
Variational inference via Wasserstein gradient flows

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Along with Markov chain Monte Carlo (MCMC) methods, variational inference (VI)
2 has emerged as a central computational approach to large-scale Bayesian inference.
3 Rather than sampling from the true posterior π , VI aims at producing a simple but
4 effective approximation $\hat{\pi}$ to π for which summary statistics are easy to compute.
5 However, unlike the well-studied MCMC methodology, VI is still poorly understood
6 and dominated by heuristics. In this work, we propose principled methods for VI,
7 in which $\hat{\pi}$ is taken to be a Gaussian or a mixture of Gaussians, which rest upon
8 the theory of gradient flows on the Bures–Wasserstein space of Gaussian measures.
9 Akin to MCMC, it comes with strong theoretical guarantees when π is log-concave.

10 1 Introduction

11 This work brings together three active research areas: variational inference, variational Kalman
12 filtering, and gradient flows on the Wasserstein space.

13 **Variational inference.** The development of large-scale Bayesian methods has fueled the need for
14 fast and scalable methods to approximate complex distributions. More specifically, Bayesian method-
15 ology typically generates a high-dimensional posterior distribution $\pi \propto \exp(-V)$ that is known
16 only up to normalizing constants, making the computation even of simple summary statistics such as
17 the mean and covariance a major computational hurdle. To overcome this limitation, two distinct
18 computational approaches are largely favored. The first approach consists of Markov chain Monte
19 Carlo (MCMC) methods that rely on carefully constructed Markov chains which (approximately)
20 converge to π . For example, the *Langevin diffusion*

$$dX_t = -\nabla V(X_t) dt + \sqrt{2} dB_t, \quad (1)$$

21 where $(B_t)_{t \geq 0}$ denotes standard Brownian motion on \mathbb{R}^d , admits π as a stationary distribution.
22 Crucially, the Langevin diffusion can be discretized and implemented without knowledge of the
23 normalizing constant of π , leading to practical algorithms for Bayesian inference. Recent theoretical
24 efforts have produced sharp non-asymptotic convergence guarantees for algorithms based on the
25 Langevin diffusion (or variants thereof), with many results known when π is strongly log-concave or
26 satisfies isoperimetric assumptions [see, e.g., [Durmus et al., 2019](#), [Shen and Lee, 2019](#), [Vempala and
27 Wibisono, 2019](#), [Chen et al., 2020](#), [Dalalyan and Riou-Durand, 2020](#), [Chewi et al., 2021](#), [Lee et al.,
28 2021](#), [Ma et al., 2021](#), [Wu et al., 2021](#)].

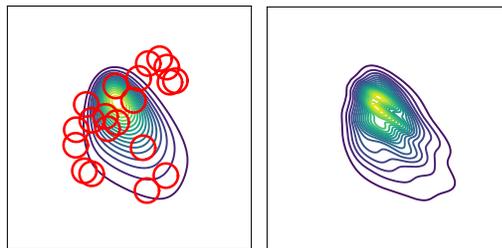
29 More recently, Variational Inference (VI) has emerged as a viable alternative to MCMC [[Jordan et al.,
30 1999](#), [Wainwright and Jordan, 2008](#), [Blei et al., 2017](#)]. The goal of VI is to approximate the posterior π
31 by a more tractable distribution $\hat{\pi} \in \mathcal{P}$ such that

$$\hat{\pi} \in \arg \min_{p \in \mathcal{P}} \text{KL}(p \parallel \pi). \quad (2)$$

32 A common example arises when \mathcal{P} is the class of product distributions, in which case $\hat{\pi}$ is called
 33 the *mean-field* approximation of \mathcal{P} . Unfortunately, by definition, mean-field approximations fail to
 34 capture important correlations present in the posterior π , and various remedies have been proposed,
 35 with varied levels of success. In this paper, we largely focus on obtaining a Gaussian approximation
 36 to π , that is, we take \mathcal{P} to be the class of non-degenerate Gaussian distributions on \mathbb{R}^d [Barber and
 37 Bishop, 1997, Seeger, 1999, Honkela and Valpola, 2004, Opper and Archambeau, 2009, Zhang et al.,
 38 2018]. The expressive power of the variational model may then be further increased by considering
 39 mixture distributions [Lin et al., 2019, Daudel and Douc, 2021, Daudel et al., 2021].

40 Although the solution $\hat{\pi}$ of (2) is no longer equal to the true posterior, variational inference remains
 41 heavily used in practice because the problem (2) can be solved for simple models \mathcal{P} via scalable
 42 optimization algorithms. In particular, VI avoids many of the practical hurdles associated with MCMC
 43 methods—such as the potentially long “burn-in” period of samplers and the lack of effective stopping
 44 criteria for the algorithm—while still producing informative summary statistics. In this regard, we
 45 highlight the fact that obtaining an approximation for the covariance matrix of π via MCMC methods
 46 requires drawing potentially many samples, whereas for many choices of \mathcal{P} (e.g., the Gaussian
 47 approximation) the covariance matrix of $\hat{\pi}$ can be directly obtained from the solution to the VI
 48 problem (2).

49 However, in contrast with MCMC methods, to
 50 date there have not been many theoretical guar-
 51 antees for VI, even when π is strongly log-
 52 concave and \mathcal{P} is taken to be the class of Gaus-
 53 sians $\mathcal{N}(m, \Sigma)$. The problem stems from the
 54 fact that the objective in (2) is typically non-
 55 convex in the pair (m, Σ) . Obtaining such guar-
 56 antees remains a pressing challenge for the field.



57 **Variational Kalman filtering.** There is also
 58 considerable interest in extending ideas behind
 59 variational inference to dynamical settings of
 60 Bayesian inference. Consider a general frame-
 61 work where $(\pi_t)_t$ represents the marginal laws
 62 of a stochastic process indexed by time t , which
 63 can be discrete or continuous. The goal is to recursively build a Gaussian approximation to $(\pi_t)_t$.

Figure 1: Left: randomly initialized mixture of 20 Gaussians (the initial covariances are depicted as red circles) and contour plot of a logistic target π . Right: contour lines of a mixture of Gaussians approximation $\hat{\pi}$ obtained from the gradient flow in Section 5.

64 As a concrete example, suppose that $(\pi_t)_{t \geq 0}$ denotes the marginal law of the solution to the Langevin
 65 diffusion (1). In the context of Bayesian optimal filtering and smoothing, Särkkä [2007] proposed
 66 the following heuristic. Let (m_t, Σ_t) denote the mean and covariance matrix of π_t . Then, it can be
 67 checked (see Section B.4) that

$$\begin{aligned} \dot{m}_t &= -\mathbb{E} \nabla V(X_t) \\ \dot{\Sigma}_t &= 2I - \mathbb{E}[\nabla V(X_t) \otimes (X_t - m_t) + (X_t - m_t) \otimes \nabla V(X_t)] \end{aligned} \quad (3)$$

68 where $X_t \sim \pi_t$. These ordinary differential equations (ODEs) are intractable because they involve
 69 expectations under the law of $X_t \sim \pi_t$, which is not available to the practitioner. However, if we
 70 replace $X_t \sim \pi_t$ with a Gaussian $Y_t \sim p_t = \mathcal{N}(m_t, \Sigma_t)$ with the same mean and covariance as X_t ,
 71 then the system of ODEs

$$\begin{aligned} \dot{m}_t &= -\mathbb{E} \nabla V(Y_t) \\ \dot{\Sigma}_t &= 2I - \mathbb{E}[\nabla V(Y_t) \otimes (Y_t - m_t) + (Y_t - m_t) \otimes \nabla V(Y_t)] \end{aligned} \quad (4)$$

72 yields a well-defined evolution of Gaussian distributions $(p_t)_{t \geq 0}$, which we may optimistically believe
 73 to be a good approximation of $(\pi_t)_{t \geq 0}$. Moreover, the system of ODEs can be numerically approxi-
 74 mated efficiently in practice using Gaussian quadrature rules to compute the above expectations. This
 75 is the principle behind the unscented Kalman filter [Julier et al., 2000].

