
Sliced-Wasserstein Flows: Nonparametric Generative Modeling via Optimal Transport and Diffusions

Antoine Liutkus¹ Umut Şimşekli² Szymon Majewski³ Alain Durmus⁴ Fabian-Robert Stöter¹

Abstract

By building upon the recent theory that established the connection between implicit generative modeling (IGM) and optimal transport, in this study, we propose a novel parameter-free algorithm for learning the underlying distributions of complicated datasets and sampling from them. The proposed algorithm is based on a functional optimization problem, which aims at finding a measure that is close to the data distribution as much as possible and also expressive enough for generative modeling purposes. We formulate the problem as a gradient flow in the space of probability measures. The connections between gradient flows and stochastic differential equations let us develop a computationally efficient algorithm for solving the optimization problem. We provide formal theoretical analysis where we prove finite-time error guarantees for the proposed algorithm. To the best of our knowledge, the proposed algorithm is the first nonparametric IGM algorithm with explicit theoretical guarantees. Our experimental results support our theory and show that our algorithm is able to successfully capture the structure of different types of data distributions.

1. Introduction

Implicit generative modeling (IGM) (Diggle & Gratton, 1984; Mohamed & Lakshminarayanan, 2016) has become very popular recently and has proven successful in various fields; variational auto-encoders (VAE) (Kingma & Welling,

2013) and generative adversarial networks (GAN) (Goodfellow et al., 2014) being its two well-known examples. The goal in IGM can be briefly described as learning the underlying probability measure of a given dataset, denoted as $\nu \in \mathcal{P}(\Omega)$, where \mathcal{P} is the space of probability measures on the measurable space (Ω, \mathcal{A}) , $\Omega \subset \mathbb{R}^d$ is a domain and \mathcal{A} is the associated Borel σ -field.

Given a set of data points $\{y_1, \dots, y_P\}$ that are assumed to be independent and identically distributed (i.i.d.) samples drawn from ν , the implicit generative framework models them as the output of a measurable map, i.e. $y = T(x)$, with $T : \Omega_\mu \mapsto \Omega$. Here, the inputs x are generated from a known and easy to sample source measure μ on Ω_μ (e.g. Gaussian or uniform measures), and the outputs $T(x)$ should match the unknown target measure ν on Ω .

Learning generative networks have witnessed several groundbreaking contributions in recent years. Motivated by this fact, there has been an interest in illuminating the theoretical foundations of VAEs and GANs (Bousquet et al., 2017; Liu et al., 2017). It has been shown that these implicit models have close connections with the theory of Optimal Transport (OT) (Villani, 2008). As it turns out, OT brings new light on the generative modeling problem: there have been several extensions of VAEs (Tolstikhin et al., 2017; Kolouri et al., 2018) and GANs (Arjovsky et al., 2017; Gulrajani et al., 2017; Guo et al., 2017; Lei et al., 2017), which exploit the links between OT and IGM.

OT studies whether it is possible to transform samples from a source distribution μ to a target distribution ν . From this perspective, an ideal generative model is simply a transport map from μ to ν . This can be written by using some ‘push-forward operators’: we seek a mapping T that ‘pushes μ onto ν ’, and is formally defined as $\nu(A) = \mu(T^{-1}(A))$ for all Borel sets $A \subset \mathcal{A}$. If this relation holds, we denote the push-forward operator $T_\#$, such that $T_\#\mu = \nu$. Provided mild conditions on these distributions hold (notably μ is nonatomic (Villani, 2008)), existence of such a transport map is guaranteed; however, it remains a challenge to construct it in practice.

One common point between VAE and GAN is to adopt an approximate strategy and consider transport maps that

¹Inria and LIRMM, Univ. of Montpellier, France
²LTCI, Télécom Paristech, Université Paris-Saclay, Paris, France
³Institute of Mathematics, Polish Academy of Sciences, Warsaw, Poland
⁴CNRS, ENS Paris-Saclay, Université Paris-Saclay, Cachan, France. Correspondence to: Antoine Liutkus <antoine.liutkus@inria.fr>, Umut Şimşekli <umut.simsekli@telecom-paristech.fr>.

belong to a *parametric* family T_ϕ with $\phi \in \Phi$. Then, they aim at finding the best parameter ϕ^* that would give $T_{\phi^* \# \mu} \approx \nu$. This is typically achieved by attempting to minimize the following optimization problem: $\phi^* = \arg \min_{\phi \in \Phi} \mathcal{W}_2(T_{\phi \# \mu}, \nu)$, where \mathcal{W}_2 denotes the Wasserstein distance that will be properly defined in Section 2. It has been shown that (Genevay et al., 2017) OT-based GANs (Arjovsky et al., 2017) and VAEs (Tolstikhin et al., 2017) both use this formulation with different parameterizations and different equivalent definitions of \mathcal{W}_2 . However, their resulting algorithms still lack theoretical understanding.

In this study, we follow a completely different approach for IGM, where we aim at developing an algorithm with explicit theoretical guarantees for estimating a transport map between source μ and target ν . The generated transport map will be *nonparametric* (in the sense that it does not belong to some family of functions, like a neural network), and it will be iteratively augmented: always increasing the quality of the fit along iterations. Formally, we take T_t as the constructed transport map at time $t \in [0, \infty)$, and define $\mu_t = T_t \# \mu$ as the corresponding output distribution. Our objective is to build the maps so that μ_t will converge to the solution of a functional optimization problem, defined through a gradient flow in the Wasserstein space. Informally, we will consider a gradient flow that has the following form:

$$\partial_t \mu_t = -\nabla_{\mathcal{W}_2} \left\{ \text{Cost}(\mu_t, \nu) + \text{Reg}(\mu_t) \right\}, \quad \mu_0 = \mu, \quad (1)$$

where the functional Cost computes a discrepancy between μ_t and ν , Reg denotes a regularization functional, and $\nabla_{\mathcal{W}_2}$ denotes a notion of gradient with respect to a probability measure in the \mathcal{W}_2 metric for probability measures¹. If this flow can be simulated, one would hope for $\mu_t = (T_t) \# \mu$ to converge to the minimum of the functional optimization problem: $\min_{\mu} (\text{Cost}(\mu, \nu) + \text{Reg}(\mu))$ (Ambrosio et al., 2008; Santambrogio, 2017).

We construct a gradient flow where we choose the Cost functional as the *sliced Wasserstein distance* (\mathcal{SW}_2) (Rabin et al., 2012; Bonneel et al., 2015) and the Reg functional as the negative entropy. The \mathcal{SW}_2 distance is equivalent to the \mathcal{W}_2 distance (Bonnotte, 2013) and has important computational implications since it can be expressed as an average of (one-dimensional) projected optimal transportation costs whose analytical expressions are available.

We first show that, with the choice of \mathcal{SW}_2 and the negative-entropy functionals as the overall objective, we obtain a valid gradient flow that has a solution path $(\mu_t)_t$, and the probability density functions of this path solve a particular

¹This gradient flow is similar to the usual Euclidean gradient flows, i.e. $\partial_t x_t = -\nabla(f(x_t) + r(x_t))$, where f is typically the data-dependent cost function and r is a regularization term. The (explicit) Euler discretization of this flow results in the well-known gradient descent algorithm for solving $\min_x (f(x) + r(x))$.

partial differential equation, which has close connections with stochastic differential equations. Even though gradient flows in Wasserstein spaces cannot be solved in general, by exploiting this connection, we are able to develop a practical algorithm that provides approximate solutions to the gradient flow and is algorithmically similar to stochastic gradient Markov Chain Monte Carlo (MCMC) methods² (Welling & Teh, 2011; Ma et al., 2015; Durmus et al., 2016; Şimşekli, 2017; Şimşekli et al., 2018). We provide finite-time error guarantees for the proposed algorithm and show explicit dependence of the error to the algorithm parameters.

To the best of our knowledge, the proposed algorithm is the first nonparametric IGM algorithm that has explicit theoretical guarantees. In addition to its nice theoretical properties, the proposed algorithm has also significant practical importance: it has low computational requirements and can be easily run on an everyday laptop CPU. Our experiments on both synthetic and real datasets support our theory and illustrate the advantages of the algorithm in several scenarios.

2. Technical Background

2.1. Wasserstein distance, optimal transport maps and Kantorovich potentials

For two probability measures $\mu, \nu \in \mathcal{P}_2(\Omega)$, $\mathcal{P}_2(\Omega) = \{\mu \in \mathcal{P}(\Omega) : \int_{\Omega} \|x\|^2 \mu(dx) < +\infty\}$, the 2-Wasserstein distance is defined as follows:

$$\mathcal{W}_2(\mu, \nu) \triangleq \left\{ \inf_{\gamma \in \mathcal{C}(\mu, \nu)} \int_{\Omega \times \Omega} \|x - y\|^2 \gamma(dx, dy) \right\}^{1/2}, \quad (2)$$

where $\mathcal{C}(\mu, \nu)$ is called the set of *transportation plans* and defined as the set of probability measures γ on $\Omega \times \Omega$ satisfying for all $A \in \mathcal{A}$, $\gamma(A \times \Omega) = \mu(A)$ and $\gamma(\Omega \times A) = \nu(A)$, i.e. the marginals of γ coincide with μ and ν . From now on, we will assume that Ω is a compact subset of \mathbb{R}^d .

In the case where Ω is finite, computing the Wasserstein distance between two probability measures turns out to be a linear program with linear constraints, and has therefore a dual formulation. Since Ω is a Polish space (i.e. a complete and separable metric space), this dual formulation can be generalized as follows (Villani, 2008)[Theorem 5.10]:

$$\mathcal{W}_2(\mu, \nu) = \sup_{\psi \in L^1(\mu)} \left\{ \int_{\Omega} \psi(x) \mu(dx) + \int_{\Omega} \psi^c(x) \nu(dx) \right\}^{1/2} \quad (3)$$

where $L^1(\mu)$ denotes the class of functions that are absolutely integrable under μ and ψ^c denotes the c-conjugate of ψ and is defined as follows: $\psi^c(y) \triangleq \{\inf_{x \in \Omega} \|x -$

²We note that, despite the algorithmic similarities, the proposed algorithm is not a Bayesian posterior sampling algorithm.

$y\|^2 - \psi(x)\}$. The functions ψ that realize the supremum in (3) are called the Kantorovich potentials between μ and ν . Provided that μ satisfies a mild condition, we have the following uniqueness result.

Theorem 1 ((Santambrogio, 2014)[Theorem 1.4]). *Assume that $\mu \in \mathcal{P}_2(\Omega)$ is absolutely continuous with respect to the Lebesgue measure. Then, there exists a unique optimal transport plan γ^* that realizes the infimum in (2) and it is of the form $(Id \times T)_{\#}\mu$, for a measurable function $T : \Omega \rightarrow \Omega$. Furthermore, there exists at least a Kantorovich potential ψ whose gradient $\nabla\psi$ is uniquely determined μ -almost everywhere. The function T and the potential ψ are linked by $T(x) = x - \nabla\psi(x)$.*

The measurable function $T : \Omega \rightarrow \Omega$ is referred to as the optimal transport map from μ to ν . This result implies that there exists a solution for transporting samples from μ to samples from ν and this solution is optimal in the sense that it minimizes the ℓ_2 displacement. However, identifying this solution is highly non-trivial. In the discrete case, effective solutions have been proposed (Cuturi, 2013). However, for continuous and high-dimensional probability measures, constructing an actual transport plan remains a challenge. Even if recent contributions (Genevay et al., 2016) have made it possible to rapidly compute \mathcal{W}_2 , they do so without constructing the optimal map T , which is our objective here.

2.2. Wasserstein spaces and gradient flows

By (Ambrosio et al., 2008)[Proposition 7.1.5], \mathcal{W}_2 is a distance over $\mathcal{P}(\Omega)$. In addition, if $\Omega \subset \mathbb{R}^d$ is compact, the topology associated with \mathcal{W}_2 is equivalent to the weak convergence of probability measures and $(\mathcal{P}(\Omega), \mathcal{W}_2)^3$ is compact. The metric space $(\mathcal{P}_2(\Omega), \mathcal{W}_2)$ is called the *Wasserstein space*.

In this study, we are interested in functional optimization problems in $(\mathcal{P}_2(\Omega), \mathcal{W}_2)$, such as $\min_{\mu \in \mathcal{P}_2(\Omega)} \mathcal{F}(\mu)$, where \mathcal{F} is the functional that we would like to minimize. Similar to Euclidean spaces, one way to formulate this optimization problem is to construct a gradient flow of the form $\partial_t \mu_t = -\nabla_{\mathcal{W}_2} \mathcal{F}(\mu_t)$ (Benamou & Brenier, 2000; Lavenant et al., 2018), where $\nabla_{\mathcal{W}_2}$ denotes a notion of gradient in $(\mathcal{P}_2(\Omega), \mathcal{W}_2)$. If such a flow can be constructed, one can utilize it both for practical algorithms and theoretical analysis.

Gradient flows $\partial_t \mu_t = \nabla_{\mathcal{W}_2} \mathcal{F}(\mu_t)$ with respect to a functional \mathcal{F} in $(\mathcal{P}_2(\Omega), \mathcal{W}_2)$ have strong connections with partial differential equations (PDE) that are of the form of a *continuity equation* (Santambrogio, 2017). Indeed, it is shown that under appropriate conditions on \mathcal{F} (see e.g. (Ambrosio et al., 2008)), $(\mu_t)_t$ is a solution of the gradient flow if and only if it admits a density ρ_t with respect to the Lebesgue measure for all $t \geq 0$, and solves the continuity equation

given by: $\partial_t \rho_t + \text{div}(v \rho_t) = 0$, where v denotes a vector field and div denotes the divergence operator. Then, for a given gradient flow in $(\mathcal{P}_2(\Omega), \mathcal{W}_2)$, we are interested in the evolution of the densities ρ_t , i.e. the PDEs which they solve. Such PDEs are of our particular interest since they have a key role for building practical algorithms.

2.3. Sliced-Wasserstein distance

In the one-dimensional case, i.e. $\mu, \nu \in \mathcal{P}_2(\mathbb{R})$, \mathcal{W}_2 has an analytical form, given as follows: $\mathcal{W}_2(\mu, \nu) = \int_0^1 |F_\mu^{-1}(\tau) - F_\nu^{-1}(\tau)|^2 d\tau$, where F_μ and F_ν denote the cumulative distribution functions (CDF) of μ and ν , respectively, and F_μ^{-1}, F_ν^{-1} denote the inverse CDFs, also called quantile functions (QF). In this case, the optimal transport map from μ to ν has a closed-form formula as well, given as follows: $T(x) = (F_\nu^{-1} \circ F_\mu)(x)$ (Villani, 2008). The optimal map T is also known as the *increasing arrangement*, which maps each quantile of μ to the same quantile of ν , e.g. minimum to minimum, median to median, maximum to maximum (Villani, 2008). Due to Theorem 1, the derivative of the corresponding Kantorovich potential is given as:

$$\psi'(x) \triangleq \partial_x \psi(x) = x - (F_\nu^{-1} \circ F_\mu)(x).$$

In the multidimensional case $d > 1$, building a transport map is much more difficult. The nice properties of the one-dimensional Wasserstein distance motivate the usage of *sliced-Wasserstein distance* (\mathcal{SW}_2) for practical applications. Before formally defining \mathcal{SW}_2 , let us first define the orthogonal projection $\theta^*(x) \triangleq \langle \theta, x \rangle$ for any direction $\theta \in \mathbb{S}^{d-1}$ and $x \in \mathbb{R}^d$, where $\langle \cdot, \cdot \rangle$ denotes the Euclidean inner-product and $\mathbb{S}^{d-1} \subset \mathbb{R}^d$ denotes the d -dimensional unit sphere. Then, the \mathcal{SW}_2 distance is formally defined as follows:

$$\mathcal{SW}_2(\mu, \nu) \triangleq \int_{\mathbb{S}^{d-1}} \mathcal{W}_2(\theta_{\#}^* \mu, \theta_{\#}^* \nu) d\theta, \quad (4)$$

where $d\theta$ represents the uniform probability measure on \mathbb{S}^{d-1} . As shown in (Bonnotte, 2013), \mathcal{SW}_2 is indeed a distance metric and induces the same topology as \mathcal{W}_2 for compact domains.

