Information Newton's flow: second-order optimization method in probability space

Yifei Wang

WANGYF18@STANFORD.EDU

Department of Electrical Engineering Stanford University Stanford, CA 94305-9505, USA

Wuchen Li

Department of Mathematics University of California Los Angeles, CA 90095-1555, USA WCLI@MATH.UCLA.EDU

Abstract

We introduce a framework for Newton's flows in probability space with information metrics, named information Newton's flows. Here two information metrics are considered, including both the Fisher-Rao metric and the Wasserstein-2 metric. A known fact is that overdamped Langevin dynamics correspond to Wasserstein gradient flows of Kullback-Leibler (KL) divergence. Extending this fact to Wasserstein Newton's flows, we derive Newton's Langevin dynamics. We provide examples of Newton's Langevin dynamics in both one-dimensional space and Gaussian families. For the numerical implementation, we design sampling efficient variational methods in affine models and reproducing kernel Hilbert space (RKHS) to approximate Wasserstein Newton's directions. We also establish convergence results of the proposed information Newton's method with approximated directions. Several numerical examples from Bayesian sampling problems are shown to demonstrate the effectiveness of the proposed method.

Keywords: Optimal transport; Information geometry; Langvien dynamics; Information Newton's flow; Newton's Langvien dynamics.

1. Introduction

Optimization problems in probability space are of great interest in inverse problems, information science, physics, and scientific computing, with applications in machine learning (Amari, 2016; Stuart, 2010; Liu, 2017; Amari, 1998; Villani, 2003). One typical problem here comes from Bayesian inference, which provides an optimal probability formulation for learning models from observed data. Given a prior distribution, the problem is to generate samples from a (target) posterior distribution (Stuart, 2010). From an optimization perspective, such a problem often refers to minimizing an objective function, such as the Kullback-Leibler (KL) divergence, in the probability space. The update relates to finding a sampling representation for the evolution of the probability.

In practice, one often needs to transfer probability optimization problems into samplingbased formulations, and then design efficient updates in the form of samples. Here first-order methods, such as gradient descent methods, play essential roles. We notice that gradient directions for samples rely on the metric over the probability space, which reflects the change of objective/loss functions. In practice, there are several important metrics, often

^{©2020} Yifei Wang and Wuchen Li.

License: CC-BY 4.0, see https://creativecommons.org/licenses/by/4.0/.

named information metrics from information geometry and optimal transport, including the Fisher-Rao metric (Amari, 1998) and the Wasserstein-2 metric (in short, Wasserstein metric) (Lafferty, 1988; Otto, 2001). In literature, along with a given information metric, the probability space can be viewed as a Riemannian manifold, named density manifold (Lafferty, 1988).

For the Fisher-Rao metric, its gradient flow, known as birth-death dynamics, are important in modeling population games and designing evolutionary dynamics (Amari, 2016). It is also important for optimization problems in discrete probability (Malagò and Pistone, 2014) and machine learning (Ollivier et al., 2017). Recently, the Fisher-Rao gradient has also been applied for accelerating Bayesian sampling problems in continuous sample space (Lu et al., 2019). The Fisher-Rao gradient direction also inspires the design of learning algorithms for probability models. Several optimization methods in machine learning approximate the Fisher-Rao gradient direction, including the Kronecker-factored Approximate Curvature (K-FAC) (Martens and Grosse, 2015) method and adaptive estimates of lower-order moments (Adam) method (Kingma and Ba, 2014).

For the Wasserstein metric, its gradient direction deeply connects with stochastic differential equations and the associated Markov chain Monte Carlo methods (MCMC). An important fact is that the Wasserstein gradient of KL divergence forms the Kolmogorov forward generator of overdamped Langevin dynamics (Jordan et al., 1998). Hence, many MCMC methods can be viewed as Wasserstein gradient descent methods. In recent years, there are also several generalized Wasserstein metrics, such as Stein metric (Liu and Wang, 2016; Liu, 2017), Hessian transport (mobility) metrics (Carrillo et al., 2010; Dolbeault et al., 2009; Li and Ying, 2019) and Kalman-Wasserstein metric (Garbuno-Inigo et al., 2019). These metrics introduce various first-order methods with sampling efficient properties. For instance, the Stein variational gradient descent (Liu and Wang, 2016, SVGD) introduces a kernelized interacting Langevin dynamics. The Kalman-Wasserstein metric introduces a particular mean-field interacting Langevin dynamics (Garbuno-Inigo et al., 2019), known as ensemble Kalman sampling. On the other hand, many approaches design fast algorithms on modified Langevin dynamics. These methods can also be viewed and analyzed by the modified Wasserstein gradient descent, see details in (Ma et al., 2019; Simsekli et al., 2016; Li, 2019). By viewing sampling as optimization problems in the probability space, many efficient sampling algorithms are inspired by classical optimization methods. E.g., Bernton (2018); Wibisono (2019) apply the operator splitting technique to improve the unadjusted Langevin algorithm. Liu et al. (2018); Taghvaei and Mehta (2019); Wang and Li (2019) study Nesterov's accelerated gradient methods in probability space.

In optimization, the Newton's method is a fundamental second-order method to accelerate optimization computations. For optimization problems in probability space, several natural questions arise: Can we systematically design Newton's methods to accelerate sampling related optimization problems? What is the Newton's flow in probability space under information metrics? Focusing on the Wasserstein metric, can we extend the relation between Wasserstein gradient flow of KL divergence and Langevin dynamics? In other words, what is the Wasserstein Newton's flow of KL divergence and which Langevin dynamics does it corresponds to?

In this paper, following (Li, 2018; Wang and Li, 2019), we complete these questions. We derive Newton's flows in probability space with general information metrics. By studying

these Newton's flows, we provide the convergence analysis.Focusing on Wasserstein Newton's flows of KL divergence, we derive several analytical examples in one-dimensional space and Gaussian families. Besides, we design two algorithms as particle implementations of Wasserstein Newton's flows in high dimensional sample space. This is to restrict the dual variable (cotangent vector) associated with Newton's direction into either finite-dimensional affine function space or RKHS. A hybrid update of Newton's direction and gradient direction is also introduced. For the concreteness of presentation, we demonstrate the Wasserstein Newton's flow of KL divergence in Theorem 1.

Theorem 1 (Wasserstein Newton's flow of KL divergence) For a density $\rho^*(x) \propto \exp(-f(x))$, where f is a given function, denote the KL divergence between ρ and ρ^* by

$$D_{KL}(\rho \| \rho^*) = \int \rho \log \frac{\rho}{e^{-f}} dx - \log Z, \qquad (1)$$

where $Z = \int \exp(-f(x)) dx$. Then the Wasserstein Newton's flow of KL divergence follows

$$\partial_t \rho_t + \nabla \cdot (\rho_t \nabla \Phi_t^{\text{Newton}}) = 0, \qquad (2)$$

where Φ_t^{Newton} satisfies the following equation

$$\nabla^2 : (\rho_t \nabla^2 \Phi_t) - \nabla \cdot (\rho_t \nabla^2 f \nabla \Phi_t) - \nabla \cdot (\rho_t \nabla f) - \Delta \rho_t = 0.$$
(3)

Here we notice that Φ_t^{Newton} is the solution to the Wasserstein Newton's direction equation (3). In Figure 1, we provide a sampling (particle) formulation of Wasserstein Newton's flows. We compare formulations among Wasserstein Newton's flows, Wasserstein gradient flows and overdamped Langevin dynamics.

Density formulation $\partial_t \rho_t = \nabla \cdot (\rho_t \nabla f) + \Delta \rho_t$ $\partial_t \rho_t = -\nabla \cdot (\rho_t \nabla \Phi_t^{\text{Newton}})$ \downarrow \downarrow Particle formulation $dX_t = -\nabla f(X_t)dt - \nabla \log \rho_t(X_t)dt$ $dX_t = \nabla \Phi_t^{\text{Newton}}(X_t)dt$

Langevin dynamics $dX_t = -\nabla f(X_t)dt + \sqrt{2}dB_t$

Figure 1: The relation among Wasserstein gradient flow, Newton's flow and Langevin dynamics. Our approach derive the particle formulation of Wasserstein Newton's flow of KL divergence.

In literature, second-order methods are developed for optimization problems on Riemannian manifold, see (Smith, 1994; Yang, 2007). Here we are interested in density manifolds,

WANG AND LI

i.e., probability space with information metrics. Compared to known results in Riemannian optimization, we not only develop methods in probability space but also find efficient sampling representations of the algorithms. In discrete probability simplex with the Fisher-Rao metric and exponential family models, the Newton's method has also been studied by Malagò and Pistone (2014), known as the second order method in information geometry. Also, Detommaso et al. (2018); Chen et al. (2019) design second-order methods for the Stein variational gradient descent direction. Our approach generalizes these results to information metrics, especially for the Wasserstein metric. On the other hand, the Newton-type MCMC method has been studied in (Simsekli et al., 2016), known as Hessian Approximated MCMC (HAMCMC) method. The differences between HAMCMC and our proposed Newton's Langevin dynamics can be observed from evolutions in probability space. HAMCMC utilizes the Hessian matrix of logarithm of target density function and derives the associated drift-diffusion process. In density space, it is still a linear local partial differential equation (PDE). Newton's Langevin dynamics apply the Hessian operator of KL divergence based on the Wasserstein metric. In density space, the Wasserstein Newton's flow is a nonlocal PDE. A careful comparison of all related Langevin dynamics in analytical (Appendix C.3) and numerical examples are provided.

We organize this paper as follows. In section 2, we briefly review information metrics and corresponding gradient operators in probability space. We introduce properties of Hessian operators and derive information Newton's flows in section 3. Focusing on Wasserstein Newton's flows of KL divergence, we derive Newton's Langevin dynamics in section 4. Two sampling efficient numerical algorithms of Wasserstein Newton's method are presented in section 5. In section 6, we prove the asymptotic convergence rate of information Newton's method with approximated Newton's direction. Several numerical examples for sampling problems are provided in section 7.

2. Review on Newton's flows and information metrics

In this section, we briefly review Newton's methods and Newton's flows in Euclidean spaces and Riemannian manifolds. Then, we focus on a probability space, in which we introduce information metrics with the associated gradient and Hessian operators. Based on them, we will derive the Newton's flow under information metrics later on. Throughout this paper, we use $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$ to denote the Euclidean inner product and norm in \mathbb{R}^d .

2.1 Finite dimensional Newton's flow

We first briefly review Newton's methods and Newton's flows in Euclidean spaces. Given an objective function $f: \mathbb{R}^d \to \mathbb{R}$, consider an optimization problem:

$$\min_{x \in \mathbb{R}^d} f(x).$$

The update rule of the (damped) Newton's method follows

$$x_{k+1} = x_k + \alpha_k p_k, \quad p_k = -(\nabla^2 f(x_k))^{-1} \nabla f(x_k).$$

Here $\alpha_k > 0$ is a step size and p_k is called the Newton's direction. With $\alpha_k = 1$, we recover the classical Newton's methods. By taking a limit $\alpha_k \to 0$, the Newton's method in

continuous-time, namely Newton's flow, writes

$$\dot{x} = -\left(\nabla^2 f(x)\right)^{-1} \nabla f(x).$$
 (Euclidean Newton's flow)

We next consider an optimization problem on a Riemannian manifold $\mathcal{M} \subset \mathbb{R}^d$. Given an objective function $f: \mathcal{M} \to \mathbb{R}$, consider

$$\min_{x \in \mathcal{M}} f(x).$$

The tangent space $T_x \mathcal{M}$ and the cotangent space $T_x^* \mathcal{M}$ at x are identical to a linear subspace of \mathbb{R}^d . For $p, q \in T_x \mathcal{M}$, let $\langle p, q \rangle_x = p^T \mathcal{G}(x) q$ denote an inner product in tangent space $T_x \mathcal{M}$ at x. Here $\mathcal{G}(x)$ is called the metric tensor, which corresponds to a symmetric semi-positive definite matrix in $\mathbb{R}^{d \times d}$. For the Euclidean case, we can view $T_x \mathcal{M} = T_x^* \mathcal{M} = \mathbb{R}^d$ and $\mathcal{G}(x) = I$, where I is an identity matrix. The Riemannian gradient of f at x is the unique tangent vector v such that the following equality holds for all $p \in T_x \mathcal{M}$.

$$\langle \operatorname{grad} f(x), p \rangle_x = \lim_{\epsilon \to 0} \frac{f(x + \epsilon p) - f(x)}{\epsilon}.$$

The Riemannian Hessian of f at x is a linear mapping from $T_x \mathcal{M}$ to $T_x \mathcal{M}$ defined by

Hess
$$f(x)p = \nabla_p \operatorname{grad} f(x), \quad \forall p \in T_x \mathcal{M}.$$

Here $\nabla_p \operatorname{grad} f(x)$ is the covariant derivative of $\operatorname{grad} f(x)$ w.r.t. the tangent vector p. Detailed definitions of gradient and Hessian operators on a Riemannian manifold can be found in (Huang, 2013, Chapter 1). The update rule of the Newton's method writes

$$x_{k+1} = R_{x_k}(\alpha_k p_k), \quad p_k = -(\text{Hess } f(x_k))^{-1} \operatorname{grad} f(x_k).$$

Here R_{x_k} can be the exponential mapping or the retraction (first-order approximation of the exponential mapping) at x_k . Based on the Riemannian metric of \mathcal{M} , the exponential mapping uniquely maps a tangent vector to a point in \mathcal{M} along the geodesic curve. Different from the Euclidean case, the update of x_{k+1} is based on the (approximated) geodesic curve of \mathcal{M} . In continuous time, the Newton's flow follows

$$\dot{x} = -(\text{Hess } f(x))^{-1} \operatorname{grad} f(x).$$
 (Riemannian Newton's flow)

From now on, we consider optimization problems in probability space. Suppose that sample space Ω is a region in \mathbb{R}^d . Let $\mathcal{F}(\Omega)$ represent the set of smooth functions on Ω . Denote the set of probability density

$$\mathcal{P}(\Omega) = \Big\{ \rho \in \mathcal{F}(\Omega) \colon \int_{\Omega} \rho dx = 1, \quad \rho \ge 0 \Big\}.$$

The optimization problem in $\mathcal{P}(\Omega)$ takes the form:

$$\min_{\rho \in \mathcal{P}(\Omega)} E(\rho)$$

Here $E(\rho)$ is the objective or loss functional. It evaluates certain divergence or metric functional between ρ and a target density $\rho^* \in \mathcal{P}(\Omega)$. In machine learning problems,

typical examples of $E(\rho)$ include the KL divergence, Maximum mean discrepancy (MMD), cross entropy, etc. Similar to (Euclidean Newton's flow) and (Riemannian Newton's flow), the Newton's flow in probability space (density manifold) takes the form

$$\partial_t \rho_t = -(\text{Hess } E(\rho_t))^{-1} \operatorname{grad} E(\rho_t).$$
 (Information Newton's flow)

Here grad and Hess represent the gradient and the Hessian operator with respect to certain information metric, respectively. To understand (Information Newton's flow), we briefly review the information metrics with the associated gradient operators.

2.2 Information metrics

We first define the tangent space and the cotangent space in probability space. The tangent space at $\rho \in \mathcal{P}(\Omega)$ is defined by

$$T_{\rho}\mathcal{P}(\Omega) = \left\{ \sigma \in \mathcal{F}(\Omega) : \int \sigma dx = 0 \right\}.$$

The cotangent space $T^*_{\rho}\mathcal{P}(\Omega)$ is equivalent to $\mathcal{F}(\Omega)/\mathbb{R}$, which represents the set of functions in $\mathcal{F}(\Omega)$ defined up to addition of constants.

Definition 2 (Metric in probability space) For a given $\rho \in \mathcal{P}(\Omega)$, a metric tensor $\mathcal{G}(\rho) : T_{\rho}\mathcal{P}(\Omega) \to T_{\rho}^*\mathcal{P}(\Omega)$ is an invertible mapping from the tangent space $T_{\rho}\mathcal{P}(\Omega)$ to the cotangent space $T_{\rho}^*\mathcal{P}(\Omega)$. This metric tensor defines the metric (inner product) on the tangent space $T_{\rho}\mathcal{P}(\Omega)$. Namely, for $\sigma_1, \sigma_2 \in T_{\rho}\mathcal{P}(\Omega)$, we define the inner product $g_{\rho}: T_{\rho}\mathcal{P}(\Omega) \times T_{\rho}\mathcal{P}(\Omega) \to \mathbb{R}$ by

$$g_{\rho}(\sigma_1, \sigma_2) = \int \sigma_1 \mathcal{G}(\rho) \sigma_2 dx = \int \Phi_1 \mathcal{G}(\rho)^{-1} \Phi_2 dx,$$

where Φ_i is the solution to $\sigma_i = \mathcal{G}(\rho)^{-1}\Phi_i$, i = 1, 2.

We present two essential examples of metrics in probability space $\mathcal{P}(\Omega)$: Fisher-Rao metric and Wasserstein metric.

Example 1 (Fisher-Rao metric) The inverse of the Fisher-Rao metric tensor follows

$$\mathcal{G}^F(\rho)^{-1}\Phi = \rho\left(\Phi - \int \Phi \rho dx\right), \quad \Phi \in T^*_{\rho}\mathcal{P}(\Omega).$$

The Fisher-Rao metric is defined by

$$g_{\rho}^{F}(\sigma_{1},\sigma_{2}) = \int \Phi_{1}\Phi_{2}\rho dx - \left(\int \Phi_{1}\rho dx\right) \left(\int \Phi_{2}\rho dx\right), \quad \sigma_{1},\sigma_{2} \in T_{\rho}\mathcal{P}(\Omega),$$

where Φ_i satisfies $\sigma_i = \rho \left(\Phi_i - \int \Phi_i \rho dx \right), i = 1, 2.$

Example 2 (Wasserstein metric) The inverse of the Wasserstein metric tensor satisfies

$$\mathcal{G}^W(\rho)^{-1}\Phi = -\nabla \cdot (\rho \nabla \Phi), \quad \Phi \in T^*_{\rho} \mathcal{P}(\Omega).$$

The Wasserstein metric is given by

$$g_{\rho}^{W}(\sigma_{1},\sigma_{2}) = \int \rho \left\langle \nabla \Phi_{1}, \nabla \Phi_{2} \right\rangle dx, \quad \sigma_{1},\sigma_{2} \in T_{\rho}\mathcal{P}(\Omega),$$

where Φ_i is the solution to $\sigma_i = -\nabla \cdot (\rho \nabla \Phi_i), i = 1, 2.$

2.3 Gradient operators

The gradient operator for the objective functional $E(\rho)$ in $(\mathcal{P}(\Omega), \mathcal{G}(\rho))$ satisfies

grad
$$E(\rho) = -\mathcal{G}(\rho)^{-1} \frac{\delta E}{\delta \rho}.$$

Here $\frac{\delta E}{\delta \rho}$ is the L^2 first variation w.r.t. ρ . The gradient flow follows

$$\partial_t \rho_t = -\operatorname{grad} E(\rho_t) = -\mathcal{G}(\rho)^{-1} \frac{\delta E}{\delta \rho_t}.$$

We present gradient operators under either Fisher-Rao metric or Wasserstein metric.

Example 3 (Fisher-Rao gradient operator) The Fisher-Rao gradient operator satisfies

$$\operatorname{grad}^{F} E(\rho) = \rho \left(\frac{\delta E}{\delta \rho} - \int \frac{\delta E}{\delta \rho} \rho dx \right).$$

Example 4 (Wasserstein gradient operator) The Wasserstein gradient operator writes

$$\operatorname{grad}^{W} E(\rho) = -\nabla \cdot \left(\rho \nabla \frac{\delta E}{\delta \rho}\right).$$

3. Information Newton's flow

In this section, we introduce and discuss properties of Hessian operators in probability space. Then, we formulate Newton's flows under information metrics. This is based on the previous definition of gradient operators and the inverse of Hessian operators.

3.1 Information Hessian operators

In this subsection, we review the definition of Hessian operators in probability space and provide the exact formulations of Hessian operators.

For $\sigma \in T_{\rho}\mathcal{P}(\Omega)$, there exists a unique geodesic curve $\hat{\rho}_s$, which satisfies $\hat{\rho}_s|_{s=0} = \rho$ and $\hat{\partial}_s \rho_s|_{s=0} = \sigma$. The Hessian operator of $E(\rho)$ w.r.t. metric tensor $\mathcal{G}(\rho)$ is a mapping Hess $E(\rho) : T_{\rho}\mathcal{P}(\Omega) \to T_{\rho}\mathcal{P}(\Omega)$, which is defined by

$$g_{\rho}(\text{Hess } E(\rho)\sigma, \sigma) = g_{\rho}(\sigma, \text{Hess } E(\rho)\sigma) = \left. \frac{d^2}{ds^2} E(\hat{\rho}_s) \right|_{s=0}$$

Combining with the metric tensor, the Hessian operator uniquely defines a self-adjoint mapping $\mathcal{H}_E(\rho): T^*_{\rho}\mathcal{P}(\Omega) \to T_{\rho}\mathcal{P}(\Omega)$, which satisfies

$$\int \Phi \mathcal{H}_E(\rho) \Phi dx = g_\rho(\sigma, \operatorname{Hess} E(\rho)\sigma), \quad \Phi = \mathcal{G}(\rho)\sigma$$

In Proposition 3, we give an exact formulation of $\int \Phi \mathcal{H}_E(\rho) \Phi dx$ and a relationship between $\mathcal{H}_E(\rho)$ and Hess $E(\rho)$.

Proposition 3 The quantity $g_{\rho}(\sigma, \text{Hess } E(\rho)\sigma)$ is a bi-linear form of Φ :

$$\int \Phi \mathcal{H}_E(\rho) \Phi dx = g_\rho(\sigma, \text{Hess } E(\rho)\sigma)$$

= $-\frac{1}{2} \int \mathcal{A}(\rho)(\Phi, \Phi) \mathcal{G}(\rho)^{-1} \frac{\delta E}{\delta \rho} dx + \int \mathcal{A}(\rho) \left(\Phi, \frac{\delta E}{\delta \rho}\right) \mathcal{G}(\rho)^{-1} \Phi dx \qquad (4)$
+ $\int \int \left(\mathcal{G}(\rho)^{-1}\Phi\right) (y) \frac{\delta^2 E}{\delta \rho^2} (x, y) dy \left(\mathcal{G}(\rho)^{-1}\Phi\right) (x) dx.$

Here $\frac{\delta^2 E}{\delta \rho^2}(x, y)$ is defined by

$$\frac{\delta^2 E}{\delta \rho^2}(x,y) = \frac{\delta}{\delta \rho} \left(\int \frac{\delta E}{\delta \rho}(y) \delta(x-y) dy \right),$$

where $\delta(x)$ is the Dirac delta function. Here $\mathcal{A}(\rho) : T^*_{\rho}\mathcal{P}(\Omega) \times T^*_{\rho}\mathcal{P}(\Omega) \to T^*_{\rho}\mathcal{P}(\Omega)$ is a bi-linear operator which satisfies

$$\mathcal{A}(\rho)(\Phi_1, \Phi_2) = \frac{\delta}{\delta\rho} \int \Phi_1 \mathcal{G}(\rho)^{-1} \Phi_2 dx, \quad \forall \Phi_1, \Phi_2 \in T^*_{\rho} \mathcal{P}(\Omega).$$

Moreover, the operator $\mathcal{H}_E(\rho)$ satisfies

$$\mathcal{H}_E(\rho) = \operatorname{Hess} E(\rho)\mathcal{G}(\rho)^{-1}.$$
(5)

Now, we are ready to present the information Newton's flow in probability space.

Proposition 4 (Information Newton's flow) The Newton's flow of $E(\rho)$ in $(\mathcal{P}(\Omega), \mathcal{G}(\rho))$ satisfies

$$\partial_t \rho_t + (\operatorname{Hess} E(\rho_t))^{-1} \mathcal{G}(\rho_t)^{-1} \frac{\delta E}{\delta \rho_t} = 0.$$

This is equivalent to

$$\begin{cases} \partial_t \rho_t - \mathcal{G}(\rho_t)^{-1} \Phi_t = 0, \\ \mathcal{H}_E(\rho_t) \Phi_t + \mathcal{G}(\rho_t)^{-1} \frac{\delta}{\delta \rho_t} E(\rho_t) = 0. \end{cases}$$
(6)

In particular, we focus on Wasserstein Newton's flow of KL divergence. Other examples of Newton's flows of different objective functions under either Fisher-Rao metric or Wasserstein metric are presented in Appendix B.2 and B.3.

Example 5 (Wasserstein Newton's flow of KL divergence) In this example we prove Theorem 1. As a known fact in (Otto and Villani, 2000) and Gamma calculus (Bakry and Émery, 1985; Li, 2018), the Hessian operator of KL divergence under the Wasserstein metric follows

$$g_{\rho}^{W}(\sigma, \operatorname{Hess}^{W} E(\rho)\sigma) = \int \left(\|\nabla^{2}\Phi\|_{F}^{2} + (\nabla\Phi)^{T}\nabla^{2}f\nabla\Phi) \right) \rho dx,$$

where $\sigma = -\nabla \cdot (\rho \nabla \Phi)$ and $\|\cdot\|_F$ is the Frobenius norm of a matrix in $\mathbb{R}^{n \times n}$. Via integration by parts, we validate that the operator $\mathcal{H}^W_E(\rho)$ follows

$$\mathcal{H}_{E}^{W}(\rho)\Phi = \nabla^{2}: (\rho\nabla^{2}\Phi) - \nabla \cdot (\rho\nabla^{2}f\nabla\Phi).$$
(7)

We also present the Wasserstein Newton's flow of KL divergence in Gaussian families. Proposition 5 ensures the existence of information Newton's flows in Gaussian families.

