

Minimizing f -Divergences by Interpolating Velocity Fields

Song Liu¹ Jiahao Yu¹ Jack Simons¹ Mingxuan Yi¹ Mark Beaumont¹

Abstract

Many machine learning problems can be seen as approximating a *target* distribution using a *particle* distribution by minimizing their statistical discrepancy. Wasserstein Gradient Flow can move particles along a path that minimizes the f -divergence between the target and particle distributions. To move particles, we need to calculate the corresponding velocity fields derived from a density ratio function between these two distributions. Previous works estimated such density ratio functions and then differentiated the estimated ratios. These approaches may suffer from overfitting, leading to a less accurate estimate of the velocity fields. Inspired by non-parametric curve fitting, we directly estimate these velocity fields using interpolation techniques. We prove that our estimators are consistent under mild conditions. We validate their effectiveness using novel applications on domain adaptation and missing data imputation. The code for reproducing our results can be found at <https://github.com/aneewgithubname/gradest2>. This manuscript is an extended version of the ICML2024 version.

1. Introduction

Many machine learning problems can be formulated as minimizing statistical divergences between distributions, e.g., Variational Inference (Blei et al., 2017), Generative Modeling (Nowozin et al., 2016; Yi et al., 2023), Domain Adaptation (Courty et al., 2017a; Yu et al., 2021), and Data Imputation (Muzellec et al., 2020). Among many divergence minimization techniques, Particle-based gradient descent reduces the divergence between a set of particles (which define a distribution) and the target distribution by iteratively moving particles according to an update rule.

¹University of Bristol, Bristol, UK. Correspondence to: Song Liu <song.liu@bristol.ac.uk>.

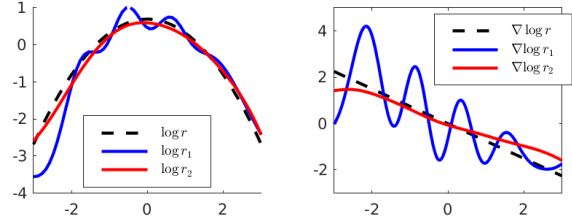


Figure 1. Estimating a log density ratio $\log r$ using a flexible model (RBF kernel) leads to an overfitted estimate ($\log r_1$). The overfitting consequently causes huge fluctuations in the derivative $(\log r_1)'$. Our proposed method provides a much more stable estimate $\log r_2$ and a more accurate estimate of $(\log r_2)'$.

One such algorithm is Stein Variational Gradient Descent (SVGD) (Liu and Wang, 2016; Liu, 2017). It minimizes the Kullback-Leibler (KL) divergence using the steepest descent algorithm and has achieved promising results in Bayesian inference. However, to compute the SVGD updates, we need unnormalized target density functions. If we only have samples from the target distribution—as is often the case in applications like domain adaptation and generative model training—we cannot directly apply SVGD to minimize the f -divergences.

Wasserstein Gradient Flow (WGF) describes the evolution of a marginal measure q_t along the steepest descent direction of a functional objective in Wasserstein geometry where $x_t \sim q_t$ follows a probability flow ODE. It can be employed to minimize an f -divergence between a particle distribution and the target distribution. Particularly, such WGFs characterize the following ODE (Yi et al., 2023; Gao et al., 2019)

$$dx_t = \nabla(h \circ r_t)(x_t)dt, \quad t \in [0, \infty).$$

Here $r_t := \frac{p}{q_t}$ is the ratio between the target density p and particle density q_t at time t , and h is a known function which depends on the f -divergence. The gradient operator ∇ computes the gradient of the composite function $h \circ r_t$. If we know r_t , simulating the above ODE would be straightforward using a discrete-time Euler method. However, the main challenge is that we do not know the ratio r_t in practice. In the context we consider, we know neither the particle nor target density, which constitutes r_t .

To overcome this issue, recent works (Gao et al., 2019;

Ansari et al., 2021; Simons et al., 2021) first obtain an estimate \hat{r}_t using a density ratio estimator (Sugiyama et al., 2012), then differentiate $h \circ \hat{r}_t$ to obtain $\nabla(h \circ \hat{r}_t)(\mathbf{x}_t)$. However, like other estimation tasks, density ratio estimation can be prone to overfitting. The risk of overfitting is further exacerbated when employing flexible models such as kernel models or neural networks. Overfitting can be disastrous for gradient estimation: A wiggly fit of r_t will cause huge fluctuations in gradient (see the blue fit in Figure 1). Moreover, density ratio estimation lacks the inductive bias for the density ratio gradient estimation. In other words, while it might be good at estimating the ratio itself, it doesn't have any built-in assumptions to capture the gradient of the density ratio function accurately.

In this paper, we **directly approximate the velocity fields induced by WGF**, i.e., $\nabla(h \circ r_t)(\mathbf{x}_t)$. We show that the backward KL velocity field, where $h = \log$, can be effectively estimated using Nadaraya-Watson (NW) interpolation if we know $\nabla \log p$, and this estimator is closely related to SVGD. We prove the estimation error of $\nabla \log r_t(\mathbf{x}_t)$ vanishes as the kernel bandwidth approaches to zero. This finding motivates us to propose a more general linear interpolation method to approximate $\nabla(h \circ r_t)(\mathbf{x}_t)$ for any general h functions using only samples from p and q_t . Our estimators are based on the idea that, within the neighbourhood of a given point, the best linear approximation of $h \circ r_t$ has slope $\nabla(h \circ r_t)$. Under mild conditions, we show that our estimators are also consistent for estimating $\nabla(h \circ r_t)(\mathbf{x}_t)$ and achieve the optimal non-parametric regression rate. Finally, equipped with our proposed gradient estimator, we test WGF on two novel applications: *domain adaptation* and *missing data imputation* and achieve promising performance.

2. Background

Notation: \mathbb{R}^d is the d -dimensional real domain. Vectors are lowercase bold letters, e.g., $\mathbf{x} := [x_1, \dots, x_d]^\top$. \mathbf{x}_{-j} is a subvector of \mathbf{x} obtained by excluding the j -th dimension. Matrices are uppercase bold letters, e.g., \mathbf{X}, \mathbf{Y} . $f \circ g$ is the composite function $f(g(\cdot))$. f' is the derivative of a univariate function. $\partial_i f$ means the partial derivative with respect to the i -th input of f and $\nabla f := [\partial_1 f, \partial_2 f, \dots, \partial_d f]^\top$. $\nabla^\top f$ represents its transpose. $\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$ represents the gradient of f with respect to the input \mathbf{x} . $\nabla \mathbf{f} \in \mathbb{R}^{d \times m}$ represents the Jacobian of a vector-valued function $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^d$. $\partial_i^2 f$ means the second-order partial derivative with respect to the i -th input of f . $\nabla^2 f$ is the Hessian of f . $\|\cdot\|$ is the ℓ_2 norm of a vector or the spectral norm of a matrix. $\lambda_{\min}[\mathbf{X}]$ represents the smallest eigenvalue of a matrix \mathbf{X} . $a \vee b$ is the greater value between two scalars a and b . $\mathcal{P}(\mathbb{R}^d)$ is the space of probability measures defined on \mathbb{R}^d equipped with Wasserstein-2 metric or Wasserstein space for short. $\hat{\mathbb{E}}[\cdot]$ is

Name	Notation	$f(r_t)$	$h(r_t)$
Forw. KL	$\text{KL}[p, q_t]$	$r_t \log(r_t)$	r_t
Back. KL	$\text{KL}[q_t, p]$	$-\log(r_t)$	$\log r_t - 1$
Pearson's χ^2	$\chi_p^2[p, q_t]$	$\frac{1}{2}(r_t - 1)^2$	$\frac{1}{2}r_t^2 - \frac{1}{2}$
Neyman's χ^2	$\chi_n^2[p, q_t]$	$\frac{1}{2r_t} - \frac{1}{2}$	$-\frac{1}{r_t} + \frac{1}{2}$

Table 1. Some f -divergences, their definitions and h functions computed according to Theorem 2.1.

the sample approximation of an expectation $\mathbb{E}[\cdot]$.

We begin by introducing WGF of f -divergence, an effective technique for minimizing f -divergence between the particle and target distribution.

2.1. Wasserstein Gradient Flows of f -divergence

In general terms, a Wasserstein Gradient Flow is a curve in probability space (Ambrosio et al., 2005). By moving a probability measure along this curve, a functional objective (such as a statistical divergence) is reduced. In this work, we focus solely on using f -divergences as the functional objective. Let $q_t : \mathbb{R}_+ \rightarrow \mathcal{P}(\mathbb{R}^d)$ be a curve in Wasserstein space. Consider an f -divergence defined by $D_f[p, q_t] := \int q_t(\mathbf{x}) f(r_t(\mathbf{x})) d\mathbf{x}$, where $r_t := \frac{p}{q_t}$. f is a twice differentiable convex function with $f(1) = 0$.

Theorem 2.1 (Corollary 3.3 in (Yi et al., 2023)). *The Wasserstein gradient flow of $D_f[p, q_t]$ characterizes the particle evolution via the ODE:*

$$d\mathbf{x}_t = \nabla(h \circ r_t)(\mathbf{x}_t) dt, \quad h(r_t) = r_t f'(r_t) - f(r_t).$$

Simply speaking, particles evolve in Euclidean space according to the above ODE moves the corresponding q_t along a curve where $D_f[p, q_t]$ always decreases with time.

Theorem 2.1 establishes a relationship between f and a function $h : \mathbb{R}_+ \mapsto \mathbb{R}$. The gradient field of $h \circ r_t$ over time as $t \rightarrow \infty$ is referred to as the *WGF velocity field*. Note that, in Gao et al. (2019) and Ansari et al. (2021), authors provided a similar theorem for the ‘‘reversed’’ f -divergence $\int p(\mathbf{x}) f(\frac{q_t(\mathbf{x})}{p(\mathbf{x})}) d\mathbf{x}$, which is different from the definition of f -divergence that we use in our paper. Some frequently used f -divergences and their corresponding h functions are listed in Table 1. Specifically, for the backward KL divergence we have $\nabla(h \circ r_t)(\mathbf{x}_t)|_{h=\log(\cdot)} = \nabla \log r_t(\mathbf{x}_t)$.

In reality, we move particles by simulating the above ODE using the forward Euler method: We draw particles from an initial distribution q_0 and iteratively update them for time $t = 0, 1 \dots T$ according to the following rule:

$$\mathbf{x}_{t+1} := \mathbf{x}_t + \eta \nabla(h \circ r_t)(\mathbf{x}_t) \quad (1)$$

where η is a small step size. There is a slight abuse of notation, and we reuse t for discrete-time indices ¹.

¹From now on, we will only discuss discrete-time algorithms.

Although (1) seems straightforward, we normally do not have access to r_t , so the update in (1) cannot be readily performed. In previous works, such as (Gao et al., 2019; Simons et al., 2021), r_t is estimated using density ratio estimators and the WGF is simulated using the estimated ratio. Although these estimators achieved promising results, the density ratio estimators are not designed for usage in WGF algorithms. For example, a small density ratio estimation error could lead to huge deviations in gradient estimation, as we demonstrated in Figure 1. Others (Wang et al., 2022) propose to estimate the gradient flow using Kernel Density Estimation (KDE) on densities p and q_t separately (See Section J), then compute the log ratio and its gradient. However, KDE tends to perform poorly in high dimensional settings (see e.g., (Scott, 1991)).

3. Direct Velocity Field Estimation by Interpolation

In this work, we consider directly estimating the velocity field, i.e., directly modelling and estimating $\nabla(h \circ r_t)$. We are encouraged by the recent successes in Score Matching (Hyvärinen, 2005; 2007; Vincent, 2011; Song et al., 2020), which is a direct estimator of a log density gradient. It works by minimizing the squared differences between the true log density gradient and the model gradient. However, such a technique cannot be easily adapted to estimate $\nabla(h \circ r_t)$, even for $h(r) = \log r$ (See Appendix K).

We start by looking at a simpler setting where $\nabla \log p$ is known. In fact, this setting itself has many interesting applications such as Bayesian inference. The solution we derive using interpolation will serve as a motivation for other interpolation based approaches in later sections.

3.1. Nadaraya-Watson (NW) Interpolation of Backward KL Velocity Field

Define a local weighting function with a parameter $\sigma > 0$, $k_\sigma(\mathbf{x}, \mathbf{x}^*) := \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}^*\|^2}{2\sigma^2}\right)$. Nadaraya-Watson (NW) estimator (Nadaraya, 1964; Watson, 1964) interpolates a function g at a **fixed point** \mathbf{x}^* . Suppose that we observe $g(\mathbf{x})$ at a set of sample points $\{\mathbf{x}_i\}_{i=1}^n \sim q$, NW interpolates $g(\mathbf{x}^*)$ by computing

$$\hat{g}(\mathbf{x}^*) := \widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)g(\mathbf{x})] / \widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)]. \quad (2)$$

Thus, the NW interpolation of the backward KL field² is

$$\hat{\mathbf{u}}_t(\mathbf{x}^*) := \widehat{\mathbb{E}}_{q_t}[k_\sigma(\mathbf{x}, \mathbf{x}^*)\nabla \log r_t(\mathbf{x})] / \widehat{\mathbb{E}}_{q_t}[k_\sigma(\mathbf{x}, \mathbf{x}^*)]. \quad (3)$$

Since we cannot evaluate $\nabla \log r_t(\mathbf{x})$, (3) is intractable. However, assuming that $\lim_{\|\mathbf{x}\| \rightarrow \infty} q_t(\mathbf{x})k_\sigma(\mathbf{x}, \mathbf{x}^*) = 0$, us-

ing integration by parts³, the expectation in the numerator of (3) can be rewritten as:

$$\begin{aligned} \mathbb{E}_{q_t}[k_\sigma^* \nabla \log r_t(\mathbf{x})] &= \mathbb{E}_{q_t}[k_\sigma^* \nabla \log p(\mathbf{x})] - \mathbb{E}_{q_t}[\nabla \log q_t(\mathbf{x})] \\ &= \mathbb{E}_{q_t}[k_\sigma^* \nabla \log p(\mathbf{x}) + \nabla k_\sigma^*], \end{aligned} \quad (4)$$

where we shortened the kernel $k_\sigma(\mathbf{x}, \mathbf{x}^*)$ as k_σ^* . Since we can evaluate $\nabla \log p$, (4) can be approximated using samples from the particle distribution q_t . Thus, the NW estimator of the backward KL field can be approximated by

$$\hat{\mathbf{u}}_t(\mathbf{x}^*) \approx \widehat{\mathbb{E}}_{q_t}[k_\sigma^* \nabla \log p(\mathbf{x}) + \nabla k_\sigma^*] / \widehat{\mathbb{E}}_{q_t}[k_\sigma^*]. \quad (5)$$

Interestingly, the numerator of (5) is exactly the particle update of the SVGD algorithm (Liu and Wang, 2016) for an RKHS induced by a Gaussian kernel (See Appendix N for details on SVGD), and the equality (4) has been noticed by Chewi et al. (2020). Note that for different \mathbf{x}^* , the denominator in (5) is different and thus cannot be combined into the overall learning rate of SVGD.

3.2. Effectiveness of NW Estimator

For simplicity, we drop t from \mathbf{u}_t , r_t , \mathbf{x}_t and q_t when our analysis holds the same for all t .

Although there have been theoretical justifications for the convergence analysis of WGF given the ground truth velocity fields such as Langevin dynamics (Wibisono, 2018). Few theories have been dedicated to the estimation of velocity fields themselves. One of the contributions of this paper is that we study the statistical theory of the velocity field estimation through the lenses of non-parametric regression/curve approximation.

Now, we prove the convergence rate of the NW estimator under the assumption that the second-order derivative of $\log r$ is well-behaved. Although $\hat{\mathbf{u}}(\mathbf{x}^*)$ cannot be directly computed, assuming $\nabla \log p$, k and ∇k are well-behaved, using concentration inequalities (such as Hoeffding's inequality (Hoeffding, 1963)), the difference between $\hat{\mathbf{u}}(\mathbf{x}^*)$ and its approximation (5) can be easily bounded. Thus, we focus on the classical NW estimator in (3).