The \mathcal{SW}_2 distance has important practical implications: provided that the projected distributions $\theta_{\#}^* \mu$ and $\theta_{\#}^* \nu$ can be computed, then for any $\theta \in \mathbb{S}^{d-1}$, the distance $\mathcal{W}_2(\theta_{\#}^* \mu, \theta_{\#}^* \nu)$, as well as its optimal transport map and the corresponding Kantorovich potential can be analytically computed (since the projected measures are one-dimensional). Therefore, one can easily approximate (4) by using a simple Monte Carlo scheme that draws uniform random samples from \mathbb{S}^{d-1} and replaces the integral in (4) with a finite-sample average. Thanks to its computational benefits, \mathcal{SW}_2 was very recently considered for OT-based VAEs and GANs (Deshpande et al., 2018; Wu et al., 2018;

³Note that in that case, $\mathcal{P}_2(\Omega) = \mathcal{P}(\Omega)$

Kolouri et al., 2018), appearing as a stable alternative to the adversarial methods.

3. Regularized Sliced-Wasserstein Flows for Generative Modeling

3.1. Construction of the gradient flow

In this paper, we propose the following functional minimization problem on $\mathcal{P}_2(\Omega)$ for implicit generative modeling:

$$\min_{\mu} \left\{ \mathcal{F}_{\lambda}^{\nu}(\mu) \triangleq \frac{1}{2} \mathcal{SW}_2^2(\mu, \nu) + \lambda \mathcal{H}(\mu) \right\}, \quad (5)$$

where $\lambda > 0$ is a regularization parameter and \mathcal{H} denotes the negative entropy defined by $\mathcal{H}(\mu) \triangleq \int_{\Omega} \rho(x) \log \rho(x) dx$ if μ has density ρ with respect to the Lebesgue measure and $\mathcal{H}(\mu) = +\infty$ otherwise. Note that the case $\lambda = 0$ has been already proposed and studied in (Bonnotte, 2013) in a more general OT context. Here, in order to introduce the necessary noise inherent to generative model, we suggest to penalize the slice-Wasserstein distance using \mathcal{H} . In other words, the main idea is to find a measure μ^* that is close to ν as much as possible and also has a certain amount of entropy to make sure that it is sufficiently expressive for generative modeling purposes. The importance of the entropy regularization becomes prominent in practical applications where we have finitely many data samples that are assumed to be drawn from ν . In such a circumstance, the regularization would prevent μ^* to collapse on the data points and therefore avoid ‘over-fitting’ to the data distribution. Note that this regularization is fundamentally different from the one used in Sinkhorn distances (Genevay et al., 2018).

In our first result, we show that there exists a flow $(\mu_t)_{t \geq 0}$ in $(\mathcal{P}(\bar{\mathbb{B}}(0, r)), \mathcal{W}_2)$ which decreases along $\mathcal{F}_{\lambda}^{\nu}$, where $\bar{\mathbb{B}}(0, a)$ denotes the closed unit ball centered at 0 and radius a . This flow will be referred to as a generalized minimizing movement scheme (see Definition 1 in the supplementary document). In addition, the flow $(\mu_t)_{t \geq 0}$ admits a density ρ_t with respect to the Lebesgue measure for all $t > 0$ and $(\rho_t)_{t \geq 0}$ is solution of a non-linear PDE (in the weak sense).

Theorem 2. *Let ν be a probability measure on $\bar{\mathbb{B}}(0, 1)$ with a strictly positive smooth density. Choose a regularization constant $\lambda > 0$ and radius $r > \sqrt{d}$, where d is the data dimension. Assume that $\mu_0 \in \mathcal{P}(\bar{\mathbb{B}}(0, r))$ is absolutely continuous with respect to the Lebesgue measure with density $\rho_0 \in L^{\infty}(\bar{\mathbb{B}}(0, r))$. There exists a generalized minimizing movement scheme $(\mu_t)_{t \geq 0}$ associated to (5) and if ρ_t stands for the density of μ_t for all $t \geq 0$, then $(\rho_t)_t$ satisfies the following continuity equation:*

$$\frac{\partial \rho_t}{\partial t} = -\operatorname{div}(v_t \rho_t) + \lambda \Delta \rho_t, \quad (6)$$

$$v_t(x) \triangleq v(x, \mu_t) = - \int_{\mathbb{S}^{d-1}} \psi'_{t, \theta}(\langle x, \theta \rangle) \theta d\theta \quad (7)$$

in a weak sense. Here, Δ denotes the Laplacian operator, div the divergence operator, and $\psi_{t, \theta}$ denotes the Kantorovich potential between $\theta_{\#}^* \mu_t$ and $\theta_{\#}^* \nu$.

The precise statement of this Theorem, related results and its proof are postponed to the supplementary document. For its proof, we use the technique introduced in (Jordan et al., 1998): we first prove the existence of a generalized minimizing movement scheme by showing that the solution curve $(\mu_t)_t$ is a limit of the solution of a time-discretized problem. Then we prove that the curve $(\rho_t)_t$ solves the PDE given in (6).

3.2. Connection with stochastic differential equations

As a consequence of the entropy regularization, we obtain the Laplacian operator Δ in the PDE given in (6). We therefore observe that the overall PDE is a Fokker-Planck-type equation (Bogachev et al., 2015) that has a well-known probabilistic counterpart, which can be expressed as a stochastic differential equation (SDE). More precisely, let us consider a stochastic process $(X_t)_t$, that is the solution of the following SDE starting at $X_0 \sim \mu_0$:

$$dX_t = v(X_t, \mu_t) dt + \sqrt{2\lambda} dW_t, \quad (8)$$

where $(W_t)_t$ denotes a standard Brownian motion. Then, the probability distribution of X_t at time t solves the PDE given in (6) (Bogachev et al., 2015). This informally means that, if we could simulate (8), then the distribution of X_t would converge to the solution of (5), therefore, we could use the sample paths $(X_t)_t$ as samples drawn from $(\mu_t)_t$. However, in practice this is not possible due to two reasons: (i) the drift v_t cannot be computed analytically since it depends on the probability distribution of X_t , (ii) the SDE (8) is a continuous-time process, it needs to be discretized.

We now focus on the first issue. We observe that the SDE (8) is similar to McKean-Vlasov SDEs (Veretennikov, 2006; Mishura & Veretennikov, 2016), a family of SDEs whose drift depends on the distribution of X_t . By using this connection, we can borrow tools from the relevant SDE literature (Malrieu, 2003; Cattiaux et al., 2008) for developing an approximate simulation method for (8).

Our approach is based on defining a *particle system* that serves as an approximation to the original SDE (8). The particle system can be written as a collection of SDEs, given as follows (Bossy & Talay, 1997):

$$dX_t^i = v(X_t^i, \mu_t^N) dt + \sqrt{2\lambda} dW_t^i, \quad i = 1, \dots, N, \quad (9)$$

where i denotes the particle index, $N \in \mathbb{N}_+$ denotes the total number of particles, and $\mu_t^N = (1/N) \sum_{j=1}^N \delta_{X_t^j}$ denotes the empirical distribution of the particles $\{X_t^j\}_{j=1}^N$. This particle system is particularly interesting, since (i) one

typically has $\lim_{N \rightarrow \infty} \mu_t^N = \mu_t$ with a rate of convergence of order $\mathcal{O}(1/\sqrt{N})$ for all t (Malrieu, 2003; Cattiaux et al., 2008), and (ii) each of the particle systems in (9) can be simulated by using an Euler-Maruyama discretization scheme. We note that the existing theoretical results in (Veretennikov, 2006; Mishura & Veretennikov, 2016) do not directly apply to our case due to the non-standard form of our drift. However, we conjecture that a similar result holds for our problem as well. Such a result would be proven by using the techniques given in (Zhang et al., 2018); however, it is out of the scope of this study.

3.3. Approximate Euler-Maruyama discretization

In order to be able to simulate the particle SDEs (9) in practice, we propose an approximate Euler-Maruyama discretization for each particle SDE. The algorithm iteratively applies the following update equation: ($\forall i \in \{1, \dots, N\}$)

$$\bar{X}_0^i \stackrel{\text{i.i.d.}}{\sim} \mu_0, \quad \bar{X}_{k+1}^i = \bar{X}_k^i + h \hat{v}_k(\bar{X}_k^i) + \sqrt{2\lambda h} Z_{k+1}^i, \quad (10)$$

where $k \in \mathbb{N}_+$ denotes the iteration number, Z_k^i is a standard Gaussian random vector in \mathbb{R}^d , h denotes the step-size, and \hat{v}_k is a short-hand notation for a computationally tractable estimator of the original drift $v(\cdot, \bar{\mu}_{kh}^N)$, with $\bar{\mu}_{kh}^N = (1/N) \sum_{j=1}^N \delta_{\bar{X}_k^j}$ being the empirical distribution of $\{\bar{X}_k^j\}_{j=1}^N$. A question of fundamental practical importance is how to compute this function \hat{v} .

We propose to approximate the integral in (7) via a simple Monte Carlo estimate. This is done by first drawing N_θ uniform i.i.d. samples from the sphere \mathbb{S}^{d-1} , $\{\theta_n\}_{n=1}^{N_\theta}$. Then, at each iteration k , we compute:

$$\hat{v}_k(x) \triangleq -(1/N_\theta) \sum_{n=1}^{N_\theta} \psi'_{k,\theta_n}(\langle \theta_n, x \rangle) \theta_n, \quad (11)$$

where for any θ , $\psi'_{k,\theta}$ is the derivative of the Kantorovich potential (cf. Section 2) that is applied to the OT problem from $\theta_{\#}^* \bar{\mu}_{kh}^N$ to $\theta_{\#}^* \nu$: i.e.

$$\psi'_{k,\theta}(z) = [z - (F_{\theta_{\#}^* \nu}^{-1} \circ F_{\theta_{\#}^* \bar{\mu}_{kh}^N})(z)]. \quad (12)$$

For any particular $\theta \in \mathbb{S}^{d-1}$, the QF, $F_{\theta_{\#}^* \nu}^{-1}$ for the projection of the target distribution ν on θ can be easily computed from the data. This is done by first computing the projections $\langle \theta, y_i \rangle$ for all data points y_i , and then computing the empirical quantile function for this set of P scalars. Similarly, $F_{\theta_{\#}^* \bar{\mu}_{kh}^N}$, the CDF of the particles at iteration k , is easy to compute: we first project all particles \bar{X}_k^i to get $\langle \theta, \bar{X}_k^i \rangle$, and then compute the empirical CDF of this set of N scalar values.

In both cases, the true CDF and quantile functions are approximated as a linear interpolation between a set of the

Algorithm 1: Sliced-Wasserstein Flow (SWF)

```

input :  $\mathcal{D} \equiv \{y_i\}_{i=1}^P, \mu_0, N, N_\theta, h, \lambda$ 
output :  $\{\bar{X}_K^i\}_{i=1}^N$ 
// Initialize the particles
 $\bar{X}_0^i \stackrel{\text{i.i.d.}}{\sim} \mu_0, \quad i = 1, \dots, N$ 
// Generate random directions
 $\theta_n \sim \text{Uniform}(\mathbb{S}^{d-1}), \quad n = 1, \dots, N_\theta$ 
// Quantiles of projected target
for  $\theta \in \{\theta_n\}_{n=1}^{N_\theta}$  do
     $F_{\theta_{\#}^* \nu}^{-1} = \text{QF}\{\langle \theta, y_i \rangle\}_{i=1}^P$ 
// Iterations
for  $k = 0, \dots, K-1$  do
    for  $\theta \in \{\theta_n\}_{n=1}^{N_\theta}$  do
        // CDF of projected particles
         $F_{\theta_{\#}^* \bar{\mu}_{kh}^N} = \text{CDF}\{\langle \theta, \bar{X}_k^i \rangle\}_{i=1}^N$ 
        // Update the particles
         $\bar{X}_{k+1}^i = \bar{X}_k^i - h \hat{v}_k(\bar{X}_k^i) + \sqrt{2\lambda h} Z_{k+1}^i$ 
         $i = 1, \dots, N$ 

```

computed $Q \in \mathbb{N}_+$ empirical quantiles. Another source of approximation here comes from the fact that the target ν will in practice be a collection of Dirac measures on the observations y_i . Since it is currently common to have a very large dataset, we believe this approximation to be accurate in practice for the target. Finally, yet another source of approximation comes from the error induced by using a finite number of θ_n instead of a sum over \mathbb{S}^{d-1} in (12).

Even though the error induced by these approximation schemes can be incorporated into our current analysis framework, we choose to neglect it for now, because (i) all of these one-dimensional computations can be done very accurately and (ii) the quantization of the empirical CDF and QF can be modeled as additive Gaussian noise that enters our discretization scheme (10) (Van der Vaart, 1998). Therefore, we will assume that \hat{v}_k is an *unbiased* estimator of v , i.e. $\mathbb{E}[\hat{v}(x, \mu)] = v(x, \mu)$, for any x and μ , where the expectation is taken over θ_n .

The overall algorithm is illustrated in Algorithm 1. It is remarkable that the updates of the particles only involves the learning data $\{y_i\}$ through the CDFs of its projections on the many $\theta_n \in \mathbb{S}^{d-1}$. This has a fundamental consequence of high practical interest: these CDF may be computed beforehand in a massively distributed manner that is independent of the sliced Wasserstein flow. This aspect is reminiscent of the *compressive learning* methodology (Gribonval et al., 2017), except we exploit quantiles of random projections here, instead of random generalized moments as done there.

Besides, we can obtain further reductions in the computing time if the CDF, $F_{\theta_{\#}^* \nu}$ for the target is computed on random

mini-batches of the data, instead of the whole dataset of size P . This simplified procedure might also have some interesting consequences in privacy-preserving settings: since we can vary the number of projection directions N_θ for each data point y_i , we may guarantee that y_i cannot be recovered via these projections, by picking fewer than necessary for reconstruction using, e.g. compressed sensing (Donoho & Tanner, 2009).

