

GUIDING TIME-VARYING GENERATIVE MODELLING WITH NATURAL GRADIENTS ON EXPONENTIAL FAMILY MANIFOLD

Song Liu¹, Leyang Wang², Yakun Wang¹

¹School of Mathematics, University of Bristol, UK

²Department of Computer Science, University College London, UK

{song.liu, yakun.wang}@bristol.ac.uk

leyang.wang.24@ucl.ac.uk

ABSTRACT

Optimising probabilistic models is a well-studied field in statistics. However, its connection with the training of generative models remains largely under-explored. In this paper, we show that the evolution of time-varying generative models can be projected onto an exponential family manifold, naturally creating a link between the parameters of a generative model and those of a probabilistic model. We then train the generative model by moving its projection on the manifold according to the natural gradient descent scheme. This approach also allows us to approximate the natural gradient of the KL divergence efficiently without relying on MCMC for intractable models. Furthermore, we propose particle versions of the algorithm, which feature closed-form update rules for any parametric model within the exponential family. Through toy and real-world experiments, we validate the effectiveness of the proposed algorithms. The code of the proposed method could be found at <https://github.com/anewgithubname/iNGD>.

1 INTRODUCTION

Modern generative models (Goodfellow et al., 2014; Ho et al., 2020; Song et al., 2021b) have become indispensable tools in modern machine learning, achieving remarkable success in applications (Rombach et al., 2022; Gu et al., 2022; Li et al., 2019a; Tan et al., 2024). These generative models are neural networks transforming a latent variable to a higher dimensional sample. They overcome *classic restrictions imposed on probability density models*, such as positivity, normalization, and encoding of conditional independence via factorization; thus, they can be designed freely to capture complex, intricate patterns from the high-dimensional data.

Albeit these models produce highly realistic outputs, they can be hard to train. The model training requires massive data and maintains a tricky balance between “generators” and “discriminators” (Goodfellow et al., 2014; Arjovsky et al., 2017; Wang et al., 2023) or building effective “bridges” between the reference and target dataset (Ho et al., 2020; Song et al., 2021a;b; Lipman et al., 2023; Liu et al., 2023; Ou et al., 2025), which are quite undesirable in the context of specific domain applications when these techniques cannot be applied.

Parametric probabilistic models, particularly those in the exponential family (Casella & Berger, 2024; Wainwright et al., 2008), play a central role in modern day’s statistical inference, and have well-established theoretical framework and training algorithms. These models, characterised by their sufficient statistics and natural parameters, define a *probability distributions manifold*. The geometric structure of which inspired efficient optimisation methods such as natural gradient descent (NGD) Amari (1998), ensuring stable and efficient parameter estimation (Amari & Nagaoka, 2000). However, while exponential family models are versatile *in theory*, their use may be limited in practice: the hand-crafted sufficient statistic may fail to capture complex relationships in data; more flexible sufficient statistics (such as neural nets) result in intractable likelihoods (see Section 2.2). Thus, parameter estimation that requires a likelihood, such as NGD, cannot be easily applied to fit the model.

We aim to unlock the power of modern generative models through the principled training of probabilistic models.

Our work is inspired by two distinct research directions developed in recent years: time-varying generative models (Ho et al., 2020; Song et al., 2021b; Liu et al., 2023; Lipman et al., 2023) and Time Score Matching (TSM) (Choi et al., 2022). Time-varying generative models generate samples progressively, evolving them over time until they match the target distribution. Meanwhile, TSM learns the instantaneous change of a time-varying distribution from data. Recent work demonstrates that TSM can “project” temporal variations in a dataset onto the parameter space of exponential family distributions (Williams et al., 2025).

The core idea of this paper is to evolve a generative model such that its *projected trajectory* on an exponential family manifold aligns with the trajectory induced by NGD.

This way, we get the expressiveness of a modern generative model, and the training efficiency of NGD. The exponential family model acts as a guiding framework for the generative model throughout the training process.

An illustration of this idea could be seen in Fig. 1. We align the projected changes of the generative model with the NGD update on a parameter manifold (shrinking the length of the dotted line in Fig. 1). First, we apply TSM to project the temporal evolution of the generative model onto the manifold. Second, we match this projected evolution to the NGD update by updating the parameters of the generative model. Then generating samples from the updated generator. This process repeats until convergence.

Although this technique applies to all time-varying generative models, in this paper, we focus on drift-based generative models, where samples are iteratively perturbed by vector-valued functions. Specifically, we develop two NGD-guided drift-based generative models: kernel NGD and neural tangent kernel NGD, both of which admit closed-form expressions for sample updates.

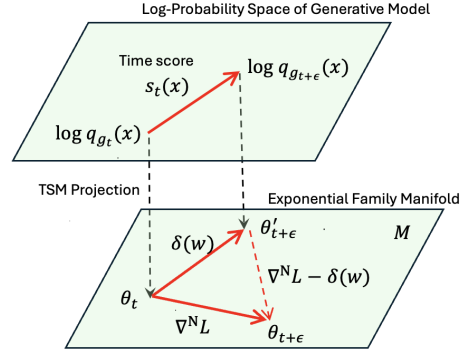


Figure 1: Technical notions illustrated. Matching the projected generative model change $\delta(w)$ to the natural gradient descent meaning reducing the length of the red dotted line. Symbols are defined in Section 3.

2 BACKGROUNDS

In recent years, it has been recognized that there are many similarities between sampling and optimization (Wibisono, 2018; Chewi, 2024; Cheng et al., 2018; He et al., 2025). Encouraged by these results, we ask: Given a time-varying generative model g_t , where t is time, how can we move its output distribution q_{g_t} toward an optimal distribution q^* using NGD? In contemporary machine learning literature, generative models are often referred to as implicit models and generative distribution q_{g_t} normally don’t have a parametric density¹, thus optimization designed for parametric densities cannot be directly applied.

In this paper, we propose an alternative: rather than directly optimizing q_{g_t} , we move its projection on a chosen manifold of exponential family distributions toward optimality. We select this exponential family to be sufficiently expressive, aiming to approximate the full range of distributions that the generative model can produce.

Before stating our ideas, we first formally introduce the time-varying exponential family, NGD, and Time Score Matching (TSM), which are essential for understanding the proposed algorithm. A summary of notation used in this paper can be found in Appendix A.1

¹Some generative models, such as normalizing flows (Rezende & Mohamed, 2015), neural ODEs (Chen et al., 2018), and probability flow models (Song et al., 2021b), do admit explicit parametric forms. However, in this paper, we consider a more general class of generative models.

2.1 EXPONENTIAL FAMILY DISTRIBUTION

Definition 2.1. (See, e.g., Section 3.4, Casella & Berger (2024)) A distribution belongs to the **exponential family** with sufficient statistic T if and only if its density q can be expressed as:

$$q(\mathbf{x}; \boldsymbol{\theta}) = \exp(\langle \boldsymbol{\theta}, T(\mathbf{x}) \rangle - A(\boldsymbol{\theta})),$$

where $\boldsymbol{\theta}$ is the natural parameter, and $A(\boldsymbol{\theta})$ is the log-normalisation function, defined as:

$$A(\boldsymbol{\theta}) = \log \int \exp(\langle \boldsymbol{\theta}, T(\mathbf{x}) \rangle) d\mathbf{x}.$$

This family includes many known distributions, such as Gaussian, Gamma and Exponential distributions.

Definition 2.2. The parameter manifold of exponential family distributions with a chosen sufficient statistic T is defined as:

$$\mathcal{M}(T) = \left\{ \boldsymbol{\theta} \in \mathbb{R}^d \mid q(\mathbf{x}; \boldsymbol{\theta}) > 0, \int q(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} = 1 \right\},$$

where $q(\mathbf{x}; \boldsymbol{\theta})$ depends on T as in Definition 2.1.

More details on the time-varying exponential family can be found in Appendix A.2.

2.2 NATURAL GRADIENT DESCENT AND NATURAL GRADIENT FLOW

Definition 2.3. The natural gradient of a loss function $\mathcal{L}(\boldsymbol{\theta})$ is the gradient of the loss function scaled by the inverse Fisher information matrix \mathcal{F} .

$$\nabla^N \mathcal{L}(\boldsymbol{\theta}) := \mathcal{F}^{-1} \nabla \mathcal{L}(\boldsymbol{\theta}).$$

Assuming $\boldsymbol{\theta} \in \mathcal{M}(T)$, $\mathcal{F} = \mathbb{E}_{q_{\boldsymbol{\theta}}} [-\nabla_{\boldsymbol{\theta}}^2 \log q(\mathbf{x}; \boldsymbol{\theta})] = \text{Cov}_{q_{\boldsymbol{\theta}}} [T(\mathbf{x})]$, describing the curvature of $\mathcal{M}(T)$ at the neighbourhood of $\boldsymbol{\theta}$ (See Section 2.2 in (Amari & Nagaoka, 2000)).

The continuous-time limit of NGD, known as Natural Gradient Flow, describes a time-varying distribution $q_{\boldsymbol{\theta}_t}$, whose parameters trace a trajectory on $\mathcal{M}(T)$ following the steepest descent direction to minimize \mathcal{L} . More discussions can be found in Appendix A.3.