Proposition 3.1. *Suppose $\sup_{\mathbf{x} \in \mathbb{R}^d} \|\nabla^2 \log r(\mathbf{x})\| \leq \kappa < \infty$. Define $k(\mathbf{y}) := \exp[-\|\mathbf{y}\|^2/2]$. Assume that there*

³

$$\begin{aligned} &\int \int q(\mathbf{x})k(\mathbf{x}, \mathbf{x}^*)\partial_i \log q(\mathbf{x})d\mathbf{x}_i d\mathbf{x}_{-i} \\ &= \int \left[[q(\mathbf{x})k(\mathbf{x}, \mathbf{x}^*)]_{\mathbf{x}_i \rightarrow -\infty}^{\mathbf{x}_i \rightarrow \infty} - \int q(\mathbf{x})\partial_i k(\mathbf{x}, \mathbf{x}^*)d\mathbf{x}_i \right] d\mathbf{x}_{-i} \end{aligned}$$

²“backward KL field” is short for backward KL velocity field.

exist constants C_k, K that are independent of σ , such that

$$\frac{\int q(\sigma \mathbf{y} + \mathbf{x}^*)k(\mathbf{y})\|\mathbf{y}\|d\mathbf{y}}{\int q(\sigma \mathbf{y} + \mathbf{x}^*)k(\mathbf{y})d\mathbf{y}} \leq C_k, \mathbb{E}_q \left[\frac{1}{\sigma^d} k_\sigma^* \right] \geq K > 0 \quad (6)$$

$$\begin{aligned} \text{Var}_q \left[\frac{1}{\sigma^d} k_\sigma^* \|\mathbf{x} - \mathbf{x}^*\| \right] &= O \left(\frac{1}{\sigma^d} \right), \\ \text{Var}_q \left[\frac{1}{\sigma^d} k_\sigma^* \right] &= O \left(\frac{1}{\sigma^d} \right), \text{ as } \sigma \rightarrow 0. \end{aligned} \quad (7)$$

Then, with high probability, there exists a constant K' :

$$\|\hat{\mathbf{u}}(\mathbf{x}^*) - \nabla \log r(\mathbf{x}^*)\| \leq \sqrt{d}\kappa \left[\frac{K'}{\sqrt{n\sigma^d}} + \sigma C_k \right]$$

holds for all $\sigma > 0$.

We can also deduce that, with an optimal choice of $\sigma \sim n^{-1/(d+2)}$, the estimation error is $O_p(n^{-1/(d+2)})$.

See Appendix A for the proof. $q(\sigma \mathbf{y} + \mathbf{x}^*)$ in the first inequality (6) can be further expanded using Taylor expansion. Provided that the kernel k is well-behaved, this condition becomes a regularity condition on $q(\mathbf{x}^*)$ and its higher order moments. The second inequality in (6) means that there should be enough mass around \mathbf{x}^* under the distribution q , which is a key assumption in classical nonparametric curve estimation (See, e.g., Chapter 20 in (Wasserman, 2010)). (7) is required to ensure the empirical quantities in the NW estimator converge to their population counterparts in probability.

It can be seen that the estimation error is bounded by the sample approximation error $\frac{K'}{\sqrt{n\sigma^d}}$ and a bias depending on by σ . Interestingly, the bias term decreases at the rate of σ , slower than the classical rate σ^2 for the non-parametric regression of a second-order differentiable function (Theorem 20.21 in (Wasserman, 2010)). This is expected as NW estimates the *gradient* of $h \circ r$, not $h \circ r$. To achieve a faster σ^2 rate, one needs to assume conditions on $\nabla^3 \log r$. This also highlights a slight downside of using NW to estimate the gradient. However, in the following section, we show that another interpolator achieves the superior rate when using the same type of assumption on $\nabla^2(h \circ r)$.

4. Velocity Field Interpolation from Samples

Although we have seen that $\hat{\mathbf{u}}$ is an effective estimator of the backward KL velocity field, it can only be approximated when we can evaluate $\nabla \log p$. In some applications, such as domain adaptation or generative modelling, we only have samples from the target distribution, and $\nabla \log p$ is unavailable. Moreover, since the f -divergence family consists of a wide variety of divergences, we hope to provide a *general computational framework* to estimate different velocity fields that minimize different f -divergences.

Nonetheless, the success of NW motivates us to look for other interpolators to approximate $\nabla(h \circ r)(\mathbf{x}^*)$.

Another common interpolation technique is *local linear regression* (See e.g., (Gasser and Müller, 1979; Fan, 1993) or Chapter 6, (Hastie et al., 2001)). It approximates an unknown function g at \mathbf{x}^* by using a linear function: $\hat{g}(\mathbf{x}) := \langle \boldsymbol{\beta}(\mathbf{x}^*), \mathbf{x} \rangle + \beta_0(\mathbf{x}^*)$. $\boldsymbol{\beta}(\mathbf{x}^*)$ and $\beta_0(\mathbf{x}^*)$ are the minimizer of the following weighted least squares objective:

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \widehat{\mathbb{E}}_q \left[k_\sigma^* (g(\mathbf{x}) - \langle \boldsymbol{\beta}, \mathbf{x} \rangle - \beta_0)^2 \right]. \quad (8)$$

A key insight is, since *the gradient of a function is the slope of its best local linear approximation*, it is reasonable to just use the slope of the fitted linear model, a.k.a., $\boldsymbol{\beta}(\mathbf{x}^*)$, to approximate the gradient $\nabla g(\mathbf{x}^*)$. See Figure 6 in Appendix for an illustration.

We apply the same rationale to estimate $\nabla(h \circ r)(\mathbf{x}^*)$. However, unlike local linear interpolation, we cannot evaluate $h \circ r$ at any input. Thus, we cannot directly use the least squares objective (8) to obtain a local linear interpolation. Similar to what we have done in Section 3.1, we look for a tractable population estimator for estimating $h \circ r$, which can be approximated using samples from p and q . Then, we “convert it” into a local linear objective.

In the following section, we derive an objective for estimating $h \circ r$ by maximizing a variational lower bound of a *mirror divergence*.

4.1. Mirror Divergence

Definition 4.1. Let $D_\phi[p, q]$ and $D_\psi[p, q]$ denote two f -divergences with f being ϕ and ψ respectively. D_ψ is the mirror of D_ϕ if and only if $\psi'(r) \triangleq r\phi'(r) - \phi(r)$, where \triangleq means equal up to a constant.

For example, let $D_\phi[p, q]$ and $D_\psi[p, q]$ be $\text{KL}[p, q]$ and $\chi_p^2[p, q]$ respectively. From Table 1, we can see that $\phi(r) = r \log r$ and $\psi(r) = \frac{1}{2}(r-1)^2$. Thus $r\phi'(r) - \phi(r) = r(1 + \log r) - r \log r = r \triangleq \psi'(r) = r - 1$. Therefore, $\chi_p^2[p, q]$ is the mirror of $\text{KL}[p, q]$. Similarly, we can verify that $\text{KL}[p, q]$ is the mirror of $\text{KL}[q, p]$ and $\text{KL}[q, p]$ is the mirror of $\chi_n^2[p, q]$. In general, $D_\psi[p, q]$ is the mirror of $D_\phi[p, q]$ does not imply the other direction. The mirror divergence is also not unique.

Here, we list a few more examples of f -divergences and their mirror divergences:

- Jensen-Shannon Divergence:

$$\phi = \frac{1}{2}r \log(r) - (r+1) \log \left(\frac{r+1}{2} \right).$$

– Mirror Divergence:

$$\psi = \frac{\log(r+1)}{2} - \frac{\log(2)}{2} + r \left(\frac{\log\left(\frac{r}{2} + \frac{1}{2}\right)}{2} - \frac{1}{2} \right) + \frac{1}{2}.$$

- Squared Hellinger Distance: $\phi = \frac{(\sqrt{r}-1)^2}{2}$.
 - Mirror Divergence: $\psi = \frac{r^{3/2}}{3} - \frac{r}{2}$.
- Total Variational Distance: $\phi = \frac{1}{2}|r-1|$.
 - Mirror Divergence: $\psi = \frac{1}{2}|r-1|$.

We also provide a MATLAB script in the GitHub repo to compute the mirror divergence of an f -divergence.

4.2. Gradient Estimator using Local Linear Interpolation

The key observation that helps derive a tractable objective is that $h \circ r$ is the argmax of “the mirror variational lowerbound”. Suppose h is associated with an f -divergence D_ϕ as per Theorem 2.1 and D_ψ is the mirror of D_ϕ . Then $h \circ r$ is the argmax of the following objective:

$$D_\psi[p, q] = \max_d \mathbb{E}_p[d(\mathbf{x})] - \mathbb{E}_q[\psi_{\text{con}}(d(\mathbf{x}))], \quad (9)$$

where ψ_{con} is the *convex conjugate* of ψ . The formal statement and its proof can be found in Appendix C. The equality in (9) is known in previous literature (Nguyen et al., 2010; Nowozin et al., 2016) and the objective in (9) is commonly referred to as the variational lowerbound of $D_\psi[p, q]$.

Notice that the expectations in (9) can be approximated by $\widehat{\mathbb{E}}_p[\cdot]$ and $\widehat{\mathbb{E}}_q[\cdot]$ using samples from p and q respectively.

This is a *surprising result*. Since h is related to D_ϕ ’s field (as per Theorem 2.1), one may associate maximizing D_ϕ ’s variational lowerbound with its velocity field estimation. However, the above observation shows that, to approximate D_ϕ ’s field, one should maximize the variational lowerbound of its *mirror divergence* D_ψ ! To our best knowledge, this “mirror structure” in the context of WGF has never been studied before.

We then localize (9) to obtain a local linear estimator of $h \circ r$ at a fixed point \mathbf{x}^* . First, we parameterize the function d using a linear model $d_{\mathbf{w}, b}(\mathbf{x}) := \langle \mathbf{w}, \mathbf{x} \rangle + b$. Second, we weight the objective using $k_\sigma(\mathbf{x}, \mathbf{x}^*)$, which leads to the following local linear objective:

$$(\mathbf{w}(\mathbf{x}^*), b(\mathbf{x}^*)) = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \ell(\mathbf{w}, b; \mathbf{x}^*),$$

$$\text{where } \ell(\mathbf{w}, b; \mathbf{x}^*) := \widehat{\mathbb{E}}_p[k_\sigma^* d_{\mathbf{w}, b}(\mathbf{x})] - \widehat{\mathbb{E}}_q[k_\sigma^* \psi_{\text{con}}(d_{\mathbf{w}, b}(\mathbf{x}))] \quad (10)$$

The above transformation is similar to how the local linear regression “localizes” the ordinary least squares objective.

Solving (10), we get a linear approximation of $h \circ r$ at \mathbf{x}^* :

$$h(r(\mathbf{x}^*)) \approx \langle \mathbf{w}(\mathbf{x}^*), \mathbf{x}^* \rangle + b(\mathbf{x}^*).$$

Following the intuition that $\nabla(h \circ r)(\mathbf{x}^*)$ is the slope of the best local linear fit of $h(r(\mathbf{x}^*))$, we use $\mathbf{w}(\mathbf{x}^*)$ to approximate $\nabla(h \circ r)(\mathbf{x}^*)$. We will theoretically justify this approximation in Section 4.3. Now let us study two examples:

Example 4.2. Suppose we would like to estimate $\text{KL}[p, q]$ ’s field at \mathbf{x}^* , which is $\nabla r(\mathbf{x}^*)$. Using Definition 4.1, we can verify that the mirror of $\text{KL}[p, q]$ is $D_\psi = \chi_p^2[p, q]$, in which case $\psi = \frac{1}{2}(r-1)^2$. The convex conjugate of ψ is $\psi_{\text{con}}(d) = d^2/2 + d$. Substituting ψ_{con} in (10), the gradient estimator $\mathbf{w}_\rightarrow(\mathbf{x}^*) \approx \nabla r(\mathbf{x}^*)$ is obtained by the following objective:

$$(\mathbf{w}_\rightarrow(\mathbf{x}^*), b_\rightarrow(\mathbf{x}^*)) := \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \ell_\rightarrow(\mathbf{w}, b; \mathbf{x}^*)$$

where $\ell_\rightarrow(\mathbf{w}, b; \mathbf{x}^*) := \widehat{\mathbb{E}}_p[k_\sigma^* \cdot d_{\mathbf{w}, b}(\mathbf{x})] -$

$$\widehat{\mathbb{E}}_q \left[k_\sigma^* \cdot \left(\frac{d_{\mathbf{w}, b}(\mathbf{x})^2}{2} + d_{\mathbf{w}, b}(\mathbf{x}) \right) \right]. \quad (11)$$

Example 4.3. Suppose we would like to estimate $\text{KL}[q, p]$ ’s field at \mathbf{x}^* , which is $\nabla \log r(\mathbf{x}^*)$. Using Definition 4.1, we can verify that the mirror of $\text{KL}[q, p]$ is $D_\psi = \text{KL}[p, q]$, in which case, $\psi(r) = r \log r$. The convex conjugate of ψ is $\psi_{\text{con}}(d) = \exp(d-1)$. The gradient estimator $\mathbf{w}_\leftarrow(\mathbf{x}^*) \approx \nabla \log r(\mathbf{x}^*)$ is obtained by the following objective:

$$(\mathbf{w}_\leftarrow(\mathbf{x}^*), b_\leftarrow(\mathbf{x}^*)) := \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \ell_\leftarrow(\mathbf{w}, b; \mathbf{x}^*)$$

where $\ell_\leftarrow(\mathbf{w}, b; \mathbf{x}^*) := \widehat{\mathbb{E}}_p[k_\sigma^* \cdot d_{\mathbf{w}, b}(\mathbf{x})] -$

$$\widehat{\mathbb{E}}_q[k_\sigma^* \cdot \exp(d_{\mathbf{w}, b}(\mathbf{x}) - 1)]. \quad (12)$$

In the following section, we show that the estimation error $\|\mathbf{w}(\mathbf{x}^*) - \nabla(h \circ r)(\mathbf{x}^*)\|$ vanishes as $\sigma \rightarrow 0$ and $n \rightarrow \infty$.

4.3. Effectiveness of Local Linear Interpolation

In this section, we state our main theoretical result. Let $(\mathbf{w}(\mathbf{x}^*), b(\mathbf{x}^*))$ be a stationary point of $\ell(\mathbf{w}, b; \mathbf{x}^*)$. We denote the domain of \mathbf{x} as \mathcal{X} (not necessarily \mathbb{R}^d). Without loss of generality, we also assume all sample averages $\widehat{\mathbb{E}}[\cdot]$ are averaged over n samples. We prove that, $\mathbf{w}(\mathbf{x}^*)$ is a consistent estimate of $\nabla(h \circ r)(\mathbf{x}^*)$ assuming the change rate of the flow is bounded.

Assumption 4.4. The change rate of the velocity fields is well-behaved, i.e.,

$$\sup_{\mathbf{x} \in \mathcal{X}} \|\nabla^2(h \circ r)(\mathbf{x})\| \leq \kappa.$$

This is an analogue of the assumption on $\nabla^2 \log r$ in Proposition 3.1.

Assumption 4.5. There exists a constant $C_k > 0$ independent of σ ,

$$\frac{1}{2} \int q(\sigma \mathbf{y} + \mathbf{x}^*) k(\mathbf{y}) \cdot \|\mathbf{y}\|^2 \cdot \|\llbracket \sigma \mathbf{y} + \mathbf{x}^*, 1 \rrbracket\| d\mathbf{y} \leq C_k.$$

This assumption is similar to the first inequality in (6). Expanding $q(\sigma \mathbf{y} + \mathbf{x}^*)$ using Taylor expansion and providing that our kernel is well behaved, this assumption essentially implies the boundedness of the $q(\mathbf{x}^*)$ and $\|\nabla^2 q(\mathbf{x}^*)\|$.

Define two shorthands: $\mathbf{w}^* := \nabla(h \circ r)(\mathbf{x}^*)$ and $b^* := h(r(\mathbf{x}^*)) - \langle \nabla(h \circ r)(\mathbf{x}^*), \mathbf{x}^* \rangle$.

Assumption 4.6. Let $\tilde{\mathbf{x}} := [\mathbf{x}^\top, 1]^\top$. As $\sigma \rightarrow 0$,

$$\begin{aligned} \text{tr} \left[\text{Cov}_p \left[\frac{1}{\sigma^d} k_\sigma^* \cdot \tilde{\mathbf{x}} \right] \right] &= O\left(\frac{1}{\sigma^d}\right), \\ \text{tr} \left[\text{Cov}_q \left[\frac{1}{\sigma^d} k_\sigma^* \cdot \psi'_{\text{con}}(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*) \tilde{\mathbf{x}} \right] \right] &= O\left(\frac{1}{\sigma^d}\right). \end{aligned}$$

These two are analogues of (7). They are required so that our sample approximation of the objective is valid and concentration inequalities can be applied.

Assumption 4.7. For all $a \in [0, 1]$ and $\mathbf{x} \in \mathcal{X}$, $\psi''_{\text{con}}[ah(r(\mathbf{x})) + (1-a)(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*)] \leq C_{\psi''_{\text{con}}}$.

