(\mathbf{x}) = K_i(-\mathbf{x})$ for all $i = 1, 2$ and $\mathbf{x} \in \mathcal{X}$. We use $B > 0$ as a bound for all derivatives in the proofs.

(B2) There exists a constant $D > 0$ such that both k_1 and ∇k_2 and are D -Lipschitz, and $\nabla V(\cdot)k_1(\cdot, \mathbf{z})$ is D -Lipschitz for each \mathbf{z} . That is,

$$\begin{aligned} |k_1(\mathbf{x}, \mathbf{x}') - k_1(\mathbf{y}, \mathbf{y}')| &\leq D(\|\mathbf{x} - \mathbf{y}\|_2 + \|\mathbf{x}' - \mathbf{y}'\|_2), \\ \|\nabla_{\mathbf{x}} k_2(\mathbf{x}, \mathbf{x}') - \nabla_{\mathbf{y}} k_2(\mathbf{y}, \mathbf{y}')\| &\leq D(\|\mathbf{x} - \mathbf{y}\|_2 + \|\mathbf{x}' - \mathbf{y}'\|_2), \\ \|\nabla V(\mathbf{x})k_1(\mathbf{x}, \mathbf{z}) - \nabla V(\mathbf{y})k_1(\mathbf{y}, \mathbf{z})\| &\leq D(\|\mathbf{x} - \mathbf{y}\|_2) \end{aligned}$$

for all $\mathbf{x}, \mathbf{x}', \mathbf{y}, \mathbf{y}', \mathbf{z} \in \mathcal{X}$.

Assumption (B1) is a slight relaxation of Assumption 2.1 in (Lu et al., 2019), and contains Assumption 2.6 in (Salim et al., 2022). The first two parts of Assumption (B2) are hybrid kernel versions of Assumption (B2) from (Korba et al., 2020), and the third part of Assumption (B2) relaxes the restrictive Assumption (B1) from (Korba et al., 2020).

3.2 LARGE PARTICLE LIMIT

We begin our theoretical study with a result ensuring convergence of h-SVGD in the large particle limit. The single kernel version (Korba et al., 2020, Proposition 7) of the following result uses an assumption that is quite restrictive. It requires $|V(\mathbf{x})| \leq C_V$ for some constant $C_V > 0$, which rules out even a normal target distribution. We relax this with the third part of Assumption (B2) and provide an updated proof in the appendix.

Proposition 3.1. *Assume (A1), (A4), (B1) and (B2), and let $T > 0$. For any $0 \leq \ell \leq \frac{T}{\epsilon_\ell}$, there exists a constant L depending on k_1, k_2 and π such that*

$$\mathbb{E}[W_2^2(\mu_\ell^N, \mu_\ell^\infty)] \leq \frac{1}{2\sqrt{N}} \sqrt{\text{var}(\mu_0^\infty)} e^{LT} (e^{2LT} - 1).$$

3.3 HYBRID STEIN PDE

In the continuous time limit, equation (5) becomes a coupled system of differential equations,

$$\frac{d\mathbf{x}_i}{dt} = \frac{1}{N} \sum_{j=1}^N (k_1(\mathbf{x}_j, \mathbf{x}_i) \nabla_{\mathbf{x}_j} \log \pi(\mathbf{x}_j) + \nabla_{\mathbf{x}_j} k_2(\mathbf{x}_j, \mathbf{x}_i)), \quad i = 1, \dots, N. \quad (9)$$

In the mean field limit, integration by parts yields

$$\frac{d\mathbf{x}}{dt} = \int (k_1(\mathbf{x}', \mathbf{x}) \nabla_{\mathbf{x}'} \log \pi(\mathbf{x}') + \nabla_{\mathbf{x}'} k_2(\mathbf{x}', \mathbf{x})) \rho(\mathbf{x}') d\mathbf{x}' \quad (10)$$

$$= T_{k_1, \rho}(\nabla \log \pi)(\mathbf{x}) - T_{k_2, \rho}(\nabla \log \rho)(\mathbf{x}) \quad (11)$$

where $T_{k, \rho} : L^2(\rho)^d \rightarrow \mathcal{H}_k^d$ is the Hilbert-Schmidt operator $(T_{k, \rho} f)(\cdot) = \int k(\cdot, \mathbf{x}) f(\mathbf{x}) d\rho(\mathbf{x})$ for a kernel k . Recalling that V is the potential of π and so $V = -\log \pi$, this mean field limit can be described by the hybrid Stein partial differential equation

$$\partial_t \rho_t = \nabla \cdot (\rho_t (K_1 * \nabla V \rho_t + K_2 * \nabla \rho_t)). \quad (12)$$

Definition 3.1. Given a probability measure ν on \mathbb{R}^d , a map $X(t, \mathbf{x}; \nu) : [0, \infty) \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ that is C^1 with respect to t and satisfies

$$\begin{aligned}\partial_t X(t, \mathbf{x}; \nu) &= -(K_1 * (\nabla V \rho_t))(X(t, \mathbf{x}; \nu)) - (\nabla K_2 * \rho_t)(X(t, \mathbf{x}; \nu)) \\ \rho_t &= X(t, \cdot, \nu)_\# \nu \\ X(0, \mathbf{x}; \nu) &= \mathbf{x}\end{aligned}\tag{13}$$

is called a mean field characteristic flow of (10) or of (12).

We now generalise Theorem 2.4 from (Lu et al., 2019), ensuring the existence of a solution to the hybrid Stein PDE. First, define the set $Y := \{u \in C(\mathcal{X}, \mathcal{X}) : \sup_{\mathbf{x} \in \mathcal{X}} |u(\mathbf{x}) - \mathbf{x}| < \infty\}$ with $d_Y(u, v) = \sup_{\mathbf{x} \in \mathcal{X}} |u(\mathbf{x}) - v(\mathbf{x})|$ and note that (Y, d_Y) is a complete metric space.

Proposition 3.2. Assume (A1), (A2), (A3) and (B1), and let $T > 0$. Then there is a unique solution $X(\cdot, \cdot, \nu) \in C^1([0, T]; Y)$ to (13) and the corresponding ρ_t is a weak solution to (12) satisfying

$$\|\rho_t\|_{\mathcal{P}_V} \leq \|\pi\|_{\mathcal{P}_V} \exp(C \min(\|\nabla K_1\|_\infty, \|\nabla K_2\|_\infty) t)\tag{14}$$

for some constant $C > 0$ depending on K_1 , K_2 and V .

The second kernel enables a stronger bound than Theorem 2.4 from (Lu et al., 2019) by careful modification of the proof (see Appendix A). In particular, ensuring that $\|\nabla K_1\|_\infty < \|\nabla K_2\|_\infty$ when choosing K_2 yields a stronger bound in (14) than if K_1 were used for both kernels. We remark that this bound describes regularity of the solution to the PDE, not a rate of convergence.

3.4 KERNELISED WASSERSTEIN GRADIENT FLOW AND ASYMPTOTIC DENSITY: $k_2 = ck_1$

Zhuo et al. (2018) uncovered a correlation between dimension and the magnitude of the repulsive force. Under some technical conditions, for any $\alpha, \delta \in (0, 1)$, they show that SVGD under an RBF kernel yields $\|\mathbf{R}(\cdot; k_2, \mu)\|_\infty = O(d^{-\alpha})$ with probability at least $1 - \delta$. This suggests that simply scaling the repulsive force by d^α for some $\alpha \in (0, 1)$ should offset the decrease in $\|\mathbf{R}(\cdot; k_2, \mu)\|_\infty$ in high dimensions, thereby alleviating variance collapse at negligible additional computational cost. Scaling the repulsive kernel in this way corresponds to h-SVGd where k_1 is an RBF kernel and $k_2 = d^\alpha k_1$. This motivates our study of the h-SVGd gradient flow under the special case $k_2 = ck_1$.