2.3 TIME-VARYING GENERATIVE MODEL

In recent years, there has been a growing trend of designing generative models as functions of time, in contrast to the classic generative models where the sample generating mechanism is independent of time. For example, the diffusion generative model (Song et al., 2021b) can be interpreted as the solution to a Stochastic Differential Equation at time t , with initial samples drawn from a reference distribution. Similarly, the rectified flow generative model (Liu et al., 2023) can be viewed as the solution to an Ordinary Differential Equation at time t .

Definition 2.4. A **generative model** is a data-generating mechanism defined as $X = g(Z; \mathbf{w})$, where Z is a random variable sampled from a latent distribution p_Z , and X is the output of the deterministic function g applied to Z . The parameters $\mathbf{w} \in \mathcal{W}$ parametrise the generative model g . A **time-varying generative model** is a generative model with time dependency, defined as $X_t = g(Z, t; \mathbf{w})$, where the function g explicitly depends on time t .

2.4 TIME SCORE MATCHING

Now we introduce how to measure the change of distributions over time using samples.

Definition 2.5. **Time score** is the time derivative of the log density of a time-varying distribution. Given a time-varying distribution q_t , its time score is $s_t := \partial_t (\log q_t)$.

Given a parametric time score model $v(\mathbf{x}; t)$ and a time-varying sample $X_t \sim q_t$, the time score can be learned by the **Time Score Matching (TSM)** which minimizes the objective: $\int_t \mathbb{E} [\lambda(t) \|s_t(X_t) - v(X_t; t)\|^2] dt$, where $\lambda(t)$ is a weighting function.

The time score plays an important role in the density ratio estimation as $\log(q_{t_1}/q_{t_0}) = \int_{t_0}^{t_1} s_t dt$. Choi et al. leverages this fact and estimates the time score at intermediate distributions and aggregates them to obtain the overall ratio. Moreover, Williams et al. (2025) shows that for the special case where $q_t(\mathbf{x}) \in \mathcal{M}(T)$, TSM can directly learn the time differential natural parameter $\partial_t \boldsymbol{\theta}(t)$.

3 PROPOSED ALGORITHM: EVOLUTION PROJECTION

Now, we combine the concepts introduced in the earlier sections and propose a novel generative training method.

Suppose we only have access to a target distribution p through samples $Y \sim p$ and a latent variable Z . Our goal is to find a time-varying generative model g_t that progressively approximates the target distribution as $t \rightarrow \infty$. Informally speaking, we seek a model such that $Y \stackrel{d}{\approx} g_\infty(Z)$.

Note that we do not want to train a generative model using GANs or diffusion models, as these models require a large number of samples and are hard to train. Instead, given a loss function $\mathcal{L}(p, q_{\theta_t})$ that measures the difference between p and a time-varying probabilistic model q_{θ_t} , we aim to guide the training of the generative model by minimizing the loss \mathcal{L} over time.

The key idea of this paper is to align the evolution of the generative model with $\nabla^N \mathcal{L}(\theta)$ on the manifold $\mathcal{M}(T)$. Consequently, minimizing the loss for q_θ drives the generative model toward matching p .

In the section below, we establish a direct correspondence between the instantaneous change in the generative model and the parametric update on $\mathcal{M}(T)$.

3.1 PROJECTING THE CHANGE

Denote the sample of a time-varying generative model g_t as $X_t \sim q_{g_t}$. We can measure the instantaneous change of the generative model via the its time score $s_t := \partial_t(\log q_{g_t})$.

At a fixed time t_0 , we can “project” the time score s_t onto the manifold $\mathcal{M}(T)$ by minimizing the squared difference between s_t and the time score of an exponential family distribution, $\partial_t(\log q_{\theta_t}), q_{\theta_t} \in \mathcal{M}(T)$,

$$\begin{aligned} & \int \lambda_{t_0}(t) \mathbb{E} \left[(s_t(X_t) - \partial_t(\log q_{\theta_t})(X_t))^2 \right] dt, \\ &= \int \lambda_{t_0}(t) \mathbb{E} \left[(s_t(X_t) - \langle \partial_t \theta(t), T(X_t) - \mathbb{E}[T(X'_t)] \rangle)^2 \right] dt, \end{aligned} \quad (1)$$

where $\lambda_{t_0} = \exp(-(t - t_0)^2/\sigma^2)$ and X'_t is an independent copy of X_t . σ is a hyperparameter fixed in advance. The second line is due to Eq. (10). The integration is over the entire real number domain.

Introducing a linear-in-time model $\theta(t; \delta) = t\delta$, the above objective becomes

$$J(\delta) := \int \lambda_{t_0}(t) \mathbb{E} \left[(s_t(X_t) - \langle \delta, T(X_t) - \mathbb{E}[T(X'_t)] \rangle)^2 \right] dt, \quad (2)$$

which now depends on the parameter δ . One can view (2) as a local regression at the fixed time point t_0 : λ is a time smoothing kernel, Eq. (2) finds the best score model that approximates the time-varying function $s_t(x)$ at time t_0 .

The minimizer to the above objective is a vector in \mathbb{R}^d , the tangent vector that best describes the instantaneous change in the generative model. Now, we provide a closed form expression of such tangent vector.

Theorem 3.1. *Let $\delta_{t_0} := \operatorname{argmin} J(\delta)$. Then δ_{t_0} is unique if the Fisher information matrix $\mathcal{F}_t = \operatorname{Cov}[T(X_t)]$ is invertible and has a closed form expression*

$$\delta_{t_0} = - \left(\int \lambda_{t_0}(t) \operatorname{Cov}[T(X_t)] dt \right)^{-1} \int \partial_t \lambda_{t_0}(t) \mathbb{E}[T(X_t)] dt.$$

The proof can be found in Appendix B. Notice that X_t is generated through $g(Z, t; \mathbf{w})$, we can express δ_{t_0} as a function of \mathbf{w} using the reparameterization trick.

$$\delta_{t_0}(\mathbf{w}) = -C^{-1} \int \partial_t \lambda_{t_0}(t) \mathbb{E}[T(g(Z, t; \mathbf{w}))] dt, C = \int \lambda_{t_0}(t) \operatorname{Cov}[T(g(Z, t; \mathbf{w}))] dt. \quad (3)$$

Algorithm 1 Generative Model Training Guided by Natural Gradient Descent**Require:** Target samples $Y \sim p$, latent samples $Z \sim p_Z$, step size ϵ , number of iterations n

```

1:  $t[0] \leftarrow 0$ , initialize  $\mathbf{w}$ .
2: for  $i \leftarrow 1$  to  $n$  do
3:   Sample  $X_{t[i]} \leftarrow g(Z, t[i]; \mathbf{w})$ 
4:   Approximate  $\nabla^N \mathcal{L}(\boldsymbol{\theta}_{t[i]})$  with  $X_{t[i]}$  and  $Y$  using Monte Carlo.
5:   Approximate  $\boldsymbol{\delta}(\mathbf{w})$  with  $X_{t[i]}$  using Monte Carlo.
6:    $\mathbf{w} \leftarrow \underset{\mathbf{w}}{\operatorname{argmin}} \|\nabla^N \mathcal{L}(\boldsymbol{\theta}_{t[i]}) - \boldsymbol{\delta}(\mathbf{w})\|^2$ 
7:    $t[i+1] \leftarrow t[i] + \epsilon$ 
8: end for
9: Return: Samples from  $g(Z, t[n]; \mathbf{w})$ 

```

This expression allows $\mathbb{E}[\cdot]$ and $\operatorname{Cov}[\cdot]$ to be approximated using samples of Z . The projection process could be seen from Figure 1 where the downward dotted arrow represents the projection by TSM.

Remarks: Different time scores can be mapped to the same $\boldsymbol{\delta}_{t_0}$, especially when the sufficient statistic T is restrictive. However, if T is chosen such that the exponential family is expressive enough, we expect to avoid such information collapse. The rigorous proof of this claim is left as a future work.

The hyperparameter σ in Eq. (1) will introduce additional biases to the estimation, similar to how a non-zero bandwidth in local regression introduces biases to the estimate. In experiments, we observe that reasonable σ choices (e.g., 0.1) work well. Moreover, in Section 4, we show how to obtain an unbiased estimator for a special type of time-varying generative models.

3.2 MATCHING TO NGD

We can see the projection $\boldsymbol{\delta}_{t_0}(\mathbf{w})$ creates a connection between the evolution of a generative model and those of a probabilistic model. The key idea of this paper is to align the evolution of the generative model with that of the probabilistic model. In particular, we match $\boldsymbol{\delta}_{t_0}(\mathbf{w})$ to the NGD update using the following objective

$$\mathbf{w}_{t_0} = \underset{\mathbf{w} \in \mathcal{W}}{\operatorname{argmin}} \|\mathcal{F}_{t_0}^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_{t_0}) - \boldsymbol{\delta}_{t_0}(\mathbf{w})\|^2. \quad (4)$$

In words, we find the generative model parameter \mathbf{w} that results in the projected update $\boldsymbol{\delta}_{t_0}(\mathbf{w})$ closest to the natural gradient update of the loss function.