3.4. Finite-time analysis for the infinite particle regime

In this section we will analyze the behavior of the proposed algorithm in the asymptotic regime where the number of particles $N \rightarrow \infty$. Within this regime, we will assume that the original SDE (8) can be directly simulated by using an approximate Euler-Maruyama scheme, defined starting at $\bar{X}_0 \sim \mu_0$ as follows:

$$\bar{X}_{k+1} = \bar{X}_k + h\hat{v}(\bar{X}_k^i, \bar{\mu}_{kh}) + \sqrt{2\lambda h}Z_{k+1}, \quad (13)$$

where $\bar{\mu}_{kh}$ denotes the law of \bar{X}_k with step size h and $\{Z_k\}_k$ denotes a collection of standard Gaussian random variables. Apart from its theoretical significance, this scheme is also practically relevant, since one would expect that it captures the behavior of the particle method (10) with large number of particles.

In practice, we would like to approximate the measure sequence $(\mu_t)_t$ as accurate as possible, where μ_t denotes the law of X_t . Therefore, we are interested in analyzing the distance $\|\bar{\mu}_{Kh} - \mu_T\|_{\text{TV}}$, where K denotes the total number of iterations, $T = Kh$ is called the horizon, and $\|\mu - \nu\|_{\text{TV}}$ denotes the total variation distance between two probability measures μ and ν : $\|\mu - \nu\|_{\text{TV}} \triangleq \sup_{A \in \mathcal{B}(\Omega)} |\mu(A) - \nu(A)|$.

In order to analyze this distance, we exploit the algorithmic similarities between (13) and the stochastic gradient Langevin dynamics (SGLD) algorithm (Welling & Teh, 2011), which is a Bayesian posterior sampling method having a completely different goal, and is obtained as a discretization of an SDE whose drift has a much simpler form. We then bound the distance by extending the recent results on SGLD (Raginsky et al., 2017) to time- and measure-dependent drifts, that are of our interest in the paper.

We now present our second main theoretical result. We present all our assumptions and the explicit forms of the constants in the supplementary document.

Theorem 3. *Assume that the conditions given in the supplementary document hold. Then, the following bound holds for $T = Kh$:*

$$\|\bar{\mu}_{Kh} - \mu_T\|_{\text{TV}}^2 \leq \delta_\lambda \left\{ \frac{L^2 K}{2\lambda} \left(\frac{C_1 h^3}{3} + 3\lambda d h^2 \right) + \frac{C_2 \delta K h}{4\lambda} \right\}, \quad (14)$$

for some $C_1, C_2, L > 0$, $\delta \in (0, 1)$, and $\delta_\lambda > 1$.

Here, the constants C_1, C_2, L are related to the regularity and smoothness of the functions v and \hat{v} ; δ is directly proportional to the variance of \hat{v} , and δ_λ is inversely proportional to λ . The theorem shows that if we choose h small enough, we can have a non-asymptotic error guarantee, which is formally shown in the following corollary.

Corollary 1. *Assume that the conditions of Theorem 3 hold. Then for all $\varepsilon > 0$, $K \in \mathbb{N}_+$, setting*

$$h = (3/C_1) \wedge \left(\frac{2\varepsilon^2 \lambda}{\delta_\lambda L^2 T} (1 + 3\lambda d)^{-1} \right)^{1/2}, \quad (15)$$

we have

$$\|\bar{\mu}_{Kh} - \mu_T\|_{\text{TV}} \leq \varepsilon + \left(\frac{C_2 \delta_\lambda \delta T}{4\lambda} \right)^{1/2} \quad (16)$$

for $T = Kh$.

This corollary shows that for a large horizon T , the approximate drift \hat{v} should have a small variance in order to obtain accurate estimations. This result is similar to (Raginsky et al., 2017) and (Nguyen et al., 2019): for small ε the variance of the approximate drift should be small as well. On the other hand, we observe that the error decreases as λ increases. This behavior is expected since for large λ , the Brownian term in (8) dominates the drift, which makes the simulation easier.

We note that these results establish the explicit dependency of the error with respect to the algorithm parameters (e.g. step-size, gradient noise) for a fixed number of iterations, rather than explaining the asymptotic behavior of the algorithm when K goes to infinity.

4. Experiments

In this section, we evaluate the SWF algorithm on a synthetic and a real data setting. Our primary goal is to validate our theory and illustrate the behavior of our non-standard approach, rather than to obtain the state-of-the-art results in IGM. In all our experiments, the initial distribution μ_0 is selected as the standard Gaussian distribution on \mathbb{R}^d , we take $Q = 100$ quantiles and $N = 5000$ particles, which proved sufficient to approximate the quantile functions accurately.

4.1. Gaussian Mixture Model

We perform the first set of experiments on synthetic data where we consider a standard Gaussian mixture model (GMM) with 10 components and random parameters. Centroids are taken as sufficiently distant from each other to make the problem more challenging. We generate $P = 50000$ data samples in each experiment.

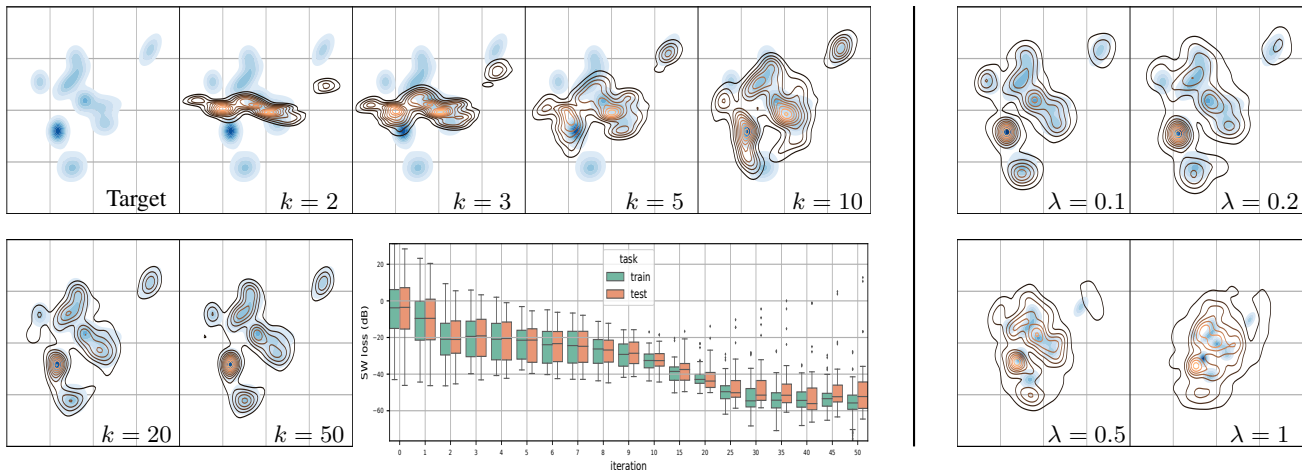


Figure 1. SWF on toy 2D data. **Left:** Target distribution (shaded contour plot) and distribution of particles (lines) during SWF. (bottom) SW cost over iterations during training (left) and test (right) stages. **Right:** Influence of the regularization parameter λ .

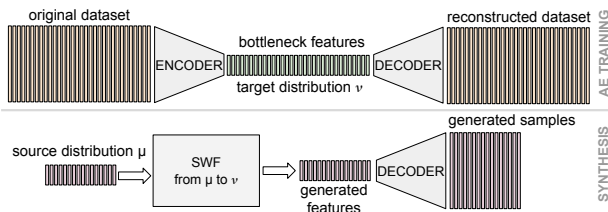


Figure 2. First, we learn an autoencoder (AE). Then, we use SWF to transport random vectors to the distribution of the bottleneck features of the training set. The trained decoder is used for visualization.

In our first experiment, we set $d = 2$ for visualization purposes and illustrate the general behavior of the algorithm. Figure 1 shows the evolution of the particles through the iterations. Here, we set $N_\theta = 30$, $h = 1$ and $\lambda = 10^{-4}$. We first observe that the SW cost between the empirical distributions of training data and particles is steadily decreasing along the SW flow. Furthermore, we see that the QFs, $F_{\theta^*}^{-1}$ that are computed with the initial set of particles (the *training* stage) can be perfectly re-used for new unseen particles in a subsequent *test* stage, yielding similar — yet slightly higher — SW cost.

In our second experiment on Figure 1, we investigate the effect of the level of the regularization λ . The distribution of the particles becomes more spread with increasing λ . This is due to the increment of the entropy, as expected.

4.2. Experiments on real data

In the second set of experiments, we test the SWF algorithm on two real datasets. (i) The traditional MNIST dataset that contains 70K binary images corresponding to different digits. (ii) The popular CelebA dataset (Liu et al., 2015), that



Figure 3. Samples generated after 200 iterations of SWF to match the distribution of bottleneck features for the training dataset. Visualization is done with the pre-trained decoder.

contains 202K color-scale images. This dataset is advocated as more challenging than MNIST. Images were interpolated as 32×32 for MNIST, and 64×64 for CelebA.

In experiments reported in the supplementary document, we found out that directly applying SWF to such high-dimensional data yielded noisy results, possibly due to the insufficient sampling of S^{d-1} . To reduce the dimensionality, we trained a standard convolutional autoencoder (AE) on the training set of both datasets (see Figure 2 and the supplementary document), and the target distribution ν considered becomes the distribution of the resulting bottleneck features, with dimension d . Particles can be visualized with the pre-trained decoder. Our goal is to show that SWF permits to directly sample from the distribution of bottleneck features, as an alternative to enforcing this distribution to

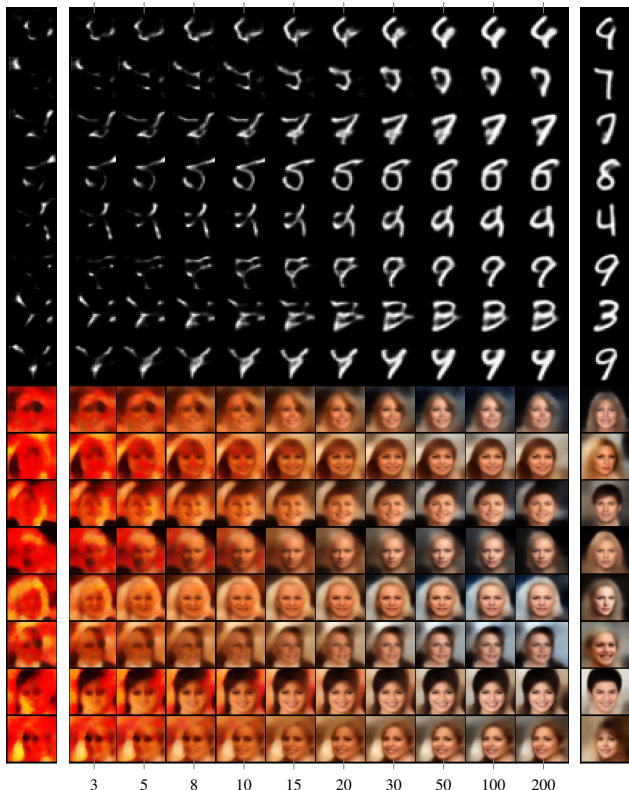


Figure 4. Initial random particles (left), particles through iterations (middle, from 1 to 200 iterations) and closest sample from the training dataset (right), for both MNIST and CelebA.

match some prior, as in VAE. In the following, we set $\lambda = 0$, $N_\theta = 40000$, $d = 32$ for MNIST and $d = 64$ for CelebA.

Assessing the validity of IGM algorithms is generally done by visualizing the generated samples. Figure 3 shows some particles after 500 iterations of SWF. We can observe they are considerably accurate. Interestingly, the generated samples gradually take the form of either digits or faces along the iterations, as seen on Figure 4. In this figure, we also display the closest sample from the original database to check we are not just reproducing training data.

For a visual comparison, we provide the results presented in (Deshpande et al., 2018) in Figure 5. These results are obtained by running different IGM approaches on the MNIST



Figure 5. Performance of GAN (left), W-GAN (middle), SWG (right) on MNIST. (The figure is directly taken from (Deshpande et al., 2018).)



Figure 6. Applying a pre-trained SWF on new samples located in-between the ones used for training. Visualization is done with the pre-trained decoder.

dataset, namely GAN (Goodfellow et al., 2014), Wasserstein GAN (W-GAN) (Arjovsky et al., 2017) and the Sliced-Wasserstein Generator (SWG) (Deshpande et al., 2018). The visual comparison suggests that the samples generated by SWF are of slightly better quality than those, although research must still be undertaken to scale up to high dimensions without an AE.

We also provide the outcome of the pre-trained SWF with samples that are regularly spaced in between those used for training. The result is shown in Figure 4.2. This plot suggests that SWF is a way to interpolate non-parametrically in between latent spaces of regular AE.

5. Conclusion and Future Directions

In this study, we proposed SWF, an efficient, nonparametric IGM algorithm. SWF is based on formulating IGM as a functional optimization problem in Wasserstein spaces, where the aim is to find a probability measure that is close to the data distribution as much as possible while maintaining the expressiveness at a certain level. SWF lies in the intersection of OT, gradient flows, and SDEs, which allowed us to convert the IGM problem to an SDE simulation problem. We provided finite-time bounds for the infinite-particle regime and established explicit links between the algorithm parameters and the overall error. We conducted several experiments, where we showed that the results support our theory: SWF is able to generate samples from non-trivial distributions with low computational requirements.

The SWF algorithm opens up interesting future directions: (i) extension to differentially private settings (Dwork & Roth, 2014) by exploiting the fact that it only requires random projections of the data, (ii) showing the convergence scheme of the particle system (9) to the original SDE (8), (iii) providing bounds directly for the particle scheme (10).

Acknowledgments

This work is partly supported by the French National Research Agency (ANR) as a part of the FBIMATRIX (ANR-16-CE23-0014) and KAMoulox (ANR-15-CE38-0003-01) projects. Szymon Majewski is partially supported by Polish National Science Center grant number 2016/23/B/ST1/00454.