This is a unique assumption to our estimator, where we assume that the convex conjugate ψ_{con} has a bounded second-order derivative.

Theorem 4.8. Suppose Assumption 4.4, 4.5, 4.6 and 4.7 holds. If there exist strictly positive constants W, B, Λ_{\min} that are independent of σ and n such that,

$$\|\mathbf{w}^*\| \leq W, \quad |b^*| \leq B \quad (13)$$

and for all $\mathbf{w} \in \{\mathbf{w} \mid \|\mathbf{w}\| < 2W\}$ and $b \in \{b \mid |b| < 2B\}$,

$$\lambda_{\min} \left\{ \widehat{\mathbb{E}}_q \left[k_\sigma^* \nabla_{[\mathbf{w}, b]}^2 \psi_{\text{con}}(\langle \mathbf{w}, \mathbf{x} \rangle + b) \right] \right\} \geq \sigma^d \Lambda_{\min}, \quad (14)$$

holds with high probability. Then there exists $\sigma_0, N, K > 0$ such that for all $0 < \sigma < \sigma_0, n > N$,

$$\|\mathbf{w}(\mathbf{x}^*) - \mathbf{w}^*\| \leq \frac{\frac{K}{\sqrt{n\sigma^d}} + \kappa C_k C_{\psi''_{\text{con}}} \sigma^2}{\Lambda_{\min}},$$

with high probability.

The proof can be found in Appendix D.

Similar to Proposition 3.1, the estimation error is upper-bounded by the sample approximation error that reduces with the ‘‘effective sample size’’ $n\sigma^d$, and a bias term that

reduces with σ . Interestingly, although we make the smoothness assumption on $\nabla^2(h \circ r)$, similar to Proposition 3.1, the bias vanishes at a *quadratic rate* σ^2 , unlike the linear rate obtained in Proposition 3.1.

We can deduce that, with an optimal choice of $\sigma \sim n^{-1/(d+4)}$, the estimation error is $O_p(n^{-2/(d+4)})$, a rate faster than the rate implied by Proposition 3.1.

Using Theorem 4.8, we can prove the consistency of various velocity field estimators for different f -divergences.

Corollary 4.9. Suppose $\sup_{\mathbf{x} \in \mathcal{X}} \|\nabla^2 r(\mathbf{x})\| \leq \kappa_{\rightarrow}$, Assumption 4.5, 4.6 holds and

$$\lambda_{\min} \left\{ \widehat{\mathbb{E}}_q \left[k_\sigma(\mathbf{x}, \mathbf{x}^*) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \right] \right\} \geq \sigma^d \cdot \Lambda_{\rightarrow} > 0, \quad (15)$$

holds with high probability. Then $\exists \sigma_0, N, K > 0$,

$$\|\mathbf{w}_{\rightarrow}(\mathbf{x}^*) - \nabla r(\mathbf{x}^*)\| \leq \frac{\frac{K}{\sqrt{n\sigma^d}} + \kappa_{\rightarrow} \cdot C_k \cdot \sigma^2}{\Lambda_{\rightarrow}}$$

for all $0 < \sigma < \sigma_0, n > N$ holds with high probability.

See Appendix E for the proof.

Corollary 4.10. Suppose $\sup_{\mathbf{x} \in \mathcal{X}} \|\nabla^2 \log r(\mathbf{x})\| \leq \kappa_{\leftarrow}$, and Assumption 4.5, 4.6 holds. Assume

$$\sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\| \leq C_{\mathcal{X}}, \quad \sup_{\mathbf{x} \in \mathcal{X}} \log r(\mathbf{x}) - 1 \leq C_{\log r},$$

and there exists $B, W < \infty$, so that

$$\|\nabla \log r(\mathbf{x}^*)\| < W, \quad |\log r(\mathbf{x}^*) - \langle \nabla \log r(\mathbf{x}^*), \mathbf{x}^* \rangle| < B.$$

Additionally,

$$\lambda_{\min} \left[\widehat{\mathbb{E}}_q [k_\sigma(\mathbf{x}, \mathbf{x}^*) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right] \geq \sigma^d \cdot \Lambda_{\leftarrow} > 0, \quad (16)$$

holds with high probability. Then $\exists \sigma_0, N, K > 0$ and for all $0 < \sigma < \sigma_0, n > N$,

$$\begin{aligned} &\|\mathbf{w}_{\leftarrow}(\mathbf{x}^*) - \nabla \log r(\mathbf{x}^*)\| \\ &\leq \frac{1}{(\Lambda_{\leftarrow} \cdot \exp(-2WC_{\mathcal{X}} - 2B - 1)) \cdot \sqrt{n\sigma^d}} + \\ &\quad \frac{\kappa_{\leftarrow} \cdot C_k \cdot [\exp(WC_{\mathcal{X}} + B - 1) \vee \exp(C_{\log r} - 1)] \cdot \sigma^2}{\Lambda_{\leftarrow} \cdot \exp(-2WC_{\mathcal{X}} - 2B - 1)}, \end{aligned}$$

holds with high probability.

See Appendix F for the proof. Note that due to the assumption that $\|\mathbf{x}\|$ is bounded, Corollary 4.10 can only be applied to density ratio functions with bounded input domains. However, this does include important examples such as images, where pixel brightnesses are bounded within $[0, 1]$. Using the same proof techniques, it is possible to derive more Corollaries for other f -divergence velocity fields. We leave them as a future investigation.

4.4. Model Selection via Local Linear Interpolation

Although Theorem 4.8 says the estimation bias disappears as $\sigma \rightarrow 0$, when we only have a finite number of samples, the choice of the kernel bandwidth σ controls the bias-variance trade-off of the local estimation. Thus we propose a model selection criterion. The details of the procedure is provided in Appendix I.1.

The high level idea is: Suppose we have testing samples from p and q . A good choice of σ would result in a good approximation of $h \circ r$ on *testing points*, thus the best approximation of $h \circ r$ would maximize the variational lower bound (9). Therefore, we only need to evaluate (9) on testing samples to determine the optimality of σ .

Our local linear estimator offers a unique advantage of model selection because it is formulated as a non-parametric curve fitting problem. In contrast, SVGD lacks a systematic approach and has to resort to the ‘‘median trick’’.

5. Experiments

5.1. Reducing KL Divergence: SVGD vs. NW vs. Local Linear Estimator

In this experiment, we investigate the performance of SVGD, NW and Local Linear (LL) estimator through the task of minimizing $\text{KL}[q_t, p]$. We let SVGD, NW and LL fit the target distribution $p = \mathcal{N}(0, 1)$. 500 iid initial particles are drawn from $q_0 = \mathcal{N}(-1, 0.25^2)$. For all methods, we use naive gradient descent to update particles with a fixed step size 0.1. We also consider a variant of SVGD where AdaGrad (Duchi et al., 2011) is applied to adjust the step size dynamically. For SVGD and SVGD with AdaGrad, we use the MATLAB code provided by Liu and Wang (2016) with its default heuristics. We plot the trajectories of particles of all three methods in Figure 2.

Although all three algorithms move particles toward the target distribution, the naive SVGD does not spread the particle mass quickly enough to cover the target distribution when using the same step size. This situation is much improved by applying the adaptive learning rate method (AdaGrad). In comparison, NW and LL both converge fast. After 20 iterations, all particles have arrived at the target positions. Since all methods are motivated by minimizing $\text{KL}[q_t, p]$, we plot the $\text{KL}[q_t, p]$ approximated by Donsker and Varadhan Lower Bound (Donsker and Varadhan, 1976) for all four methods. The plot of $\text{KL}[q_t, p]$ agrees with our qualitative assessments: AdaGrad SVGD can reduce the KL significantly faster than the vanilla SVGD with naive gradient descent. After 20 iterations, the KL divergence for both NW and LL particles reaches zero, indicating that the particles have fully converged to the target distribution. LL achieves a performance comparable to NW. This is a

remarkable result as NW, and SVGD has access to the true $\nabla \log p$, but LL only has samples from p . We also compare NW, LL and SVGD performance at different sample sizes ($n = 100, n = 250$). The results can be found in the Appendix O.4.

In the next sections, we will showcase the performance of LL in forward/backward KL minimization problems.

5.2. Joint Domain Adaptation

In domain adaptation, we want to use source domain samples to help a prediction task in a target domain. This addresses situations where the training data for a method may differ from the real data when deployed. We assume that source samples $\mathcal{D}_q := \{(\mathbf{x}_q^{(i)}, y_q^{(i)})\}_{i=1}^{n_q}$ are drawn from a joint distribution \mathbb{Q}_{XY} and target samples $\mathcal{D}_p := \{(\mathbf{x}_p^{(i)}, y_p^{(i)})\}_{i=1}^{n_p}$ are drawn from a different joint distribution \mathbb{P}_{XY} . However, y_p is missing from the target set. Thus, we want to predict missing labels in \mathcal{D}_p with the help of \mathcal{D}_q . Courty et al. (2017b;a) propose to find an optimal map that aligns the distribution \mathbb{Q}_{XY} and \mathbb{P}_{XY} , then train a classifier on the aligned source samples. Inspired by this method, we propose to align samples by minimizing $\text{KL}[q_t, p]$, where p is the density of the target \mathbb{P}_{XY} and q_t is a particle-label pair distribution whose samples are $\mathcal{D}_{q_t} := \{(\mathbf{x}_t^{(i)}, y_q^{(i)})\}_{i=1}^{n_q}$. To minimize $\text{KL}[q_t, p]$, we evolve \mathbf{x}_t according to the backward KL field and $\mathbf{x}_{t=0}$ is initialized to be the source input \mathbf{x}_q . In words, we transport source input samples so that the transported and target samples are aligned in terms of minimizing the backward KL divergence. After T iterations, we can train a classifier using transported source samples $\{(\mathbf{x}_T, y_q)\}$ to predict target labels.

One slight issue is that we do not have labels in the target domain but performing WGF requires joint samples (X, Y) . To solve this, we adopt the same approach used in Courty et al. (2017a), replacing y_p with a proxy $\hat{g}(\mathbf{x}_p)$, where \hat{g} is a prediction function trained to minimize an empirical transportation cost (See Section 2.2 in Courty et al. (2017a))⁴. We demonstrate our approach in a toy example, in Figure 3.

⁴If we know some labels y_p from the target domain, instead of using the proxy $\hat{g}(\mathbf{x}_p)$, we can use pairs (\mathbf{x}_p, y_p) to align the source and target samples.

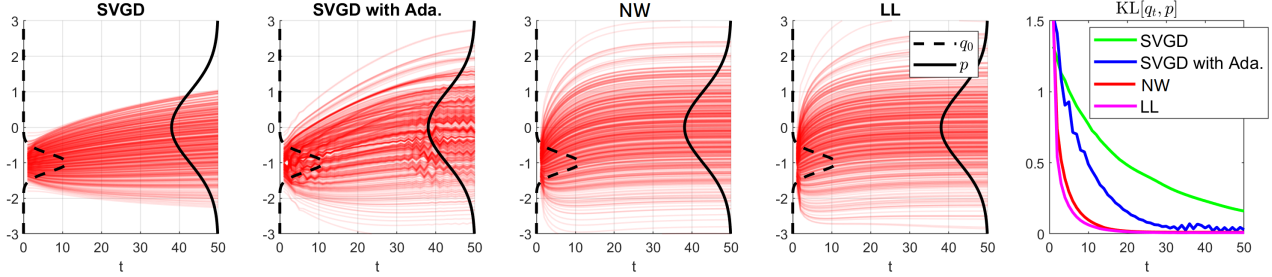


Figure 2. Particle Trajectories of SVGD, SVGD with AdaGrad, NW, LL. Approximated $\text{KL}[q_t, p]$ with different methods.

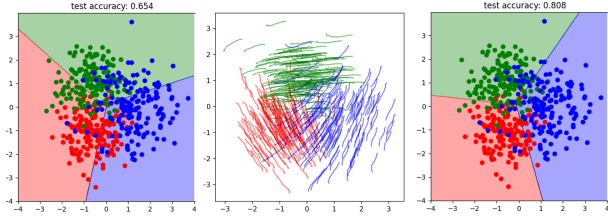


Figure 3. Left: the source classifier (represented by colored areas) misclassifies many testing points (colored dots). Middle: WGF moves particles to align the source and target samples. Lines are trajectories of sample movements in each class. Right: the retrained classifier on the transported source samples gives a much better prediction.

Table 2 compares the performance of adapted classifiers on a real-world 10-class classification dataset named “officecaltech-10”, where images of the same objects are taken from four different domains (amazon, caltech, dsrlr and webcam). We reduce the dimensionality by projecting all samples to a 100-dimensional subspace using PCA. We compare the performance of the **base** (the source RBF kernel SVM classifier), the Joint Distribution Optimal Transport (Courty et al., 2017a) (**JDOT**), an RBF kernel SVM trained on the WGF transported source samples $\{(x_T, y_q)\}$ (**WGF**) and an SVM trained on MMDFlow (Hagemann et al., 2024) transported source samples (**MMD**). The classification accuracy on the entire target sets are reported. It can be seen that in some cases, reusing the source classifiers in the target domain does lead to catastrophic results (e.g. amazon to dsrlr, caltech to dsrlr). However, we can avoid such performance decline by using any joint distribution-based domain adaptation. It can be seen that both **WGF** and **MMD** achieve superior performance compared to **JDOT** while **WGF** has the best performance in most adaptation cases.

5.3. Missing Data Imputation

In missing data imputation, we are given a joint dataset $\tilde{\mathcal{D}} := \{(\tilde{x}^{(i)}, \mathbf{m}^{(i)})\}$, where $\mathbf{m} \in \{0, 1\}^d$ is a mask vector and $m_j^{(i)} = 0$ indicates the j -th dimension of $\tilde{x}^{(i)}$ is missing. The task is to “guess” the missing values in \tilde{x} vector. In

$\mathcal{D}_q \rightarrow \mathcal{D}_p$	base	JDOT	WGF	MMD
amz. \rightarrow dsrlr	27.39%	65.61%	78.34%	78.34%
amz. \rightarrow web.	61.69%	67.80%	84.07%	89.15%
amz. \rightarrow cal.	81.66%	63.58%	82.72%	82.19%
dsrlr \rightarrow amz.	70.35%	72.96%	85.91%	76.30%
dsrlr \rightarrow web.	94.92%	76.61%	95.25%	86.10%
dsrlr \rightarrow cal.	58.95%	72.31%	79.25%	69.01%
web. \rightarrow amz.	75.78%	75.37%	91.34%	89.46%
web. \rightarrow dsrlr	94.27%	73.25%	98.73%	100.00%
web. \rightarrow cal.	67.05%	63.49%	78.09%	75.16%
cal. \rightarrow amz.	83.40%	84.55%	91.13%	89.04%
cal. \rightarrow dsrlr	26.11%	69.43%	84.71%	84.71%
cal. \rightarrow web.	63.39%	74.58%	80.34%	79.32%

Table 2. Domain adaptation of different domains in Office-Caltech-10 dataset.

recent years, GAN-based missing value imputation, e.g., **GAIN** (Yoon et al., 2018), has gained significant attention (Zhang et al., 2023). Let \mathbf{x} be an imputation of $\tilde{\mathbf{x}}$. The basic idea is that if \mathbf{x} is a perfect imputation of $\tilde{\mathbf{x}}$, then a classifier cannot predict m_j given $(\mathbf{x}, \mathbf{m}_{-j})$. For example, given a perfectly imputed image, one cannot tell which pixels are imputed and which pixels are observed. Therefore, in their approach, a generative network is trained to minimize the aforementioned classification accuracy. In fact, this method teaches the generator to *break the dependency between \mathbf{x} and \mathbf{m}* . Inspired by this idea, we propose to impute $\{\tilde{\mathbf{x}}\}$ by iteratively updating particles $\{\mathbf{x}_t\}$. The initial particle $\mathbf{x}_{t=0}$ is set to be

$$x_{0,j} = \begin{cases} \tilde{x}_j & m_j = 1 \\ \text{Sample from e.g., } \mathcal{N}(0, 1) & m_j = 0. \end{cases}$$

After that, the particles $\{\mathbf{x}_t\}$ are evolved according to the forward KL field that minimizes $\text{KL}[p_{X_t M}, p_{X_t} p_M]$, i.e., the mutual information between X_t (particles) and M (mask). Note that we only update missing dimensions, i.e.,

$$x_{t+1,j} := \begin{cases} x_{t,j} & m_j = 1 \\ x_{t,j} + \eta \partial_{x_j} r_t(\mathbf{x}_t, \mathbf{m}) & m_j = 0. \end{cases}$$

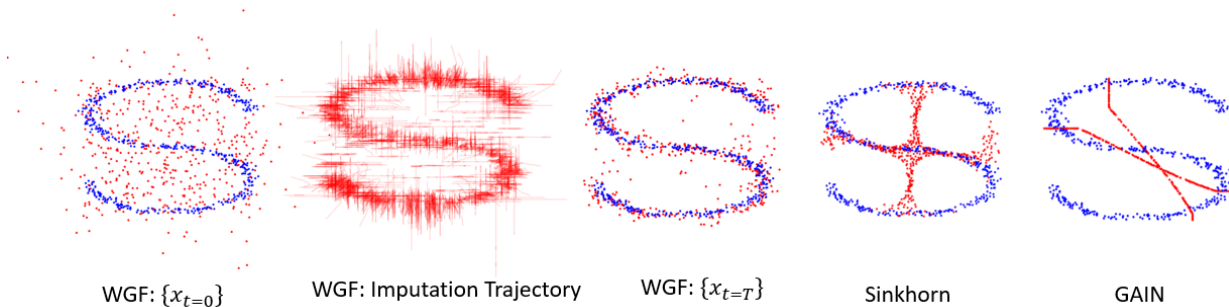


Figure 4. Comparison of imputation methods. Fully observed samples are plotted in blue, and imputed samples in red. The leftmost plot shows the initial particles in the WGF impute. The second left plot visualizes the imputation trajectories of different particles. The third left plot is the final output after 100 WGF iterations.