76 In the context of the Langevin diffusion, Särkkä’s heuristic (4) provides a promising avenue towards
 77 computational VI. Indeed, since $\pi = \exp(-V)$ is the unique stationary distribution of the Langevin
 78 diffusion (1), an algorithm to approximate $(\pi_t)_{t \geq 0}$ is expected to furnish an algorithm to solve the
 79 VI problem (2). However, at present there is little theoretical understanding of how the system (4)
 80 approximates (3); moreover, Särkkä’s heuristic only provides Gaussian approximations, and it is
 81 unclear how to extend the system (4) to more complex models (e.g., mixtures of Gaussians).

82 **Our contributions: bridging the gap via Wasserstein gradient flows.** We show that the approxi-
 83 mation $(p_t)_{t \geq 0}$ in Särkkä’s heuristic (4) arises precisely as the gradient flow of the Kullback–Leibler
 84 (KL) divergence $\text{KL}(\cdot \parallel \pi)$ on the Bures–Wasserstein space of Gaussian distributions on \mathbb{R}^d endowed
 85 with the 2-Wasserstein distance from optimal transport [Villani, 2003]. This perspective allows us to
 86 not only understand its convergence but also to extend it to the richer space of mixtures of Gaussian
 87 distributions, and propose an implementation as a novel system of interacting “Gaussian particles”.
 88 Below, we proceed to describe our contributions in greater detail.

89 Our framework builds upon the seminal work of Jordan et al. [1998], which introduced the celebrated
 90 *JKO scheme* in order to give meaning to the idea that the evolving marginal law of the Langevin
 91 diffusion (1) is a gradient flow of $\text{KL}(\cdot \parallel \pi)$ on the Wasserstein space $\mathcal{P}_2(\mathbb{R}^d)$ of probability measures
 92 with finite second moments. Subsequently, in order to emphasize the Riemannian geometry underlying
 93 this result, Otto [2001] developed his eponymous calculus on $\mathcal{P}_2(\mathbb{R}^d)$, a framework which has had
 94 tremendous impact in analysis, geometry, PDE, probability, and statistics.

95 Inspired by this perspective, we show in Theorem 1 that Särkkä’s approximation $(p_t)_{t \geq 0}$ is also a gra-
 96 dient flow of $\text{KL}(\cdot \parallel \pi)$, with the main difference being that it is *constrained* to lie on the submanifold
 97 $\text{BW}(\mathbb{R}^d)$ of $\mathcal{P}_2(\mathbb{R}^d)$ consisting of Gaussian distributions, known as the Bures–Wasserstein manifold.
 98 In turn, our result paves the way for new theoretical understanding via the powerful theory of gradient
 99 flows. As a first step, using well-known results about convex functionals on the Wasserstein space,
 100 we show in Corollary 1 that $(p_t)_{t \geq 0}$ converges rapidly to the solution of the VI problem (2) with
 101 $\mathcal{P} = \text{BW}(\mathbb{R}^d)$ as soon as V is convex. Moreover, in Section 4.1, we apply numerical integration
 102 based on cubature rules for Gaussian integrals to the system of ODEs (4), thus arriving at a fast
 103 method with robust empirical performance (details in Sections H and I).

104 This combination of results brings VI closer to Langevin-based MCMC both on the practical and
 105 theoretical fronts, but still falls short of achieving non-asymptotic discretization guarantees as
 106 pioneered by Dalalyan [2017] for MCMC. To further close the theoretical gap between VI and the
 107 state of the art for MCMC, we propose in Section 4.2 a stochastic gradient descent (SGD) algorithm as
 108 a time discretization of the Bures–Wasserstein gradient flow. This algorithm comes with convergence
 109 guarantees that establish VI as a solid competitor to MCMC not only from a practical standpoint but
 110 also from a theoretical one. Both have their relative merits; whereas MCMC targets the true posterior,
 111 VI leads to fast computation of summary statistics of the approximation $\hat{\pi}$ to π .

112 In Section 5, we consider an extension of these ideas to the substantially more flexible class of
 113 mixtures of Gaussians. Namely, the space of mixtures of Gaussians can be identified as a Wasserstein
 114 space over $\text{BW}(\mathbb{R}^d)$ and hence inherits Otto’s differential calculus. Leveraging this viewpoint, in
 115 Theorem 3 we derive the gradient flow of $\text{KL}(\cdot \parallel \pi)$ over the space of mixtures of Gaussians and
 116 propose to implement it via a system of interacting particles. Unlike typical particle-based algorithms,
 117 here our particles correspond to Gaussian distributions, and the collection thereof to a Gaussian
 118 mixture which is better equipped to approximate a continuous measure. We validate the empirical
 119 performance of our method with promising experimental results (see Section I). Although we focus
 120 on the VI problem in this work, we anticipate that our notion of “Gaussian particles” may be a broadly
 121 useful extension of classical particle methods for PDEs.

122 **Related work.** Classical VI methods define a parametric family $\mathcal{P} = \{p_\theta : \theta \in \Theta\}$ and minimize
 123 $\theta \mapsto \text{KL}(p_\theta \parallel \pi)$ over $\theta \in \Theta$ using off-the-shelf optimization algorithms [Paisley et al., 2012, Ran-
 124 ganath et al., 2014]. Since (2) is an optimization problem over the space of probability distributions,
 125 we argue for methods that respect a natural geometric structure on this space. In this regard, previous

126 approaches to VI using natural gradients implicitly employ a different geometry [Wu et al., 2019,
 127 Khan and Håvard, 2022], namely the Fisher–Rao geometry [Amari and Nagaoka, 2000]. The appli-
 128 cation of Wasserstein gradient flows to VI was introduced earlier in work on normalizing flows and
 129 Stein Variational Gradient Descent (SVGD) [Liu and Wang, 2016, Liu, 2017].

130 The connection between VI and Kalman filtering was studied in the static case by Lambert et al. [2021,
 131 2022a], and extended to the dynamical case by Lambert et al. [2022b], providing a first justification
 132 of Särkkä’s heuristic in terms of local variational Gaussian approximation. In particular, the closest
 133 linear process to the Langevin diffusion (1) is a Gaussian process governed by a McKean–Vlasov
 134 equation whose Gaussian marginals have parameters evolving according to Särkkä’s ODEs.

135 Constrained gradient flows on the Wasserstein space have also been extensively studied [Carlen and
 136 Gangbo, 2003, Caglioti et al., 2009, Tudorascu and Wunsch, 2011, Eberle et al., 2017], although our
 137 interpretation of Särkkä’s heuristic is, to the best of our knowledge, new.