Proposition 5 Suppose that ρ_0, ρ^* are Gaussian distributions with zero means and their covariance matrices are Σ_0 and Σ^* . $E(\Sigma)$ evaluates the KL divergence from ρ to ρ^* :

$$E(\Sigma) = \frac{1}{2} \left(\operatorname{tr}(\Sigma(\Sigma^*)^{-1}) - d - \log \det \left(\Sigma(\Sigma^*)^{-1} \right) \right).$$
(8)

Let (Σ_t, S_t) satisfy

$$\begin{cases} \dot{\Sigma}_t - 2(S\Sigma_t + \Sigma S_t) = 0, \\ 2\Sigma_t S_t (\Sigma^*)^{-1} + 2(\Sigma^*)^{-1} S_t \Sigma_t + 4S_t = -(\Sigma_t (\Sigma^*)^{-1} + (\Sigma^*)^{-1} \Sigma_t - 2I). \end{cases}$$
(9)

with initial values $\Sigma_t|_{t=0} = \Sigma_0$ and $S_t|_{t=0} = 0$. Thus, for any $t \ge 0$, Σ_t is well-defined and stays positive definite. We denote

$$\rho_t(x) = \frac{(2\pi)^{-n/2}}{\sqrt{\det(\Sigma_t)}} \exp\left(-\frac{1}{2}x^T \Sigma_t^{-1} x\right), \quad \Phi_t(x) = x^T S_t x + C(t),$$

where $C(t) = -t + \frac{1}{2} \int_0^t \log \det(\Sigma_s(\Sigma^*)^{-1}) ds$. Then, ρ_t and Φ_t follow the information Newton's flow (3) with initial values $\rho_t|_{t=0} = \rho_0$ and $\Phi_t|_{t=0} = 0$.

4. Newton's Langevin dynamics

In this section, we primarily focus on the Wasserstein Newton's flow of KL divergence. We formulate it into the Newton's Langevin dynamics for Bayesian sampling problems. The connection and difference with

Let the objective functional $E(\rho) = D_{KL}(\rho || \rho^*)$ evaluate the KL divergence from ρ to a target density $\rho^*(x) \propto \exp(-f(x))$ with $\int \exp(-f(x)) dx < \infty$. This specific optimization problem is important since it corresponds to sampling from the target density ρ^* . Classical Langevin MCMC algorithms evolves samples following overdamped Langevin dynamics (OLD), which satisfies

$$dX_t = -\nabla f(X_t)dt + \sqrt{2}dB_t,$$

where B_t is the standard Brownian motion. Denote ρ_t as the density function of the distribution of X_t . The evolution of ρ_t satisfies the Fokker-Planck equation

$$\partial_t \rho_t = \nabla \cdot (\rho_t \nabla f) + \Delta \rho_t.$$

A known fact is that the Fokker-Planck equation is the Wasserstein gradient flow (WGF) of KL divergence, i.e.

$$\partial_t \rho_t = -\operatorname{grad}^W \operatorname{D}_{\operatorname{KL}}(\rho_t \| \rho^*)$$

= $\mathcal{G}^W(\rho_t)^{-1} \frac{\delta}{\delta \rho_t} \operatorname{D}_{\operatorname{KL}}(\rho_t \| \rho^*)$
= $\nabla \cdot (\rho_t \nabla (f + \log \rho_t + 1))$
= $\nabla \cdot (\rho_t \nabla f) + \Delta \rho_t.$ (10)

where we use the fact that $\frac{\delta}{\delta\rho} D_{\text{KL}}(\rho_t \| \rho^*) = \log \rho + t + f + 1$ and $\rho \nabla \log \rho = \nabla \rho$.

It is worth mentioning that OLD can be viewed as particle implementations of WGF

(10). From the viewpoint of fluid dynamics, WGF also has a Lagrangian formulation

$$dX_t = -\nabla f(X_t)dt - \nabla \log \rho_t(X_t)dt$$

We name above dynamics by the Lagrangian Langevin Dynamics (LLD). Here 'Lagrangian' refers to the Lagrangian coordinates (flow map) in fluid dynamics (Villani, 2008).

Overall, many sampling algorithms follow OLD or LLD. The evolution of corresponding density follows the Wasserstein gradient flow (10). E.g. the classical Langevin MCMC (unadjusted Langevin algorithm) is the time discretization of OLD. The Particle-based Variational Inference methods (ParVI), (Liu et al., 2019) can be viewed as the discrete-time approximation of LLD.

In short, we notice that the Langevein dynamics can be viewed as first-order methods for Bayesian sampling problems. Analogously, the Wasserstein Newton's flow of KL divergence derived in Example 5 corresponds to certain Langevin dynamics of particle systems, named *Newton's Langevin dynamics*.

Theorem 6 Consider the Newton's Langevin dynamics

$$dX_t = \nabla \Phi_t^{\text{Newton}}(X_t) dt, \tag{11}$$

where $\Phi_t^{\text{Newton}}(x)$ is the solution to Wasserstein Newton's direction equation (3):

$$\nabla^2 : (\rho_t \nabla^2 \Phi_t) - \nabla \cdot (\rho_t \nabla^2 f \nabla \Phi_t) - \nabla \cdot (\rho_t \nabla f) - \Delta \rho_t = 0.$$

Here X_0 follows an initial distribution ρ^0 and ρ_t is the distribution of X_t . Then, ρ_t is the solution to Wasserstein Newton's flow with an initial value $\rho_0 = \rho^0$.

Proof Note that ρ_t is the distribution of X_t . The dynamics of X_t implies

$$\partial_t \rho_t + \nabla \cdot (\rho_t \nabla \Phi_t^{\text{Newton}}) = 0.$$

Because Φ_t satisfies the Wasserstein Newton's direction equation (3), ρ_t is the solution to Wasserstein Newton's flow.

Remark 7 We notice that the Newton's Langevien dynamics is different from HAMCMC (Simsekli et al., 2016). Detailed comparisons can be found in Appendix C.1.

The following proposition provide a closed-form formula for NLD in 1D Gaussian family.

Proposition 8 Assume that $f(x) = (2\Sigma^*)^{-1}(x - \mu^*)^2$, where $\Sigma^* > 0$ and μ^* are given. Suppose that the particle system X_0 follows the Gaussian distribution. Then X_t follows a Gaussian distribution with mean μ_t and variance Σ_t . The corresponding NLD satisfies

$$dX_t = \left(\frac{\Sigma^* - \Sigma}{\Sigma^* + \Sigma_t} X_t - \frac{2\Sigma^*}{\Sigma^* + \Sigma_t} \mu_t + \mu^*\right) dt.$$

And the evolution of μ_t and Σ_t satisfies

$$d\mu_t = (-\mu_t + \mu^*)dt, \quad d\Sigma_t = 2\frac{\Sigma^* - \Sigma_t}{\Sigma^* + \Sigma_t}\Sigma_t dt.$$

The explicit solutions of μ_t and Σ_t satisfy

$$\mu_t = e^{-t}(\mu_0 - \mu^*) + \mu^*, \quad \Sigma_t = \Sigma^* + (\Sigma_0 - \Sigma^*)e^{-t}\sqrt{\frac{e^{-2t}(\Sigma_0 - \Sigma^*)^2}{4\Sigma_0^2}} + \frac{1}{\Sigma_0\Sigma^*}.$$

We present discrete-time particle implementations of Newton's Langevin dynamics in section 5 and numerical examples in section 7.

5. Particle implementation of Wasserstein Newton's method

In this section, we design sampling efficient implementations of Wasserstein Newton's method. Focusing on Wasserstein Newton's flow of KL divergence, we introduce a variational formulation for computing the Wasserstein Newton's direction. By restricting the domain of the variational problem in a linear subspace or reproducing kernel Hilbert space (RKHS), we derive sampling efficient algorithms. Besides, a hybrid method between Newton's Langevin dynamics and overdamped Langevin dynamics is provided.

We briefly review update rules of Newton's methods and hybrid methods in Euclidean space. In each iteration, the update rule of Newton's method follows

$$x_{k+1} = x_k + \alpha_k p_k, \quad p_k = -\nabla^2 f(x_k)^{-1} \nabla f(x),$$

Suppose that f(x) is strictly convex. Namely, $\nabla^2 f(x)$ is positive definite for all $x \in \mathbb{R}^d$. To compute the Newton's direction p_k , it is equivalent to solve the following variational problem

$$\min_{p \in \mathbb{R}^n} p^T \nabla^2 f(x_k) p + 2 \nabla f(x_k)^T p$$

In practice, the Newton's direction may not lead to the decrease in the objective function, especially when f(x) is non-convex. Nevertheless, the Newton's method often converges when the update is close to the minimizer. One way to overcome this problem is the hybrid method. Consider a hybrid update of the Newton's direction and the gradient's direction

$$x_{k+1} = x_k + \alpha_k p_k - \alpha_k \gamma \nabla f(x_k),$$

where $\gamma > 0$ is a parameter.

Following above ideas in Euclidean space, we present a particle implementation of information Newton's method. Here we connect density $\rho_k \in \mathcal{P}(\Omega)$ with a particle system $\{x_k^n\}_{i=1}^N$. Namely, we assume that the distribution $\{x_k^n\}_{n=1}^N$ follows $\rho_k(x)$. We update each particle by

$$x_{k+1}^n = x_k^n + \alpha_k \nabla \Phi_k(x_k^n), \quad i = 1, 2 \dots N.$$

Here $\hat{\Phi}_k$ is an approximated solution to the Wasserstein Newton's direction equation (3). The details on obtaining $\hat{\Phi}_k$ is left in subsection 5.1.

WANG AND LI

In practice, the Wasserstein Newton's direction may not be a descent direction if the update is far away from the target distribution. To overcome this issue, we propose a hybrid update of the Wasserstein Newton's direction and the Wasserstein gradient direction.

Let $\gamma \ge 0$ be a parameter. Here we recall that there are two choices for using the gradient direction. Namely, if we use overdamped Langevin dynamics as the gradient direction, the hybrid update rule follows

$$x_{k+1}^n = x_k^n + \alpha_k \nabla \hat{\Phi}_k(x_k^n) - \gamma \alpha_k \nabla f(x_k^n) + \sqrt{2\gamma \alpha_k} z_k,$$
(12)

where $z_k \sim \mathcal{N}(0, I)$. If we use Lagrangian Langevin dynamics as the gradient direction, the hybrid update rule satisfies

$$x_{k+1}^n = x_k^n + \alpha_k \nabla \hat{\Phi}_k(x_k^n) - \gamma \alpha_k (\nabla f(x_k^n) + \xi_k(x_k^n)).$$
(13)

Here ξ_k is an approximation of $\nabla \log \rho_k$. For general ρ_k and ρ^* , we can approximate $\nabla \log \rho_k$ via kernel density estimation (KDE) (Gretton et al., 2012). Namely, we approximate $\nabla \log \rho_k$ by

$$\xi_k(x) = \frac{\sum_{n=1}^N \nabla_y k(x, x_k^n)}{\sum_{n=1}^N \nabla k(x, x_k^n)}.$$

Here k(x, y) is a given positive kernel. A typical choice of k(x, y) is a Gaussian kernel with a bandwidth h > 0, such that

$$k(x,y) = (2\pi h)^{-n/2} \exp\left(-\frac{\|x-y\|^2}{2h}\right).$$

The overall algorithm is summarized in Algorithm 1.

Algorithm 1 Wasserstein Newton's method with hybrid update

Require: initial positions $\{x_0^n\}_{n=1}^N$, $\epsilon \ge 0$, step sizes α_k , parameters $\lambda_k \ge 0$, maximum iteration K.

1: Set k = 0.

- 2: while k < K and the convergence criterion is not met do
- 3: Compute an approximate solution Φ_k to (3).
- 4: Update particle positions by (12) or (13).
- 5: Set k = k + 1.
- 6: end while

Remark 9 It worths mentioning that our algorithm corresponds to the following hybrid Langvien dynamics

$$dX_t = (\nabla \Phi_t - \gamma \nabla f)dt + \sqrt{2\gamma}dB_t,$$

where B_t is the standard Brownian motion, $\gamma \geq 0$ is a parameter and Φ_t satisfies (3).

5.1 Variational formulation for Wasserstein Newton's direction

Similar to the Euclidean case, we derive a variational formulation for estimating Wasserstein Newton's direction, and provide the associated particle formulations.

Proposition 10 Suppose that $\mathcal{H}: T^*_{\rho}\mathcal{P}(\Omega) \to T_{\rho}\mathcal{P}$ is a linear self-adjoint operator and \mathcal{H} is positive definite. Let $u \in T_{\rho}\mathcal{P}$. Then the minimizer of variational problem

$$\min_{\Phi \in T^*_{\rho} \mathcal{P}(\Omega)} J(\Phi) = \int \left(\Phi \mathcal{H} \Phi - 2u\Phi\right) dx,$$

satisfies $\mathcal{H}\Phi = u$, where $\Phi \in T^*_{\rho}\mathcal{P}(\Omega)$.

Proof Since \mathcal{H} is linear and self-adjoint, the optimal solution of satisfies

$$0 = \frac{\delta J}{\delta \Phi} = 2\mathcal{H}\Phi - 2u$$

Hence, Φ satisfies $\mathcal{H}\Phi = u$. On the other hand, let Φ satisfy $\mathcal{H}\Phi = u$. Then, for any $\Psi \in T^*_{\rho}\mathcal{P}(\Omega)$, it follows

$$J(\Phi + \Psi) = \int \left((\Phi + \Psi) \mathcal{H}(\Phi + \Psi) - 2u(\Phi + \Psi) \right) dx$$
$$= \int \left(\Phi \mathcal{H}\Phi - 2u\Phi \right) dx + \int \left(\Psi \mathcal{H}\Psi - 2u\Psi - 2\Psi \mathcal{H}\Phi \right) dx$$
$$= J(\Phi) + \int \Psi \mathcal{H}\Psi dx \ge J(\Phi).$$

The last inequality is based on the fact that \mathcal{H} is positive definite. Hence, Φ is the optimal solution to the proposed variational problem. This completes the proof.

Suppose that f is strongly convex, or equivalent, $\nabla^2 f(x)$ is positive definite for $x \in \Omega$. Then, the operator $\mathcal{H}_E(\rho)$ defined in (7) is positive definite. In this case, proposition 10 indicates that solving Wasserstein Newton's direction equation (3) is equivalent to optimizing the following variational problem.

$$\min_{\Phi \in T^*_{\rho_k} \mathcal{P}(\Omega)} J(\Phi) = \int \left(\|\nabla^2 \Phi\|_F^2 + \|\nabla \Phi\|_{\nabla^2 f}^2 + 2\left\langle \nabla f + \nabla \log \rho_k, \nabla \Phi \right\rangle \right) \rho_k dx.$$

Here we denote $||v||_A^2 = v^T A v$. For possibly non-convex f, we consider a regularized problem

$$\min_{\Phi \in T^*_{\rho_k} \mathcal{P}(\Omega)} J^{\epsilon}(\Phi) = \int \left(\|\nabla^2 \Phi\|_F^2 + \|\nabla \Phi\|_{\nabla^2 f + \epsilon I}^2 + 2\left\langle \nabla f + \nabla \log \rho_k, \nabla \Phi \right\rangle \right) \rho_k dx.$$
(14)

Here $\epsilon \geq 0$ is a regularization parameter to ensure that $\nabla^2 f(x) + \epsilon I$ is positive definite for $x \in \Omega$.

Remark 11 Namely, we penalize the objective function by adding the squared norm of Φ induced by the Wasserstein metric. In other words,

$$\min_{\Phi \in T^*_{\rho_k} \mathcal{P}(\Omega)} J(\Phi) + \epsilon \int \|\nabla \Phi\|^2 \rho_k dx$$

In terms of samples, we can rewrite (14) into

$$\min_{\Phi \in T^*_{\rho_k} \mathcal{P}(\Omega)} J^{\epsilon}(\Phi) = \frac{1}{N} \sum_{n=1}^{N} \left(\|\nabla^2 \Phi(x_k^n)\|_F^2 + \|\nabla \Phi(x_k^n)\|_{\nabla^2 f(x_k^n) + \epsilon I}^2 + 2 \left\langle \nabla f(x_k^n) + \nabla \log \rho_k(x_k^n), \nabla \Phi(x_k^n) \right\rangle \right).$$
(15)

In high dimensional sample space, directly solving (15) for $\Phi \in T^*_{\rho_k} \mathcal{P}(\Omega)$ can be difficult. To deal with this issue, we restrict the functional space of Φ into a linear subspace $\mathcal{S} \subseteq T^*_{\rho_k} \mathcal{P}(\Omega)$. An appropriately chosen \mathcal{S} can lead to a closed-form solution to (14). For the rest of this section, we discuss two choices of \mathcal{S} , including finite dimensional affine subspace and reproducing kernel Hilbert space (RKHS).

5.2 Affine models

Consider $S = \text{span}\{\psi_i\}_{i=1}^m$, where $\psi_i : \Omega \to R$ are given basis functions. Namely, we assume that $\Phi(x)$ is a linear combination of ψ_1, \ldots, ψ_m , such that

$$\Phi(x) = \langle \mathbf{a}, \psi(x) \rangle = \sum_{i=1}^{m} a_i \psi_i(x),$$

where $\mathbf{a} \in \mathbb{R}^m$ and $\psi(x) = [\psi_1(x), \psi_2(x), \dots, \psi_m(x)].$

Proposition 12 Suppose that $\Phi(x) = \langle \mathbf{a}, \psi(x) \rangle$. Then, the optimization problem (15) with the constraint $\Phi \in S$ is equivalent to

$$\min_{\mathbf{a}\in\mathbb{R}^m} J^{\epsilon}(\mathbf{a}) = \mathbf{a}^T (\mathbf{B}_k + \mathbf{D}_k) \mathbf{a} + 2\mathbf{c}_k^T \mathbf{a},$$

where $\mathbf{B}_k, \mathbf{D}_k \in \mathbb{R}^{m \times m}$ and $\mathbf{c}_k \in \mathbb{R}^m$. The detailed formulations of $\mathbf{B}_k, \mathbf{D}_k$ and \mathbf{c}_k are provided as follows.

$$\mathbf{B}_{k} = \frac{1}{N} \sum_{n=1}^{N} \nabla \psi(x_{k}^{n}) (\nabla^{2} f(x_{k}^{n}) + \epsilon I) (\nabla \psi(x_{k}^{n}))^{T},$$
$$\mathbf{D}(x)_{j_{1},j_{2}} = \frac{1}{N} \sum_{n=1}^{N} \operatorname{tr}(\nabla^{2} \psi_{j_{1}}(x_{k}^{n}) \nabla^{2} \psi_{j_{2}}(x_{k}^{n})),$$
$$\mathbf{c}(x) = \frac{1}{N} \sum_{n=1}^{N} \nabla \psi(x_{k}^{n}) (\nabla f(x_{k}^{n}) + \xi_{k}(x_{k}^{n})).$$

If $\mathbf{B}_k + \mathbf{D}_k$ is positive definite, the optimal solution follows $\mathbf{a} = -(\mathbf{B}_k + \mathbf{D}_k)^{-1}c_k$. The optimal solution $\hat{\Phi}$ follows $\hat{\Phi}(x) = \langle \mathbf{a}, \psi(x) \rangle$.

Proof We denote the Jacobian $\nabla \psi(x) \in \mathbb{R}^{n \times m}$. As a result, $J(\mathbf{a})$ turns to be

$$J^{\epsilon}(\mathbf{a}) = \left\{ \frac{1}{N} \sum_{n=1}^{N} \left\| \sum_{j=1}^{m} a_{j} \nabla^{2} \psi_{j}(x_{k}^{n}) \right\|_{F}^{2} + \mathbf{a}^{T} \mathbf{B}(x_{k}^{n}) \mathbf{a} + 2\mathbf{a}^{T} \mathbf{c}(x_{k}^{n}) \right\}$$

We can further compute that

$$\left\|\sum_{j=1}^{m} \mathbf{a}_{j} \nabla^{2} \psi_{j}(x_{k}^{n})\right\|_{F}^{2} = \sum_{j_{1}=1}^{m} \sum_{j_{2}=1}^{m} \mathbf{a}_{j_{1}} \nabla^{2} \psi_{j_{1}}(x_{k}^{n}) \nabla^{2} \psi_{j_{2}}(x_{k}^{n}) \mathbf{a}_{j_{2}} = \mathbf{a}^{T} \mathbf{D}(x_{k}^{n}) \mathbf{a}.$$

This completes the proof.

This affine approximation technique has been used in approximating natural gradient direction in (Li et al., 2019). Hence, we call our method affine information Newton's method. In particular, we set m = 2d and consider the basis

In particular, we set m = 2d and consider the basis

$$\psi_i(x) = x_i, \quad \psi_{i+d}(x) = x_i^2, \quad 1 \le i \le d.$$

In other words, we assume that $\Phi(x)$ takes the form $\Phi(x) = \frac{1}{2}x \operatorname{diag}(s)x + b^T x$, where $s, b \in \mathbb{R}^d$. For simplicity, we denote $v_k^n = \nabla f(x_k^n) + \xi_k(x_k^n)$.

$$J^{\epsilon}(s,b) = \begin{bmatrix} s \\ b \end{bmatrix}^{T} \mathbf{H}_{k} \begin{bmatrix} s \\ b \end{bmatrix} + 2 \begin{bmatrix} s \\ b \end{bmatrix}^{T} u_{k}.$$

where we denote $\mathbf{H}_k \in \mathbb{R}^{2d \times 2d}$ via

$$\mathbf{H}_{k} = \begin{bmatrix} I + \frac{1}{N} \sum_{n=1}^{N} \operatorname{diag}(x_{k}^{n}) (\nabla^{2} f(x_{k}^{n}) + \epsilon I) \operatorname{diag}(x_{k}^{n}) & \frac{1}{N} \sum_{n=1}^{N} \operatorname{diag}(x_{k}^{n}) (\nabla^{2} f(x_{k}^{n}) + \epsilon I) \\ \frac{1}{N} \sum_{n=1}^{N} (\nabla^{2} f(x_{k}^{n}) + \epsilon I) \operatorname{diag}(x_{k}^{n}) & \frac{1}{N} \sum_{n=1}^{N} (\nabla^{2} f(x_{k}^{n}) + \epsilon I) \end{bmatrix},$$

and $u_k \in \mathbb{R}^{2d}$ via

$$u_k = \begin{bmatrix} \frac{1}{N} \sum_{n=1}^N \operatorname{diag}(x_k^n) v_k^n \\ \frac{1}{N} \sum_{n=1}^N v_k^n \end{bmatrix}.$$

Hence, the optimal solution for minimizing J(s, b) follows

$$\begin{bmatrix} s_k \\ b_k \end{bmatrix} = -(\mathbf{H}_k)^{-1} u_k$$

Hence, the approximate solution $\hat{\Phi}_k$ computed via the affine method follows

$$\nabla \hat{\Phi}_k(x) = \operatorname{diag}(s_k)x + b_k. \tag{16}$$

The overall algorithm are summarized in Algorithm 2. For simplicity, we do not mention the hybrid update.

When the optimal solution Φ_k to (15) is highly non-linear, S in affine methods may not be large enough to approximate Φ_k well.

Algorithm 2 Wasserstein Newton's method with affine models.

Require: initial positions $\{x_0^i\}_{i=1}^N$, $\epsilon \ge 0$, step sizes α_k , maximum iteration K. 1: Set k = 0.

- 2: while k < K and the convergence criterion is not met do
- Compute $v_k^n = \nabla f(x_k^n) + \xi_k(x_k^n)$. Here ξ_k is an approximation of $\nabla \log \rho_k$. Calculate \mathbf{H}_k by 3:
- 4:

$$\mathbf{H}_{k} = \begin{bmatrix} I + \frac{1}{N} \sum_{n=1}^{N} \operatorname{diag}(x_{k}^{n})(\nabla^{2}f(x_{k}^{n}) + \epsilon I) \operatorname{diag}(x_{k}^{n}) & \frac{1}{N} \sum_{n=1}^{N} \operatorname{diag}(x_{k}^{n})(\nabla^{2}f(x_{k}^{n}) + \epsilon I) \\ \frac{1}{N} \sum_{n=1}^{N} (\nabla^{2}f(x_{k}^{n}) + \epsilon I) \operatorname{diag}(x_{k}^{n}) & \frac{1}{N} \sum_{n=1}^{N} (\nabla^{2}f(x_{k}^{n}) + \epsilon I) \end{bmatrix},$$

and formulate u_k by

$$u_k = \begin{bmatrix} \frac{1}{N} \sum_{n=1}^N \operatorname{diag}(x_k^n) v_k^n \\ \frac{1}{N} \sum_{n=1}^N v_k^n \end{bmatrix}.$$

5:Compute s_k and b_k by

$$\begin{bmatrix} s_k \\ b_k \end{bmatrix} = -(\mathbf{H}_k)^{-1} u_k.$$

Update particle positions by 6:

$$x_{k+1}^n = x_k^n + \alpha_k (\operatorname{diag}(s_k) x_k^n + b_k).$$

Set k = k + 1. 7:

8: end while

5.3 Kernel models

In this subsection, we approximate the Wasserstein Newton's direction in kernel models. Specifically, we consider S as the RKHS with an associated kernel function $k(x, y) : \mathbb{R}^d \times$ $\mathbb{R}^d \to \mathbb{R}$. Compared to finite-dimensional linear subspace, RKHS can be viewed as with infinitely many feature functions. Detailed description about RKHS and the related norm can be found in (Berlinet and Thomas-Agnan, 2011).

