

Distributional Sliced-Wasserstein and Applications to Generative Modeling

Khai Nguyen[◊] Nhat Ho[†] Tung Pham[◊] Hung Bui[◊]

VinAI Research[◊], University of Texas, Austin[†]

October 6, 2020

Abstract

Sliced-Wasserstein distance (SW) and its variant, Max Sliced-Wasserstein distance (Max-SW), have been used widely in the recent years due to their fast computation and scalability even when the probability measures lie in a very high dimensional space. However, SW requires many unnecessary projection samples to approximate its value while Max-SW only uses the most important projection, which ignores the information of other useful directions. In order to account for these weaknesses, we propose a novel distance, named *Distributional Sliced-Wasserstein* distance (DSW), that finds an *optimal* distribution over projections that can balance between exploring distinctive projecting directions and the informativeness of projections themselves. We show that the DSW is a generalization of Max-SW, and it can be computed efficiently by searching for the optimal push-forward measure over a set of probability measures over the unit sphere satisfying certain regularizing constraints that favor distinct directions. Finally, we conduct extensive experiments with large-scale datasets to demonstrate the favorable performances of the proposed distances over the previous sliced-based distances in generative modeling applications.

1 Introduction

Optimal transport (OT) is a classical problem in mathematics and operation research. Due to its appealing theoretical properties and flexibility in practical applications, it has recently become an important tool in the machine learning and statistics community; see for example, [12, 3, 40, 19] and references therein. The main usage of OT is to provide a distance named Wasserstein distance, to measure the discrepancy between two probability distributions. However, that distance suffers from expensive computational complexity, which is the main obstacle to using OT in practical applications.

There have been two main approaches to overcome the high computational complexity problem: either approximate the value of OT or apply the OT adaptively to specific situations. The first approach was initiated by Cuturi [14] using an entropic regularizer to speed up the computation of the OT [39, 23]. The entropic regularization approach has demonstrated its usefulness in several application domains [13, 18, 9]. Along this direction, several works proposed efficient algorithms for solving the entropic OT [1, 31, 30] as well as methods to stabilize these algorithms [11, 36, 11, 38]. However, these algorithms have complexities of the order $\mathcal{O}(k^2)$, where k is the number of atoms. It is expensive when we need to compute the OT repeatedly, especially in learning the data distribution.

The second approach, known as "slicing", takes a rather different perspective. It leverages two key ideas: the OT closed-form expression for two distributions in one-dimensional space, and the transformation of a distribution into a set of projected one-dimensional distributions by the Radon transform (RT) [20]. The popular proposal along this direction is Sliced-Wasserstein (SW) distance [7],

which samples the projecting directions uniformly over a unit sphere in the data ambient space and takes the expectation of the resulting one-dimensional OT distance. The SW distance hence requires a significantly lower computation cost than the original Wasserstein distance and is more scalable than the first approach. Due to its solid statistical guarantees and efficient computation, the SW distance has been successfully applied to a variety of practical tasks [16, 33, 26, 43, 15] where it has been shown to have comparative performances to other distances and divergences between probability distributions. However, there is an inevitable bottleneck of computing the SW distance. Specifically, the expectation with respect to the uniform distribution of projections in SW is intractable to compute; therefore, the Monte Carlo method is employed to approximate it. Nevertheless, drawing from a uniform distribution of directions in high-dimension can result in an overwhelming number of irrelevant directions, especially when the actual data lies in a low-dimensional manifold. Hence, SW typically needs to have a large number of samples to yield an accurate estimation of the discrepancy. Alternatively, in the other extreme, Max Sliced-Wasserstein (Max-SW) distance [15]) uses only one important direction to distinguish the probability distributions. However, other potentially relevant directions are ignored in Max-SW. Therefore, Max-SW can miss some important differences between the two distributions in high dimension. We note that the linear projections in the Radon transform can be replaced by non-linear projections resulting in the generalized sliced-Wasserstein distance and its variants [6, 25].

Apart from these main directions, there are also few proposals that try either to modify them or to combine the advantages of the above-mentioned approaches. In particular, Paty et al. [35] extended the idea of the max-sliced distance to the max-subspace distance by considering finding an optimal orthogonal subspace. However, this approach is computationally expensive, since it could not exploit the closed-form of the one-dimensional Wasserstein distance. Another approach named the Projected Wasserstein distance (PWD), which was proposed in [37], uses sliced decomposition to find multiple one-dimension optimal transport maps. Then, it computes the average cost of those maps equally in the original dimension.

Our contributions. Our paper also follows the slicing approach. However, we address key friction in this general line of work: how to obtain a relatively small number of slices simultaneously to maintain the computational efficiency, but at the same time, cover the major differences between two high-dimensional distributions. We take a probabilistic view of slicing by using a probability measure on the unit sphere to represent how important each direction is. From this viewpoint, SW uses the uniform distribution while Max-SW searches for the best delta-Dirac distribution over the projections, both can be considered as special cases. In this paper, we propose to search for an optimal distribution of important directions. We regularize this distribution such that it prefers directions that are far away from one another, hence encouraging an efficient exploration of the space of directions. In the case of no regularization, our proposed method recovers max-(generalized) SW as a special case. In summary, our main contributions are two-fold:

1. First, we introduce a novel distance, named *Distributional Sliced-Wasserstein distance* (DSW), to account for the issues of previous sliced distances. Our main idea is to search for not just a single most important projection, but an *optimal* distribution over projections that could balance between an expansion of the area around important projections and the informativeness of projections themselves, i.e., how well they can distinguish the two target probability measures. We show that DSW is a proper metric in the probability space and possesses appealing statistical and computational properties as the previous sliced distances.

- Second, we apply the DSW distance to generative modeling tasks based on the generative adversarial framework. The extensive experiments on real and large-scale datasets show that DSW distance significantly outperforms the SW and Max-SW distances under similar computational time on these tasks. Furthermore, the DSW distance helps model distribution converge to the data distribution faster and provides more realistic generated images than the SW and Max-SW distances.

Organization. The remainder of the paper is organized as follows. In Section 2, we provide backgrounds for Wasserstein distance and its slice-based versions. In Section 3, we propose distributional (generalized) sliced-Wasserstein distance and analyze some of its theoretical properties. Section 4 includes extensive experiment results followed by discussions in Section 5. Finally, we defer the proofs of key results and extra materials in the Appendices.

Notation. For any $\theta, \theta' \in \mathbb{R}^d$, $\cos(\theta, \theta') = \frac{\theta^\top \theta'}{\|\theta\| \|\theta'\|}$, where $\|\cdot\|$ is ℓ_2 norm. For any $d \geq 2$, \mathbb{S}^{d-1} denotes the unit sphere in d dimension in ℓ_2 norm. Furthermore, δ denotes the Dirac delta function, and $\langle \cdot, \cdot \rangle$ is the Euclidean inner-product. For any $p \geq 1$, $\mathbb{L}^p(\mathbb{R}^d)$ is the set of real-valued functions on \mathbb{R}^d with finite p -th moment.

2 Background

In this section, we provide necessary backgrounds for the (generalized) Radon transform, the Wasserstein, and sliced-Wasserstein distances.

2.1 Wasserstein distance

We start with a formal definition of Wasserstein distance. For any $p \geq 1$, we define $\mathcal{P}_p(\mathbb{R}^d)$ as the set of Borel probability measures with finite p -th moment defined on a given metric space $(\mathbb{R}^d, \|\cdot\|)$. For any probability measures μ, ν defined on $\mathcal{X}, \mathcal{Y} \subseteq \mathbb{R}^d$, we denote their corresponding probability density functions as I_μ and I_ν . The Wasserstein distance of order p between μ and ν is given by [41, 36]:

$$W_p(\mu, \nu) := \left(\inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} \|x - y\|^p d\pi(x, y) \right)^{\frac{1}{p}},$$

where $\Pi(\mu, \nu)$ is a set of all transportation plans π such that the marginal distributions of π are μ and ν , respectively. In order to simplify the presentation, we abuse the notation by using both $W_p(\mu, \nu)$ and $W_p(I_\mu, I_\nu)$ interchangeably for the Wasserstein distance between μ and ν .

When μ and ν are *one-dimension* measures, the Wasserstein distance between μ and ν has a closed-form expression $W_p(\mu, \nu) = (\int_0^1 |F_\mu^{-1}(z) - F_\nu^{-1}(z)|^p dz)^{1/p}$ where F_μ and F_ν are the cumulative distribution function (CDF) of I_μ and I_ν , respectively.

2.2 (Generalized) Radon transforms

Now, we review (generalized) Radon transform maps, which are key to the notion of (generalized) sliced-Wasserstein distance and its variants. The *Radon transform* (RT) maps a function $I \in \mathbb{L}^1(\mathbb{R}^d)$ to the space of functions defined over space of lines in \mathbb{R}^d . In particular, for any $t \in \mathbb{R}$ and direction $\theta \in \mathbb{S}^{d-1}$, the RT is defined as follows [20]: $\mathcal{R}I(t, \theta) := \int_{\mathbb{R}^d} I(x) \delta(t - \langle x, \theta \rangle) dx$.

The *generalized* Radon transform (GRT) [6] extends the original one from integration over hyperplanes of \mathbb{R}^d to integration over hypersurfaces. In particular, it is defined as: $\mathcal{G}I(t, \theta) :=$

$\int_{\mathbb{R}^d} I(x)\delta(t - g(x, \theta))dx$, where $t \in \mathbb{R}$ and $\theta \in \Omega_\theta$. Here, Ω_θ is a compact subset of \mathbb{R}^d and $g : \mathbb{R}^d \times \mathbb{S}^{d-1} \mapsto \mathbb{R}$ is a defining function (cf. Assumptions H1-H4 in [25] for the definition of defining function) inducing the hypersurfaces. When $g(x, \theta) = \langle x, \theta \rangle$ and $\Omega_\theta = \mathbb{S}^{d-1}$, the generalized Radon transform becomes the standard Radon transform.

2.3 (Generalized) sliced-Wasserstein distances

The sliced-Wasserstein distance (SW) between two probability measures μ and ν is defined as [7] :

$$SW_p(\mu, \nu) := \left(\int_{\mathbb{S}^{d-1}} W_p^p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta))d\theta \right)^{1/p}.$$

Similarly, the generalized sliced-Wasserstein distance [25] (GSW) is given by $GSW_p(\mu, \nu) := \left(\int_{\Omega_\theta} W_p^p(\mathcal{G}I_\mu(\cdot, \theta), \mathcal{G}I_\nu(\cdot, \theta))d\theta \right)^{1/p}$, where Ω_θ is the compact set of feasible parameter. However, these integrals are usually intractable. Thus, they are often approximated by using Monte Carlo scheme to draw uniform samples $\{\theta_i\}_{i=1}^N$ from \mathbb{S}^{d-1} and Ω_θ . In particular, $SW_p^p(\mu, \nu) \approx \frac{1}{N} \sum_{i=1}^N W_p^p(\mathcal{R}I_\mu(\cdot, \theta_i), \mathcal{R}I_\nu(\cdot, \theta_i))$ and $GSW_p^p(\mu, \nu) \approx \frac{1}{N} \sum_{i=1}^N W_p^p(\mathcal{G}I_\mu(\cdot, \theta_i), \mathcal{G}I_\nu(\cdot, \theta_i))$. In order to obtain a good approximation of (generalized) SW distances, N needs to be sufficiently large. However, important directions are not distributed uniformly over the sphere. Thus, this approach will draw potentially many unimportant projections that are not only expensive but also greatly reduce the effect of the SW distance.

2.4 Max (generalized) sliced-Wasserstein distances

An approach to using only informative directions is to simply take the best slice in discriminating two given probability distributions. That distance is max sliced-Wasserstein distance (Max-SW) [15], which is given by:

$$\max SW_p(\mu, \nu) := \max_{\theta \in \mathbb{S}^{d-1}} W_p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta)).$$

By combining this idea with non-linear projections from generalized Radon transform, we obtain max generalized sliced-Wasserstein distance (Max-GSW) [25]. The formal definition of that distance is: $\max GSW_p(\mu, \nu) := \max_{\theta \in \Omega_\theta} W_p(\mathcal{G}I_\mu(\cdot, \theta), \mathcal{G}I_\nu(\cdot, \theta))$. The (generalized) Max-SW distances focus on finding only the most important direction. Meanwhile, other informative directions play no role in the distance. Therefore, (generalized) Max-SW distances can ignore useful information about the structure of high dimensional probability measures.

3 Distributional Sliced-Wasserstein Distance

With the aim of improving the limitations of the previous sliced distances, we propose a novel distance, named *Distributional Sliced-Wasserstein distance* (DSW), that can search for not just a single but a distribution of important directions on the unit sphere. We prove that it is a well-defined metric and discuss its connection to the existing sliced-based distances in Section 3.1. Then, we provide a procedure to approximate DSW based on its dual form in Section 3.2.

3.1 Definition and metricity

We first start with a definition of distributional sliced-Wasserstein distance. We say $C > 0$ *admissible* if the set \mathbb{M}_C of probability measures σ on \mathbb{S}^{d-1} satisfying $\mathbb{E}_{\theta, \theta' \sim \sigma} [|\theta^\top \theta'|] \leq C$ is not empty.