138 2 Background

139 In order to define gradient flows on the space of probability measures, we must first endow this space
 140 with a geometry; see Appendix B for more details. Given probability measures μ and ν on \mathbb{R}^d , define
 141 the 2-Wasserstein distance

$$W_2(\mu, \nu) = \left[\inf_{\gamma \in \mathcal{C}(\mu, \nu)} \int \|x - y\|^2 d\gamma(x, y) \right]^{1/2},$$

142 where $\mathcal{C}(\mu, \nu)$ is the set of *couplings* of μ and ν , that is, joint distributions on $\mathbb{R}^d \times \mathbb{R}^d$ whose
 143 marginals are μ and ν respectively. This quantity is finite as long as μ and ν belong to the space
 144 $\mathcal{P}_2(\mathbb{R}^d)$ of probability measures over \mathbb{R}^d with finite second moments. The 2-Wasserstein distance has
 145 the interpretation of measuring the smallest possible mean squared displacement of mass required
 146 to *transport* μ to ν ; we refer to Villani [2003, 2009], Santambrogio [2015] for textbook treatments
 147 on optimal transport. Unlike other notions of distance between probability measures, such as the
 148 total variation distance, the 2-Wasserstein distance respects the geometry of the underlying space \mathbb{R}^d ,
 149 leading to numerous applications in modern data science [see, e.g., Peyré and Cuturi, 2019].

150 The space $(\mathcal{P}_2(\mathbb{R}^d), W_2)$ is a metric space [Villani, 2003, Theorem 7.3], and we refer to it as the
 151 *Wasserstein space*. However, as shown by Otto [Otto, 2001], it has a far richer geometric structure:
 152 formally, $(\mathcal{P}_2(\mathbb{R}^d), W_2)$ can be viewed as a Riemannian manifold, a fact which allows for considering
 153 gradient flows of functionals on $\mathcal{P}_2(\mathbb{R}^d)$. A fundamental example of such a functional is the KL
 154 divergence $\text{KL}(\cdot \parallel \pi)$ to a target density $\pi = \exp(-V)$ on \mathbb{R}^d , for which Jordan et al. [1998] showed
 155 that the Wasserstein gradient flow is the same as the evolution of the marginal law of the Langevin
 156 diffusion (1). This optimization perspective has had tremendous impact on our understanding and
 157 development of MCMC algorithms [Wibisono, 2018].

158 3 Variational inference with Gaussians

159 In this section we describe our problem using two equivalent approaches: a variational approach based
 160 on a modified version of the JKO scheme of Jordan et al. [1998] (Section 3.1), and a Wasserstein
 161 gradient flow approach based on Otto calculus (Section 3.2). Both lead to the same result (Section
 162 3.3). While the former is more accessible to readers who are unfamiliar with gradient flows on the
 163 Wasserstein space, the latter leads to strong convergence guarantees (Section 3.4).

164 3.1 Variational approach: the Bures–JKO scheme

165 The space of non-degenerate Gaussian distributions on \mathbb{R}^d equipped with the W_2 distance forms the
 166 *Bures–Wasserstein space* $\text{BW}(\mathbb{R}^d) \subseteq \mathcal{P}_2(\mathbb{R}^d)$. On $\text{BW}(\mathbb{R}^d)$, the Wasserstein distance $W_2^2(p_0, p_1)$
 167 between two Gaussians $p_0 = \mathcal{N}(m_0, \Sigma_0)$ and $p_1 = \mathcal{N}(m_1, \Sigma_1)$ admits the following closed form:

$$W_2^2(p_0, p_1) = \|m_0 - m_1\|^2 + \mathcal{B}^2(\Sigma_0, \Sigma_1), \quad (5)$$

168 where $\mathcal{B}^2(\Sigma_0, \Sigma_1) = \text{tr}(\Sigma_0 + \Sigma_1 - 2(\Sigma_0^{\frac{1}{2}}\Sigma_1\Sigma_0^{\frac{1}{2}})^{\frac{1}{2}})$ is the squared Bures metric [Bures, 1969].

169 Given a target density $\pi = \exp(-V)$ on \mathbb{R}^d , and with a step size $h > 0$, we may define the iterates
170 of the proximal point algorithm

$$p_{k+1,h} := \arg \min_{p \in \text{BW}(\mathbb{R}^d)} \left\{ \text{KL}(p \parallel \pi) + \frac{1}{2h} W_2^2(p, p_{k,h}) \right\}. \quad (6)$$

171 Using (5), this is an explicit optimization problem involving the mean and covariance matrix of p .
172 Although (6) is not solvable in closed form, by letting $h \searrow 0$ we obtain a limiting curve $(p_t)_{t \geq 0}$
173 via $p_t = \lim_{h \searrow 0} p_{\lfloor t/h \rfloor, h}$, which can be interpreted as the Bures–Wasserstein gradient flow of the
174 KL divergence $\text{KL}(\cdot \parallel \pi)$. This procedure mimics the JKO scheme [Jordan et al., 1998] with the
175 additional constraint that the iterates lie in $\text{BW}(\mathbb{R}^d)$, and we therefore call it the Bures–JKO scheme.

176 3.2 Geometric approach: the Bures–Wasserstein gradient flow of the KL divergence

177 In the formal sense of Otto described above, $\text{BW}(\mathbb{R}^d)$ is a submanifold of $\mathcal{P}_2(\mathbb{R}^d)$. Moreover, since
178 Gaussians can be parameterized by their mean and covariance, $\text{BW}(\mathbb{R}^d)$ can be identified with the
179 manifold $\mathbb{R}^d \times \mathbf{S}_{++}^d$, where \mathbf{S}_{++}^d is the cone of symmetric positive definite $d \times d$ matrices. Hence,
180 $\text{BW}(\mathbb{R}^d)$ is a genuine Riemannian manifold in its own right [see Bhatia et al., 2019], and gradient
181 flows can be defined using Riemannian geometry [do Carmo, 1992]. See Section B.3 for more details.
182 Since the functional $\mu \mapsto \mathcal{F}(\mu) = \text{KL}(\mu \parallel \pi)$ defined over $\mathcal{P}_2(\mathbb{R}^d)$ restricts to a functional over
183 $\text{BW}(\mathbb{R}^d)$, we can also consider the gradient flow of \mathcal{F} over the Bures–Wasserstein space; note that
184 this latter gradient flow is necessarily a curve $(p_t)_{t \geq 0}$ such that each p_t is a Gaussian measure.

185 3.3 Variational inference via the Bures–Wasserstein gradient flow

186 Using either approach, we can prove the following theorem.

187 **Theorem 1.** *Let $\pi = \exp(-V)$ be the target density on \mathbb{R}^d . Then, the limiting curve $(p_t)_{t \geq 0}$ where
188 $p_t = \mathcal{N}(m_t, \Sigma_t)$ is obtained via the Bures–JKO scheme (6), or equivalently, the Bures–Wasserstein
189 gradient flow $(p_t)_{t \geq 0}$ of the KL divergence $\text{KL}(\cdot \parallel \pi)$, satisfies Särkkä’s system of ODEs (4).*

190 *Proof.* The proof using the Bures–JKO scheme is given in Section A.1 and the proof using Otto
191 calculus is presented in Section C. \square

192 This theorem shows that Särkkä’s heuristic (4) precisely yields the Wasserstein gradient flow of the
193 KL divergence over the submanifold $\text{BW}(\mathbb{R}^d)$. Equipped with this interpretation, we are now able
194 to obtain information about the asymptotic behavior of the approximation $(p_t)_{t \geq 0}$. Namely, we can
195 hope that it converges to constrained minimizer $\hat{\pi} = \arg \min_{p \in \text{BW}(\mathbb{R}^d)} \text{KL}(p \parallel \pi)$, i.e., precisely the
196 solution to the VI problem (2). In the next section, we show that this convergence in fact holds as
197 soon as V is convex, and moreover with quantitative rates.