To ensure the well-posedness of the optimal solution, we penalize the objective function using the RKHS norm $\|\cdot\|_{\mathcal{S}}$. Hence, we consider a regularized variational problem based on (14)

$$\min_{\Phi \in \mathcal{H}} \int \left(\|\nabla^2 \Phi\|_F^2 + \|\nabla \Phi\|_{\nabla^2 f + \epsilon I}^2 + 2 \langle \nabla f, \nabla \Phi + \nabla \log \rho_k \rangle \right) \rho_k dx + \lambda \|\Phi\|_{\mathcal{S}}^2
= \int \left(\|\nabla^2 \Phi\|_F^2 + \|\nabla \Phi\|_{\nabla^2 f + \epsilon I}^2 + 2 \langle \nabla f, \nabla \Phi \rangle - 2\Delta \Phi \right) \rho_k dx + \lambda \|\Phi\|_{\mathcal{S}}^2.$$
(17)

In terms of samples, this variational problem becomes

$$\min_{\Phi \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^{N} \left(\|\nabla^2 \Phi(x_k^n)\|_F^2 + \|\nabla \Phi(x_k^n)\|_{\nabla^2 f(x_k^n) + \epsilon I}^2 + 2 \left\langle \nabla f(x_k^n), \nabla \Phi(x_k^n) \right\rangle - 2\Delta \Phi(x_k^n) \right) + \lambda \|\Phi\|_{\mathcal{S}}^2.$$
(18)

From the general representation theorem (Schölkopf et al., 2001), the minimizer of (18) can take the form

$$\Phi(x) = \sum_{n=1}^{N} \left(\sum_{i=1}^{d} \alpha_{i,n} \partial_i k(x_k^n, x) + \sum_{j_1=1}^{d} \sum_{j_2=1}^{d} \beta_{j_1, j_2, n} \partial_{j_1, j_2} k(x_k^n, x) \right).$$
(19)

Proposition 13 Let Φ take the form (19). Then, (18) is equivalent to

$$\inf_{\alpha \in \mathbb{R}^{Nd}, \beta \in \mathbb{R}^{Nd^2}} \begin{bmatrix} \alpha \\ \beta \end{bmatrix}^T \begin{bmatrix} K^{1,2} \\ K^{2,2} \end{bmatrix} \begin{bmatrix} K^{1,2} \\ K^{2,2} \end{bmatrix}^T \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + \begin{bmatrix} \alpha \\ \beta \end{bmatrix}^T \begin{bmatrix} K^{1,1} \\ K^{2,1} \end{bmatrix} H \begin{bmatrix} K^{1,1} \\ K^{2,1} \end{bmatrix}^T \begin{bmatrix} \alpha \\ \beta \end{bmatrix} \\
+ N\lambda \begin{bmatrix} \alpha \\ \beta \end{bmatrix}^T \begin{bmatrix} K^{1,1} \\ K^{2,1} \end{bmatrix} \begin{bmatrix} K^{1,2} \\ K^{2,2} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} - 2 \begin{bmatrix} \alpha \\ \beta \end{bmatrix}^T \begin{bmatrix} K^{1,1} \\ K^{2,1} \end{bmatrix} \begin{bmatrix} v \\ e \end{bmatrix}.$$
(20)

 $Here \ we \ denote$

$$\begin{aligned} v &= \begin{bmatrix} -\nabla f(x_k^1) \\ \vdots \\ -\nabla f(x_k^N) \end{bmatrix} \in \mathbb{R}^{Nd}, \quad e = \begin{bmatrix} vec(I_d) \\ \vdots \\ vec(I_d) \end{bmatrix} \in \mathbb{R}^{Nd^2}, \\ H &= \begin{bmatrix} \nabla^2 f(x_k^1) + \epsilon I & 0 & \dots & 0 \\ 0 & \nabla^2 f(x_k^2) + \epsilon I & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \nabla^2 f(x_k^N) + \epsilon I \end{bmatrix} \in \mathbb{R}^{Nd \times Nd}, \end{aligned}$$

and

$$K^{p,q} = \begin{bmatrix} K_{1,1}^{p,q} & \dots & K_{1,N}^{p,q} \\ \vdots & \ddots & \vdots \\ K_{N,1}^{p,q} & \dots & K_{N,N}^{p,q} \end{bmatrix}, \quad p,q \in \{1,2\}.$$

Each $K^{p,q}_{n,n'}$ are defined by

$$\begin{pmatrix} K_{n,n'}^{1,1} \end{pmatrix}_{i,j} = \partial_{i,j+d} k(x_k^n, x_k^{n'}), \quad K_{n,n'}^{1,1} \in \mathbb{R}^{d \times d},$$

$$\begin{pmatrix} K_{n,n'}^{1,2} \end{pmatrix}_{i,(j_1-1)d+j_2} = \partial_{i,j_1+d,j_2+d} k(x_k^n, x_k^{n'}), \quad K_{n,n'}^{1,2} \in \mathbb{R}^{d \times d^2},$$

$$\begin{pmatrix} K_{n,n'}^{2,1} \end{pmatrix}_{(j_1-1)d+j_2,i} = \partial_{j_1,j_2,i+d} k(x_k^n, x_k^{n'}), \quad K_{n,n'}^{2,1} \in \mathbb{R}^{d^2 \times d},$$

$$\begin{pmatrix} K_{n,n'}^{2,2} \\ \dots \end{pmatrix}_{(i_1-1)d+i_2,(j_1-1)d+j_2} = \partial_{i_1,i_2,j_1+d,j_2+d} k(x_k^n, x_k^{n'}), \quad K_{n,n'}^{2,2} \in \mathbb{R}^{d^2 \times d^2}.$$

Here we use the notation $\partial_i k(x,y) = \partial_{x_i} k(x,y)$ and $\partial_{j+d} = \partial_{y_j} k(x,y)$. The optimal solution follows

$$\begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \left(\begin{bmatrix} K^{1,2} \\ K^{2,2} \end{bmatrix} \begin{bmatrix} K^{1,2} \\ K^{2,2} \end{bmatrix}^T + \begin{bmatrix} K^{1,1} \\ K^{2,1} \end{bmatrix} H \begin{bmatrix} K^{1,1} \\ K^{2,1} \end{bmatrix}^T + N\lambda \begin{bmatrix} \alpha \\ \beta \end{bmatrix}^T \begin{bmatrix} K^{1,1} & K^{1,2} \\ K^{2,1} & K^{2,2} \end{bmatrix} \right)^{\dagger} \begin{bmatrix} K^{1,1} & K^{1,2} \\ K^{2,1} & K^{2,2} \end{bmatrix} \begin{bmatrix} v \\ e \end{bmatrix}$$

Here \dagger denotes the Moore pseudo-inverse. Hence the approximated solution $\hat{\Phi}_k$ satisfies

$$\begin{bmatrix} \nabla \hat{\Phi}_k(x_k^1) \\ \vdots \\ \nabla \hat{\Phi}_k(x_k^N) \end{bmatrix} = K^{1,1} \alpha + K^{1,2} \beta.$$

To solve (20) is equivalent to solve a $N(d + d^2) \times N(d + d^2)$ linear system. Moreover, this linear system is potentially to be ill-posed, especially for large N and d. Hence, we further restrict $\beta = 0$ in (20) (this is equivalent to choose a smaller basis in representing $\Phi(x)$). Then, (20) reduces to

$$\inf_{\alpha \in \mathbb{R}^{Nd}} \quad \alpha^T K^{1,2} K^{2,1} \alpha + \alpha^T K^{1,1} H K^{1,1} \alpha + N \lambda \alpha^T K^{1,1} \alpha - 2\alpha^T \begin{bmatrix} K^{1,1} & K^{1,2} \end{bmatrix} \begin{bmatrix} v \\ e \end{bmatrix}.$$
(21)

The optimal solution follows

$$\alpha = (K^{1,2}K^{2,1} + K^{1,1}HK^{1,1} + N\lambda K^{1,1})^{-1} \begin{bmatrix} K^{1,1} & K^{1,2} \end{bmatrix} \begin{bmatrix} v \\ e \end{bmatrix}$$

Denote $\mathbf{C} = K^{1,2}K^{2,1} + K^{1,1}HK^{1,1} + N\lambda K^{1,1}$. Hence, the approximate solution $\hat{\Phi}_k(x_k^n)$ satisfies

$$\begin{bmatrix} \nabla \Phi_k(x_k^1) \\ \vdots \\ \nabla \hat{\Phi}_k(x_k^N) \end{bmatrix} = K^{1,1} \alpha = K^{1,1} \mathbf{C}^{-1} (K^{1,1} v + K^{1,2} e).$$
(22)

In practice, when N, d are large, the computation cost of $K^{1,2}K^{2,1}$ is quite heavy, which is of order $O(N^3d^4)$. Hence, we consider a block-diagonal approximation \mathbf{C}_{bd} of \mathbf{C} , which is defined by

$$\mathbf{C}_{\rm bd} = \begin{bmatrix} C_{1,1} & 0 & \dots & 0 \\ 0 & C_{2,2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & C_{N,N} \end{bmatrix}$$

Here each block $C_{i,i} \in \mathbb{R}^{d \times d}$ can be computed by

$$C_{i,i} = N\lambda K_{i,i}^{1,1} + \sum_{j=1}^{N} \left(K_{i,j}^{1,2} K_{j,i}^{2,1} + K_{i,j}^{1,1} \nabla^2 f(x_k^j) K_{j,i}^{1,1} \right).$$

The computational cost of \mathbf{C}_{bd} is $O(N^2 d^4)$. We also note that for Gaussian kernel, with $\lambda > 0$, $C_{i,i}$ is invertible. Hence, we can compute the approximate solution $\hat{\Phi}_k(x_k^n)$ by

$$\begin{bmatrix} \nabla \hat{\Phi}_k(x_k^1) \\ \vdots \\ \nabla \hat{\Phi}_k(x_k^N) \end{bmatrix} = K^{1,1} \mathbf{C}_{\mathrm{bd}}^{-1}(K^{1,1}v + K^{1,2}e).$$
(23)

The overall algorithm is summarized in Algorithm 3.

Algorithm 3 Wasserstein Newton's method with RKHS.

Require: initial positions $\{x_0^n\}_{n=1}^N$, $\epsilon \ge 0$, step sizes α_k , maximum iteration K. 1: Set k = 0.

- 2: while k < K and the convergence criterion is not met do
- 3: Calculate $H, v, e, K^{1,1}, K^{1,2}$ and $K^{2,1}$ in Proposition 13 based on $\{x_k^n\}_{n=1}^N$.
- 4: Formulate $\hat{\Phi}_k(x_k^n)$ via (22) or (23).
- 5: Update particle positions by

$$x_{k+1}^n = x_k^n + \alpha_k \nabla \hat{\Phi}_k(x_k^n).$$

6: Set k = k + 1. 7: **end while**

Besides, we can use a sparse kernel approximation (Arbel et al., 2019; Maoutsa et al., 2020) to further reduce the computational cost. Namely, we assume that $\Phi(x)$ takes the form

$$\Phi(x) = \sum_{m=1}^{M} \sum_{i=1}^{d} \alpha_{i,m} \partial_i k(z^m, x).$$
(24)

Here $M \ll N$ and $\{z^m\}_{m=1}^M$ are randomly sampled from $\{x_k^n\}_{n=1}^N$. This can reduce the computational cost to $O(MN^2d^4)$ (or $O(MNd^4)$ if we apply the block-diagonal approximation).

Remark 14 In future works, we expect to find efficient methods to approximate the solution to (20) with low computational cost in terms of N and d.

Remark 15 We notice that our Wasserstein Newton's method with RKHS is related to Stein variational Newton's method (SVN) (Detommaso et al., 2018). Here SVN restricts the Newton's direction of general transformation map in RKHS, while our method restricts the potential function of gradient transportation map in RKHS. See details in the appendix. We also provide detailed numerical comparison of these methods in section 7.

6. Convergence analysis of Information Newton's method

In this section, we introduce general update rules of information Newton's method in terms of probability densities and analyze their convergence rates in both distance and objective function value.

We briefly review the Riemannian structure of probability space as follows. Given a metric tensor $\mathcal{G}(\rho)$ and two probability densities $\rho_0, \rho_1 \in \mathcal{P}(\Omega)$, we denote the distance $\mathcal{D}(\rho_0, \rho_1)$ as follows

$$\mathcal{D}(\rho_0, \rho_1)^2 = \inf_{\hat{\rho}_s, s \in [0,1]} \left\{ \int_0^1 \int \partial_s \hat{\rho}_s \mathcal{G}(\hat{\rho}_s)^{-1} \partial_s \hat{\rho}_s dx ds : \hat{\rho}_s|_{s=0} = \rho_0, \hat{\rho}_s|_{s=1} = \rho_1 \right\}.$$

For the Wasserstein metric, $\mathcal{D}(\rho_0, \rho_1)$ is the Wasserstein-2 distance between ρ_0 and ρ_1 . Denote the inner product on cotangent space $T^*_{\rho}\mathcal{P}(\Omega)$ by

$$\langle \Phi_1, \Phi_2 \rangle_{\rho} = \int \Phi_1 \mathcal{G}(\rho)^{-1} \Phi_2 dx, \quad \Phi_1, \Phi_2 \in T^*_{\rho} \mathcal{P}(\Omega),$$

and $\|\Phi\|_{\rho}^2 = \langle \Phi, \Phi \rangle_{\rho}$. And we introduce the definition of the parallelism.

Definition 16 (Parallelism) We say that $\tau : T_{\rho_0}\mathcal{P}(\Omega) \to T_{\rho_1}\mathcal{P}(\Omega)$ is a parallelism from ρ_0 to ρ_1 , if for all $\Phi_1, \Phi_2 \in T_{\rho_0}\mathcal{P}(\Omega)$, it follows

$$\left\langle \Phi_{1},\Phi_{2}
ight
angle _{
ho_{0}}=\left\langle au\Phi_{1}, au\Phi_{2}
ight
angle _{
ho_{1}}.$$

To analyze the convergence rate, we introduce $\nabla^n E(\rho)$. This is a *n*-form on the cotangent space $T^*_{\rho}\mathcal{P}(\Omega)$, which is recursively defined by

$$\nabla^{n} E(\rho)(\Phi_{1},\ldots,\Phi_{n}) = \left. \frac{\partial}{\partial s} \nabla^{n-1} E(\operatorname{Exp}_{\rho}(s\Phi_{n}))(\tau_{s}\Phi_{1},\ldots,\tau_{s}\Phi_{n-1}) \right|_{s=0},$$

where τ_s is the parallelism from ρ to $\text{Exp}_{\rho}(s\Phi_n)$.

6.1 Convergence analysis in distance

The general update rule of the information Newton's method follows

$$\rho_{k+1} = \operatorname{Exp}_{\rho_k}(\alpha_k \Phi_k), \quad \mathcal{H}_E(\rho_k) \Phi_k + \mathcal{G}(\rho_k)^{-1} \frac{\delta E}{\delta \rho_k} = 0.$$
(25)

Here $\alpha_k > 0$ is a step size and $\operatorname{Exp}_{\rho_k}(\cdot)$ is the exponential map at ρ_k .

Recall that in the convergence proof of Euclidean Newton methods, it is assumed that $\nabla^2 f(x)$ is positive definite around a small neighbour of the optimal solution x^* . In the probability space, we assume that the following assumption holds analogously.

Assumption 1 Assume that there exists $\zeta, \delta_1, \delta_2, \delta_3 > 0$, such that for all ρ satisfying $\mathcal{D}(\rho, \rho^*) < \zeta$ and $\Phi_1, \Phi_2 \in T^*_{\rho}\mathcal{P}(\Omega)$, the following statements hold.

$$\nabla^2 E(\rho)(\Phi_1, \Phi_1) \ge \delta_1 \|\Phi_1\|_{\rho}^2, \tag{A1}$$

$$\nabla^2 E(\rho)(\Phi_1, \Phi_1) \le \delta_2 \|\Phi_1\|_{\rho}^2,$$
 (A2)

$$|\nabla^{3} E(\rho)(\Phi_{1}, \Phi_{1}, \Phi_{2})| \le \delta_{3} \|\Phi_{1}\|_{\rho}^{2} \|\Phi_{2}\|_{\rho}.$$
(A3)

Relying on Assumption 1, Theorem 17 shows the quadratic convergence rate of the Newton's method in the probability space.

Theorem 17 Suppose that Assumption 1 holds, ρ_k satisfies $\mathcal{D}(\rho_k, \rho^*) < \zeta$ and the step size $\alpha_k = 1$. Then, we have

$$\mathcal{D}(\rho_{k+1}, \rho^*) = O(\mathcal{D}(\rho_k, \rho^*)^2).$$

We present a sketch of the proof. For simplicity, we denote $T_k = \text{Exp}_{\rho_k}^{-1}(\rho^*)$.

Proposition 18 Suppose that Assumption 1 holds. Let τ be the parallelism from ρ_k to ρ_{k+1} . There exists a unique $R_k \in T^*_{\rho_k} \mathcal{P}(\Omega)$ such that

$$T_k = \tau^{-1} T_{k+1} + \Phi_k + R_k$$

Then, we have

$$||T_{k+1}||_{\rho_{k+1}} \le \frac{\delta_3}{\delta_1} ||T_k||_{\rho_k}^2 + \frac{\delta_2}{\delta_1} ||R_k||_{\rho_k}.$$

In order to provide an estimation on $||R_k||_{\rho_k}$, we introduce Lemma 19.

Lemma 19 For all $\Psi \in \mathcal{T}_{\rho_k}^* \mathcal{P}(\Omega)$, it follows

$$\int \Psi \mathcal{G}(\rho_k)^{-1} R_k dx = O(\|\Psi\|_{\rho_k} \|T_k\|_{\rho_k}^2).$$

Taking $\Psi = R_k$ in Lemma 19 yields $||R_k||_{\rho_k} = O(||T_k||_{\rho_k}^2)$. Because the geodesic curve has constant speed (Boothby, 1986), $||T_k||_{\rho_k}^2 = \mathcal{D}(\rho_k, \rho^*)^2$. As a result, we have

$$\mathcal{D}(\rho_{k+1},\rho^*) \leq \frac{\delta_2}{\delta_1} \mathcal{D}(\rho_k,\rho^*)^2 + \frac{\delta_3}{\delta_1} \|R_k\|_{\rho_k} = O(\mathcal{D}(\rho_k,\rho^*)^2).$$

6.2 Convergence analysis in objective function value

We next analyze the convergence rate based on our approximation methods in section 5. In practice, we use the approximated solution Φ_k to update ρ_k . Here Φ_k is the solution to the variational problem

$$\inf_{\Phi \in \mathcal{S}} \int \Phi \mathcal{H}_E(\rho_k) \Phi dx + 2 \int \Phi \mathcal{G}(\rho_k)^{-1} \frac{\delta E}{\delta \rho_k} dx + \lambda \int \Phi \mathcal{R}_{\mathcal{S}} \Phi dx.$$
(26)

Here \mathcal{H} is a linear subspace of $\mathcal{F}(\Omega)$, $\lambda \geq 0$ is a regularization parameter and $\int \Phi \mathcal{R}_{\mathcal{S}} \Phi dx$ is a regularization term in \mathcal{S} . For instance, if \mathcal{S} is an RKHS, then $\int \Phi \mathcal{R}_{\mathcal{H}} \Phi dx$ can be the squared norm of RKHS, i.e., $\|\Phi\|_{\mathcal{S}}^2$.

Suppose that $P: T^*_{\rho_k}\mathcal{P}(\Omega) \to \mathcal{S}$ is a projection operator from $T^*_{\rho_k}\mathcal{P}(\Omega)$ to \mathcal{S} and $P^*: \mathcal{S} \to T_{\rho_k}\mathcal{P}(\Omega)$ is its adjoint operator. Then, we can write $\hat{\Phi}_k$ in the closed-form formulation:

$$\Phi_k = -P(P^*\mathcal{H}_E(\rho_k)P + \mathcal{R}_S)^{-1}P^*\mathcal{G}(\rho_k)^{-1}\frac{\delta E}{\delta\rho_k}.$$
(27)

For simplicity, we use the following notations.

$$g_k = \mathcal{G}(\rho_k)^{-1} \frac{\delta E}{\delta \rho_k}, \quad \mathcal{H}_{E,P} = P(P^* \mathcal{H}_E(\rho_k) P + \mathcal{R}_S)^{-1} P^*.$$
(28)

For the subspace S and the regularization term $\lambda \int \Phi \mathcal{R}_S \Phi dx$, we further assume that the following three statements hold.

Assumption 2 There exists $\epsilon_1 \geq 0$, for all ρ_k satisfying $\mathcal{D}(\rho_k, \rho^*) < \zeta$, such that

$$\left| \int g_k (\mathcal{H}_{E,P} - \mathcal{H}_E(\rho_k)^{-1}) g_k dx \right| \le \epsilon_1 \int g_k \mathcal{H}_E(\rho_k)^{-1} g_k dx.$$
(A4)

There exists $\epsilon_2 \geq 0$, for all ρ_k satisfying $\mathcal{D}(\rho_k, \rho^*) < \zeta$, such that

$$\left| \int g_k (\mathcal{H}_{E,P} \mathcal{H}_E(\rho_k) \mathcal{H}_{E,P} - \mathcal{H}_{E,P}) g_k dx \right| \le \epsilon_2 \int g_k \mathcal{H}_{E,P} g_k dx.$$
(A5)

There exists $\delta_4 \geq 0$, for all ρ_k satisfying $\mathcal{D}(\rho_k, \rho^*) < \zeta$, such that

$$\left\|\mathcal{H}_{E,P}\mathcal{G}(\rho)^{-1}\Phi\right\|_{\rho_k} \le \delta_4 \left\|\Phi\right\|_{\rho_k}.$$
(A6)

The update rule in terms of density follows

$$\rho_{k+1} = \operatorname{Exp}_{\rho_k}(\alpha_k \Phi_k).$$

Theorem 20 Under Assumption 1 and 2, for ρ_k satisfying $\mathcal{D}(\rho_k, \rho^*) < \zeta$, with $\alpha_k = 1$, we have the linear convergence rate

$$E(\rho_{k+1}) - E(\rho^*) \le (\epsilon_1 + \epsilon_2 + \epsilon_1 \epsilon_2)(E(\rho_k) - E(\rho^*)) + \mathcal{O}((E(\rho_k) - E(\rho^*))^{3/2}).$$

From Theorem 20, we note that if the linear subspace S is appropriately chosen such that $\mathcal{H}_{E,P}$ is close to $\mathcal{H}_E(\rho_k)^{-1}$ in the sense of (A4) and (A5), then ϵ_1, ϵ_2 will be close to 0. This yields a sharp asymptotic convergence rate in terms of optimality gap, i.e., $E(\rho_k) - E(\rho^*)$.

Remark 21 We note that $\epsilon_2 = \mathcal{O}(\lambda)$. This comes from the following identity.

$$\mathcal{H}_{E,P}\mathcal{H}_{E}(\rho_{k})\mathcal{H}_{E,P} - \mathcal{H}_{E,P}$$

$$=P(P^{*}\mathcal{H}_{E}(\rho_{k})P + \lambda\mathcal{R}_{\mathcal{H}})^{-1}P^{*}\mathcal{H}_{E}(\rho_{k})P(P^{*}\mathcal{H}_{E}(\rho_{k})P + \lambda\mathcal{R}_{\mathcal{H}})^{-1}P^{*}$$

$$-P(P^{*}\mathcal{H}_{E}(\rho_{k})P + \lambda\mathcal{R}_{\mathcal{H}})^{-1}P^{*}$$

$$=\lambda P(P^{*}\mathcal{H}_{E}(\rho_{k})P + \lambda\mathcal{R}_{\mathcal{H}})^{-1}\mathcal{R}_{\mathcal{S}}(P^{*}\mathcal{H}_{E}(\rho_{k})P + \lambda\mathcal{R}_{\mathcal{H}})^{-1}P^{*}.$$
(29)

6.3 Convergence analysis in terms of samples

In practice, we replace ρ_k in the variational problem (26) by $\hat{\rho}_k(x) = \frac{1}{N} \sum_{n=1}^N \delta(x - x_k^n)$ to solve $\hat{\Phi}_k$. Here $x_k^n \sim \rho_k$. A natural question arises: with increasing sample numbers N, does $\hat{\Phi}_k$ from samples converge to Φ_k from distribution? Under further assumptions, the answer is yes and we postpone the justification in the appendix.

To establish the convergence rate, we further assume that the following statements hold.

Assumption 3 There exists $\epsilon_3 \geq 0$, for all ρ_k satisfying $\mathcal{D}(\rho_k, \rho^*) < \zeta$, such that

$$\left| \int (\hat{\Phi}_k - \Phi_k) g_k dx - \frac{1}{2} \int (\hat{\Phi}_k - \Phi_k) (\mathcal{H}_{E,P} \mathcal{H}_E(\rho_k) + \mathcal{H}_E(\rho_k) \mathcal{H}_{E,P}) g_k dx \right|$$

$$\leq \frac{\epsilon_3}{2} \int g_k \mathcal{H}_E(\rho_k)^{-1} g_k dx.$$
(A7)

There exists $\epsilon_4 \geq 0$, for all ρ_k satisfying $\mathcal{D}(\rho_k, \rho^*) < \zeta$, such that

$$\left| \int (\hat{\Phi}_k - \Phi_k) \mathcal{H}_E(\rho_k) (\hat{\Phi}_k - \Phi_k) dx \right| \le \epsilon_4 \int g_k \mathcal{H}_E(\rho_k)^{-1} g_k dx.$$
(A8)

The update rule in terms of density follows

$$\rho_{k+1} = \operatorname{Exp}_{\rho_k}(\alpha_k \hat{\Phi}_k).$$

Theorem 22 Under Assumption 1, 2 and 3, for ρ_k satisfying $\mathcal{D}(\rho_k, \rho^*) < \zeta$, with $\alpha_k = 1$, we have the linear convergence rate

$$E(\rho_{k+1}) - E(\rho^*) \le (\epsilon_1 + \epsilon_2 + \epsilon_1 \epsilon_2 + \epsilon_3 + \epsilon_4)(E(\rho_k) - E(\rho^*)) + \mathcal{O}((E(\rho_k) - E(\rho^*))^{3/2}).$$

7. Numerical experiments

In this section, we present numerical experiments to demonstrate the strength of information Newton's methods.

7.1 Toy examples

We compare particle implementations among Wasserstein Newton's methods with affine models 2/RKHS 3 (WNewton-a/WNewton-k), Wasserstein gradient flow (WGF), Hessian Approximated Lagrangian Langevin dynamics (HALLD) and Stein variational Newton's method with the scaled Hessian kernel (SVN-H) (Detommaso et al., 2018). We note that the update rule of WGF satisfies

$$x_{k+1}^n = x_k^n - \alpha_k (\nabla f(x_k^n) + \xi_k(x_k^n)).$$

The update rule of HALLD follows

$$x_{k+1}^{n} = x_{k}^{n} - \alpha_{k} \nabla^{2} f(x_{k}^{n})^{-1} (\nabla f(x_{k}^{n}) + \xi_{k}(x_{k}^{n})).$$

We note that the density evolution of HALLD and HAMCMC are identical to each other. In other words, we replace the Brownian motion in HAMCMC by ξ_k in HALLD. Here ξ_k is an approximation of $\nabla \log \rho_k$. For all compared methods, we use constant step sizes. For the calculation of ξ_k , we apply KDE with Gaussian kernels and the kernel bandwidth is selected by the Brownian Motion method (Wang and Li, 2019)[section 5.1]. This method adaptively learns the bandwidth from samples generated by Brownian motions.