Recall that the Wasserstein gradient flow on a functional \mathcal{F} is the time evolution of the probability measure ρ_t that minimises $\mathcal{F}(\rho)$. This is described by a PDE and the Wasserstein gradient is defined through the Wasserstein distance W_2 . We refer the reader to (Ambrosio et al., 2005; Santambrogio, 2015) for further details.

In the case where $k = k_1 = k_2$, (10) and (12) describe a kernelised Wasserstein gradient flow of the form $\partial_t \rho_t = \nabla \cdot (\rho_t T_{k, \rho_t} \nabla_W \mathcal{F}(\rho_t))$, where

$$\mathcal{F}(\rho) = \text{KL}(\rho \parallel \pi) = \mathbb{E}_{\mathbf{x} \sim \rho} [\log \rho(\mathbf{x}) - \log \pi(\mathbf{x})]$$

is the KL divergence functional (Liu, 2017). This has functional derivative $\frac{\delta \mathcal{F}}{\delta \rho} = \log \rho - \log \pi + 1$ up to a constant, and Wasserstein gradient

$$\nabla_W \mathcal{F}(\rho) = \nabla \frac{\delta \mathcal{F}}{\delta \rho}(\rho) = \nabla \log \rho - \nabla \log \pi.\tag{15}$$

Then (11) can be written as

$$\frac{d\mathbf{x}}{dt} = -T_{k, \rho} \nabla_W \mathcal{F}(\rho),$$

and the corresponding Fokker-Planck equation is

$$\partial_t \rho_t = \nabla \cdot (\rho_t T_{k, \rho_t} (\nabla \log \rho_t - \nabla \log \pi))\tag{16}$$

where $\{\rho_t : t \geq 0\}$ is a curve of probability densities. The following result generalises this gradient flow interpretation to the case where $k_2 = ck_1$. Note that it applies to any positive definite kernel k_1 , not just the RBF kernel as discussed earlier.

Proposition 3.3. *Given a positive definite kernel k_1 and constant $c > 0$, let $k_2 = ck_1$. Then the mean field dynamics of h-SVGD describe a kernelised Wasserstein gradient flow on the functional*

$$\mathcal{F}(\rho) = \mathbb{E}_{\mathbf{x} \sim \rho} [c \log \rho - \log \pi(\mathbf{x})], \quad (17)$$

whose Wasserstein gradient is

$$\nabla_W \mathcal{F}(\rho) = c \nabla \log \rho - \nabla \log \pi. \quad (18)$$

The corresponding continuity equation is

$$\partial_t \rho_t = \nabla \cdot (\rho_t T_{k_1, \rho_t} (c \nabla \log \rho_t - \nabla \log \pi)). \quad (19)$$

Even in this simple hybrid kernel setting, the following result establishes that the limiting distribution ρ^* of the mean field regime is not equal to the target distribution π .

Corollary 3.4. *If $k_2 = ck_1$ for some constant $c > 0$ where k_1 is a positive definite kernel, then the mean field h-SVGD has a fixed point $\rho^*(\mathbf{x}) \propto \pi(\mathbf{x})^{1/c}$.*

Although Corollary 3.4 applies to any target density π satisfying (A1), it is especially insightful to consider a normal target. If $\pi(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$, then $\rho^*(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, c\boldsymbol{\Sigma})$. So scaling the repulsive kernel k_2 will adjust the variance of the target by the same factor. This supports the motivation that scaling the repulsive kernel should offset the variance underestimation in high dimensions at a negligible additional cost. This idea will be revisited in Section 4.

We now generalise an existing result (Korba et al., 2020, Proposition 1) that describes the dissipation of the KL divergence along the SVGD gradient flow. The result below describes the dissipation of the functional in (17) along the h-SVGD gradient flow, ensuring that the functional always decreases. It also describes the dissipation of the KL divergence with respect to the mean field limiting distribution ρ^* , which we emphasise is not equal to the target distribution π , as per Corollary 3.4.

Proposition 3.5. *Under the assumptions of Proposition 3.3,*

$$\frac{d}{dt} \mathcal{F}(\rho_t) = -c^2 \|T_{k_1, \rho_t} (\nabla \log \rho - \nabla \log \rho^*)\|_{\mathcal{H}_1^d}^2 + c \int \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x}, \quad (20)$$

$$\frac{d}{dt} \text{KL}(\rho_t \parallel \rho^*) = -c \|T_{k_1, \rho_t} (\nabla \log \rho - \nabla \log \rho^*)\|_{\mathcal{H}_1^d}^2 + \int \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x}, \quad (21)$$

where $\rho^(\mathbf{x}) \propto \pi(\mathbf{x})^{1/c}$ is the mean field fixed point. Furthermore,*

$$\int \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x} \leq 0,$$

so $\frac{d}{dt} \mathcal{F}(\rho_t) \leq 0$ and $\frac{d}{dt} \text{KL}(\rho \parallel \rho^) \leq 0$ for $t \geq 0$, and $\frac{d}{dt} \mathcal{F}(\rho^*) = 0$.*

The following descent lemma is adapted from (Liu, 2017, Theorem 3.3) and it provides a discrete time version of Proposition 3.5. We use μ_ℓ to denote discrete time steps of the algorithm, as defined in (8), as opposed to ρ_t for the continuous time analysis. We remark that other descent lemmas for SVGD have been established (Korba et al., 2020; Salim et al., 2022).

Proposition 3.6. *Set $\epsilon_\ell = (2 \sup_{\mathbf{x}} \sigma(\nabla \phi^* + \nabla \phi^{*\top}))^{-1}$ where $\sigma(\cdot)$ denotes the spectral radius of a matrix. Define the quantity $R = \sup_{\mathbf{x}} \{\frac{1}{2} \|\nabla \log \pi\|_{\text{Lip}} k(\mathbf{x}, \mathbf{x}) + 2 \nabla_{\mathbf{x}, \mathbf{y}} k(\mathbf{x}, \mathbf{x})\}$ where $\nabla_{\mathbf{x}, \mathbf{y}} k(\mathbf{x}, \mathbf{x}) = \sum_i \partial_{x_i} \partial_{y_i} k(\mathbf{x}, \mathbf{y})|_{\mathbf{x}=\mathbf{y}}$. Then*

$$\frac{1}{\epsilon_\ell} (\mathcal{F}(\mu_{\ell+1}^\infty) - \mathcal{F}(\mu_\ell^\infty)) \leq -(1 - \epsilon_\ell R) c \mathbb{S}(\mu_\ell, \rho^*).$$

This result ensures that for a sufficiently small step size ϵ_ℓ , the functional \mathcal{F} will decrease at each step of the algorithm, until μ_ℓ approaches ρ^* .

3.5 KERNELISED WASSERSTEIN GRADIENT FLOW: THE GENERAL CASE

In this section, we present a generalisation of Proposition (3.3) and discuss some difficulties in finding kernels that satisfy the required conditions. For ease of presentation, we restrict our attention

to $d = 1$. We assume **(B1)** throughout, and also assume that any required Fourier transforms exist. Define the function $r : \mathcal{X} \rightarrow \mathbb{R}$ by

$$r(x; \rho) := \mathcal{F}^{-1} \left(\frac{\mathcal{F}(K_2)}{\mathcal{F}(K_1)} \cdot \mathcal{F}(\nabla \rho) \right) (x) \quad (22)$$

and let $R : \mathcal{X} \rightarrow \mathbb{R}$ denote a function satisfying $\nabla R(x; \rho) = r(x; \rho)/\rho(x)$.