198 The solution $\hat{\pi}$ to (2), and consequently the limit point of Särkkä’s approximation, is well-studied in
199 the variational inference literature [see, e.g., Opper and Archambeau, 2009], and we recall standard
200 facts about $\hat{\pi}$ here for completeness. It is known that $\hat{\pi}$ satisfies the equations

$$\mathbb{E}_{\hat{\pi}} \nabla V = 0 \quad \text{and} \quad \mathbb{E}_{\hat{\pi}} \nabla^2 V = \hat{\Sigma}^{-1}, \quad (7)$$

201 where $\hat{\Sigma}$ is the covariance matrix of $\hat{\pi}$ (these equations can also be derived as first-order necessary
202 conditions by setting the Bures–Wasserstein gradient derived in Section C to zero). In particular, it
203 follows from (7) that if $\nabla^2 V$ enjoys the bounds $\alpha I \preceq \nabla^2 V \preceq \beta I$ for some $-\infty \leq \alpha \leq \beta \leq \infty$,
204 then any solution $\hat{\pi}$ to the constrained problem also satisfies $\beta^{-1} I \preceq \hat{\Sigma} \preceq (\alpha \vee 0)^{-1} I$.

205 3.4 Continuous-time convergence

206 Besides providing an intuitive interpretation of Särkkä’s heuristic, Theorem 1 readily yields conver-
207 gence criteria for the system (4) which rest upon general principles for gradient flows. We begin with

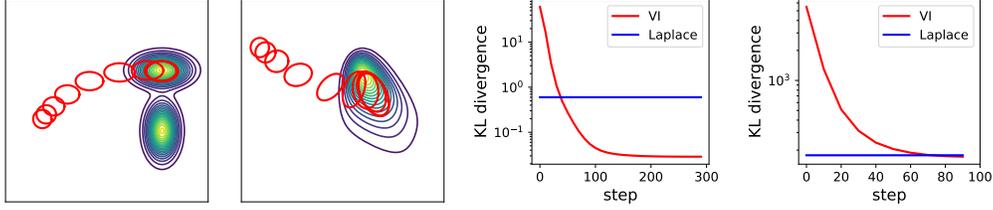


Figure 2: Two left plots: approximation of a bimodal target and a logistic target. Two right plots: convergence of the KL in dimension 2 and 100 for the logistic target. Our algorithm yields better approximation in KL than the Laplace approximation (see Appendix H.4 for details).

208 a key observation. For a functional $\mathcal{F} : \text{BW}(\mathbb{R}^d) \rightarrow \mathbb{R} \cup \{\infty\}$ and $\alpha \in \mathbb{R}$, we say that \mathcal{F} is α -convex
 209 if for all constant-speed geodesics $(p_t)_{t \in [0,1]}$ in $\text{BW}(\mathbb{R}^d)$,

$$\mathcal{F}(p_t) \leq (1-t)\mathcal{F}(p_0) + t\mathcal{F}(p_1) - \frac{\alpha t(1-t)}{2} W_2^2(p_0, p_1), \quad t \in [0, 1].$$

210 **Lemma 1.** For any $\alpha \in \mathbb{R}$, if $\nabla^2 V \succeq \alpha I$, then $\text{KL}(\cdot \parallel \pi)$ is α -convex on $\text{BW}(\mathbb{R}^d)$.

211 *Proof.* The assumption that $\nabla^2 V \succeq \alpha I$ entails that the functional $\text{KL}(\cdot \parallel \pi)$ is α -convex on the
 212 entire Wasserstein space $(\mathcal{P}_2(\mathbb{R}^d), W_2)$ [see, e.g., Villani, 2009, Theorem 17.15]. Since $\text{BW}(\mathbb{R}^d)$ is a
 213 geodesically convex subset of $\mathcal{P}_2(\mathbb{R}^d)$ (see Section B.3), then the geodesics in $\text{BW}(\mathbb{R}^d)$ agree with
 214 the geodesics in $\mathcal{P}_2(\mathbb{R}^d)$, from which it follows that $\text{KL}(\cdot \parallel \pi)$ is α -convex on $\text{BW}(\mathbb{R}^d)$. \square

215 Consequently, we obtain the following corollary. Its proof is postponed to Section D.

216 **Corollary 1.** Suppose that $\nabla^2 V \succeq \alpha I$ for some $\alpha \in \mathbb{R}$. Then, for any $p_0 \in \text{BW}(\mathbb{R}^d)$, there is a
 217 unique solution to the $\text{BW}(\mathbb{R}^d)$ gradient flow of $\text{KL}(\cdot \parallel \pi)$ started at p_0 . Moreover:

- 218 1. If $\alpha > 0$, then for all $t \geq 0$, $W_2^2(p_t, \hat{\pi}) \leq \exp(-2\alpha t) W_2^2(p_0, \hat{\pi})$.
- 219 2. If $\alpha > 0$, then for all $t \geq 0$, $\text{KL}(p_t \parallel \pi) - \text{KL}(\hat{\pi} \parallel \pi) \leq \exp(-2\alpha t) \{\text{KL}(p_0 \parallel \pi) - \text{KL}(\hat{\pi} \parallel \pi)\}$.
- 220 3. If $\alpha = 0$, then for all $t > 0$, $\text{KL}(p_t \parallel \pi) - \text{KL}(\hat{\pi} \parallel \pi) \leq \frac{1}{2t} W_2^2(p_0, \hat{\pi})$.

221 The assumption that $\nabla^2 V \succeq \alpha I$ for some $\alpha > 0$, i.e., that π is *strongly log-concave*, is a standard
 222 assumption in the MCMC literature. Under this same assumption, Corollary 1 yields convergence
 223 for the Bures–Wasserstein gradient flow of $\text{KL}(\cdot \parallel \pi)$; however, the flow must first be discretized in
 224 time for implementation. If we assume additionally that the smoothness condition $\nabla^2 V \preceq \beta I$ holds,
 225 then a surge of recent research has succeeded in obtaining precise non-asymptotic guarantees for
 226 discretized MCMC algorithms. In Section 4.2 below, we will show how to do the same for VI.

227 4 Time discretization of the Bures–Wasserstein gradient flow

228 We are now equipped with dual perspectives on a dynamical solution to Gaussian VI: ODE and
 229 gradient flow. Each perspective leads to a different implementation. On the one hand, we discretize
 230 the system of ODEs defined in (4) using numerical integration. On the other, we discretize the
 231 gradient flow using stochastic gradient descent in the Bures–Wasserstein space.

232 4.1 Numerical integration of the ODEs

233 The system of ODEs (4) can be integrated in time using a classical Runge–Kutta scheme. The
 234 expectations under a Gaussian support are approximated by cubature rules used in Kalman filter-
 235 ing [Arasaratnam and Haykin, 2009]. Moreover, a square root version of the ODE is also considered
 236 to ensure that covariance matrices remain symmetric and positive. See Appendix H.2 for more details.
 237 We have tested our method on a bimodal distribution and on a posterior distribution arising from a
 238 logistic regression problem. We observe fast convergence as shown in Figure 2.

239 **4.2 Bures–Wasserstein SGD and theoretical guarantees for VI**

240 Although the ODE discretization proposed in the preceding section enjoys strong empirical perfor-
 241 mance, it is unclear how to quantify its impact on the convergence rates established in Corollary 1.
 242 Therefore, we now propose a stochastic gradient descent algorithm over the Bures–Wasserstein space,
 243 for which useful analysis tools have been developed [Chewi et al., 2020, Altschuler et al., 2021]. This
 244 approach bypasses the use of the system of ODEs (4), and instead discretizes the Bures–Wasserstein
 245 gradient flow directly. Under the standard assumption of strong log-concavity and log-smoothness, it
 246 leads to an algorithm (Algorithm 1) for approximating $\hat{\pi}$ with provable convergence guarantees.