We first consider a 1D target density $\rho^*(x) \propto \exp(-f(x))$, where $f(x) = \frac{1}{2}(x^2-1)^2$. For WGF, we set $\alpha_k = 0.01$. For SVGD, we set $\alpha_1 = 0.1$ and adjust the step size by Adagrad (Duchi et al., 2011). For WNewton-a and WNewton-k, we let $\alpha_k = 1$, $\epsilon = 0$ and $\gamma = 0$. Namely, we do not apply the hybrid update. For HALLD and SVN-H, we set $\alpha_k = 1$.

The sample number follows N = 100. The initial distribution follows $\mathcal{N}(0, 0.01)$. We plot the distribution after 2, 5, 10, 20 iterations in Figure 2. Although we use affine/kernel approximations to compute the Newton's direction, WNewton-a and WNewton-k tend to



Figure 2: Comparison among WGF, SVGD, WNewton-a, WNewton-k, HALLD and SVN-H in 1D toy example. Left to right: sample distribution after 2, 5, 10, 20 iterations.

converge to the target density and they are faster than WGF. SVGD has similar performance with WGF. HALLD and SVN-H have some particle which tend to diverge. This may result from that the target density is not log-concave.

Then, we let the target density ρ^* to be a 2D bimodal distribution (Rezende and Mohamed, 2015). For WGF, we set $\alpha_k = 0.1$. For SVGD, we set $\alpha_1 = 1$ and adjust the step size via Adagrad. For WNewton-a, we apply the hybrid update and set $\alpha_k = 0.2, \epsilon = 0$ and $\gamma = 0.5$. For WNewton-k, we set $\alpha_k = 1, \epsilon = 0, \gamma = 0$. For HALLD, we set $\alpha_k = 0.2$. For SVN-H, we set $\alpha = 1$.

The initial distribution follows $\mathcal{N}([0, 10]', I)$. We plot the distribution after 2, 5, 10, 20 iterations with N = 100 samples in Figure 3. WNewton-k converges rapidly toward the target density. HALLD fails to converge because $\nabla^2 f$ becomes singular on certain sample points. SVN-H barely moves because the initial distribution is not close enough to the target distribution. SVGD converges slower than WGF. The Wasserstein Newton's direction helps samples to converge faster towards the target density with robustness.

Next, we present numerical results on a 2D double-banana shape posterior density in (Detommaso et al., 2018). For WGF, we set $\alpha_k = 0.002$. For SVGD, we set $\alpha_1 = 0.1$ and adjust the step size via Adagrad. For WNewton-a, we apply the hybrid update and set $\alpha_k = 0.2, \epsilon = 0$ and $\gamma = 0.001$. For WNewton-k, we set $\alpha_k = 1, \epsilon = 0, \gamma = 0$. For HALLD and SVN-H, we set $\alpha_k = 1$.

Similarly, we plot the distribution after 2, 5, 10, 20 iterations with N = 100 samples in Figure 4. WNewton-k and SVN-H converges toward the posterior distribution in no more than 5 iterations. WNewton-a collapses around the center of the lower banana. WGF and SVGD take nearly 20 iterations to converge. HALLD converges rapidly but it diverge at iteration 20. Here we notice that WNewton does not require heavy tunes of step sizes. The step size $\alpha_k = 1$ usually leads to robust performance.

7.2 Conditioned diffusion

The conditioned diffusion example is a 100-dimensional model from a Langevin SDE, with state $u_t : [0, T] \to \mathbb{R}$ and dynamics give by

$$du_t = \frac{\beta u(1-u^2)}{1+u^2}dt + dx_t, \quad u_0 = 0.$$

Here $x = (x_t)_{t\geq 0}$ is the standard Brownian motion. The goal is to infer the driving process x_t and its pushfoward to the state u. Detailed setup of this test case can be found in (Detommaso et al., 2018).

We compare WNewton-a with WGF, SVGD, SVN-H and HALLD. We do not compare WNewton-k because per-iteration computation cost in the current implementation is too heavy on this test case with N = 1000 and d = 100. For WGF, we set $\alpha_k = 0.01$. For SVGD, we set $\alpha_1 = 0.1$ and adjust step sizes via Adagrad. For WNewton-a, SVN-H and HALLD, we set $\alpha_k = 1$. From Figure 5, we note that the posterior mean (which captures the trends of true path) from WNewton-a, SVN-H and HALLD almost converge in approximately 10 iterations. Meanwhile, the posterior mean from WGF and SVGD takes 50-100 iterations to converge. Compared to SVN-H, WNewton-a tends to have narrower credible interval. The credible interval of HALLD in [0, 0.5] after 100 iterations has larger fluctuation.



Figure 3: Comparison among WGF, SVGD, WNewton-a, WNewton-k, HALLD and SVN-H in 2D toy example. Left to right: sample distribution after 2, 5, 10, 20 iterations.



Figure 4: Comparison among WGF, SVGD, WNewton-a, WNewton-k, HALLD and SVN-H in 2D double banana example. Left to right: sample distribution after 2, 5, 10, 20 iterations.



Figure 5: Comparison among WGF, SVGD, WNewton-a, HALLD and SVN-H in 100D conditioned diffusion example. Left to right: sample distribution after 10, 50, 100 iterations. Red dots: noisy observations. Purple line: ground truth. Blue line: posterior mean. Shaded area: 90% credible interval.

7.3 Bayesian logistic regression

We perform the standard Bayesian logistic regression experiment on the Covertype dataset, following the settings in (Liu and Wang, 2016). We compare WNewton-a and WNewton-k with MCMC, SVGD (Liu and Wang, 2016), and WGF. The performances of SVN-H and HALLD on this test example are not ideal. For the calculation of ξ_k in WGF and WNewtona, we use KDE with Gaussian kernel and the bandwidth is selected by the median method, which is the same as (Liu and Wang, 2016). The sample number follows N = 50. The mini-batch size for stochastic gradient and Hessian evaluations in each iteration is 100.

We first discuss the choice of step sizes. The initial step sizes for the compared methods are given in Table 1. Except for SVGD, the initial step sizes are selected from $\{i \cdot 10^j | i \in \{1, 2, 5\}, j \in \{-3, \ldots, -7\}\}$ to ensure the best performance. For SVGD, we use the initial step size in (Liu and Wang, 2016) and adjust step sizes by Adagrad. For MCMC, WGF and WNewton-k, the step size is multiplied by 0.9 every 100 iterations. For WNewton-a, the step size is multiplied by 0.82 every 100 iterations.

Method	MCMC	SVGD	WGF	WNewton-a	WNewton-k
Step size α_1	1e-5	0.05	1e-5	2e-3	2e-3

Table 1: Initial step sizes for algorithms in comparison.

We then elaborate on the implementation details of compared methods. For WNewtonk, we apply the block-diagonal approximation to accelerate the computation. For WNewtona and WNewton-k, we set $\epsilon = 1$ and use the hybrid update with $\gamma = 5 \times 10^{-3}$ and $\gamma = 10^{-3}$ respectively.

From Figure 6, we observe that WNewton-k has the best performance in terms of test accuracy and test log-likelihood and it converges much faster compared to other methods. Namely, WNewton-k has ideal performance on test test tests in less than 200 iterations. WNewton-a and WNewton-k achieves higher test log-likelihood. This indicates that the approximated Wasserstein Newton's direction leads to better generalization on the test set.



Figure 6: Comparison of different methods on Bayesian logistic regression, averaged over 10 independent trials. The shaded areas show the variance over 10 trials. Left: Test accuracy; Right: Test log-likelihood.

8. Conclusion

In this paper, we introduce information Newton's flows (second-order optimization methods) for optimization problems in probability space arising from Bayesian statistics, inverse problems, and machine learning. Here two information metrics, such as Fisher-Rao metric and Wasserstein-2 metric, are considered. Several examples and convergence analysis of the proposed second-order methods are provided. Following the fact that the Wasserstein gradient flow of KL divergence formulates the Langevin dynamics, we derive the Wasserstein Newton's flow of KL divergence as Newton's Langevin dynamics. Focusing on Newton's Langevin dynamics, we study analytical examples in one-dimensional sample space and Gaussian families. We further propose practical sampling efficient algorithms, in affine models and RKHS, to implement Newton's Langevin dynamics. We show the convergence rate of information Newton's method with approximated solutions. The numerical examples in Bayesian sampling problems demonstrate the effectiveness of the proposed method.

Appendix A. Definitions and notations

In this section, we present several definitions and notations used in this paper. We briefly review the concept of self-adjoint operator.

Definition 23 (Self-adjoint) Suppose that V is a Hilbert space and let $\mathcal{H} : V \to V^*$ be a linear operator. V^* is the adjoint space of V, which consists of all linear functionals on V. Let (f, v) = (v, f) = f(v) denote the coupling of $v \in V$ and $f \in V^*$. The adjoint operator of \mathcal{H} is the unique linear operator $\mathcal{H}^* : V \to V^*$, which satisfies

$$(\mathcal{H}v_1, v_2) = (v_1, \mathcal{H}^* v_2), \quad \forall v_1, v_2 \in V.$$

We say that \mathcal{H} is self-adjoint if $\mathcal{H} = \mathcal{H}^*$.

Remark 24 If $V = \mathbb{R}^d$ is the Euclidean space, then the linear operator \mathcal{H} can be viewed as a matrix in $\mathbb{R}^{d \times d}$. Then, to say that \mathcal{H} is self-adjoint operator is equivalent to say that \mathcal{H} is a symmetric matrix.

We define positive definite operators as follows.

Definition 25 Suppose that V is a Hilbert space and let $\mathcal{H} : V \to V^*$ be a self-adjoint linear operator. We say that \mathcal{H} is positive definite, if $(\mathcal{H}v, v) > 0$ for all $v \in V$, $v \neq 0$.

Appendix B. Proofs in section 3

In this section, we present details and proofs for propositions in section 3. Proposition 26 provides a sufficient condition to ensure that the Hessian operator is injective (invertible).

Proposition 26 Suppose that $g_{\rho}(\text{Hess } E(\rho)\sigma, \sigma) > 0$ for all $\sigma \neq 0, \sigma \in T_{\rho}\mathcal{P}(\Omega)$. Namely, $\mathcal{H}_{E}(\rho)$ is positive definite. Then, $\text{Hess } E(\rho)$ is injective.

Proof If there exist $\sigma_1, \sigma_2 \in T_{\rho}\mathcal{P}(\Omega)$ such that $\operatorname{Hess} E(\rho)\sigma_1 = \operatorname{Hess} E(\rho)\sigma_2$. Then,

$$g_{\rho}((\sigma_1 - \sigma_2), \operatorname{Hess} E(\rho)(\sigma_1 - \sigma_2)) = \int (\sigma_1 - \sigma_2)G(\rho)^{-1} \operatorname{Hess} E(\rho)(\sigma_1 - \sigma_2)dx = 0.$$

By our assumption $g_{\rho}(\text{Hess } E(\rho)\sigma, \sigma) > 0$ for all $\sigma \neq 0$, we have $\sigma_1 = \sigma_2$.

B.1 Proof of Proposition 3

The geodesic curve $\hat{\rho}_s$ satisfies geodesic equation

$$\begin{cases} \partial_s \hat{\rho}_s - \mathcal{G}(\hat{\rho}_s)^{-1} \Phi_s = 0, \\ \partial_s \Phi_s + \frac{1}{2} \frac{\delta}{\delta \hat{\rho}_s} \left(\int \Phi_s \mathcal{G}(\hat{\rho}_s)^{-1} \Phi_s dx \right) = 0, \end{cases}$$
(30)

with initial values $\hat{\rho}_s|_{s=0} = \rho$ and $\Phi_s|_{s=0} = \Phi$. For the first-order derivative, it follows

$$\frac{d}{ds}E(\hat{\rho}_s) = \int \partial_s \hat{\rho}_s \frac{\delta E}{\delta \hat{\rho}_s} dx = \int \Phi_s \mathcal{G}(\hat{\rho}_s)^{-1} \frac{\delta E}{\delta \hat{\rho}_s} dx,$$

where we utilize the fact that $\mathcal{G}(\hat{\rho}_s)$ is self-adjoint. For the second-order derivative,

$$\begin{split} \frac{d^2}{ds^2} E(\hat{\rho}_s) &= \int \partial_s \Phi_s \mathcal{G}(\hat{\rho}_s)^{-1} \frac{\delta E}{\delta \hat{\rho}_s} dx + \int \partial_s \hat{\rho}_s \frac{\delta}{\delta \hat{\rho}_s} \left(\frac{d}{ds} E(\hat{\rho}_s) \right) dx \\ &= -\frac{1}{2} \int \mathcal{A}(\hat{\rho}_s) (\Phi_s, \Phi_s) \mathcal{G}(\hat{\rho}_s)^{-1} \frac{\delta E}{\delta \hat{\rho}_s} dx + \int \mathcal{A}(\hat{\rho}_s) \left(\Phi_s, \frac{\delta E}{\delta \hat{\rho}_s} \right) \mathcal{G}(\hat{\rho}_s)^{-1} \Phi_s dx \\ &+ \int \int \left(\mathcal{G}(\hat{\rho}_s)^{-1} \Phi_s \right) (y) \frac{\delta^2 E}{\delta \hat{\rho}_s^2} (x, y) \left(\mathcal{G}(\hat{\rho}_s)^{-1} \Phi_s \right) (x) dx dy. \end{split}$$

Based on the definition of $\mathcal{H}_E(\rho)$, (4) is proved by setting s = 0 in the above formula. To prove (5), we introduce Lemma 27.

Lemma 27 Let \mathcal{H} be a self-adjoint linear operator from $T^*_{\rho}\mathcal{P}(\Omega) \to T_{\rho}\mathcal{P}(\Omega)$. Namely $\mathcal{H}^* = \mathcal{H}$. Suppose that $\int \Phi \mathcal{H} \Phi dx = 0$ for all $\Phi \in T^*_{\rho}\mathcal{P}(\Omega)$. Then, $\mathcal{H} = 0$.

Proof Because \mathcal{H} is self-adjoint and linear, for any $\Phi \in T^*_{\rho}\mathcal{P}(\Omega)$, it follows

$$\mathcal{H}\Phi = \frac{1}{2}\frac{\delta}{\delta\Phi}\int \Phi\mathcal{H}\Phi dx = 0.$$

This completes the proof.

Note that Hess $E(\rho)$ is self-adjoint w.r.t. the metric tensor $G(\rho)$, namely

$$(\operatorname{Hess} E(\rho))^* \mathcal{G}(\rho) = \mathcal{G}(\rho) \operatorname{Hess} E(\rho), \quad \mathcal{G}(\rho)^{-1} (\operatorname{Hess} E(\rho))^* = \operatorname{Hess} E(\rho) \mathcal{G}(\rho)^{-1}.$$

where $(\text{Hess } E(\rho))^*$ is the adjoint operator of $\text{Hess } E(\rho)$. This tells that $\text{Hess } E(\rho)\mathcal{G}(\rho)^{-1}$ is self-adjoint. We have the following relationship.

$$\int \Phi \mathcal{H}_E(\rho) \Phi dx = g_\rho(\operatorname{Hess} E(\rho)\sigma, \sigma) = \int \Phi \mathcal{G}(\rho)^{-1} \operatorname{Hess} E(\rho) \Phi dx.$$

As a direct result of Proposition 26, it follows $\mathcal{H}_E(\rho) = \operatorname{Hess} E(\rho) \mathcal{G}(\rho)^{-1}$.

B.2 Newton's flows under Fisher-Rao metric

For Fisher-Rao metric, the geodesic curve $\hat{\rho}_s$ satisfies

$$\begin{cases} \partial_s \hat{\rho}_s - \rho_s \left(\Phi_s - \int \Phi_s \hat{\rho}_s dy \right) = 0, \\ \partial_s \Phi_s + \frac{1}{2} \Phi_s^2 - \left(\int \Phi_s \hat{\rho}_s dy \right) \Phi_s = 0 \end{cases}$$

And the bi-linear operator $\mathcal{A}^F(\rho)$ follows

$$\mathcal{A}^{F}(\rho)(\Phi_{1},\Phi_{2}) = \Phi_{1}\Phi_{2} - \left(\int \Phi_{2}\rho dy\right)\Phi_{1} - \left(\int \Phi_{1}\rho dy\right)\Phi_{2}.$$
(31)

For simplicity, we let $\mathbb{E}_{\rho}[\Phi] = \int \Phi \rho dx$, where $\Phi \in T^*_{\rho} \mathcal{P}(\Omega)$.

Proposition 28 (Fisher-Rao Newton's flow) For an objective function $E : \mathcal{P}(\Omega) \to \mathbb{R}$, the Fisher-Rao Newton's flow follows

$$\begin{cases} \partial_t \rho_t - \rho_t (\Phi_t - \mathbb{E}_{\rho_t} [\Phi_t]) = 0, \\ \mathcal{H}_E^F(\rho_t) \Phi_t - \rho_t \left(\frac{\delta E}{\delta \rho_t} - \mathbb{E}_{\rho_t} \left[\frac{\delta E}{\delta \rho_t} \right] \right) = 0, \end{cases}$$
(32)

where $\mathcal{H}_E^F(\rho): T^*_{\rho}\mathcal{P}(\Omega) \to T_{\rho}\mathcal{P}(\Omega)$ defines a bi-linear form: for $\Phi \in T^*_{\rho}\mathbb{P}(\Omega)$,

$$\int \Phi \mathcal{H}_{E}^{F}(\rho) \Phi dx = \frac{1}{2} \int \mathcal{A}^{F}(\rho) \left(\Phi, \frac{\delta E}{\delta \rho}\right) (\Phi - \mathbb{E}_{\rho}[\Phi]) \rho dx + \int \int \rho(y) (\Phi(y) - \mathbb{E}_{\rho}[\Phi]) \frac{\delta^{2} E}{\delta \rho^{2}}(x, y) dy \rho(x) (\Phi(x) - \mathbb{E}_{\rho}[\Phi]) dx.$$
(33)

Proof Based on Proposition 3, we only need to prove that

$$\int \mathcal{A}^{F}(\rho)(\Phi,\Phi)\mathcal{G}^{F}(\rho)^{-1}\frac{\delta E}{\delta\rho}dx = \int \mathcal{A}^{F}(\rho)\left(\Phi,\frac{\delta E}{\delta\rho}\right)\mathcal{G}^{F}(\rho)^{-1}\Phi dx$$

The left hand side follows

$$\int \mathcal{A}^{F}(\rho)(\Phi,\Phi)\mathcal{G}^{F}(\rho)^{-1}\frac{\delta E}{\delta\rho}dx$$

= $\int \left(\Phi^{2} - 2\mathbb{E}_{\rho}[\Phi]\Phi\right)\left(\frac{\delta E}{\delta\rho} - \mathbb{E}_{\rho}\left[\frac{\delta E}{\delta\rho}\right]\right)\rho dx$
= $\int \left(\Phi - \mathbb{E}_{\rho}[\Phi]\right)\left(\frac{\delta E}{\delta\rho} - \mathbb{E}_{\rho}\left[\frac{\delta E}{\delta\rho}\right]\right)\Phi\rho dx - \mathbb{E}_{\rho}[\Phi]\int \left(\frac{\delta E}{\delta\rho} - \mathbb{E}_{\rho}\left[\frac{\delta E}{\delta\rho}\right]\right)\Phi\rho dx.$

The right hand side satisfies

$$\int \mathcal{A}^{F}(\rho) \left(\Phi, \frac{\delta E}{\delta \rho}\right) \mathcal{G}^{F}(\rho)^{-1} \Phi dx$$

=
$$\int \left(\Phi \frac{\delta E}{\delta \rho} - \mathbb{E}_{\rho} \left[\frac{\delta E}{\delta \rho}\right] \Phi - \mathbb{E}_{\rho}[\Phi] \frac{\delta E}{\delta \rho}\right) \left(\Phi - \mathbb{E}_{\rho}[\Phi]\right) \rho dx$$

=
$$\int \left(\Phi - \mathbb{E}_{\rho}[\Phi]\right) \left(\frac{\delta E}{\delta \rho} - \mathbb{E}_{\rho} \left[\frac{\delta E}{\delta \rho}\right]\right) \Phi \rho dx - \mathbb{E}_{\rho}[\Phi] \int \frac{\delta E}{\delta \rho} \left(\Phi - \mathbb{E}_{\rho}[\Phi]\right) \rho dx$$

We also observe that

$$\int \frac{\delta E}{\delta \rho} \left(\Phi - \mathbb{E}_{\rho}[\Phi] \right) \rho dx = \mathbb{E}_{\rho} \left[\Phi \frac{\delta E}{\delta \rho} \right] - \mathbb{E}_{\rho}[\Phi] \mathbb{E}_{\rho} \left[\frac{\delta E}{\delta \rho} \right] = \int \left(\frac{\delta E}{\delta \rho} - \mathbb{E}_{\rho} \left[\frac{\delta E}{\delta \rho} \right] \right) \Phi \rho dx.$$

Hence, the left hand side is equal to the right hand side.

Example 6 (Fisher-Rao Newton's flow of KL divergence) Suppose that $E(\rho)$ evaluates the KL divergence from ρ to $\rho^* \sim \exp(-f)$. This objective functional also writes

$$E(\rho) = \int (\rho \log \rho + f\rho) dx.$$

We derive that

$$\frac{\delta E}{\delta \rho}(x) = \log \rho(x) + f + 1, \quad \frac{\delta^2 E}{\delta \rho^2}(x, y) = \frac{\delta(x - y)}{\rho(y)}.$$

Based on Proposition 28, we can compute that (4) is equivalent to

$$\begin{split} \int \Phi \mathcal{H}_E(\rho) \Phi dx &= \frac{1}{2} \int \left(\Phi^2 - 2\mathbb{E}_{\rho}[\Phi] \Phi \right) \left(\log \rho + f - \mathbb{E}_{\rho}[\log \rho + f] \right) \rho dx \\ &+ \int \left(\Phi(x) - \mathbb{E}_{\rho}[\Phi] \right) \rho(x) \int \frac{\delta(y - x)}{\rho(y)} \left(\Phi(y) - \mathbb{E}_{\rho}[\Phi] \right) \rho(y) dy dx \\ &= \frac{1}{2} \int \left(\log \rho + f - \mathbb{E}_{\rho}[\log \rho + f] \right) \Phi^2 \rho dx \\ &- \mathbb{E}_{\rho}[\Phi] \int \left(\log \rho + f - \mathbb{E}_{\rho}[\log \rho + f] \right) \Phi \rho dx \\ &+ \int \Phi^2 \rho dx - \left(\int \Phi \rho dx \right)^2. \end{split}$$

Hence, the operator $\mathcal{H}^F_E(\rho)$ follows

$$\begin{aligned} \mathcal{H}_{E}^{F}(\rho)\Phi =& \frac{1}{2}\left(\log\rho + f - \mathbb{E}_{\rho}[\log\rho + f]\right)\Phi\rho - \frac{1}{2}\left(\int\left(\log\rho + f - \mathbb{E}_{\rho}[\log\rho + f]\right)\Phi\rho dy\right)\rho \\ &- \frac{1}{2}\mathbb{E}_{\rho}[\Phi]\left(\log\rho + f - \mathbb{E}_{\rho}[\log\rho + f]\right)\rho + \Phi\rho - \mathbb{E}_{\rho}[\Phi]\rho \\ =& \frac{1}{2}\left(2 + \log\rho + f - \mathbb{E}_{\rho}[\log\rho + f]\right)\left(\Phi - \mathbb{E}_{\rho}[\Phi]\right)\rho \\ &- \frac{1}{2}\left(\mathbb{E}_{\rho}[\Phi(\log\rho + f)] - \mathbb{E}_{\rho}[\Phi]\mathbb{E}_{\rho}[(\log\rho + f)]\right)\rho.\end{aligned}$$

Example 7 (Fisher-Rao Newton's flow of interaction energy) Consider an interaction energy

$$E(\rho) = \frac{1}{2} \int \int \rho(x) W(x, y) \rho(y) dx dy,$$

where W(x, y) = W(y, x) is a kernel function. The interaction energy also formulates the MMD, see details in (Gretton et al., 2012). We can compute that

$$\frac{\delta E}{\delta \rho}(x) = \int W(x,y)\rho(y)dy, \quad \frac{\delta^2 E}{\delta \rho^2}(x,y) = W(x,y).$$

We denote $(W * \rho)(x) = \int W(x, y)\rho(y)dx$. Based on Proposition 28, it follows

$$\begin{split} &\int \Phi \mathcal{H}_E^F(\rho) \Phi dx \\ = &\frac{1}{2} \int (\Phi^2 - 2\mathbb{E}_{\rho}[\Phi] \Phi) (W * \rho - \mathbb{E}_{\rho}[W * \rho]) \rho dx \\ &+ \int \int (\Phi(y) - \mathbb{E}_{\rho}[\Phi]) W(x, y) \rho(y) \rho(x) (\Phi(x) - \mathbb{E}_{\rho}[\Phi]) dy dx \\ = &\frac{1}{2} \int \Phi^2 (W * \rho - \mathbb{E}_{\rho}[W * \rho]) \rho dx - \mathbb{E}_{\rho}[\Phi] \left(\int \Phi (W * \rho - \mathbb{E}_{\rho}[W * \rho]) \rho dx \right) \\ &+ \int \int \Phi(x) \rho(x) W(x, y) \Phi(y) \rho(y) dx dy + (\mathbb{E}_{\rho}[\Phi])^2 \left(\int \int \rho(x) W(x, y) \rho(y) dx dy \right) \\ &- 2\mathbb{E}_{\rho}[\Phi] \left(\int \int \rho(x) W(x, y) \Phi(x) \rho(y) dx dy \right). \end{split}$$

Hence, the operator $\mathcal{H}^F_E(\rho)$ satisfies

$$\begin{aligned} \mathcal{H}_{E}^{F}(\rho)\Phi(x) =& \frac{1}{2}(W*\rho - \mathbb{E}_{\rho}[W*\rho])\rho\Phi - \frac{1}{2}\left(\int\Phi(W*\rho - \mathbb{E}_{\rho}[W*\rho])\rho dy\right)\rho \\ &\quad -\frac{1}{2}\mathbb{E}_{\rho}[\Phi](W*\rho - \mathbb{E}_{\rho}[W*\rho])\rho + (W*(\rho\Phi))\rho \\ &\quad +\mathbb{E}_{\rho}[W*\rho]\mathbb{E}_{\rho}[\Phi]\rho - \mathbb{E}_{\rho}[W*(\rho\Phi)]\rho - \mathbb{E}_{\rho}[\Phi](W*\rho)\rho \\ &\quad = \frac{1}{2}(W*\rho - \mathbb{E}_{\rho}[W*\rho])(\Phi - \mathbb{E}_{\rho}[\Phi])\rho - \frac{1}{2}\left(\mathbb{E}_{\rho}[\Phi(W*\rho)] - \mathbb{E}_{\rho}[\Phi]\mathbb{E}_{\rho}[W*\rho]\right)\rho \\ &\quad + (W*(\rho\Phi) - \mathbb{E}_{\rho}[W*(\rho\Phi)])\rho - \mathbb{E}_{\rho}[\Phi]\left((W*\rho) - \mathbb{E}_{\rho}[W*\rho]\right)\rho. \end{aligned}$$