Proposition 3.7. *Assume that both r and R exist. Then the corresponding continuity equation is*

$$\partial_t \rho_t = \nabla \cdot (\rho_t T_{k_1, \rho_t} (\nabla R(\cdot; \rho_t) - \nabla \log \pi)). \quad (23)$$

If in addition

$$\int \frac{\partial}{\partial \epsilon} R(x; \rho + \epsilon \chi) d(\rho + \epsilon \chi)(x) \Big|_{\epsilon=0} = 0 \quad (24)$$

for any measure $\chi = \tilde{\rho} - \rho$ with $\tilde{\rho} \in L_c^\infty(\mathcal{X}) \cap \mathcal{P}(\mathcal{X})$, then the mean field dynamics of h -SVGD describe a kernelised Wasserstein gradient flow on the functional

$$\mathcal{F}(\rho) = \mathbb{E}_{x \sim \rho} [R(x; \rho) - \log \pi(x)], \quad (25)$$

whose Wasserstein gradient is

$$\nabla_W \mathcal{F}(\rho) = r(x; \rho)/\rho(x) - \nabla \log \pi(x). \quad (26)$$

Note that when $k_2 = ck_1$, we have $\mathcal{F}(K_2)/\mathcal{F}(K_1) = c$ and so $r(x; \rho) = c\nabla \rho(x)$. Therefore, $R(x; \rho) = c \log \rho(x)$, and so (23) reduces to (19). Furthermore, the left hand side of (24) is equal to $c \int d\chi(x) = c \int d\tilde{\rho}(x) - c \int d\rho(x) = 0$ because $\tilde{\rho}, \rho$ are both probability measures. So (24) is satisfied in this special case.

We remark that verifying (24) remains a challenge in general hybrid kernel settings. However, the derivation of the continuity equation does not rely on (24), so we can gain some insights into the steady state by studying when (23) is equal to zero. We present a specific example below to demonstrate the behaviour of SVGD in the general hybrid kernel setting. In particular, the form of the target π and the mean field steady state ρ^* can look quite different, even under a hybrid kernel setting as simple as two RBF kernels with different bandwidths.

Proposition 3.8. *Let $h_1, h_2, \sigma > 0$ such that $\Delta h := h_2 - h_1 \neq 0$. Let $k_1(x, y) = k_{\text{RBF}}(x, y; h_1)$ and $k_2(x, y) = k_{\text{RBF}}(x, y; h_2)$, and let $\pi(x) \propto \exp(-\alpha \exp(\frac{x^2}{2\beta}))$ be the target on $\mathcal{X} = \mathbb{R}$ where*

$$\alpha = \sqrt{\frac{h_2}{h_1}} \cdot \sqrt{\frac{\sigma^2}{\sigma^2 + \Delta h}} \cdot \frac{\sigma^2}{\Delta h}, \quad \beta = \frac{\sigma^2(\sigma^2 + \Delta h)}{\Delta h}.$$

Then $\rho^(x) = \mathcal{N}(x; 0, \sigma^2)$ is a fixed point of the h -SVGD mean field dynamics.*

We emphasise that the $k_2 = ck_1$ case studied in the previous section is the focus of this paper due to its capacity to alleviate variance collapse. We leave a more detailed study of the general hybrid setting for future work.

4 EXPERIMENTS

The problem of variance collapse in SVGD has been successfully mitigated in the setting of probabilistic graphical models where the conditional dependence structure enables π to be factorised and $\mathbf{R}(\cdot; k_2, \mu)$ to be replaced with a set of lower dimensional repulsive forces (Zhuo et al., 2018). Other methods such as S-SVG (Gong et al., 2021) and GSVG (Liu et al., 2022) have demonstrated that variance collapse can be avoided, at the expense of additional computational cost. In particular, S-SVG requires computation of the optimal test directions, and GSVG requires the projectors to be updated at each step. Our numerical experiments demonstrate that in the absence of a conditional dependence structure, and without incurring additional computational cost, h -SVGD can mitigate variance collapse while improving the inference capabilities of SVGD. In light of the discussion following Corollary 3.4, we argue that this mitigation occurs because the higher variance in the mean field limit offsets the variance underestimation in the finite particle setting. We measure variance collapse using dimension averaged marginal variance (DAMV), $\frac{1}{d} \sum_{j=1}^d \text{Var}_j(\{\mathbf{x}_i\}_{i=1}^N)$, as is standard in the literature (Zhuo et al., 2018; Ba et al., 2019; 2021; Gong et al., 2021). Further details on the the experimental details can be found in Appendix B.

4.1 MULTIVARIATE GAUSSIAN MIXTURE

In this first example, we sample from a high-dimensional Gaussian distribution, a setting which has been explored previously to study variance collapse (Gong et al., 2021; Liu et al., 2022; Zhuo et al., 2018). We use a target distribution $\pi = \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ for dimensions up to $d = 100$. We choose to sample $N = 50$ particles in order to demonstrate the performance of h-SVGD when d is much greater than N , as is often the case in high dimensional Bayesian inference. Particles are initialised from $\mathcal{N}(2\mathbf{1}_d, 2\mathbf{I}_d)$, where $\mathbf{1}_d$ is the vector of ones, and each SVGD variant is run for 2000 iterations with an initial step size of $\epsilon = 0.01$, adapted using AdaGrad. We run SVGD and S-SVGD with the RBF kernel k_{RBF} and compare against h-SVGD with kernels $k_1 = k_{\text{RBF}}$ and $k_2 = f(d)k_{\text{RBF}}$ for the following choices of factor f ,

$$f(d; \alpha) := d^\alpha, \quad \alpha \in (0, 1), \quad (27)$$

$$f(d; \alpha_1, \alpha_2, c) := \frac{cd^{\alpha_1} + d^{1+\alpha_2}}{c + d}, \quad \alpha_1, \alpha_2 \in (0, 1), c > 0. \quad (28)$$

All algorithms use an adaptive bandwidth $h = \text{med}^2 / \log(N)$ where med is the median pairwise distance between particles (Liu & Wang, 2016). For S-SVGD, there is one bandwidth per dimension, so the median distances are computed along each projection, as described in (Gong et al., 2021). All configurations are run 5 times and results for each configuration are averaged.

Figure 1 demonstrates that h-SVGD provides a noticeable uplift in marginal variance estimation under the repulsive factors $f(d; 0.5)$ and $f(d; 0.8)$ when compared to SVGD, although increasing α too much leads to overestimation, especially in lower dimensions. The motivation for (28) was the observation that the variance estimation of (27) is better in lower dimensions for smaller α , and better in higher dimensions for larger α . Since $f(d; \alpha_1, \alpha_2, c) = \frac{c}{c+d}f(d; \alpha_1) + \frac{d}{c+d}f(d; \alpha_2)$, this choice should enable the advantages of both small and large α for moderately large c . Figure 1 demonstrates that choosing $\alpha_1 = 0.5$, $\alpha_2 = 0.8$ and $c = 100$ provides consistent variance estimation up to $d = 100$. An interesting avenue for future research would be to develop a more principled method of choosing such a repulsive factor. The techniques in (Zhuo et al., 2018, Propositions 1-2) offer a promising direction. We remark that S-SVGD estimates the variance fairly consistently as the dimension increases. However, this comes at a significant increase in runtime, whereas there is no noticeable difference between the runtimes of SVGD and h-SVGD.