In the previous section, we have seen that $\boldsymbol{\delta}_{t_0}(\mathbf{w})$ can be approximated using samples Z . Assume that at the starting point, the generator $g(Z, 0; \mathbf{w}_0)$ produces an output distribution q_0 on the manifold, and its movement is precisely tracked by the trajectory of its projection, we can expect that the samples generated from $g(Z, t; \mathbf{w}_t)$ will be close to the samples generated from $q_{\boldsymbol{\theta}_t}$. Thus, we approximate $\mathcal{F}_{t_0}^{-1} \nabla_{\boldsymbol{\theta}}$ using samples from $g(Z, t; \mathbf{w}_t)$.

After solving Eq. (4), instead of taking an actual natural gradient step, we directly sample from $g(Z, t_0 + \epsilon; \mathbf{w}_{t_0})$ using a small $\epsilon > 0$. Since we have already aligned the projected change of our time-varying generator with the natural gradient step, we can expect that the samples generated from $g(Z, t_0 + \epsilon; \mathbf{w}_{t_0})$ will be close to the samples generated from $q(\mathbf{x}; \boldsymbol{\theta}_{t_0} - \epsilon \nabla^N \mathcal{L})$ by actually taking a natural gradient step with step size ϵ . We summarize the entire algorithm in Algorithm 1 and name it implicit NGD (iNGD). A visualization of iNGD is also provided in Appendix E.

4 SPECIAL CASE: DRIFT MODEL

One popular class of generative models is the *drift-based generative model*. This model iteratively perturbs samples using a vector-valued function until convergence.

Algorithm 2 RKHS/Neural Tangent Kernel Implicit Natural Gradient Descent**Require:** Target samples $Y \sim p$, initial particles $X_0 \sim q_0$, step size ϵ , number of iterations n

```

1:  $t[0] \leftarrow 0$ 
2: for  $i \leftarrow 1$  to  $n$  do
3:   Approximate  $\mathbf{h}$  with  $Y$  and  $X_{t[i]}$  using (7) or (9).
4:    $t[i+1] = t[i] + \epsilon$ 
5:    $X_{t[i+1]} = X_{t[i]} + \epsilon \mathbf{h}(X_{t[i]})$ 
6: end for
7: Return:  $X_{t[n]} \sim q_{t[n]}$ .
```

Definition 4.1. At a fixed time t_0 , suppose we have samples X_{t_0} , a drift generative model generates samples at a new time point t via the following scheme:

$$g(X_{t_0}, t; \mathbf{w}) = X_{t_0} + (t - t_0) \mathbf{h}(X_{t_0}; \mathbf{w}).$$

This model could be seen as a local linear model and the amount of update $\mathbf{h}(X_{t_0}; \mathbf{w})$ depends linearly on $t - t_0$. The input of the model is a sample at the time point t_0 . To generate samples, we need to draw a batch of samples from an initial distribution q_0 and then successively apply the drift generative model until convergence. An Euler solver of a flow-based generative model is an example of a drift generative model, in which case, $t - t_0$ is the step size of the Euler solver. We can prove that when using a model introduced in Definition 4.1, Eq. (3) has a limiting solution as $\sigma \rightarrow 0$, which eliminates the bias caused by the time-smoothing kernel λ .

Theorem 4.2. *The projection of the time score s_t of a drift generative model onto the manifold $\mathcal{M}(T)$ has a closed form expression at the limit of $\sigma \rightarrow 0$, i.e.,*

$$\lim_{\sigma \rightarrow 0} \delta_{t_0}(\mathbf{w}) = \text{Cov}[T(X_{t_0})]^{-1} \mathbb{E} [\nabla T(X_{t_0}) \mathbf{h}(X_{t_0}; \mathbf{w})] = \mathcal{F}_{t_0}^{-1} \mathbb{E} [\nabla T(X_{t_0}) \mathbf{h}(X_{t_0}; \mathbf{w})].$$

The proof can be found in Appendix C. Using this limiting solution, the objective of Eq. (4) can be rewritten as:

$$\mathbf{w}_{t_0} = \underset{\mathbf{w} \in \mathcal{W}}{\text{argmin}} \|\nabla \mathcal{L}(\boldsymbol{\theta}_{t_0}) - \mathbb{E} [\nabla^\top T(X_{t_0}) \mathbf{h}(X_{t_0}; \mathbf{w})]\|_{\mathcal{F}_{t_0}^{-1}}^2. \quad (5)$$

In our experiments, this objective function is more stable and computationally efficient than (4), since it does not require back-propagating through a matrix inversion.

4.1 KERNEL NGD

Now let us consider an example of the drift model where the drift function \mathbf{h} is defined as the gradient of a Reproducing Kernel Hilbert Space (RKHS) function.

Example 1 (RKHS Drift Model). A kernel drift function is defined as

$$\mathbf{h}(\mathbf{x}; \mathbf{w}) := \langle \mathbf{w}, \nabla_{\mathbf{x}} k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}, \mathbf{w} \in \mathcal{H},$$

where \mathcal{H} is an RKHS with a kernel function k .

To align this generative model with the NGD, we introduce a regularized version of Eq. (5)

$$\mathbf{w}_{t_0} = \underset{\mathbf{w} \in \mathcal{H}}{\text{argmin}} \|\nabla \mathcal{L}(\boldsymbol{\theta}_{t_0}) - \mathbb{E} [\nabla T(X_{t_0}) \mathbf{h}_w(X_{t_0})]\|_{\mathcal{F}_{t_0}^{-1}}^2 + \lambda \|\mathbf{w}\|_{\mathcal{H}}^2. \quad (6)$$

Theorem 4.3. *The optimal drift function that minimizes Eq. (6) can be found as*

$$\begin{aligned} \mathbf{h}_{\mathbf{w}_{t_0}}(\mathbf{x}) &= \mathbb{E} [\nabla T(X_{t_0}) \nabla \nabla k(X_{t_0}, \mathbf{x})]^\top \boldsymbol{\Gamma}^{-1} \nabla \mathcal{L}(\boldsymbol{\theta}_{t_0}) \\ \boldsymbol{\Gamma} &= \lambda \mathcal{F}_{t_0} + \mathbb{E} [\nabla T(X_{t_0}) \nabla \nabla k(X_{t_0}, X'_{t_0}) \nabla^\top T(X'_{t_0})] \end{aligned} \quad (7)$$

where X'_{t_0} is an independent copy of X_{t_0} and $\nabla \nabla k(\mathbf{x}, \mathbf{y})$ is the short hand for $\nabla_{\mathbf{x}} \nabla_{\mathbf{y}} k(\mathbf{x}, \mathbf{y})$.

The proof can be found in Appendix D. This result enables the *direct calculation* of particle updates without fitting a generative model first. This inspires us to build a *particle evolution strategy* guided by NGD: First, we sample X_0 from an initial distribution q_0 , then we iteratively update each sample X_{t_0} using the formula given by Theorem 4.3 until they converge. This algorithm is summarized in Algorithm 2 and we name this algorithm Kernel implicit NGD (KiNG).

4.2 NEURAL TANGENT KERNEL NGD

KiNG can be generalized to other types of kernels.

Example 2 (Neural Tangent Drift Model). Given a neural network $\phi : \mathbb{R}^{\dim(\mathbf{x})} \rightarrow \mathbb{R}^{\dim(\mathbf{x})}$, a neural tangent drift function is defined as

$$\mathbf{h}(\mathbf{x}) := \nabla_{\beta} \phi(\mathbf{x}; \beta_0) \mathbf{w},$$

where β_0 are initial weights, $\nabla_{\beta} \phi(\mathbf{x}; \beta_0)$ is called a neural tangent.

Similar to Eq. (6), we can solve the following regularized objective to find the update of the particles

$$\mathbf{w}_{t_0} = \underset{\mathbf{w} \in \mathcal{W}}{\operatorname{argmin}} \|\nabla \mathcal{L}(\boldsymbol{\theta}_{t_0}) - \mathbb{E}[\nabla T(X_{t_0}) \mathbf{h}_{\mathbf{w}}(X_{t_0})]\|_{\mathcal{F}_0^{-1}}^2 + \lambda \|\mathbf{w}\|^2. \quad (8)$$

The optimal drift that minimizes Eq. (8) can be expressed using the *neural tangent kernel* (Jacot et al., 2018):

$$\begin{aligned} \mathbf{h}_{\mathbf{w}_{t_0}}(\mathbf{x}) &= \mathbb{E}[\nabla T(X_{t_0}) \mathbf{K}_{\text{NTK}}(X_{t_0}, \mathbf{x})]^\top \mathbf{\Gamma}^{-1} \nabla \mathcal{L}(\boldsymbol{\theta}_{t_0}) \\ \mathbf{\Gamma} &:= \lambda \mathcal{F}_{t_0} + \mathbb{E}[\nabla T(X_{t_0}) \mathbf{K}_{\text{NTK}}(X_{t_0}, X'_{t_0}) \nabla^\top T(X'_{t_0})], \end{aligned} \quad (9)$$

where \mathbf{K}_{NTK} is the matrix-valued *neural tangent kernel*, defined as $\mathbf{K}_{\text{NTK}}(\mathbf{x}, \mathbf{y}) := \nabla_{\beta} \phi(\mathbf{x}) \nabla_{\beta}^\top \phi(\mathbf{y})$. In this paper, we use empirical and a finite-width NTK for simplicity, but NTKs that are infinitely wide can be efficiently computed using off-the-shelf package such as `neural-tangents` (Novak et al., 2020) for a variety of neural network architectures. We name this variant of KiNG as ntKiNG.