Example 8 (Fisher-Rao Newton's flow of cross entropy) Suppose that $E(\rho)$ is the cross entropy, i.e., reverse KL divergence. It evaluates the KL divergence from a given density ρ^* to ρ

$$E(\rho) = -\int \log\left(\frac{\rho}{\rho^*}\right)\rho^* dx = -\int (\log\rho)\rho^* dx + \int (\log\rho^*)\rho^* dx.$$

It is equivalent to optimize $E(\rho) = -\int (\log \rho) \rho^* dx$. We compute that

$$\frac{\delta E}{\delta \rho}(x) = -\frac{\rho^*(x)}{\rho(x)}, \quad \frac{\delta^2 E}{\delta \rho^2}(x,y) = \frac{\rho^*(y)}{\rho^2(y)}\delta(x-y).$$

Proposition 28 indicates that

$$\begin{split} \int \Phi \mathcal{H}_{E}^{F}(\rho) \Phi dx &= \frac{1}{2} \int (\Phi^{2} - 2\mathbb{E}_{\rho}[\Phi] \Phi) (-\rho^{*}/\rho + \mathbb{E}_{\rho}[\rho^{*}/\rho]) \rho dx \\ &+ \int (\Phi(x) - \mathbb{E}_{\rho}[\Phi]) \rho(x) \int \frac{\rho^{*}(y)}{\rho^{2}(y)} \delta(x - y) (\Phi(y) - \mathbb{E}_{\rho}[\Phi]) \rho(y) dy dx \\ &= \frac{1}{2} \int \Phi^{2}(\rho - \rho^{*}) dx - \mathbb{E}_{\rho}[\Phi] \int \Phi(\rho - \rho^{*}) dx + \int (\Phi - \mathbb{E}_{\rho}[\Phi])^{2} \rho^{*} dx \\ &= \frac{1}{2} (\mathbb{E}_{\rho}[\Phi^{2}] - \mathbb{E}_{\rho^{*}}[\Phi^{2}]) - \mathbb{E}_{\rho}[\Phi] (\mathbb{E}_{\rho}[\Phi] - \mathbb{E}_{\rho^{*}}[\Phi]) \\ &+ \mathbb{E}_{\rho^{*}}[\Phi^{2}] - 2\mathbb{E}_{\rho}[\Phi] \mathbb{E}_{\rho^{*}}[\Phi] + (\mathbb{E}_{\rho}[\Phi])^{2} \\ &= \frac{1}{2} (\mathbb{E}_{\rho}[\Phi^{2}] + \mathbb{E}_{\rho^{*}}[\Phi^{2}]) - \mathbb{E}_{\rho}[\Phi] \mathbb{E}_{\rho^{*}}[\Phi]. \end{split}$$

Hence, the operator $\mathcal{H}^F_E(\rho)$ follows

$$\mathcal{H}_E^F(\rho)\Phi = \frac{1}{2}((\Phi - \mathbb{E}_{\rho^*}[\Phi])\rho + (\Phi - \mathbb{E}_{\rho}[\Phi])\rho^*).$$

B.3 Newton's flows under Wasserstein metric

For Wasserstein metric, the geodesic curve $\hat{\rho}_s$ satisfies

$$\begin{cases} \partial_s \hat{\rho}_s + \nabla \cdot (\hat{\rho}_s \nabla \Phi_s) = 0, \\ \partial_s \Phi_s + \frac{1}{2} \| \nabla \Phi_s \|^2 = 0. \end{cases}$$

The bi-linear operator $\mathcal{A}^W(\rho)$ follows

$$\mathcal{A}^{W}(\rho)(\Phi_{1},\Phi_{2}) = \langle \nabla \Phi_{1}, \nabla \Phi_{2} \rangle.$$

Proposition 29 (Wasserstein Newton's flow) For an objective functional $E: \mathcal{P}(\Omega) \rightarrow \mathbb{R}$, the Wasserstein Newton's flow follows

$$\begin{cases} \partial_t \rho_t + \nabla \cdot (\rho \nabla \Phi_t) = 0, \\ \mathcal{H}_E^W(\rho_t) \Phi_t - \nabla \cdot \left(\rho_t \nabla \frac{\delta E}{\delta \rho_t}\right) = 0. \end{cases}$$
(34)

Here $\mathcal{H}_E(\rho): T^*_{\rho}\mathcal{P}(\Omega) \to T_{\rho}\mathcal{P}(\Omega)$ defines a bi-linear form: for $\Phi \in T^*_{\rho}\mathcal{P}(\Omega)$,

$$\int \Phi \mathcal{H}_{E}^{W}(\rho) \Phi dx = \int \int \left\langle \nabla \Phi(x), \nabla_{x} \nabla_{y} \frac{\delta^{2} E}{\delta \rho^{2}}(x, y) \nabla \Phi(y) \right\rangle \rho(x) \rho(y) dx dy + \int \left\langle \nabla \Phi, \nabla^{2} \frac{\delta E}{\delta \rho} \nabla \Phi \right\rangle \rho dx.$$
(35)

Proof Based on integration by parts, we observe that

$$\int \mathcal{A}^{W}(\rho)(\Phi, \Phi)\mathcal{G}^{W}(\rho)^{-1}\frac{\delta E}{\delta\rho}dx$$
$$= -\int \|\nabla\Phi\|^{2}\nabla \cdot \left(\rho\nabla\frac{\delta E}{\delta\rho}\right)dx$$
$$= \int \left\langle\nabla\frac{\delta E}{\delta\rho}, \nabla\|\nabla\Phi\|^{2}\right\rangle\rho dx$$
$$= 2\int \left\langle\nabla\frac{\delta E}{\delta\rho}, \nabla^{2}\Phi\nabla\Phi\right\rangle\rho dx,$$

and

$$\begin{split} &\int \mathcal{A}^{W}(\rho) \left(\Phi, \frac{\delta E}{\delta \rho}\right) \mathcal{G}^{W}(\rho)^{-1} \Phi dx \\ &= -\int \left\langle \nabla \Phi, \nabla \frac{\delta E}{\delta \rho} \right\rangle \nabla \cdot (\rho \nabla \Phi) dx \\ &= \int \left\langle \nabla \left\langle \nabla \Phi, \nabla \frac{\delta E}{\delta \rho} \right\rangle, \nabla \Phi \right\rangle \rho dx \\ &= \int \left\langle \nabla \Phi, \nabla^{2} \Phi \nabla \frac{\delta E}{\delta \rho} \right\rangle \rho dx + \int \left\langle \nabla \Phi, \nabla^{2} \frac{\delta E}{\delta \rho} \nabla \Phi \right\rangle \rho dx. \end{split}$$

Combining above two observations with Proposition 3, we derive

$$\begin{split} \int \Phi \mathcal{H}_E^W(\rho) \Phi dx &= \int \left\langle \nabla \Phi, \nabla^2 \frac{\delta E}{\delta \rho} \nabla \Phi \right\rangle \rho dx \\ &+ \int \int \nabla \cdot (\rho \nabla \Phi)(y) \frac{\delta^2 E}{\delta \rho^2}(x,y) \nabla \cdot (\rho \nabla \Phi)(x) dx dy \\ &= \int \int \left\langle \nabla \Phi(x), \nabla_x \nabla_y \frac{\delta^2 E}{\delta \rho^2}(x,y) \nabla \Phi(y) \right\rangle \rho(x) \rho(y) dx dy \\ &+ \int \left\langle \nabla \Phi, \nabla^2 \frac{\delta E}{\delta \rho} \nabla \Phi \right\rangle \rho dx. \end{split}$$

This proves Proposition 29.

Example 9 (Wasserstein Newton's flow of interaction energy) Consider an interaction energy

$$E(\rho) = \frac{1}{2} \int \int \rho(x) W(x,y) \rho(y) dx dy.$$

Combining with previous computations, Proposition 29 yields that

$$\begin{split} \int \Phi \mathcal{H}_E^W(\rho) \Phi dx &= \int \int \left\langle \nabla \Phi(x), \nabla_x \nabla_y W(x, y) \nabla \Phi(y) \right\rangle \rho(x) \rho(y) dx dy \\ &+ \int \left\langle \nabla \Phi(x), \int \nabla_x^2 W(x, y) \rho(y) dy \nabla \Phi(x) \right\rangle \rho(x) dx \\ &= \frac{1}{2} \mathbb{E}_{x, y \sim \rho} \begin{bmatrix} \nabla \Phi(x) \\ \nabla \Phi(y) \end{bmatrix}^T \begin{bmatrix} \nabla_{xx}^2 W(x, y) & \nabla_{xy}^2 W(x, y) \\ \nabla_{yx}^2 W(x, y) & \nabla_{yy}^2 W(x, y) \end{bmatrix} \begin{bmatrix} \nabla \Phi(x) \\ \nabla \Phi(y) \end{bmatrix}. \end{split}$$

Based on integration by parts, the operator $\mathcal{H}^W_E(\rho)$ is given by

$$\mathcal{H}^W_E(\rho)\Phi = -\nabla \cdot \left(\rho(\nabla^2_{xy}W * (\rho\nabla\Phi))\right) - \nabla \cdot \left(\rho(\nabla^2_{xx}W * \rho)\nabla\Phi\right)$$

Example 10 (Wasserstein Newton's flow of cross entropy) Suppose that $E(\rho)$ evaluates the KL divergence from a given density ρ^* to ρ

$$E(\rho) = -\int \log\left(\frac{\rho}{\rho^*}\right)\rho^* dx = -\int (\log\rho)\rho^* dx + \int (\log\rho^*)\rho^* dx.$$

It is equivalent to optimize $E(\rho) = -\int (\log \rho) \rho^* dx$. Proposition 29 yields

$$\int \Phi \mathcal{H}_E^W(\rho) \Phi dx = \int \nabla \cdot (\rho(x) \nabla \Phi(x)) \int \frac{\rho^*(y)}{\rho^2(y)} \delta(x-y) \nabla \cdot (\rho(y) \nabla \Phi(y)) dy dx$$
$$- \int \left\langle \nabla \Phi(x), \nabla_x^2 \left(\frac{\rho^*(x)}{\rho(x)}\right) \nabla \Phi(x) \right\rangle \rho(x) dx$$
$$= \int (\rho^{-1} \nabla \cdot (\rho \nabla \Phi))^2 \rho^* dx - \int \left\langle \nabla \Phi, \nabla^2 \left(\frac{\rho^*}{\rho}\right) \nabla \Phi \right\rangle \rho dx.$$

Hence, the operator $\mathcal{H}^W_E(\rho)$ satisfies

$$\mathcal{H}_{E}^{W}(\rho)\Phi = \nabla \cdot \left(\rho \nabla \left(\frac{\rho^{*}}{\rho^{2}} \nabla \cdot (\rho \nabla \Phi)\right)\right) + \nabla \cdot \left(\rho \nabla^{2} \left(\frac{\rho^{*}}{\rho}\right) \nabla \Phi\right).$$

Remark 30 For simplicity of presentations, we only present the Hessian formulas for Fisher-Rao and Wasserstein information metrics. In fact, there are many interesting generalized Hessian formulas in Li (2019) from Hessian transport metrics. We leave systematic studies of Newton's flows for general metrics in future works.

We summarize formulations of Hessian-related operators $\mathcal{H}_E(\rho)$ under both Fisher-Rao metric and Wasserstein metric.

Objective functional $E(\rho)$	$\mathcal{H}^F_E(ho)\Phi$		
KL divergence: $\int (\rho \log \rho + f \rho) dx.$	$\frac{1}{2} \left(2 + \log \rho + f - \mathbb{E}_{\rho}[\log \rho + f]\right) \left(\Phi - \mathbb{E}_{\rho}[\Phi]\right) \rho$ $-\frac{1}{2} \left(\mathbb{E}_{\rho}[\Phi(\log \rho + f)] - \mathbb{E}_{\rho}[\Phi]\mathbb{E}_{\rho}[(\log \rho + f)]\right) \rho.$		
Interaction energy: $\frac{1}{2} \int \int \rho(x) W(x,y) \rho(y) dx dy$	$ \frac{1}{2}(W*\rho - \mathbb{E}_{\rho}[W*\rho])(\Phi - \mathbb{E}_{\rho}[\Phi])\rho -\frac{1}{2}(\mathbb{E}_{\rho}[\Phi(W*\rho)] - \mathbb{E}_{\rho}[\Phi]\mathbb{E}_{\rho}[W*\rho])\rho +(W*(\rho\Phi) - \mathbb{E}_{\rho}[W*(\rho\Phi)])\rho -\mathbb{E}_{\rho}[\Phi]((W*\rho) - \mathbb{E}_{\rho}[W*\rho])\rho. $		
Reverse KL divergence: $\int (\log \rho^* - \log \rho) \rho^* dx$	$\frac{1}{2}(\Phi - \mathbb{E}_{\rho^*}[\Phi])\rho + \frac{1}{2}(\Phi - \mathbb{E}_{\rho}[\Phi])\rho^*.$		

Table 2: The formulation of $\mathcal{H}^F_E(\rho)$ under the Fisher-Rao metric.

Objective functional $E(\rho)$	$\mathcal{H}^W_E(ho)\Phi$	
KL divergence: $\int (\rho \log \rho + f \rho) dx.$	$ abla^2: (ho abla^2 \Phi) - abla \cdot (ho abla^2 f abla \Phi).$	
Interaction energy: $\frac{1}{2} \int \int \rho(x) W(x,y) \rho(y) dx dy$	$-\nabla \cdot ((\nabla_{xy}^2 W * (\nabla \Phi \rho))\rho) - \nabla \cdot ((\nabla_{xx}^2 W * \rho_t)\rho \nabla \Phi).$	
Reverse KL divergence: $\int (\log \rho^* - \log \rho) \mu dx$	$\nabla \cdot \left(\rho \nabla \left(\frac{\rho^*}{\rho^2} \nabla \cdot (\rho \nabla \Phi) \right) \right) + \nabla \cdot \left(\rho \nabla^2 \left(\frac{\rho^*}{\rho} \right) \nabla \Phi \right).$	

Table 3: The formulation of $\mathcal{H}_{E}^{W}(\rho)$ under the Wasserstein metric.

B.4 Wasserstein Newton's flows in Gaussian families

In this subsection, we study information Newton's flows in Gaussian families with respect to Wasserstein metric. We leave the proof of Proposition 5 in next subsection. Let \mathbb{P}^n and \mathbb{S}^n represent the space of symmetric positive definite matrices and symmetric matrices with size $n \times n$ respectively.

We let \mathcal{N}_n^0 denote multivariate Gaussian densities with zero means. Each $\rho \in \mathcal{N}_n^0$ is uniquely determined by its covariance matrix $\Sigma \in \mathbb{P}^n$. So we can view $\mathcal{N}_n^0 \simeq \mathbb{P}^n$. The Wasserstein metric $\mathcal{G}^W(\rho)$ on $\mathcal{P}(\mathbb{R}^d)$ induces the Wasserstein metric $\mathcal{G}^W(\Sigma)$ on \mathbb{P}^n , see (Takatsu, 2008; Modin, 2017; Malagò et al., 2018). For $\Sigma \in \mathbb{P}^n$, tangent space and cotangent space follow

$$T_{\Sigma}\mathbb{P}^n \simeq T_{\Sigma}^*\mathbb{P}^n \simeq \mathbb{S}^n.$$

Definition 31 (Wasserstein metric in Gaussian families) Given $\Sigma \in \mathbb{P}^n$, the Wasserstein metric tensor $\mathcal{G}^W(\Sigma) : \mathbb{S}^n \to \mathbb{S}^n$ is defined by

$$\mathcal{G}^W(\Sigma)^{-1}S = 2(\Sigma S + S\Sigma).$$

It defines an inner product on the tangent space $T_{\Sigma}\mathbb{P}^n$. Namely, for $A_1, A_2 \in T_{\Sigma}\mathbb{P}^n \simeq \mathbb{S}^n$

$$g_{\Sigma}^{W}(A_1, A_2) = \operatorname{tr}(A_1 \mathcal{G}^{W}(\Sigma) A_2) = \operatorname{tr}(S_1 \mathcal{G}^{W}(\Sigma)^{-1} S_2) = 4 \operatorname{tr}(S_1 \Sigma S_2).$$

Here $S_i \in T^*_{\Sigma} \mathbb{P}^n \simeq \mathbb{S}^n$ is the solution to discrete Lyapunov equation

$$A_i = 2(\Sigma S_i + S_i \Sigma), \quad i = 1, 2.$$

For $\Sigma \in \mathbb{P}^n$, there exits a unique solution to discrete Lyapunov equation. Again, we focus on the case where the objective functional $E(\Sigma)$ evaluates the KL divergence from ρ with covariance matrix Σ to a target Gaussian density ρ^* with covariance matrix Σ^* . Then, $E(\Sigma)$ satisfies (8).

Proposition 32 (Gradient and Hessian operators in \mathbb{P}^n) The gradient operator follows

$$\operatorname{grad}^{W} E(\Sigma) = \mathcal{G}^{W}(\Sigma)^{-1} \nabla E(\Sigma) = \Sigma(\Sigma^{*})^{-1} + (\Sigma^{*})^{-1} \Sigma - 2I.$$

And the Hessian operator satisfies that for all $A \in \mathbb{S}^n$,

$$g_{\Sigma}^{W}(A, \operatorname{Hess}^{W} E(\Sigma)A) = 4\operatorname{tr}(S\Sigma S(\Sigma^{*})^{-1}) + 4\operatorname{tr}(S^{2}),$$

where S is the unique solution to $A = 2(\Sigma S + S\Sigma)$.

Given $A \in \mathbb{S}^n$, the geodesic curve $\hat{\Sigma}_s$ with $\hat{\Sigma}_s|_{s=0} = \Sigma$ and $\partial_s \hat{\Sigma}_s|_{s=0} = A$ follows $\hat{\Sigma}_s = (I + 2sS)\Sigma(I + 2sS)$, where $S = \mathcal{G}(\Sigma)^{-1}A$ is the solution to $A = 2(\Sigma S + S\Sigma)$. We can compute that

$$E(\hat{\Sigma}_s) = \frac{1}{2} (\operatorname{tr}((I+2sS)\Sigma(I+2sS)(\Sigma^*)^{-1}) - n - \log \det((I+2sS)\Sigma(I+2sS)(\Sigma^*)^{-1})).$$

The Taylor expansion of $\log \det(I + sS)$ w.r.t. s satisfies

$$\log \det(I + sS) = s \operatorname{tr}(S) - \frac{s^2}{2} \operatorname{tr}(S^2) + o(s^2).$$

Hence, the first-order derivative follows

$$\frac{\partial}{\partial s} E(\Sigma(s)) \bigg|_{s=0} = \operatorname{tr}(S\Sigma(\Sigma^*)^{-1}) + \operatorname{tr}(\Sigma S(\Sigma^*)^{-1}) - 2\operatorname{tr}(S)$$
$$= \operatorname{tr}\left(S\left(\Sigma(\Sigma^*)^{-1} + (\Sigma^*)^{-1}\Sigma - 2I\right)\right).$$

By the definition $\frac{\partial}{\partial s} E(\Sigma(s))|_{s=0} = \operatorname{tr}(S \operatorname{grad} E(\Sigma))$, this yields $\operatorname{grad} E(\Sigma) = \Sigma(\Sigma^*)^{-1} + (\Sigma^*)^{-1}\Sigma - 2I$ and the second-order derivative follows

$$\left. \frac{\partial^2}{\partial s^2} E(\Sigma(s)) \right|_{s=0} = 4 \operatorname{tr}(S\Sigma S(\Sigma^*)^{-1}) + 4 \operatorname{tr}(S^2).$$

This completes the proof.

Similarly, let us consider the linear self-adjoint operator $\mathcal{H}_E^W(\Sigma) : \mathbb{S}^n \to \mathbb{S}^n$, which defines a bi-linear form

$$\operatorname{tr}(S\mathcal{H}_E^W(\Sigma)S) = g_{\Sigma}^W(A, \operatorname{Hess}^W E(\Sigma)A) = 4\operatorname{tr}(S\Sigma S(\Sigma^*)^{-1}) + 4\operatorname{tr}(S^2).$$

We can compute that $\mathcal{H}_E(\Sigma)$ is uniquely defined by

$$\mathcal{H}_E(\Sigma)S = 2\Sigma S(\Sigma^*)^{-1} + 2(\Sigma^*)^{-1}S\Sigma + 4S, \quad \forall S \in \mathbb{S}^n.$$

Because $\operatorname{tr}(S\mathcal{H}_E(\Sigma)S) = 4\operatorname{tr}(S\Sigma S(\Sigma^*)^{-1}) + 4\operatorname{tr}(S^2) > 0$ for $S \neq 0, S \in \mathbb{S}^n$, \mathcal{H}_E is injective and invertible. Now, we are ready to present the Newton's flow of KL divergence in Gaussian families.

Proposition 33 The Newton's flow of KL divergence in Gaussian families follows

$$\begin{cases} \dot{\Sigma}_t - 2(S\Sigma_t + \Sigma S_t) = 0, \\ 2\Sigma_t S_t (\Sigma^*)^{-1} + 2(\Sigma^*)^{-1} S_t \Sigma_t + 4S_t = -(\Sigma_t (\Sigma^*)^{-1} + (\Sigma^*)^{-1} \Sigma_t - 2I). \end{cases}$$
(36)

Proof The Newton's flow follows

$$\dot{\Sigma}_t - (\operatorname{Hess}^W E(\Sigma_t))^{-1} \operatorname{grad}^W E(\Sigma_t) = 0.$$

We note that $\operatorname{Hess}^W E(\Sigma)\mathcal{G}^W(\Sigma)^{-1} = \mathcal{H}^W_E(\Sigma)$, which implies

$$(\operatorname{Hess}^{W} E(\Sigma))^{-1} = \mathcal{G}^{W}(\Sigma)^{-1} \mathcal{H}_{E}^{W}(\Sigma)^{-1}.$$

Hence, we can reformulate the Newton's flow by

$$\begin{cases} \dot{\Sigma}_t - \mathcal{G}^W(\Sigma_t)^{-1} S_t = 0, \\ \mathcal{H}_E^W(\Sigma_t) S_t = -\operatorname{grad}^W E(\Sigma_t) \end{cases}$$

From the formulations of $\mathcal{G}(\Sigma)^{-1}$, $\operatorname{grad}^W E(\Sigma)$ and $\mathcal{H}_E(\Sigma)$, we obtain (36).

Example 11 In one dimensional case, the second equation in (36) has an explicit solution $S_t = -\frac{(\Sigma^*)^{-1}\Sigma_t - 1}{2((\Sigma^*)^{-1}\Sigma_t + 1)}$. Let $\Sigma_t = Y_t^2$, where $Y_t > 0$. Then, the first equation in (36) turns to

$$2Y_t \dot{Y}_t + 4Y_t^2 \frac{(\Sigma^*)^{-1} Y_t^2 - 1}{2((\Sigma^*)^{-1} Y_t^2 + 1)} = 0,$$

or equivalently,

$$\dot{Y}_t + \frac{(\Sigma^*)^{-1}Y_t - Y_t^{-1}}{(\Sigma^*)^{-1} + Y_t^{-2}} = 0.$$
(37)

Let $f(Y) = \frac{1}{2}((\Sigma^*)^{-1}Y^2 - 1 - \log((\Sigma^*)^{-1}Y^2))$. Then, we have $\nabla f(Y) = (\Sigma^*)^{-1}Y - Y^{-1}$ and $\nabla^2 f(Y) = (\Sigma^*)^{-1} + Y^{-2}$. Hence, the Newton's flow (37) coincides with Newton's flow of f(X) in Euclidean space. We also note that (37) is identical to the evolution of Σ_t in Proposition 8 by substituting $\Sigma_t = Y_t^2$.

B.5 Proof of Proposition 5

We first prove that Σ_t is positive definite. We formulate that

$$\partial_t E(\Sigma_t) = \operatorname{tr}(\partial_t \Sigma_t \nabla E(\Sigma_t)) = 2 \operatorname{tr}(S_t \Sigma_t ((\Sigma^*)^{-1} - \Sigma_t^{-1})) \\ = \operatorname{tr}(S_t (\Sigma_t (\Sigma^*)^{-1} + (\Sigma^*)^{-1} \Sigma_t - 2I)) = -\operatorname{tr}(S(2\Sigma S_t (\Sigma^*)^{-1} + 2(\Sigma^*)^{-1} S_t \Sigma_t + 4S_t)) \\ = -4 \operatorname{tr}(S_t \Sigma_t S_t (\Sigma^*)^{-1}) - 4 \operatorname{tr}(S_t^2) \le 0.$$

As a result, $E(\Sigma_t)$ is non-increasing. Applying the idea of proof in (Wang and Li, 2019, Theorem 1), we can establish that Σ_t is positive definite. Then, we examine that Φ_t satisfies (3). We observe that

$$\nabla^2 : (\rho_t \nabla^2 \Phi_t) - \nabla \cdot (\rho_t \nabla^2 f \nabla \Phi_t) - \nabla \cdot (\rho_t \nabla f) - \Delta \rho_t$$

=2\sum 2 : (S_t\rho_t(x)) - 2\sum \cdot (\rho_t(x)(\Sigma^*)^{-1}S_tx) - \sum \cdot (\rho_t(x)(\Sigma^*)^{-1}x) - \Delta\rho_t.