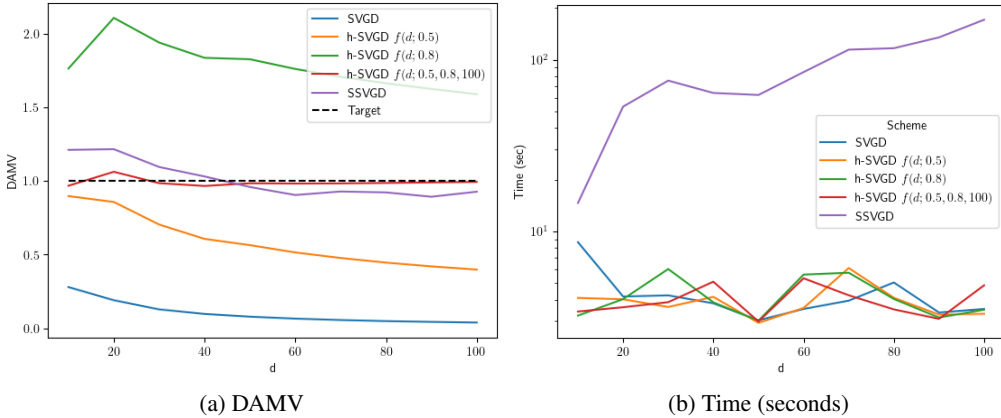


Figure 1: DAMV and runtime for different SVGD variants sampling from $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.

4.2 BAYESIAN NEURAL NETWORK

In this section, we sample weights from a Bayesian neural network. Aside from scaling the repulsive kernel by $f(d; 0.5, 0.8, 100)$, as defined in (28) during the previous experiment, our setup is identical to (Liu & Wang, 2016). For completeness, details are included in Appendix B. Table 1 shows that with no additional computational cost, the problem of variance collapse is mitigated under h-SVGD through an increased DAMV. Following Corollary 3.4, one may interpret the increased

variance in the mean field limit as an offset to the variance underestimation in finite dimensions, leading to improved variance estimation under h-SVGD. Table 2 demonstrates that it remains competitive at inference through improved test log-likelihood (LL) for all datasets and improved root mean squared error (RMSE) for all but one dataset. We observe that for two datasets, Kin8nm and Protein, h-SVGD enjoys an improvement in RMSE and LL despite a small decrease in the DAMV. This suggests that the improved performance may be linked to another property of h-SVGD. We leave further exploration of such properties for future work. Appendix B includes a comparison of the same metrics between h-SVGD and S-SVGD to emphasise that the improvement in variance estimation comes at a lower cost.

Table 1: DAMV and runtime in seconds of SVGD and h-SVGD.

Dataset	DAMV		Runtime (seconds)	
	SVGD	h-SVGD	SVGD	h-SVGD
Boston	0.051 ± 0.011	0.087 ± 0.010	28.9 ± 1.4	28.5 ± 0.9
Concrete	0.084 ± 0.01	0.102 ± 0.006	28.7 ± 1.3	28.7 ± 1.3
Energy	0.065 ± 0.015	0.106 ± 0.011	30.8 ± 2.8	30.7 ± 2.2
Kin8nm	0.105 ± 0.003	0.093 ± 0.003	35.9 ± 1.3	36.0 ± 1.3
Naval	0.059 ± 0.004	0.068 ± 0.011	33.4 ± 0.8	34.5 ± 0.9
Combined	0.128 ± 0.008	0.138 ± 0.006	36.3 ± 1.5	36.7 ± 2.4
Protein	0.089 ± 0.001	0.084 ± 0.001	72.7 ± 1.0	72.8 ± 1.7
Wine	0.068 ± 0.005	0.075 ± 0.005	29.5 ± 1.4	29.8 ± 1.1
Yacht	0.060 ± 0.020	0.121 ± 0.012	29.3 ± 0.4	29.8 ± 1.3
Year	$0.011 \pm \text{NA}$	$0.012 \pm \text{NA}$	$673 \pm \text{NA}$	$666 \pm \text{NA}$

Table 2: Average RMSE and LL of SGVD and h-SVGD evaluated on the test dataset.

Dataset	Test RMSE		Test LL	
	SVGD	h-SVGD	SVGD	h-SVGD
Boston	3.094 ± 0.579	3.001 ± 0.584	-2.123 ± 0.116	-1.988 ± 0.221
Concrete	5.857 ± 0.468	5.210 ± 0.529	-2.616 ± 0.099	-2.535 ± 0.179
Energy	1.528 ± 0.169	1.040 ± 0.128	-1.702 ± 0.094	-0.805 ± 0.104
Kin8nm	0.124 ± 0.005	0.090 ± 0.003	-1.293 ± 0.108	0.468 ± 0.090
Naval	0.006 ± 0.000	0.004 ± 0.000	-1.353 ± 0.161	-0.090 ± 0.105
Combined	4.105 ± 0.220	4.057 ± 0.218	-2.459 ± 0.051	-2.354 ± 0.052
Protein	4.791 ± 0.025	4.600 ± 0.026	-2.633 ± 0.035	-2.456 ± 0.017
Wine	0.637 ± 0.044	0.626 ± 0.045	-1.463 ± 0.120	-0.750 ± 0.097
Yacht	1.677 ± 0.571	1.861 ± 0.662	-1.587 ± 0.120	-0.813 ± 0.227
Year	$8.837 \pm \text{NA}$	$8.689 \pm \text{NA}$	$-2.883 \pm \text{NA}$	$-2.872 \pm \text{NA}$

5 CONCLUSION

We developed the mean field theory of h-SVGD by proving the existence of a solution to the hybrid Stein PDE and identifying it as a gradient flow on a functional other than the KL divergence. We characterised the mean field fixed point for the special case $k_2 = ck_1$ and demonstrated that h-SVGD does not converge to the target in the mean field limit. We provided a result on the dissipation of the new functional, as well as a discrete time version, otherwise known as a descent lemma. We also highlighted the complexities of the gradient flow in the general hybrid kernel setting. Experimental results demonstrated that h-SVGD can alleviate variance collapse in high dimensions at a much lower cost than other SVGD variants. We also showed that h-SVGD maintains its performance on high dimensional inference tasks, whilst improving variance estimation without the additional computational cost required of other SVGD variants. One interesting direction for future research is to find a principled method of scaling the repulsive kernel, using the theoretical analysis in (Zhuo et al., 2018) as a starting point. Another avenue is to further develop the theory of h-SVGD in the general hybrid kernel setting.