Remark: Eq. (9) requires computing a matrix-valued kernel, which may be computationally demanding if the dimensionality of X_{t_0} is high. However, in experiments, we observe that the formulation works well with a diagonalized scalar kernel, i.e.,

$$[\mathbf{K}(X_{t_0}, X'_{t_0})]_{l,m \in [1, \dim(\mathbf{x})]} := \begin{cases} k(X_{t_0}, X'_{t_0}), & l = m \\ 0, & l \neq m, \end{cases}$$

where k is any scalar kernel function. As both KiNG and ntKiNG involve matrix inversion of $\mathbf{\Gamma}$, the computation complexity of Algorithm 2 is $\mathcal{O}(\dim(T)^3 \# \text{particles}^2)$.

5 EXPERIMENTS

5.1 ILLUSTRATIVE EXAMPLES

In Fig. 2a, we plot the trajectories of KiNG with different one-dimensional initial and target distributions. We choose $\mathcal{M}(T)$ to be a Gaussian manifold, i.e., $T(x) = [x, x^2]^\top$. Note that in the right plot, the particles do not converge to the bi-modal target distribution, since the movements of our particles are restricted by the Gaussian manifold, and in this case, the best approximation is a Gaussian with a larger standard deviation. This example shows that the particles sticks to the Gaussian manifold throughout the generative process. This phenomenon could be beneficial in some applications, if the target is to find the best approximation within a given family.

This behaviour could be changed by replacing the Gaussian manifold with a more expressive manifold. In the left plot of Fig. 2b, we run a similar experiments by letting T be the Radial Basis Functions (RBFs) and we can see that indeed the particles bifurcate and converge to both modes. The right plot of Fig. 2c shows the particle trajectory of a neural tangent KiNG with T as RBF basis.

5.2 COMPARISON WITH REVERSE KL WASSERSTEIN GRADIENT FLOW AND MMD FLOW

In this experiment, we compare KiNG, ntKiNG, with reverse KL Wasserstein Gradient Flow (WGF) (Gao et al., 2019; Liu et al., 2024) and Maximum Mean Discrepancy (MMD) flow (Hagemann et al., 2024) on small datasets with different dimensions. See Appendix F for more explanations.

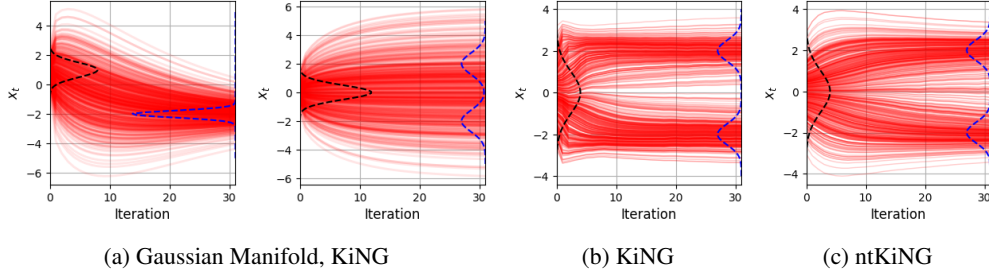


Figure 2: (a) The evolution of particles X_t under KiNG on the manifold of Gaussian distributions. Each red line is a trajectory of a particle in the space-time. The initial distribution q_0 is plotted on the left as black dotted lines, while the target distribution p is plotted on the right as blue dotted lines. (b, c) **Particle trajectories when using more expressive manifold**, i.e., $T(\mathbf{x}) := [k(x_1, b_1), k(x_1, b_2), \dots]^\top$, where k is a RBF basis function and b_i are kernel basis randomly chosen from samples of X_{t_0} . (b) KiNG, (c) ntKiNG.

Since they all minimize different divergences, we measure their performance using MMD (Gretton et al., 2012) between a fresh batch of target samples and X_t . Let $p = 0.5\mathcal{N}(-2, \mathbf{I}) + 0.5\mathcal{N}(2, \mathbf{I})$. We draw 100 samples from p as Y , 100 samples from $\mathcal{N}(0, \mathbf{I})$ as X_0 , and run Algorithm 2, WGF, MMD flow to evolve particles X_t . We plot $\text{MMD}(Y, X_t)$ over iterations.

For all methods, we set the learning rate to be 1, which is the largest learning rate without causing numerical instability. It can be seen that when dimension is small (5), all methods work relatively well and ntKiNG and KiNG can reduce MMD faster, but when we increase the dimension to 20 the performance gap widens. However, ntKiNG and KiNG still have a commanding lead.

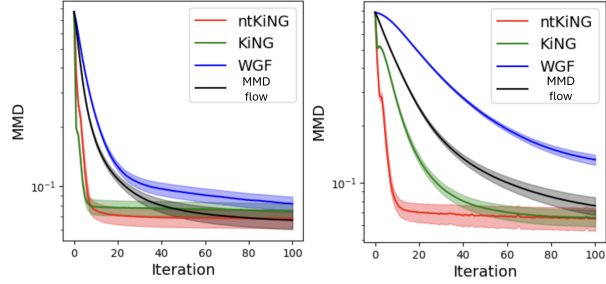


Figure 3: $\text{MMD}[Y, X_t]$ over iterations, the lower the better. Left: 5 dimensions, Right: 20 dimensions. The error bar indicates the standard error.

5.3 GRAPHICAL MODEL RECOVERY USING INFORMATIVE SUFFICIENT STATISTICS

In this experiment, we demonstrate how an informative sufficient statistic could help accelerate the training of the generative model using ntKiNG. Details can be found in Appendix G.

5.4 COVARIATE SHIFT BY DISTRIBUTION MATCHING

We further test our algorithm on the Office+Caltech dataset (Gong et al., 2012) which is an object recognition dataset with photos collected from four different places: amazon, ds1r, webcam, caltech. The task is to train a source classifier using one of the places, and test it on samples from another place. We test the performance of ntKiNG against two other particle-based transport methods WGF and MMD flow by matching q_t with p_X . Further details can be found in Appendix H.

6 RELATED WORKS

We compare our proposal with some related works in Appendix J.

REFERENCES

- Shun-Ichi Amari. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276, 1998.
- Shun-ichi Amari and Hiroshi Nagaoka. *Methods of information geometry*, volume 191. American Mathematical Soc., 2000.
- Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2008.
- M Arbel, A Gretton, W Li, and G Montufar. Kernelized wasserstein natural gradient. In *International Conference on Learning Representations*, 2020.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, 2017.
- George Casella and Roger Berger. *Statistical inference, 2nd Edition*. CRC press, 2024.
- Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in Neural Information Processing Systems*, 2018.
- Zhichao Chen, Haoxuan Li, Fangyikang Wang, Odin Zhang, Hu Xu, Xiaoyu Jiang, Zhihuan Song, and Hao Wang. Rethinking the diffusion models for missing data imputation: A gradient flow perspective. In *Advances in Neural Information Processing Systems*, 2024.
- Xiang Cheng, Niladri S Chatterji, Peter L Bartlett, and Michael I Jordan. Underdamped langevin mcmc: A non-asymptotic analysis. In *Conference on learning theory*, 2018.
- Sinho Chewi. *Log-Concave Sampling*. 2024.
- Sinho Chewi, Thibaut Le Gouic, Chen Lu, Tyler Maunu, and Philippe Rigollet. Svgd as a kernelized wasserstein gradient flow of the chi-squared divergence. In *Advances in Neural Information Processing Systems*, 2020.
- Jaemoo Choi, Jaewoong Choi, and Myungjoo Kang. Scalable wasserstein gradient flow for generative modeling through unbalanced optimal transport. In *International Conference on Machine Learning*, 2024.
- Kristy Choi, Chenlin Meng, Yang Song, and Stefano Ermon. Density ratio estimation via infinitesimal classification. In *International Conference on Artificial Intelligence and Statistics*, 2022.
- Nicolas Courty, Rémi Flamary, Devis Tuia, and Alain Rakotomamonjy. Optimal transport for domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):1853–1865, 2016.
- Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2007.
- Yuan Gao, Yuling Jiao, Yang Wang, Yao Wang, Can Yang, and Shunkang Zhang. Deep generative learning via variational gradient flow. In *International Conference on Machine Learning*, 2019.
- Thomas George, César Laurent, Xavier Bouthillier, Nicolas Ballas, and Pascal Vincent. Fast approximate natural gradient descent in a kronecker factored eigenbasis. *Advances in Neural Information Processing Systems*, 2018.
- Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *Conference on Computer Vision and Pattern Recognition*, 2012.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in Neural Information Processing Systems*, 2014.
- Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(25):723–773, 2012.