We note that $\nabla \rho_t(x) = -\Sigma_t^{-1} x \rho_t(x)$ and $\nabla^2 \rho_t(x) = -\Sigma_t^{-1} \rho_t(x) + \Sigma_t^{-1} x x^T \Sigma_t^{-1} \rho_t(x)$. Hence, we derive all four terms in the above equation as follows. First, it is easy to observe that

$$\nabla^2 : (\rho_t(x)S_t) = \operatorname{tr}(S_t \nabla^2 \rho_t(x)), \quad -\Delta \rho_t = -\nabla^2 : (\rho_t I) = -\operatorname{tr}(\nabla^2 \rho_t(x)).$$

We can also compute that

$$-\nabla \cdot (\rho_t(\Sigma^*)^{-1}S_tx)$$

$$= -\sum_{i=1}^n \partial_i (\rho_t(x)WS_tx)_i$$

$$= -\sum_{i=1}^n \left[\rho_t(x)\partial_i ((\Sigma^*)^{-1}S_tx)_i + (WS_tx)_i\partial_i\rho_t(x) \right]$$

$$= -\rho_t(x) \left[\operatorname{tr}((\Sigma^*)^{-1}S_t) + ((\Sigma^*)^{-1}S_tx)^T (-\Sigma_t^{-1}x) \right]$$

$$= -\rho_t(x) \operatorname{tr}(S_t(\Sigma^*)^{-1}(I - \Sigma_t^{-1}xx^T))$$

$$= \frac{1}{2} \operatorname{tr}((\Sigma_t S_t(\Sigma^*)^{-1} + (\Sigma^*)^{-1}S_t\Sigma_t)\nabla^2\rho_t(x)).$$

Taking $S_t = I$ into the above equation yields

$$-\nabla \cdot (\rho_t(\Sigma^*)^{-1}x) = \frac{1}{2} \operatorname{tr}((\Sigma_t(\Sigma^*)^{-1} + (\Sigma^*)^{-1}\Sigma_t)\nabla^2 \rho_t(x)).$$

Because (Σ_t, S_t) satisfies (36), we have

$$2\nabla^{2} : (S_{t}\rho_{t}(x)) - 2\nabla \cdot (\rho_{t}(x)(\Sigma^{*})^{-1}S_{t}x) - \nabla \cdot (\rho_{t}(x)(\Sigma^{*})^{-1}x) - \Delta\rho_{t}$$

= tr((2S_{t} + \Sigma_{t}S_{t}(\Sigma^{*})^{-1} + (\Sigma^{*})^{-1}S_{t}\Sigma_{t} + \Sigma_{t}(\Sigma^{*})^{-1} + (\Sigma^{*})^{-1}\Sigma_{t} - 2I)\nabla^{2}\rho_{t}(x))
=0.

This completes the proof.

Appendix C. Details in section 4

In this section, we present detailed discussion of Wasserstein Newton's flow and Newton's Langevin dynamics with particular examples.

C.1 Connections and differences with HAMCMC

HAMCMC approximates the dynamics of

$$dX_t = -(\nabla^2 f(X_t))^{-1} \left(\nabla f(X_t) + \Gamma(X_t)\right) dt + \sqrt{2\nabla^2 f(X_t)^{-1}} dB_t,$$

where $\Gamma_i(x) = \sum_{j=1} \frac{\partial}{\partial x_j} \left(\left(\nabla^2 f(X_t) \right)^{-1} \right)_{i,j}$. Here $\Gamma(x)$ is a correction term to ensure that ρ_t converges to ρ^* . The evolution of ρ_t follows

$$\partial_t \rho_t = \nabla \cdot \left(\left((\nabla^2 f)^{-1} \nabla f + \Gamma \right) \rho_t \right) + \nabla^2 : \left((\nabla^2 f)^{-1} \rho_t \right).$$

We formulate the above equation as

$$\partial_t \rho_t = \nabla \cdot \left(\rho_t (\nabla^2 f)^{-1} (\nabla f + \nabla \log \rho_t)\right) = \nabla \cdot \left(\rho_t \mathbf{v}_t\right),\tag{38}$$

where we denote $\mathbf{v}_t = (\nabla^2 f)^{-1} (\nabla f + \nabla \log \rho_t)$. Moreover, \mathbf{v}_t satisfies

$$-\nabla \cdot (\rho_t \nabla^2 f \mathbf{v}_t) - \nabla \cdot (\rho_t \nabla f) - \Delta \rho_t = 0.$$

On the other hand, replacing $\nabla \Phi_t^{\text{Newton}}$ by $\mathbf{v}_t^{\text{Newton}}$ in the Newton's direction equation (3) yields

$$\nabla^2 : (\rho_t \nabla \mathbf{v}_t^{\text{Newton}}) - \nabla \cdot (\rho_t \nabla^2 f \mathbf{v}_t^{\text{Newton}}) - \nabla \cdot (\rho_t \nabla f) - \Delta \rho_t = 0$$

Hence, $\mathbf{v}_t \neq \mathbf{v}_t^{\text{Newton}}$. Then, (38) is different from information Newton's flow because the term $\nabla^2 : (\rho_t \nabla \mathbf{v}_t^{\text{Newton}})$ is not considered.

C.2 Connections and differences with Newton's flows in Euclidean space

We recall that the density evolution of particle's gradient flow in Euclidean space corresponds to the Wasserstein gradient flow (Villani, 2008). We notice that this relationship does not hold for the Wasserstein Newton's flow.

Consider an objective function:

$$E(\rho) = \int \rho(x) f(x) dx$$

where f(x) is a given smooth function. Here we notice that minimize ρ for $E(\rho)$ in probability space is equivalent to minimize x for f(x) in Euclidean space. Namely, the support of the optimal solution ρ contains all global minimizers of f(x). The gradient flow in Euclidean space of each particle follows

$$dX_t = -\nabla f(X_t)dt,$$

A known fact is that the density evolution of particles satisfies the following continuity equation

$$\partial_t \rho_t = \nabla \cdot (\rho_t \nabla f) = -\operatorname{grad}^W E(\rho_t),$$

which is the Wasserstein gradient flow of $E(\rho)$ in probability space.

We next show that Newton's flow in Euclidean space of each particle does not coincide with the Wasserstein Newton's flow in probability space. For simplicity, we assume that f(x) is strictly convex so $\nabla^2 f(x)$ is invertible for all x. Here, the Euclidean Newton's flow of each particle follows

$$dX_t = -(\nabla^2 f(X_t))^{-1} \nabla f(X_t) dt$$

The density evolution of particles satisfies the continuity equation

$$\partial_t \rho_t = \nabla \cdot (\rho_t (\nabla^2 f)^{-1} \nabla f). \tag{39}$$

On the other hand, the Wasserstein Newton's flow writes

$$\partial_t \rho_t + \nabla \cdot (\rho_t \nabla \Phi_t^{\text{Newton}}) = 0, \qquad (40)$$

where Φ_t^{Newton} is the unique solution to

$$-\nabla \cdot (\rho_t \nabla^2 f \nabla \Phi) - \nabla \cdot (\rho_t \nabla f) = 0.$$
(41)

We note that in general equation (39) can be different from equation (40). Later on in Lemma 35, we formulate the following Hodge decomposition of the Euclidean Newton's direction

$$-(\nabla^2 f)^{-1} \nabla f = \nabla \Phi_t^{\text{Newton}} + \boldsymbol{\xi}_t,$$

where $\nabla \cdot (\rho_t \nabla^2 f \boldsymbol{\xi}_t) = 0$. Here, the constraint on $\boldsymbol{\xi}_t$ does not necessarily ensure that $\nabla \cdot (\rho_t \boldsymbol{\xi}_t) = 0$. Hence, equation (39) can be different from equation (40).

Remark 34 In one dimensional case or f is a quadratic function, there exists Φ^{Newton} , such that $-(\nabla^2 f)^{-1}\nabla f = \nabla \Phi^{\text{Newton}}$. Hence equation (39) is same as equation (40). We also show an example of $\boldsymbol{\xi} \neq 0$. Let $\Omega = \mathbb{R}^2$ and we define

$$f(x) = \log(\exp(x_1) + \exp(x_2)) + \frac{\lambda}{2}(x_1^2 + x_2^2),$$

where $\lambda > 0$ is a parameter. For simplicity, we denote $p_1 = \exp(x_1)/(\exp(x_1) + \exp(x_2))$ and $p_2 = \exp(x_1)/(\exp(x_1) + \exp(x_2))$. Then, we can compute that the gradient and Hessian of f(x) follows

$$\nabla f(x) = \begin{bmatrix} p_1 + \lambda x_1 \\ p_2 + \lambda x_2 \end{bmatrix}, \quad \nabla^2 f(x) = \begin{bmatrix} p_1 p_2 + \lambda & -p_1 p_2 \\ -p_1 p_2 & p_1 p_2 + \lambda \end{bmatrix}.$$

Because $p_1p_2 + \lambda > 0$ and $det(\nabla^2 f(x)) = \lambda^2 + 2\lambda p_1p_2 > 0$, $\nabla^2 f(x)$ is positive definite. We note that

$$(\nabla^2 f(x))^{-1} \nabla f(x) = \frac{1}{\lambda^2 + 2\lambda p_1 p_2} \begin{bmatrix} p_1 p_2 + \lambda & p_1 p_2 \\ p_1 p_2 & p_1 p_2 + \lambda \end{bmatrix} \begin{bmatrix} p_1 + \lambda x_1 \\ p_2 + \lambda x_2 \end{bmatrix}$$
$$= \frac{1}{\lambda^2 + 2\lambda p_1 p_2} \begin{bmatrix} p_1 p_2 (1 + \lambda (x_1 + x_2)) + \lambda (p_1 + \lambda x_1) \\ p_1 p_2 (1 + \lambda (x_1 + x_2)) + \lambda (p_2 + \lambda x_2) \end{bmatrix}$$
$$= : \begin{bmatrix} F_1(x) \\ F_2(x) \end{bmatrix}.$$

If $(\nabla^2 f(x))^{-1} \nabla f(x)$ is a gradient vector field, we shall have

$$\partial_{x_2} F_1(x) = \partial_{x_1} F_2(x).$$

However, we can examine that

$$\partial_{x_2} F_1(x) = \frac{p_1 p_2}{\lambda + 2p_1 p_2} \left(1 + p_1 (1 + \lambda(x_1 + x_2)) + \frac{2\lambda p_1 (\lambda(p_1 + \lambda x_1) + p_1 p_2 (1 + \lambda(x_1 + x_2)))}{\lambda^2 + 2\lambda p_1 p_2} \right).$$

$$\partial_{x_1} F_2(x) = \frac{p_1 p_2}{\lambda + 2p_1 p_2} \left(1 + p_2 (1 + \lambda(x_1 + x_2)) + \frac{2\lambda p_2 (\lambda(p_2 + \lambda x_2) + p_1 p_2 (1 + \lambda(x_1 + x_2)))}{\lambda^2 + 2\lambda p_1 p_2} \right).$$

This indicates that $(\nabla^2 f(x))^{-1} \nabla f(x)$ is not a gradient vector field. Hence, $\boldsymbol{\xi} \neq 0$.

Lemma 35 For given $\rho \in \mathcal{P}(\mathbb{R}^d)$, there exists a unique $\Phi \in T^*_{\rho}\mathcal{P}(\mathbb{R}^d)$ (up to a constant shrift) and a vector field $\boldsymbol{\xi} : \mathbb{R}^d \to \mathbb{R}^d$ satisfying $\nabla \cdot (\rho \nabla^2 f \boldsymbol{\xi}) = 0$ such that

$$-(\nabla^2 f(x))^{-1} \nabla f(x) = \nabla \Phi(x) + \boldsymbol{\xi}(x).$$

Proof We first show the existence of $\Phi \in T^*_{\rho}\mathcal{P}(\mathbb{R}^d)$ and $\boldsymbol{\xi}$. Note that Φ is the solution to

$$-\nabla \cdot (\rho \nabla^2 f \nabla \Phi) = \nabla \cdot (\rho \nabla f).$$

Denote $\mathcal{H}\Phi = -\nabla \cdot (\rho \nabla^2 f \nabla \Phi)$. Then, for $\Phi \neq 0$, we have

$$\int \Phi \mathcal{H} \Phi dx = \int \nabla \Phi^T \nabla^2 f \nabla \Phi \rho dx > 0$$

Hence, \mathcal{H} is a positive definite operator and it is invertible. Thus $\Phi = \mathcal{H}^{-1} \left(\nabla \cdot (\rho \nabla f) \right)$ exists. Because $\nabla^2 f \boldsymbol{\xi} = \nabla f - \nabla^2 f \nabla \Phi$, it follows

$$\nabla \cdot (\nabla^2 f \boldsymbol{\xi}) = \nabla \cdot (\rho \nabla f) - \nabla \cdot (\rho \nabla^2 f \nabla \Phi) = 0.$$

Hence, $\boldsymbol{\xi}$ also exists. We next prove the uniqueness. Suppose that $\nabla^2 f(x)^{-1} \nabla f(x) = \nabla \Phi_1(x) + \boldsymbol{\xi}_1(x) = \nabla \Phi_2(x) + \boldsymbol{\xi}_2(x)$. Then, we have $\nabla \Phi_1 - \nabla \Phi_2 = \boldsymbol{\xi}_2 - \boldsymbol{\xi}_1$. Hence

$$\int (\Phi_1 - \Phi_2) \mathcal{H}(\Phi_1 - \Phi_2) dx = \int (\nabla \Phi_1 - \nabla \Phi_2)^T \nabla^2 f(\nabla \Phi_1 - \nabla \Phi_2) \rho dx$$
$$= \int (\nabla \Phi_1 - \nabla \Phi_2)^T \nabla^2 f(\boldsymbol{\xi}_2 - \boldsymbol{\xi}_1) \rho dx = -\int (\Phi_1 - \Phi_2) \nabla \cdot (\rho \nabla^2 f(\boldsymbol{\xi}_2 - \boldsymbol{\xi}_1)) dx = 0.$$

Because \mathcal{H} is positive definite, this yields that $\Phi_1 - \Phi_2 = 0$ (up to a spatial constant).

C.3 Newton's Langevin dynamics in one dimensional sample space

In this subsection, we provide examples of Newton's Langevin dynamics in one dimensional sample space. In particular, similar to the Ornstein–Uhlenbeck (OU) process in classical Langevin dynamics, we derive a closed form solution to Newton's OU process.

Here we assume that $\Omega = \mathbb{R}$ and f is strictly convex. The essence of Newton's Langevin dynamics is to compute Φ_t^{Newton} from the Wasserstein Newton's direction equation (3). Proposition 26 ensures the uniqueness of the solution to (3). For the simplicity of notations, we neglect the subscript t.

Proposition 36 Suppose that $\rho > 0$ and let $u = \nabla \Phi$. Then, the Newton's direction equation (3) reduces to an ODE

$$u'' + u'(\log \rho)' - f''u - f' - (\log \rho)' = 0.$$
(42)

Proof In 1-dimensional case, the equation (3) follows

$$\nabla^2(\rho\nabla^2\Phi) - \nabla(\rho\nabla^2f\nabla\Phi) - \nabla(\rho\nabla f) - \nabla^2\rho = 0$$

The above equation is equivalent to

$$\rho \nabla^3 \Phi + \nabla \rho \nabla^2 \Phi - \rho \nabla^2 f \nabla \Phi - \rho \nabla f - \nabla \rho + C = 0,$$

where C is a constant. Because $\rho \in \mathcal{P}(\mathbb{R}) \subset L^1(\mathbb{R})$. Hence $\lim_{|x|\to\infty} \rho(x) = 0$, which indicates C = 0. Suppose that $\rho > 0$ and let $u = \nabla \Phi$. Dividing both sides by ρ , we obtain

$$u'' + u'\rho'/\rho - f''u - f' - \rho'/\rho = 0.$$

By the fact that $\rho'/\rho = (\log \rho)'$, we derive (42).

We consider the case where f'(x) and $(\log \rho)'(x)$ are affine functions. Then, ODE (42) has a closed-form solution. Applying ODE (42), we obtain the exact formulation of Newton's

Langevin dynamics in Proposition 8. For the rest of this section, we present the proof of Proposition 8.

Proof In section 3 Proposition 5, we show that if the evolution of X_t follows NLD, then X_t follows the Gaussian distribution. We first solve the Newton's direction from ODE (42). Suppose that $(\log \rho)'(x) = \Sigma^{-1}(x - \mu)$. The ODE turns to be

$$u'' - u'\Sigma^{-1}(x-\mu) - (\Sigma^*)^{-1}u - (\Sigma^*)^{-1}(x-\mu^*) + \Sigma^{-1}(x-\mu) = 0.$$

We can examine that the following u is a solution to the above ODE.

$$u(x) = \frac{\Sigma^{-1} - (\Sigma^*)^{-1}}{\Sigma^{-1} + (\Sigma^*)^{-1}} x - \frac{2\Sigma^{-1}}{\Sigma^{-1} + (\Sigma^*)^{-1}} \mu + \mu^*.$$

Hence, we have $\Phi^{\text{Newton}}(x) = \frac{\Sigma^* - \Sigma}{2(\Sigma^* + \Sigma)} x^2 - \frac{2\Sigma^*}{\Sigma^* + \Sigma} \mu x + \mu^* x$. As a result, NLD follows

$$dX_t = \left(\frac{\Sigma^* - \Sigma_t}{\Sigma^* + \Sigma_t} X_t - \frac{2\Sigma^*}{\Sigma^* + \Sigma_t} \mu_t + \mu^*\right) dt.$$

The dynamics of μ_t satisfies

$$d\mu_t = d\mathbb{E}[X_t] = \mathbb{E}[dX_t] = \left(\frac{\Sigma^* - \Sigma_t}{\Sigma^* + \Sigma_t}\mu_t - \frac{2\Sigma^*}{\Sigma^* + \Sigma_t}\mu_t + \mu^*\right)dt = (-\mu_t + \mu^*)dt.$$

This indicates that $\mu_t = \mu^* + e^{-t}(\mu_0 - \mu^*)$. The dynamics of Σ_t follows

$$d\Sigma_t = d(\mathbb{E}[X_t^2] - \mu_t^2) = 2\mathbb{E}[X_t dX_t] - 2\mu_t d\mu_t$$

=2\left[\frac{\Sigma^* - \Sigma_t}{\Sigma^* + \Sigma_t} \left(\Sigma_t + \mu_t^2\right) - \frac{2\Sigma^*}{\Sigma^* + \Sigma_t} \mu_t^2 + \mu^* \mu_t - \mu_t (-\mu_t + \mu^*)\right] dt = 2\frac{\Sigma^* - \Sigma_t}{\Sigma^* + \Sigma_t} \Sigma_t dt.

We can rewrite that

$$dt = \frac{(\Sigma^* + \Sigma_t)d\Sigma_t}{2(\Sigma^* - \Sigma_t)\Sigma_t} = \left(\frac{1}{(\Sigma^* - \Sigma_t)} + \frac{1}{2\Sigma_t}\right)d\Sigma_t.$$

Integrating both sides of the above equation yields

$$t - \log|\Sigma^* - \Sigma_0| + \frac{1}{2}\log\Sigma_0 = -\log|\Sigma^* - \Sigma_t| + \frac{1}{2}\log\Sigma_t, \quad (\Sigma_t - \Sigma^*)^2 = \frac{(\Sigma_0 - \Sigma^*)^2}{\Sigma_0}e^{-2t}\Sigma_t.$$

Hence, the solution Σ_t follows

$$\Sigma_t = \Sigma^* + \frac{e^{-2t}(\Sigma_0 - \Sigma^*)^2}{2\Sigma_0} + (\Sigma_0 - \Sigma^*)e^{-t}\sqrt{\frac{e^{-2t}(\Sigma_0 - \Sigma^*)^2}{4\Sigma_0^2} + \frac{\Sigma^*}{\Sigma_0}}.$$

Now, we are ready to compare the NLD with OLD, LLD and HAMCMC. Here we consider $f(x) = (2\Sigma^*)^{-1}(x - \mu^*)^2$, where $\Sigma^* > 0$ and μ^* are given. The OLD satisfies

$$dX_t = -(\Sigma^*)^{-1}(X_t - \mu^*)dt + \sqrt{2}dB_t,$$

which is also known as the Ornstein-Uhlenbeck process. And LLD writes

$$dX_t = -(\Sigma^*)^{-1}(X_t - \mu^*)dt + \Sigma_t^{-1}(X_t - \mu_t)dt.$$

The mean μ_t and variance Σ_t of OLD and LLD both satisfy

$$\mu_t = \mu^* + e^{-(\Sigma^*)^{-1}t} (\mu_0 - \mu^*), \quad \Sigma_t = \Sigma^* + e^{-2(\Sigma^*)^{-1}t} (\Sigma_0 - \Sigma^*).$$

On the other hand, HAMCMC follows the dynamics

$$dX_t = -(X_t - \mu^*)dt + \sqrt{2\Sigma^*}dB_t.$$

For HAMCMC, the evolution of mean μ_t follows

$$d\mu_t = d\mathbb{E}[X_t] = -(\mu_t - \mu^*)dt,$$

and the evolution of variance Σ_t satisfies

$$d\Sigma_t = d(\mathbb{E}[X_t^2] - \mu_t^2) = 2\mathbb{E}[X_t dX_t] - 2\mu_t d\mu_t$$

=2\left[-\left(\Sigma_t + \mu_t^2\right) + \mu^*\mu_t + \Sigma^* + \mu_t(\mu_t - \mu^*)\right] dt = 2\left(\Sigma^* - \Sigma_t) dt.

The mean μ_t and variance Σ_t of HAMCMC follows

$$\mu_t = \mu^* + e^{-t}(\mu_0 - \mu^*), \quad \Sigma_t = \Sigma^* + e^{-2t}(\Sigma_0 - \Sigma^*).$$

We summarize our results in Table 4.

Dynamics	Particle	Mean and variance
NLD	$dX_t = \left(\frac{\Sigma^* - \Sigma_t}{\Sigma^* + \Sigma_t} X_t - \frac{2\Sigma^*}{\Sigma^* + \Sigma_t} \mu_t + \mu^*\right) dt$	$\mu_t = \mu^* + e^{-t}(\mu_0 - \mu^*)$ $\frac{\Sigma_t - \Sigma^*}{\Sigma_0 - \Sigma_t^*} = \frac{e^{-2t}(\Sigma_0 - \Sigma^*)}{2\Sigma_0}$
		$+e^{-t}\sqrt{\frac{e^{-2t}(\Sigma_0-\Sigma^*)^2}{4\Sigma_0^2}+\frac{\Sigma^*}{\Sigma_0}}$
OLD	$dX_t = -(\Sigma^*)^{-1}(X_t - \mu^*)dt + \sqrt{2}dB_t$	$\mu_t = \mu^* + e^{-(\Sigma^*)^{-1}t}(\mu_0 - \mu^*)$
LLD	$dX_t = -(\Sigma^*)^{-1}(X_t - \mu^*)dt + \Sigma_t^{-1}(X_t - \mu_t)dt$	$\Sigma_t = \Sigma^* + e^{-2(\Sigma^*)^{-1}t} (\Sigma_0 - \Sigma^*)$
HAMCMC	$dX_t = -(X_t - \mu^*)dt + \sqrt{2\Sigma^*}dB_t.$	$\mu_t = \mu^* + e^{-t}(\mu_0 - \mu^*)$ $\Sigma_t = \Sigma^* + e^{-2t}(\Sigma_0 - \Sigma^*)$

Table 4: Comparison among different Langevin dynamics on 1D Gaussian family.

Compared to OLD and LLD, the exponential convergence rate of μ_t and Σ_t in NLD does not depend on Σ^* . This fact shows that the NLD is the Newton's flow for both the evolution of mean and variance in Gaussian process. We also note that the convergence rates of mean and variance are different in HAMCMC, while they are same in NLD. In section 7, we use numerical examples to further demonstrate the differences between NLD and HAMCMC.

Appendix D. Connection with Stein variational Newton's method

The Stein variational Newton's method (SVN) is also a second-order method for sampling. It aims to minimize $J_{\rho}[\phi]$, which evaluates the change of $E(\rho)$ along the transformation map $\phi : \mathbb{R}^d \to \mathbb{R}^d$.

$$J_{\rho}[\boldsymbol{\phi}] = E((I + \boldsymbol{\phi}) \# \rho). \tag{43}$$

Here $(I + \phi) \# \rho$ denotes the pushforward density of ρ along the map $I(x) + \phi(x)$ and I(x) is the identity map. In each iteration, SVN solves $\phi \in S^d$ via the following equation:

$$D^2 J_{\rho}[0](\boldsymbol{\psi}, \boldsymbol{\phi}) = -D J_{\rho}[0](\boldsymbol{\psi}), \quad \boldsymbol{\psi} \in \mathcal{S}^d.$$
(44)

Here S is the RKHS related to a kernel function k(x, y) and $S^d = S \times \cdots \times S$. Besides, DJ_{ρ} and D^2J_{ρ} denote the first and second variation of J_{ρ} .

We note that the following relationships hold

$$DJ_{\rho}[0][\psi] = \int \psi^{T} (\nabla f + \nabla \log \rho) \rho dx.$$
$$D^{2}J_{\rho}[0](\psi, \phi) = \mathbb{E}_{x \sim \rho}[\phi(x)^{T} \nabla^{2} f(x) \psi(x) + \operatorname{tr}(\nabla \phi(x) \nabla \psi(x))].$$

If we restrict ψ and ϕ to be gradient vector fields. Namely, there exists $\Psi(x), \Phi(x) : \mathbb{R}^d \to \mathbb{R}$ such that $\psi(x) = \nabla \Psi(x)$ and $\phi(x) = \nabla \Phi(x)$. Then, we recover the gradient and Hessian operators in probability space with Wasserstein-2 metric.

$$DJ_{\rho}[0][\nabla\Psi] = \int (\langle \nabla\psi, \nabla f \rangle + \Delta\Psi)\rho dx = \int \Psi \operatorname{grad}^{W} E(\rho) dx.$$
$$D^{2}J_{\rho}[0](\nabla\Psi, \nabla\Phi) = \int \left(\langle \nabla^{2}\Psi, \nabla^{2}\Phi \rangle + \nabla\Psi^{T}\nabla^{2}f\nabla\Phi\right)\rho dx$$
$$= \int \Psi \mathcal{H}_{E}^{W}(\rho)\Phi dx.$$

On the other hand, the kernelized Wasserstein Newton's method in each step solves $\Phi \in S$ from (3). Because S is a Hilbert space, this is equivalent to find $\Phi \in S$ such that

$$\int \Psi \operatorname{Hess}^{W} E(\rho)[\Phi] dx = -\int \Psi \operatorname{grad}^{W} E(\rho) dx, \quad \forall \Psi \in \mathcal{S},$$

or equivalently,

$$D^2 J_{\rho}[0](\nabla \Phi, \nabla \Psi) = -D J_{\rho}[0](\nabla \Psi), \quad \forall \Psi \in \mathcal{S}.$$

This can be viewed as a restriction on (44). Namely, we solve $D^2 J_{\rho}[0](\psi, \phi) = -D J_{\rho}[0](\psi)$ in the space $\{\phi = \nabla \Phi | \Phi \in S\}$ instead of S^d .