Definition 1. Given two probability measures μ and ν on \mathbb{R}^d with finite p -th moments where $p \geq 1$ and an admissible regularizing constant $C > 0$. The distributional sliced-Wasserstein distance (DSW) of order p between μ and ν is given by:

$$DSW_p(\mu, \nu; C) := \sup_{\sigma \in \mathbb{M}_C} \left(\mathbb{E}_{\theta \sim \sigma} \left[W_p^p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta)) \right] \right)^{\frac{1}{p}}, \quad (1)$$

where \mathcal{R} is the Radon transform operator.

The DSW aims to find the optimal probability measure of slices on the unit sphere \mathbb{S}^{d-1} . Note that, the Max-SW distance is equivalent to searching for the best Dirac measure on a single point in \mathbb{S}^{d-1} , which puts all weights in only one direction. Meanwhile, the uniform measure in the formulation of SW distance distributes the same weights in all directions. Indeed, the uniform and Dirac measures are two special cases, because they view that either all directions are equally important or only one direction is important. That view is too restricted if the data actually lie on low dimensional space. Thus, we aim to find a probability measure which concentrates only on areas around important directions. Furthermore, we do not want these directions to lie in only one small area, because under the orthogonal projection of RT, their corresponding one-dimensional distributions will become similar. In order to achieve this, we search for an optimal measure σ that satisfies the regularization constraint $\mathbb{E}_{\theta, \theta' \sim \sigma} [|\theta^\top \theta'|] \leq C$. By Cauchy-Schwarz inequality, C is no greater than 1, thus \mathbb{M}_1 contains all probability measures on the unit sphere. Optimizing over \mathbb{M}_1 simply returns the best Dirac measure corresponding to the Max-SW distance. When C is small, the constraint forces the measure σ to distribute more weights to directions that are far from each other (in terms of their angles). Thus, a small appropriate value of C will help to balance between the distinctiveness and informativeness of these targeted directions. For further discussion about C , see Appendix B.1.

Next, we show that DSW is a well-defined metric on the probability space.

Theorem 1. For any $p \geq 1$ and admissible $C > 0$, $DSW_p(\cdot, \cdot; C)$ is a well-defined metric in the space of Borel probability measures with finite p -th moment. In particular, it is non-negative, symmetric, identity, and satisfies the triangle inequality.

The proof of Theorem 1 is in Appendix A.1. Our next result establishes the topological equivalence between DSW distance and (max)-sliced Wasserstein and Wasserstein distances.

Theorem 2. For any $p \geq 1$ and admissible $C > 0$, the following holds

(a) $DSW_p(\mu, \nu; C) \leq \max SW_p(\mu, \nu) \leq W_p(\mu, \nu)$.

(b) If $C \geq 1/d$, we have $DSW_p(\mu, \nu; C) \geq \left(\frac{1}{d}\right)^{1/p} \max SW_p(\mu, \nu) \geq \left(\frac{1}{d}\right)^{1/p} W_p(\mu, \nu)$.

As a consequence, when $p \geq 1$ and $C \geq 1/d$, $DSW_p(\cdot, \cdot; C)$, SW_p , $\max SW_p$, and W_p are topologically equivalent, namely, the convergence of probability measures under $DSW_p(\cdot, \cdot; C)$ implies the convergence of these measures under other metrics and vice versa.

The proof of Theorem 2 is in Appendix A.2. As a consequence of Theorem 2, the statistical error of estimating the unknown distribution based on the empirical distribution of n i.i.d data under DSW distance is $C_d \cdot n^{-1/2}$ with high probability where C_d is some universal constant depending on dimension d (see Appendix B.3). Therefore, as other sliced-based Wasserstein distances, the DSW distance does not suffer from the curse of dimensionality.

3.2 Computation of DSW distance

Direct computation of DSW distance is challenging. Hence we consider a dual form of DSW distance and a reparametrization of σ as follows.

Definition 2. For any $p \geq 1$ and admissible $C > 0$, there exists a non-negative constant λ_C depending on C such that the dual form of DSW distance takes the following form

$$DSW_p^*(\mu, \nu; C) = \sup_{\sigma \in \mathbb{M}} \left\{ \left(\mathbb{E}_{\theta \sim \sigma} \left[W_p^p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta)) \right] \right)^{\frac{1}{p}} - \lambda_C \mathbb{E}_{\theta, \theta' \sim \sigma} [|\theta^\top \theta'|] \right\} + \lambda_C C,$$

where \mathbb{M} denotes the space of all probability measures on the unit sphere \mathbb{S}^{d-1} .

By the Lagrangian duality theory, $DSW_p(\mu, \nu; C) \geq DSW_p^*(\mu, \nu; C)$ for any $p \geq 1$ and admissible $C > 0$. In Definition 2, the set \mathbb{M}_C disappears and λ_C plays the tuning role for the regularized term $\mathbb{E}_{\theta, \theta' \sim \sigma} [|\theta^\top \theta'|]$. When λ_C is large, $\mathbb{E}_{\theta, \theta' \sim \sigma} [|\theta^\top \theta'|]$ needs to be small, meaning that C is small. When λ_C is small, the value of $\mathbb{E}_{\theta, \theta' \sim \sigma} [|\theta^\top \theta'|]$ becomes less important, i.e., C is large.

For reparametrizing the measure σ , we use a pushforward of uniform measure on the unit sphere through some measurable function f . In particular, let f be a Borel measurable function from \mathbb{S}^{d-1} to \mathbb{S}^{d-1} . For any Borel set $A \subset \mathbb{S}^{d-1}$, we define $\sigma(A) = \sigma^{d-1}(f^{-1}(A))$, where σ^{d-1} is the uniform probability measure on \mathbb{S}^{d-1} . Then for any Borel measurable function $g: \mathbb{S}^{d-1} \rightarrow \mathbb{R}$, we have $\int_{\theta \sim \sigma} g(\theta) d\sigma(\theta) = \int_{\theta \sim \sigma^{d-1}} (g \circ f)(\theta) d\sigma^{d-1}(\theta)$. Therefore, we obtain the equivalent dual form of DSW as follows:

$$DSW_p^*(\mu, \nu; C) = \sup_{f \in \mathcal{F}} \left\{ \left(\mathbb{E}_{\theta \sim \sigma^{d-1}} \left[W_p^p(\mathcal{R}I_\mu(\cdot, f(\theta)), \mathcal{R}I_\nu(\cdot, f(\theta))) \right] \right)^{1/p} - \lambda_C \mathbb{E}_{\theta, \theta' \sim \sigma^{d-1}} [|f(\theta)^\top f(\theta')|] \right\} + \lambda_C C := \sup_{f \in \mathcal{F}} DS(f), \quad (2)$$

where \mathcal{F} is a class of all Borel measurable functions from \mathbb{S}^{d-1} to \mathbb{S}^{d-1} .

Finding the optimal f : We parameterize f in the dual form (2) by using a deep neural network with parameter ϕ , defined as f_ϕ . Then, we estimate the gradient of the objective function $DS(f_\phi)$ in equation (2) with respect to ϕ and use stochastic gradient ascent algorithm to update ϕ . Since there are expectations over uniform distribution in the gradient of $DS(f_\phi)$, we use the Monte Carlo method to approximate these expectations. Note that, we can use the fixed point from the stochastic ascent algorithm to approximate the dual value of DSW in equation (2). A detailed argument for this point is in Appendix B.2. Finally, in generative model applications with DSW being the loss function, we only need to use the gradient of the function $DS(\cdot)$ to update the parameters of interest. Therefore, we can treat λ_C as a regularized parameter and tune it to find suitable value in these applications.

Illustration of the roles of λ_C and C : To illustrate the roles of λ_C and C in finding optimal distribution σ , we conduct a simple experiment on two Gaussian distributions with zero means and

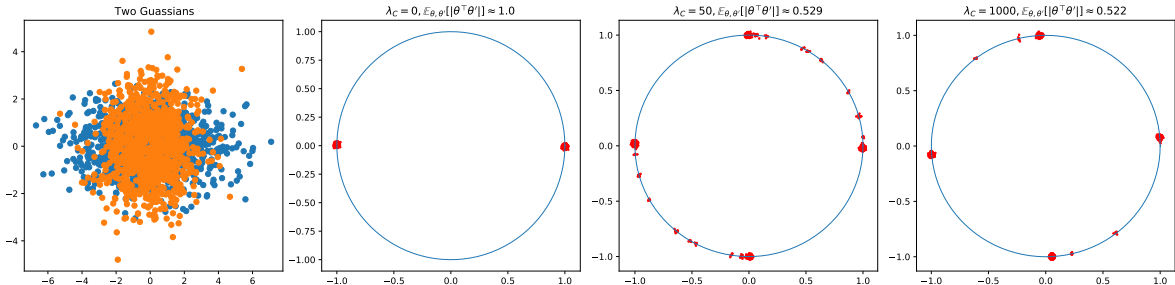


Figure 1: Empirical behavior of optimal measure σ , approximated by 1000 samples, on a circle for different values of λ_C (the constant in the dual form of DSW in Definition 2) when μ and ν are bivariate Gaussian distributions sharing the same eigenvectors. When $\lambda_C = 0$, $C = 1$. When λ_C increases, C becomes small.

covariance matrices given by $\begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$ and $\begin{pmatrix} 5 & 0 \\ 0 & 1 \end{pmatrix}$. The experiment optimizes the empirical form of Definition 2 with different choices of λ_C . The results are shown in Figure 1 with the reported value of λ_C and $\mathbb{E}_{\theta, \theta' \sim \sigma}[|\theta^\top \theta'|]$. For $\lambda_C = 0$, the obtained distribution concentrates only on one direction. When $\lambda_C = 50$, optimal σ distributes more weights to other directions on the circle. When $\lambda_C = 1000$, optimal σ is close to the discrete distribution concentrated on two eigenvectors of the covariance matrices, which are the main directions differentiating the two Gaussian distributions.

Extension of DSW and comparison of DSW to Max-GSW-NN: Similar to SW, we extend DSW to distributional generalized sliced Wasserstein (DGSW) by using the non-linear projecting operator via GRT. The definition of the DGSW and its properties are in Appendix C. Finally, in Appendix E.1, we show the distinction of the DSW to Max-GSW-NN [25] when the neural network defining function in Max-GSW-NN is $g(x, \theta) = \langle x, f(\theta) \rangle$ where $f: \mathbb{S}^{d-1} \rightarrow \mathbb{S}^{d-1}$.

4 Experiments

In this section, we conduct extensive experiments comparing the performance in both generative quality and computational speed of the proposed DSW distance with other sliced-based distances, namely the SW, Max-SW, Max-GSW-NN [25] and projected robust subspace Wasserstein (PRW) [35, 29] using the minimum expected distance estimator (MEDE) [5] on MNIST [28], CIFAR10 [27], CelebA [32] and LSUN [44] datasets. The details of the MEDE framework are described in Appendix D. On MNIST dataset, we train generative models with different distances and then evaluate their performances by comparing Wasserstein-2 distances between 10000 random generated images and all images from the MNIST test set. Due to the very large size of other datasets, e.g., 3 million images in LSUN, it is expensive to compute empirical Wasserstein-2 distance as its complexity is of order $\mathcal{O}(k^2 \log k)$ where k is the number of support points. Therefore, after we train generative models, we use FID score [21] to evaluate the generative quality of these generators. The FID score is calculated from 10000 random generated images and all training samples using precomputed statistics in [21]. Finally, for λ_C in DSW (see Definition 2), it is chosen in the set $\{1, 10, 100, 1000\}$ such that its Wasserstein-2 (FID score) (between 10000 random generated images and all images from corresponding validation set) is the lowest among the four values. Detailed experiment settings

are in Appendix F.

4.1 Results on MNIST

Generative quality and computational speed: We report the performance of the learned generative models for MNIST in Figure 2(a). To plot this figure, we vary the number of projections $N \in \{1, 10, 10^2, 5 \times 10^2, 10^3, 5 \times 10^3, 10^4\}$ for the SW, and $N \in \{1, 10, 10^2, 5 \times 10^2, 10^3, 5 \times 10^3\}$ for the DSW. Then we measure the computational time per minibatch and the Wasserstein-2 score of the learned generators for each N . We plot the Wasserstein-2 score and computational time of Max-SW and Max-GSW-NN in their standard settings [25]. Except for the regime with very fast but low-quality learned models, DSW is better than all the existing slice-based baselines in terms of both model quality and computational speed. Moreover, DSW can learn good models with very few projections, e.g., DSW-10 achieves better model quality than Max-GSW-NN and Max-SW and is one order-of-magnitude faster than these sliced distances. Finally, with a similar computational time, a learned generator by DSW has the Wasserstein-2 score that is roughly 10% lower than the one got from SW. For the qualitative comparison between these distances, we show random generated images from their generative models in Figure 7 in Appendix E.1. We observe that generated images from DSW are sharper and easier to classify into numbers than those from other baseline distances.