REFERENCES

- Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2005.
- Nachman Aronszajn. Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68(3):337–404, 1950.
- Jimmy Ba, Murat A Erdogdu, Marzyeh Ghassemi, Taiji Suzuki, Shengyang Sun, Denny Wu, and Tianzong Zhang. Towards characterizing the high-dimensional bias of kernel-based particle inference algorithms. In *Second Symposium on Advances in Approximate Bayesian Inference*, 2019.
- Jimmy Ba, Murat A Erdogdu, Marzyeh Ghassemi, Shengyang Sun, Taiji Suzuki, Denny Wu, and Tianzong Zhang. Understanding the variance collapse of SVGD in high dimensions. In *International Conference on Learning Representations*, 2021.
- Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media, 2011.
- Peng Chen and Omar Ghattas. Projected Stein variational gradient descent. *Advances in Neural Information Processing Systems*, 33:1947–1958, 2020.
- Sinho Chewi, Thibaut Le Gouic, Chen Lu, Tyler Maunu, and Philippe Rigollet. SVGD as a kernelized Wasserstein gradient flow of the chi-squared divergence. *Advances in Neural Information Processing Systems*, 33:2098–2109, 2020.
- Francesco D’Angelo and Vincent Fortuin. Annealed Stein variational gradient descent. In *Third Symposium on Advances in Approximate Bayesian Inference*, 2021.
- Francesco D’Angelo, Vincent Fortuin, and Florian Wenzel. On Stein variational neural network ensembles. *arXiv preprint arXiv:2106.10760*, 2021.
- Andrew Duncan, Nikolas Nüsken, and Lukasz Szpruch. On the geometry of Stein variational gradient descent. *Journal of Machine Learning Research*, 24:1–39, 2023.
- Wenbo Gong, Yingzhen Li, and José Miguel Hernández-Lobato. Sliced kernelized Stein discrepancy. In *International Conference on Learning Representations*, 2021.
- Anna Korba, Adil Salim, Michael Arbel, Giulia Luise, and Arthur Gretton. A non-asymptotic analysis for Stein variational gradient descent. *Advances in Neural Information Processing Systems*, 33:4672–4682, 2020.
- Chang Liu and Jun Zhu. Riemannian Stein variational gradient descent for Bayesian inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Qiang Liu. Stein variational gradient descent as gradient flow. *Advances in Neural Information Processing Systems*, 30, 2017.
- Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose Bayesian inference algorithm. *Advances in Neural Information Processing Systems*, 29, 2016.
- Qiang Liu, Jason Lee, and Michael Jordan. A kernelized Stein discrepancy for goodness-of-fit tests. In *International Conference on Machine Learning*, pp. 276–284. PMLR, 2016.
- Xing Liu, Harrison Zhu, Jean-François Ton, George Wynne, and Andrew Duncan. Grassmann Stein variational gradient descent. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151, pp. 2002–2021. PMLR, 2022.
- Jianfeng Lu, Yulong Lu, and James Nolen. Scaling limit of the Stein variational gradient descent: The mean field regime. *SIAM Journal on Mathematical Analysis*, 51(2):648–671, 2019.
- Adil Salim, Lukang Sun, and Peter Richtarik. A convergence theory for SVGD in the population limit under Talagrand’s inequality T1. In *International Conference on Machine Learning*, pp. 19139–19152. PMLR, 2022.

- Filippo Santambrogio. Optimal transport for applied mathematicians. *Birkäuser, NY*, 55(58-63):94, 2015.
- Ingo Steinwart and Andreas Christmann. *Support vector machines*. Springer Science & Business Media, 2008.
- Dilin Wang, Zhe Zeng, and Qiang Liu. Stein variational message passing for continuous graphical models. In *International Conference on Machine Learning*, pp. 5219–5227. PMLR, 2018.
- Dilin Wang, Ziyang Tang, Chandrajit Bajaj, and Qiang Liu. Stein variational gradient descent with matrix-valued kernels. *Advances in Neural Information Processing Systems*, 32, 2019.
- Jingwei Zhuo, Chang Liu, Jiaxin Shi, Jun Zhu, Ning Chen, and Bo Zhang. Message passing Stein variational gradient descent. In *International Conference on Machine Learning*, pp. 6018–6027. PMLR, 2018.

A PROOFS

Lemma A.1. *Under assumptions (A1), (A4), (B1), (B2) the map*

$$(z, \mu) \mapsto E(z, \mu) := \int_{\mathcal{X}} -k_1(x, z) \nabla V(x) + \nabla_x k_2(x, z) d\mu(x)$$

is L -Lipschitz. That is,

$$\|E(z, \mu) - E(z', \mu')\|_2 \leq L(\|z - z'\|_2 + W_2(\mu, \mu'))$$

where $L > 0$ depends on k_1 , k_2 and V .

Proof. Largely following the proof of Lemma 14 (Korba et al., 2020), choosing an optimal coupling γ of μ and μ' ,

$$\begin{aligned} \|E(z, \mu) - E(z', \mu')\|_2 &\leq \|\mathbb{E}_\gamma [\nabla V(x)(k_1(x, z) - k_1(x', z'))] \\ &\quad + \mathbb{E}_\gamma [(\nabla V(x') - \nabla V(x))k_1(x', z')] \\ &\quad + \mathbb{E}_\gamma [\nabla_x k_2(x, z) - \nabla_{x'} k_2(x', z')] \| \\ &\leq D\mathbb{E}_\gamma [\|x - x'\|_2 + \|z - z'\|_2] \\ &\quad + BM\mathbb{E}_\gamma [\|x - x'\|_2] \\ &\quad + D\mathbb{E}_\gamma [\|x - x'\|_2 + \|z - z'\|_2] \\ &\leq (2D + BM)(\|z - z'\|_2 + W_2(\mu, \mu')). \end{aligned}$$

Note that the second term is bounded using the relaxed Assumption (B2) and there is no need to require that $|V|$ is bounded by a constant. \square

Proof of Proposition 3.1. This follows identically to the proof of (Korba et al., 2020, Proposition 7) with Lemma A.1 in place of Lemma 14. \square

Proof of Proposition 3.2. The proof largely follows those of (Lu et al., 2019, Theorems 2.4 and 3.2) with some minor adjustments. Notably, after fixing $r > 0$ and defining

$$Y_r := \left\{ u \in Y : \sup_{x \in \mathcal{X}} |u(x) - x| < r \right\}$$

and the complete metric space

$$\begin{aligned} S_r &:= C([0, T_0]; Y_r), \\ d_S(u, v) &:= \sup_{t \in [0, T_0]} d_Y(u(t), v(t)) \end{aligned}$$

for some sufficiently small T_0 (to be determined later), the operator $\mathcal{G} : u(t, \cdot) \mapsto \mathcal{G}(u)(t, \cdot)$ must be modified to act on $u \in S_r$ via

$$\begin{aligned} \mathcal{G}(u)(t, x) &= x - \int_0^t \int_{\mathcal{X}} \nabla K_2(u(s, x) - u(s, x')) \nu(dx') ds \\ &\quad - \int_0^t \int_{\mathcal{X}} K_1(u(s, x) - u(s, x')) \nabla V(u(s, x')) \nu(dx') ds. \end{aligned}$$

Note that we use \mathcal{G} instead of the \mathcal{F} used in (Lu et al., 2019) to avoid confusion between the functional \mathcal{F} defined in (17). The same techniques of (Lu et al., 2019) are sufficient to establish the required bounds to show that \mathcal{G} is a contraction on S_r for sufficiently small T_0 . Note that Assumptions (B1), (A2) and (A3) are used to establish this. So the unique fixed point $X(\cdot, \cdot; \nu) \in S_r$ of \mathcal{G} solves (13) in the interval $[0, T_0]$.

The $\min(\|\nabla K_1\|_\infty, \|\nabla K_2\|_\infty)$ term emerges because the telescoping in (Lu et al., 2019, Equation (3.8)) can be performed with either kernel. The remainder of the proof follows (Lu et al., 2019, Theorems 2.4 and 3.2). \square

Proof of Proposition 3.3. Following (Santambrogio, 2015, Definition 7.12), a functional derivative of \mathcal{F} is a measurable function $\frac{\delta \mathcal{F}}{\delta \rho}(\rho)$ satisfying

$$\left. \frac{d}{d\epsilon} \mathcal{F}(\rho + \epsilon \chi) \right|_{\epsilon=0} = \int \frac{\delta \mathcal{F}}{\delta \rho}(\rho) d\chi$$

for all perturbations $\chi = \tilde{\rho} - \rho$ with $\tilde{\rho} \in L_c^\infty(\mathcal{X}) \cap \mathcal{P}(\mathcal{X})$. If it exists, $\frac{\delta \mathcal{F}}{\delta \rho}(\rho)$ is unique up to additive constants.