- Roger Grosse and James Martens. A kronecker-factored approximate fisher matrix for convolution layers. In *International Conference on Machine Learning*, 2016.
- Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Conference on Computer Vision and Pattern Recognition*, 2022.
- Paul Hagemann, Johannes Hertrich, Fabian Altekrüger, Robert Beinert, Jannis Chemseddine, and Gabriele Steidl. Posterior sampling based on gradient flows of the MMD with negative distance kernel. In *International Conference on Learning Representations*, 2024.
- Jiajun He, Wenlin Chen, Mingtian Zhang, David Barber, and José Miguel Hernández-Lobato. Training neural samplers with reverse diffusive kl divergence. *arXiv preprint arXiv:2410.12456*, 2025.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 2020.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, 2018.
- Naihan Li, Shujie Liu, Yanqing Liu, Sheng Zhao, and Ming Liu. Neural speech synthesis with transformer network. In *AAAI Conference on Artificial Intelligence*, 2019a.
- Wuchen Li, Alex Tong Lin, and Guido Montúfar. Affine natural proximal learning. In *International Conference on Geometric Science of Information*, 2019b.
- Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *International Conference on Learning Representations*, 2023.
- Qiang Liu. Stein variational gradient descent as gradient flow. In *Advances in Neural Information Processing Systems*, 2017.
- Song Liu, Jiahao Yu, Jack Simons, Mingxuan Yi, and Mark Beaumont. Minimizing f -divergences by interpolating velocity fields. In *International Conference on Machine Learning*, 2024.
- Xingchao Liu, Chengyue Gong, and qiang liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *International Conference on Learning Representations*, 2023.
- James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International Conference on Machine Learning*, 2015.
- Aimee Maurais and Youssef Marzouk. Sampling in unit time with kernel fisher-rao flow. In *International Conference on Machine Learning*, 2024.
- Roman Novak, Lechao Xiao, Jiri Hron, Jaehoon Lee, Alexander A. Alemi, Jascha Sohl-Dickstein, and Samuel S. Schoenholz. Neural tangents: Fast and easy infinite neural networks in python. In *International Conference on Learning Representations*, 2020.
- Zijing Ou, Mingtian Zhang, Andi Zhang, Tim Z Xiao, Yingzhen Li, and David Barber. Diffusion model with optimal covariance matching. In *International Conference on Learning Representations*, 2025.
- Joaquin Quiñero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset shift in machine learning*. Mit Press, 2009.
- Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International Conference on Machine Learning*, 2015.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Conference on Computer Vision and Pattern Recognition*, 2022.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021a.

- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021b.
- Masashi Sugiyama, Shinichi Nakajima, Hisashi Kashima, Paul Buenau, and Motoaki Kawanabe. Direct importance estimation with model selection and its application to covariate shift adaptation. In *Advances in Neural Information Processing Systems*, 2008.
- Xu Tan, Jiawei Chen, Haohe Liu, Jian Cong, Chen Zhang, Yanqing Liu, Xi Wang, Yichong Leng, Yuanhao Yi, Lei He, et al. Naturalspeech: End-to-end text-to-speech synthesis with human-level quality. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- Zhendong Wang, Huangjie Zheng, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. Diffusion-gan: Training gans with diffusion. In *International Conference on Learning Representations*, 2023.
- Max Welling. Kernel ridge regression. <https://web2.qatar.cmu.edu/~gdicaro/10315-Fall19/additional/welling-notes-on-kernel-ridge.pdf>, 2019. Accessed: 2025-02-10.
- Andre Wibisono. Sampling as optimization in the space of measures: The langevin dynamics as a composite optimization problem. In *Conference on Learning Theory*, 2018.
- Daniel J Williams, Leyang Wang, Qizhen Ying, Song Liu, and Mladen Kolar. High-dimensional differential parameter inference in exponential family using time score matching. *International conference on artificial intelligence and statistics*, 2025.

A PRELIMINARIES

A.1 NOTATIONS

Vectors (\mathbf{x}) and matrices (\mathbf{X}) are denoted in bold. Random variables (X) are non-bold, capital letters. Euclidean norms are denoted as $\|\cdot\|$, and for elements in a Hilbert space, the Hilbert norm is written as $\|\cdot\|_{\mathcal{H}}$. Expectations and covariances under q are denoted as $\mathbb{E}_q[\cdot]$ and $\text{Cov}_q[\cdot]$, respectively. $\nabla f(\mathbf{x}) = [\partial_1 f(\mathbf{x}), \partial_2 f(\mathbf{x}), \dots]^\top$. $\nabla f(\mathbf{x}_0)$ means the gradient of f evaluated at \mathbf{x}_0 . Suppose $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n$, $\nabla \mathbf{f}(\mathbf{x})$ denotes the Jacobian of \mathbf{f} and is a matrix of size $n \times m$.

A.2 TIME-VARYING EXPONENTIAL FAMILY

Time-varying exponential family distributions refers to exponential family distributions whose natural parameter is a continuous function of time, i.e., $q_{\theta_t} = q(\mathbf{x}; \theta(t))$. Moreover, the time derivative of $\log q_t$ is:

$$\partial_t \log q_{\theta_t} = \langle \partial_t \theta(t), T(\mathbf{x}) \rangle - \partial_t A(\theta(t)).$$

Equivalently, this can be expressed as the inner product between the rate of change of the natural parameter and the centered sufficient statistic (Proposition 3.1 in Williams et al. (2025)):

$$\partial_t \log q_{\theta_t} = \langle \partial_t \theta(t), T(\mathbf{x}) - \mathbb{E}_{q_{\theta_t}}[T(\mathbf{x})] \rangle. \quad (10)$$

By definition, the time-varying process $\theta(t)$ is a curve on $\mathcal{M}(T)$ and $\partial_t \theta(t)$ is the tangent vector of such curve.

A.2.1 INFINITE-DIMENSIONAL EXPONENTIAL FAMILY

One important class of exponential family is infinite dimensional exponential family.

Example 3. An infinite-dimensional exponential family is an exponential family with a sufficient statistic $T(\mathbf{x}) := k(\mathbf{x}, \cdot)$:

$$q(\mathbf{x}; \theta) = \exp(\langle \theta, k(\mathbf{x}, \cdot) \rangle - A(\theta)),$$

where $A(\theta) := \int \exp(\langle \theta, k(\mathbf{x}, \cdot) \rangle - A(\theta)) d\mathbf{x}$ and θ is in an RKHS \mathcal{H} with kernel k .

Given a dataset of n data points $\{\mathbf{x}_i\}_{i=1}^n \sim q_\theta$, the empirical natural gradient is defined as

Definition A.1.

$$\text{NGD} := \underset{f \in \mathcal{H}}{\text{argmax}} \langle f, \nabla \mathcal{L}(\theta) \rangle_{\mathcal{H}} - \frac{1}{2} \langle f, \hat{\Sigma} f \rangle_{\mathcal{H}} - \lambda \|f\|_{\mathcal{H}}^2, \quad (11)$$

where $\hat{\Sigma}$ is the empirical covariance operator defined as

$$\hat{\Sigma} := \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i)^\top \mathbf{H} \phi(\mathbf{x}_i),$$

and \mathbf{H} is a centering matrix and $\lambda > 0$.

Due to the Representer Theorem, we can see that the optimal solution takes the form $\text{NGD} := \sum_i \alpha_i \phi(\mathbf{x}_i, \cdot)$.

Since our aim is to match the projection $\delta \in \mathcal{H}$ to the natural gradient, we let $\delta = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \cdot)$. The regularized objective of Eq. (2) takes the form

$$J(\boldsymbol{\alpha}) := \int \lambda_{t_0}(t) \mathbb{E} \left[(s_t(X_t) - \langle \boldsymbol{\alpha}, \mathbf{k}(X_t) - \mathbb{E}[\mathbf{k}(X'_t)] \rangle)^2 \right] dt + \lambda \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \quad (12)$$

where $\mathbf{k}(X) := [k(\mathbf{x}_1, X), k(\mathbf{x}_2, X), \dots]$. The limiting solution of Eq. (12) is

$$\lim_{\sigma \rightarrow 0} \boldsymbol{\alpha}_{t_0}(\mathbf{w}) = (\text{Cov}[\mathbf{k}(X_{t_0})] + \lambda \mathbf{K})^{-1} \mathbb{E} \left[\nabla \mathbf{k}(X_{t_0}) \mathbf{h}^\top(X_{t_0}; \mathbf{w}) \right].$$

Substituting the form of NGD and solving (11) for $\boldsymbol{\alpha}$, we get

$$\boldsymbol{\alpha}^* = (\text{Cov}[\mathbf{k}(X_{t_0})] + \lambda \mathbf{K})^{-1} (\mathbb{E}_p[\mathbf{k}(\mathbf{x})] - \mathbb{E}_q[\mathbf{k}(\mathbf{x})]).$$

Thus, the optimal drift parameter could be obtained via $\mathbf{w}^* := \underset{\mathbf{w}}{\text{argmin}} \|\boldsymbol{\alpha}^* - \lim_{\sigma \rightarrow 0} \boldsymbol{\alpha}_{t_0}(\mathbf{w})\|^2$, which has a similar solution to Eq. (7) or Eq. (9) depending on which model is used for the drift.