247 Algorithm 1 maintains a sequence of Gaussian
 248 distributions $(p_k)_{k \in \mathbb{N}}$; here (m_k, Σ_k) denote the
 249 mean vector and covariance matrix at iteration k
 250 (see Section E for a derivation of the algorithm
 251 as SGD in the Bures–Wasserstein space). The
 252 clipping operator clip^τ , which is introduced
 253 purely for the purpose of theoretical analysis,
 254 simply truncates the eigenvalues from above;
 255 see Section E. Our theoretical result for VI is
 256 given as the following theorem, whose proof is
 257 deferred to Section E.

Algorithm 1 Bures–Wasserstein SGD

Require: strong convexity parameter $\alpha > 0$; step
 size $h > 0$; mean m_0 and covariance matrix Σ_0
for $k = 1, \dots, N$ **do**
 draw a sample $\hat{X}_k \sim p_k$
 set $m_{k+1} \leftarrow m_k - h \nabla V(\hat{X}_k)$
 set $M_k \leftarrow I - h (\nabla^2 V(\hat{X}_k) - \Sigma_k^{-1})$
 set $\Sigma_k^+ \leftarrow M_k \Sigma_k M_k$
 set $\Sigma_{k+1} \leftarrow \text{clip}^{1/\alpha} \Sigma_k$
end for

258 **Theorem 2.** Assume that $0 \prec \alpha I \preceq \nabla^2 V \preceq I$. Also, assume that $h \leq \frac{\alpha}{6}$ and that we initialize
 259 Algorithm 1 at a matrix satisfying $\frac{\alpha}{4} I \preceq \Sigma_{\mu_0} \preceq \frac{1}{\alpha} I$. Then, for all $k \in \mathbb{N}$,

$$\mathbb{E} W_2^2(p_k, \hat{\pi}) \leq \exp(-\alpha k h) W_2^2(p_0, \hat{\pi}) + \frac{21dh}{\alpha^2}.$$

260 In particular, we obtain $\mathbb{E} W_2^2(p_k, \hat{\pi}) \leq \varepsilon^2$ provided we set $h \asymp \frac{\alpha^2 \varepsilon^2}{d}$ and the number of iterations to
 261 be $k \gtrsim \frac{d}{\alpha^3 \varepsilon^2} \log(W_2(p_0, \hat{\pi})/\varepsilon)$.

262 The upper bound $\nabla^2 V \preceq I$ is notationally convenient for our proof but not necessary; in any case,
 263 any strongly log-concave and log-smooth density π can be rescaled so that the assumption holds.

264 Theorem 2 is similar in flavor to modern results for MCMC, both in terms of the assumptions (Hessian
 265 bounds and query access to the derivatives¹ of V) and the conclusion (a non-asymptotic polynomial-
 266 time algorithmic guarantee). We hope that such an encouraging result for VI will prompt more
 267 theoretical studies aimed at closing the gap between the two approaches.

268 **5 Variational inference with mixtures of Gaussians**

269 Thus far, we have shown that the tractability of Gaussians can be readily exploited in the context of
 270 Bures–Wasserstein gradient flows and translated into useful results for variation inference. Never-
 271 theless, these results are limited by the lack of expressivity of Gaussians, namely their inability to
 272 capture complex features such as multimodality and, more generally, heterogeneity. To overcome
 273 this limitation, mixtures of Gaussians arise as a natural and powerful alternative; indeed, universal
 274 approximation of arbitrary probability measures by mixtures of Gaussians is well-known [see, e.g.,
 275 Delon and Desolneux, 2020]. As we show next, the space of mixtures of Gaussians can also be
 276 equipped with a Wasserstein structure which gives rise to implementable gradient flows.

277 **5.1 Geometry of the space of mixtures of Gaussians**

278 We begin with the key observation already made by Chen et al. [2019], that any mixture of Gaussians
 279 can be canonically identified with a probability distribution (the mixing distribution) over the param-
 280 eter space $\Theta = \mathbb{R}^d \times \mathbf{S}_{++}^d$ (the space of means and covariance matrices). Explicitly a probability

¹A notable downside of Algorithm 1 is the requirement of a Hessian oracle for V , which results in a higher per-iteration cost than typical MCMC samplers.

281 measure $\mu \in \mathcal{P}(\Theta)$ corresponds to a Gaussian mixture as follows:

$$\mu \quad \leftrightarrow \quad \mathfrak{p}_\mu := \int p_\theta \, d\mu(\theta), \quad (8)$$

282 where p_θ is the Gaussian distribution with parameters $\theta \in \Theta$. Equivalently, μ can be thought of as a
 283 probability measure over $\text{BW}(\mathbb{R}^d)$, and hence the space of Gaussian mixtures on \mathbb{R}^d can be identified
 284 with the Wasserstein space $\mathcal{P}_2(\text{BW}(\mathbb{R}^d))$ over the Bures–Wasserstein space which is endowed with
 285 the distance (5) between Gaussian measures. Indeed, the theory of optimal transport can be developed
 286 with any Riemannian manifold (rather than \mathbb{R}^d) as the base space [Villani, 2009]. As before, the
 287 space $\mathcal{P}_2(\text{BW}(\mathbb{R}^d))$ is endowed with a formal Riemannian structure, which respects the geometry of
 288 the base space $\text{BW}(\mathbb{R}^d)$, and we can consider Wasserstein gradient flows over $\mathcal{P}_2(\text{BW}(\mathbb{R}^d))$.

289 Note that this framework encompasses both discrete mixtures of Gaussians (when μ is a discrete
 290 measure) and continuous mixtures of Gaussians. In the case when the mixing distribution μ is
 291 discrete, the geometry of $\mathcal{P}_2(\text{BW}(\mathbb{R}^d))$ was studied by Chen et al. [2019], Delon and Desolneux
 292 [2020]. An important insight of our work, however, is that it is fruitful to consider the full space
 293 $\mathcal{P}_2(\text{BW}(\mathbb{R}^d))$ for deriving gradient flows, even if we eventually develop algorithms which propagate
 294 a finite number of mixture components.

295 5.2 Gradient flow of the KL divergence and particle discretization

296 We consider the gradient flow of the KL divergence functional

$$\mu \mapsto \mathcal{F}(\mu) := \text{KL}(\mathfrak{p}_\mu \parallel \pi) \quad (9)$$

297 over the space $\mathcal{P}_2(\text{BW}(\mathbb{R}^d))$. The proof of the following theorem is given in Section F.

298 **Theorem 3.** *The gradient flow $(\mu_t)_{t \geq 0}$ of the functional \mathcal{F} defined in (9) over $\mathcal{P}_2(\text{BW}(\mathbb{R}^d))$ can be
 299 described as follows. Let $\theta_0 = (m_0, \Sigma_0) \sim \mu_0$, and let $\theta_t = (m_t, \Sigma_t)$ evolve according to the ODE*

$$\begin{cases} \dot{m}_t = -\mathbb{E} \nabla \ln \frac{\mathfrak{p}_{\mu_t}}{\pi}(Y_t) \\ \dot{\Sigma}_t = -\mathbb{E} \nabla^2 \ln \frac{\mathfrak{p}_{\mu_t}}{\pi}(Y_t) \Sigma_t - \Sigma_t \mathbb{E} \nabla^2 \ln \frac{\mathfrak{p}_{\mu_t}}{\pi}(Y_t) \end{cases} \quad (10)$$

300 where $Y_t \sim \mathcal{N}(m_t, \Sigma_t)$. Then $\theta_t \sim \mu_t$.