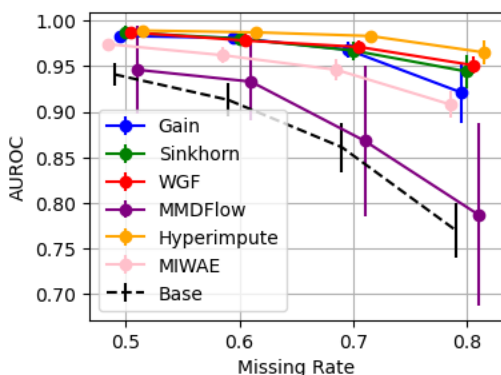


Figure 5. AUROC of a linear SVM classifier on the imputed Breast Cancer dataset. **Base** indicates the performance of a baseline imputer where we impute the missing values with Gaussian noises.

Samples from $p_{X_t M}$ are available to us since we observe the pairs $\{(x, m)\}$ and samples from $p_{X_t} p_M$ can be constructed as $\{(x, m')\}$, where m' is a random sample of M given the missing pattern.

In the first experiment, we test the performance of various imputers on an “S”-shaped dataset, where samples are Missing Completely at Random (MCAR) (Rubin, 1976). The results are plotted in Figure 4. We compare our imputed results (**WGF**) with Optimal Transport-based imputation method (**Sinkhorn**) (Muzellec et al., 2020), and GAN-based imputation method (**GAIN**) (Yoon et al., 2018). σ is chosen by automatic model selection described in Section I.1. It can be seen that our imputer, based on minimizing $\text{KL}[p_{X_t M}, p_{X_t} p_M]$ nicely recovers the “S”-shape after 100 particle update iterations. However, Sinkhorn and GAIN imputation methods struggle to accurately restore the “S”-shaped pattern.

In the second experiment, we test the performance of our algorithm on a real-world Breast Cancer classification dataset (Zwitter and Soklic, 1988) in Figure 5. This is a 30-

dimensional binary classification dataset and we artificially create missing values by following the MCAR paradigm with different missing rates. Since the dataset is a binary classification dataset, we compare the performance of linear SVM classifiers trained on imputed datasets. The performance is measured by the Area Under the ROC Curve (AUROC) on hold-out testing sets. In addition to **Sinkhorn** and **GAIN**, we validate the performance of our method with three additional algorithms: MMD flow using negative distance kernel (**MMDFlow**) (Hagemann et al., 2024), a model selection-based method **HyperImpute** (Jarrett et al., 2022), and an auto encoder-based approach (**MIWAE**) (Mattei and Frellsen, 2019). The result shows that SVM trained on the dataset imputed by our method achieves comparable performance to datasets imputed by the other benchmark methods. Our method is robust against the choice of σ . Details and the selection of hyperparameters can be found in Section I.2.

6. Limitations

All f -divergence fields are defined using the density ratio function p/q_t . However, the ratio function is not defined when the input domains of p and q_t are non-overlapping. This problem is particularly noticeable when working on high-dimensional, real-world datasets (such as images). This so-called “density chasm problem” (Rhodes et al., 2020) will also affect the velocity field estimation as assumptions required by our estimators e.g., Assumption 4.4 or (13) depends on the density ratio. One possible solution to this issue is to build a “bridge” using interpolations of two distributions, estimate the density ratios of “neighbouring” interpolations, and then combine these estimates. However, this will significantly increase the computational cost as we need to estimate the density ratio gradients of many pairs of distributions. Developing an efficient gradient estimator using target and particle distribution interpolation is an interesting future task. Our assumptions in Proposition 3.1 and Theorem 4.8 indicate that the algorithm’s effective-

ness depends on the boundedness of $\sup \|\nabla^2(h \circ r)\|$. This assumption is not necessarily true for WGF of some divergence over the entire real domain (e.g. See Appendix G.). This restriction suggests in some applications, these flows are non-ideal choices.

Another potential limitation of the proposed estimator is its applicability to high-dimensional datasets (such as the ImageNet dataset (Deng et al., 2009)). Although some preliminary investigations show that the proposed estimator does work reasonably well on some high-dimensional generative tasks (see Section O.2), it is known that local regression tends to perform poorly on high-dimensional datasets (see, e.g., (Stone, 1980; 1982)). Instead of performing WGF updates directly on the high-dimensional datasets, we can consider performing the updates in a low-dimensional feature space (see Section L). Some preliminary investigations have shown promising results (see Section O.3).

Finally, our method does not have a generator in the form of a neural network (like the one used by (Gao et al., 2019)), so we can only rely on the current particle set to estimate the vector field. This means that the sample size must be fixed before the flow. Thus, our method is better suited for applications such as domain adaptation and data imputation, where the number of samples to be transported is fixed before running the flow.

7. Conclusion

In recent years, it has been discovered that by iteratively updating a set of particles, WGF can approximate a target distribution by reducing the f -divergences between corresponding distributions. This paper addresses the important problem of estimating the velocity fields induced by various f -divergence WGFs. We propose novel interpolation-based estimators for different f -divergences fields and prove their validity. We demonstrate the effectiveness of these estimators through two novel applications: domain adaptation and missing data imputation. Our results show that even without access to the density ratio function, the velocity fields can be efficiently estimated from samples of the target and particle distributions.

Acknowledgments

We thank four anonymous reviewers for their feedback on our paper. JS was supported by the EPSRC Centre for Doctoral Training in Computational Statistics and Data Science, grant number EP/S023569/1.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal

consequences of our work, none of which we feel must be specifically highlighted here.

References

- L. Ambrosio, N. Gigli, and G. Savaré. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2005.
- A. F. Ansari, M. L. Ang, and H. Soh. Refining deep generative models via discriminator gradient flow. In *International Conference on Learning Representations (ICLR 2021)*, 2021.
- D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- S. Chewi, T. Le Gouic, C. Lu, T. Maunu, and P. Rigollet. Svdg as a kernelized wasserstein gradient flow of the chi-squared divergence. In *Advances in Neural Information Processing Systems (NeurIPS 2020)*, volume 33, pages 2098–2109, 2020.
- N. Courty, R. Flamary, A. Habrard, and A. Rakotomamonjy. Joint distribution optimal transportation for domain adaptation. In *Advances in Neural Information Processing Systems (NeurIPS 2017)*, volume 30, 2017a.
- N. Courty, R. Flamary, D. Tuia, and A. Rakotomamonjy. Optimal transport for domain adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(9):1853–1865, 2017b.
- J. Deng, W. Dong, R. Socher, L-J Li, K. Li, and F-F Li. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- M. D. Donsker and S. R. S. Varadhan. Asymptotic evaluation of certain markov process expectations for large time—iii. *Communications on Pure and Applied Mathematics*, 29(4):389–461, 1976.
- J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(7), 2011.
- J. Fan. Local Linear Regression Smoothers and Their Minimax Efficiencies. *The Annals of Statistics*, 21(1):196–216, 1993.
- Y. Gao, Y. Jiao, Y. Wang, Y. Wang, C. Yang, and S. Zhang. Deep generative learning via variational gradient flow. In *International Conference on Machine Learning (ICML 2019)*, pages 2093–2101, 2019.

- T. Gasser and H-G Müller. Kernel estimation of regression functions. In T. Gasser and M. Rosenblatt, editors, *Smoothing Techniques for Curve Estimation*, pages 23–68, Berlin, Heidelberg, 1979. Springer Berlin Heidelberg.
- J. Givens, S. Liu, and H. W. J. Reeve. Density ratio estimation and neyman pearson classification with missing data. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics (AISTATS 2023)*, volume 206, pages 8645–8681, 2023.
- P. Hagemann, J. Hertrich, F. Altekrüger, R. Beinert, J. Chemseddine, and G. Steidl. Posterior sampling based on gradient flows of the MMD with negative distance kernel. In *International Conference on Learning Representations (ICLR 2024)*, 2024.
- T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2001.
- W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American statistical association*, 58(301):13–30, 1963.
- A. Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6:695–709, 2005.
- A. Hyvärinen. Some extensions of score matching. *Computational Statistics & Data Analysis*, 51(5):2499–2512, 2007.
- D. Jarrett, B. Cebere, T. Liu, A. Curth, and M. van der Schaar. Hyperimpute: Generalized iterative imputation with automatic model selection. In *International Conference on Machine Learning (ICML 2022)*, volume 162, page 9916–9937, 2022.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR 2015)*, 2015.
- Q. Liu. Stein variational gradient descent as gradient flow. In *Advances in Neural Information Processing Systems (NeurIPS 2017)*, volume 30, pages 3118–3126, 2017.
- Q. Liu and D. Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In *Advances in Neural Information Processing Systems (NeurIPS 2016)*, volume 29, pages 2378–2386, 2016.
- D. Maoutsa, S. Reich, and M. Opper. Interacting particle solutions of fokker–planck equations through gradient–log–density estimation. *Entropy*, 22(8):802, 2020.
- P-A Mattei and J Frellsen. MIWAE: Deep generative modelling and imputation of incomplete data sets. In *International Conference on Machine Learning (ICML 2019)*, volume 97, pages 4413–4423, 2019.
- B. Muzellec, J. Josse, C. Boyer, and M. Cuturi. Missing data imputation using optimal transport. In *International Conference on Machine Learning (ICML 2020)*, pages 7130–7140, 2020.
- E. A. Nadaraya. On estimating regression. *Theory of Probability & Its Applications*, 9(1):141–142, 1964.
- X. Nguyen, M. J. Wainwright, and M. I. Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861, 2010.
- S. Nowozin, B. Cseke, and R. Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems (NeurIPS 2016)*, volume 29, 2016.
- B. Rhodes, K. Xu, and M. U. Gutmann. Telescoping density-ratio estimation. *Advances in Neural Information Processing Systems (NeurIPS 2020)*, 33:4905–4916, 2020.
- D. B. Rubin. Inference and Missing Data. *Biometrika*, 63(3):581–592, 1976. Publisher: Oxford University Press.
- D. W. Scott. Feasibility of multivariate density estimates. *Biometrika*, 78(1):197–205, 1991.
- J. Simons, S. Liu, and M. Beaumont. Variational likelihood-free gradient descent. In *Fourth Symposium on Advances in Approximate Bayesian Inference (AABI 2021)*, 2021.
- Y. Song, S. Garg, J. Shi, and S. Ermon. Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in Artificial Intelligence (UAI 2020)*, pages 574–584, 2020.
- C. J. Stone. Optimal rates of convergence for nonparametric estimators. *The Annals of Statistics*, pages 1348–1360, 1980.
- C. J. Stone. Optimal Global Rates of Convergence for Nonparametric Regression. *The Annals of Statistics*, 10(4):1040 – 1053, 1982.
- M. Sugiyama, T. Suzuki, and T. Kanamori. *Density Ratio Estimation in Machine Learning*. Cambridge University Press, 2012.
- P. Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661–1674, 2011.
- M. J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, 2019.
- Y. Wang, P. Chen, and W. Li. Projected wasserstein gradient descent for high-dimensional bayesian inference. *SIAM/ASA Journal on Uncertainty Quantification*, 10(4): 1513–1532, 2022.

- L. Wasserman. *All of Statistics: A Concise Course in Statistical Inference*. Springer Publishing Company, Incorporated, 2010.
- G. S. Watson. Smooth regression analysis. *Sankhyā: The Indian Journal of Statistics, Series A (1961-2002)*, 26(4): 359–372, 1964.
- A. Wibisono. Sampling as optimization in the space of measures: The langevin dynamics as a composite optimization problem. In *Conference on Learning Theory (COLT 2018)*, pages 2093–3027, 2018.
- M. Yi, Z. Zhu, and S. Liu. Monoflow: Rethinking divergence gans via the perspective of wasserstein gradient flows. In *International Conference on Machine Learning (ICML 2023)*, pages 39984–40000, 2023.
- J. Yoon, J. Jordon, and M. Schaar. Gain: Missing data imputation using generative adversarial nets. In *International Conference on Machine Learning (ICML 2018)*, pages 5689–5698, 2018.
- Q. Yu, A. Hashimoto, and Y. Ushiku. Divergence optimization for noisy universal domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR 2021)*, pages 2515–2524, 2021.
- Y. Zhang, R. Zhang, and B. Zhao. A systematic review of generative adversarial imputation network in missing data imputation. *Neural Computing and Applications*, 35 (27):19685–19705, September 2023.
- M. Zwitter and M. Soklic. Breast Cancer. UCI Machine Learning Repository, 1988.

A. Proof of Proposition 3.1

Proof.

$$\begin{aligned}
 |\hat{u}_i(\mathbf{x}^*) - \partial_i \log r(\mathbf{x}^*)| &= \left| \frac{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*) \partial_i \log r(\mathbf{x})]}{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)]} - \partial_i \log r(\mathbf{x}^*) \right| \\
 &= \left| \frac{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*) (\partial_i \log r(\mathbf{x}^*) + \langle \nabla \partial_i \log r(\bar{\mathbf{x}}), \mathbf{x} - \mathbf{x}^* \rangle)]}{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)]} - \partial_i \log r(\mathbf{x}^*) \right| \\
 &= \left| \frac{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*) \langle \nabla \partial_i \log r(\bar{\mathbf{x}}), \mathbf{x} - \mathbf{x}^* \rangle]}{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)]} \right| \\
 &\leq \frac{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*) \|\nabla \partial_i \log r(\bar{\mathbf{x}})\| \cdot \|\mathbf{x} - \mathbf{x}^*\|]}{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)]} \\
 &\leq \kappa \cdot \frac{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*) \|\mathbf{x} - \mathbf{x}^*\|]}{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)]} \\
 &\leq \kappa \cdot \left(\frac{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*) \|\mathbf{x} - \mathbf{x}^*\|]}{\widehat{\mathbb{E}}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)]} - \frac{\mathbb{E}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*) \|\mathbf{x} - \mathbf{x}^*\|]}{\mathbb{E}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)]} + \frac{\mathbb{E}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*) \|\mathbf{x} - \mathbf{x}^*\|]}{\mathbb{E}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)]} \right) \quad (17)
 \end{aligned}$$

The second line is due to the mean value theorem and $\bar{\mathbf{x}}$ is a point in between \mathbf{x} and \mathbf{x}^* in a coordinate-wise fashion. The second inequality is due to the operator norm of a matrix is always greater than a row/column norm.