References

- Ambrosio, L., Gigli, N., and Savaré, G. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2008.
- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pp. 214–223, 2017.
- Benamou, J.-D. and Brenier, Y. A computational fluid mechanics solution to the Monge-Kantorovich mass transfer problem. *Numerische Mathematik*, 84(3):375–393, 2000.
- Bogachev, V. I., Krylov, N. V., Röckner, M., and Shaposhnikov, S. V. *Fokker-Planck-Kolmogorov Equations*, volume 207. American Mathematical Soc., 2015.
- Bonneel, N., Rabin, J., Peyré, G., and Pfister, H. Sliced and Radon Wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51(1):22–45, 2015.
- Bonnotte, N. *Unidimensional and evolution methods for optimal transportation*. PhD thesis, Paris 11, 2013.
- Bossy, M. and Talay, D. A stochastic particle method for the McKean-Vlasov and the Burgers equation. *Mathematics of Computation of the American Mathematical Society*, 66(217):157–192, 1997.
- Bousquet, O., Gelly, S., Tolstikhin, I., Simon-Gabriel, C.-J., and Schoelkopf, B. From optimal transport to generative modeling: the vegan cookbook. *arXiv preprint arXiv:1705.07642*, 2017.
- Cattiaux, P., Guillin, A., and Malrieu, F. Probabilistic approach for granular media equations in the non uniformly convex case. *Prob. Theor. Rel. Fields*, 140(1-2):19–40, 2008.
- Şimşekli, U., Yildiz, C., Nguyen, T. H., Cemgil, A. T., and Richard, G. Asynchronous stochastic quasi-Newton MCMC for non-convex optimization. In *ICML*, pp. 4674–4683, 2018.
- Cuturi, M. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pp. 2292–2300, 2013.
- Dalalyan, A. S. Theoretical guarantees for approximate sampling from smooth and log-concave densities. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(3):651–676, 2017.
- Deshpande, I., Zhang, Z., and Schwing, A. Generative modeling using the sliced wasserstein distance. *arXiv preprint arXiv:1803.11188*, 2018.
- Diggle, P. J. and Gratton, R. J. Monte carlo methods of inference for implicit statistical models. *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 193–227, 1984.
- Donoho, D. and Tanner, J. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 367(1906):4273–4293, 2009.
- Durmus, A., Şimşekli, U., Moulines, E., Badeau, R., and Richard, G. Stochastic gradient Richardson-Romberg Markov Chain Monte Carlo. In *NIPS*, 2016.
- Dwork, C. and Roth, A. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- Genevay, A., Cuturi, M., Peyré, G., and Bach, F. Stochastic optimization for large-scale optimal transport. In *Advances in Neural Information Processing Systems*, pp. 3440–3448, 2016.
- Genevay, A., Peyré, G., and Cuturi, M. Gan and vae from an optimal transport point of view. *arXiv preprint arXiv:1706.01807*, 2017.
- Genevay, A., Peyré, G., and Cuturi, M. Learning generative models with Sinkhorn divergences. In *International Conference on Artificial Intelligence and Statistics*, pp. 1608–1617, 2018.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- Gribonval, R., Blanchard, G., Keriven, N., and Traonmilin, Y. Compressive statistical learning with random feature moments. *arXiv preprint arXiv:1706.07180*, 2017.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. C. Improved training of Wasserstein GANs. In *Advances in Neural Information Processing Systems*, pp. 5769–5779, 2017.

- Guo, X., Hong, J., Lin, T., and Yang, N. Relaxed Wasserstein with applications to GANs. *arXiv preprint arXiv:1705.07164*, 2017.
- Jordan, R., Kinderlehrer, D., and Otto, F. The variational formulation of the Fokker–Planck equation. *SIAM journal on mathematical analysis*, 29(1):1–17, 1998.
- Kingma, D. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Kolouri, S., Martin, C. E., and Rohde, G. K. Sliced-wasserstein autoencoder: An embarrassingly simple generative model. *arXiv preprint arXiv:1804.01947*, 2018.
- Lavenant, H., Claiici, S., Chien, E., and Solomon, J. Dynamical optimal transport on discrete surfaces. In *SIGGRAPH Asia 2018 Technical Papers*, pp. 250. ACM, 2018.
- Lei, N., Su, K., Cui, L., Yau, S.-T., and Gu, D. X. A geometric view of optimal transportation and generative model. *arXiv preprint arXiv:1710.05488*, 2017.
- Liu, S., Bousquet, O., and Chaudhuri, K. Approximation and convergence properties of generative adversarial learning. In *Advances in Neural Information Processing Systems*, pp. 5551–5559, 2017.
- Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- Ma, Y. A., Chen, T., and Fox, E. A complete recipe for stochastic gradient MCMC. In *Advances in Neural Information Processing Systems*, pp. 2899–2907, 2015.
- Malrieu, F. Convergence to equilibrium for granular media equations and their Euler schemes. *Ann. Appl. Probab.*, 13(2):540–560, 2003.
- Mishura, Y. S. and Veretennikov, A. Y. Existence and uniqueness theorems for solutions of McKean–Vlasov stochastic equations. *arXiv preprint arXiv:1603.02212*, 2016.
- Mohamed, S. and Lakshminarayanan, B. Learning in implicit generative models. *arXiv preprint arXiv:1610.03483*, 2016.
- Nguyen, T. H., Şimşekli, U., , and Richard, G. Non-asymptotic analysis of fractional Langevin Monte Carlo for non-convex optimization. In *ICML*, 2019.
- Rabin, J., Peyré, G., Delon, J., and Bernot, M. Wasserstein barycenter and its application to texture mixing. In Bruckstein, A. M., ter Haar Romeny, B. M., Bronstein, A. M., and Bronstein, M. M. (eds.), *Scale Space and Variational Methods in Computer Vision*, pp. 435–446. Berlin, Heidelberg, 2012. Springer Berlin Heidelberg. ISBN 978-3-642-24785-9.
- Raginsky, M., Rakhlin, A., and Telgarsky, M. Non-convex learning via stochastic gradient Langevin dynamics: a nonasymptotic analysis. In *Proceedings of the 2017 Conference on Learning Theory*, volume 65, pp. 1674–1703, 2017.
- Samangouei, P., Kabkab, M., and Chellappa, R. Defense-GAN: Protecting classifiers against adversarial attacks using generative models. In *International Conference on Learning Representations*, 2018.
- Santambrogio, F. Introduction to optimal transport theory. In Pajot, H., Ollivier, Y., and Villani, C. (eds.), *Optimal Transportation: Theory and Applications*, chapter 1. Cambridge University Press, 2014.
- Santambrogio, F. {Euclidean, metric, and Wasserstein} gradient flows: an overview. *Bulletin of Mathematical Sciences*, 7(1):87–154, 2017.
- Şimşekli, U. Fractional Langevin Monte Carlo: Exploring Lévy Driven Stochastic Differential Equations for Markov Chain Monte Carlo. In *International Conference on Machine Learning*, 2017.
- Tolstikhin, I., Bousquet, O., Gelly, S., and Schoelkopf, B. Wasserstein auto-encoders. *arXiv preprint arXiv:1711.01558*, 2017.
- Van der Vaart, A. W. *Asymptotic statistics*, volume 3. Cambridge university press, 1998.
- Veretennikov, A. Y. On ergodic measures for McKean–Vlasov stochastic equations. In *Monte Carlo and Quasi-Monte Carlo Methods 2004*, pp. 471–486. Springer, 2006.
- Villani, C. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- Welling, M. and Teh, Y. W. Bayesian learning via stochastic gradient Langevin dynamics. In *International Conference on Machine Learning*, pp. 681–688, 2011.
- Wu, J., Huang, Z., Li, W., and Gool, L. V. Sliced wasserstein generative models. *arXiv preprint arXiv:1706.02631*, abs/1706.02631, 2018.
- Zhang, J., Zhang, R., and Chen, C. Stochastic particle-optimization sampling and the non-asymptotic convergence theory. *arXiv preprint arXiv:1809.01293*, 2018.

Sliced-Wasserstein Flows: Nonparametric Generative Modeling via Optimal Transport and Diffusions

SUPPLEMENTARY DOCUMENT

Antoine Liutkus¹ Umut Şimşekli² Szymon Majewski³ Alain Durmus⁴ Fabian-Robert Stöter¹

1. Proof of Theorem 2

We first need to generalize (Bonnotte, 2013)[Lemma 5.4.3] to distribution $\rho \in L^\infty(\bar{B}(0, r))$, $r > 0$.

Theorem S4. *Let ν be a probability measure on $\bar{B}(0, 1)$ with a strictly positive smooth density. Fix a time step $h > 0$, regularization constant $\lambda > 0$ and a radius $r > \sqrt{d}$. For any probability measure μ_0 on $\bar{B}(0, r)$ with density $\rho_0 \in L^\infty(\bar{B}(0, r))$, there is a probability measure μ on $\bar{B}(0, r)$ minimizing:*

$$\mathcal{G}(\mu) = \mathcal{F}_\lambda^\nu(\mu) + \frac{1}{2h} \mathcal{W}_2^2(\mu, \mu_0),$$

where \mathcal{F}_λ^ν is given by (5). Moreover the optimal μ has a density ρ on $\bar{B}(0, r)$ and:

$$\|\rho\|_{L^\infty} \leq (1 + h/\sqrt{d})^d \|\rho_0\|_{L^\infty}. \quad (\text{S1})$$

Proof. The set of measures supported on $\bar{B}(0, r)$ is compact in the topology given by \mathcal{W}_2 metric. Furthermore by (Ambrosio et al., 2008)[Lemma 9.4.3] \mathcal{H} is lower semicontinuous on $(\mathcal{P}(\bar{B}(0, r)), \mathcal{W}_2)$. Since by (Bonnotte, 2013)[Proposition 5.1.2, Proposition 5.1.3], \mathcal{SW}_2 is a distance on $\mathcal{P}(\bar{B}(0, r))$, dominated by $d^{-1/2}\mathcal{W}_2$, we have:

$$|\mathcal{SW}_2(\pi_0, \nu) - \mathcal{SW}_2(\pi_1, \nu)| \leq \mathcal{SW}_2(\pi_0, \pi_1) \leq \frac{1}{\sqrt{d}} \mathcal{W}_2(\pi_0, \pi_1).$$

The above means that $\mathcal{SW}_2(\cdot, \nu)$ is continuous with respect to topology given by \mathcal{W}_2 , which implies that $\mathcal{SW}_2^2(\cdot, \nu)$ is continuous in this topology as well. Therefore $\mathcal{G} : \mathcal{P}(\bar{B}(0, r)) \rightarrow (-\infty, +\infty]$ is a lower semicontinuous function on the compact set $(\mathcal{P}(\bar{B}(0, r)), \mathcal{W}_2)$. Hence there exists a minimum μ of \mathcal{G} on $\mathcal{P}(\bar{B}(0, r))$. Furthermore, since $\mathcal{H}(\pi) = +\infty$ for measures π that do not admit a density with respect to Lebesgue measure, the measure μ must admit a density ρ .

If ρ_0 is smooth and positive on $\bar{B}(0, r)$, the inequality S1 is true by (Bonnotte, 2013)[Lemma 5.4.3.] When ρ_0 is just in $L^\infty(\bar{B}(0, r))$, we proceed by smoothing. For $t \in (0, 1]$, let ρ_t be a function obtained by convolution of ρ_0 with a Gaussian kernel $(t, x, y) \mapsto (2\pi)^{d/2} \exp(-\|x - y\|^2 / 2t)$, restricting the result to $\bar{B}(0, r)$ and normalizing to obtain a probability density. Then $(\rho_t)_t$ are smooth positive densities, and it is easy to see that $\lim_{t \rightarrow 0} \|\rho_t\|_{L^\infty} \leq \|\rho_0\|_{L^\infty}$. Furthermore, if we denote by μ_t the measure on $\bar{B}(0, r)$ with density ρ_t , then μ_t converge weakly to μ_0 . For $t \in (0, 1]$ let $\hat{\mu}_t$ be the minimum of $\mathcal{F}_\lambda^\nu(\cdot) + \frac{1}{2h} \mathcal{W}_2^2(\cdot, \mu_t)$, and let $\hat{\rho}_t$ be the density of $\hat{\mu}_t$. Using (Bonnotte, 2013)[Lemma 5.4.3.] we get

$$\|\hat{\rho}_t\|_{L^\infty} \leq (1 + h\sqrt{d})^d \|\rho_t\|_{L^\infty}.$$

so $\hat{\rho}_t$ lies in a ball of finite radius in L^∞ . Using compactness of $\mathcal{P}(\bar{B}(0, r))$ in weak topology and compactness of closed ball in $L^\infty(\bar{B}(0, r))$ in weak star topology, we can choose a subsequence $\hat{\mu}_{t_k}, \hat{\rho}_{t_k}$, $\lim_{k \rightarrow +\infty} t_k = 0$, that converges along that subsequence to limits $\hat{\mu}, \hat{\rho}$. Obviously $\hat{\rho}$ is the density of $\hat{\mu}$, since for any continuous function f on $\bar{B}(0, r)$ we have:

$$\int \hat{\rho} f dx = \lim_{k \rightarrow \infty} \int \rho_{t_k} f dx = \lim_{k \rightarrow \infty} \int f d\mu_{t_k} = \int f d\mu.$$

Furthermore, since $\hat{\rho}$ is the weak star limit of a bounded subsequence, we have:

$$\|\hat{\rho}\|_{L^\infty} \leq \limsup_{k \rightarrow \infty} (1 + h\sqrt{d})^d \|\rho_{t_k}\|_{L^\infty} \leq (1 + h\sqrt{d})^d \|\rho_0\|_{L^\infty}.$$

To finish, we just need to prove that $\hat{\mu}$ is a minimum of \mathcal{G} . We remind our reader, that we already established existence of some minimum μ (that might be different from $\hat{\mu}$). Since $\hat{\mu}_{t_k}$ converges weakly to $\hat{\mu}$ in $\mathcal{P}(\overline{\mathbb{B}}(0, r))$, it implies convergence in \mathcal{W}_2 as well since $\overline{\mathbb{B}}(0, r)$ is compact. Similarly μ_{t_k} converges to μ_0 in \mathcal{W}_2 . Using the lower semicontinuity of \mathcal{G} we now have:

$$\begin{aligned} \mathcal{F}_\lambda^\nu(\hat{\mu}) + \frac{1}{2h} \mathcal{W}_2^2(\hat{\mu}, \mu_0) &\leq \liminf_{k \rightarrow \infty} \left(\mathcal{F}_\lambda^\nu(\hat{\mu}_{t_k}) + \frac{1}{2h} \mathcal{W}_2^2(\hat{\mu}_{t_k}, \mu_0) \right) \\ &\leq \liminf_{k \rightarrow \infty} \mathcal{F}_\lambda^\nu(\mu) + \frac{1}{2h} \mathcal{W}_2^2(\mu, \mu_{t_k}) \\ &\quad + \frac{1}{2h} \mathcal{W}_2^2(\hat{\mu}_{t_k}, \mu_0) - \frac{1}{2h} \mathcal{W}_2^2(\hat{\mu}_{t_k}, \mu_{t_k}) \\ &= \mathcal{F}_\lambda^\nu(\mu) + \frac{1}{2h} \mathcal{W}_2^2(\mu, \mu_0), \end{aligned}$$

where the second inequality comes from the fact, that $\hat{\mu}_{t_k}$ minimizes $\mathcal{F}_\lambda^\nu(\cdot) + \frac{1}{2h} \mathcal{W}_2^2(\cdot, \mu_{t_k})$. From the above inequality and previously established facts, it follows that $\hat{\mu}$ is a minimum of \mathcal{G} with density satisfying **S1**. \square

Definition 1. Minimizing movement scheme Let $r > 0$ and $\mathcal{F} : \mathbb{R}_+ \times \mathcal{P}(\overline{\mathbb{B}}(0, r)) \times \mathcal{P}(\overline{\mathbb{B}}(0, r)) \rightarrow \mathbb{R}$ be a functional. Let $\mu_0 \in \mathcal{P}(\overline{\mathbb{B}}(0, r))$ be a starting point. For $h > 0$ a piecewise constant trajectory $\mu^h : [0, \infty) \rightarrow \mathcal{P}(\overline{\mathbb{B}}(0, r))$ for \mathcal{F} starting at μ_0 is a function such that:

- $\mu^h(0) = \mu_0$.
- μ^h is constant on each interval $[nh, (n+1)h)$, so $\mu^h(t) = \mu^h(nh)$ with $n = \lfloor t/h \rfloor$.
- $\mu^h((n+1)h)$ minimizes the functional $\zeta \mapsto \mathcal{F}(h, \zeta, \mu^h(nh))$, for all $n \in \mathbb{N}$.