We first compute

$$\begin{aligned} & \left. \frac{d}{d\epsilon} \mathcal{F}(\rho + \epsilon \chi) \right|_{\epsilon=0} \\ &= \left. \frac{d}{d\epsilon} \left(\int c \log(\rho(\mathbf{x}) + \epsilon \chi(\mathbf{x})) \rho(\mathbf{x}) d\mathbf{x} + \epsilon \int c \log(\rho(\mathbf{x}) + \epsilon \chi(\mathbf{x})) \chi(\mathbf{x}) d\mathbf{x} \right. \right. \\ & \quad \left. \left. - \int \log \pi(\mathbf{x}) \rho(\mathbf{x}) d\mathbf{x} - \epsilon \int \log \pi(\mathbf{x}) \chi(\mathbf{x}) d\mathbf{x} \right) \right|_{\epsilon=0} \\ &= \left(\int c \frac{\chi(\mathbf{x})}{\rho(\mathbf{x}) + \epsilon \chi(\mathbf{x})} \rho(\mathbf{x}) d\mathbf{x} + \int c \log(\rho(\mathbf{x}) + \epsilon \chi(\mathbf{x})) \chi(\mathbf{x}) d\mathbf{x} \right. \\ & \quad \left. + \epsilon \int c \frac{\chi(\mathbf{x})}{\rho(\mathbf{x}) + \epsilon \chi(\mathbf{x})} \chi(\mathbf{x}) d\mathbf{x} - \int \log \pi(\mathbf{x}) \chi(\mathbf{x}) d\mathbf{x} \right) \Big|_{\epsilon=0} \\ &= \int c \log \rho(\mathbf{x}) - \log \pi(\mathbf{x}) \chi(\mathbf{x}) d\mathbf{x} + c \int \chi(\mathbf{x}) d\mathbf{x}. \end{aligned} \tag{29}$$

Since $\rho, \tilde{\rho} \in \mathcal{P}(\mathcal{X})$, the final integral is zero. So the functional gradient of \mathcal{F} is

$$\frac{\delta \mathcal{F}}{\delta \rho}(\rho) = c \log \rho - \log \pi.$$

Its Wasserstein gradient is then

$$\nabla_W \mathcal{F}(\rho) = c \nabla \log \rho - \nabla \log \pi.$$

Since $k_2 = ck_1$, equation (10) can be written as

$$\begin{aligned} \frac{d\mathbf{x}}{dt} &= \int (k_1(\mathbf{x}', \mathbf{x}) \nabla_{\mathbf{x}'} \log \pi(\mathbf{x}') + c \nabla_{\mathbf{x}'} k_1(\mathbf{x}', \mathbf{x})) \rho(\mathbf{x}') d\mathbf{x}' \\ &= \int k_1(\mathbf{x}', \mathbf{x}) \nabla_{\mathbf{x}'} \log \pi(\mathbf{x}') \rho(\mathbf{x}') - ck_1(\mathbf{x}', \mathbf{x}) \nabla_{\mathbf{x}'} \rho(\mathbf{x}') d\mathbf{x}' \\ &= \int k_1(\mathbf{x}', \mathbf{x}) \nabla_{\mathbf{x}'} \log \pi(\mathbf{x}') \rho(\mathbf{x}') - ck_1(\mathbf{x}', \mathbf{x}) \nabla_{\mathbf{x}'} \log \rho(\mathbf{x}') \rho(\mathbf{x}') d\mathbf{x}' \\ &= T_{k_1, \rho}(\nabla \log \pi - c \nabla \log \rho)(\mathbf{x}) \\ &= -T_{k_1, \rho}(\nabla_W \mathcal{F}(\rho))(\mathbf{x}). \end{aligned}$$

The continuity equation $\partial_t \rho_t + \nabla \cdot (\rho_t \frac{d\mathbf{x}}{dt}) = 0$ then becomes

$$\begin{aligned} \partial_t \rho_t &= \nabla \cdot (\rho_t T_{k_1, \rho}(\nabla_W \mathcal{F}(\rho))(\mathbf{x})) \\ &= \nabla \cdot (\rho_t T_{k_1, \rho_t}(c \nabla \log \rho_t - \nabla \log \pi)). \end{aligned}$$

□

Proof of Corollary (3.4). A measure ρ^* satisfying $c \nabla \log \rho^* = \nabla \log \pi$ will be a fixed point because this would imply $\nabla_W \mathcal{F}(\rho) = 0$. Solving this for ρ^* , we have

$$\begin{aligned} c \nabla \log \rho^*(\mathbf{x}) &= \nabla \log \pi(\mathbf{x}) \\ c \log \rho^*(\mathbf{x}) &= \log \pi(\mathbf{x}) + A \\ \log(\rho^*(\mathbf{x})^c) &= \log(e^A \pi(\mathbf{x})) \\ \rho^*(\mathbf{x})^c &= e^A \pi(\mathbf{x}) \\ \rho^*(\mathbf{x}) &= e^{A/c} \pi(\mathbf{x})^{1/c} \end{aligned}$$

for some $A \in \mathbb{R}$.

□

Proof of Proposition 3.5. Using (19), integration by parts, and the fact that $T_{k_1, \rho} : L^2(\rho_t)^d \rightarrow \mathcal{H}_1^d$ is the adjoint of the inclusion $\iota : \mathcal{H}_1^d \rightarrow L^2(\rho_t)^d$,

$$\begin{aligned}
\frac{d}{dt} \mathcal{F}(\rho_t) &= \frac{d}{dt} \int \rho_t(\mathbf{x}) (c \log \rho_t(\mathbf{x}) - \log \pi(\mathbf{x})) d\mathbf{x} \\
&= \int (c \log \rho_t(\mathbf{x}) - \log \pi(\mathbf{x})) \frac{\partial}{\partial t} \rho_t(\mathbf{x}) + c \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x} \\
&= \int (c \log \rho_t(\mathbf{x}) - \log \pi(\mathbf{x})) \nabla \cdot (\rho_t(\mathbf{x}) T_{k_1, \rho_t} (c \nabla \log \rho_t(\mathbf{x}) - \nabla \log \pi(\mathbf{x}))) d\mathbf{x} \\
&\quad + \int c \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x} \\
&= - \int \nabla (c \log \rho_t(\mathbf{x}) - \log \pi(\mathbf{x})) \cdot T_{k_1, \rho_t} (c \nabla \log \rho_t - \nabla \log \pi)(\mathbf{x}) \rho_t(\mathbf{x}) d\mathbf{x} \\
&\quad + \int c \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x} \\
&= - \langle c \nabla \log \rho_t - \nabla \log \pi, T_{k_1, \rho_t} (c \nabla \log \rho_t - \nabla \log \pi) \rangle_{L^2(\rho_t)^d} \\
&\quad + \int c \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x} \\
&= - \|T_{k_1, \rho_t} (c \nabla \log \rho_t - \nabla \log \pi)\|_{\mathcal{H}_1^d}^2 + c \int \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x} \\
&= -c^2 \|T_{k_1, \rho_t} (\nabla \log \rho_t - \nabla \log \rho^*)\|_{\mathcal{H}_1^d}^2 + c \int \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x}. \tag{30}
\end{aligned}$$