Since $\text{Var}_q[\frac{1}{\sigma^d} k_\sigma(\mathbf{x}, \mathbf{x}^*) \|\mathbf{x} - \mathbf{x}^*\|] = O(\frac{1}{\sigma^d})$ due to the assumption, using Chebyshev inequality

$$\widehat{\mathbb{E}}_q \left[\frac{1}{\sigma^d} k_\sigma(\mathbf{x}, \mathbf{x}^*) \|\mathbf{x} - \mathbf{x}^*\| \right] - \mathbb{E}_q \left[\frac{1}{\sigma^d} k_\sigma(\mathbf{x}, \mathbf{x}^*) \|\mathbf{x} - \mathbf{x}^*\| \right] = O_p\left(\frac{1}{\sqrt{n\sigma^d}}\right).$$

Similarly, due to the assumption that $\text{Var}_q[\frac{1}{\sigma^d} k_\sigma(\mathbf{x}, \mathbf{x}^*)] = O(\frac{1}{\sigma^d})$, we have

$$\widehat{\mathbb{E}}_q \left[\frac{1}{\sigma^d} k_\sigma(\mathbf{x}, \mathbf{x}^*) \right] - \mathbb{E}_q \left[\frac{1}{\sigma^d} k_\sigma(\mathbf{x}, \mathbf{x}^*) \right] = O_p\left(\frac{1}{\sqrt{n\sigma^d}}\right).$$

Lemma A.1. Suppose $|\mathbb{E}A| \leq A_1 < \infty$, $|\mathbb{E}B| > B_0 > 0$. $\widehat{\mathbb{E}}A - \mathbb{E}A = O_p(\frac{1}{\sqrt{n\sigma^d}})$ and $\widehat{\mathbb{E}}B - \mathbb{E}B = O_p(\frac{1}{\sqrt{n\sigma^d}})$

$$\frac{\widehat{\mathbb{E}}A}{\widehat{\mathbb{E}}B} - \frac{\mathbb{E}A}{\mathbb{E}B} = O_p\left(\frac{1}{\sqrt{n\sigma^d}}\right).$$

Sketch of proof

The argument below resembles the proof technique used in A.2 in Givens et al. (2023). Here, we provide a sketch argument.

$$\left| \frac{\widehat{\mathbb{E}}A}{\widehat{\mathbb{E}}B} - \frac{\mathbb{E}A}{\mathbb{E}B} \right| = \left| \frac{\widehat{\mathbb{E}}A\mathbb{E}B - \mathbb{E}A\widehat{\mathbb{E}}B}{\widehat{\mathbb{E}}B\mathbb{E}B} \right| \leq \left| \frac{\widehat{\mathbb{E}}A - \mathbb{E}A}{\widehat{\mathbb{E}}B} \right| + \left| \frac{\widehat{\mathbb{E}}B - \mathbb{E}B}{\widehat{\mathbb{E}}B\mathbb{E}B} \right| \cdot |\mathbb{E}A|.$$

First, due to the boundedness of $|\mathbb{E}B|$, for any ϵ , $\exists N, \forall n > N$, $|\widehat{\mathbb{E}}B| > |\mathbb{E}B/2| > B_0/2$ with a probability $1 - \epsilon$. Thus, $\widehat{\mathbb{E}}A - \mathbb{E}A = O_p(\frac{1}{\sqrt{n\sigma^d}})$ implies $\left| \frac{\widehat{\mathbb{E}}A - \mathbb{E}A}{\widehat{\mathbb{E}}B} \right| = O_p(\frac{1}{\sqrt{n\sigma^d}})$.

Second, due to the boundedness of $|\mathbb{E}A|$, and the boundedness of $|\widehat{\mathbb{E}}B|$ argued above, $\widehat{\mathbb{E}}A - \mathbb{E}A = O_p(\frac{1}{\sqrt{n\sigma^d}})$ implies

$$\left| \frac{\widehat{\mathbb{E}}B - \mathbb{E}B}{\widehat{\mathbb{E}}B\mathbb{E}B} \right| \cdot |\mathbb{E}A| = O_p\left(\frac{1}{\sqrt{n\sigma^d}}\right).$$

Hence, $\frac{\widehat{\mathbb{E}}A}{\widehat{\mathbb{E}}B} - \frac{\mathbb{E}A}{\mathbb{E}B} = O_p\left(\frac{1}{\sqrt{n\sigma^d}}\right)$.

Due to Lemma A.1 and (17), we have with high probability that

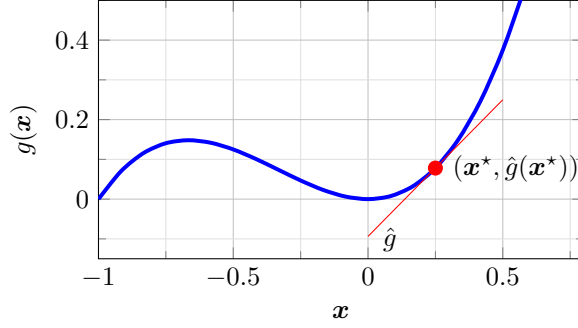


Figure 6. Fit g locally at \mathbf{x}^* using a linear function \hat{g} . Its “slope” is a natural estimator of $\nabla_{\mathbf{x}}g(\mathbf{x}^*)$.

$$\begin{aligned}
 |\hat{u}_i(\mathbf{x}^*) - \partial_i \log r(\mathbf{x}^*)| &\leq \kappa \cdot \left(\frac{K}{\sqrt{n\sigma^d}} + \frac{\mathbb{E}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)\|\mathbf{x} - \mathbf{x}^*\|]}{\mathbb{E}_q[k_\sigma(\mathbf{x}, \mathbf{x}^*)]} \right) \\
 &\leq \kappa \cdot \left(\frac{K}{\sqrt{n\sigma^d}} + \frac{\sigma \int q(\mathbf{x}^* + \sigma\mathbf{y})k(\mathbf{y})\|\mathbf{y}\|d\mathbf{y}}{\int q(\mathbf{x}^* + \sigma\mathbf{y})k(\mathbf{y})d\mathbf{y}} \right) & \mathbf{y} &:= (\mathbf{x} - \mathbf{x}^*)/\sigma \\
 &\leq \kappa \cdot \left(\frac{K}{\sqrt{n\sigma^d}} + \sigma C_k \right),
 \end{aligned}$$

where K is constant. □

B. Visualization of Gradient Estimation using Local Linear Fitting

C. Variational Objective for Estimating $h \circ r$

Proposition C.1. *The maximum in (9) is attained if and only if $d = h \circ r$.*

Proof. Due to the maximizing argument, the maximum of (9) is attained if and only if $d = \psi'$. The definition of mirror divergence indicates that $\psi' = r\phi'(r) - \phi(r)$. Theorem 2.1 states that $h \circ r = r\phi'(r) - \phi(r) = \psi'$, thus the maximum is attained if and only if $d = h \circ r$. □

D. Proof of Theorem 4.8, Estimation Error Bound of $w(\mathbf{x}^*)$

Proof. In this section, to simplify notations, we denote $\boldsymbol{\theta}$ as the parameter vector that combines both w and b , i.e., $\boldsymbol{\theta} := [w, b]^\top$. Specifically, we define

$$\boldsymbol{\theta}^* = [w^*, b^*]^\top := [\nabla^\top(h \circ r)(\mathbf{x}^*), h(r(\mathbf{x}^*)) - \langle \nabla(h \circ r)(\mathbf{x}^*), \mathbf{x}^* \rangle]^\top.$$

Let us denote the negative objective function in (10) as $\ell(\boldsymbol{\theta})$ and consider a constrained optimization problem:

$$\min_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}) \text{ subject to: } \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|^2 \leq \min(W, B)^2. \quad (18)$$

This convex optimization has a Lagrangian $\ell(\boldsymbol{\theta}) + \lambda(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|^2 - \min(W, B)^2)$, where $\lambda \geq 0$ is the Lagrangian multiplier. According to KKT condition, the optimal solution $\boldsymbol{\theta}_0$ of (18) satisfies $\nabla \ell(\boldsymbol{\theta}_0) + 2\lambda(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) = \mathbf{0}$.

We apply mean value theorem to $g(\boldsymbol{\theta}) := \langle \boldsymbol{\theta}_0 - \boldsymbol{\theta}^*, \nabla \ell(\boldsymbol{\theta}) + 2\lambda(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \rangle$:

$$\underbrace{\langle \boldsymbol{\theta}_0 - \boldsymbol{\theta}^*, \nabla \ell(\boldsymbol{\theta}_0) + 2\lambda(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) \rangle}_{g(\boldsymbol{\theta}_0)=0, \text{ KKT condition}} = \underbrace{\langle \boldsymbol{\theta}_0 - \boldsymbol{\theta}^*, \nabla \ell(\boldsymbol{\theta}^*) \rangle}_{g(\boldsymbol{\theta}^*)} + \underbrace{\langle (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)^\top [\nabla^2 \ell(\bar{\boldsymbol{\theta}}) + 2\lambda \mathbf{I}], (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) \rangle}_{\nabla g(\bar{\boldsymbol{\theta}})},$$

where $\bar{\theta}$ is a point between θ^* and θ_0 in an elementwise fashion. Since $\bar{\theta}$ is in a hyper cube with θ_0 and θ^* as opposite corners, it is in the constrain set of (18). Let us rearrange terms:

$$-\langle \theta_0 - \theta^*, \nabla \ell(\theta^*) \rangle = \langle \theta_0 - \theta^*, [\nabla^2 \ell(\bar{\theta}) + 2\lambda \mathbf{I}] (\theta_0 - \theta^*) \rangle \geq \langle \theta_0 - \theta^*, \nabla^2 \ell(\bar{\theta})(\theta_0 - \theta^*) \rangle,$$

For all θ in the constraint set, $\|\mathbf{w}\| = \|\mathbf{w} - \mathbf{w}^* + \mathbf{w}^*\| \leq \|\mathbf{w} - \mathbf{w}^*\| + \|\mathbf{w}^*\| \leq \|\theta - \theta^*\| + \|\mathbf{w}^*\| \leq 2W$ and $|b| = |b - b^* + b^*| \leq |b - b^*| + |b^*| \leq \|\theta - \theta^*\| + |b^*| \leq 2B$. Under our assumption (14), the lowest eigenvalue of $\nabla^2 \ell(\theta)$ for all θ in the constrain set of (18) is always lower bounded by $\sigma^d \Lambda_{\min}$. Therefore,

$$-\langle \theta_0 - \theta^*, \nabla \ell(\theta^*) \rangle \geq \sigma^d \cdot \Lambda_{\min} \|\theta_0 - \theta^*\|^2.$$

Using Cauchy–Schwarz inequality

$$\|\theta_0 - \theta^*\| \|\nabla \ell(\theta^*)\| \geq \sigma^d \cdot \Lambda_{\min} \|\theta_0 - \theta^*\|^2.$$

Assume $\|\theta_0 - \theta^*\|$ is not zero (if it is, our estimator is already consistent).

$$\begin{aligned} \|\theta_0 - \theta^*\| &\leq \frac{1}{\sigma^d \Lambda_{\min}} \|\nabla \ell(\theta^*)\| \\ &\leq \frac{1}{\sigma^d \Lambda_{\min}} \|\nabla \ell(\theta^*) - \mathbb{E} \nabla \ell(\theta^*) + \mathbb{E} \nabla \ell(\theta^*)\| \\ &\leq \frac{1}{\sigma^d \Lambda_{\min}} (\|\nabla \ell(\theta^*) - \mathbb{E} \nabla \ell(\theta^*)\| + \|\mathbb{E} \nabla \ell(\theta^*)\|) \end{aligned} \quad (19)$$

Since $\text{tr} [\text{Cov}_p [\frac{1}{\sigma^d} k(\mathbf{x}, \mathbf{x}^*) \cdot \tilde{\mathbf{x}}]] = O(\frac{1}{\sigma^d})$ and $\text{tr} [\text{Cov}_q [\frac{1}{\sigma^d} k(\mathbf{x}, \mathbf{x}^*) \cdot \psi_{\text{con}}(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*) \tilde{\mathbf{x}}]] = O(\frac{1}{\sigma^d})$ by assumption, due to multi-dimensional version of Chebyshev's inequality, $\sigma^d \|\frac{1}{\sigma^d} (\nabla \ell(\theta^*) - \mathbb{E} \nabla \ell(\theta^*))\| = \sigma^d O_p(\frac{1}{\sqrt{n\sigma^d}})$. Therefore,

$$\|\theta_0 - \theta^*\| \leq \frac{1}{\sigma^d \Lambda_{\min}} \left(\sigma^d \cdot \frac{K}{\sqrt{n\sigma^d}} + \|\mathbb{E} \nabla \ell(\theta^*)\| \right),$$

with high probability and K is a constant.

Now we proceed to bound $\|\mathbb{E} \nabla \ell(\theta^*)\|$.

Lemma D.1. $\|\mathbb{E} \nabla \ell(\theta^*)\| \leq \mathbb{E}_q \left[k_\sigma(\mathbf{x}, \mathbf{x}^*) \frac{1}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 \|\tilde{\mathbf{x}}\| \right] \kappa \cdot C_{\psi'_{\text{con}}}$

Proof. The expression of $\mathbb{E} \nabla \ell(\theta^*)$ is:

$$\mathbb{E} \nabla \ell(\theta^*) := -\mathbb{E}_p [k_\sigma(\mathbf{x}, \mathbf{x}^*) \cdot \tilde{\mathbf{x}}] + \mathbb{E}_q [k_\sigma(\mathbf{x}, \mathbf{x}^*) \psi'_{\text{con}}(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*) \cdot \tilde{\mathbf{x}}]. \quad (20)$$

Due to Taylor's theorem, $\langle \mathbf{w}^*, \mathbf{x} \rangle + b^* = h(r(\mathbf{x})) - \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \nabla^2 (h \circ r)(\bar{\mathbf{x}}) (\mathbf{x} - \mathbf{x}^*)$ where $\bar{\mathbf{x}}$ is a point in between \mathbf{x} and \mathbf{x}^* in an elementwise fashion. Thus, applying the mean value theorem on ψ'_{con} ,

$$\begin{aligned} \psi'_{\text{con}}(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*) &= \psi'_{\text{con}} \left[h(r(\mathbf{x})) - \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \nabla^2 (h \circ r)(\bar{\mathbf{x}}) (\mathbf{x} - \mathbf{x}^*) \right] \\ &= \psi'_{\text{con}} [h(r(\mathbf{x}))] - \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \nabla^2 (h \circ r)(\bar{\mathbf{x}}) (\mathbf{x} - \mathbf{x}^*) \psi''_{\text{con}}(y), \end{aligned}$$

where y is a scalar in between $h(r(\mathbf{x}))$ and $h(r(\mathbf{x})) - \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \nabla^2 (h \circ r)(\bar{\mathbf{x}}) (\mathbf{x} - \mathbf{x}^*)$ or equivalently, in between $h(r(\mathbf{x}))$ and $\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*$.

Theorem 2.1 states that $h(r) = r\phi' + \phi$, then, by the definition of the mirror divergence, $\psi' = h(r)$. Moreover, due to the maximizing argument, ψ'_{con} is the input argument of ψ (i.e., r) and ψ' is the input argument of ψ_{con} . Thus, $\psi'_{\text{con}}(h(r)) = \psi'_{\text{con}}(\psi'(r)) = r$. Let us write

$$\begin{aligned} \mathbb{E} \nabla \ell(\theta^*) &= -\mathbb{E}_p [k_\sigma(\mathbf{x}, \mathbf{x}^*) \tilde{\mathbf{x}}] + \mathbb{E}_q [k_\sigma(\mathbf{x}, \mathbf{x}^*) \underbrace{\psi'_{\text{con}}(h(r(\mathbf{x})))}_r \tilde{\mathbf{x}}] \\ &\quad - \mathbb{E}_q \left[k_\sigma(\mathbf{x}, \mathbf{x}^*) \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \nabla^2 (h \circ r)(\bar{\mathbf{x}}) (\mathbf{x} - \mathbf{x}^*) \psi''_{\text{con}}(y) \tilde{\mathbf{x}} \right] \\ &= -\mathbb{E}_q \left[k_\sigma(\mathbf{x}, \mathbf{x}^*) \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \nabla^2 (h \circ r)(\bar{\mathbf{x}}) (\mathbf{x} - \mathbf{x}^*) \psi''_{\text{con}}(y) \tilde{\mathbf{x}} \right], \end{aligned}$$

We can derive a bound for $\|\mathbb{E}\nabla\ell(\boldsymbol{\theta}^*)\|$:

$$\begin{aligned}
 \|\mathbb{E}\nabla\ell(\boldsymbol{\theta}^*)\| &\leq \mathbb{E}_q \left[k_\sigma(\mathbf{x}, \mathbf{x}^*) \cdot \frac{1}{2} \left| (\mathbf{x} - \mathbf{x}^*)^\top \nabla^2(h \circ r)(\bar{\mathbf{x}}) (\mathbf{x} - \mathbf{x}^*) \right| \cdot |\psi''_{\text{con}}(y)| \cdot \|\tilde{\mathbf{x}}\| \right] \\
 &\leq \mathbb{E}_q \left[k_\sigma(\mathbf{x}, \mathbf{x}^*) \frac{1}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 \|\nabla^2(h \circ r)(\bar{\mathbf{x}})\| \|\tilde{\mathbf{x}}\| \right] C_{\psi''_{\text{con}}} \\
 &\leq \mathbb{E}_q \left[k_\sigma(\mathbf{x}, \mathbf{x}^*) \frac{1}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 \|\tilde{\mathbf{x}}\| \right] \kappa \cdot C_{\psi''_{\text{con}}} \\
 &\leq \sigma^{d+2} \cdot \kappa \cdot C_{\psi''_{\text{con}}} \frac{1}{2} \int q(\sigma\mathbf{y} + \mathbf{x}^*) \cdot k(\mathbf{y}) \cdot \|\mathbf{y}\|^2 \cdot \|[\sigma\mathbf{y} + \mathbf{x}^*, 1]\| d\mathbf{y} \leq \sigma^{d+2} \cdot \kappa \cdot C_{\psi''_{\text{con}}} C_k
 \end{aligned}$$

□

Finally, due to Lemma D.1 and Assumption 4.5, we can see that with high probability

$$\|\boldsymbol{\theta}^* - \boldsymbol{\theta}_0\| \leq \frac{\frac{K\sigma^d}{\sqrt{n\sigma^d}} + \kappa \cdot \sigma^{d+2} \cdot C_k \cdot C_{\psi''_{\text{con}}}}{\sigma^d \Lambda_{\min}} \leq \frac{\frac{K}{\sqrt{n\sigma^d}} + \kappa \cdot \sigma^2 \cdot C_k \cdot C_{\psi''_{\text{con}}}}{\Lambda_{\min}}.$$

Since $n\sigma^d \rightarrow \infty, \sigma \rightarrow 0, \|\boldsymbol{\theta}^* - \boldsymbol{\theta}_0\| \rightarrow 0$ with high probability. There always exists σ_0 and N , such that for all $\sigma < \sigma_0$ and $n > N$, $\min(W, B) > \frac{\frac{K}{\sqrt{n\sigma^d}} + \kappa \cdot C_k \cdot C_{\psi''_{\text{con}}} \cdot \sigma^2}{\Lambda_{\min}}$. When it happens, $\boldsymbol{\theta}_0$ must be the interior of the constrain set of (18). i.e., the constraints in (18) are not active. It implies $\boldsymbol{\theta}_0$ must be the stationary point of $\ell(\boldsymbol{\theta})$ as long as σ is sufficiently small and n is sufficiently large.