Remark 37 We notice the differences between Wasserstein Newton and Stein variational Newton in formulations. SVN studies the second order variations w.r.t. transportation maps, while we focus on these variations w.r.t. densities. Besides, we benefit from the utilization of gradient and Hessian operators in probability space with Wasserstein-2 metric. This allows us to to prove the convergence rate of information Newton's method in the sense of density.

Appendix E. Proofs in Section 6

In this section, we provide convergence proofs of information Newton's method with approximated Newton's direction in section 6.

E.1 Riemannian structure of probability space

We first provide some background knowledge for the Riemannian structure of probability space. For simplicity, we define the exponential map and other Riemannian operators on cotangent space.

Definition 38 (Exponential map on cotangent space and its inverse) The exponential map $\operatorname{Exp}_{\rho_0}$ is a mapping from the cotangent space $T^*_{\rho_0}\mathcal{P}(\Omega)$ to $\mathcal{P}(\Omega)$. Namely, $\operatorname{Exp}_{\rho_0}(\Phi) = \hat{\rho}_s|_{s=1}$. Here $\hat{\rho}_s, s \in [0, 1]$ is the solution to geodesic equation (30) with initial conditions $\hat{\rho}_s|_{s=0} = \rho_0, \Phi_s|_{s=0} = \Phi$.

The inverse of the exponential map $\operatorname{Exp}_{\rho_0}(\rho_1)$ follows $\operatorname{Exp}_{\rho_0}^{-1}(\rho_1) = \mathcal{G}(\hat{\rho}_s)\partial_s \hat{\rho}_s|_{s=0}$. Here $\hat{\rho}_s, s \in [0, 1]$ is the solution to geodesic equation (30) with boundary conditions $\hat{\rho}_s|_{s=0} = \rho_0$ and $\hat{\rho}_s|_{s=1} = \rho_1$.

We also denote $\text{Exp}_{\rho}^{\alpha}(\Phi)$ to be the solution at time $t = \alpha$ to the geodesic equation (30) with initial values $\hat{\rho}_0 = \rho$ and $\Phi_0 = \Phi$. As a known result of Riemannian geometry, the geodesic curve has constant speed (Boothby, 1986). Namely, for $\Phi \in T_{\rho}^* \mathcal{P}(\Omega)$ and $\alpha > 0$, we have

$$\operatorname{Exp}_{\rho}^{\alpha}(\Phi) = \operatorname{Exp}_{\rho}(\alpha \Phi).$$

And for $\rho_0, \rho_1 \in \mathcal{P}(\Omega)$, it follows

$$\|\operatorname{Exp}_{\rho_0}^{-1}(\rho_1)\|_{\rho_0}^2 = \mathcal{D}(\rho_0, \rho_1)^2.$$

We define high-order derivatives on the cotangent-space in Proposition 39.

Proposition 39 For all $\Phi \in T^*_{\rho}\mathcal{P}(\Omega)$, it follows

$$E(\operatorname{Exp}^{s}_{\rho}(\Phi)) = E(\rho) + s\nabla E(\rho)(\Phi) + \dots \frac{s^{n-1}}{(n-1)!} \nabla^{n-1} E(\rho)(\Phi, \dots, \Phi) + \frac{s^{n}}{n!} \nabla^{n} E(\operatorname{Exp}_{\rho}(\lambda \Phi))(\tau_{\lambda} \Phi, \dots, \tau_{\lambda} \Phi),$$

where τ_{λ} is the parallelism from ρ to $\operatorname{Exp}_{\rho}^{\lambda}(\Phi)$ and $\lambda \in (0, s)$. Here $\nabla^{n} E(\rho)$ defines a n-form on the cotangent space $T_{\rho}^{*}\mathcal{P}(\Omega)$. Namely, it is recursively defined by

$$\nabla^{n} E(\rho)(\Phi_{1},\ldots,\Phi_{n}) = \left. \frac{\partial}{\partial s} \nabla^{n-1} E(\operatorname{Exp}_{\rho}(s\Phi_{n}))(\tau_{s}\Phi_{1},\ldots,\tau_{s}\Phi_{n-1}) \right|_{s=0},$$

where τ_s is the parallelism from ρ to $\operatorname{Exp}_{\rho}(s\Phi_n)$.

Proof We first show that

$$\frac{\partial}{\partial s} \nabla^{n-1} E(\operatorname{Exp}^{s}_{\rho}(\Phi_{n}))(\tau_{s}\Phi_{1},\ldots,\tau_{s}\Phi_{n-1}) = \nabla^{n} E(\operatorname{Exp}^{s}_{\rho}(\Phi_{n}))(\tau_{s}\Phi_{1},\ldots,\tau_{s}\Phi_{n}).$$
(45)

From the definition, it follows that

$$\begin{aligned} &\frac{\partial}{\partial s} \nabla^{n-1} E(\operatorname{Exp}_{\rho}^{s}(\Phi_{n}))(\tau_{s}\Phi_{1},\ldots,\tau_{s}\Phi_{n-1}) \\ &= \frac{\partial}{\partial t} \nabla^{n-1} E(\operatorname{Exp}_{\rho}^{s+t}(\Phi_{n}))(\tau_{s+t}\Phi_{1},\ldots,\tau_{s+t}\Phi_{n-1})\Big|_{t=0} \\ &= \frac{\partial}{\partial t} \nabla^{n-1} E(\operatorname{Exp}_{\operatorname{Exp}_{\rho}^{s}(\Phi_{n})}^{t}(\tau_{s}\Phi_{n}))(\tau_{t}\tau_{s}\Phi_{1},\ldots,\tau_{t}\tau_{s}\Phi_{n-1})\Big|_{t=0} \\ &= \nabla^{n} E(\operatorname{Exp}_{\rho}^{s}(\Phi_{n}))(\tau_{s}\Phi_{1},\ldots,\tau_{s}\Phi_{n}). \end{aligned}$$

From (45), we can recursively compute that

$$\frac{\partial^n}{(\partial s)^n} E(\operatorname{Exp}^s_{\rho}(\Phi)) = \nabla^n E(\operatorname{Exp}^s_{\rho}(\Phi))(\tau_s \Phi, \dots \tau_s \Phi).$$

The Taylor expansion of $E(\text{Exp}_{\rho}^{s}(\Phi))$ w.r.t. s completes the proof.

E.2 Cauchy-Schwarz inequality

To complete proofs in section 6, we introduce Lemma 40.

Lemma 40 (Cauchy-Schwarz inequality) Suppose that $\mathcal{H} : T^*_{\rho}\mathcal{P}(\Omega) \to T_{\rho}\mathcal{P}(\Omega)$ is a self-adjoint linear operator and \mathcal{H} is positive definite. Then, for $\Phi_1, \Phi_2 \in T^*_{\rho}\mathcal{P}(\Omega)$, we have

$$\left(\int \Phi_1 \mathcal{H} \Phi_2 dx\right)^2 \leq \left(\int \Phi_1 \mathcal{H} \Phi_1 dx\right) \left(\int \Phi_2 \mathcal{H} \Phi_2 dx\right).$$

Proof The proof is quite similar to the Euclidean space. For all $s \in \mathbb{R}$, we have

$$0 \leq \int (\Phi_1 + s\Phi_2) \mathcal{H}(\Phi_1 + s\Phi_2) dx$$
$$= s^2 \int \Phi_2 \mathcal{H}\Phi_2 dx + 2s \int \Phi_1 \mathcal{H}\Phi_2 dx + \int \Phi_1 \mathcal{H}\Phi_1 dx$$

Because the arbitrary choice of s, it follows that

$$\left(2\int\Phi_1\mathcal{H}_E(\rho)\Phi_2dx\right)^2 - 4\left(\int\Phi_1\mathcal{H}_E(\rho)\Phi_1dx\right)\left(\int\Phi_2\mathcal{H}_E(\rho)\Phi_2dx\right) \ge 0.$$

This completes the proof.

E.3 Proofs of Proposition 18 and Lemma 19

To prove Proposition 18, we introduce Lemma 41.

Lemma 41 For all $\Phi \in T^*_{\rho_k} \mathcal{P}(\Omega)$, it follows

$$\nabla E(\rho_k)(\Phi) + \nabla^2 E(\rho_k)(T_k, \Phi) = -\frac{1}{2} \nabla^3 E(\operatorname{Exp}_{\rho_k}^{\lambda})(\tau_{\lambda} T_k, \tau_{\lambda} T_k, \tau_{\lambda} \Phi),$$

where τ_{λ} is the parallelism from ρ_k to $\operatorname{Exp}_{\rho_k}^{\lambda}(T_k)$ and $\lambda \in (0,1)$.

Proof Consider an auxiliary function

$$A(s) = \nabla E(\operatorname{Exp}_{\rho_k}^s(T_k))(\tau_s \Phi).$$

Directly from the definition of high-order derivatives, it follows

$$\frac{\partial}{\partial s}A(s) = \nabla^2 E(\operatorname{Exp}_{\rho_k}^s(T_k))(\tau_s T_k, \tau_s \Phi),$$
$$\frac{\partial^2}{\partial s^2}A(s) = \nabla^3 E(\operatorname{Exp}_{\rho_k}^s(T_k))(\tau_s T_k, \tau_s T_k, \tau_s \Phi).$$

Hence, we can compute the Taylor expansion

$$\nabla E(\operatorname{Exp}_{\rho_k}^1(T_k))(\tau_1\Phi) = \nabla E(\rho_k)(\Phi) + \nabla^2 E(\rho_k)(T_k,\Phi) + \frac{1}{2}\nabla^3 E(\operatorname{Exp}_{\rho_k}^\lambda)(\tau_\lambda T_k,\tau_\lambda T_k,\tau_\lambda\Phi).$$

On the other hand, we notice that

$$\nabla E(\operatorname{Exp}^{1}_{\rho_{k}}(T_{k}))(\tau_{1}\Phi) = \nabla E(\rho^{*})(\tau_{1}\Phi) = \int \tau_{1}\Phi \mathcal{G}(\rho)^{-1} \frac{\delta E}{\delta \rho^{*}} dx = 0.$$

This completes the proof.

Based on Lemma 41, Note that $\Phi_k = -\mathcal{H}_E(\rho_k)^{-1}\mathcal{G}(\rho_k)^{-1}\frac{\delta E}{\delta\rho_k}$. Hence, it follows

$$\mathcal{H}_E(\rho_k)\tau^{-1}T_{k+1} = \mathcal{H}_E(\rho_k)T_k + \mathcal{G}(\rho_k)^{-1}\frac{\delta E}{\delta\rho_k} - \mathcal{H}_E(\rho_k)R_k.$$

For arbitrary $\Psi \in T^*_{\rho_k} \mathcal{P}(\Omega)$, we have

$$\nabla^{2} E(\rho_{k})(\Psi, \tau^{-1}T_{k+1})$$

$$= \int \Psi \mathcal{H}_{E}(\rho_{k})\tau^{-1}T_{k+1}dx$$

$$= \int \Psi(\mathcal{H}_{E}(\rho_{k})T_{k} + \mathcal{G}(\rho_{k})^{-1}\frac{\delta E}{\delta\rho_{k}} - \mathcal{H}_{E}(\rho_{k})R_{k})dx$$

$$= \nabla^{2} E(\rho_{k})(\Psi, T_{k}) + \nabla E(\rho_{k})(\Psi) - \nabla^{2} E(\rho_{k})(\Psi, R_{k})$$

$$= -\frac{1}{2}\nabla^{3} E(\operatorname{Exp}_{\rho_{k}}^{\lambda})(\tau_{\lambda}\Psi, \tau_{\lambda}T_{k}, \tau_{\lambda}T_{k}) - \nabla^{2} E(\rho_{k})(\Psi, R_{k}).$$
(46)

Here the last equality comes from Lemma 41. Based on the definition of parallelism, we notice the fact

$$\| au_{\lambda}\Psi\|_{\operatorname{Exp}_{\rho_{k}}^{\lambda}(\Phi_{k})} = \|\Psi\|_{\rho_{k}}, \quad \forall \Psi \in T^{*}_{\rho_{k}}\mathcal{P}(\Omega).$$

Taking $\Psi = \tau^{-1}T_{k+1}$ in (46), applying Assumption 1 and utilizing Lemma 40 yields

$$\begin{split} \delta_{1} \| \tau^{-1} T_{k+1} \|_{\rho_{k}}^{2} &\leq \left| \nabla^{2} E(\rho_{k}) (\tau^{-1} T_{k+1}, \tau^{-1} T_{k+1}) \right| \\ &\leq \frac{1}{2} \left| \nabla^{3} E(\operatorname{Exp}_{\rho_{k}}^{\lambda}) (\tau_{\lambda} \tau^{-1} T_{k+1}, \tau_{\lambda} T_{k}, \tau_{\lambda} T_{k}) \right| + \left| \nabla^{2} E(\rho_{k}) (\tau^{-1} T_{k+1}, R_{k}) \right| \\ &\leq \frac{1}{2} \left| \nabla^{3} E(\operatorname{Exp}_{\rho_{k}}^{\lambda}) (\tau_{\lambda} \tau^{-1} T_{k+1}, \tau_{\lambda} T_{k}, \tau_{\lambda} T_{k}) \right| \\ &+ \sqrt{\left| \nabla^{2} E(\rho_{k}) (R_{k}, R_{k}) \right| \left| \nabla^{2} E(\rho_{k}) (\tau_{\lambda} \tau^{-1} T_{k+1}, \tau_{\lambda} \tau^{-1} T_{k+1}) \right|} \\ &\leq \delta_{3} \| \tau_{\lambda} T_{k} \|_{\operatorname{Exp}_{\rho_{k}}^{\lambda}(\Phi_{k})}^{2} \| \tau_{\lambda} \tau^{-1} T_{k+1} \|_{\operatorname{Exp}_{\rho_{k}}^{\lambda}(\Phi_{k})} + \delta_{2} \| \tau^{-1} T_{k+1} \|_{\rho_{k}} \| R_{k} \|_{\rho_{k}} \\ &= \delta_{3} \| T_{k} \|_{\rho_{k}}^{2} \| \tau^{-1} T_{k+1} \|_{\rho_{k}} + \delta_{2} \| \tau^{-1} T_{k+1} \|_{\rho_{k}} \| R_{k} \|_{\rho_{k}}. \end{split}$$

Hence, it follows

$$||T_{k+1}||_{\rho_{k+1}} = ||\tau^{-1}T_{k+1}||_{\rho_k} \le \frac{\delta_3}{\delta_1} ||T_k||_{\rho_k}^2 + \frac{\delta_2}{\delta_1} ||R_k||_{\rho_k}.$$

To prove Lemma 19, we introduce the following Lemma 42.

Lemma 42 We have following estimations

$$\|\Phi_k\|_{\rho_k} = O(\|T_k\|_{\rho_k}), \quad \|T_{k+1}\|_{\rho_{k+1}} = O(\|T_k\|_{\rho_k}).$$

Proof From Assumption 1 and Cauchy-Swarz inequality, it follows that

$$\begin{split} \|\Phi_k\|_{\rho_k}^2 &= \int \Phi_k \mathcal{G}(\rho_k)^{-1} \Phi_k dx \le \delta_1^{-1} \int \Phi_k \mathcal{H}_E(\rho_k) \Phi_k dx \\ &= \delta_1^{-1} \int \Phi_k \mathcal{G}(\rho_k)^{-1} \frac{\delta E}{\delta \rho_k} dx \le \delta_1^{-1} \|\Phi_k\|_{\rho_k} \left\| \frac{\delta E}{\delta \rho_k} \right\|_{\rho_k} \end{split}$$

We also notice that from Lemma 41,

$$\begin{split} & \left\| \frac{\delta E}{\delta \rho_k} \right\|_{\rho_k}^2 = \nabla E(\rho_k) \left(\frac{\delta E}{\delta \rho_k} \right) \\ = \nabla^2 E(\rho_k) \left(T_k, \frac{\delta E}{\delta \rho_k} \right) + O\left(\|T_k\|_{\rho_k}^2 \left\| \frac{\delta E}{\delta \rho_k} \right\|_{\rho_k} \right) \\ = O\left(\|T_k\|_{\rho_k} \left\| \frac{\delta E}{\delta \rho_k} \right\|_{\rho_k} \right). \end{split}$$

As a result, we have $\|\Phi_k\|_{\rho_k} = O\left(\left\|\frac{\delta E}{\delta \rho_k}\right\|_{\rho_k}\right) = O\left(\|T_k\|_{\rho_k}\right)$. We also note the triangle inequality

$$|||T_k||_{\rho_k} - ||\Phi_k||_{\rho_k}| \le ||T_{k+1}||_{\rho_{k+1}} \le ||T_k||_{\rho_k} + ||\Phi_k||_{\rho_k}.$$

This yields $||T_{k+1}||_{\rho_{k+1}} = O(||T_k||_{\rho_k}).$

We finally show the estimation of $||R_k||_{\rho_k}$. Based on the first-order approximation of the exponential map and the parallelsim, we have the following estimations

$$\int \Psi(\rho^* - \rho_k) dx = \int \Psi \mathcal{G}(\rho_k)^{-1} T_k dx + O(\|\Psi\|_{\rho_k} \|T_k\|_{\rho_k}^2),$$

$$\int \Psi(\rho_{k+1} - \rho_k) dx = \int \Psi \mathcal{G}(\rho_k)^{-1} \Phi_k dx + O(\|\Psi\|_{\rho_k} \|\Phi_k\|_{\rho_k}^2)$$
$$= \int \Psi \mathcal{G}(\rho_k)^{-1} \Phi_k dx + O(\|\Psi\|_{\rho_k} \|T_k\|_{\rho_k}^2),$$

and

$$\int \Psi(\rho^* - \rho_{k+1}) dx = \int \Psi \mathcal{G}(\rho_{k+1})^{-1} T_{k+1} dx + O(\|\Psi\|_{\rho_{k+1}} \|T_{k+1}\|_{\rho_{k+1}}^2)$$

$$= \int \tau^{-1} \Psi \mathcal{G}(\rho_k)^{-1} \tau^{-1} T_{k+1} dx + O(\|\Psi\|_{\rho_k} \|T_{k+1}\|_{\rho_{k+1}}^2 + \|\Psi - \tau^{-1}\Psi\|_{\rho_k} \|T_{k+1}\|_{\rho_{k+1}}^2)$$

$$= \int \tau^{-1} \Psi \mathcal{G}(\rho_k)^{-1} \tau^{-1} T_{k+1} dx + O(\|\Psi\|_{\rho_k} \|T_{k+1}\|_{\rho_{k+1}}^2 + \|\Psi\|_{\rho_k} \|\Phi_k\|_{\rho_k} \|T_{k+1}\|_{\rho_{k+1}}^2)$$

$$= \int \Psi \mathcal{G}(\rho_k)^{-1} \tau^{-1} T_{k+1} dx + O(\|\Psi\|_{\rho_k} \|T_{k+1}\|_{\rho_{k+1}}^2 + \|\Psi - \tau^{-1}\Psi\|_{\rho_k} \|\tau^{-1} T_{k+1}\|_{\rho_k})$$

$$= \int \Psi \mathcal{G}(\rho_k)^{-1} \tau^{-1} T_{k+1} dx + O(\|\Psi\|_{\rho_k} \|T_{k+1}\|_{\rho_{k+1}}^2 + \|\Psi\|_{\rho_k} \|\Phi_k\|_{\rho_k} \|T_{k+1}\|_{\rho_{k+1}})$$

$$= \int \Psi \mathcal{G}(\rho_k)^{-1} \tau^{-1} T_{k+1} dx + O(\|\Psi\|_{\rho_k} \|T_{k+1}\|_{\rho_{k+1}}^2).$$

Furthermore, we have $R_k = T_k - \tau^{-1}T_{k+1} - \Phi_k$ and

$$\int \Psi(\rho^* - \rho_k) dx - \int \Psi(\rho^* - \rho_{k+1}) dx - \int \Psi(\rho_{k+1} - \rho_k) dx = 0.$$

This completes the proof.

E.4 Proof of Theorem 20

We first notice that

$$\nabla E(\rho)(\Phi) = \int \Phi \mathcal{G}(\rho)^{-1} \frac{\delta E}{\delta \rho} dx, \quad \nabla^2 E(\rho)(\Phi_1, \Phi_1) = \int \Phi \mathcal{H}_E(\rho) \Phi dx.$$
(47)

By taking (47) into Lemma 19 and utilizing (A3), we note that for $\sigma \in T^*_{\rho_k}\mathcal{P}(\Omega)$,

$$\int g_k \sigma dx = -\int \mathcal{H}_E(\rho_k) T_k \sigma dx + \mathcal{O}(\|\sigma\|_{\rho_k} \|T_k\|_{\rho_k}^2).$$
(48)

Based on the Taylor expansion on the Riemannian manifold with (A3), it follows

$$E(\rho_{k+1}) = E(\rho_k) + \alpha_k \int \Phi_k \mathcal{G}(\rho_k)^{-1} \frac{\delta E}{\delta \rho_k} dx + \frac{\alpha_k^2}{2} \int \Phi_k \mathcal{H}_E(\rho_k) \Phi_k dx + \mathcal{O}(\|\Phi_k\|_{\rho_k}^3).$$

Following (27) and (A5), this yields

$$E(\rho_{k+1}) - E(\rho_{k})$$

$$= -\alpha_{k} \int g_{k} \mathcal{H}_{E,P} g_{k} dx + \frac{\alpha_{k}^{2}}{2} \int g_{k} \mathcal{H}_{E,P} \mathcal{H}_{E}(\rho_{k}) \mathcal{H}_{E,P} g_{k} dx + \mathcal{O}(\|\Phi_{k}\|_{\rho_{k}}^{3})$$

$$= \frac{\alpha_{k}^{2} - 2\alpha_{k}}{2} \int g_{k} \mathcal{H}_{E,P} g_{k} dx + \frac{\alpha_{k}^{2}}{2} \int g_{k} (\mathcal{H}_{E,P} \mathcal{H}_{E}(\rho_{k}) \mathcal{H}_{E,P} - \mathcal{H}_{E,P}) g_{k} dx + \mathcal{O}(\|\Phi_{k}\|_{\rho_{k}}^{3})$$

$$\leq \frac{\alpha_{k}^{2} - 2\alpha_{k}}{2} \int g_{k} \mathcal{H}_{E,P} g_{k} dx + \frac{\epsilon_{2} \alpha_{k}^{2}}{2} \int g_{k} \mathcal{H}_{E,P}(\rho_{k}) g_{k} dx + \mathcal{O}(\|\Phi_{k}\|_{\rho_{k}}^{3}).$$

$$(49)$$

Similarly, by the Taylor expansion along with (A3), we have

$$E(\rho^{*}) - E(\rho_{k}) = \int g_{k}T_{k}dx + \frac{1}{2}\int T_{k}\mathcal{H}_{E}(\rho_{k})T_{k}dx + \mathcal{O}(||T_{k}||_{\rho_{k}}^{3})$$

$$= -\frac{1}{2}\int T_{k}\mathcal{H}_{E}(\rho_{k})T_{k}dx + \mathcal{O}(||T_{k}||_{\rho_{k}}^{3}).$$
(50)

According to (A1), (A2) and Cauchy-Schwartz inequality, we have

$$\begin{aligned} \|\mathcal{H}_{E}(\rho_{k})^{-1}g_{k}\|_{\rho_{k}}^{2} &= \int \mathcal{H}_{E}(\rho_{k})^{-1}g_{k}\mathcal{G}(\rho_{k})^{-1}\mathcal{H}_{E}(\rho_{k})^{-1}g_{k}dx\\ &\leq \delta_{1}^{-1}\int \mathcal{H}_{E}(\rho_{k})^{-1}g_{k}\mathcal{H}_{E}(\rho_{k})\mathcal{H}_{E}(\rho_{k})^{-1}g_{k}dx\\ &= \delta_{1}^{-1}\int g_{k}\mathcal{H}_{E}(\rho_{k})^{-1}\mathcal{G}(\rho_{k})^{-1}\mathcal{G}(\rho_{k})g_{k}dx\\ &\leq \delta_{1}^{-1}\|\mathcal{H}_{E}(\rho_{k})^{-1}g_{k}\|_{\rho_{k}}\|\mathcal{G}(\rho_{k})g_{k}\|_{\rho_{k}}.\end{aligned}$$

Besides, from the proof of Lemma 42, we have

$$\|\mathcal{G}(\rho_k)g_k\|_{\rho_k} = \left\|\frac{\delta E}{\delta\rho_k}\right\|_{\rho_k} = O(\|T_k\|_{\rho_k}).$$

This tells $\|\mathcal{H}_E(\rho_k)^{-1}g_k\|_{\rho_k} = O(\|\mathcal{G}(\rho_k)g_k\|_{\rho_k}) = O(\|T_k\|_{\rho_k})$. Hence, by utilizing (48) two times, we have

$$\int g_k \mathcal{H}_E(\rho_k)^{-1} g_k dx$$

= $-\int \mathcal{H}_E(\rho_k)^{-1} T_k \mathcal{H}_E(\rho_k) g_k dx + \mathcal{O}(\|T_k\|_{\rho_k}^2 \|\mathcal{H}_E(\rho_k)^{-1} g_k\|_{\rho_k})$
= $-\int T_k g_k dx + \mathcal{O}(\|T_k\|_{\rho_k}^3)$
= $\int T_k \mathcal{H}_E(\rho_k) T_k dx + \mathcal{O}(\|T_k\|_{\rho_k}^3).$

This indicates

$$E(\rho^*) - E(\rho_k) = -\frac{1}{2} \int T_k \mathcal{H}_E(\rho_k) T_k dx + \mathcal{O}(||T_k||^3_{\rho_k})$$

= $-\frac{1}{2} \int g_k \mathcal{H}_E(\rho_k)^{-1} g_k dx + \mathcal{O}(||T_k||^3_{\rho_k}).$ (51)