Comparison with projected robust subspace Wasserstein (PRW): In Figure 2(a), we plot the Wasserstein-2 score and computational time of PRW, where the subspace dimension of PRW varies in the range $\{2, 5, 10, 50\}$. PRW is able to improve upon the model quality of slice-based methods including DSW, however at the cost of being an order of magnitude slower than DSW with 10 projections (DSW-10). We observe that DSW-10 obtains a better Wasserstein-2 score than PRW with 5-dimensional subspace, while its corresponding computational time is 30 times faster than that of PRW-5. Using 50 dimension, PRW’s Wasserstein-2 score improves about 29% to that of DSW-10 but the computational cost is also around 40 times slower. The main computational advantage of DSW comes from the exact calculation of Wasserstein distance in one-dimension. The visual comparison between PRW and DSW based on their generated images is in Figure 12 in Appendix E.2.

Convergence behavior: Figure 2(b) shows that DSW learns better models at a faster speed of convergence than other baseline distances with a very small number of projections, e.g., DSW-10 is the second lowest curve compared to curves from other sliced-based distances.

Scalability over sample size of minibatch: Results in Figure 2(c) show that DSW has a computational complexity of the order $\mathcal{O}(k \log k)$, which is similar to those of other sliced-based distances, where k is the number of samples per batch.

Effect of the regularization parameter λ_C : For each value of $\lambda_C \in \{1, 10, 100, 1000\}$, we find the optimal distribution σ of DSW with $N = 10$ projections, and then calculate $A_N = \frac{1}{N^2} \sum_{i,j=1}^N |\theta_i^\top \theta_j|$, an approximation of the regularized term $\mathbb{E}_{\theta, \theta' \sim \sigma} [|\theta^\top \theta'|]$ in the dual form of DSW in equation (2), where $\{\theta_i\}_{i=1}^N \sim \sigma$. The results are shown in Figure 2(d). We observe that when λ_C increases, A_N goes down. When $\lambda_C = 0$, i.e., no regularization, A_N gets close to 1, meaning that all projected directions collapse to one direction. When $\lambda_C = 1000$, A_N is close to 0.1, suggesting that all projected directions are nearly orthogonal.

Additional experiments: We also investigate how the number of gradient-steps used for updating distribution of directions σ , and how the size of minibatches affects the quality of DSW (see

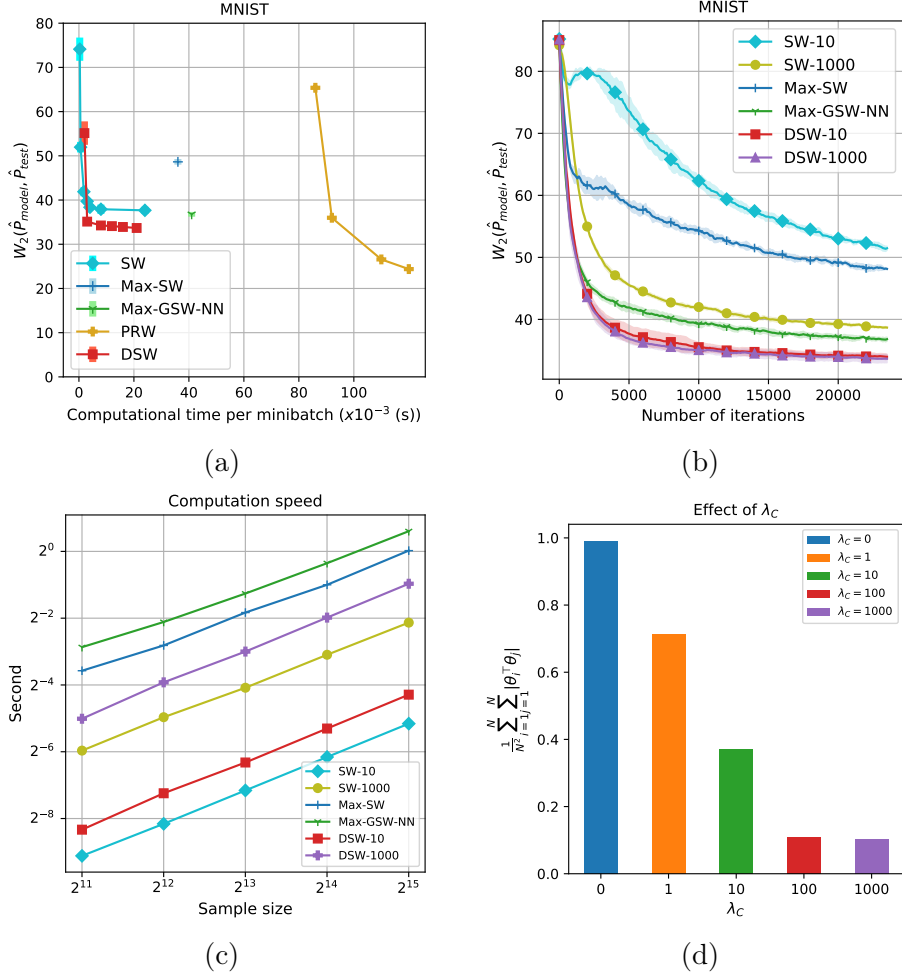


Figure 2: (a) Comparison between DSW, SW, Max-SW and Max-GSW-NN based on execution time and performance. Here, each dot of SW and DSW corresponds to the number of projections chosen in $\{1, 10, 10^2, 5 \times 10^2, 10^3, 5 \times 10^3, 10^4\}$. Each dot of PRW corresponds to the dimension of the subspace chosen in $\{2, 5, 10, 50\}$; (b) Comparison between SW, DSW, Max-SW and Max-GSW-NN based on Wasserstein-2 distance between distributions of learned model and test set over iterations; (c) Computation speed of distances based on the number of minibatch’s samples (log-log scale); (d) Effect of λ_C on the mean of absolute values of pairwise cosine similarity between 10 random directions from the found distribution σ of DSW.

Appendix E.1). The results show that an increasing number of gradient steps to update σ leads to better performance of DSW but also slows down the computation speed. Furthermore, we carry out experiments with DGSW, an extension of DSW to non-linear projections, and test the new proposed distances in training encoder-generator models on MNIST using joint contrastive inference (JCI) in Appendices E.1 and E.3. The description of these models is in Appendix D.

4.2 Results on Large-scale Datasets

Next, we conduct large-scale experiments on a range of more realistic image datasets. We train generative models using CIFAR10, CelebA, and LSUN datasets (all these datasets are rescaled

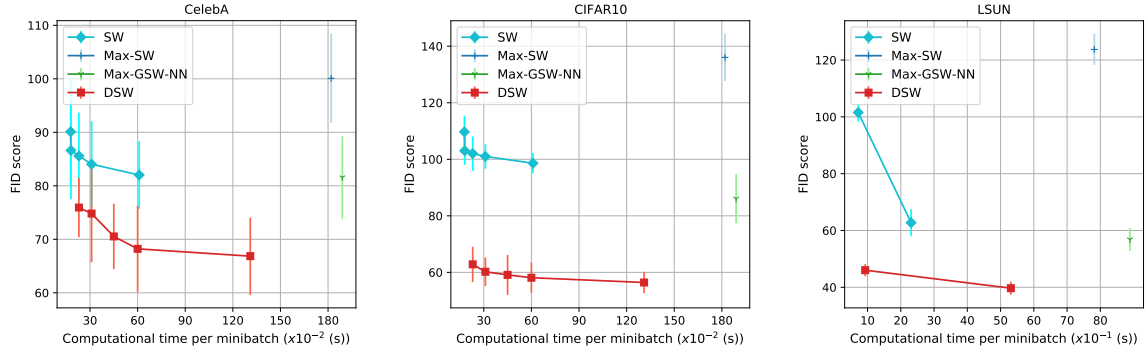


Figure 3: Comparison between DSW, SW, Max-SW and Max-GSW-NN in terms of execution time and performance. Here, each dot of SW and DSW corresponds to the number of projections chosen in $\{10^2, 5 \times 10^2, 10^3, 5 \times 10^3, 10^4\}$. We set the minibatch size be 512 on CelebA and CIFAR, and be 4096 on LSUN.

to 64x64 resolution). When working with high dimensional distributions, [16] proposed a trick to improve the quality of the generator by learning a feature function which maps data to a new feature space that is more manageable in size. When the feature function is fixed, the generator is trained to match the distribution of features. When the generator is fixed, the feature function tries to tease apart the data empirical features from the generated feature distribution. For the experiments in this section, we use the same technique with DSW and all other baseline distances.

We compare DSW with SW, Max-SW, and Max-GSW-NN in both generative quality (FID score) and computational time in Figure 3. We could not compare DSW with PRW on the large-scale datasets since PRW is computationally expensive to train to obtain good generated images. On CelebA and CIFAR10, we let N , the number of projections of both DSW and SW, vary in the set $\{10^2, 5 \times 10^2, 10^3, 5 \times 10^3, 10^4\}$. For LSUN, since it takes considerably longer time to train each model, we only vary N in the set $\{10^2, 10^4\}$. On all these large datasets, DSW outperforms all the other baselines in both FID score of the learned model and computational efficiency. The gap of FID scores between DSW and other methods is especially large on CIFAR10 and LSUN. For example, on CIFAR10, with the same computational time, FID scores of DSW are always lower than those of SW about 20 units. On LSUN, with 100 projections, DSW can achieve an FID score of 46 while SW with 10000 projections still has a worse FID score of over 60. It is interesting to note that on these high-dimensional datasets, Max-SW performs rather poorly: it obtains the highest FID scores among all distances while requires heavy computation. Max-GSW-NN has better FID scores than (Max)-SW; however, it is still worse than DSW and while being slower. This is consistent with the intuition that as the number of dimension of the data grows, the use of a single important slice in Max-SW becomes a less efficient approximation. DSW, on the other hand, is able to make use of more important slices, and at the same time avoids SW’s inefficiency of uniform slice-sampling.

Generated images from CelebA, CIFAR10 and LSUN are deferred to Appendix E.1. Comparing to other sliced-Wasserstein distances, generated samples obtained from the DSW’s generative model are also more visually realistic. Further experiments to compare DGSW with GSW, Max-GSW, and Max-GSW-NN are also given in the Appendix E.1. Based on these experiments, we can conclude that the distributional approach also improves the generative quality of non-linear slicing distances.

5 Conclusion

In this paper, we have presented the novel distributional sliced-Wasserstein (DSW) distances between two probability measures. Our main idea is to search for the best distribution of important directions while regularizing towards orthogonal directions. We prove that they are well-defined metrics and provide their theoretical and computational properties. We compare our proposed distances to other sliced-based distances in a variety of generative modeling tasks, including estimating generative models and jointly estimating both generators and inference models. Extensive experiments demonstrate that our new distances yield significantly better models and convergence behaviors during training than the previous sliced-based distances. One important future direction is to investigate theoretically the optimal choice of the regularization parameter λ_C such that the DSW distance can capture all the important directions that can distinguish two target probability measures well.

References

- [1] J. Altschuler, J. Niles-Weed, and P. Rigollet. Near-linear time approximation algorithms for optimal transport via sinkhorn iteration. In *Advances in neural information processing systems*, pages 1964–1974, 2017. (Cited on page 1.)
- [2] L. Ambrogioni, U. Güçlü, Y. Güçlütürk, M. Hinne, M. A. van Gerven, and E. Maris. Wasserstein variational inference. In *Advances in Neural Information Processing Systems*, pages 2473–2482, 2018. (Cited on page 22.)
- [3] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pages 214–223, 2017. (Cited on pages 1 and 22.)
- [4] E. Bayraktar and G. Guo. Strong equivalence between metrics of Wasserstein type. *arXiv preprint arXiv:1912.08247*, 2019. (Cited on pages 16 and 22.)
- [5] E. Bernton, P. E. Jacob, M. Gerber, and C. P. Robert. On parameter estimation with the Wasserstein distance. *Information and Inference: A Journal of the IMA*, 8(4):657–676, 2019. (Cited on pages 7 and 22.)
- [6] G. Beylkin. The inversion problem and applications of the generalized Radon transform. *Communications on pure and applied mathematics*, 37(5):579–599, 1984. (Cited on pages 2 and 3.)
- [7] N. Bonneel, J. Rabin, G. Peyré, and H. Pfister. Sliced and Radon Wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 1(51):22–45, 2015. (Cited on pages 1 and 4.)
- [8] N. Bonneotte. *Unidimensional and evolution methods for optimal transportation*. PhD thesis, Paris 11, 2013. (Cited on page 15.)
- [9] C. Bunne, D. Alvarez-Melis, A. Krause, and S. Jegelka. Learning generative models across incomparable spaces. In *International Conference on Machine Learning*, 2019. (Cited on page 1.)
- [10] X. Chen, Y. Yang, and Y. Li. Augmented sliced Wasserstein distances. *arXiv preprint arXiv:2006.08812*, 2020. (Cited on pages 23 and 24.)