□

E. Proof of Corollary 4.9

Proof. Since (11) has a unconstrained quadratic objective, its maximizers are stationary points.

We can see that Assumption 4.4 holds. To apply Theorem 4.8, we still need to show that Assumption 4.7 holds and W and B exist. In this case, $\psi_{\text{con}}(d) = d^2/2 + d$, so $\psi''_{\text{con}} = 1$. Thus Assumption 4.7 holds automatically for every $C_{\psi''_{\text{con}}} \geq 1$. Additionally,

$$\nabla_{[\mathbf{w}, b]}^2 \psi_{\text{con}}(\langle \mathbf{w}, \mathbf{x} \rangle + b) = \widehat{\mathbb{E}}_q [k_\sigma(\mathbf{x}, \mathbf{x}^*) \psi''_{\text{con}}(\langle \mathbf{w}, \mathbf{x} \rangle + b) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] = \widehat{\mathbb{E}}_q [k_\sigma(\mathbf{x}, \mathbf{x}^*) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top].$$

Therefore, (15) implies the minimum eigenvalue assumption (14) holds for every $W > 0, B > 0$. Thus, we can choose any $W > 0$ and $B > 0$ that satisfies (13). Noticing that $h(r(\mathbf{x})) = r(\mathbf{x})$, applying Theorem 4.8 gives the desired result. □

F. Proof of Corollary 4.10

Proof. Assumption 4.4, 4.5 and (13) are already satisfied. Let us verify the eigenvalue condition (14). In this case, $\psi''_{\text{con}}(d) = \exp(d - 1)$. Thus $\exp(\langle \mathbf{w}, \mathbf{x} \rangle + b - 1) \leq \exp(2WC_{\mathcal{X}} + 2B - 1) < \infty$ for $\|\mathbf{w}\| \leq 2W, |b| \leq 2B$. Moreover, because $\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top$ is positive semi-definite, due to (16),

$$\begin{aligned}
 \lambda_{\min} \left[\widehat{\mathbb{E}}_q [k_\sigma(\mathbf{x}, \mathbf{x}^*) \exp(\langle \mathbf{w}, \mathbf{x} \rangle + b - 1) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right] &\geq \sigma^2 \cdot \lambda_{\min} \left[\widehat{\mathbb{E}}_q [k_\sigma(\mathbf{x}, \mathbf{x}^*) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right] \cdot \exp(-2WC_{\mathcal{X}} - 2B - 1) \\
 &> \sigma^2 \cdot \Lambda_{\leftarrow} \cdot \exp(-2WC_{\mathcal{X}} - 2B - 1) > 0,
 \end{aligned}$$

for all \mathbf{w}, b that $\|\mathbf{w}\| \leq 2W, |b| \leq 2B$. So (14) holds. Finally, let us verify Assumption 4.7. Since ψ_{con} is a strictly monotone increasing function, $\sup_{a \in [0, 1]} \psi_{\text{con}}(ax_0 + (1-a)y_0)$ is obtained either at x_0 or y_0 . We only need to verify that $\psi_{\text{con}}(h(r(\mathbf{x})))$ and $\psi_{\text{con}}(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*)$ are both bounded for all \mathbf{x} . Both $h(r(\mathbf{x}))$ and $\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*$ can be bounded using our assumptions. Thus, for a

$$C_{\psi''_{\text{con}}} = \exp(W \cdot C_{\mathcal{X}} + B - 1) \vee \exp(C_{\log r} - 1)$$

Assumption 4.7 holds. Applying Theorem 4.8 completes the proof. □

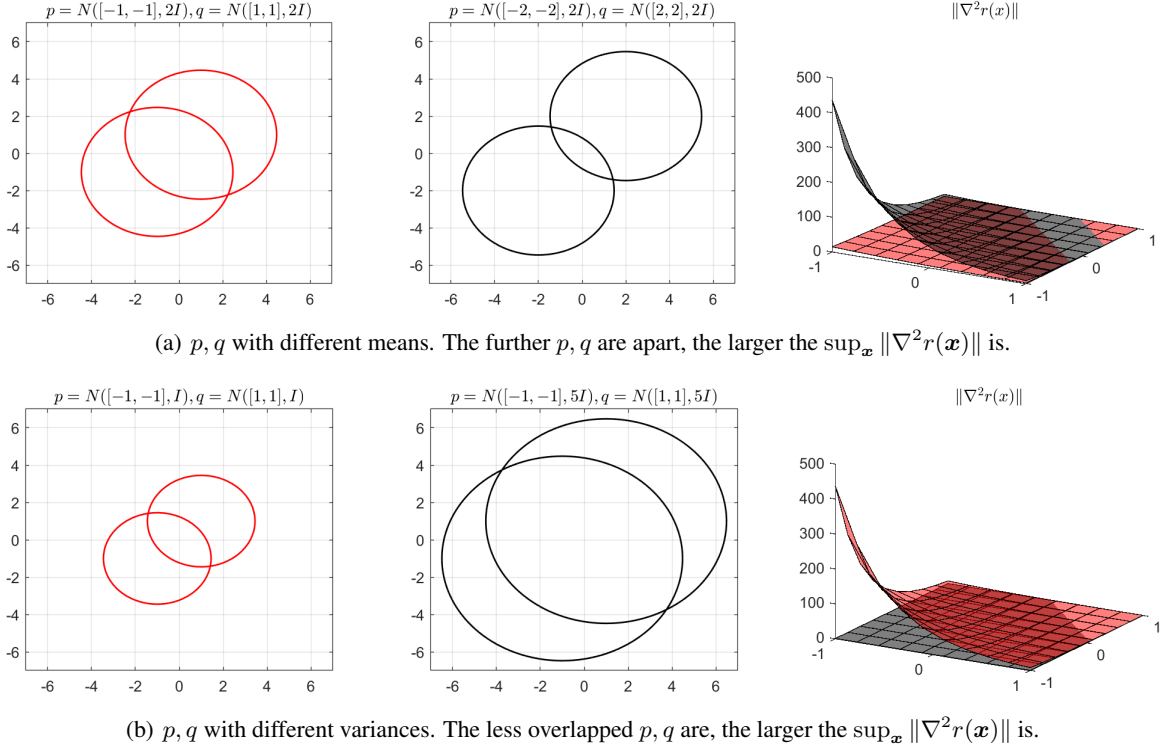


Figure 7. Visualizing $\nabla^2(h \circ r)(\mathbf{x})$ for $D_\phi = \text{KL}[p, q]$ with two different settings of p and q .

G. $\|\nabla^2 r(\mathbf{x})\|$ for Different p and q

See Figure 7.

H. Finite-sample Objectives

H.1. Practical Implementation

In our experiments, we observe that (10) can be efficiently minimized by using gradient descent with adaptive learning rate schemes, e.g., Adam (Kingma and Ba, 2015). One computational advantage of the local linear model is that the computation for each \mathbf{x}^* is independent from the others. This property allows us to parallelize the optimization. Even using a single CPU/GPU, we can easily write highly vectorized code to compute the gradient of (10) with respect to \mathbf{W} and \mathbf{b} for a large particle set $\{\mathbf{x}_i^*\}_{i=1}^n$.

Suppose $\mathbf{X}_p \in \mathbb{R}^{n_p \times d}$ and $\mathbf{X}_q \in \mathbb{R}^{n_q \times d}$ are the matrices whose rows are $\mathbf{x}_p^{(i)}$ and $\mathbf{x}_q^{(i)}$ respectively. $\mathbf{K}_p \in \mathbb{R}^{n_p \times n_p}$, $\mathbf{K}_q \in \mathbb{R}^{n_q \times n_q}$ are the kernel matrices between $\{\mathbf{x}^*\}$ and \mathbf{X}_p , $\{\mathbf{x}^*\}$ and \mathbf{X}_q respectively. $\mathbf{W} \in \mathbb{R}^{n \times d}$ and $\mathbf{b} \in \mathbb{R}^n$ are the parameters whose rows are $\mathbf{w}(\mathbf{x}^*)$ and $\mathbf{b}(\mathbf{x}^*)$ respectively. Then the gradient of (10) with respect to \mathbf{W} , \mathbf{b} can be expressed as

$$\begin{aligned} \mathbf{K}_p \mathbf{X}_p / n_p - \mathbf{K}_q \odot \psi'_{\text{con}}(\mathbf{W} \mathbf{X}_q^\top + \mathbf{b}) \mathbf{X}_q / n_q, \\ \mathbf{K}_p \mathbf{1}_{n_p} / n_p - \mathbf{K}_q \odot \psi'_{\text{con}}(\mathbf{W} \mathbf{X}_q^\top + \mathbf{b}) \mathbf{1}_{n_q} / n_q, \end{aligned}$$

where $\mathbf{1}_{n_p}$ and $\mathbf{1}_{n_q}$ are the vectors of ones with length n_p and n_q respectively. ψ'_{con} is evaluated element-wise. \odot is the element-wise product and the vector \mathbf{b} is broadcast to a matrix with n_q columns.

I. Experiment Details in Section 5

I.1. Model Selection

Let $\mathcal{D}_p := \{\mathbf{x}_p^{(i)}\}_{i=1}^{n_p}$, $\mathcal{D}_q := \{\mathbf{x}_q^{(i)}\}_{i=1}^{n_q}$ be training sets from p and q respectively and $\tilde{\mathcal{D}}_p = \{\tilde{\mathbf{x}}_p^{(i)}\}_{i=1}^{\tilde{n}_p}$ and $\tilde{\mathcal{D}}_q = \{\tilde{\mathbf{x}}_q^{(i)}\}_{i=1}^{\tilde{n}_q}$ be testing sets. We can fit a local linear model *at each testing point using the training sets*, i.e.,

$$\left(\hat{\mathbf{w}}_\sigma(\tilde{\mathbf{x}}), \hat{b}_\sigma(\tilde{\mathbf{x}})\right) := \underset{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}}{\operatorname{argmax}} \underbrace{\ell(\mathbf{w}, b; \tilde{\mathbf{x}}, \mathcal{D}_p, \mathcal{D}_q)}_{\text{training objective}}.$$

The dependency on σ comes from the smoothing kernel in the training objective. We can tune σ by evaluating the variational lower bound (9) approximated using testing samples:

$$\tilde{\ell}(\sigma; \tilde{\mathcal{D}}_p, \tilde{\mathcal{D}}_q) := \underbrace{\tilde{\mathbb{E}}_p[\hat{d}_\sigma(\mathbf{x})] - \tilde{\mathbb{E}}_q[\psi_{\text{con}}(\hat{d}_\sigma(\mathbf{x}))]}_{\text{testing criterion}}, \quad (21)$$

where $\tilde{\mathbb{E}}_p[\hat{d}(\mathbf{x})] := \frac{1}{\tilde{n}_p} \sum_{i=1}^{\tilde{n}_p} \hat{d}(\tilde{\mathbf{x}}_p^{(i)})$, i.e., the sample average over the testing points. $\hat{d}_\sigma(\tilde{\mathbf{x}}) := \langle \hat{\mathbf{w}}_\sigma(\tilde{\mathbf{x}}), \tilde{\mathbf{x}} \rangle + \hat{b}_\sigma(\tilde{\mathbf{x}})$, is the interpolation of $d(\tilde{\mathbf{x}})$ using training samples. The best choice of σ should maximize the above testing criterion.

In our experiments, we construct training and testing sets using cross validation and choose a list of candidate σ for the model selection. This procedure is parallel to selecting k in k -nearest neighbors to minimize the testing error. In our case, (21) is the “negative testing error”.

I.2. Missing Data Imputation

Before running the experiments, we first pre-process data in the following way:

1. Suppose \mathbf{X}_{true} is the original data matrix, i.e. without missing values. We introduce missingness to \mathbf{X}_{true} , and call the matrix with missing values $\tilde{\mathbf{X}}$, following MCAR paradigm. Denote the corresponding mask matrix as \mathbf{M} , where $m_j^{(i)} = 0$ if $\tilde{x}_j^{(i)}$ is missing, and $m_j^{(i)} = 1$ otherwise.
2. Calculate column-wise mean $\tilde{\mathbf{x}}$ and standard deviation $\tilde{\mathbf{s}}$ (excluding missing values) of $\tilde{\mathbf{X}}$.
3. Standardize $\tilde{\mathbf{X}}$ by taking $\tilde{\mathbf{X}} = \frac{\tilde{\mathbf{X}} - \tilde{\mathbf{x}}}{\tilde{\mathbf{s}}}$, where the vectors $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{s}}$ are broadcasted to the same dimensions as the matrix $\tilde{\mathbf{X}}$. Note that the division here is element-wise.

Denote \mathbf{X}^t as the imputed data of $\tilde{\mathbf{X}}$ at iteration i , where $\mathbf{X}^0 = \tilde{\mathbf{X}} \odot \mathbf{M} + \mathbf{Z} \odot (1 - \mathbf{M})$, and $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \text{diag}(\mathbf{1}))$.

We performed two experiments on both toy data (“S”-shape) and real world data (UCI Breast Cancer⁵ data).

Let N_{WGF} be the number of iterations WGF is performed. In each iteration, let N_{GradEst} be the number gradient descent steps for gradient estimation. In the missing data experiments, we set the hyperparameters to be:

- “S”-shape data: $T_{\text{WGF}} = 100$, $T_{\text{GradEst}} = 2000$, σ is chosen by model selection described in Section I.1.
- UCI Breast Cancer data: $T_{\text{WGF}} = 1000$, $T_{\text{GradEst}} = 100$, $\sigma = \text{median}\left(\sqrt{\frac{\text{pairwise distance of } \mathbf{X}}{2}}\right)$.

Across experiments, we observe that our method is robust against the choice of the bandwidth σ . We demonstrate this via a new experiment that tests the imputation performance with different selections of σ . Figure 8 demonstrates this experiment on the UCI Breast Cancer data. We can see that the imputation performance is still comparable with the original results (and baseline methods) when the scale of σ varies from 0.01 to 1.

⁵Available at <https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/wdbc.data>.

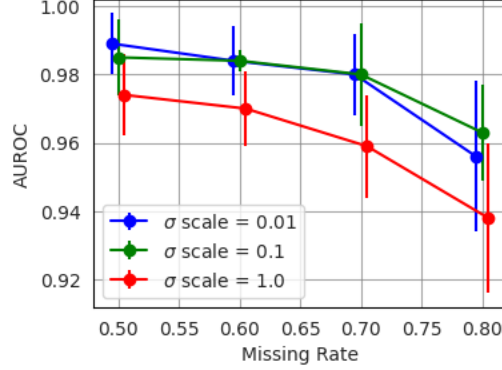


Figure 8. AUROC of a linear SVM classifier on the imputed Breast Cancer dataset, with various scales of σ from 0.01 to 1.