Following (A6), we note that

$$\|\Phi_k\|_{\rho_k} = \left\|\mathcal{H}_{E,P}\mathcal{G}(\rho_k)^{-1}\frac{\delta E}{\delta\rho_k}\right\|_{\rho_k} = \mathcal{O}\left(\left\|\frac{\delta E}{\delta\rho_k}\right\|_{\rho_k}\right) = \mathcal{O}(\|T_k\|_{\rho_k}).$$

In summary, combining (A4), (49) and (51), we have

$$\begin{split} E(\rho_{k+1}) &- E(\rho^{*}) \\ \leq & E(\rho_{k}) - E(\rho^{*}) + \frac{\alpha_{k}^{2} - 2\alpha_{k}}{2} \int g_{k} \mathcal{H}_{E,P} g_{k} dx \\ &+ \frac{\epsilon_{1} \alpha_{k}^{2}}{2} \int g_{k} \mathcal{H}_{E,P}(\rho_{k}) g_{k} dx + \mathcal{O}(\|\Phi_{k}\|_{\rho_{k}}^{3}) \\ \leq & \frac{1}{2} \int g_{k} \mathcal{H}_{E}(\rho_{k})^{-1} g_{k} dx + \frac{\alpha_{k}^{2} - 2\alpha_{k}}{2} \int g_{k} \mathcal{H}_{E}(\rho_{k}) g_{k} dx \\ &+ \frac{|\alpha_{k}^{2} - 2\alpha_{k}|\epsilon_{1}}{2} \int g_{k} \mathcal{H}_{E}(\rho_{k})^{-1} g_{k} dx + \frac{\epsilon_{2}(1 + \epsilon_{1})\alpha_{k}^{2}}{2} \int g_{k} \mathcal{H}_{E}(\rho_{k})^{-1} g_{k} dx + \mathcal{O}(\|T_{k}\|_{\rho_{k}}^{3}) \\ &= \left(\frac{(\alpha_{k} - 1)^{2}}{2} + \frac{|\alpha_{k}^{2} - 2\alpha_{k}|\epsilon_{1}}{2} + \frac{\epsilon_{2}(1 + \epsilon_{1})\alpha_{k}^{2}}{2}\right) \int g_{k} \mathcal{H}_{E}(\rho_{k})^{-1} g_{k} dx + \mathcal{O}(\|T_{k}\|_{\rho_{k}}^{3}). \end{split}$$

By taking $\alpha_k = 1$ and utilizing (51), we have

$$E(\rho_{k+1}) - E(\rho^*) \leq \frac{\epsilon_1 + \epsilon_2 + \epsilon_1 \epsilon_2}{2} \int g_k \mathcal{H}_E(\rho_k)^{-1} g_k dx + \mathcal{O}(\|T_k\|_{\rho_k}^3)$$

= $(\epsilon_1 + \epsilon_2 + \epsilon_1 \epsilon_2) (E(\rho_k) - E(\rho^*)) + \mathcal{O}((E(\rho_k) - E(\rho^*))^{3/2}).$

The last equality comes from $||T_k||_{\rho_k}^2 = O\left(\int T_k \mathcal{H}_E(\rho_k) T_k dx\right) = O(E(\rho_k) - E(\rho^*)).$

E.5 Proof of Theorem 22

For simplicity, denote $p_k = \hat{\Phi}_k - \Phi_k$. From the previous derivation, with $\alpha_k = 1$, we note that

$$\begin{split} E(\rho_{k+1}) - E(\rho_k) \\ &= -\int g_k (\mathcal{H}_{E,P}g_k + p_k) dx + \frac{1}{2} \int (p_k + \mathcal{H}_{E,P}g_k) \mathcal{H}_E(\rho_k) (\mathcal{H}_{E,P}g_k + p_k) dx + \mathcal{O}(\|\hat{\Phi}_k\|_{\rho_k}^3) \\ &= \frac{1}{2} \int g_k \mathcal{H}_{E,P}g_k dx + \frac{1}{2} \int g_k (\mathcal{H}_{E,P}\mathcal{H}_E(\rho_k)\mathcal{H}_{E,P} - \mathcal{H}_{E,P})g_k dx \\ &- \int \left(g_k p_k - \frac{1}{2} p_k (\mathcal{H}_E(\rho_k)\mathcal{H}_{E,P} + \mathcal{H}_{E,P}\mathcal{H}_E(\rho_k))g_k \right) dx \\ &+ \frac{1}{2} \int p_k \mathcal{H}_E(\rho_k)p_k dx + \mathcal{O}(\|\hat{\Phi}_k\|_{\rho_k}^3) \\ &\leq \frac{1}{2} \int g_k \mathcal{H}_{E,P}g_k dx + \frac{\epsilon_2}{2} \int g_k \mathcal{H}_{E,P}(\rho_k)g_k dx + \frac{\epsilon_3 + \epsilon_4}{2} \int g_k \mathcal{H}_E^{-1}(\rho_k)g_k dx + \mathcal{O}(\|\hat{\Phi}_k\|_{\rho_k}^3). \end{split}$$

The last inequality further utilizes (A7) and (A8). We also note that

$$\|\hat{\Phi}_k\|_{\rho_k} \le \|p_k\|_{\rho_k} + \|\Phi_k\|_{\rho_k}.$$

And we have

$$\begin{aligned} \|p_k\|_{\rho_k}^2 &= \int p_k \mathcal{G}(\rho_k)^{-1} p_k dx \le \frac{1}{\delta_1} \int p_k \mathcal{H}_E(\rho_k) p_k dx \\ &\le \frac{\epsilon_4}{\delta_1} \int g_k \mathcal{H}_E(\rho_k)^{-1} g_k dx = \mathcal{O}(\|T_k\|_{\rho_k}^2). \end{aligned}$$

Hence, $\|\hat{\Phi}_k\|_{\rho_k} = \mathcal{OO}(\|T_k\|_{\rho_k})$. As a result, by utilizing (51), we complete the proof.

E.6 Justification of Assumption 3

To justify Assumption 3, we first introduce some definitions.

For an energy function $E(\rho)$, we call it *well-defined w.r.t. samples* if $E(\hat{\rho})$ is well-defined for $\hat{\rho} = \frac{1}{N} \sum_{i=1}^{N} \delta(x - x_i)$, where δ is the Dirac-delta distribution. We denote

$$\hat{\mathcal{P}}(\Omega) = \mathcal{P}(\Omega) \cup \left\{ \hat{\rho} = \frac{1}{N} \sum_{i=1}^{N} \delta(x - x_i) | x_i \sim \rho, \rho \in \mathcal{P}(\Omega) \right\}.$$

Remark 43 Typical examples of such energy functions include

$$E(\rho) = \int f(x)\rho(x)dx,$$

where f(x) is a smooth function. Or

$$E(\rho) = \int f(x;\rho)\rho(x)dx.$$

Here $f(x; \rho)$ is well-defined w.r.t. samples for fixed x. For instance, $f(x; \rho) = \int w(x, y)\rho(y)dy$ for some smooth function w(x, y).

We say that $\{\hat{\rho}_n\} \subseteq \hat{\mathcal{P}}(\Omega)$ weakly converges to $\rho \in \mathcal{P}(\Omega)$ if for any smooth (test) function f,

$$\lim_{N \to \infty} \int f(x) \hat{\rho}_N(x) dx = \int f(x) \rho(x) dx.$$

We say that $E(\rho)$ is convergent w.r.t. samples if $E(\rho)$ is well-defined w.r.t. samples and

$$\lim_{n \to \infty} E(\hat{\rho}_n) = E(\rho),$$

for any $\{\hat{\rho}_n\} \subseteq \hat{\mathcal{P}}(\Omega)$ weakly converges to $\rho \in \mathcal{P}(\Omega)$.

For $\rho \in \hat{\mathcal{P}}(\Omega)$, we define the variational problem

$$J(\rho, \Phi) = \int \Phi \mathcal{H}_E(\rho) \Phi dx + 2 \int \Phi \mathcal{G}(\rho)^{-1} \frac{\delta E}{\delta \rho} dx + \lambda \int \Phi \mathcal{R}_S \Phi dx.$$

Suppose that $\|\Phi\|_{\mathcal{S}}$ is a norm in \mathcal{S} , which is independent of ρ . We further assumes that $\|\Phi\|_{\mathcal{S}}$ and the regularization term $\int \Phi \mathcal{R}_{\mathcal{S}} \Phi dx$ satisfy Assumption 4.

Assumption 4 There exists $\delta_5, \delta_6 > 0$ such that for all $\mathcal{D}(\rho, \rho^*) < \zeta$,

$$\delta_6 \|\Phi\|_{\mathcal{S}}^2 \le \|\Phi\|_{\rho}^2 \le \delta_5 \|\Phi\|_{\mathcal{S}}^2. \tag{A9}$$

There exists $\delta_7 \geq 0$, such that

$$\int \Phi \mathcal{R}_{\mathcal{S}} \Phi dx \le \delta_7 \|\Phi\|_{\mathcal{S}}^2.$$
(A10)

Suppose that for fixed $\Phi \in S$, $J(\rho, \Phi)$ is convergent w.r.t. samples. Then, for fixed $\rho \in \hat{\mathcal{P}}(\Omega)$, $J(\rho, \Phi)$ is well-defined and we denote $\Phi(\rho)$ as the minimizer of $\min_{\Phi \in S} J(\rho, \Phi)$. Then, $\Phi(\rho)$ is well-defined w.r.t. samples. We then show that $\Phi(\rho)$ is convergent w.r.t. samples.

For $\rho \in \hat{\mathcal{P}}(\Omega)$ satisfying (A1), we note that

$$J(\rho, \Phi) \ge \delta_1 \|\Phi\|_{\rho}^2 + 2 \int \Phi \mathcal{G}(\rho)^{-1} \frac{\delta E}{\delta \rho} dx + \lambda \int \Phi \mathcal{R}_{\mathcal{S}} \Phi dx$$

As a result, for fixed ρ , $J(\rho, \Phi)$ is δ_1 -strictly convex in Φ w.r.t. the norm $\|\cdot\|_{\rho}$, i.e.,

$$J(\rho, \Phi_1) - J(\rho, \Phi_2) \ge \int (\Phi_1 - \Phi_2) \left. \frac{\delta J(\rho, \Phi)}{\delta \Phi} \right|_{\Phi = \Phi_2} dx + \delta_1 \|\Phi_1 - \Phi_2\|_{\rho}^2.$$
(52)

Similarly, for $\rho \in \hat{\mathcal{P}}(\Omega)$ satisfying (A2), we note that

$$J(\rho, \Phi) \le \delta_2 \|\Phi\|_{\rho}^2 + 2 \int \Phi \mathcal{G}(\rho)^{-1} \frac{\delta E}{\delta \rho} dx + \lambda \delta_7 \|\Phi\|_{\mathcal{S}}^2.$$

Hence, this yields

$$J(\rho, \Phi_1) - J(\rho, \Phi_2) \le \int (\Phi_1 - \Phi_2) \left. \frac{\delta J(\rho, \Phi)}{\delta \Phi} \right|_{\Phi = \Phi_2} dx + \delta_2 \|\Phi_1 - \Phi_2\|_{\rho}^2 + \lambda \delta_7 \|\Phi_1 - \Phi_2\|_{\mathcal{S}}^2.$$
(53)

Lemma 44 Suppose that $S \subseteq \mathcal{F}(\Omega)/\mathbb{R}$ is a Hilbert space. $J(\Phi)$ is strictly convex in Φ w.r.t. some norm. For a variational problem $\min_{\Phi \in S} J(\Phi)$, the unique minimizer Φ^* satisfies

$$\int (\Psi - \Phi^*) \left. \frac{\delta J}{\delta \Phi} \right|_{\Phi = \Phi^*} dx = 0, \quad \forall \Psi \in \mathcal{S}.$$

Proof The variational problem $\min_{\Phi \in \mathcal{S}} J(\Phi)$ is equivalent to

$$\min_{\Phi \in \mathcal{F}(\Omega)/\mathbb{R}, \Psi \in \mathcal{S}} J(\Phi), \quad \text{s.t.} \quad \Phi = \Psi.$$

Consider the Lagrangian $\mathcal{L}(\Phi, \Psi, \lambda) = J(\Phi) + \int \lambda(\Phi - \Psi) dx$. The KKT conditions include:

$$\frac{\delta J}{\delta \Phi} + \lambda = 0, \quad \Psi = \Phi, \quad \int \lambda \tilde{\Psi} dx = 0, \quad \forall \tilde{\Psi} \in \mathcal{S}.$$

Here the equality holds up to a spatial-shift. As a result, for the minimizer Φ^* , we have

$$\int \left. \frac{\delta J}{\delta \Phi} \right|_{\Phi = \Phi^*} \tilde{\Psi} dx = 0, \quad \forall \tilde{\Psi} \in \mathcal{S}.$$

Because $\Psi - \Phi^* \in \mathcal{S}$, this completes the proof.

Proposition 45 S is a Hilbert space. Suppose that (A1) and (A2) in Assumption 1 further holds for $\rho \in \hat{\mathcal{P}}(\Omega)$. We assume the following statements hold.

- For $\rho \in \hat{\mathcal{P}}(\Omega)$, $\|\Phi(\rho)\|_{\mathcal{S}}$ is bounded.
- For fixed $\Phi \in S$, $J(\rho; \Phi)$ is convergent w.r.t. samples.
- For fixed $\Phi \in S$, $\|\Phi\|_{\rho}^2$ is well-defined w.r.t. samples.

Then, under Assumption 4, $\Phi(\rho)$ is convergent w.r.t. samples.

Proof Suppose that $\Phi(\rho)$ is not convergent w.r.t. samples. Then, there exists $\{\hat{\rho}_n\}_{n=1}^{\infty} \subseteq$ $\hat{\mathcal{P}}(\Omega)$ and $\epsilon > 0$ such that $\hat{\rho}_n$ weakly converges to $\rho \in \mathcal{P}(\Omega)$, while $\|\Phi(\hat{\rho}_n) - \Phi(\rho)\|_{\rho} > \epsilon$. We note that

$$J(\hat{\rho}_n, \Phi(\hat{\rho}_n)) - J(\rho, \Phi(\rho)) = J(\hat{\rho}_n, \Phi(\hat{\rho}_n)) - J(\rho, \Phi(\hat{\rho}_n)) + J(\rho, \Phi(\hat{\rho}_n)) - J(\rho, \Phi(\rho)) = J(\hat{\rho}_n, \Phi(\hat{\rho}_n)) - J(\hat{\rho}_n, \Phi(\rho)) + J(\hat{\rho}_n, \Phi(\rho)) - J(\rho, \Phi(\rho)).$$

Because $\Phi(\hat{\rho}_n)$ is the minimizer of $J(\hat{\rho}_n, \Phi)$, by applying (52) and Lemma 44, we have

$$J(\hat{\rho}_{n}, \Phi(\hat{\rho}_{n})) - J(\hat{\rho}_{n}, \Phi(\rho)) \leq -\delta_{1} \|\Phi(\hat{\rho}_{n}) - \Phi(\rho)\|_{\hat{\rho}_{n}}^{2}$$

$$\leq -\delta_{1}\delta_{6} \|\Phi(\hat{\rho}_{n}) - \Phi(\rho)\|_{\mathcal{S}}^{2} \leq -\frac{\delta_{1}\delta_{6}}{\delta_{5}} \|\Phi(\hat{\rho}_{n}) - \Phi(\rho)\|_{\rho}^{2} \leq -\frac{\delta_{1}\delta_{6}\epsilon^{2}}{\delta_{5}}.$$

Similarly, because $\Phi(\rho)$ is the minimizer of $J(\rho, \Phi)$, we have

- ())

$$J(\hat{\rho}_n, \Phi(\rho)) - J(\rho, \Phi(\rho)) \ge \delta_1 \|\Phi(\hat{\rho}_n) - \Phi(\rho)\|_{\rho}^2 \ge \frac{\delta_1 \epsilon^2}{2}$$

Because S is a Hilbert space and $\{\Phi(\hat{\rho}_n)\}$ is bounded, according to the Banach-Alaoglu theorem, $\{\Phi(\hat{\rho}_n)\}$ is weakly sequentially compact. Namely, there exists a weakly convergent subsequent $\{\Phi(\hat{\rho}_{n_k})\}$ (which is also convergent because S is a Hilbert space). Suppose that this sequence converges to Φ^* . As a result,

$$\lim_{k \to \infty} J(\rho, \Phi(\hat{\rho}_{n_k})) = J(\rho, \Phi^*).$$

From (53) and Assumption 4, we have

_ / .

$$J(\hat{\rho}_{n_k}, \Phi(\hat{\rho}_{n_k})) - J(\hat{\rho}_{n_k}, \Phi^*) \ge -\delta_2 \|\Phi(\hat{\rho}_{n_k}) - \Phi^*\|^2_{\hat{\rho}_{n_k}} - \lambda\delta_7 \|\Phi(\hat{\rho}_{n_k}) - \Phi^*\|^2_{\hat{\rho}_{n_k}} \\\ge -(\delta_2\delta_5 + \lambda\delta_7) \|\Phi(\hat{\rho}_{n_k}) - \Phi^*\|^2_{\mathcal{S}}.$$

Hence, we have

$$\lim_{k \to \infty} J(\hat{\rho}_{n_k}, \Phi(\hat{\rho}_{n_k})) = J(\rho, \Phi^*).$$

On the other hand, because $J(\rho, \Phi)$ is convergent w.r.t. samples for fixed Φ , $\lim_{k\to\infty} J(\hat{\rho}_{n_k}, \Phi(\rho)) J(\rho, \Phi(\rho)) = 0$. Hence, for sufficiently large k, we have

$$J(\hat{\rho}_{n_{k}}, \Phi(\hat{\rho}_{n_{k}})) - J(\rho, \Phi(\rho)) = J(\hat{\rho}_{n_{k}}, \Phi(\hat{\rho}_{n_{k}})) - J(\rho, \Phi(\hat{\rho}_{n_{k}})) + J(\rho, \Phi(\hat{\rho}_{n_{k}})) - J(\rho, \Phi(\rho)) \le -\frac{\delta_{1}\delta_{6}\epsilon^{2}}{2\delta_{5}},$$

and

$$J(\hat{\rho}_{n_{k}}, \Phi(\hat{\rho}_{n_{k}})) - J(\rho, \Phi(\rho)) = J(\hat{\rho}_{n_{k}}, \Phi(\hat{\rho}_{n_{k}})) - J(\hat{\rho}_{n_{k}}, \Phi(\rho)) + J(\hat{\rho}_{n_{k}}, \Phi(\rho)) - J(\rho, \Phi(\rho)) \ge \frac{\delta_{1}}{2}\epsilon^{2}.$$

This leads to a contradiction.

I(ô.

References

- Shun-Ichi Amari. Natural gradient works efficiently in learning. *Neural computation*, 10(2): 251–276, 1998.
- Shun-ichi Amari. Information geometry and its applications, volume 194. Springer, 2016.
- Michael Arbel, Arthur Gretton, Wuchen Li, and Guido Montúfar. Kernelized wasserstein natural gradient. arXiv preprint arXiv:1910.09652, 2019.
- Dominique Bakry and Michel Émery. Diffusions hypercontractives. In Séminaire de Probabilités XIX 1983/84, pages 177–206. Springer, 1985.
- Alain Berlinet and Christine Thomas-Agnan. Reproducing kernel Hilbert spaces in probability and statistics. Springer Science & Business Media, 2011.
- Espen Bernton. Langevin monte carlo and jko splitting. In Conference On Learning Theory, pages 1777–1798, 2018.
- William M Boothby. An introduction to differentiable manifolds and Riemannian geometry, volume 120. Academic press, 1986.
- J. A. Carrillo, S. Lisini, G. Savare, and D. Slepcev. Nonlinear mobility continuity equations and generalized displacement convexity. J. Funct. Anal., 258(4):1273-1309, 2010. ISSN 0022-1236. doi: 10.1016/j.jfa.2009.10.016. URL https://doi-org.stanford.idm.oclc. org/10.1016/j.jfa.2009.10.016.
- Peng Chen, Keyi Wu, Joshua Chen, Tom O'Leary-Roseberry, and Omar Ghattas. Projected stein variational newton: A fast and scalable bayesian inference method in high dimensions. In Advances in Neural Information Processing Systems, pages 15104–15113, 2019.
- Gianluca Detommaso, Tiangang Cui, Youssef Marzouk, Alessio Spantini, and Robert Scheichl. A stein variational newton method. In Advances in Neural Information Processing Systems, pages 9169–9179, 2018.
- Jean Dolbeault, Bruno Nazaret, and Giuseppe Savaré. A new class of transport distances between measures. *Calculus of Variations and Partial Differential Equations*, 34(2):193– 231, Feb 2009. ISSN 1432-0835. doi: 10.1007/s00526-008-0182-5. URL https://doi. org/10.1007/s00526-008-0182-5.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. Journal of Machine Learning Research, 12(Jul): 2121–2159, 2011.
- Alfredo Garbuno-Inigo, Franca Hoffmann, Wuchen Li, and Andrew M Stuart. Interacting langevin diffusions: Gradient structure and ensemble kalman sampler. arXiv preprint arXiv:1903.08866, 2019.

- Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. Journal of Machine Learning Research, 13(Mar):723– 773, 2012.
- Wen Huang. Optimization algorithms on riemannian manifolds with applications. 2013.
- Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the fokker-planck equation. SIAM journal on mathematical analysis, 29(1):1–17, 1998.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- John D Lafferty. The density manifold and configuration space quantization. *Transactions* of the American Mathematical Society, 305(2):699–741, 1988.
- Wuchen Li. Geometry of probability simplex via optimal transport. arXiv preprint arXiv:1803.06360, 2018.
- Wuchen Li. Diffusion hypercontractivity via generalized density manifold. *CoRR*, abs/1907.12546, 2019. URL http://arxiv.org/abs/1907.12546.
- Wuchen Li and Lexing Ying. Hessian transport gradient flows. Research in the Mathematical Sciences, 6(4):34, Oct 2019. ISSN 2197-9847. doi: 10.1007/s40687-019-0198-9. URL https://doi.org/10.1007/s40687-019-0198-9.
- Wuchen Li, Alex Tong Lin, and Guido Montúfar. Affine natural proximal learning. *Geo*metric science of information, 2019.
- Chang Liu, Jingwei Zhuo, Pengyu Cheng, Ruiyi Zhang, Jun Zhu, and Lawrence Carin. Accelerated first-order methods on the Wasserstein space for Bayesian inference. arXiv preprint arXiv:1807.01750, 2018.
- Chang Liu, Jingwei Zhuo, Pengyu Cheng, Ruiyi Zhang, and Jun Zhu. Understanding and accelerating particle-based variational inference. In *International Conference on Machine Learning*, pages 4082–4092, 2019.
- Qiang Liu. Stein variational gradient descent as gradient flow. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 3115-3123. Curran Associates, Inc., 2017. URL http://papers.nips.cc/paper/ 6904-stein-variational-gradient-descent-as-gradient-flow.pdf.
- Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In Advances in neural information processing systems, pages 2378– 2386, 2016.
- Yulong Lu, Jianfeng Lu, and James Nolen. Accelerating langevin sampling with birth-death. arXiv preprint arXiv:1905.09863, 2019.

- Yi-An Ma, Niladri Chatterji, Xiang Cheng, Nicolas Flammarion, Peter Bartlett, and Michael I Jordan. Is there an analog of nesterov acceleration for mcmc? arXiv preprint arXiv:1902.00996, 2019.
- Luigi Malagò and Giovanni Pistone. Combinatorial optimization with information geometry: The newton method. *Entropy*, 16(8):4260–4289, 2014.
- Luigi Malagò, Luigi Montrucchio, and Giovanni Pistone. Wasserstein riemannian geometry of positive definite matrices. arXiv preprint arXiv:1801.09269, 2018.
- Dimitra Maoutsa, Sebastian Reich, and Manfred Opper. Interacting particle solutions of fokker-planck equations through gradient-log-density estimation. *arXiv preprint* arXiv:2006.00702, 2020.
- James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pages 2408–2417, 2015.
- Klas Modin. Geometry of matrix decompositions seen through optimal transport and information geometry. *Journal of Geometric Mechanics*, 9(3), 2017.
- Yann Ollivier, Ludovic Arnold, Anne Auger, and Nikolaus Hansen. Information-geometric optimization algorithms: A unifying picture via invariance principles. *The Journal of Machine Learning Research*, 18(1):564–628, 2017.
- Felix Otto. The geometry of dissipative evolution equations: the porous medium equation. Communications in Partial Differential Equations, 26(1-2):101–174, 2001.
- Felix Otto and Cédric Villani. Generalization of an inequality by talagrand and links with the logarithmic sobolev inequality. *Journal of Functional Analysis*, 173(2):361–400, 2000.
- Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In International Conference on Machine Learning, pages 1530–1538, 2015.
- Bernhard Schölkopf, Ralf Herbrich, and Alex J Smola. A generalized representer theorem. In *International conference on computational learning theory*, pages 416–426. Springer, 2001.
- Umut Simsekli, Roland Badeau, Taylan Cemgil, and Gaël Richard. Stochastic quasi-newton langevin monte carlo. In International Conference on Machine Learning (ICML), 2016.
- Steven T Smith. Optimization techniques on riemannian manifolds. Fields institute communications, 3(3):113–135, 1994.
- Andrew M Stuart. Inverse problems: a Bayesian perspective. Acta numerica, 19:451–559, 2010.
- Amirhossein Taghvaei and Prashant G Mehta. Accelerated flow for probability distributions. arXiv preprint arXiv:1901.03317, 2019.

- Asuka Takatsu. On Wasserstein geometry of the space of gaussian measures. arXiv preprint arXiv:0801.2250, 2008.
- Cédric Villani. Topics in optimal transportation. American Mathematical Soc., 2003.
- Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- Yifei Wang and Wuchen Li. Accelerated information gradient flow. arXiv preprint arXiv:1909.02102, 2019.
- Andre Wibisono. Proximal langevin algorithm: Rapid convergence under isoperimetry. arXiv preprint arXiv:1911.01469, 2019.
- Yaguang Yang. Globally convergent optimization algorithms on riemannian manifolds: Uniform framework for unconstrained and constrained optimization. Journal of Optimization Theory and Applications, 132(2):245–265, 2007.