We say $\hat{\mu}$ is a minimizing movement scheme for \mathcal{F} starting at μ_0 , if there exists a family of piecewise constant trajectory $(\mu^h)_{h>0}$ for \mathcal{F} such that $\hat{\mu}$ is a pointwise limit of μ^h as h goes to 0, i.e. for all $t \in \mathbb{R}_+$, $\lim_{h \rightarrow 0} \mu^h(t) = \mu(t)$ in $\mathcal{P}(\overline{\mathbb{B}}(0, r))$. We say that $\tilde{\mu}$ is a generalized minimizing movement for \mathcal{F} starting at μ_0 , if there exists a family of piecewise constant trajectory $(\mu^h)_{h>0}$ for \mathcal{F} and a sequence $(h_n)_n$, $\lim_{n \rightarrow \infty} h_n = 0$, such that μ^{h_n} converges pointwise to $\tilde{\mu}$.

Theorem S5. Let ν be a probability measure on $\overline{\mathbb{B}}(0, 1)$ with a strictly positive smooth density. Fix a regularization constant $\lambda > 0$ and radius $r > \sqrt{d}$. Given an absolutely continuous measure $\mu_0 \in \mathcal{P}(\overline{\mathbb{B}}(0, r))$ with density $\rho_0 \in L^\infty(\overline{\mathbb{B}}(0, r))$, there is a generalized minimizing movement scheme $(\mu_t)_t$ in $\mathcal{P}(\overline{\mathbb{B}}(0, r))$ starting from μ_0 for the functional defined by

$$\mathcal{F}^\nu(h, \mu_+, \mu_-) = \mathcal{F}_\lambda^\nu(\mu_+) + \frac{1}{2h} \mathcal{W}_2^2(\mu_+, \mu_-). \quad (\text{S2})$$

Moreover for any time $t > 0$, the probability measure $\mu_t = \mu(t)$ has density ρ_t with respect to the Lebesgue measure and:

$$\|\rho_t\|_{L^\infty} \leq e^{dt\sqrt{d}} \|\rho_0\|_{L^\infty}. \quad (\text{S3})$$

Proof. We start by noting, that by **S4** for any $h > 0$ there exists a piecewise constant trajectory μ^h for **S2** starting at μ_0 . Furthermore for $t \geq 0$ measure $\mu_t^h = \mu^h(t)$ has density ρ_t^h , and:

$$\|\rho_t^h\|_{L^\infty} \leq e^{d\sqrt{d}(t+h)} \|\rho_0\|_{L^\infty}. \quad (\text{S4})$$

Let us choose $T > 0$. We denote $\rho^h(t, x) = \rho_t^h(x)$. For $h \leq 1$, the functions ρ^h lie in a ball in $L^\infty([0, T] \times \overline{\mathbb{B}}(0, r))$, so from Banach-Alaoglu theorem there is a sequence h_n converging to 0, such that ρ^{h_n} converges in weak-star topology in $L^\infty([0, T] \times \overline{\mathbb{B}}(0, r))$ to a certain limit ρ . Since ρ has to be nonnegative except for a set of measure zero, we assume ρ is nonnegative. We denote $\rho_t(x) = \rho(t, x)$. We will prove that for almost all t , ρ_t is a probability density and $\mu_t^{h_n}$ converges in \mathcal{W}_2 to a measure μ_t with density ρ_t .

First of all, for almost all $t \in [0, T]$, ρ_t is a probability density, since for any Borel set $A \subseteq [0, T]$ the indicator of set $A \times \overline{\mathbb{B}}(0, r)$ is integrable, and hence by definition of the weak-star topology:

$$\int_A \int_{\overline{\mathbb{B}}(0, r)} \rho_t(x) dx dt = \lim_{n \rightarrow \infty} \int_A \int_{\overline{\mathbb{B}}(0, r)} \rho_t^{h_n}(x) dx dt,$$

and so we have to have $\int \rho_t(x)dx = 1$ for almost all $t \in [0, T]$. Nonnegativity of ρ_t follows from nonnegativity of ρ .

We will now prove, that for almost all $t \in [0, T]$ the measures $\mu_t^{h_n}$ converge to a measure with density ρ_t . Let $t \in (0, T)$, take $\delta < \min(T - t, t)$ and $\zeta \in C^1(\overline{\mathbb{B}}(0, r))$. We have:

$$\begin{aligned} & \left| \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_t^{h_n} - \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_t^{h_m} \right| \leq \\ & \left| \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_t^{h_n} - \frac{1}{2\delta} \int_{t-\delta}^{t+\delta} \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_s^{h_n} ds \right| + \left| \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_t^{h_m} - \frac{1}{2\delta} \int_{t-\delta}^{t+\delta} \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_s^{h_m} ds \right| + \\ & \left| \frac{1}{2\delta} \int_{t-\delta}^{t+\delta} \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_s^{h_m} ds - \frac{1}{2\delta} \int_{t-\delta}^{t+\delta} \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_s^{h_n} ds \right|. \quad (\text{S5}) \end{aligned}$$

Because $\mu_t^{h_n}$ have densities $\rho_t^{h_n}$ and both ρ^{h_n}, ρ^{h_m} converge to ρ in weak-star topology, the last element of the sum on the right hand side converges to zero, as $n, m \rightarrow \infty$. Next, we get a bound on the other two terms.

First, if we denote by γ the optimal transport plan between $\mu_t^{h_n}$ and $\mu_s^{h_n}$, we have:

$$\left| \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_t^{h_n} - \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_s^{h_n} \right|^2 \leq \int_{\overline{\mathbb{B}}(0,r) \times \overline{\mathbb{B}}(0,r)} |\zeta(x) - \zeta(y)|^2 d\gamma(x, y) \leq \|\nabla \zeta\|_\infty^2 \mathcal{W}_2^2(\mu_t^{h_n}, \mu_s^{h_n}). \quad (\text{S6})$$

In addition, for $n_t = \lfloor t/h_n \rfloor$ and $n_s = \lfloor s/h_n \rfloor$ we have $\mu_t^{h_n} = \mu_{n_t h_n}^{h_n}$ and $\mu_s^{h_n} = \mu_{n_s h_n}^{h_n}$. For all $k \geq 0$ we have:

$$\mathcal{W}_2^2(\mu_{k h_n}^{h_n}, \mu_{(k+1)h_n}^{h_n}) \leq 2h_n (\mathcal{F}_\lambda^\nu(\mu_{k h_n}^{h_n}) - \mathcal{F}_\lambda^\nu(\mu_{(k+1)h_n}^{h_n})). \quad (\text{S7})$$

Using this result and (S6) and assuming without loss of generality $n_t \leq n_s$, from the Cauchy-Schwartz inequality we get:

$$\begin{aligned} \mathcal{W}_2^2(\mu_t^{h_n}, \mu_s^{h_n}) & \leq \left(\sum_{k=n_t}^{n_s-1} \mathcal{W}_2(\mu_{k h_n}^{h_n}, \mu_{(k+1)h_n}^{h_n}) \right)^2 \\ & \leq |n_t - n_s| \sum_{k=n_t}^{n_s-1} \mathcal{W}_2^2(\mu_{k h_n}^{h_n}, \mu_{(k+1)h_n}^{h_n}) \\ & \leq 2h_n |n_t - n_s| (\mathcal{F}_\lambda^\nu(\mu_{n_t h_n}^{h_n}) - \mathcal{F}_\lambda^\nu(\mu_{n_s h_n}^{h_n})) \leq 2C(|t - s| + h_n), \quad (\text{S8}) \end{aligned}$$

where we used for the last inequality, denoting $C = \mathcal{F}_\lambda^\nu(\mu_0) - \min_{\mathcal{P}(\overline{\mathbb{B}}(0,r))} \mathcal{F}_\lambda^\nu$, that $(\mathcal{F}_\lambda^\nu(\mu_{k h_n}^{h_n}))_n$ is non-increasing by (S7) and $\min_{\mathcal{P}(\overline{\mathbb{B}}(0,r))} \mathcal{F}_\lambda^\nu$ is finite since \mathcal{F}_λ^ν is lower semi-continuous. Finally, using Jensen's inequality, the above bound and S6 we get:

$$\begin{aligned} & \left| \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_t^{h_n} - \frac{1}{2\delta} \int_{t-\delta}^{t+\delta} \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_s^{h_n} ds \right|^2 \leq \frac{1}{2\delta} \int_{t-\delta}^{t+\delta} \left| \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_t^{h_n} - \int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_s^{h_n} \right|^2 ds \\ & \leq \frac{C \|\nabla \zeta\|_\infty^2}{\delta} \int_{t-\delta}^{t+\delta} (|t - s| + h_n) ds \\ & \leq 2C \|\nabla \zeta\|_\infty^2 (h_n + \delta). \end{aligned}$$

Together with (S5), when taking $\delta = h_n$, this result means that $\int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_t^{h_n}$ is a Cauchy sequence for all $t \in (0, T)$. On the other hand, since ρ^{h_n} converges to ρ in weak-star topology on L^∞ , the limit of $\int_{\overline{\mathbb{B}}(0,r)} \zeta d\mu_t^{h_n}$ has to be $\int_{\overline{\mathbb{B}}(0,r)} \zeta(x) \rho_t(x) dx$ for almost all $t \in (0, T)$. This means that for almost all $t \in [0, T]$ sequence $\mu_t^{h_n}$ converges to a measure μ_t with density ρ_t .

Let $S \in [0, T]$ be the set of times such that for $t \in S$ sequence $\mu_t^{h_n}$ converges to μ_t . As we established almost all points from $[0, T]$ belong to S . Let $t \in [0, T] \setminus S$. Then, there exists a sequence of times $t_k \in S$ converging to t , such that μ_{t_k} converge to some limit μ_t . We have:

$$\mathcal{W}_2(\mu_t^{h_n}, \mu_t) \leq \mathcal{W}_2(\mu_t^{h_n}, \mu_{t_k}^{h_n}) + \mathcal{W}_2(\mu_{t_k}^{h_n}, \mu_{t_k}) + \mathcal{W}_2(\mu_{t_k}, \mu_t).$$

From which we have for all $k \geq 1$:

$$\limsup_{n \rightarrow \infty} \mathcal{W}_2(\mu_t^{h_n}, \mu_t) \leq \mathcal{W}_2(\mu_{t_k}, \mu_t) + \limsup_{n \rightarrow \infty} \mathcal{W}_2(\mu_t^{h_n}, \mu_{t_k}),$$

and using (S8), we get $\mu_t^{h_n} \rightarrow \mu_t$. Furthermore, the measure μ_t has to have density, since $\rho_t^{h_n}$ lie in a ball in $L^\infty(\overline{\mathbb{B}}(0, r))$, so we can choose a subsequence of $\rho_t^{h_n}$ converging in weak-star topology to a certain limit $\hat{\rho}_t$, which is the density of μ_t .

We use now the diagonal argument to get convergence for all $t > 0$. Let $(T_k)_{k=1}^\infty$ be a sequence of times increasing to infinity. Let h_n^1 be a sequence converging to 0, such that $\mu_t^{h_n^1}$ converge to μ_t for all $t \in [0, T_1]$. Using the same arguments as above, we can choose a subsequence h_n^2 of h_n^1 , such that $\mu_t^{h_n^2}$ converges to a limit μ_t for all $t \in [0, T_2]$. Inductively, we construct subsequences h_n^k , and in the end take $h_n = h_n^n$. For this subsequence we have that $\mu_t^{h_n}$ converges to μ_t for all $t > 0$, and μ_t has a density satisfying the bound from the statement of the theorem.

Finally, note that (S5) follows from (S4). \square

Theorem S6. *Let $(\mu_t)_{t \geq 0}$ be a generalized minimizing movement scheme given by Theorem S5 with initial distribution μ_0 with density $\rho_0 \in L(\overline{\mathbb{B}}(0, r))$. We denote by ρ_t the density of μ_t for all $t \geq 0$. Then ρ_t satisfies the continuity equation:*

$$\frac{\partial \rho_t}{\partial t} + \operatorname{div}(v_t \rho_t) + \lambda \Delta \rho_t = 0, \quad v_t(x) = - \int_{\mathbb{S}^{d-1}} \psi'_{t,\theta}(\langle x, \theta \rangle) \theta d\theta,$$

in a weak sense, that is for all $\xi \in C_c^\infty([0, \infty) \times \overline{\mathbb{B}}(0, r))$ we have:

$$\int_0^\infty \int_{\overline{\mathbb{B}}(0,r)} \left[\frac{\partial \xi}{\partial t}(t, x) - v_t \nabla \xi(t, x) - \lambda \Delta \xi(t, x) \right] \rho_t(x) dx dt = - \int_{\overline{\mathbb{B}}(0,r)} \xi(0, x) \rho_0(x) dx.$$

Proof. Our proof is based on the proof of (Bonnotte, 2013)[Theorem 5.6.1]. We proceed in five steps.