A.3 THE NATURAL GRADIENT OF KL DIVERGENCE

For easier conceptualization, we here provide a concrete example of NGD for KL divergence.

Example 4. Suppose \mathcal{L} is a KL divergence from p to q .

$$\mathcal{L} = \text{KL}[p, q] = \mathbb{E}_p[\log p(\mathbf{x}) - \log q(\mathbf{x})].$$

Suppose q is an exponential family distribution with sufficient statistic T . Then the natural gradient of \mathcal{L} is given by

$$\begin{aligned} \text{NGD}_{\text{KL}} &:= -\mathcal{F}^{-1} \mathbb{E}_p[\nabla_{\theta} \log q] \\ &= -\text{Cov}_{q_{\theta}}[T(\mathbf{x})]^{-1} \{ \mathbb{E}_p[T(\mathbf{x})] - \mathbb{E}_{q_{\theta}}[T(\mathbf{x})] \}. \end{aligned} \quad (13)$$

Except for certain specific choices of T , NGD_{KL} does not admit a closed-form solution, as neither $\mathbb{E}_{q_{\theta}}[T(\mathbf{x})]$ nor $\text{Cov}_{q_{\theta}}[T(\mathbf{x})]$ can be expressed in closed form for a general T . If we could sample from q_{θ} , these expectations and covariances could be approximated using Monte Carlo methods. However, generating samples from a complex distribution q_{θ} itself remains a challenging problem.

B PROOF OF THEOREM 3.1

Theorem B.1. (Theorem 4.1 in (Williams et al., 2025)) Eq. (2) can be rewritten as the following form

$$\mathcal{L}(\delta) = \int_0^1 \lambda_{t_0}(t) \mathbb{E}[\langle \delta, T(X_t) - \mathbb{E}[T(X_t')] \rangle^2] dt + 2 \int_0^1 \partial_t \lambda_{t_0}(t) \mathbb{E}[\langle \delta, T(X_t) \rangle] dt + \text{const.} \quad (14)$$

Let $\mathbf{B}_t = T(X_t) - \mathbb{E}[T(X_t')]$. Firstly, we find the derivative of the quadratic term. For each t , since $\frac{d}{d\delta} (\langle \delta, \mathbf{B}_t \rangle^2) = 2 \langle \delta, \mathbf{B}_t \rangle \mathbf{B}_t$, it follows that

$$\frac{\partial}{\partial \delta} \int_t \lambda_{t_0}(t) \mathbb{E}[\langle \delta, \mathbf{B}_t \rangle^2] dt = 2 \int_t \lambda_{t_0}(t) \mathbb{E}[\langle \delta, \mathbf{B}_t \rangle \mathbf{B}_t] dt.$$

Since $\mathbb{E}[\mathbf{B}_t] = 0$, one has $\mathbb{E}[\langle \delta, \mathbf{B}_t \rangle \mathbf{B}_t] = \text{Cov}[\mathbf{B}_t] \delta = \text{Cov}[T(X_t)] \delta$.

Hence this part becomes

$$2 \int_t \lambda_{t_0}(t) \text{Cov}[T(X_t)] \delta dt.$$

Derivative of the linear-in- δ term. The term $\int_t 2 \partial_t \lambda_{t_0}(t) \mathbb{E}[\langle \delta, T(X_t) \rangle] dt$ is linear in δ . Its gradient w.r.t. δ is simply

$$2 \int_t \partial_t \lambda_{t_0}(t) \mathbb{E}[T(X_t)] dt.$$

Putting these together, the gradient of $\mathcal{L}(\delta)$ is

$$\nabla_{\delta} \mathcal{L}(\delta) = 2 \int_t \lambda_{t_0}(t) \text{Cov}[T(X_t)] \delta dt + 2 \int_t \partial_t \lambda_{t_0}(t) \mathbb{E}[T(X_t)] dt.$$

To find the minimizer, set this gradient to zero:

$$0 = 2 \int_t \lambda_{t_0}(t) \text{Cov}[T(X_t)] \delta dt + 2 \int_t \partial_t \lambda_{t_0}(t) \mathbb{E}[T(X_t)] dt.$$

which can be rewritten as

$$\left(\int_t \lambda_{t_0}(t) \text{Cov}[T(X_t)] dt \right) \delta = - \int_t \partial_t \lambda_{t_0}(t) \mathbb{E}[T(X_t)] dt.$$

If the matrix $\int_0^1 \lambda_{t_0}(t) \text{Cov}[T(X_t)] dt$ is invertible—which is precisely the non-degeneracy (invertibility) of the Fisher information—then we can solve uniquely for δ :

$$\delta_{t_0} = \left(\int_t \lambda_{t_0}(t) \text{Cov}[T(X_t)] dt \right)^{-1} \left(- \int_t \partial_t \lambda_{t_0}(t) \mathbb{E}[T(X_t)] dt \right).$$

We can conclude

$$\delta_{t_0} = \left(\int_t \lambda_{t_0}(t) \text{Cov}[T(X_t)] dt \right)^{-1} \int_t \partial_t \lambda_{t_0}(t) \mathbb{E}[T(X_t)] dt. \quad (15)$$

C PROOF OF THEOREM 4.2

Recall:

$$\delta_{t_0}(\mathbf{w}) = C^{-1} \int_{-\infty}^{\infty} \partial_t \lambda(t, t_0) \mathbb{E}[T(g(Z, t; \mathbf{w}))] dt, \quad C = \int_{-\infty}^{\infty} \lambda(t, t_0) \text{Cov}[T(g(Z, t; \mathbf{w}))] dt.$$

The first lemma states a few properties of the Gaussian kernel.

Lemma C.1. $\int_{-\infty}^{\infty} \partial_t \lambda(t, t_0) dt = 0$. $\int_{-\infty}^{\infty} (t - t_0) \partial_t \lambda_\sigma(t, t_0) dt = -1$.

Proof. The first result is due to the Fundamental Theorem of Calculus and the fact that $\lambda(t, t_0) \rightarrow 0$ as $|t| \rightarrow \infty$. Now, we prove the second statement. Since

$$\lambda_\sigma(t, t_0) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(t - t_0)^2}{2\sigma^2}\right),$$

we have

$$\int_{-\infty}^{\infty} \lambda_\sigma(t, t_0) dt = 1, \quad \text{and} \quad \partial_t \lambda_\sigma(t, t_0) = -\frac{(t - t_0)}{\sigma^2} \lambda_\sigma(t, t_0).$$

We then have with integration by parts

$$\int_{-\infty}^{\infty} (t - t_0) \partial_t \lambda_\sigma(t, t_0) dt = (t - t_0) \lambda_\sigma(t, t_0) \Big|_{-\infty}^{\infty} - \int_{-\infty}^{\infty} \lambda_\sigma(t, t_0) dt = 0 - 1 = -1. \quad (16)$$

the last equality is due to $\lim_{|t| \rightarrow \infty} (t - t_0) \lambda_\sigma(t, t_0) = 0$. \square

First, we inspect $\int_{-\infty}^{\infty} \partial_t \lambda(t, t_0) \mathbb{E}[T(g(Z, t; \mathbf{w}))] dt$. Using the Taylor expansion on $\mathbb{E}[T(g(Z, t; \mathbf{w}))]$ at t_0 , we obtain

$$\mathbb{E}[T(g(Z, t; \mathbf{w}))] = \mathbb{E}[T(X_{t_0})] + (t - t_0) \mathbb{E}[\nabla T(X_{t_0})^\top \mathbf{h}(X_{t_0}; \mathbf{w})].$$

Note that we don't have higher order terms as the drift model $g(Z, t; \mathbf{w})$ is a linear function of t by definition (see Definition 4.1).

Thus, due to Lemma C.1, we have

$$\int_{-\infty}^{\infty} \partial_t \lambda_\sigma(t, t_0) \left[\mathbb{E}(T(X_{t_0})) + (t - t_0) \mathbb{E}(\nabla T(X_{t_0})^\top \mathbf{h}(X_{t_0}; \mathbf{w})) \right] dt = -\mathbb{E}[\nabla T(X_{t_0})^\top \mathbf{h}(X_{t_0}; \mathbf{w})]. \quad (17)$$

Now we shift our focus on

$$C = \int_{-\infty}^{\infty} \lambda_\sigma(t, t_0) \text{Cov}[T(g(Z, t; \mathbf{w}))] dt.$$

As $\sigma \rightarrow 0$, $\lambda_\sigma(t, t_0)$ converges to $\delta(t - t_0)$, so $\lim_{\sigma \rightarrow 0} \left(\int \lambda(t, t_0) f(t) dt \right) = f(t_0)$. Hence

$$\lim_{\sigma \rightarrow 0} C = \lim_{\sigma \rightarrow 0} \int_{-\infty}^{\infty} \lambda_\sigma(t, t_0) \text{Cov}[T(g(Z, t; \mathbf{w}))] dt = \text{Cov}[T(X_{t_0})]. \quad (18)$$

Finally, combining Eq. (17) and Eq. (18) we have the desired result.