301 The gradient flow in Theorem 3 describes the evolution of a particle θ_t which describes the parameters
 302 of a Gaussian measure, hence the name *Gaussian particle*. The intuition behind this evolution is
 303 as follows. Suppose we draw infinitely many initial particles (each being a Gaussian) from μ_0 .
 304 By evolving all those particles through (10), which interact with each other via the term \mathfrak{p}_{μ_t} , they
 305 tend to aggregate in some parts of the space of Gaussian parameters and spread out in others. This
 306 distribution of Gaussian particles is precisely the mixing measure μ_t , which, in turn, corresponds to a
 307 Gaussian mixture. Since an infinite number of Gaussian particles is impractical, consider initializing
 308 this evolution at a finitely supported distribution μ_0 , thus corresponding to a more familiar Gaussian
 309 mixture model with a finite number of components:

$$\mu_0 = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_0^{(i)}} = \frac{1}{N} \sum_{i=1}^N \delta_{(m_0^{(i)}, \Sigma_0^{(i)})} \quad \leftrightarrow \quad \mathfrak{p}_{\mu_0} := \frac{1}{N} \sum_{i=1}^N \mathfrak{p}_{(m_0^{(i)}, \Sigma_0^{(i)})}.$$

Interestingly, it can be readily checked that the system of ODEs (10) thus initialized maintains a finite
 mixture distribution:

$$\mu_t = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_t^{(i)}} = \frac{1}{N} \sum_{i=1}^N \delta_{(m_t^{(i)}, \Sigma_t^{(i)})},$$

310 where the parameters $\theta_t^{(i)} = (m_t^{(i)}, \Sigma_t^{(i)})$ evolve according to the following interacting particle
 311 system, for $i \in [N]$

$$\dot{m}_t^{(i)} = -\mathbb{E} \nabla \ln \frac{\mathfrak{p}_{\mu_t}}{\pi}(Y_t^{(i)}), \quad (11)$$

$$\dot{\Sigma}_t^{(i)} = -\mathbb{E} \nabla^2 \ln \frac{\mathfrak{p}_{\mu_t}}{\pi}(Y_t^{(i)}) \Sigma_t^{(i)} - \Sigma_t^{(i)} \mathbb{E} \nabla^2 \ln \frac{\mathfrak{p}_{\mu_t}}{\pi}(Y_t^{(i)}), \quad (12)$$

312 where $Y_t^{(i)} \sim p_{\theta_t^{(i)}}$. This finite system of particles can now be implemented using the same numerical
 313 tools as for Gaussian VI, see Section I. Note that due to this property of the dynamics, we can hope at
 314 best to converge to the best mixture of N Gaussians approximating π , but this approximation error is
 315 expected to vanish as $N \rightarrow \infty$.

316 The above system of particles may also be derived using a proximal point method similar to the
 317 Bures–JKO scheme, see Section A.2. Indeed, infinitesimally, it has the variational interpretation

$$(\theta_{t+h}^{(1)}, \dots, \theta_{t+h}^{(N)}) \approx \arg \min_{\theta^{(1)}, \dots, \theta^{(N)} \in \Theta} \left\{ \text{KL} \left(\frac{1}{N} \sum_{i=1}^N p_{\theta^{(i)}} \parallel \pi \right) + \frac{1}{2Nh} \sum_{i=1}^N W_2^2(p_{\theta^{(i)}}, p_{\theta_t^{(i)}}) \right\}.$$

318 Reassuringly, Equations (11)-(12) reduce to (4) when $\mu_0 = \delta_{(m_0, \Sigma_0)}$ is a point mass, indicating that
 319 the theorem provides a natural extension of our previous results. However, although the model (8)
 320 is substantially more expressive than the Gaussian VI considered in Section 3, it has the downside
 321 that we lose many of the theoretical guarantees. For example, even when V is convex, the objective
 322 functional \mathcal{F} considered here need not be convex; see Section G. We nevertheless validate the practical
 323 utility of our approach in experiments (see Figure 3 and Section I).

324 Unlike typical interacting particle systems which arise from discretizations of Wasserstein gradient
 325 flows, at each time t , the distribution p_{μ_t} is continuous. This extension provides considerably more
 326 flexibility—from a mixture of point masses to a mixture of Gaussians—compared to interacting
 327 particle-based algorithms hitherto considered for either sampling [Liu and Wang, 2016, Liu, 2017,
 328 Duncan et al., 2019, Chewi et al., 2020], or solving partial differential equations [Carrillo et al., 2011,
 329 2012, Bonaschi et al., 2015, Craig and Bertozzi, 2016, Carrillo et al., 2019, Craig et al., 2022].

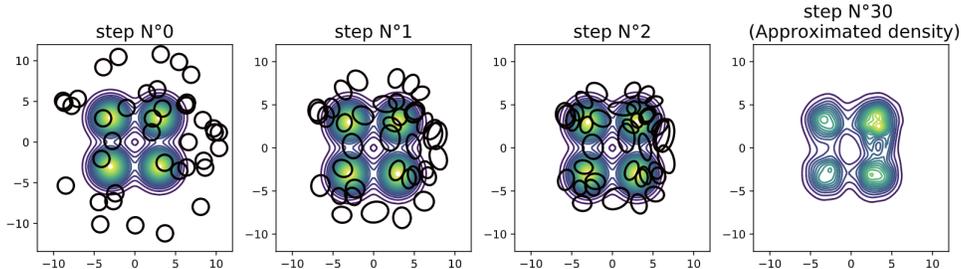


Figure 3: Approximation of a Gaussian mixture target π with 40 Gaussian particles. The particles are represented by their covariance ellipsoids shown at Steps 0, 1, and 2. The right figure shows the final step with the approximated density in contour-lines. More figures are available in Appendix I.

330 6 Conclusion

331 Using the powerful theory of Wasserstein gradient flows, we derived new algorithms for VI using either
 332 Gaussians or mixtures of Gaussians as approximating distributions. The consequences are twofold.
 333 On the one hand, strong convergence guarantees under classical conditions contribute markedly to
 334 closing the theoretical gap between MCMC and Gaussian VI. On the other hand, discretization of the
 335 Wasserstein gradient flow for mixtures of Gaussians yields a new *Gaussian particle method* which
 336 appears to be significantly more powerful than classical particle methods.

337 We conclude by briefly listing some possible directions for future study. For Gaussian variational
 338 inference, our theoretical result (Theorem 2) can be strengthened by weakening the assumption that π
 339 is strongly log-concave, or by developing algorithms which do not require Hessian information for V .
 340 For mixtures of Gaussians, it is desirable to design a principled algorithm which also allows for the
 341 mixture weights to be updated.