(1) Let $h_n \rightarrow 0$ be a sequence given by Theorem S5, such that $\mu_t^{h_n}$ converges to μ_t pointwise. Furthermore we know that $\mu_t^{h_n}$ have densities $\rho_t^{h_n}$ that converge to ρ in L^r , for $r \geq 1$, and in weak-star topology in L^∞ . Let $\xi \in C_c^\infty([0, \infty) \times \overline{\mathbb{B}}(0, r))$. We denote $\xi_k^n(x) = \xi(kh_n, x)$. Using part 1 of the proof of (Bonnotte, 2013)[Theorem 5.6.1], we obtain:

$$\begin{aligned} \int_{\overline{\mathbb{B}}(0,r)} \xi(0, x) \rho_0(x) dx + \int_0^\infty \int_{\overline{\mathbb{B}}(0,r)} \frac{\partial \xi}{\partial t}(t, x) \rho_t(x) dx dt \\ = \lim_{n \rightarrow \infty} -h_n \sum_{k=1}^\infty \int_{\overline{\mathbb{B}}(0,r)} \xi_k^n(x) \frac{\rho_{kh_n}^{h_n}(x) - \rho_{(k-1)h_n}^{h_n}(x)}{h_n} dx. \end{aligned} \quad (\text{S9})$$

(2) Again, this part is the same as part 2 of the proof of (Bonnotte, 2013)[Theorem 5.6.1]. For any $\theta \in \mathbb{S}^{d-1}$ we denote by $\psi_{t,\theta}$ the unique Kantorovich potential from $\theta_{\#}^* \mu_t$ to $\theta_{\#}^* \nu$, and by $\psi_{t,\theta}^{h_n}$ the unique Kantorovich potential from $\theta_{\#}^* \mu_t^{h_n}$ to $\theta_{\#}^* \nu$. Then, by the same reasoning as part 2 of the proof of (Bonnotte, 2013)[Theorem 5.6.1], we get:

$$\begin{aligned} \int_0^\infty \int_{\overline{\mathbb{B}}(0,r)} \int_{\mathbb{S}^{d-1}} (\psi_{t,\theta})'(\langle \theta, x \rangle) \langle \theta, \nabla \xi(x, t) \rangle d\theta d\mu_t(x) dt \\ = \lim_{n \rightarrow \infty} h_n \sum_{k=1}^\infty \int_{\overline{\mathbb{B}}(0,r)} \int_{\mathbb{S}^{d-1}} \psi_{kh_n,\theta}^{h_n}(\theta^*) \langle \theta, \nabla \xi_k^n \rangle d\theta d\mu_{kh_n}^{h_n}. \end{aligned} \quad (\text{S10})$$

(3) Since ξ is compactly supported and smooth, $\Delta \xi$ is Lipschitz, and so for any $t \geq 0$ if we take $k = \lfloor t/h_n \rfloor$ we get $|\Delta \xi_k^n(x) - \Delta \xi(t, x)| \leq Ch_n$ for some constant C . Let $T > 0$ be such that $\xi(t, x) = 0$ for $t > T$. We have:

$$\left| \sum_{k=1}^\infty h_n \int_{\overline{\mathbb{B}}(0,r)} \Delta \xi_k^n(x) \rho_{kh_n}^{h_n}(x) dx - \int_0^{+\infty} \int_{\overline{\mathbb{B}}(0,r)} \Delta \xi(t, x) \rho_t^{h_n}(x) dx dt \right| \leq CTh_n.$$

On the other hand, we know, that ρ^{h_n} converges to ρ in weak star topology on $L^\infty([0, T] \times \overline{B}(0, r))$, and $\Delta\xi$ is bounded, so:

$$\lim_{n \rightarrow +\infty} \left| \int_0^{+\infty} \int_{\overline{B}(0, r)} \Delta\xi(t, x) \rho_t^{h_n}(x) dx dt - \int_0^{+\infty} \int_{\overline{B}(0, r)} \Delta\xi(t, x) \rho_t(x) dx dt \right| = 0.$$

Combining those two results give:

$$\lim_{n \rightarrow \infty} h_n \sum_{k=1}^{\infty} \int_{\overline{B}(0, r)} \Delta\xi_k^n(x) \rho_{kh_n}^{h_n}(x) dx = \int_0^{+\infty} \int_{\overline{B}(0, r)} \Delta\xi(t, x) \rho_t(x) dx dt. \quad (\text{S11})$$

(4) Let $\phi_k^{h_n}$ denote the unique Kantorovich potential from $\mu_{kh_n}^{h_n}$ to $\mu_{(k-1)h_n}^{h_n}$. Using (Bonnotte, 2013)[Propositions 1.5.7 and 5.1.7], as well as (Jordan et al., 1998)[Equation (38)] with $\Psi = 0$, and optimality of $\mu_{kh_n}^{h_n}$, we get:

$$\begin{aligned} \frac{1}{h_n} \int_{\overline{B}(0, r)} \langle \nabla \phi_k^{h_n}(x), \nabla \xi_k^n(x) \rangle d\mu_{kh_n}^{h_n}(x) - \int_{\overline{B}(0, r)} \int_{\mathbb{S}^{d-1}} (\psi_{kh_n}^{h_n})'(\theta^*) \langle \theta, \nabla \xi_k^n(x) \rangle d\theta d\mu_{kh_n}^{h_n}(x) \\ - \lambda \int_{\overline{B}(0, r)} \Delta\xi_k^n(x) d\mu_{kh_n}^{h_n}(x), \end{aligned} \quad (\text{S12})$$

which is the derivative of $\mathcal{F}_\lambda^\nu(\cdot) + \frac{1}{2h_n} \mathcal{W}_2^2(\cdot, \mu_{(k-1)h_n}^{h_n})$ in the direction given by vector field $\nabla \xi_k^n$ is zero.

Let γ be the optimal transport between $\mu_{kh_n}^{h_n}$ and $\mu_{(k-1)h_n}^{h_n}$. Then:

$$\int_{\overline{B}(0, r)} \xi_k^n(x) \frac{\rho_{kh_n}^{h_n}(x) - \rho_{(k-1)h_n}^{h_n}(x)}{h_n} dx = \frac{1}{h_n} \int_{\overline{B}(0, r)} (\xi_k^n(y) - \xi_k^n(x)) d\gamma(x, y). \quad (\text{S13})$$

$$\frac{1}{h_n} \int_{\overline{B}(0, r)} \langle \nabla \phi_k^{h_n}(x), \nabla \xi_k^n(x) \rangle d\mu_{kh_n}^{h_n}(x) = \frac{1}{h_n} \int_{\overline{B}(0, r)} \langle \nabla \xi_k^n(x), y - x \rangle d\gamma(x, y). \quad (\text{S14})$$

Since ξ is C_c^∞ , it has Lipschitz gradient. Let C be twice the Lipschitz constant of $\nabla \xi$. Then we have $|\xi(y) - \xi(x) - \langle \nabla \xi(x), y - x \rangle| \leq C|x - y|^2$, and hence:

$$\int_{\overline{B}(0, r)} |\xi_k^n(y) - \xi_k^n(x) - \langle \nabla \xi_k^n(x), y - x \rangle| d\gamma(x, y) \leq C \mathcal{W}_2^2(\mu_{(k-1)h_n}^{h_n}, \mu_{kh_n}^{h_n}). \quad (\text{S15})$$

Combining (S13), (S14) and (S15), we get:

$$\begin{aligned} \left| \sum_{k=1}^{\infty} h_n \int_{\overline{B}(0, r)} \xi_k^n(x) \frac{\rho_{kh_n}^{h_n} - \rho_{(k-1)h_n}^{h_n}}{h_n} dx + \sum_{k=1}^{\infty} h_n \int_{\overline{B}(0, r)} \langle \nabla \phi_k^{h_n}, \nabla \xi_k^n \rangle d\mu_{kh_n}^{h_n} \right| \\ \leq C \sum_{k=1}^{\infty} \mathcal{W}_2^2(\mu_{(k-1)h_n}^{h_n}, \mu_{kh_n}^{h_n}). \end{aligned} \quad (\text{S16})$$

As some \mathcal{F}_λ^ν have a finite minimum on $\mathcal{P}(\overline{B}(0, r))$, we have:

$$\begin{aligned} \sum_{k=1}^{\infty} \mathcal{W}_2^2(\mu_{(k-1)h_n}^{h_n}, \mu_{kh_n}^{h_n}) &\leq 2h_n \sum_{k=1}^{\infty} \mathcal{F}_\lambda^\nu(\mu_{(k-1)h_n}^{h_n}) - \mathcal{F}_\lambda^\nu(\mu_{kh_n}^{h_n}) \\ &\leq 2h_n \left(\mathcal{F}_\lambda^\nu(\mu_0) - \min_{\mathcal{P}(\overline{B}(0, r))} \mathcal{F}_\lambda^\nu \right). \end{aligned} \quad (\text{S17})$$

and so the sum on the right hand side of the equation goes to zero as n goes to infinity.

From (S16), (S17) and (S12) we conclude:

$$\lim_{n \rightarrow \infty} -h_n \sum_{k=1}^{\infty} \xi_k^n(x) \frac{\rho_{kh_n}^{h_n} - \rho_{(k-1)h_n}^{h_n}}{h_n} dx = \lim_{n \rightarrow \infty} \left(h_n \sum_{k=1}^{\infty} \int_{\overline{B}(0,r)} \int_{\mathbb{S}^{d-1}} \psi_{kh_n, \theta}^{h_n}(\theta^*) \langle \theta, \nabla \xi_k^n \rangle d\theta d\mu_{kh_n}^{h_n} + h_n \sum_{k=1}^{\infty} \int_{\overline{B}(0,r)} \Delta \xi_k^n(x) \rho_{kh_n}^{h_n}(x) dx \right), \quad (\text{S18})$$

where both limits exist, since the difference of left hand side and right hand side of the equation goes to zero, while the left hand side converges to a finite value by (S9).

(5) Combining (S9), (S10), (S11) and (S18) we get the result. □

2. Proof of Theorem 3

Before proceeding to the proof, let us first define the following Euler-Maruyama scheme which will be useful for our analysis:

$$\hat{X}_{k+1} = \hat{X}_k + h\hat{v}(\hat{X}_k, \mu_{kh}) + \sqrt{2\lambda h}Z_{n+1}, \quad (\text{S19})$$

where μ_t denotes the probability distribution of X_t with $(X_t)_t$ being the solution of the original SDE (8). Now, consider the probability distribution of \hat{X}_k as $\hat{\mu}_{kh}$. Starting from the discrete-time process $(\hat{X}_k)_{k \in \mathbb{N}_+}$, we first define a continuous-time process $(Y_t)_{t \geq 0}$ that linearly interpolates $(\hat{X}_k)_{k \in \mathbb{N}_+}$, given as follows:

$$dY_t = \tilde{v}_t(Y)dt + \sqrt{2\lambda}dW_t, \quad (\text{S20})$$

where $\tilde{v}_t(Y) \triangleq -\sum_{k=0}^{\infty} \hat{v}_{kh}(Y_{kh})\mathbb{1}_{[kh, (k+1)h)}(t)$ and $\mathbb{1}$ denotes the indicator function. Similarly, we define a continuous-time process $(U_t)_{t \geq 0}$ that linearly interpolates $(\bar{X}_k)_{k \in \mathbb{N}_+}$, defined by (13), given as follows:

$$dU_t = \bar{v}_t(U)dt + \sqrt{2\lambda}dW_t, \quad (\text{S21})$$

where $\bar{v}_t(U) \triangleq -\sum_{k=0}^{\infty} \hat{v}(U_{kh}, \bar{\mu}_{kh})\mathbb{1}_{[kh, (k+1)h)}(t)$ and $\bar{\mu}_{kh}$ denotes the probability distribution of \bar{X}_k . Let us denote the distributions of $(X_t)_{t \in [0, T]}$, $(Y_t)_{t \in [0, T]}$ and $(U_t)_{t \in [0, T]}$ as π_X^T , π_Y^T and π_U^T respectively with $T = Kh$.

We consider the following assumptions:

HS1. For all $\lambda > 0$, the SDE (8) has a unique strong solution denoted by $(X_t)_{t \geq 0}$ for any starting point $x \in \mathbb{R}^d$.

HS2. There exists $L < \infty$ such that

$$\|v_t(x) - v_{t'}(x')\| \leq L(\|x - x'\| + |t - t'|), \quad (\text{S22})$$

where $v_t(x) = v(x, \mu_t)$ and

$$\|\hat{v}(x, \mu) - \hat{v}(x', \mu')\| \leq L(\|x - x'\| + \|\mu - \mu'\|_{\text{TV}}). \quad (\text{S23})$$

HS3. For all $t \geq 0$, v_t is dissipative, i.e. for all $x \in \mathbb{R}^d$,

$$\langle x, v_t(x) \rangle \geq m\|x\|^2 - b, \quad (\text{S24})$$

for some $m, b > 0$.

HS4. The estimator of the drift satisfies the following conditions: $\mathbb{E}[\hat{v}_t] = v_t$ for all $t \geq 0$, and for all $t \geq 0$, $x \in \mathbb{R}^d$,

$$\mathbb{E}[\|\hat{v}(x, \mu_t) - v(x, \mu_t)\|^2] \leq 2\delta(L^2\|x\|^2 + B^2), \quad (\text{S25})$$

for some $\delta \in (0, 1)$.

HS5. For all $t \geq 0$: $|\Psi_t(0)| \leq A$ and $\|v_t(0)\| \leq B$, for $A, B \geq 0$, where $\Psi_t = \int_{\mathbb{S}^{d-1}} \psi_t(\langle \theta, \cdot \rangle) d\theta$.

We start by upper-bounding $\|\hat{\mu}_{Kh} - \mu_T\|_{\text{TV}}$.

Lemma S1. Assume that the conditions **HS2** to **S5** hold. Then, the following bound holds:

$$\|\hat{\mu}_{Kh} - \mu_T\|_{\text{TV}}^2 \leq \|\pi_Y^T - \pi_X^T\|_{\text{TV}}^2 \leq \frac{L^2 K}{4\lambda} \left(\frac{C_1 h^3}{3} + 3\lambda d h^2 \right) + \frac{C_2 \delta K h}{8\lambda}, \quad (\text{S26})$$

where $C_1 \triangleq 12(L^2 C_0 + B^2) + 1$, $C_2 \triangleq 2(L^2 C_0 + B^2)$, $C_0 \triangleq C_e + 2(1 \vee \frac{1}{m})(b + 2B^2 + d\lambda)$, and C_e denotes the entropy of μ_0 .

Proof. We use the proof technique presented in (Dalalyan, 2017; Raginsky et al., 2017). It is easy to verify that for all $k \in \mathbb{N}_+$, we have $Y_{kh} = \hat{X}_k$.

By Girsanov's theorem to express the Kullback-Leibler (KL) divergence between these two distributions, given as follows:

$$\text{KL}(\pi_X^T \| \pi_Y^T) = \frac{1}{4\lambda} \int_0^{Kh} \mathbb{E}[\|v_t(Y_t) + \tilde{v}_t(Y)\|^2] dt \quad (\text{S27})$$

$$= \frac{1}{4\lambda} \sum_{k=0}^{K-1} \int_{kh}^{(k+1)h} \mathbb{E}[\|v_t(Y_t) + \tilde{v}_t(Y)\|^2] dt \quad (\text{S28})$$

$$= \frac{1}{4\lambda} \sum_{k=0}^{K-1} \int_{kh}^{(k+1)h} \mathbb{E}[\|v_t(Y_t) - \hat{v}_{kh}(Y_{kh})\|^2] dt. \quad (\text{S29})$$

By using $v_t(Y_t) - \hat{v}_{kh}(Y_{kh}) = (v_t(Y_t) - v_{kh}(Y_{kh})) + (v_{kh}(Y_{kh}) - \hat{v}_{kh}(Y_{kh}))$, we obtain

$$\begin{aligned} \text{KL}(\pi_X^T \| \pi_Y^T) &\leq \frac{1}{2\lambda} \sum_{k=0}^{K-1} \int_{kh}^{(k+1)h} \mathbb{E}[\|v_t(Y_t) - v_{kh}(Y_{kh})\|^2] dt \\ &\quad + \frac{1}{2\lambda} \sum_{k=0}^{K-1} \int_{kh}^{(k+1)h} \mathbb{E}[\|v_{kh}(Y_{kh}) - \hat{v}_{kh}(Y_{kh})\|^2] dt \end{aligned} \quad (\text{S30})$$

$$\begin{aligned} &\leq \frac{L^2}{\lambda} \sum_{k=0}^{K-1} \int_{kh}^{(k+1)h} (\mathbb{E}[\|Y_t - Y_{kh}\|^2] + (t - kh)^2) dt \\ &\quad + \frac{1}{2\lambda} \sum_{k=0}^{K-1} \int_{kh}^{(k+1)h} \mathbb{E}[\|v_{kh}(Y_{kh}) - \hat{v}_{kh}(Y_{kh})\|^2] dt. \end{aligned} \quad (\text{S31})$$

The last inequality is due to the Lipschitz condition **HS2**.