D PROOF OF THEOREM 4.3

Proof. First, we introduce Welling's Woodbury identity (Welling, 2019):

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}.$$

Recall that we try to minimize Eq. (6)

$$\|\nabla \mathcal{L}(\boldsymbol{\theta}_{t_0}) - \mathbb{E}[\nabla T(X_{t_0}) \nabla \mathbf{h}_w(X_{t_0})]\|_{\mathcal{F}_{t_0}^{-1}}^2 + \lambda \|\mathbf{w}\|_{\mathcal{H}}^2. \quad (19)$$

Expanding the first square, up to a constant that does not depend on w , we obtain

$$\mathbb{E}[\nabla T(x_{t_0}) \mathbf{h}_w(X_{t_0})]^\top \mathcal{F}_{t_0}^{-1} \mathbb{E}[\nabla T(x_{t_0}) \mathbf{h}_w(X_{t_0})] - 2\nabla^\top \mathcal{L}(\boldsymbol{\theta}_{t_0}) \mathcal{F}_{t_0}^{-1} \mathbb{E}[\nabla T(x_{t_0}) \mathbf{h}_w(X_{t_0})] + \lambda \|w\|_{\mathcal{H}}^2. \quad (20)$$

By definition, $\mathbf{h}_w(X_{t_0}) = \langle w, \nabla k(X_{t_0}, \cdot) \rangle_{\mathcal{H}}$, we obtain a quadratic form with respect to w ,

$$\begin{aligned} & \langle w, \mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)]^\top \mathcal{F}_{t_0}^{-1} \mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)] w \\ & - \langle w, 2\mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)]^\top \mathcal{F}_{t_0}^{-1} \nabla \mathcal{L}(\boldsymbol{\theta}_{t_0}) \rangle + \lambda \|w\|_{\mathcal{H}}^2, \end{aligned} \quad (21)$$

where we used $a^\top B C d = (a^\top B C d)^\top = (B C d)^\top a = d^\top (B C)^\top a = \langle d, C^\top B^\top a \rangle$ and the fact that the inverse of Fisher Information Matrix \mathcal{F}_{t_0} is a symmetric matrix. Differentiating both sides by w and setting the derivative to zero, we obtain the following optimality condition of the least squares:

$$\begin{aligned} & 2\mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)]^\top \mathcal{F}_{t_0}^{-1} \mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)] w - \\ & 2\mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)]^\top \mathcal{F}_{t_0}^{-1} \nabla \mathcal{L}(\boldsymbol{\theta}_{t_0}) + 2\lambda w = 0. \end{aligned} \quad (22)$$

Thus, the closed-form solution of the optimal solution w^* is

$$\left(\underbrace{\mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)]^\top \mathcal{F}_{t_0}^{-1} \mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)]}_{R^{-1}} + \underbrace{\lambda \mathbf{I}}_{P^{-1}} \right)^{-1} \mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)]^\top \mathcal{F}_{t_0}^{-1} \nabla \mathcal{L}(\boldsymbol{\theta}_{t_0}). \quad (23)$$

Applying Woodbury's identity, we get:

$$\mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)]^\top (\lambda \mathcal{F}_{t_0} + \mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)] \mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)]^\top)^{-1} \nabla \mathcal{L}(\boldsymbol{\theta}_{t_0}). \quad (24)$$

Note that product

$$\mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)] \mathbb{E}[\nabla T(X_{t_0}) \nabla k(X_{t_0}, \cdot)]^\top = \mathbb{E}[\nabla T(X_{t_0}) \nabla \nabla k(X_{t_0}, X'_{t_0}) \nabla^\top T(X'_{t_0})],$$

we obtain the desired result. \square

E VISUALIZATION OF INGD ON GAUSSIAN MANIFOLD

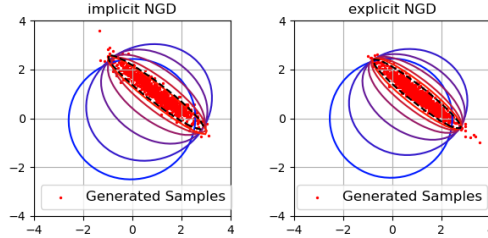


Figure 4: The evolution of parametric distributions under iNGD and classic NGD on the manifold of Gaussian distributions, i.e., $\mathcal{M}([x, x^\top x])$, starting from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Each ellipse represents the 95% confidence interval of the parametric distribution at the current step. As t increases, the ellipse turns from blue to red. Black dotted line marks the confidence interval the target distribution p . For iNGD, the ellipses are approximated by fitting a Gaussian model to samples generated by $g(Z, t; w)$.

In Fig. 4, we show an example of the iNGD and compare it with the actual NGD on a Gaussian manifold $T(x) = [x, x x^\top]$, $x \in \mathbb{R}^2$. In this example, the generative model is an MLP with one hidden layer consisting of 67 neurons. We can see that the samples generated from iNGD accurately retrace the steps of the classic NGD, ultimately producing a set of samples (red dots) that resemble the target distribution (black dotted line).

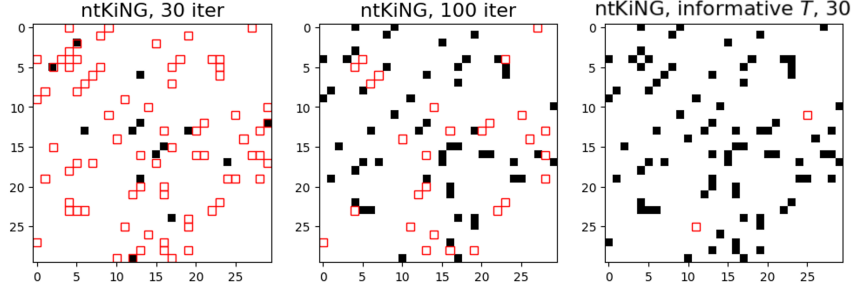


Figure 5: Sparsity pattern of Θ recovered by graphical lasso, using samples trained by ntKiNG with different sufficient statistics and number of iterations (red boxes indicate missing edges, the less boxes the better)

F REVERSE KL WASSERSTEIN GRADIENT FLOW AND MMD FLOW

The WGF dynamics that minimize $\text{KL}[q_t, p]$ move particles X_t using a simple velocity field:

$$\frac{dX_t}{dt} = \nabla (\log p)(X_t) - \nabla (\log q_t)(X_t),$$

where the gradient of log density could be easily estimated via kernel density estimation and the MMD flow minimizes $\text{MMD}[Y, X_t]$ using the following velocity field:

$$\frac{dX_t}{dt} = N \nabla \left(\frac{1}{N} \sum_{i=1}^n \|X_t - X_t^{(i)}\| - \frac{1}{M} \sum_{j=1}^M \|X_t - Y^{(j)}\| \right).$$

G GRAPHICAL MODEL RECOVERY WITH INFORMATIVE SUFFICIENT STATISTICS

In this experiment, we showcase how much improvement we can get when using informative sufficient statistics. Exponential family are commonly used to encode graphical models. For example, a Gaussian graphical model is a Gaussian density $p = \mathcal{N}(\mathbf{0}, \Theta^{-1})$, where the sparsity pattern of Θ encodes an undirected graph, describing the interactions of the random variables. One can imagine that if the generated samples approximate p well, we should recover the correct graphical model from these samples. In this experiment, we let p be a 30-dimensional Gaussian graphical model, and draw 200 samples $Y \sim p$, 200 samples $X_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and move the X_0 toward p using ntKiNG algorithm. Finally we apply the graphical lasso (Friedman et al., 2007) to estimate graphical models displayed in the left and middle plots of Fig. 5. Here $T(\mathbf{x}) := [\text{RBF basis}]$.

Since our methods use a probabilistic model to guide its training process, one may wonder if knowing the graphical model structure would improve the performance of the algorithm. To test this, we design a new sufficient statistic $T(\mathbf{x}) := [\text{RBF basis}, \forall (i, j) \in \{(i, j) | \Theta_{i,j} \neq 0\}, x_i x_j]$, i.e., we added pairwise potential functions that corresponds to pairwise factors in this graphical model. We ran the ntKiNG again with this new, better informed sufficient statistic for 30 runs. The graphical lasso estimate is shown in the right plot of Fig. 5. It can be seen that, when using the informed sufficient statistics, ntKiNG could recover the almost-correct graphical structure in only 30 iterations, while it takes the regular ntKiNG much longer. This suggests, our methods can indeed use a pre-existing probabilistic model to accelerate its generative model training process.

H COVARIATE SHIFT BY DISTRIBUTION MATCHING

In domain adaptation tasks, samples are drawn from the source distribution p_{XY} and the target distribution q_{XY} where X, Y are covariates and label respectively. The problem is that a classifier trained on source distribution samples may not work on target distribution samples. Covariate shift (Sugiyama et al., 2008; Quiñero-Candela et al., 2009) refers to a special case where $p_X \neq q_X$

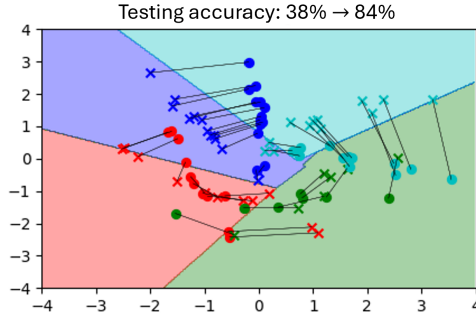


Figure 6: The inverse mapping ψ^{-1} found by KiNG. The classification boundary of a source classifier trained on p is depicted in four distinct colors. Target domain samples are marked with \bullet , and KiNG transforms these samples to new positions indicated by \times . Notably, many samples—especially those in the cyan and blue classes—are transported from the incorrect side of the classification boundary to the correct side and the test accuracy increases as a result.