- [11] L. Chizat, G. Peyré, B. Schmitzer, and F.-X. Vialard. Scaling algorithms for unbalanced optimal transport problems. *Mathematics of Computation*, 87(314):2563–2609, 2018. (Cited on page 1.)
- [12] N. Courty, R. Flamary, A. Habrard, and A. Rakotomamonjy. Joint distribution optimal transportation for domain adaptation. In *Advances in Neural Information Processing Systems*, pages 3730–3739, 2017. (Cited on page 1.)
- [13] N. Courty, R. Flamary, and D. Tuia. Domain adaptation with regularized optimal transport. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 274–289. Springer, 2014. (Cited on page 1.)
- [14] M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pages 2292–2300, 2013. (Cited on page 1.)
- [15] I. Deshpande, Y.-T. Hu, R. Sun, A. Pyrros, N. Siddiqui, S. Koyejo, Z. Zhao, D. Forsyth, and A. G. Schwing. Max-sliced Wasserstein distance and its use for GANs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10648–10656, 2019. (Cited on pages 2, 4, and 22.)
- [16] I. Deshpande, Z. Zhang, and A. G. Schwing. Generative modeling using the sliced Wasserstein distance. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3483–3491, 2018. (Cited on pages 2, 10, and 22.)
- [17] V. Dumoulin, I. Belghazi, B. Poole, O. Mastropietro, A. Lamb, M. Arjovsky, and A. Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016. (Cited on page 22.)
- [18] A. Genevay, G. Peyre, and M. Cuturi. Learning generative models with Sinkhorn divergences. In *International Conference on Artificial Intelligence and Statistics*, pages 1608–1617, 2018. (Cited on pages 1 and 22.)
- [19] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville. Improved training of Wasserstein Gans. In *Advances in neural information processing systems*, pages 5767–5777, 2017. (Cited on page 1.)
- [20] S. Helgason. *Integral geometry and Radon transforms*. Springer Science & Business Media, 2010. (Cited on pages 1 and 3.)
- [21] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, pages 6626–6637, 2017. (Cited on page 7.)
- [22] D. P. Kingma and M. Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013. (Cited on page 22.)
- [23] P. A. Knight. The Sinkhorn–Knopp algorithm: convergence and applications. *SIAM Journal on Matrix Analysis and Applications*, 30(1):261–275, 2008. (Cited on page 1.)
- [24] M. Kochurov, R. Karimov, and S. Kozlukov. Geopt: Riemannian optimization in pytorch, 2020. (Cited on page 28.)

- [25] S. Kolouri, K. Nadjahi, U. Simsekli, R. Badeau, and G. Rohde. Generalized sliced Wasserstein distances. In *Advances in Neural Information Processing Systems*, pages 261–272, 2019. (Cited on pages 2, 4, 7, 8, and 26.)
- [26] S. Kolouri, P. E. Pope, C. E. Martin, and G. K. Rohde. Sliced Wasserstein auto-encoders. In *International Conference on Learning Representations*, 2018. (Cited on page 2.)
- [27] A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009. (Cited on page 7.)
- [28] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. (Cited on page 7.)
- [29] T. Lin, C. Fan, N. Ho, M. Cuturi, and M. I. Jordan. Projection robust Wasserstein distance and Riemannian optimization. *arXiv preprint arXiv:2006.07458*, 2020. (Cited on pages 7, 23, and 27.)
- [30] T. Lin, N. Ho, and M. Jordan. On efficient optimal transport: An analysis of greedy and accelerated mirror descent algorithms. In *International Conference on Machine Learning*, pages 3982–3991, 2019. (Cited on page 1.)
- [31] T. Lin, N. Ho, and M. I. Jordan. On the acceleration of the Sinkhorn and Greenkhorn algorithms for optimal transport. *arXiv preprint arXiv:1906.01437*, 2019. (Cited on page 1.)
- [32] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015. (Cited on page 7.)
- [33] A. Liutkus, U. Simsekli, S. Majewski, A. Durmus, and F.-R. Stöter. Sliced-Wasserstein flows: Nonparametric generative modeling via optimal transport and diffusions. In *International Conference on Machine Learning*, pages 4104–4113, 2019. (Cited on page 2.)
- [34] K. Nadjahi, A. Durmus, U. Simsekli, and R. Badeau. Asymptotic guarantees for learning generative models with the sliced-Wasserstein distance. In *Advances in Neural Information Processing Systems*, pages 250–260, 2019. (Cited on page 22.)
- [35] F.-P. Paty and M. Cuturi. Subspace robust Wasserstein distances. In *International Conference on Machine Learning*, pages 5072–5081, 2019. (Cited on pages 2, 7, 23, and 27.)
- [36] G. Peyré and M. Cuturi. Computational optimal transport. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019. (Cited on pages 1 and 3.)
- [37] M. Rowland, J. Hron, Y. Tang, K. Choromanski, T. Sarlos, and A. Weller. Orthogonal estimation of Wasserstein distances. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 186–195, 2019. (Cited on page 2.)
- [38] B. Schmitzer. Stabilized sparse scaling algorithms for entropy regularized transport problems. *SIAM Journal on Scientific Computing*, 41(3):A1443–A1481, 2019. (Cited on page 1.)
- [39] R. Sinkhorn. Diagonal equivalence to matrices with prescribed row and column sums. *The American Mathematical Monthly*, 74(4):402–405, 1967. (Cited on page 1.)

- [40] I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schoelkopf. Wasserstein auto-encoders. In *International Conference on Learning Representations*, 2018. (Cited on pages 1 and 22.)
- [41] C. Villani. *Optimal transport: Old and New*. Springer, 2008. (Cited on page 3.)
- [42] M. J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, 2019. (Cited on page 20.)
- [43] J. Wu, Z. Huang, D. Acharya, W. Li, J. Thoma, D. P. Paudel, and L. V. Gool. Sliced Wasserstein generative models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3713–3722, 2019. (Cited on page 2.)
- [44] F. Yu, Y. Zhang, S. Song, A. Seff, and J. Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. (Cited on page 7.)

Supplement to “Distributional Sliced-Wasserstein and Applications to Generative Modeling”

In this supplementary material, we collect several proofs and remaining materials that were deferred from the main paper. In Appendix A, we provide the proofs of the main results in the paper. In Appendix B, additional properties of distributional sliced-Wasserstein (DSW) distance are provided. In Appendix C, we discuss distributional generalized sliced-Wasserstein distance (DGSW) and its dual form and properties. We describe in detail the applications of DSW and DGSW to generative modelings in Appendix D. Furthermore, we provide additional experiments and experiment settings in Appendices E and F.

A Proofs

In this appendix, we collect the proofs for all the results in the main text.

A.1 Proof of Theorem 1

We first show that the distributional sliced-Wasserstein distance satisfies the triangle inequality property for any three probability measures μ_1, μ_2 , and μ_3 . In fact, from the definition of distributional sliced-Wasserstein distance for admissible $C > 0$, for any $\epsilon > 0$ we find that

$$\begin{aligned}
 \text{DSW}_p(\mu_1, \mu_2; C) &\stackrel{(i)}{\leq} \left\{ \mathbb{E}_{\theta \sim \sigma_\epsilon^*} W_p^p(\mathcal{R}I_{\mu_1}(\cdot, \theta), \mathcal{R}I_{\mu_2}(\cdot, \theta)) \right\}^{\frac{1}{p}} + \epsilon \\
 &\stackrel{(ii)}{\leq} \left\{ \mathbb{E}_{\theta \sim \sigma_\epsilon^*} [W_p(\mathcal{R}I_{\mu_1}(\cdot, \theta), \mathcal{R}I_{\mu_3}(\cdot, \theta)) + W_p^p(\mathcal{R}I_{\mu_3}(\cdot, \theta), \mathcal{R}I_{\mu_2}(\cdot, \theta))] \right\}^{\frac{1}{p}} + \epsilon \\
 &\stackrel{(iii)}{\leq} \left\{ \mathbb{E}_{\theta \sim \sigma_\epsilon^*} W_p^p(\mathcal{R}I_{\mu_1}(\cdot, \theta), \mathcal{R}I_{\mu_3}(\cdot, \theta)) \right\}^{\frac{1}{p}} \\
 &\quad + \left\{ \mathbb{E}_{\theta \sim \sigma_\epsilon^*} W_p^p(\mathcal{R}I_{\mu_3}(\cdot, \theta), \mathcal{R}I_{\mu_2}(\cdot, \theta)) \right\}^{\frac{1}{p}} + \epsilon \\
 &\leq \sup_{\sigma \in \mathbb{M}_C} \left\{ \mathbb{E}_{\theta \sim \sigma} [W_p^p(\mathcal{R}I_{\mu_1}(\cdot, \theta), \mathcal{R}I_{\mu_3}(\cdot, \theta))] \right\}^{\frac{1}{p}} \\
 &\quad + \sup_{\sigma \in \mathbb{M}_C} \left\{ \mathbb{E}_{\theta \sim \sigma} [W_p^p(\mathcal{R}I_{\mu_3}(\cdot, \theta), \mathcal{R}I_{\mu_2}(\cdot, \theta))] \right\}^{\frac{1}{p}} + \epsilon \\
 &= \text{DSW}_p(\mu_1, \mu_3; C) + \text{DSW}_p(\mu_2, \mu_3; C) + \epsilon,
 \end{aligned}$$

where the existence of σ_ϵ^* in (i) is from the definition of supremum; inequality in (ii) is due to the triangle inequality with Wasserstein distance of order p ; inequality in (iii) follows from the application of the Minkowski inequality. By letting $\epsilon \rightarrow 0$ in the above inequality, we obtain the conclusion with the triangle inequality of distributional sliced-Wasserstein distance.

The non-negativity and symmetry of distributional sliced-Wasserstein distance follow directly from the non-negativity and symmetry of Wasserstein distance. For the identity property, it is straight-forward that if $\mu_1 \equiv \mu_2$ then $\text{DSW}_p(\mu_1, \mu_2) = 0$. On the other hand, if $\text{DSW}_p(\mu_1, \mu_2) = 0$, an application of Fourier transform as that in [8] leads to $\mu_1 \equiv \mu_2$.

As a consequence, for any $p \geq 1$ and admissible $C > 0$, $\text{DSW}_p(\cdot, \cdot; C)$ is a well-defined metric in the space of Borel probability measures with finite p -th moment.

A.2 Proof of Theorem 2

(a) From the definition of distributional sliced-Wasserstein distance, for any $p \geq 1$ and admissible $C > 0$ we find that

$$\text{DSW}_p(\mu, \nu; C) \leq \sup_{\sigma \in \mathbb{M}} \left(\mathbb{E}_{\theta \sim \sigma} \left[W_p^p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta)) \right] \right)^{\frac{1}{p}} = \text{maxSW}_p(\mu, \nu),$$

where \mathbb{M} is the space of all probability measures. The inequality is due to the fact that $\mathbb{M}_C \subseteq \mathbb{M}$ for all admissible $C > 0$. The second equality is true because $W_p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta)) \leq \text{maxSW}_p(\mu, \nu)$ for all $\theta \in \mathbb{S}^{d-1}$, which leads to $\mathbb{E}_{\theta \sim \sigma} [W_p^p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta))] \leq \text{maxSW}_p^p(\mu, \nu)$. The inequality becomes equality when σ is the Dirac measure at θ^* that maximizes the value of $W_p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta))$.

Furthermore, we have

$$\begin{aligned} W_p^p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta)) &= \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} |x^\top \theta - y^\top \theta|^p d\pi(x, y) \\ &\leq \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} \|x - y\|^p d\pi(x, y) = W_p^p(\mu, \nu), \end{aligned}$$

where the last inequality is due to the fact that length of side of right triangle $|(x - y)^\top \theta|$ is less than length of its hypotenuse $\|x - y\|$ for all $\theta \in \mathbb{S}^{d-1}$. Therefore, $\text{maxSW}_p(\mu, \nu) \leq W_p(\mu, \nu)$ for any $p \geq 1$.

Putting the above results together, we obtain the conclusion of part (a) of the theorem.

(b) Denote $\bar{\sigma} = \sum_{i=1}^d \frac{1}{d} \delta_{\theta_i}$ where $\theta_1 = \theta^*$, which maximizes the value of $W_p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta))$ and $\theta_1, \dots, \theta_d$ form an orthonormal basis in \mathbb{R}^d . Simple algebra shows that

$$\mathbb{E}_{\theta, \theta' \sim \bar{\sigma}} [|\theta^\top \theta'|] = \sum_{1 \leq i, j \leq d} \left(\frac{1}{d}\right)^2 |\theta_i^\top \theta_j| = \frac{1}{d}.$$

Since $C \geq \frac{1}{d}$, the above equation indicates that $\bar{\sigma} \in \mathbb{M}_C$. Therefore, we find that

$$\begin{aligned} \text{DSW}_p(\mu, \nu; C) &\geq \left(\mathbb{E}_{\theta \sim \bar{\sigma}} \left[W_p^p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta)) \right] \right)^{\frac{1}{p}} \\ &= \left(\sum_{i=1}^d \frac{1}{d} W_p^p(\mathcal{R}I_\mu(\cdot, \theta_i), \mathcal{R}I_\nu(\cdot, \theta_i)) \right)^{\frac{1}{p}} \\ &\geq \left(\frac{1}{d}\right)^{\frac{1}{p}} W_p(\mathcal{R}I_\mu(\cdot, \theta_1), \mathcal{R}I_\nu(\cdot, \theta_1)) = \left(\frac{1}{d}\right)^{\frac{1}{p}} \text{maxSW}_p(\mu, \nu). \end{aligned}$$

Moreover, for any $p \geq 1$, $\text{SW}_p(\mu, \nu) \leq \text{maxSW}_p(\mu, \nu)$. Collecting the previous results, we reach the conclusion of part (b).