I.3. Wasserstein Gradient Flow

In this experiment, we first expand MNIST digits into 32×32 pictures then adds a small random noise $\epsilon \sim \mathcal{N}(\mathbf{0}, 0.001^2 \mathbf{I})$ to each picture so that computing the sample mean and covariance will not cause numerical issues. For both forward and backward KL WGF, we use a kernel bandwidth that equals to $1/5$ of the pairwise distances in the particle dataset, as it is too computationally expensive to perform cross validation at each iteration. After each update, we clip pixel values so that they are in between $[0, 1]$. It is done using pyTorch `torch.clamp` function.

To reduce computational cost, at each iteration, we randomly select 4000 samples from the original dataset and 4000 particles from the particle set. We use these samples to estimate the WGF updates.

J. Discussion: Kernel Density Gradient Estimation

The Kernel Density Estimator (KDE) of p is

$$\hat{p}(\mathbf{y}) := \frac{1}{n_p} \sum_{i=1}^{n_p} k(\mathbf{x}_p^{(i)}, \mathbf{y}) / Z,$$

where Z is a normalization constant to ensure that $\int \hat{p}(\mathbf{x}) d\mathbf{x} = 1$. Thus,

$$\nabla_{\mathbf{y}} \log \hat{p}(\mathbf{y}) := \frac{1}{n_p} \sum_{i=1}^{n_p} \nabla_{\mathbf{y}} k(\mathbf{x}_p^{(i)}, \mathbf{y}) / \frac{1}{n_p} \sum_{i=1}^{n_p} k(\mathbf{x}_p^{(i)}, \mathbf{y}).$$

The normalizing constant Z is cancelled.

K. Discussion: Why Score Matching does not Work on Log Ratio Gradient Estimation

For Score Matching (SM), the estimator of $\hat{p} := \operatorname{argmin}_f \int p \|\nabla \log p - \nabla \log f\|^2 d\mathbf{x}$, where the objective function is commonly referred to as *Fisher Divergence*. To use SM in practice, the objective function is further broken down to

$$\int p \|\nabla \log p - \nabla \log f\|^2 d\mathbf{x} = \int p \|\nabla \log f\|^2 d\mathbf{x} + 2 \sum_{i=1}^d \int p \partial_i^2 \log f d\mathbf{x} + C, \quad (22)$$

where we used the dimension-wise integration by parts and C is a constant.

Since our target is to estimate $\nabla \log p$, we can directly model $\nabla \log p$ as $\mathbf{g} : \mathbb{R}^d \rightarrow \mathbb{R}^d$. The objective becomes

$$\int p \|\mathbf{g}\|^2 d\mathbf{x} + 2 \sum_{i=1}^d \int p \partial_i g_i d\mathbf{x} + C.$$

SM can be used to estimate $\nabla \log p$. One might assume that SM can also be used for estimating $\nabla \log r$, where $r := \frac{p}{q}$. Let us replace $\nabla \log p$ with $\nabla \log r$ in (22),

$$\begin{aligned} \int p \|\nabla \log r - \nabla \log f\|^2 d\mathbf{x} &= \int p \|\nabla \log f\|^2 d\mathbf{x} - 2 \sum_{i=1}^d \int p \partial_i \log r \partial_i \log f d\mathbf{x} + C \\ &= \int p \|\nabla \log f\|^2 d\mathbf{x} + 2 \sum_{i=1}^d \int q \partial_i r \partial_i \log f d\mathbf{x} + C \\ &= \int p \|\nabla \log f\|^2 d\mathbf{x} + 2 \sum_{i=1}^d \int r \cdot \partial_i (q \partial_i \log f) d\mathbf{x} + C, \end{aligned} \quad (23)$$

$$= \int p \|\nabla \log f\|^2 d\mathbf{x} + 2 \sum_{i=1}^d \int r \cdot q \cdot \partial_i^2 \log f d\mathbf{x} + 2 \sum_{i=1}^d \int r \cdot \partial_i q \cdot \partial_i \log f d\mathbf{x} + C \quad (24)$$

and to get (23) we applied integration by parts, where we assumed $q \cdot r \cdot \partial_i \log f \rightarrow 0$ as $x_i \rightarrow \infty$. In (24), the third term is not tractable due to the lack of information about r and q . Changing the objective to $\int p \|\nabla \log r - \nabla \log f\|^2 d\mathbf{x}$ would also yield an intractable objective for a similar reason.

L. Discussion: Gradient Flow Estimation in Feature Space

One of the issues of local estimation is the *curse of dimensionality*: Local approximation does not work well in high dimensional spaces. However, since the f -divergence gradient flow is always associated with the density ratio function, we can utilize special structures in density ratio functions to estimate $\nabla(h \circ r)(\mathbf{x}^*)$ more effectively.

L.1. Density Ratio Preserving Map

Let $\mathbf{s}(\mathbf{x})$ be a measurable function, where $\mathbf{s} : \mathbb{R}^d \rightarrow \mathbb{R}^m, m \leq d$. Consider two random variables, X_p and X_q , each associated with probability density functions p and q , respectively. Define p° and q° as the probability density functions of the random variables $S_p := \mathbf{s}(X_p)$ and $S_q := \mathbf{s}(X_q)$.

Definition L.1. $\mathbf{s}(\mathbf{x})$ is a density ratio preserving map if and only if it satisfies the following equality

$$r(\mathbf{x}) = r^\circ(\mathbf{s}(\mathbf{x})), \text{ where } r^\circ(\mathbf{s}(\mathbf{x})) := \frac{p^\circ(\mathbf{s}(\mathbf{x}))}{q^\circ(\mathbf{s}(\mathbf{x}))}, \quad \forall \mathbf{x} \in \mathcal{X}.$$

We can leverage the density ratio preserving map to reduce the dimensionality of gradient flow estimation. Suppose $\mathbf{s}(\mathbf{x})$ is a known density ratio preserving map. Define $\mathbf{z} := \mathbf{s}(\mathbf{x})$ and $\mathbf{z}^* := \mathbf{s}(\mathbf{x}^*)$. We can see that

$$\underbrace{\nabla(h \circ r)(\mathbf{x}^*)}_{\mathbb{R}^d \mapsto \mathbb{R}^d} = \nabla(h \circ r^\circ)(\mathbf{s}(\mathbf{x}^*)) = \underbrace{\nabla^\top \mathbf{s}(\mathbf{x}^*)}_{\mathbb{R}^d \mapsto \mathbb{R}^{d \times m}, \text{ known}} \underbrace{\nabla(h \circ r^\circ)(\mathbf{z}^*)}_{\mathbb{R}^d \mapsto \mathbb{R}^m}. \quad (25)$$

If we can evaluate $\nabla \mathbf{s}(\mathbf{x}^*)$, we only need to estimate an m -dimensional gradient $\nabla(h \circ r^\circ)(\mathbf{z}^*)$, which is potentially easier than estimating the original d -dimensional gradient $\nabla(h \circ r)(\mathbf{x}^*)$ using a local linear model.

While Definition L.1 might suggest that \mathbf{s} is a very specific function, the requirement for \mathbf{s} to preserve the density ratio is quite straightforward. Specifically, \mathbf{s} must be *sufficient in expressing the density ratio function*. This requirement is formalized in the following proposition:

Proposition L.2. Consider a function $\mathbf{s} : \mathbb{R}^d \rightarrow \mathbb{R}^m$. If there exists a function $g : \mathbb{R}^m \rightarrow \mathbb{R}_+$ such that $r(\mathbf{x}) = g(\mathbf{s}(\mathbf{x})), \forall \mathbf{x} \in \mathcal{X}$ holds, then \mathbf{s} is a density ratio preserving map. Additionally, it follows that $g \circ \mathbf{s} = r = r^\circ \circ \mathbf{s} \implies g = r^\circ$.

The proof can be found in Section M.

Proposition L.2 implies that we can identify the density ratio preserving map \mathbf{s} by simply learning the ratio function r and using the trained feature transform function as \mathbf{s} . For instance, in the context of a neural network used to estimate r , \mathbf{s} could correspond to the functions represented by the penultimate layer of the network. After identifying \mathbf{s} , we can simply translate a high dimensional gradient flow estimation into a low dimensional problem according to (25).

Algorithm 1 Searching for a Density Ratio Preserving Map \mathbf{s}

```

1: Inputs:  $\mathcal{D}_p, \mathcal{D}_q$  and an initial guess of  $\hat{\mathbf{s}}$ .
2: while  $\hat{\mathbf{s}}$  not converged do
3:   for each  $\mathbf{x} \in \mathcal{D}_p \cup \mathcal{D}_q$  do
4:      $\mathbf{z} := \hat{\mathbf{s}}(\mathbf{x})$ 
5:      $(\hat{\mathbf{w}}(\mathbf{z}), \hat{\mathbf{b}}(\mathbf{z})) := \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^m, \mathbf{b} \in \mathbb{R}} \ell(\mathbf{w}, \mathbf{b}; \mathbf{z}, \mathcal{D}_{p^\circ}, \mathcal{D}_{q^\circ}) \mid_{\mathbf{s}=\hat{\mathbf{s}}}$ 
6:      $\hat{d}(\mathbf{z}) := \langle \hat{\mathbf{w}}(\mathbf{z}), \mathbf{z} \rangle + \hat{\mathbf{b}}(\mathbf{z})$ 
7:   end for
8:    $\hat{\mathbf{s}} := \operatorname{argmax}_{\mathbf{s} \in \mathbb{S}} \mathbb{E}_p[\hat{d}(\mathbf{s}(\mathbf{x}))] - \mathbb{E}_q[\psi_{\text{con}}(\hat{d}(\mathbf{s}(\mathbf{x})))]$ 
9: end while
10: Output:  $\hat{\mathbf{s}}$ .

```

In practice, we find this method works well. However, this approach still requires us to estimate a high dimensional density ratio function r to obtain a feature map \mathbf{s} .

In the next section, we propose an algorithm of learning \mathbf{s} from data without estimating a high dimensional density ratio function.

L.2. Finding Density Ratio Preserving Map

Theorem L.3. *Suppose h is associated with an f -divergence D_ϕ according to Theorem 2.1 and D_ψ is the mirror of D_ϕ . If $r(\mathbf{x}) = g^*(\mathbf{s}^*(\mathbf{x}))$, then \mathbf{s}^* must be an arg sup of the following objective:*

$$\sup_{\mathbf{s}} \mathbb{E}_p[h(r^\circ(\mathbf{s}(\mathbf{x})))] - \mathbb{E}_q[\psi_{\text{con}}(h(r^\circ(\mathbf{s}(\mathbf{x}))))]. \quad (26)$$

Proof. Since $r(\mathbf{x}) = g^*(\mathbf{s}^*(\mathbf{x}))$, Proposition C.1 implies that (g^*, \mathbf{s}^*) is necessarily an arg sup to the following optimization problem:

$$\sup_{\mathbf{s}, g} \mathbb{E}_p[h(g(\mathbf{s}(\mathbf{x})))] - \mathbb{E}_q[\psi_{\text{con}}(h(g(\mathbf{s}(\mathbf{x}))))].$$

Due to the law of unconscious statistician, $\mathbb{E}_p[f(\mathbf{s}(\mathbf{x}))] = \mathbb{E}_{p^\circ}[f(\mathbf{z})]$, where $\mathbf{z} = \mathbf{s}(\mathbf{x})$. The above optimization problem can be rewritten as

$$\sup_{\mathbf{s}} \sup_g \mathbb{E}_{p^\circ}[h(g(\mathbf{z}))] - \mathbb{E}_{q^\circ}[\psi_{\text{con}}(h(g(\mathbf{z})))] ,$$

Proposition C.1 states that for all \mathbf{s} , $g = r^\circ$ is an arg sup of the inner optimization problem. Substituting this optimal solution of g and rewriting the expectation using \mathbf{x} again, we arrive

$$\sup_{\mathbf{s}} \mathbb{E}_p[h(r^\circ(\mathbf{s}(\mathbf{x})))] - \mathbb{E}_q[\psi_{\text{con}}(h(r^\circ(\mathbf{s}(\mathbf{x}))))].$$

□

Both expectations in (26) can be approximated using samples from p and q . Given a fixed \mathbf{s} , $h(r^\circ(\mathbf{z}))$ can be approximated by an m -dimensional local linear interpolation

$$h(r^\circ(\mathbf{z})) \approx \hat{d}(\mathbf{z}) := \langle \hat{\mathbf{w}}(\mathbf{z}), \mathbf{z} \rangle + \hat{\mathbf{b}}(\mathbf{z}), \quad \left(\hat{\mathbf{w}}(\mathbf{z}), \hat{\mathbf{b}}(\mathbf{z}) \right) := \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^m, \mathbf{b} \in \mathbb{R}} \ell(\mathbf{w}, \mathbf{b}; \mathbf{z}, \mathcal{D}_{p^\circ}, \mathcal{D}_{q^\circ}), \quad (27)$$

where $\mathcal{D}_{p^\circ}, \mathcal{D}_{q^\circ}$ are sets of samples from p° and q° respectively.

Approximating expectations in (26) with samples in \mathcal{D}_p and \mathcal{D}_q and replacing $h \circ r^\circ$ with \hat{d} , we solve the following optimization to obtain an estimate of \mathbf{s} :

$$\hat{\mathbf{s}} := \operatorname{argmax}_{\mathbf{s} \in \mathbb{S}} \frac{1}{n_p} \sum_{i=1}^{n_p} \hat{d} \left[\mathbf{s} \left(\mathbf{x}_p^{(i)} \right) \right] - \frac{1}{n_q} \sum_{i=1}^{n_q} \psi_{\text{con}} \left[\hat{d} \left(\mathbf{s} \left(\mathbf{x}_q^{(i)} \right) \right) \right]. \quad (28)$$

The optimization of (28) is a bi-level optimization problem as \hat{d} depends on (27). We propose to divide the whole problem into two steps: First, let $\mathbf{s} = \hat{\mathbf{s}}$ and solve for $(\hat{\mathbf{w}}, \hat{b})$. Then, with the estimated $(\hat{\mathbf{w}}, \hat{b})$, we solve for $\hat{\mathbf{s}}$. Repeat the above procedure until convergence. This algorithm is detailed in Algorithm 1.

In practice, we restrict \mathbb{S} to be the set of all linear maps via a matrix $\mathbf{S} \in \mathbb{R}^{d \times m}$ whose columns are orthonormal basis, i.e., $\mathbf{s}(\mathbf{x}) := \mathbf{S}^\top \mathbf{x}$. The Jacobian $\nabla \mathbf{s}(\mathbf{x})$ is simply \mathbf{S}^\top .

After obtaining $\hat{\mathbf{s}}$, we can approximate the gradient flow using the chain rule described in (25):

$$\nabla(h \circ r)(\mathbf{x}^*) \approx \nabla^\top \hat{\mathbf{s}}(\mathbf{x}^*) \hat{\mathbf{w}}(\mathbf{z}^*),$$

where $\hat{\mathbf{w}}(\mathbf{z}^*)$ is approximated by an m -dimensional local linear interpolation

$$\left(\hat{\mathbf{w}}(\mathbf{z}^*), \hat{b}(\mathbf{z}^*) \right) := \underset{\mathbf{w} \in \mathbb{R}^m, b \in \mathbb{R}}{\operatorname{argmin}} \ell(\mathbf{w}, b; \mathbf{z}^*, \mathcal{D}_{p^\circ}, \mathcal{D}_{q^\circ}) \Big|_{\mathbf{s}=\hat{\mathbf{s}}}.$$

M. Discussion: Sufficient Condition of Density Ratio Preserving Map

In this Section, we provide a sufficient condition for \mathbf{s} to be a density ratio preserving map.

Lemma M.1. *If there exists some $g : \mathbb{R}^m \rightarrow \mathbb{R}_+$, such that $r(\mathbf{x}) = g(\mathbf{s}(\mathbf{x}))$, then \mathbf{s} is a density ratio preserving map and $g(\mathbf{s}(\mathbf{x})) = r^\circ(\mathbf{s}(\mathbf{x}))$.*

Proof. The statement $g(\mathbf{s}(\mathbf{x}))$ being a density ratio $\frac{p(\mathbf{x})}{q(\mathbf{x})}$ is equivalent to asserting that $\operatorname{KL}[p, q \cdot (g \circ \mathbf{s})] = 0$. Since KL divergence is always non-negative, it means

$$g = \underset{g: \int q(\mathbf{x})g(\mathbf{s}(\mathbf{x}))d\mathbf{x}=1, g \geq 0}{\operatorname{argmin}} \mathbb{E}_p \left[\log \frac{p(\mathbf{x})}{q(\mathbf{x})g(\mathbf{s}(\mathbf{x}))} \right] = \mathbb{E}_q[\log g(\mathbf{s}(\mathbf{x}))] + C_1, \quad (29)$$

i.e., g is a minimizer of $\operatorname{KL}[p, q \cdot (g \circ \mathbf{s})]$, where $g \geq 0$ is constrained in a domain where $q \cdot (g \circ \mathbf{s})$ is normalized to 1.