Src. \rightarrow Tar.	base (%)	ntKiNG	WGF	MMD
amz \rightarrow dslr	69.50	+13.75	-0.50	+2.75
amz \rightarrow web	72.75	+8.25	-2.00	+0.75
amz \rightarrow cal	91.50	-0.75	-5.75	+0.00
dslr \rightarrow amz	86.00	+1.27	-6.37	+1.27
dslr \rightarrow web	98.09	+0.00	-0.64	+0.00
dslr \rightarrow cal	84.08	+5.73	-7.64	+0.64
web \rightarrow amz	77.97	+3.39	-2.71	+1.02
web \rightarrow dslr	91.19	+3.39	-4.07	+1.36
web \rightarrow cal	76.61	+3.39	-3.73	+0.34
cal \rightarrow amz	82.00	-2.50	-4.50	-0.25
cal \rightarrow dslr	58.25	+16.50	+7.00	+3.25
cal \rightarrow web	65.50	+10.25	+1.00	+2.00
Average	79.45	+5.22	-2.49	+1.09

Table 1: Comparison of Testing Accuracy Differences (in %) Relative to the Base Classifier

but $p_{Y|X} = q_{Y|X}$. We adopt a “marginal transport” assumption (Courty et al., 2016) that $X \sim q_X$ are generated as $X = \psi(X')$ where $X' \sim p_X$. It means, samples are generated from the source distribution and then “transported” to the target domain. For example, images in the source domain contains photos of objects, while in the target domain, photos contain the same objects but are filtered to reflect certain styles.

In the covariate shift setting, we observe joint samples from the source $p_{X,Y}$, but only have target covariates q_X . The goal is to find ψ^{-1} . In this paper, we propose to reverse the process by minimizing $\text{KL}[p_X, q_t]$ using Algorithm 2, where $X_0 \sim q_0$ are set to be the target domain covariates.

We demonstrate the effectiveness of this algorithm in Fig. 6, where the transfer ψ is a clockwise rotation on samples by 45 degrees. An inverse ψ^{-1} is a counter-clockwise rotation and has been correctly identified by ntKiNG.

We further test our algorithm on the Office+Caltech dataset (Gong et al., 2012) which is an object recognition dataset with photos collected from four different places: amazon, dslr, webcam, caltech. The task is to train a source classifier using one of the places, and test it on samples from another place. We test the performance of ntKiNG against two other particle-based transport methods WGF and MMD that also matches q_t with p_X . The performance is measured by the percentage gains compared with directly applying the source classifier to the target samples. The results show that our method achieves the most accuracy gains comparing to WGF and MMD. In 10 out of 12 domain adaptations settings, ntKiNG improves the testing accuracy.

I EXPERIMENT SETUP

We summarize the experiments’ setup details in this section. For each experiment, we provide details of the dataset and pre-processing procedure, as well as the details of tuning parameters.

I.1 COMPARISON WITH REVERSE KL WASSERSTEIN GRADIENT FLOW AND MMD FLOW

I.1.1 DATASET AND PRE-PROCESSING

In this experiments, we let $p = 0.5\mathcal{N}(-2, \mathbf{I}) + 0.5\mathcal{N}(2, \mathbf{I})$. We draw 100 samples from p as the target samples Y , 100 samples from $\mathcal{N}(0, \mathbf{I})$ as the initial samples X_0 . No further processing is required.

I.1.2 PARAMETER TUNING

The main tuning parameter are kernel bandwidth and step sizes.

For all methods that uses RBF kernel/basis, we set the bandwidth to be the median pairwise distance of all samples.

For all methods, we use step size 1, as any larger learning rate would result in numerical instability for each method.

For all methods, we run 100 particle updates.

The performance metric MMD uses a Gaussian kernel and the bandwidth is set as the median of pairwise distances of all samples Y and X_t .

I.2 GRAPHICAL MODEL RECOVERY

I.2.1 DATASET AND PRE-PROCESSING

In this experiment, we let p be a 30-dimensional Gaussian graphical model, and draw 200 samples $Y \sim p$, 200 samples $X_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The graphical model Θ is generated as a random graph, with edge probability 0.05. For each non-zero off-diagonal entry, we set $\Theta_{i,j} = 0.3$. No further processing is required.

I.2.2 PARAMETER TUNING

For all methods, we use the median of sample pairwise distances as the bandwidth.

For all methods, we set the step size to be 1.

Parameter tuning of Graphical Lasso is handled by `sklearn` internally using 5-fold cross validation and the sparse graph in graphical model is obtained by truncating all values smaller than 0.1. Below are the Python code.

```
# Fit with cross-validation to select alpha
model_cv = GraphicalLassoCV(alphas=10, # number of alphas or list of alphas
                             cv=5, # how many folds in cross-validation
                             max_iter=100,
                             tol=1e-4)
model_cv = model_cv.fit(x1_test.cpu().numpy())
Theta_cv = model_cv.precision_
Theta_cv = Theta_cv > 1e-1
```

I.3 EXPERIMENT 2: COVARIATE SHIFT

I.3.1 DATASET AND PRE-PROCESSING

We validate our method on the dataset Office+Caltech², which is a dataset for domain adaptation, consisting of Office 10 and Caltech 10 datasets. It contains the 10 overlapping categories between the Office dataset and Caltech256 dataset (Gong et al., 2012).

The original features are extracted using a DECAF network, and are 4096 dimensional. We apply PCA on the source and target domain to reduce the dimension to 50 with Python code

```
from sklearn.decomposition import PCA
pca = PCA(n_components=50)
pca.fit(X)
X = pca.transform(X)
X = X / 100
```

Due to memory space limit, we also randomly pick 200 samples from all target domains as X_0 .

I.3.2 PARAMETER TUNING

For all methods, we set step size to 0.1.

For ntKiNG, we run 100 steps due to reduce the computation cost.

For WGF and MMD flow, we run 1000 steps.

The source classifier is an RBF kernel Support Vector Machines with all hyper-parameters chosen by cross-validation with the following python code:

```
# Split the data into training and test sets (optional)
X_train, X_test, y_train, y_test = train_test_split(x, y,
                                                    test_size=0.3, random_state=42)

# Define parameter grid
param_grid = {
    'C': np.logspace(-3, 3, 5),
    'gamma': np.linspace(.2, 5, 5) * gamma,
    'kernel': ['rbf']
}

# Create a SVC classifier
svc = SVC()

# Initialize GridSearchCV
grid_search = GridSearchCV(svc, param_grid, refit=True, verbose=2, cv=5)

# Fit the model
grid_search.fit(X_train, y_train)
```

where gamma is the inverse of the median pairwise distances of all inputs.

J RELATED WORKS

Our methods bridge the gap between generative model training/sampling and parametric model optimization. Both domains are extensively studied in the machine learning community.

There has been a trend toward using optimization techniques to sample from unknown distributions. These methods first draw samples from an initial distribution and then move them according to a

²<https://github.com/jindongwang/transferlearning/blob/master/data/dataset.md#office+caltech>

time-dependent velocity field (Liu, 2017; Chewi et al., 2020; Maurais & Marzouk, 2024). A typical family of such gradient flows is the Wasserstein Gradient Flow (Ambrosio et al., 2008), which has found various applications outside of sampling, such as generative modelling (Gao et al., 2019; Choi et al., 2024) and missing data imputation (Chen et al., 2024). Our method falls within this family of algorithms, and we compared two of its variants in our experiments. To the best of our knowledge, none of the existing approaches could leverage a pre-existing probabilistic model to guide the flow of particles. Our framework is also more general: Algorithm 1 works for non-particle based generative models as well.

Another trend in generative modelling is “flow matching”, where one aligns the drift function with a pre-constructed flow (Lipman et al., 2023; Liu et al., 2023). In a similar spirit, our method also aligns the instantaneous change of the generative distribution with a prescribed dynamics (NGD). However, instead of directly matching the velocity field in the sample space, we match the projections of these changes in the parametric space. This approach avoids building arbitrary “bridges” between the reference and target distributions in sample space and instead leverages an effective parametric optimization algorithm to guide the training of the generative model.

In recent years, there has also been efforts to accelerate and approximate NGD using kernel methods, for example, (Arbel et al., 2020; Li et al., 2019b) propose to approximate the natural gradient by optimizing a dual formulation. However, both methods consider optimizing a probabilistic model, rather than a generative model as described in this paper. Performing NGD requires inverting a large matrix. Many research on NGD focuses on approximating the inverse the Fisher Information Matrix (Martens & Grosse, 2015; Grosse & Martens, 2016; George et al., 2018). Our particle update, e.g., Theorem 4.3 also requires us inverting a matrix with the dimension of the sufficient statistic. It would be an interesting future work to see if these techniques could be adapted to our approach.