Now, let us focus on the term $\mathbb{E}[\|Y_t - Y_{kh}\|^2]$. By using (S20), we obtain:

$$Y_t - Y_{kh} = -(t - kh)\hat{v}_{kh}(Y_{kh}) + \sqrt{2\lambda(t - kh)}Z, \quad (\text{S32})$$

where Z denotes a standard normal random variable. By adding and subtracting the term $-(t - kh)v_{kh}(Y_{kh})$, we have:

$$Y_t - Y_{kh} = -(t - kh)v_{kh}(Y_{kh}) + (t - kh)(v_{kh}(Y_{kh}) - \hat{v}_{kh}(Y_{kh})) + \sqrt{2\lambda(t - kh)}Z. \quad (\text{S33})$$

Taking the square and then the expectation of both sides yields:

$$\begin{aligned} \mathbb{E}[\|Y_t - Y_{kh}\|^2] &\leq 3(t - kh)^2 \mathbb{E}[\|v_{kh}(Y_{kh})\|^2] + 3(t - kh)^2 \mathbb{E}[\|v_{kh}(Y_{kh}) - \hat{v}_{kh}(Y_{kh})\|^2] \\ &\quad + 6\lambda(t - kh)d. \end{aligned} \quad (\text{S34})$$

As a consequence of **HS2** and **HS5**, we have $\|v_t(x)\| \leq L\|x\| + B$ for all $t \geq 0$, $x \in \mathbb{R}^d$. Combining this inequality with **H S4**, we obtain:

$$\begin{aligned} \mathbb{E}[\|Y_t - Y_{kh}\|^2] &\leq 6(t - kh)^2 (L^2 \mathbb{E}[\|Y_{kh}\|^2] + B^2) + 6(t - kh)^2 (L^2 \mathbb{E}[\|Y_{kh}\|^2] + B^2) \\ &\quad + 6\lambda(t - kh)d \end{aligned} \quad (\text{S35})$$

$$= 12(t - kh)^2 (L^2 \mathbb{E}[\|Y_{kh}\|^2] + B^2) + 6\lambda(t - kh)d. \quad (\text{S36})$$

By Lemma 3.2 of (Raginsky et al., 2017)⁴, we have $\mathbb{E}[\|Y_{kh}\|^2] \leq C_0 \triangleq C_e + 2(1 \vee \frac{1}{m})(b + 2B^2 + d\lambda)$, where C_e denotes the entropy of μ_0 . Using this result in the above equation yields:

$$\mathbb{E}[\|Y_t - Y_{kh}\|^2] \leq 12(t - kh)^2(L^2C_0 + B^2) + 6\lambda(t - kh)d. \quad (\text{S37})$$

We now focus on the term $\mathbb{E}[\|v_{kh}(Y_{kh}) - \hat{v}_{kh}(Y_{kh})\|^2]$ in (S31). Similarly to the previous term, we can upper-bound this term as follows:

$$\mathbb{E}[\|v_{kh}(Y_{kh}) - \hat{v}_{kh}(Y_{kh})\|^2] \leq 2\delta(L^2\mathbb{E}[\|Y_{kh}\|^2] + B^2) \quad (\text{S38})$$

$$\leq 2\delta(L^2C_0 + B^2). \quad (\text{S39})$$

By using (S37) and (S39) in (S31), we obtain:

$$\begin{aligned} \text{KL}(\pi_X^T \|\pi_Y^T) &\leq \frac{L^2}{\lambda} \sum_{k=0}^{K-1} \int_{kh}^{(k+1)h} (12(t - kh)^2(L^2C_0 + B^2) + 6\lambda(t - kh)d + (t - kh)^2) dt \\ &\quad + \frac{1}{2\lambda} \sum_{k=0}^{K-1} \int_{kh}^{(k+1)h} 2\delta(L^2C_0 + B^2) dt \end{aligned} \quad (\text{S40})$$

$$= \frac{L^2K}{\lambda} \left(\frac{C_1 h^3}{3} + \frac{6\lambda d h^2}{2} \right) + \frac{C_2 \delta K h}{2\lambda}, \quad (\text{S41})$$

where $C_1 = 12(L^2C_0 + B^2) + 1$ and $C_2 = 2(L^2C_0 + B^2)$.

Finally, by using the data processing and Pinsker inequalities, we obtain:

$$\|\hat{\mu}_{Kh} - \mu_T\|_{\text{TV}}^2 \leq \|\pi_X^T - \pi_Y^T\|_{\text{TV}}^2 \leq \frac{1}{4} \text{KL}(\pi_X^T \|\pi_Y^T) \quad (\text{S42})$$

$$= \frac{L^2K}{4\lambda} \left(\frac{C_1 h^3}{3} + 3\lambda d h^2 \right) + \frac{C_2 \delta K h}{8\lambda}. \quad (\text{S43})$$

This concludes the proof. \square

Now, we bound the term $\|\bar{\mu}_{Kh} - \hat{\mu}_{Kh}\|_{\text{TV}}$.

Lemma S2. *Assume that HS2 holds. Then the following bound holds:*

$$\|\pi_U^T - \pi_Y^T\|_{\text{TV}}^2 \leq \frac{L^2Kh}{16\lambda} \|\pi_X^T - \pi_U^T\|_{\text{TV}}^2. \quad (\text{S44})$$

Proof. We use that same approach than in Lemma S1. By Girsanov's theorem once again, we have

$$\text{KL}(\pi_Y^T \|\pi_U^T) = \frac{1}{4\lambda} \sum_{k=0}^{K-1} \int_{kh}^{(k+1)h} \mathbb{E}[\|\hat{v}(U_{kh}, \mu_{kh}) - \hat{v}(U_{kh}, \bar{\mu}_{kh})\|^2] dt, \quad (\text{S45})$$

where π_U^T denotes the distributions of $(U_t)_{t \in [0, T]}$ with $T = Kh$. By using HS2, we have:

$$\text{KL}(\pi_Y^T \|\pi_U^T) \leq \frac{L^2h}{4\lambda} \sum_{k=0}^{K-1} \|\mu_{kh} - \bar{\mu}_{kh}\|_{\text{TV}}^2 \quad (\text{S46})$$

$$\leq \frac{L^2Kh}{4\lambda} \|\pi_X^T - \pi_U^T\|_{\text{TV}}^2. \quad (\text{S47})$$

By applying the data processing and Pinsker inequalities, we obtain the desired result. \square

⁴Note that Lemma 3.2 of (Raginsky et al., 2017) considers the case where the drift is not time- or measure-dependent. However, with HS3 it is easy to show that the same result holds for our case as well.

2.1. Proof of Theorem 3

Here, we precise the statement of Theorem 3.

Theorem S7. *Assume that the assumptions in Lemma S1 and Lemma S2 hold. Then for $\lambda > \frac{KL^2h}{8}$, the following bound holds:*

$$\|\bar{\mu}_{Kh} - \mu_T\|_{\text{TV}}^2 \leq \delta_\lambda \left\{ \frac{L^2K}{2\lambda} \left(\frac{C_1h^3}{3} + 3\lambda dh^2 \right) + \frac{C_2\delta Kh}{4\lambda} \right\}, \quad (\text{S48})$$

where $\delta_\lambda = (1 - \frac{KL^2h}{8\lambda})^{-1}$.

Proof. We have the following decomposition: (with $T = Kh$)

$$\|\pi_X^T - \pi_U^T\|_{\text{TV}}^2 \leq 2\|\pi_X^T - \pi_Y^T\|_{\text{TV}}^2 + 2\|\pi_Y^T - \pi_U^T\|_{\text{TV}}^2 \quad (\text{S49})$$

$$\leq \frac{L^2K}{2\lambda} \left(\frac{C_1h^3}{3} + 3\lambda dh^2 \right) + \frac{C_2\delta Kh}{4\lambda} + \frac{L^2Kh}{8\lambda} \|\pi_X^T - \pi_U^T\|_{\text{TV}}^2 \quad (\text{S50})$$

$$\leq \left(1 - \frac{KL^2h}{8\lambda}\right)^{-1} \left\{ \frac{L^2K}{2\lambda} \left(\frac{C_1h^3}{3} + 3\lambda dh^2 \right) + \frac{C_2\delta Kh}{4\lambda} \right\}. \quad (\text{S51})$$

The second line follows from Lemma S1 and Lemma S2. Last line follows from the assumption that λ is large enough. This completes the proof. \square

3. Proof of Corollary 1

Proof. Considering the bound given in Theorem 3, the choice h implies that

$$\frac{\delta_\lambda L^2K}{2\lambda} \left(\frac{C_1h^3}{3} + 3\lambda dh^2 \right) \leq \varepsilon^2. \quad (\text{S52})$$

This finalizes the proof. \square

4. Additional Experimental Results

4.1. The Sliced Wasserstein Flow

The whole code for the Sliced Wasserstein Flow was implemented in Python, for use with Pytorch⁵. The code was written so as to run efficiently on GPU, and is available on the publicly available repository related to this paper⁶.

In practice, the SWF involves relatively simple operations, the most important being:

- For each random $\theta \in \{\theta_n\}_{n=1\dots N_\theta}$, compute its inner product with all items from a dataset and obtain the empirical quantiles for these *projections*.
- At each step k of the SWF, for each projection $z = \langle \theta, \bar{X}_k^i \rangle$, apply two piece-wise linear functions, corresponding to the scalar optimal transport $\psi'_{k,\theta}(z)$.

Even if such steps are conceptually simple, the quantile and required linear interpolation functions were not available on GPU for any framework we could figure out at the time of writing this paper. Hence, we implemented them ourselves for use with Pytorch, and the interested reader will find the details in the Github repository dedicated to this paper.

Given these operations, putting a SWF implementation together is straightforward. The code provided allows not only to apply it on any dataset, but also provides routines to have the computation of these sketches running in the background in a parallel manner.

⁵<http://www.pytorch.org>.

⁶<https://github.com/aliutkus/swf>.

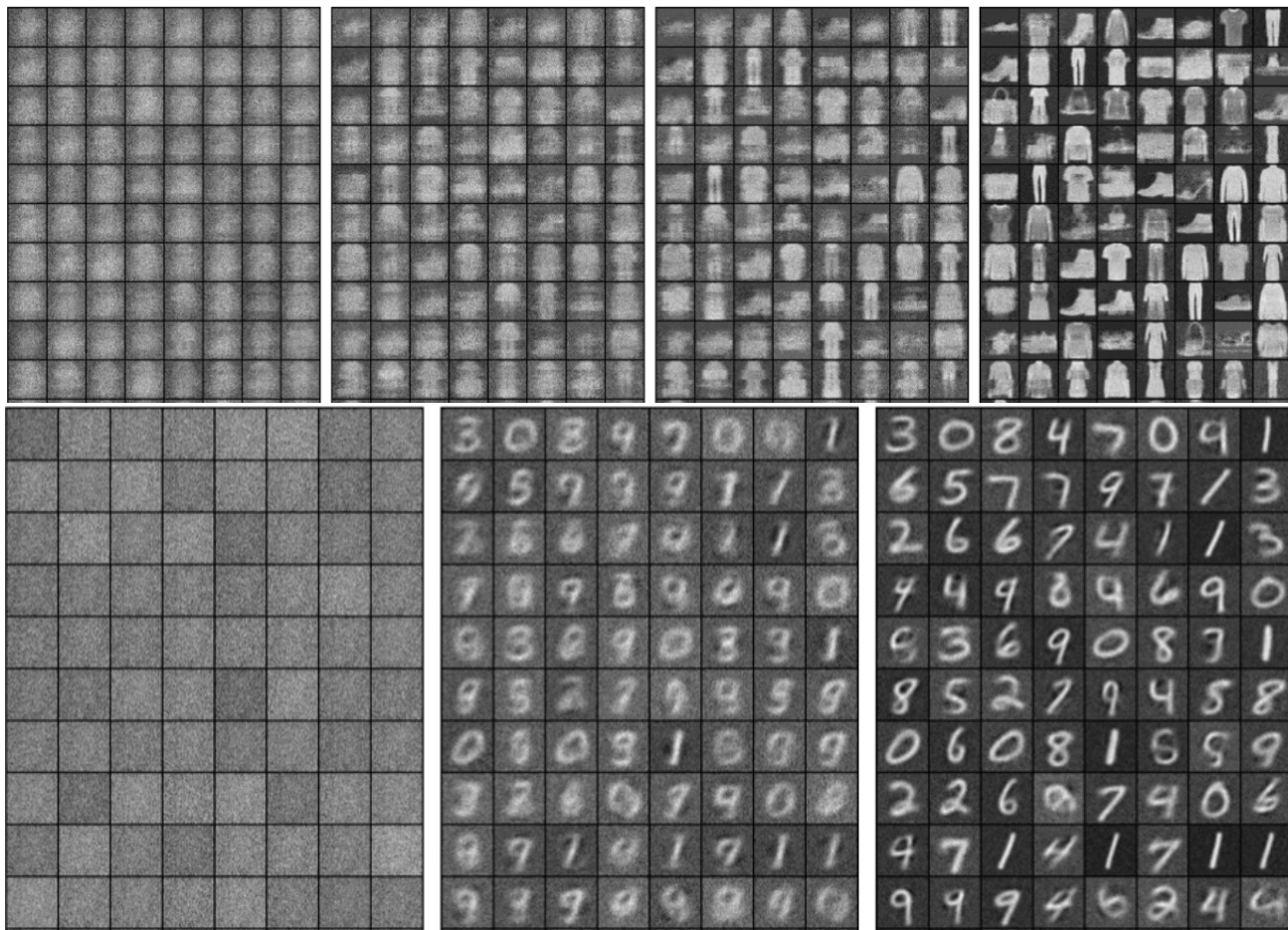


Figure S1. The evolution of SWF through 15000 iterations, when the original high-dimensional data is kept instead of working on reduced bottleneck features as done in the main document. Showing results on the MNIST and FashionMNIST datasets. For a visual comparison for FashionMNIST, we refer the reader to (Samangouei et al., 2018).

4.2. The need for dimension reduction through autoencoders

In this study, we used an autoencoder trained on the dataset as a dimension reduction technique, so that the SWF is applied to transport particles in a latent space of dimension $d \approx 50$, instead of the original $d > 1000$ of image data.

The curious reader may wonder why SWF is not applied directly to this original space, and what performances should be expected there. We have done this experiment, and we found out that SWF has much trouble rapidly converging to satisfying samples. In figure S1, we show the progressive evolution of particles undergoing SWF when the target is directly taken as the uncompressed dataset.

In this experiment, the strategy was to change the projections θ at each iteration, so that we ended up with a set of projections being $\{\theta_{n,k}\}_{n=1\dots N_\theta}^{k=1\dots K}$ instead of the fixed set of N_θ we now consider in the main document (for this, we picked $N_\theta = 200$). This strategy is motivated by the complete failure we observed whenever we picked such fixed projections throughout iterations, even for a relatively large number as $N_\theta = 16000$.

As may be seen on Figure S1, the particles definitely converge to samples from the desired datasets, and this is encouraging. However, we feel that the extreme number of iterations required to achieve such convergence comes from the fact that theory needs an integral over the d -dimensional sphere at each step of the SWF, which is clearly an issue whenever d gets too large. Although our solution of picking new samples from the sphere at each iteration alleviated this issue to some extent, the curse of dimensionality prevents us from doing much better with just thousands of *random* projections at a time.