Equivalence of $\text{DSW}_p(\cdot, \cdot; C)$ to other distances: Based on the result of Theorem 2.1 in [4], maxSW_p , SW_p , and W_p are equivalent distances for any $p \geq 1$. In particular, for any sequence $(\mu_n)_{n \geq 1} \in \mathcal{P}_p(\mathbb{R}^d)$ and $\mu \in \mathcal{P}_p(\mathbb{R}^d)$, the following holds

$$\lim_{n \rightarrow \infty} \text{maxSW}_p(\mu_n, \mu) = 0 \iff \lim_{n \rightarrow \infty} \text{SW}_p(\mu_n, \mu) = 0 \iff \lim_{n \rightarrow \infty} W_p(\mu_n, \mu) = 0. \quad (3)$$

Now, if we have $\lim_{n \rightarrow \infty} \max \text{SW}_p(\mu_n, \mu) = 0$ for $p \geq 1$, the result of part (a) shows that $\lim_{n \rightarrow \infty} \text{DSW}_p(\mu_n, \mu; C) = 0$. On the other hand, when $\lim_{n \rightarrow \infty} \text{DSW}_p(\mu_n, \mu; C) = 0$, as long as $C \geq \frac{1}{d}$ and $p \geq 1$, the result of part (b) leads to $\lim_{n \rightarrow \infty} \max \text{SW}_p(\mu_n, \mu) = 0$. As a consequence, when $C \geq \frac{1}{d}$ and $p \geq 1$ we have

$$\lim_{n \rightarrow \infty} \text{DSW}_p(\mu_n, \mu; C) = 0 \iff \lim_{n \rightarrow \infty} \max \text{SW}_p(\mu_n, \mu) = 0. \quad (4)$$

Combining the results in equations (3) and (4), we reach the conclusion that when $C \geq \frac{1}{d}$ and $p \geq 1$, $\text{DSW}_p(\cdot, \cdot; C)$, $\max \text{SW}_p$, SW_p , and W_p are equivalent distances.

B Additional Studies with Distributional Sliced-Wasserstein Distance

In this appendix, we provide further studies with distributional sliced-Wasserstein distance.

B.1 Discussion of the constraint in DSW

We first compute $\mathbb{E}_{\theta, \theta' \sim \sigma^{d-1}} [|\theta^\top \theta'|]$ where σ^{d-1} is the uniform distribution on the unit sphere \mathbb{S}^{d-1} .

Theorem 3. *For uniform measure σ^{d-1} on the unit sphere \mathbb{S}^{d-1} , we have*

$$\int_{\theta, \theta' \sim \sigma^{d-1}} |\theta^\top \theta'| d\sigma^{d-1}(\theta) d\sigma^{d-1}(\theta') = \frac{\Gamma(\frac{d}{2})}{\pi^{\frac{1}{2}} \Gamma(\frac{d+1}{2})},$$

where $\Gamma(\cdot)$ is the Gamma function.

Remark. *The result of Theorem 3 indicates that as long as $C \geq \frac{\Gamma(\frac{d}{2})}{\pi^{\frac{1}{2}} \Gamma(\frac{d+1}{2})}$, we have $\sigma^{d-1} \in \mathbb{M}_C$. Furthermore, by Gautschi's inequality for the Gamma function, we find that*

$$\frac{1}{\pi^{\frac{1}{2}} (\frac{d+1}{2})^{\frac{1}{2}}} < \frac{\Gamma(\frac{d}{2})}{\pi^{\frac{1}{2}} \Gamma(\frac{d+1}{2})} < \frac{1}{\pi^{\frac{1}{2}} (\frac{d-1}{2})^{\frac{1}{2}}}$$

For $d \geq 3$, we have $2d^2/(d+1) > \pi$. Hence, we obtain that

$$\frac{1}{\pi^{\frac{1}{2}} (\frac{d+1}{2})^{\frac{1}{2}}} > \frac{1}{d}.$$

Given the above bound, when $C \geq \frac{\Gamma(\frac{d}{2})}{\pi^{\frac{1}{2}} \Gamma(\frac{d+1}{2})}$, the set \mathbb{M}_C contains both σ^{d-1} and $\bar{\sigma} = \sum_{i=1}^d \frac{1}{d} \delta_{\theta_i}$ where $\theta_1, \dots, \theta_d$ form any orthonormal basis in \mathbb{R}^d . Furthermore, for $d = 2$, we have

$$\frac{\Gamma(1)}{\pi^{\frac{1}{2}} \Gamma(\frac{2+1}{2})} = \frac{2}{\pi} > \frac{1}{2}.$$

Therefore, when $d = 2$ and $C \geq \frac{\Gamma(\frac{d}{2})}{\pi^{\frac{1}{2}} \Gamma(\frac{d+1}{2})}$, the set \mathbb{M}_C also contains $\bar{\sigma} = \sum_{i=1}^d \frac{1}{d} \delta_{\theta_i}$.

Proof. Since σ^{d-1} is the uniform measure on the unit sphere \mathbb{S}^{d-1} , the integral

$$\int_{\theta \sim \sigma^{d-1}} |\theta^\top \theta'| d\sigma^{d-1}(\theta)$$

is the same for all fixed θ' . Hence for any fixed $\theta^* \in \mathbb{S}^{d-1}$, we obtain

$$I = \int_{\theta, \theta' \sim \sigma^{d-1}} |\theta^\top \theta'| d\sigma^{d-1}(\theta) d\sigma^{d-1}(\theta') = \int_{\theta \sim \sigma^{d-1}} |\theta^\top \theta^*| d\sigma^{d-1}(\theta).$$

Without loss of generality, we choose $\theta^* = (1, 0, \dots, 0)$, I is equal to

$$\int_{\theta \sim \sigma^{d-1}} |\theta^{(1)}| d\sigma^{d-1}(\theta),$$

where $\theta = (\theta^{(1)}, \dots, \theta^{(d)})$. For any measurable subset S of \mathbb{S}^{d-1} , let $A(S)$ be the area of S on the surface of \mathbb{S}^{d-1} and $A(\mathbb{S}^{d-1})$ be the area of the surface of \mathbb{S}^{d-1} which is equal to

$$A(\mathbb{S}^{d-1}) = \frac{d\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2} + 1)}.$$

Now, we have

$$\int_{\theta \sim \sigma^{d-1}} |\theta^{(1)}| d\sigma^{d-1}(\theta) = \frac{1}{A(\mathbb{S}^{d-1})} \int_{\theta \in \mathbb{S}^{d-1}} |\theta^{(1)}| dA(\mathbb{S}^{d-1}(\theta)).$$

Let H_1 be the hyperplane formed by $\theta^{(2)}, \dots, \theta^{(d)}$ and H_θ be the hyperplane tangent to the sphere \mathbb{S}^{d-1} at θ . Then θ is the normal vector to H_θ and $\theta^* = (1, 0, \dots, 0)$ is orthogonal to H_1 . Let α be the angle between θ and θ^* . Then

$$\begin{aligned} dA(\mathbb{S}^{d-1}(\theta)) \cos(H_1, H_\theta) &= d\theta^{(2)} \dots d\theta^{(d)} \\ dA(\mathbb{S}^{d-1}(\theta)) &= \frac{1}{\cos(\theta, \theta^*)} d\theta^{(2)} \dots d\theta^{(d)} = \frac{1}{|\theta^{(1)}|} d\theta^{(2)} \dots d\theta^{(d)}. \end{aligned}$$

Return to the integral, we find that

$$\begin{aligned} I &= \int_{\theta \sim \sigma^{d-1}} |\theta^{(1)}| d\sigma^{d-1}(\theta) = \frac{1}{A(\mathbb{S}^{d-1})} \int_{\sum_{i=1}^d (\theta^{(i)})^2 = 1} |\theta^{(1)}| \frac{1}{|\theta^{(1)}|} d\theta^{(2)} \dots d\theta^{(d)} \\ &= \frac{2}{A(\mathbb{S}^{d-1})} \int_{\theta^{(1)} > 0, \sum_{i=2}^d (\theta^{(i)})^2 \leq 1} d\theta^{(2)} \dots d\theta^{(d)} \\ &= \frac{2}{A(\mathbb{S}^{d-1})} V(\mathbb{B}^{d-1}) \\ &= \frac{2\Gamma(\frac{d}{2})}{2\pi^{\frac{d}{2}}} \times \frac{\pi^{\frac{d-1}{2}}}{\Gamma(\frac{d-1}{2} + 1)} \\ &= \frac{\Gamma(\frac{d}{2})}{\pi^{\frac{1}{2}} \Gamma(\frac{d+1}{2})}, \end{aligned}$$

where \mathbb{B}^{d-1} is the unit ball in the $d-1$ dimensional space and $V(\mathbb{B}^{d-1})$ is its corresponding volume. As a consequence, we obtain the conclusion of the theorem. \square

B.2 Approximation of dual value of DSW

Now, we give a detailed form of the objective function $\text{DS}(f_\phi)$ in the dual form of DSW in equation (2). In particular, simple calculation shows that

$$\begin{aligned} \nabla_\phi \text{DS}(f_\phi) &= \frac{1}{p} \left\{ \mathbb{E}_{\theta \sim \sigma^{d-1}} [W_p^p(\mathcal{R}I_\mu(\cdot, f_\phi(\theta)), \mathcal{R}I_\nu(\cdot, f_\phi(\theta)))] \right\}^{\frac{1}{p}-1} \\ &\quad \times \mathbb{E}_{\theta \sim \sigma^{d-1}} [\nabla_\phi W_p^p(\mathcal{R}I_\mu(\cdot, f_\phi(\theta)), \mathcal{R}I_\nu(\cdot, f_\phi(\theta)))] - \lambda_C \mathbb{E}_{\theta, \theta' \sim \sigma^{d-1}} [\nabla_\phi |f_\phi(\theta)^\top f_\phi(\theta')|]. \end{aligned} \quad (5)$$

Since the outer expectations in equation (5) are intractable to compute, we employ the standard Monte Carlo scheme to approximate these expectations. Therefore, we obtain the following approximation:

$$\begin{aligned} \nabla_\phi \text{DS}(f_\phi) &\approx \frac{1}{p} \left\{ \frac{1}{n} \sum_{i=1}^n [W_p^p(\mathcal{R}I_\mu(\cdot, f_\phi(\theta_i)), \mathcal{R}I_\nu(\cdot, f_\phi(\theta_i)))] \right\}^{\frac{1}{p}-1} \\ &\quad \times \left\{ \frac{1}{n} \sum_{i=1}^n [\nabla_\phi W_p^p(\mathcal{R}I_\mu(\cdot, f_\phi(\theta_i)), \mathcal{R}I_\nu(\cdot, f_\phi(\theta_i)))] \right\} - \frac{\lambda_C}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \nabla_\phi |(f_\phi(\theta_i))^\top f_\phi(\theta_j)|, \end{aligned}$$

where $\theta_1, \dots, \theta_n$ are i.i.d. samples from the unit sphere \mathbb{S}^{d-1} .

Denote ϕ^* as the fixed point of the stochastic gradient ascent algorithm. Then, we can use f_{ϕ^*} as the local maxima of the optimization problem (2). By using Monte Carlo method to approximate the expectation in equation (2), we obtain the following approximation:

$$\begin{aligned} \text{DSW}_p^*(\mu, \nu; C) &\approx \left\{ \frac{1}{n} \sum_{i=1}^n [W_p^p(\mathcal{R}I_\mu(\cdot, f_{\phi^*}(\theta_i)), \mathcal{R}I_\nu(\cdot, f_{\phi^*}(\theta_i)))] \right\}^{1/p} \\ &\quad - \frac{\lambda_C}{n(n-1)} \sum_{1 \leq i \neq j \leq n} |(f_{\phi^*}(\theta_i))^\top f_{\phi^*}(\theta_j)| + \lambda_C C. \end{aligned}$$

B.3 Statistical guarantee of DSW

In this appendix, we provide the statistical guarantee of DSW.

Theorem 4. *Given probability measure P supported on a compact subset $\Theta \subset \mathbb{R}^d$. Assume that X_1, \dots, X_n are i.i.d. data from P . Denote $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ the empirical measure of the data points X_1, \dots, X_n . Then, for any admissible regularizing constant $C > 0$ and for any $p \geq 1$, we obtain that*

$$\mathbb{E} \left[\text{DSW}_p(P_n, P; C) \right] \leq c \sqrt{\frac{d \log n}{n}},$$

where $c > 0$ is some universal constant.