Similarly, $r^\circ(\mathbf{z})$ being a density ratio $\frac{p^\circ(\mathbf{z})}{q^\circ(\mathbf{z})}$ is the same as asserting that $\operatorname{KL}[p^\circ, q^\circ r^\circ] = 0$ and is equivalent to

$$r^\circ = \underset{g: \int q^\circ(\mathbf{z})g(\mathbf{z})d\mathbf{z}=1, g \geq 0}{\operatorname{argmin}} \mathbb{E}_{p^\circ} \left[\log \frac{p^\circ(\mathbf{z})}{q^\circ(\mathbf{z})g(\mathbf{z})} \right] = \mathbb{E}_{q^\circ}[\log g(\mathbf{z})] + C_2, \quad (30)$$

i.e., r° is a minimizer of $\operatorname{KL}[p^\circ, q^\circ g]$.

In fact, one can see that (29) and (30) are identical optimization problems due to the law of the unconscious statistician: $\int q(\mathbf{x})g(\mathbf{s}(\mathbf{x}))d\mathbf{x} = \mathbb{E}_q[g(\mathbf{s}(\mathbf{x}))] = \mathbb{E}_{q^\circ}[(g(\mathbf{z}))] = \int q^\circ(\mathbf{z})g(\mathbf{z})d\mathbf{z}$ and $\mathbb{E}_q[\log g(\mathbf{s}(\mathbf{x}))] = \mathbb{E}_{q^\circ}[\log g(\mathbf{z})]$, which means their solution sets are the same. Therefore, for any g that minimizes (29), it must also minimize (30). Hence it satisfies the following equality $r^\circ(\mathbf{s}(\mathbf{x})) = g(\mathbf{s}(\mathbf{x})) = r(\mathbf{x})$, where the second equality is by our assumption. \square

N. Discussion: Stein Variational Gradient Descent

SVGD minimizes $\operatorname{KL}[q_{t+1}, p]$, where samples of q_{t+1} is constructed using the following deterministic rule:

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \eta \mathbf{u}_t(\mathbf{x}_t). \quad (31)$$

\mathbf{x}_t are particles at iteration t , $\mathbf{u}_t \in \mathcal{H}^d$, a d -dimensional Reproducing Kernel Hilbert Space (RKHS) with a kernel function $l(\mathbf{x}, \mathbf{x}')$. Liu and Wang (2016) shows the optimal update \mathbf{u}_t has a closed form:

$$\mathbf{u}_t^{\operatorname{svgd}} := \mathbb{E}_{q_t} [l(\mathbf{x}_t, \cdot) \nabla \log p(\mathbf{x}_t) + \nabla l(\mathbf{x}_t, \cdot)]. \quad (32)$$

In practice, expectations $\mathbb{E}_{q_t}[\cdot]$ can be approximated by $\widehat{\mathbb{E}}_{q_t}[\cdot]$, i.e., the sample average taken from the particles at time t .

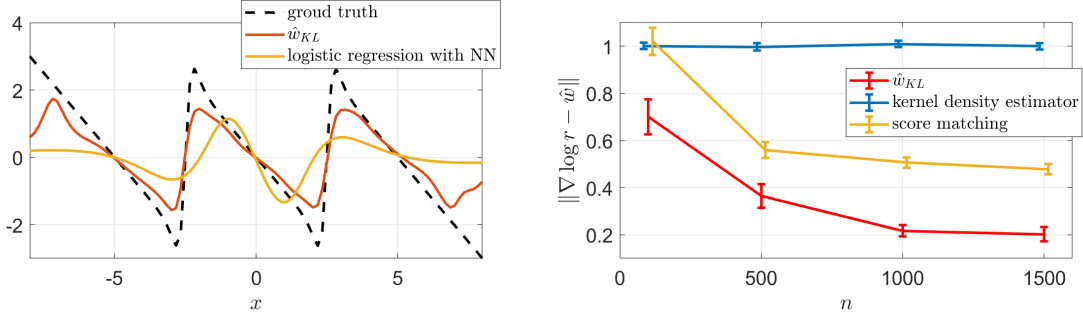


Figure 9. Left: $\nabla \log r(\mathbf{x})$ and its approximations. Right: Estimation error with standard error.

Chewi et al. (2020) links SVGD with f -divergence WGF: $\mathbf{u}_t^{\text{svgd}}$ is the backward KL divergence WGF under the coordinate-wise transform of an integral operator. Indeed, the i -th dimension of SVGD update $u_{t,i}^{\text{svgd}}$ can be expressed as

$$\begin{aligned}
 u_{t,i}^{\text{svgd}} &:= \mathbb{E}_{q_t} [l(\mathbf{x}_t, \cdot) \partial_i \log p(\mathbf{x}_t)] + \mathbb{E}_{q_t} [\partial_i l(\mathbf{x}_t, \cdot)] \\
 &= \mathbb{E}_{q_t} [l(\mathbf{x}_t, \cdot) \partial_i \log p(\mathbf{x}_t)] - \mathbb{E}_{q_t} [l(\mathbf{x}_t, \cdot) \partial_i \log q_t(\mathbf{x}_t)] \\
 &= \int q_t(\mathbf{x}_t) l(\mathbf{x}_t, \cdot) \partial_i \log \frac{p(\mathbf{x}_t)}{q_t(\mathbf{x}_t)} d\mathbf{x}_t,
 \end{aligned} \tag{33}$$

where the last line is an integral operator (Wainwright, 2019) of the functional $\partial_i \log \frac{p}{q}$, i.e., the i -th dimension of the backward KL divergence flow $\nabla \log r_t$.⁶ Due to the reproducing property of RKHS, the SVGD update at some fixed point \mathbf{x}^* can be written as

$$\mathbf{u}^{\text{svgd}}(\mathbf{x}^*) = \langle \mathbf{u}_t^{\text{svgd}}, l(\cdot, \mathbf{x}^*) \rangle_{\mathcal{H}^d} = \mathbb{E}_q [l(\mathbf{x}, \mathbf{x}^*) \nabla \log r(\mathbf{x})].$$

O. Additional Experiments

O.1. Gradient Estimation

Now we investigate the performance of estimating $\nabla \log r(\mathbf{x})$ using the proposed gradient estimator and an indirect estimator using logistic regression. For the indirect estimator, we first train a Multilayer Perceptron (MLP) using a binary logistic regression to approximate $\log r$. Then obtain $\nabla \log r$ by auto-differentiating the estimated log ratio. The kernel bandwidth in our method is tuned by using the model selection criterion described in Section 4.4.

To conduct the experiments, we let $p = \mathcal{N}(-5, .5)/3 + \mathcal{N}(0, .5)/3 + \mathcal{N}(5, .5)/3$ and $q = \mathcal{N}(-5, 1)/3 + \mathcal{N}(0, 1)/3 + \mathcal{N}(5, 1)/3$. From each distribution, 5000 samples are generated for approximating the gradients.

The left plot in Figure 9 shows the true gradient and its approximations. It can be seen that the direct gradient estimation is more accurate than estimating the log ratio first then taking the gradient.

The right plot in Figure 9 displays the estimation errors of different methods, comparing the proposed method with Kernel Density Estimation (KDE) and score matching, all applied to the same distributions in the previous experiment. KDE was previously used in approximating WGF (Wang et al., 2022). It first estimates p and q with \hat{p} and \hat{q} separately using non-parametric kernel density estimators, then approximates $\nabla \log r$ with $\nabla \log \hat{p} - \nabla \log \hat{q}$. The score matching approximates $\nabla \log \hat{p}$ and $\nabla \log \hat{q}$ with the minimizers of Fisher-divergence. It has also been used in simulating particle ODEs in a previous work (Maoutsa et al., 2020). The estimation error plot shows that the proposed estimator yields more accurate results compared to the other two kernel-based gradient estimation methods, namely KDE and score matching.

O.2. Wasserstein Gradient Flow for Generative Modelling

In this experiment, we test the performance of the proposed gradient estimators by generating samples from a high dimensional target distribution (MNIST handwritten digits). We will check whether the quality of the particles can be improved by performing WGC using our estimated updates. Note that we do not intend to compare the generated samples

⁶Note that the second equality in (33) is due to the integration by parts, and only holds under conditions that $\lim_{\|\mathbf{x}\| \rightarrow \infty} q_t(\mathbf{x}) = 0$.

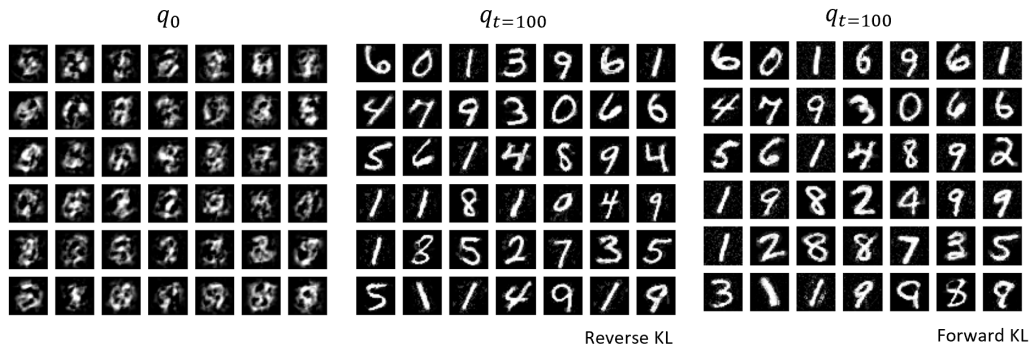


Figure 10. Samples generated using forward and backward KL velocity field.

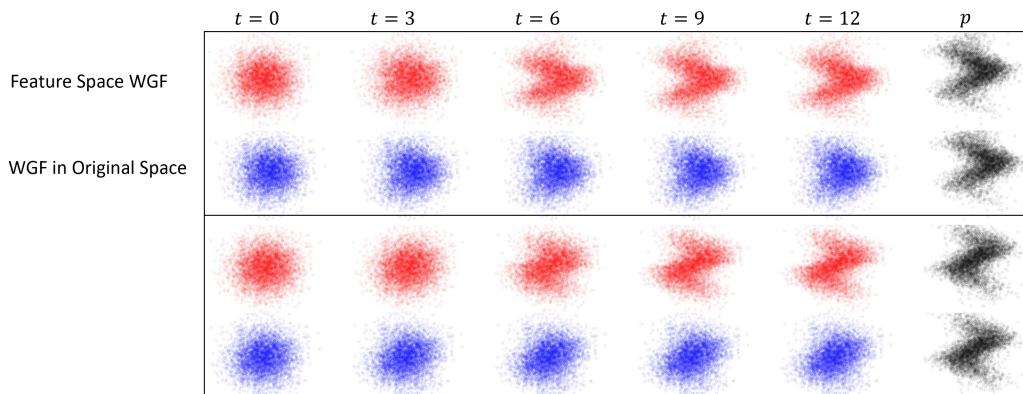


Figure 11. $\{\mathbf{x}_{q_t}\}$ in the first two dimensions x_1, x_2 . Feature space WGF converges in fewer steps compared to WGF in the original space. Above: $g = \cos$. Below: $g = \sin$.

with NN-based approaches, as the focus of our paper is on local estimation using kernel functions. We perform two different WGFs, forward KL and backward KL whose updates are approximated using \hat{w}_{\rightarrow} and \hat{w}_{\leftarrow} respectively. We let the initial particle distribution q_0 be $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ are the mean and covariance of the target dataset and fix the kernel bandwidth using “the median trick” (see the appendix for details).

The generated samples together with samples from the initial distribution q_0 are shown in Figure 10. Judging from the generated sample quality, it can be seen that both VGDs perform well and both have made significant improvements from the initial samples drawn from q_0 .

We also provide two videos showing the animation of the first 100 gradient steps:

- Forward KL: <https://youtube.com/shorts/HZcvUykrpbc>
- Backward KL: <https://youtube.com/shorts/AgN6dsDecCM>

O.3. Feature Space Wasserstein Gradient Flow

In this experiment, we run WGF in a 5-dimensional space. The target distribution is

$$p = p_{12}(\mathbf{x}_{1,2})p_{345}(\mathbf{x}_{3,4,5}), p_{345} := \mathcal{N}(\mathbf{0}, \mathbf{I})$$

and p_{12} is constructed by a generative process. We generate samples in the first two dimensions as follows

$$X_1 = g(X_2) + \epsilon, X_2 \sim \mathcal{N}(0, 1), \epsilon \sim \mathcal{N}(0, 1).$$

In our experiment, we draw 5000 samples from p , and 5000 samples from $q_0 := \mathcal{N}(\mathbf{0}, \mathbf{I})$, and run the backward KL gradient flow. Clearly, this WGF has a low dimensional structure since p and q only differs in the first two dimensions. We also run

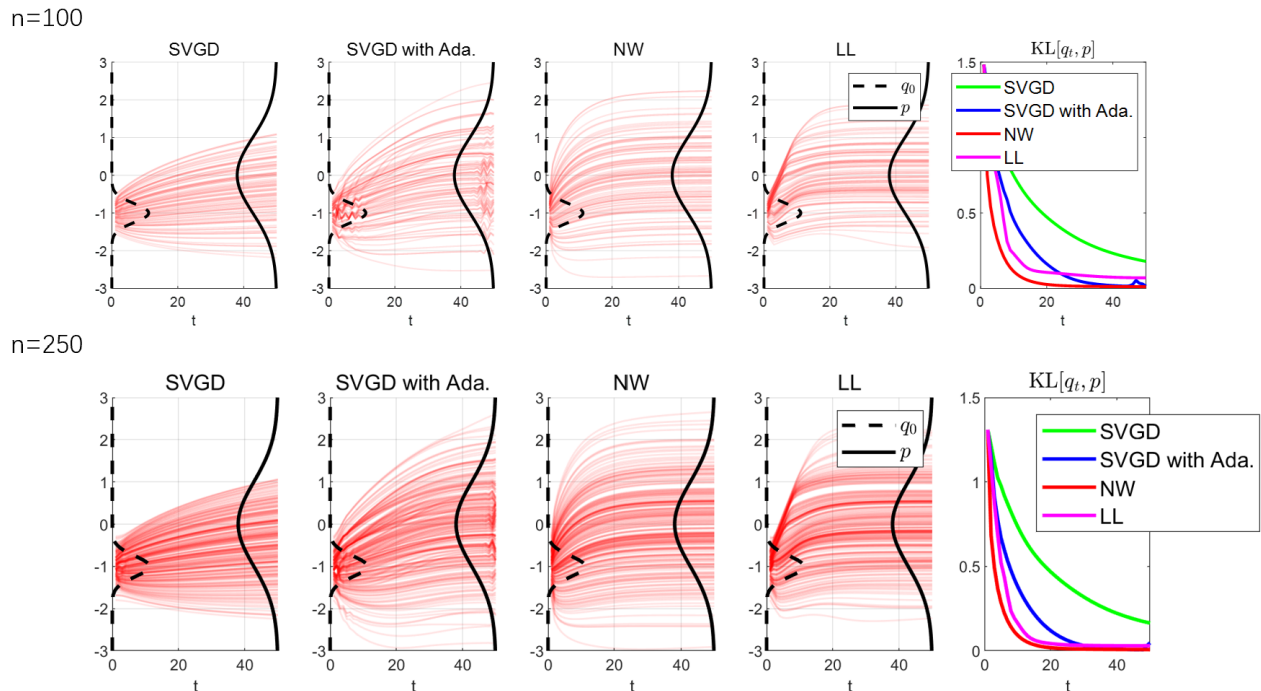


Figure 12. NW, LL, SVGD compared with sample size $n = 100, 250$.

feature space backward KL field whose updates are calculated using (25). The feature function $s(x) = S^\top x$ is learned by Algorithm 1.

The resulting particle evolution for both processes is plotted in Figure 11. For visualization purposes, we only plot the first two dimensions. It can be seen that the particles converge much faster when we explicitly exploit the subspace structure using the feature space WGF. In comparison, running WGF in the original space converges much slower.

In this experiments, we set the learning rates for both WGF and feature space WGF to be 0.1 and the kernel bandwidth σ in our local estimators is tuned using cross validation with a candidate set ranging from 0.1 to 2.

O.4. NW, LL and SVGD with Different Sample Sizes

See Figure 12.