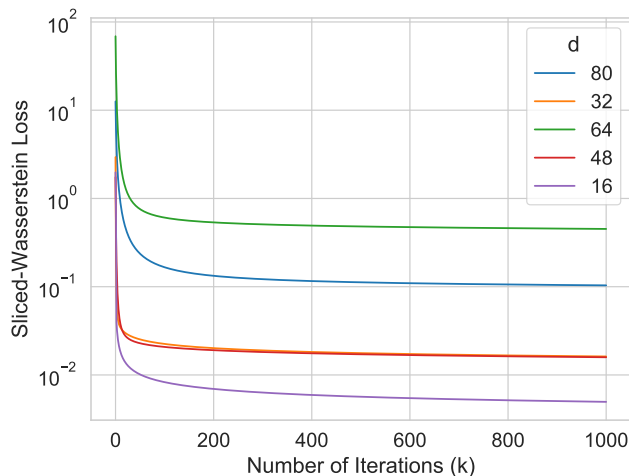


Figure S2. Approximately computed \mathcal{SW}_2 between the output $\bar{\mu}_k^N$ and data distribution ν in the MNIST experiment for different dimensions d for the bottleneck features (and the corresponding pre-trained AE).

This being said, we are confident that good performance would be obtained if millions of random projections could be considered for transporting such high dimensional data because i/ theory suggests it and ii/ we observed excellent performance on reduced dimensions.

However, we, unfortunately, did not have the computing power it takes for such large scale experiments and this is what motivated us in the first place to introduce some dimension-reduction technique through AE.

4.3. Structure of our autoencoders for reducing data dimension

As mentioned in the text, we used autoencoders to reduce the dimensionality of the transport problem. The structure of these networks is the following:

- **Encoder** Four 2d convolution layers with (num_chan_out, kernel_size, stride, padding) being (3, 3, 1, 1), (32, 2, 2, 0), (32, 3, 1, 1), (32, 3, 1, 1), each one followed by a ReLU activation. At the output, a linear layer gets the desired bottleneck size.
- **Decoder** A linear layer gets from the bottleneck features to a vector of dimension 8192, which is reshaped as (32, 16, 16). Then, three convolution layers are applied, all with 32 output channels and (kernel_size, stride, padding) being respectively (3, 1, 1), (3, 1, 1), (2, 2, 0). A 2d convolution layer is then applied with an output number of channels being that of the data (1 for black and white, 3 for color), and a (kernel_size, stride, padding) as (3, 1, 1). In any case, all layers are followed by a ReLU activation, and a sigmoid activation is applied to the very output.

Once these networks defined, these autoencoders are trained in a very simple manner by minimizing the binary cross entropy between input and output over the training set of the considered dataset (here MNIST, CelebA or FashionMNIST). This training was achieved with the Adam algorithm (Kingma & Ba, 2014) with learning rate $1e - 3$.

No additional training trick was involved as in Variational Autoencoder (Kingma & Welling, 2013) to make sure the distribution of the bottleneck features matches some prior. The core advantage of the proposed method in this respect is indeed to turn any previously learned AE as a generative model, by automatically and non-parametrically transporting particles drawn from an arbitrary prior distribution μ to the observed empirical distribution ν of the bottleneck features over the training set.

4.4. Convergence plots of SWF

In the same experimental setting as in the main document, we also illustrate the behavior of the algorithm for varying dimensionality d for the bottleneck-features. To monitor the convergence of SWF as predicted by theory, we display the

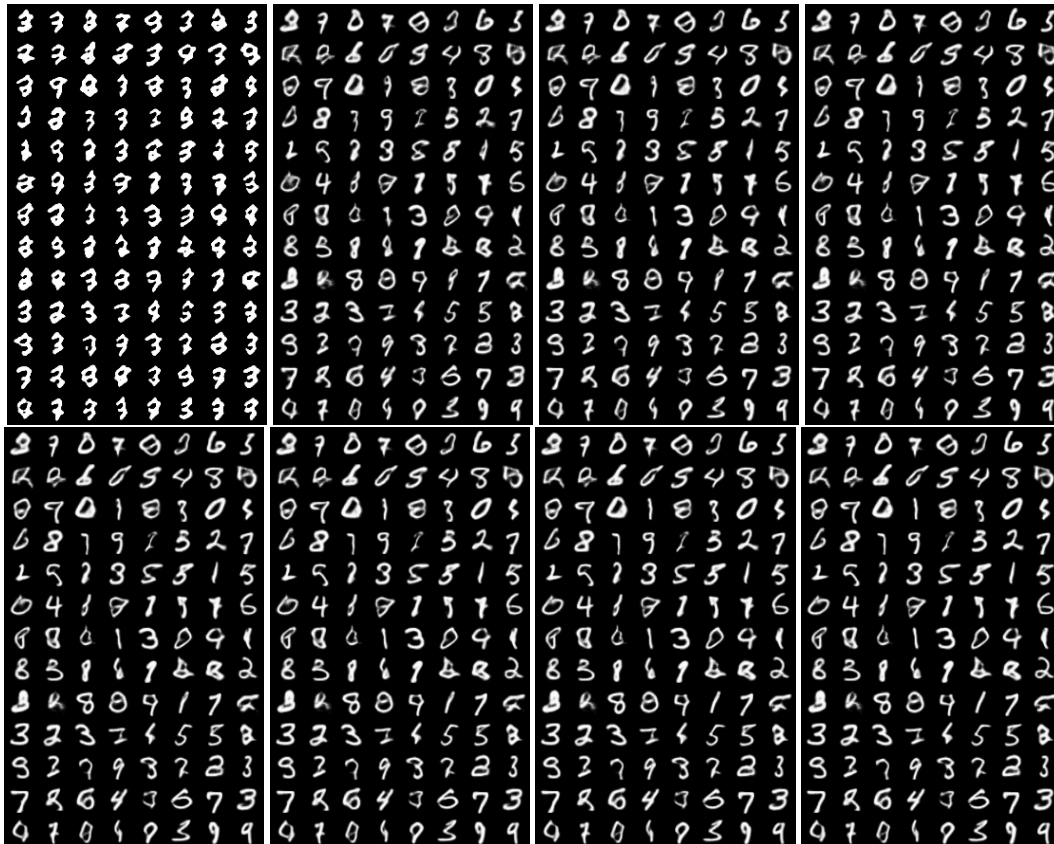


Figure S3. The evolution of SWF through 200 iterations on the MNIST dataset. Plots are for 1, 11, 21, 31, 41, 51, 101 and 201 iterations

approximately computed SW_2 distance between the distribution of the particles and the data distribution. Even though minimizing this distance is not the real objective of our method, arguably, it is still a good proxy for understanding the convergence behavior.

Figure S2 illustrates the results. We observe that, for all choices of d , we see a steady and smooth decrease in the cost for all runs, which is in line with our theory. The absolute value of the cost for varying dimensions remains hard to interpret at this stage of our investigations.

5. Additional samples

5.1. Evolution throughout iterations

In Figures S3 and S4 below, we provide the evolution of the SWF algorithm on the Fashion MNIST and the MNIST datasets in higher resolution, for an AE with $d = 48$ bottleneck features.

5.2. Training samples, interpolation and extrapolation

In Figures S5 and S6 below, we provide other examples of outcome from SWF, both for the MNIST and the FashionMNIST datasets, still with $d = 48$ bottleneck features.

The most noticeable fact we may see on these figures is that while the actual particles which went through SWF, as well as linear combinations of them, all yield very satisfying results, this is however not the case for particles that are drawn randomly and then brought through a pre-learned SWF.

Once again, we interpret this fact through the curse of dimensionality: while we saw in our toy GMM example that using a pre-trained SWF was totally working for small dimensions, it is already not so for $d = 48$ and only 3000 training samples.

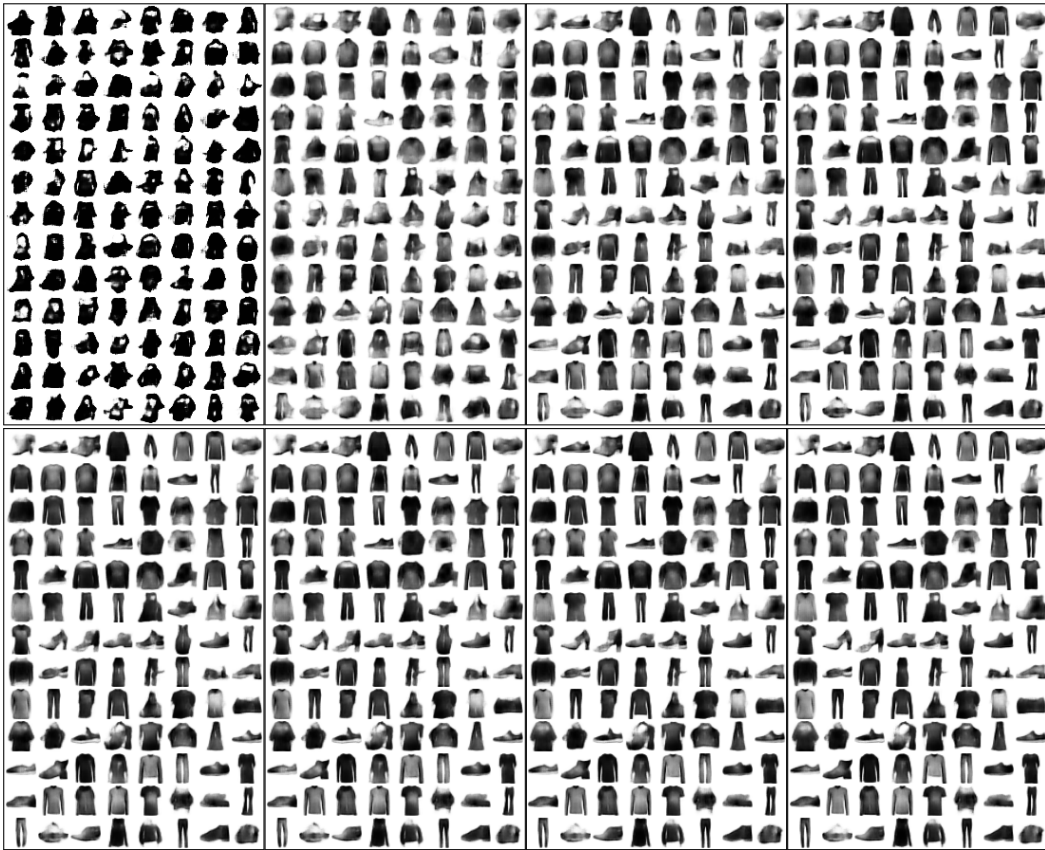
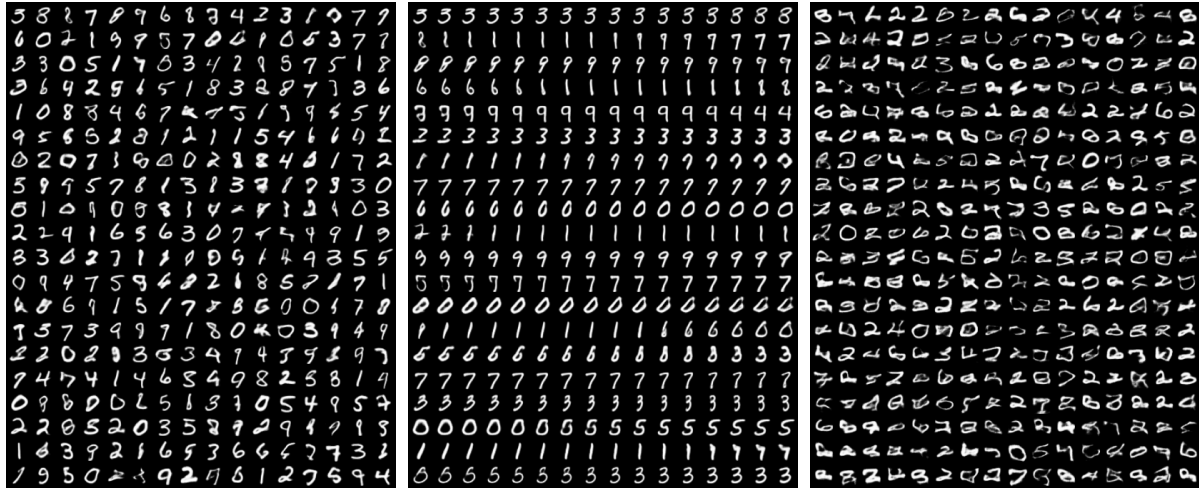


Figure S4. The evolution of SWF through 200 iterations on the FashionMNIST dataset. Plots are for 1, 11, 21, 31 (upper row) and 41, 51, 101, 201 (lower row) iterations

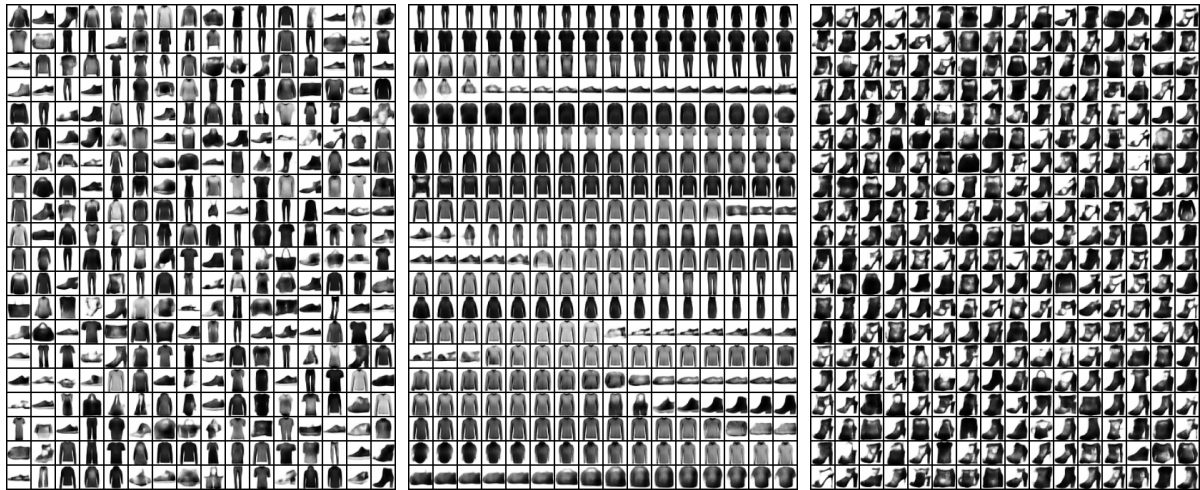


(a) particles undergoing SWF

(b) After SWF is done: applying learned map on linear combinations of train particles
(c) After SWF is done: applying learned map on random inputs.

Figure S5. SWF on MNIST: training samples, interpolation in learned mapping, extrapolation.

This noticed, we highlight that this generalization weakness of SWF for high dimensions is not really an issue, since it is always possible to i/ run SWF with more training samples if generalization is required ii/ re-run the algorithm for a set of new particles. Remember indeed that this does not require passing through the data again, since the distribution of the data projections needs to be done only once.



(a) particles undergoing SWF

(b) After SWF is done: applying learned map on linear combinations of train particles

(c) After SWF is done: applying learned map on random inputs.

Figure S6. SWF on FashionMNIST: training samples, interpolation in learned mapping, extrapolation.