Remark. *The result of Theorem 4 demonstrates that DSW has similar statistical guarantees as other sliced distances and does not suffer from the curse of dimensionality. Therefore, it is an appealing distance for applications in generative modeling.*

Proof. The proof of Theorem 4 is a direct application of Theorem 2 and statistical guarantee of max-sliced Wasserstein distance. Here, we provide the proof for the completeness. In particular, based on the result of Theorem 2, we obtain that

$$\mathbb{E} \left[\text{DSW}_p(P_n, P; C) \right] \leq \mathbb{E} \left[\max \text{SW}_p(P_n, P) \right].$$

Therefore, it is sufficient to demonstrate that $\mathbb{E} \left[\max \text{SW}_p(P_n, P) \right] \leq c \sqrt{\frac{d \log n}{n}}$ for some universal constant $c > 0$. In order to simplify the presentation, we denote a few notation. First, we define \mathcal{H} the set of half-spaces $H_{\theta, x} = \{y \in \mathbb{R}^d : \langle y, \theta \rangle \leq x\}$ for any $\theta \in \mathbb{S}^{d-1}$ and $x \in \mathbb{R}$. Then, it has been shown that \mathcal{H} has at most $d + 1$ Vapnik–Chervonenkis (VC) dimension [42]. The VC inequality implies that

$$\sup_{H \in \mathcal{H}} |P_n(H) - P(H)| \leq \sqrt{\frac{32}{n} [(d+1) \log(n+1) + \log(8/\delta)]} =: c_{n,\delta}$$

with probability at least $1 - \delta$, for any $\delta \in (0, 1)$. On the other hand, we have

$$\sup_{H \in \mathcal{H}} |P_n(H) - P(H)| = \sup_{x \in \mathbb{R}, \theta \in \mathbb{S}^{d-1}} |F_{n,\theta}(x) - F_\theta(x)|,$$

where $F_{n,\theta}$ and F_θ are respectively the cumulative distribution functions (CDF) of $\mathcal{R}I_{P_n}(\cdot, \theta)$ and $\mathcal{R}I_P(\cdot, \theta)$. Given the above equation and the close-form of Wasserstein distance in one dimension, we find that

$$\begin{aligned} \max \text{SW}_p^p(P_n, P) &= \max_{\theta \in \mathbb{S}^{d-1}} \int_0^1 |F_{n,\theta}^{-1}(u) - F_\theta^{-1}(u)|^p du \\ &= \max_{\theta \in \mathbb{S}^{d-1}} \int_{\mathbb{R}} |F_{n,\theta}(x) - F_\theta(x)|^p dx \\ &\leq \text{diam}(\Theta) \sup_{x \in \mathbb{R}, \theta \in \mathbb{S}^{d-1}} |F_{n,\theta}(x) - F_\theta(x)|^p \leq \text{diam}(\Theta) c_{n,\delta}^p. \end{aligned}$$

By using the above inequality, we obtain that $\mathbb{E} \left[\max \text{SW}_p(P_n, P) \right] \leq c \sqrt{\frac{d \log n}{n}}$ for some universal constant $c > 0$. As a consequence, we reach the conclusion of Theorem 4. \square

C An extension to distributional generalized sliced-Wasserstein distance

We now consider an extension of DSW to non-linear projections via generalized Radon transform. The constant $C > 0$ is *generalized admissible* if the set $\bar{\mathbb{M}}_C$ of probability measures σ on the compact set of feasible parameters Ω_θ satisfying $\mathbb{E}_{\theta, \theta' \sim \sigma} [|\cos(\theta, \theta')|] \leq C$ is not empty.

Definition 3. *Given two probability measures μ and ν on \mathbb{R}^d with finite p -th moments where $p \geq 1$ and a generalized admissible regularizing constant $C > 0$. The distributional generalized sliced-Wasserstein distance (DGSW) of order p between μ and ν is defined as follows:*

$$DGSW_p(\mu, \nu; C) := \sup_{\sigma \in \bar{\mathbb{M}}_C} \left\{ \mathbb{E}_{\theta \sim \sigma} W_p^p(\mathcal{G}I_\mu(\cdot, \theta), \mathcal{G}I_\nu(\cdot, \theta)) \right\}^{1/p},$$

where \mathcal{G} is generalized Radon transform defined in Section 2.2.

The DGSW distance uses the advantage of non-linear projections to capture more complex structures of the target probability measures. We show that as long as the generalized Radon transform is injective, DGSW is a proper metric in the probability space.

Theorem 5. For any $p \geq 1$ and generalized admissible $C > 0$, as long as the generalized Radon transform is injective, the distributional generalized sliced-Wasserstein is a well-defined metric in the space of Borel probability measures with finite p -th moment.

The proof of Theorem 5 simply follows the proof argument of Theorem 1 under the injectivity of GRT; thus, it is omitted. In order to compute DGSW, we also utilize the dual form of DGSW as that of DSW distance.

Dual form of distributional generalized sliced-Wasserstein distance: Similar to the distributional sliced-Wasserstein distance, we use the dual form of distributional generalized sliced-Wasserstein distance to approximate the value of distributional generalized sliced-Wasserstein distance. Recall that, for any $\theta, \theta' \in \mathbb{R}^d$, $\cos(\theta, \theta') = \frac{\theta^\top \theta'}{\|\theta\| \|\theta'\|}$.

Definition 4. For any $p \geq 1$ and generalized admissible $C > 0$, there exists a non-negative constant λ_C depending on C such that the dual form of DGSW distance takes the following form

$$\begin{aligned} \text{DGSW}_p^*(\mu, \nu; C) &:= -\sup_{\lambda \geq 0} \inf_{\sigma \in \bar{\mathbb{M}}} \left\{ -\left(\mathbb{E}_{\theta \sim \sigma} \left[W_p^p(\mathcal{G}I_\mu(\cdot, \theta), \mathcal{G}I_\nu(\cdot, \theta)) \right] \right)^{1/p} \right. \\ &\quad \left. + \lambda \left(\mathbb{E}_{\theta, \theta' \sim \sigma} \left[\frac{|\theta^\top \theta'|}{\|\theta\| \|\theta'\|} \right] - C \right) \right\} \\ &= \sup_{\sigma \in \bar{\mathbb{M}}} \left\{ \left(\mathbb{E}_{\theta \sim \sigma} \left[W_p^p(\mathcal{G}I_\mu(\cdot, \theta), \mathcal{G}I_\nu(\cdot, \theta)) \right] \right)^{1/p} - \lambda_C \mathbb{E}_{\theta, \theta' \sim \sigma} \left[\frac{|\theta^\top \theta'|}{\|\theta\| \|\theta'\|} \right] \right\} \\ &\quad + \lambda_C C, \end{aligned}$$

where $\bar{\mathbb{M}}$ denotes the space of all probability measures on the compact set of feasible parameter Ω_θ .

From the duality theory, we obtain that $\text{DGSW}_p(\mu, \nu; C) \geq \text{DGSW}_p^*(\mu, \nu; C)$ for any $p \geq 1$ and admissible $C > 0$. Similar to DSW distance, the dual form of DGSW provides an efficient way to approximate the DGSW distance. We show that when the compact set of feasible parameter $\Omega_\theta = \mathbb{S}^{d-1}$, similar reparametrization trick like that of the dual form of DSW distance can be applied to the dual form of DGSW distance. In particular, when $\Omega_\theta = \mathbb{S}^{d-1}$, we obtain the equivalent dual form of DGSW as follows:

$$\begin{aligned} \text{DGSW}_p^*(\mu, \nu; C) &= \sup_{f \in \mathcal{F}} \left\{ \left(\mathbb{E}_{\theta \sim \sigma^{d-1}} \left[W_p^p(\mathcal{G}I_\mu(\cdot, f(\theta)), \mathcal{G}I_\nu(\cdot, f(\theta))) \right] \right)^{1/p} \right. \\ &\quad \left. - \lambda_C \mathbb{E}_{\theta, \theta' \sim \sigma^{d-1}} \left[|f(\theta)^\top f(\theta')| \right] \right\} + \lambda_C C, \end{aligned} \tag{6}$$

where \mathcal{F} is a class of Borel measurable functions from \mathbb{S}^{d-1} to \mathbb{S}^{d-1} and $\lambda_C > 0$ is some positive constant given in Definition 4. Then, in order to find an optimal f , we can parameterize f as f_ϕ , which we can think as (deep) neural network. From here, with similar argument as that of equation (5), we can approximate the gradient of the objective function in equation (6) with respect to ϕ and then use stochastic gradient ascent algorithm to update ϕ . Finally, we can use the fixed point of the algorithm to approximate the dual value of DGSW in equation (6).

D Applications to Generative Modeling

The DSW and DGSW distances can potentially be applied in settings where there is a benefit of employing an optimal-transport type of distance in a computationally efficient manner. In this section, we discuss two general settings where the DSW and DGSW distances can be immediately applied. The first setting is a standard generative modeling task using the minimum expected distance estimator framework [5] where a generative model is fitted to a data distribution by minimizing an appropriate divergence. The second setting is a joint contrastive inference task where both a generative model and inference model are learned jointly, again by minimizing some divergence in the joint space of observed variable and latent variable. In each setting, we apply the DSW to these tasks as well as its generalized version, the DGSW.

D.1 Minimum expected distributional sliced-Wasserstein estimator

Minimum expected distance estimators [5] are widely used recently due to its efficiency in learning implicit generative models. Popular estimators include those based on OT distances [3, 18, 40] due to their smooth and differentiable objectives especially when the supports of the data and the generative distributions are not the same. In sliced-Wasserstein cases, SW and Max-SW have been employed with rigorous theoretical analyses in various works [4, 15, 16, 34]. They enjoy the benefits of the Wasserstein distance in one dimension and obtain fast speed in training the model. In this paper, we introduce a new novel estimator by replacing SW and Max-SW by our new DSW distance, which we refer to as *minimum expected distributional sliced-Wasserstein estimator*. The new estimator is defined as follows:

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} \mathbb{E}[\text{DSW}_p(\hat{\mu}_n, \hat{\mu}_{\theta, m}) | X_{1:n}], \quad (7)$$

where Φ is the parameter space, $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ is the empirical measure, and $\hat{\mu}_{\theta, m} = \frac{1}{m} \sum_{i=1}^m \delta_{Y_i}$ denotes the empirical distribution that is obtained by sampling i.i.d samples from μ_θ . In practice, μ_θ is created by pushing a simple distribution ϵ (such as the standard Gaussian) through a neural net, parameterized by θ , i.e., $\mu_\theta = T_\theta \# \epsilon$.

D.2 Distributional sliced-Wasserstein joint contrastive inference

Learning both a generator and an inference model, i.e., an encoder, is a central task in latent-variable modeling. A general framework for performing this task is called joint contrastive inference [17]. Let $p_\theta(z, x) = p(z)p_\theta(x|z)$ be a generative model, $q_\phi(z|x)$ be an amortized inference model and define the data-induced aggregated joint inference model as $\hat{q}_\phi(z, x) = p_{data}(x)q_\phi(z|x)$. The joint contrastive inference framework then minimizes some divergence between the two structured joint distributions $p_\theta(z, x)$ and $\hat{q}_\phi(z, x)$. This can be seen as a generalized version of amortized inference. There are some well-known examples of this kind of inference such as the Variational Autoencoder [22], Adversarially Learned Inference [17], and Wasserstein Variational Inference [2]. By using the DSW distance, we obtain a new joint contrastive inference method which inherits the benefits of optimal transport family of distances, yet remains scalable and computationally efficient. In particular, we learn both a generator and an inference model by solving:

$$(\theta_m, \phi_m) = \arg \min_{\theta \in \Theta, \phi \in \Phi} \mathbb{E}_{\hat{q}_\phi(z, x), p_\theta(z, x)}[\text{DSW}_p(\hat{q}_{\phi, m}(z, x), \hat{p}_{\theta, m}(z, x))], \quad (8)$$

where Θ, Φ are the parameter spaces, $\hat{q}_{\phi, m}(z, x)$ and $\hat{p}_{\theta, m}(z, x)$ are empirical distributions that sampled i.i.d data from $\hat{q}_\phi(z, x)$ and $p_\theta(z, x)$ respectively.

E Additional Experiments

In this appendix, we provide additional experimental results to yield more understandings about the minimum expected distance framework, which uses the new proposed distances. The appendix is divided into three parts, namely Appendices E.1, E.2 and E.3. Appendix E.1 is devoted to showing the performances of DGSW (see Appendix C for its definition) versus the generalized versions of other sliced distances on various factors which could affect the effectiveness of those methods. We also compare DSW to the recent augmented sliced Wasserstein method (ASW) [10]. Then we show the generated images from slice-based distances method for MNIST, CelebA and LSUN, when the number of projections varies. In Appendix E.2, we compare DSW to the projected robust subspcae Wasserstein (PRW) in [35, 29] on MNIST dataset . The comparison is to show Wasserstein-2 distance between the learned distribution and the data distribution versus the execution time. Finally, Appendix E.3 includes a comparison between DSW, DGSW, SW, Max-SW, Max-GSW, and Max-GSW-NN for the joint contrastive inference task on MNIST dataset.

E.1 Generative models

DGSW results on MNIST: Figure 4(a) shows the convergence of estimators of the learned distribution to the data distribution based on “generalized” sliced distances in the sense of Wasserstein-2 distance. Here, we use the circular function as the defining function for both GSW, Max-GSW, and DGSW (the polynomial function is very expensive in high-dimension). With 10 projections, DGSW produces better performance than GSW with 1000 projections, Max-GSW and Max-GSW-NN. There is a little improvement in the Wasserstein-2 score with DGSW when we increase the number of projections from 10 to 1000. For the computational speed shown in Figure 4(b), DGSW-10 is much faster than other reported methods, except the GSW-10 which has the worst Wasserstein-2 score.