342 **References**

- 343 Jason Altschuler, Sinho Chewi, Patrik Gerber, and Austin J. Stromme. Averaging on the Bures–
344 Wasserstein manifold: dimension-free convergence of gradient descent. In *Advances in Neural*
345 *Information Processing Systems*, volume 34, pages 22132–22145, 2021.
- 346 Shun-ichi Amari and Hiroshi Nagaoka. *Methods of information geometry*, volume 191 of *Translations*
347 *of Mathematical Monographs*. American Mathematical Society, Providence, RI, 2000.
- 348 Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows in metric spaces and in the*
349 *space of probability measures*. Lectures in Mathematics ETH Zürich. Birkhäuser Verlag, Basel,
350 second edition, 2008.
- 351 Ienkaran Arasaratnam and Simon Haykin. Cubature Kalman filters. *IEEE Trans. Automat. Control*,
352 54(6):1254–1269, 2009.
- 353 Dominique Bakry, Ivan Gentil, and Michel Ledoux. *Analysis and geometry of Markov diffusion oper-*
354 *ators*, volume 348 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles*
355 *of Mathematical Sciences]*. Springer, Cham, 2014.
- 356 David Barber and Christopher Bishop. Ensemble learning for multi-layer networks. In *Advances in*
357 *Neural Information Processing Systems*, volume 10, 1997.
- 358 Jean-David Benamou and Yann Brenier. A numerical method for the optimal time-continuous mass
359 transport problem and related problems. In *Monge Ampère equation: applications to geometry*
360 *and optimization (Deerfield Beach, FL, 1997)*, volume 226 of *Contemp. Math.*, pages 1–11. Amer.
361 Math. Soc., Providence, RI, 1999.
- 362 Rajendra Bhatia, Tanvi Jain, and Yongdo Lim. On the Bures–Wasserstein distance between positive
363 definite matrices. *Expo. Math.*, 37(2):165–191, 2019.
- 364 David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational inference: A review for statisticians.
365 *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- 366 Giovanni A. Bonaschi, José A. Carrillo, Marco Di Francesco, and Mark A. Peletier. Equivalence of
367 gradient flows and entropy solutions for singular nonlocal interaction equations in 1D. *ESAIM*
368 *Control Optim. Calc. Var.*, 21(2):414–441, 2015.
- 369 Donald Bures. An extension of Kakutani’s theorem on infinite product measures to the tensor product
370 of semifinite w^* -algebras. *Trans. Amer. Math. Soc.*, 135:199–212, 1969.
- 371 Emanuele Caglioti, Mario Pulvirenti, and Frédéric Rousset. On a constrained 2-D Navier–Stokes
372 equation. *Comm. Math. Phys.*, 290(2):651–677, 2009.
- 373 Eric A. Carlen and Wilfrid Gangbo. Constrained steepest descent in the 2-Wasserstein metric. *Ann.*
374 *of Math. (2)*, 157(3):807–846, 2003.
- 375 José A. Carrillo, Marco Di Francesco, Alessio Figalli, Thomas Laurent, and Dejan Slepčev. Global-
376 in-time weak measure solutions and finite-time aggregation for nonlocal interaction equations.
377 *Duke Math. J.*, 156(2):229–271, 2011.
- 378 José A. Carrillo, Marco Di Francesco, Alessio Figalli, Thomas Laurent, and Dejan Slepčev. Confine-
379 ment in nonlocal interaction equations. *Nonlinear Anal.*, 75(2):550–558, 2012.
- 380 José A. Carrillo, Katy Craig, and Francesco S. Patacchini. A blob method for diffusion. *Calc. Var.*
381 *Partial Differential Equations*, 58(2):Paper No. 53, 53, 2019.
- 382 Yongxin Chen, Tryphon T. Georgiou, and Allen Tannenbaum. Optimal transport for Gaussian mixture
383 models. *IEEE Access*, 7:6269–6278, 2019.

- 384 Yuansi Chen, Raaz Dwivedi, Martin J. Wainwright, and Bin Yu. Fast mixing of Metropolized
385 Hamiltonian Monte Carlo: benefits of multi-step gradients. *J. Mach. Learn. Res.*, 21:Paper No. 92,
386 71, 2020.
- 387 Sinho Chewi, Thibaut Le Gouic, Chen Lu, Tyler Maunu, and Philippe Rigollet. SVGD as a kernelized
388 Wasserstein gradient flow of the chi-squared divergence. In *Advances in Neural Information*
389 *Processing Systems*, volume 33, pages 2098–2109, 2020.
- 390 Sinho Chewi, Tyler Maunu, Philippe Rigollet, and Austin J. Stromme. Gradient descent algorithms
391 for Bures–Wasserstein barycenters. In *Proceedings of the Conference on Learning Theory*, volume
392 125, pages 1276–1304. PMLR, 09–12 Jul 2020.
- 393 Sinho Chewi, Murat A. Erdogdu, Mufan B. Li, Ruoqi Shen, and Matthew Zhang. Analysis of
394 Langevin Monte Carlo from Poincaré to log-Sobolev. *arXiv e-prints*, art. arXiv:2112.12662, 2021.
- 395 Katy Craig and Andrea L. Bertozzi. A blob method for the aggregation equation. *Math. Comp.*, 85
396 (300):1681–1717, 2016.
- 397 Katy Craig, Karthik Elamvazhuthi, Matt Haberland, and Olga Turanova. A blob method for inho-
398 mogeneous diffusion with applications to multi-agent control and sampling. *arXiv e-prints*, art.
399 arXiv:2202.12927, March 2022.
- 400 Arnak S. Dalalyan. Theoretical guarantees for approximate sampling from smooth and log-concave
401 densities. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 79(3):
402 651–676, 2017.
- 403 Arnak S. Dalalyan and Lionel Riou-Durand. On sampling from a log-concave density using kinetic
404 Langevin diffusions. *Bernoulli*, 26(3):1956–1988, 2020.
- 405 Kamélia Daudel and Randal Douc. Mixture weights optimisation for alpha-divergence variational
406 inference. In *Advances in Neural Information Processing Systems*, volume 34, pages 4397–4408,
407 2021.
- 408 Kamélia Daudel, Randal Douc, and François Portier. Infinite-dimensional gradient-based descent for
409 alpha-divergence minimisation. *Ann. Statist.*, 49(4):2250–2270, 2021.
- 410 Julie Delon and Agnès Desolneux. A Wasserstein-type distance in the space of Gaussian mixture
411 models. *SIAM J. Imaging Sci.*, 13(2):936–970, 2020.
- 412 Manfredo P. do Carmo. *Riemannian geometry*. Mathematics: Theory & Applications. Birkhäuser
413 Boston, Inc., Boston, MA, 1992. Translated from the second Portuguese edition by Francis
414 Flaherty.
- 415 Andrew Duncan, Nikolas Nuesken, and Lukasz Szpruch. On the geometry of Stein variational
416 gradient descent. *arXiv e-prints*, art. arXiv:1912.00894, December 2019.
- 417 Alain Durmus, Szymon Majewski, and Błażej Miasojedow. Analysis of Langevin Monte Carlo via
418 convex optimization. *J. Mach. Learn. Res.*, 20:Paper No. 73, 46, 2019.
- 419 Simon Eberle, Barbara Niethammer, and André Schlichting. Gradient flow formulation and longtime
420 behaviour of a constrained Fokker–Planck equation. *Nonlinear Anal.*, 158:142–167, 2017.
- 421 Antti Honkela and Harri Valpola. Unsupervised variational Bayesian learning of nonlinear models.
422 In *Advances in Neural Information Processing Systems*, volume 17, 2004.
- 423 Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul. An introduction
424 to variational methods for graphical models. *Mach. Learn.*, 37(2):183–233, 1999.
- 425 Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the Fokker–Planck
426 equation. *SIAM Journal on Mathematical Analysis*, 29(1):1–17, 1998.