Also,

$$\begin{aligned}
\mathcal{F}(\rho_t) &= \mathbb{E}_{\mathbf{x} \sim \rho_t} [c \log \rho_t(\mathbf{x}) - \log \pi(\mathbf{x})] \\
&= \mathbb{E}_{\mathbf{x} \sim \rho_t} [c \log \rho_t(\mathbf{x}) - \log (A \rho^*(\mathbf{x})^c)] \\
&= c \mathbb{E}_{\mathbf{x} \sim \rho_t} [\log \rho_t(\mathbf{x}) - \log \rho^*(\mathbf{x})] + \log(A) \\
&= c \text{KL}(\rho_t \parallel \rho^*) + \log(A) \tag{31}
\end{aligned}$$

for some $A \in \mathbb{R}$, we have $\frac{1}{c} \frac{d}{dt} \mathcal{F}(\rho_t) = \frac{d}{dt} \text{KL}(\rho_t \parallel \rho^*)$. Therefore,

$$\begin{aligned}
\frac{d}{dt} \text{KL}(\rho_t \parallel \rho^*) &= -\frac{1}{c} \|T_{k_1, \rho_t} (c \nabla \log \rho_t - \nabla \log \pi)\|_{\mathcal{H}_1^d}^2 + \int \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x} \\
&= -c \|T_{k_1, \rho_t} (\nabla \log \rho_t - \nabla \log \pi^{1/c})\|_{\mathcal{H}_1^d}^2 + \int \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x} \\
&= -c \|T_{k_1, \rho_t} (\nabla \log \rho_t - \nabla \log \rho^*)\|_{\mathcal{H}_1^d}^2 + \int \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x}.
\end{aligned}$$

To simplify notation in the calculations below, set $u = T_{k_1, \rho_t} (c \nabla \log \rho_t)$ and $v = T_{k_1, \rho_t} (\nabla \log \pi)$. Recall the identity

$$\langle u, u - v \rangle = \frac{1}{2} (\|u\|^2 + \|u - v\|^2 - \|v\|^2). \tag{32}$$

Now we apply the multivariate chain rule to the remainder term along with the fact about the adjoint of T_{k_1, ρ_t} , the identity in (32), and the triangle inequality. This gives

$$\begin{aligned}
\int \rho_t(\mathbf{x}) \frac{\partial}{\partial t} \log \rho_t(\mathbf{x}) d\mathbf{x} &= \int \rho_t(\mathbf{x}) \nabla \log \rho_t(\mathbf{x}) \cdot \frac{d\mathbf{x}}{dt} d\mathbf{x} \\
&= - \int \rho_t(\mathbf{x}) \nabla \log \rho_t(\mathbf{x}) \cdot T_{k_1, \rho_t}(c \nabla \log \rho_t - \nabla \log \pi)(\mathbf{x}) d\mathbf{x} \quad (33) \\
&= - \langle \nabla \log \rho_t, T_{k_1, \rho_t}(c \nabla \log \rho_t - \nabla \log \pi) \rangle_{L^2(\rho_t)^d} \\
&= - \langle T_{k_1, \rho_t}(\nabla \log \rho_t), T_{k_1, \rho_t}(c \nabla \log \rho_t - \nabla \log \pi) \rangle_{\mathcal{H}_1^d} \\
&= - \frac{1}{c} \langle u, u - v \rangle_{\mathcal{H}_1^d} \\
&= \frac{1}{2c} \left(-\|u\|_{\mathcal{H}_1^d}^2 + \|v\|_{\mathcal{H}_1^d}^2 - \|u - v\|_{\mathcal{H}_1^d}^2 \right) \\
&\leq \frac{1}{2c} \left(-\|u\|_{\mathcal{H}_1^d}^2 + \|u - v\|_{\mathcal{H}_1^d}^2 + \|u\|_{\mathcal{H}_1^d}^2 - \|u - v\|_{\mathcal{H}_1^d}^2 \right) \\
&= 0.
\end{aligned}$$

The final statement that $\frac{d}{dt} \mathcal{F}(\rho^*) = 0$ follows from equations (30) and (33) by observing that $c \nabla \log \rho^* - \nabla \log \pi = 0$. \square

Proof of Proposition 3.6. This follows from applying Theorem 3.3 (Liu, 2017) with ρ^* instead of the target, then substituting in (31). \square

Proof of Proposition 3.7. To recover the continuity equation, apply the convolution theorem to (22)

$$\begin{aligned}
\mathcal{F}^{-1} \left(\frac{\mathcal{F}(K_2)}{\mathcal{F}(K_1)} \cdot \mathcal{F}(\nabla \rho) \right) (x) &= r(x; \rho) \\
\mathcal{F}(K_2)(\omega) \cdot \mathcal{F}(\nabla \rho)(\omega) &= \mathcal{F}(K_1)(\omega) \cdot \mathcal{F}(r)(\omega) \\
\mathcal{F}(K_2 * \nabla \rho)(\omega) &= \mathcal{F}(K_1 * r)(\omega) \\
(K_2 * \nabla \rho)(x) &= (K_1 * r)(x) \\
\int K_2(x - y) \nabla \log \rho(y) \rho(y) dy &= \int K_1(x - y) \frac{r(y)}{\rho(y)} \rho(y) dy \\
(T_{k_2, \rho} \nabla \log \rho)(x) &= (T_{k_1, \rho}(r/\rho))(x)
\end{aligned}$$

Equation (10) can now be rewritten as

$$\begin{aligned}
\frac{dx}{dt} &= \int (k_1(x', x) \nabla_{x'} \log \pi(x') + \nabla_{x'} k_2(x', x)) \rho(x') dx' \\
&= \int k_1(x', x) \nabla_{x'} \log \pi(x') \rho(x') - k_2(x', x) \nabla_{x'} \rho(x') dx' \\
&= \int (k_1(x', x) \nabla_{x'} \log \pi(x') - k_2(x', x) \nabla_{x'} \log \rho(x')) \rho(x') dx' \\
&= (T_{k_1, \rho} \nabla \log \pi)(x) - (T_{k_2, \rho} \nabla \log \rho)(x) \\
&= (T_{k_1, \rho}(\nabla \log \pi - \nabla R(\cdot; \rho)))(x),
\end{aligned}$$

where $\nabla R(x; \rho) := \frac{r(x; \rho)}{\rho(x)}$, yielding the continuity equation (23).

As in the proof of Proposition 3.3, let $\chi = \tilde{\rho} - \rho$ with $\tilde{\rho} \in L_c^\infty(\mathcal{X}) \cap \mathcal{P}(\mathcal{X})$. We first compute

$$\begin{aligned}
& \left. \frac{d}{d\epsilon} \mathcal{F}(\rho + \epsilon\chi) \right|_{\epsilon=0} \\
&= \left. \frac{d}{d\epsilon} \left(\int R(x; \rho + \epsilon\chi) d\rho(x) + \epsilon \int R(x; \rho + \epsilon\chi) d\chi(x) \right. \right. \\
&\quad \left. \left. - \int \log \pi(x) d\rho(x) - \epsilon \int \log \pi(x) d\chi(x) \right) \right|_{\epsilon=0} \\
&= \left(\frac{d}{d\epsilon} \int R(x; \rho + \epsilon\chi) d\rho(x) + \int R(x; \rho + \epsilon\chi) d\chi(x) + \epsilon \frac{d}{d\epsilon} \int R(x; \rho + \epsilon\chi) d\chi(x) \right. \\
&\quad \left. - \int \log \pi(x) d\chi(x) \right) \Big|_{\epsilon=0} \\
&= \int R(x; \rho) - \log \pi(x) d\chi(x) \\
&\quad + \left(\int \frac{\partial}{\partial \epsilon} R(x; \rho + \epsilon\chi) d\rho(x) + \epsilon \int \frac{\partial}{\partial \epsilon} R(x; \rho + \epsilon\chi) d\chi(x) \right) \Big|_{\epsilon=0} \\
&= \int R(x; \rho) - \log \pi(x) d\chi(x) + \int \frac{\partial}{\partial \epsilon} R(x; \rho + \epsilon\chi) d(\rho + \epsilon\chi)(x) \Big|_{\epsilon=0}. \tag{34}
\end{aligned}$$