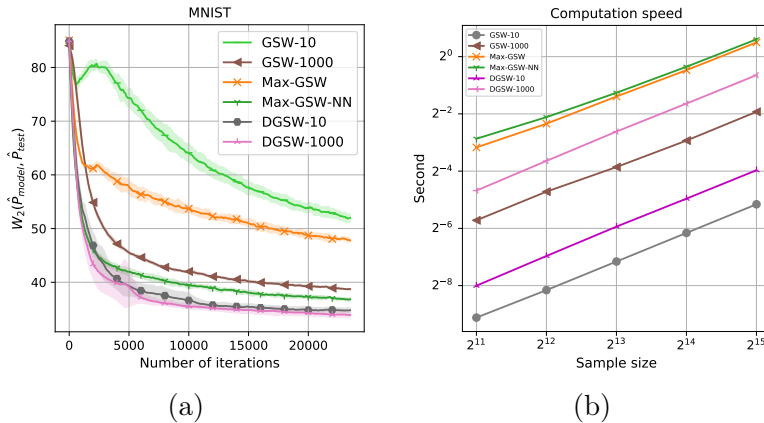


Figure 4: (a) Comparison between DGSW, GSW, Max-GSW and Max-GSW-NN using W_2 distance as metric. Here, GSW, Max-GSW and DGSW use circular function. (b) The computational speed over size of samples.

Effects of the number of samples: We conduct experiments to show how sample size (m in Appendix D.1) affects the results of DSW and DGSW in the MEDE framework. According to Figure 5(b), increasing the sample size leads to better performance of DSW. Similarly, increasing the sample size in the MEDE framework that uses DGSW (with circular defining function) helps

improve the results, see in Figure 5(d).

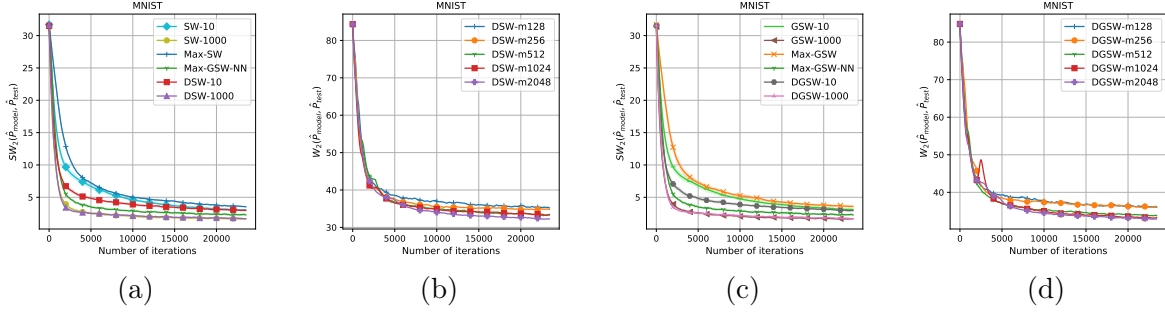


Figure 5: (a) Comparison between performances of DSW to SW, Max-SW and Max-GSW-NN using SW_2 distance as metric. (b) The effect of number of samples in minibatch to the convergence of DSW. (c) Comparing DGSW to GSW and Max-GSW-NN using SW_2 distance as metric. Here, GSW and DGSW use circular function. (d) The effect of number of samples in minibatch to the convergence of DGSW.

Effects of the number of gradient-updates: In both DSW and DGSW cases, we use a pushforward measure for the distribution over the sphere, and we use neural nets to find it. To learn these neural nets, we use gradient ascent to update their weights. In this experiment, we aim to find out how the number of iterations to update these neural net, affects the performance including the convergence behavior and computation speed. By increasing the number of updates from 1 to 10, both in DSW and DGSW, model distributions are much closer to data distribution; from 10 to 100 updates the results are improved but not too much, see the results in Figures 6(a) and 6(c). However, increasing update steps also lead to a computation problem as the al time increases considerably. When using 10 or 100 update steps, DSW and DGSW are slower than Max-SW, Max-GSW (50 gradient updates to find the max direction), and Max-GSW-NN (50 update times for the defining neural net function).

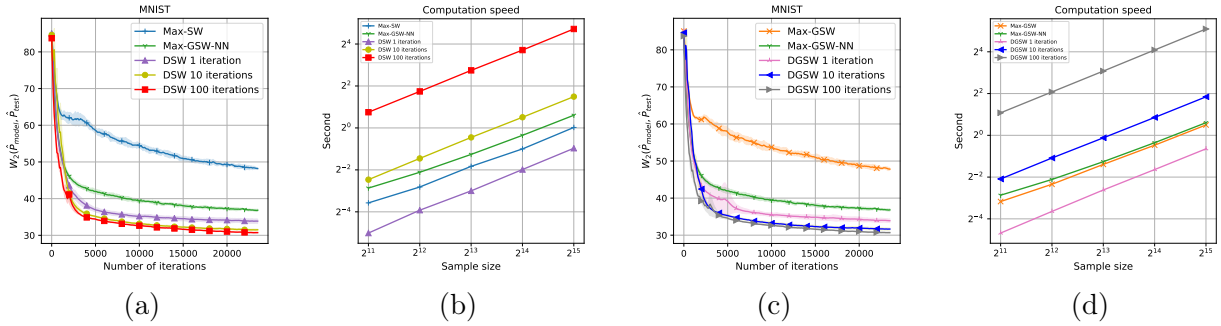


Figure 6: (a) and (c): Increasing the number of times to update push forward measure can improve the performance of both DSW and DGSW; (b) and (d): However, increasing the number of times to update push forward measure leads to much slower computation speed.

Quantitative results: We provide full FID scores of all distances mentioned in the papers and also the recent augmented sliced Wasserstein (ASW) [10] in Table 1. Based on the results in that

Table 1: FID score of generator models trained on CIFAR10 (100 epochs), CelebA (50 epochs), and LSUN (20 epochs) datasets in 64x64 resolution. Results are averaged from 5 different runs.

Model	n	CIFAR-10	CelebA	LSUN
SW	10^2	109.7 ± 5.64	90.11 ± 10.11	101.57 ± 3.24
GSW	10^2	103.11 ± 6.92	87.18 ± 8.97	92.58 ± 4.78
ASW	10^2	138.26 ± 8.31	122.11 ± 9.09	
DSW	10^2	62.83 ± 6.24	75.94 ± 5.54	46.02 ± 2.15
DGSW	10^2	68.01 ± 7.74	71.08 ± 4.24	46.91 ± 3.98
Max-SW		136.04 ± 8.35	100.09 ± 8.34	123.74 ± 5.51
Max-GSW-NN		86.04 ± 8.68	81.57 ± 7.72	56.83 ± 4.04
SW	10^4	98.61 ± 3.62	82.02 ± 6.33	62.75 ± 4.77
GSW	10^4	93.51 ± 6.12	84.22 ± 7.93	68.04 ± 2.17
ASW	10^4	121.38 ± 6.83	101 ± 7.36	
DSW	10^4	56.42 ± 3.78	66.85 ± 7.22	39.68 ± 2.33
DGSW	10^4	60.01 ± 5.58	65.8 ± 4.42	42.04 ± 4.21

Table 2: Computational speed per minibatch on CelebA and CIFAR10 dataset

Model	n	Second/Minibatch
SW	10^2	0.178
GSW	10^2	0.181
ASW	10^2	0.298
DSW	10^2	0.21
DGSW	10^2	0.212
Max-SW		1.821
Max-GSW-NN		1.895
SW	10^4	0.615
GSW	10^4	0.632
ASW	10^4	1.561
DSW	10^4	1.312
DGSW	10^4	1.384

table, DSW and DGSW (circular) achieve the best performance among all sliced distances. We also report the computational speed per minibatch in Table 2. The results show that DSW-100 is faster than DSW-10000 while its FID is lower. Regarding ASW, in our experiment, we find that the injective neural network, which is used to transform two target measures, is quite unstable to train and our obtained results with that distance are not good. Moreover, ASW is slower than DSW because ASW needs to double the dimension and still utilizes the uniform measure to slice on the new space. Note that, we use the implementation of ASW in <https://github.com/ShwanMario/ASWD>.

Qualitative results: We show random generated images from trained generators on MNIST, CelebA, CIFAR10 and LSUN datasets in Figures 7-11. Overall, we can see that the distributional approaches, i.e., DSW and DGSW distances, help to improve the quality of synthetic images in both linear and non-linear projection cases.

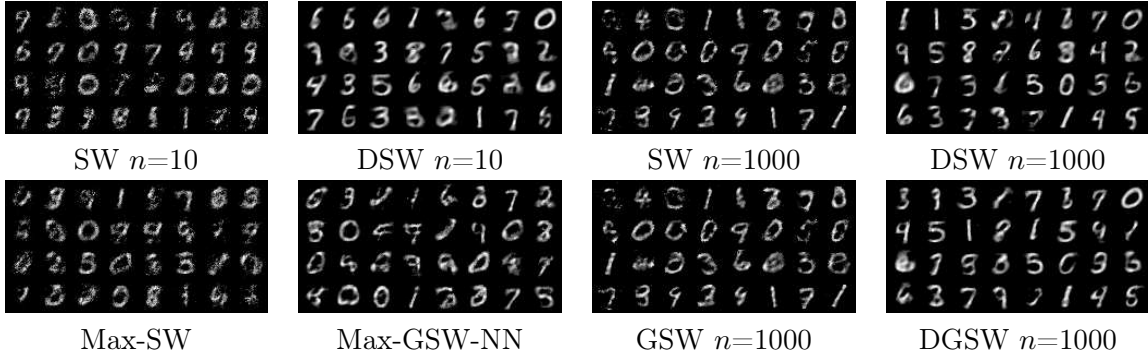


Figure 7: MNIST generated images from different generators, n is the number of projections.

Comparison with the special case of Max-GSW-NN: In Max-GSW-NN [25], one possible choice of neural network defining function is $g(x, \theta) = \langle x, f(\theta) \rangle$ where $f : \mathbb{S}^{d-1} \rightarrow \mathbb{S}^{d-1}$. That function f induces a probability measure on \mathbb{S}^{d-1} . Hence, optimizing f is equivalent to optimize over the set of probability measures without any constraints, which gives us an effect that is similar to max-SW. In contrast, the function f in our DSW is to find a push-forward probability measure that distributes high probability to informative directions, and this probability measure is regularized to avoid collapsing to a Dirac measure. To support our previous claim, we also do extra experiments on MNIST in Table 3 to clarify the role of the function f of DSW which makes DSW different from the given special case of Max-GSW-NN. The result shows that this version of Max-GSW-NN is similar to DSW when $\lambda_C = 0$ and both of them have the same performance as Max-SW.

Table 3: Comparison with the special case of Max-GSW-NN, denoting Max-GSW-NN(*) in the table, that uses the defining function $g(x, \theta) = \langle x, f(\theta) \rangle$ where $f : \mathbb{S}^{d-1} \rightarrow \mathbb{S}^{d-1}$.

Model	λ_C	Wasserstein-2
Max-SW	-	48.64
Max-GSW-NN (*)	-	49.21
DSW-10	0	49.81
DSW-10	1	38.41
DSW-10	10	33.40
DSW-10	100	40.08
DSW-10	1000	46.07

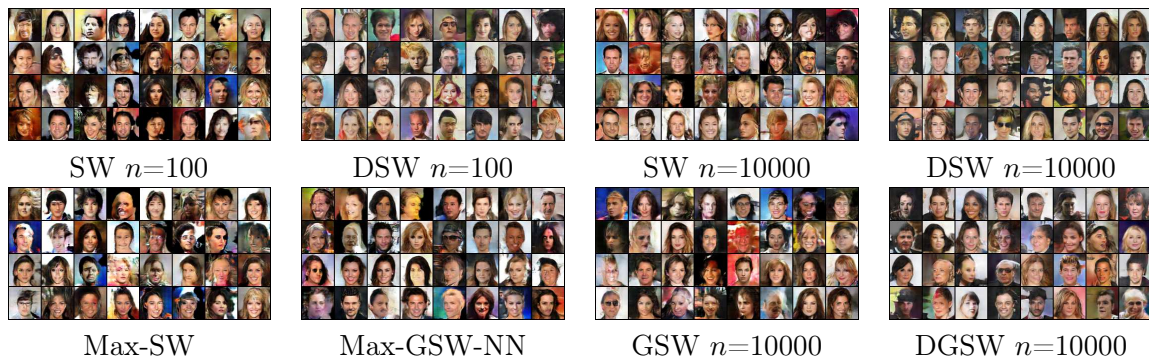


Figure 8: CelebA generated images from different generators, n is the number of projections.