- 427 Simon J. Julier and Jeffrey K. Uhlmann. Unscented filtering and nonlinear estimation. *Proceedings*
428 *of the IEEE*, 92(3):401–422, 2004.
- 429 Simon J. Julier, Jeffrey K. Uhlmann, and Hugh F. Durrant-Whyte. A new method for the nonlinear
430 transformation of means and covariances in filters and estimators. *IEEE Trans. Automat. Control*,
431 45(3):477–482, 2000.
- 432 Mohammad Emtiyaz Khan and Rue Håvard. The Bayesian learning rule. *arXiv:2107.04562*, 2022.
- 433 Marc Lambert, Silvère Bonnabel, and Francis Bach. The limited-memory recursive variational
434 Gaussian approximation (L-RVGA). *hal-03501920*, 2021.
- 435 Marc Lambert, Silvère Bonnabel, and Francis Bach. The recursive variational Gaussian approximation
436 (R-VGA). *Statistics and Computing*, 32(1):10, 2022a.
- 437 Marc Lambert, Silvère Bonnabel, and Francis Bach. The continuous-discrete variational Kalman
438 filter (CD-VKF). *hal-03665666*, 2022b.
- 439 Yin Tat Lee, Ruoqi Shen, and Kevin Tian. Structured logconcave sampling with a restricted Gaussian
440 oracle. In *Proceedings of the Conference on Learning Theory*, volume 134, pages 2993–3050,
441 15–19 Aug 2021.
- 442 Wu Lin, Mohammad E. Khan, and Mark Schmidt. Fast and simple natural-gradient variational
443 inference with mixture of exponential-family approximations. In *Proceedings of the International*
444 *Conference on Machine Learning*, volume 97, pages 3992–4002, 09–15 Jun 2019.
- 445 Qiang Liu. Stein variational gradient descent as gradient flow. In *Advances in Neural Information*
446 *Processing Systems*, volume 30, 2017.
- 447 Qiang Liu and Dilin Wang. Stein variational gradient descent: a general purpose Bayesian inference
448 algorithm. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- 449 Yi-An Ma, Niladri S. Chatterji, Xiang Cheng, Nicolas Flammarion, Peter L. Bartlett, and Michael I.
450 Jordan. Is there an analog of Nesterov acceleration for gradient-based MCMC? *Bernoulli*, 27(3):
451 1942 – 1992, 2021.
- 452 Martin Morf, Bernard Levy, and Thomas Kailath. Square-root algorithms for the continuous-time
453 linear least squares estimation problem. In *1977 IEEE Conference on Decision and Control*
454 *including the 16th Symposium on Adaptive Processes and A Special Symposium on Fuzzy Set*
455 *Theory and Applications*, pages 944–947, 1977.
- 456 Manfred Opper and Cédric Archambeau. The variational Gaussian approximation revisited. *Neural*
457 *Comput.*, 21(3):786–792, 2009.
- 458 Felix Otto. Dynamics of labyrinthine pattern formation in magnetic fluids: a mean-field theory. *Arch.*
459 *Rational Mech. Anal.*, 141(1):63–103, 1998.
- 460 Felix Otto. The geometry of dissipative evolution equations: the porous medium equation. *Comm.*
461 *Partial Differential Equations*, 26(1-2):101–174, 2001.
- 462 John Paisley, David M. Blei, and Michael I. Jordan. Variational Bayesian inference with stochastic
463 search. In *Proceedings of the International Conference on Machine Learning*, pages 1363–1370,
464 2012.
- 465 Gabriel Peyré and Marco Cuturi. *Computational optimal transport: with applications to data science*.
466 Now, 2019.
- 467 Rajesh Ranganath, Sean Gerrish, and David M. Blei. Black box variational inference. In *Proceedings*
468 *of International Conference on Artificial Intelligence and Statistics*, volume 33, pages 814–822,
469 Reykjavik, Iceland, 22–25 Apr 2014.

- 470 Filippo Santambrogio. *Optimal transport for applied mathematicians*, volume 87 of *Progress*
471 *in Nonlinear Differential Equations and their Applications*. Birkhäuser/Springer, Cham, 2015.
472 Calculus of variations, PDEs, and modeling.
- 473 Simo Särkkä. On unscented Kalman filtering for state estimation of continuous-time nonlinear
474 systems. *IEEE Trans. Automat. Control*, 52(9):1631–1641, 2007.
- 475 Matthias Seeger. Bayesian model selection for support vector machines, Gaussian processes and
476 other kernel classifiers. In *Advances in Neural Information Processing Systems*, volume 12, 1999.
- 477 Ruoqi Shen and Yin Tat Lee. The randomized midpoint method for log-concave sampling. In
478 *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 479 Adrian Tudorascu and Marcus Wunsch. On a nonlinear, nonlocal parabolic problem with conservation
480 of mass, mean and variance. *Comm. Partial Differential Equations*, 36(8):1426–1454, 2011.
- 481 Santosh Vempala and Andre Wibisono. Rapid convergence of the unadjusted Langevin algorithm:
482 isoperimetry suffices. In *Advances in Neural Information Processing Systems 32*, pages 8094–8106.
483 2019.
- 484 Cédric Villani. *Topics in optimal transportation*, volume 58 of *Graduate Studies in Mathematics*.
485 American Mathematical Society, Providence, RI, 2003.
- 486 Cédric Villani. *Optimal transport*, volume 338 of *Grundlehren der Mathematischen Wissenschaften*
487 *[Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 2009. Old and new.
- 488 Martin J. Wainwright and Michael I. Jordan. Graphical models, exponential families, and variational
489 inference. *Foundations and Trends in Machine Learning*, 1(1–2):1–305, 2008.
- 490 Andre Wibisono. Sampling as optimization in the space of measures: the Langevin dynamics as a
491 composite optimization problem. In *Proceedings of the 31st Conference On Learning Theory*,
492 volume 75, pages 2093–3027, 2018.
- 493 Keru Wu, Scott Schmidler, and Yuansi Chen. Minimax mixing time of the Metropolis-adjusted
494 Langevin algorithm for log-concave sampling. *arXiv e-prints*, art. arXiv:2109.13055, 2021.
- 495 Lin Wu, Emtiyaz Khan Mohammad, and Schmidt Mark. Stein’s lemma for the reparameterization
496 trick with exponential family mixtures. *arXiv:1910.13398*, 2019.
- 497 Guodong Zhang, Shengyang Sun, David Duvenaud, and Roger Grosse. Noisy natural gradient
498 as variational inference. In *Proceedings of the International Conference on Machine Learning*,
499 volume 80, pages 5852–5861, 2018.

500 Checklist

- 501 1. For all authors...
- 502 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
503 contributions and scope? [Yes]
- 504 (b) Did you describe the limitations of your work? [Yes]
- 505 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 506 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
507 them? [Yes]
- 508 2. If you are including theoretical results...
- 509 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 510 (b) Did you include complete proofs of all theoretical results? [Yes]
- 511 3. If you ran experiments...

- 512 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
513 mental results (either in the supplemental material or as a URL)? [No]
- 514 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
515 were chosen)? [N/A]
- 516 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
517 ments multiple times)? [N/A]
- 518 (d) Did you include the total amount of compute and the type of resources used (e.g., type
519 of GPUs, internal cluster, or cloud provider)? [No]
- 520 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 521 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 522 (b) Did you mention the license of the assets? [N/A]
- 523 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 524
- 525 (d) Did you discuss whether and how consent was obtained from people whose data you're
526 using/curating? [N/A]
- 527 (e) Did you discuss whether the data you are using/curating contains personally identifiable
528 information or offensive content? [N/A]
- 529 5. If you used crowdsourcing or conducted research with human subjects...
- 530 (a) Did you include the full text of instructions given to participants and screenshots, if
531 applicable? [N/A]
- 532 (b) Did you describe any potential participant risks, with links to Institutional Review
533 Board (IRB) approvals, if applicable? [N/A]
- 534 (c) Did you include the estimated hourly wage paid to participants and the total amount
535 spent on participant compensation? [N/A]