By assumption, the remainder term above is equal to zero and the functional derivative of \mathcal{F} is therefore

$$\frac{\delta \mathcal{F}}{\delta \rho}(\rho) = R_1(x; \rho) - \log \pi(x),$$

so its Wasserstein gradient is

$$\begin{aligned}
\nabla_W \mathcal{F}(\rho) &= \nabla \frac{\delta \mathcal{F}}{\delta \rho}(\rho) \\
&= R(x; \rho) - \nabla \log \pi(x).
\end{aligned}$$

□

Proof of Proposition 3.8. A direct computation of (22) along with the definitions of r and R yields

$$\begin{aligned}
\frac{\mathcal{F}(K_2)(\omega)}{\mathcal{F}(K_1)(\omega)} \cdot \mathcal{F}(\nabla \rho^*)(\omega) &= \sqrt{\frac{h_2}{h_1}} \cdot \exp(-2\pi^2 \omega^2 \Delta h) \cdot 2\pi i \omega \exp(-2\pi^2 \omega^2 \sigma^2) \\
&= \sqrt{\frac{h_2}{h_1}} \cdot 2\pi i \omega \cdot \exp(-2\pi^2 \omega^2 (\sigma^2 + \Delta h)) \\
r(x; \rho^*) &= -\sqrt{\frac{h_2}{h_1}} \cdot \frac{1}{\sqrt{2\pi}} \cdot \frac{x}{(\sigma^2 + \Delta h)^{3/2}} \cdot \exp\left(-\frac{x^2}{2(\sigma^2 + \Delta h)}\right) \\
\nabla R(x; \rho^*) &= -\sqrt{\frac{h_2}{h_1}} \cdot \frac{\sigma}{(\sigma^2 + \Delta h)^{3/2}} \cdot x \exp\left(-\frac{x^2}{2(\sigma^2 + \Delta h)}\right) \cdot \exp\left(\frac{x^2}{2\sigma^2}\right) \\
&= -\sqrt{\frac{h_2}{h_1}} \cdot \frac{\sigma}{(\sigma^2 + \Delta h)^{3/2}} \cdot x \exp\left(-\frac{\Delta h x^2}{2\sigma^2(\sigma^2 + \Delta h)}\right).
\end{aligned}$$

Now computing $\nabla \log \pi$, using some $A \in \mathbb{R}$ for the normalising constant, we have

$$\begin{aligned}\pi(x) &= A \exp\left(-\alpha \exp\left(\frac{x^2}{2\beta}\right)\right) \\ \log \pi(x) &= \log(A) - \alpha \exp\left(\frac{x^2}{2\beta}\right) \\ \nabla \log \pi(x) &= -\frac{\alpha x}{\beta} \exp\left(\frac{x^2}{2\beta}\right) \\ &= -\sqrt{\frac{h_2}{h_1}} \cdot \frac{\sigma}{(\sigma^2 + \Delta h)^{3/2}} \cdot x \exp\left(-\frac{\Delta h x^2}{2\sigma^2(\sigma^2 + \Delta h)}\right).\end{aligned}$$

Since $\nabla R(x; \rho^*) = \nabla \log \pi(x)$, equation (23) implies that ρ^* is a fixed point. \square

B ADDITIONAL EXPERIMENTAL RESULTS AND DETAILS

B.1 BAYESIAN NEURAL NETWORK

The results presented in Section 4.2 follow the settings of (Liu & Wang, 2016). In particular, we use normal priors for the network weights and Gamma priors for the inverse covariances. There is one hidden layer with 50 units for most datasets, Protein and Year being the exceptions with 100 units each. The datasets are randomly partitioned into 90% for training and 10% for testing with results averaged over 20 trials, Protein and Year being the exceptions with 5 trials and 1 trial respectively. The number of particles in each case is 20, the activation function is $\text{RELU}(x) = \max(0, x)$, the number of iterations is 2000, and the mini-batch size is 100 for all datasets except for Year, which uses a mini-batch size of 1000.

We recreate Tables 1 and 2 below, this time comparing h-SVGD against S-SVGD (Gong et al., 2021). Table 3 shows that h-SVGD outperforms S-SVGD in mitigation of variance collapse on all but one dataset. We note that h-SVGD achieves this with a significantly faster runtime for all datasets, which is due to S-SVGD requiring additional optimisation of the projection matrix. Table 4 shows that S-SVGD and h-SVGD are comparable in both RMSE and LL metrics. The h-SVGD algorithm achieves a better LL score on more datasets, but S-SVGD achieves a better RMSE score on more datasets.

Table 3: DAMV and runtime in seconds of S-SVGD and h-SVGD.

Dataset	DAMV		Runtime (seconds)	
	S-SVGD	h-SVGD	S-SVGD	h-SVGD
Boston	0.035 \pm 0.002	0.087 \pm 0.010	208 \pm 13	28.5 \pm 0.9
Concrete	0.070 \pm 0.004	0.102 \pm 0.006	148 \pm 56	28.7 \pm 1.3
Energy	0.053 \pm 0.005	0.106 \pm 0.011	156 \pm 28	30.7 \pm 2.2
Kin8nm	0.083 \pm 0.002	0.093 \pm 0.003	141 \pm 3.9	36.0 \pm 1.3
Naval	0.070 \pm 0.021	0.068 \pm 0.011	237 \pm 11	34.5 \pm 0.9
Combined	0.118 \pm 0.005	0.138 \pm 0.006	116 \pm 18	36.7 \pm 2.4
Protein	0.057 \pm 0.006	0.084 \pm 0.001	390 \pm 20	72.8 \pm 1.7
Wine	0.029 \pm 0.002	0.075 \pm 0.005	210 \pm 10	29.8 \pm 1.1
Yacht	0.066 \pm 0.009	0.121 \pm 0.012	97.9 \pm 1.2	30.0 \pm 1.3
Year	0.012 \pm NA	0.012 \pm NA	12488 \pm NA	666 \pm NA

Table 4: Average RMSE and LL of SSGVD and h-SVGD evaluated on the test dataset.

Dataset	Test RMSE		Test LL	
	S-SVGD	h-SVGD	S-SVGD	h-SVGD
Boston	3.024 ± 0.604	3.001 ± 0.584	-2.088 ± 0.322	-1.988 ± 0.221
Concrete	5.073 ± 0.522	5.210 ± 0.529	-2.563 ± 0.239	-2.535 ± 0.179
Energy	0.923 ± 0.123	1.040 ± 0.128	-0.631 ± 0.162	-0.805 ± 0.104
Kin8nm	0.084 ± 0.003	0.090 ± 0.003	0.232 ± 0.135	0.468 ± 0.090
Naval	0.003 ± 0.000	0.004 ± 0.000	-0.624 ± 0.161	-0.090 ± 0.105
Combined	4.028 ± 0.22	4.057 ± 0.218	-2.335 ± 0.066	-2.354 ± 0.052
Protein	4.581 ± 0.026	4.600 ± 0.026	-2.526 ± 0.045	-2.456 ± 0.017
Wine	0.676 ± 0.051	0.626 ± 0.045	-1.261 ± 0.172	-0.750 ± 0.097
Yacht	1.664 ± 0.607	1.861 ± 0.662	-0.788 ± 0.511	-0.813 ± 0.227
Year	$8.922 \pm \text{NA}$	$8.689 \pm \text{NA}$	$-2.940 \pm \text{NA}$	$-2.872 \pm \text{NA}$