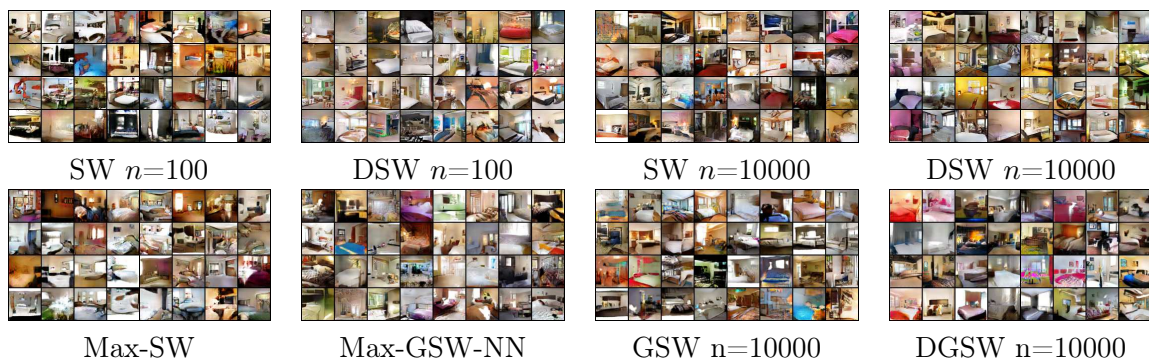


Figure 9: LSUN generated images from different generators where n is the number of projections.

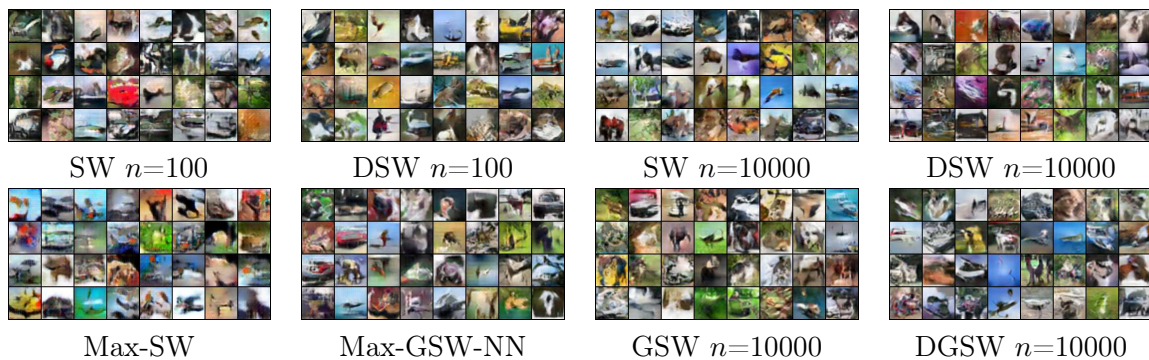


Figure 10: CIFAR10 generated images from different generators, n is the number of projections.

E.2 Comparison with Projected Robust Subspace Wasserstein

As shown in [35, 29], the main idea of projected robust subspace Wasserstein (PRW) is to find the optimal subspace (dimension ≥ 2) such that the Wasserstein-2 distance between two projected measures is maximal.

We first recall the definition of PRW in [35].

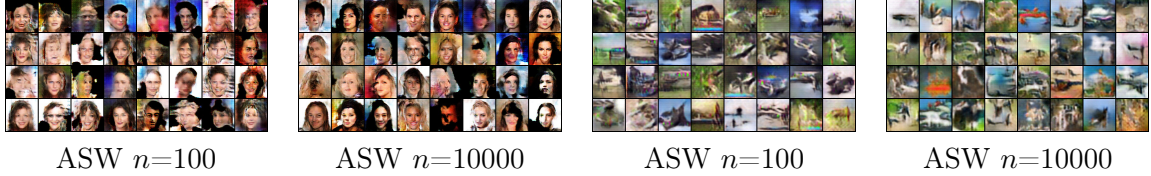


Figure 11: ASW generated images on CelebA and CIFAR10.

Definition 5. Let $\mathbb{V}_k(\mathbb{R}^d) = \{U \in \mathbb{R}^{d \times k} : U^\top U = I_k\}$ and $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$. Then, the projection robust 2-Wasserstein (PRW) distance between μ, ν is given by:

$$PRW_k(\mu, \nu) = \max_{U \in \mathbb{V}_k(\mathbb{R}^d)} W_2(U^\top \# \mu, U^\top \# \nu). \quad (9)$$

Since the projected dimension is bigger than 1, PRW does not have close-form solution on the projected space.

Experiments on generative model: We continue to use the MEDE framework on the same settings as previous experiments to compare DSW and PRW. To solve the optimization on Stiefel manifold in PRW, we use the "geopt" library [24]. We use one gradient step to solve the optimization problem of both DSW and PRW per one generator update. The experiments are carried out with both DSW and PRW on MNIST dataset. The number of projections of DSW takes value 10 and 1000 and the dimension of the subspace of PRW belongs to the set $\{2, 5, 10, 50\}$. We report the Wasserstein-2 results and the computational time in Table 4 and the generated images in Figure 12.

Table 4: Empirical Wasserstein-2 score and computation speed per minibatch on MNIST dataset.

Model	k -dimension	Wasserstein-2	Second/Minibatch
DSW-10	-	34.4	0.003
DSW-1000	-	33.11	0.018
PRW	2	65.39	0.086
PRW	5	35.99	0.092
PRW	10	26.57	0.11
PRW	50	24.38	0.12

According to Table 4, DSW with 10 projections obtains a better Wasserstein-2 score than the PRW with 5-dimensional subspace, while its corresponding computational time is 30 times faster than that of PRW. When PRW searches for the 50-dimensional subspace, the Wasserstein-2 score only improves 32.25% meanwhile the computational time increases by 10 times.

Next, we show some generated images from both DSW and PRW. We observe that these images are consistent with Wasserstein-2 score in the previous experiments.

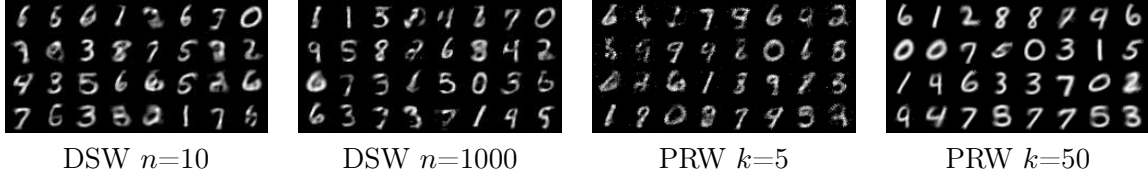


Figure 12: MNIST generated images from generators of DSW and PRW. Here, n is the number of projections of DSW and k is the projected dimension of PRW.

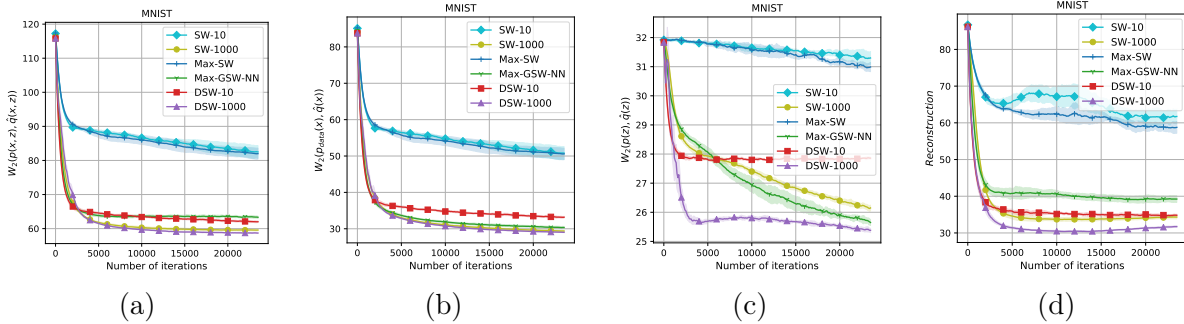


Figure 13: Joint inference model comparisons among DSW, SW, Max-SW, and Max-GSW-NN.

E.3 Joint contrastive inference

We test the performance of our distances in training encoder-generator models on MNIST using joint contrastive inference (JCI). In JCI, the joint generative latent-observed distribution $p_\theta(z, x) = p(z)p_\theta(x|z)$ is matched with the empirical joint latent-observed distribution $\hat{q}_\phi(z, x) = p_{\text{data}}(x)q_\phi(z|x)$ by minimizing a chosen distance (see Appendix D for a description of these models). We evaluate how close the two joint latent-observed distributions $p_\theta(z, x)$ and $\hat{q}_\phi(z, x)$ are, how close their corresponding marginals are (in Wasserstein-2 distance) and the ability of the encoder-generator in reconstructing images. These metrics are shown in Figures 13(a)-(d). The results show that DSW achieves better performance than SW using the same number of projections, with DSW-1000 achieves the best performance among all the other baselines in all metrics. We give experiments to compare DGSW with other non-linear sliced-Wasserstein distances in the joint contrastive inference task in Figure 14. We observed the same behavior as the linear case, the distributional version of GSW using circular function achieves better performance than the other non-linear sliced-based distances.

In order to illustrate the ability to reconstruct images of joint inference models, we show reconstructed images from the MNIST dataset. With 10 projections, SW and GSW were not able to recreate the digits; however, DSW and DGSW can recreate the digits quite correctly. Furthermore, Max-GSW-NN performs well in this task and is much better than Max-SW and Max-GSW. When having enough number of projections (for example 1000), it is very hard to compare SW, GSW, DSW, and DGSW by eyes. Nevertheless, according to reconstruction error plots in Figures 13 and 14, DSW, and DGSW distances are still better than the other sliced-based distances.

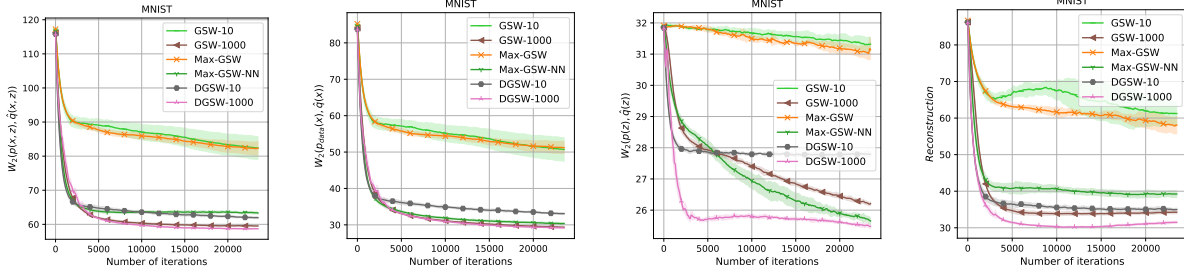


Figure 14: Joint inference model comparisons between non-linear sliced distances.

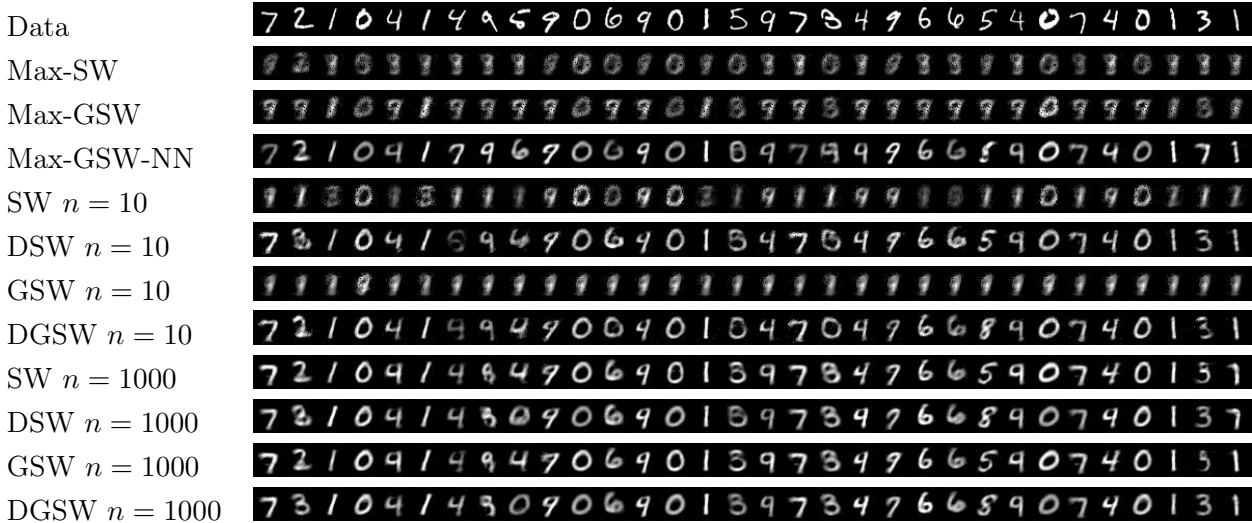


Figure 15: MNIST dataset reconstruction images (n is the number of projections).

F Experiment Settings

We use a single multi layer perceptron (MLP) layer with normalized output as the f function in the dual empirical forms of DSW and DGSW for the dual empirical forms of these distances). In all experiments, we use norm 2 as the ground metric for the Wasserstein distance. For GSW and DGSW, we use $r = 1000$ for circular function. We use the code at <https://github.com/kimiandj/gsw> for Max-SW and Max-GSW-NN (use 3 MLP layers with Leaky ReLU activation as defining function). In this implementation, Max-SW and Max-GSW-NN uses 50 gradient-update times per minibatch to find the optimal direction.

We train models on MNIST, CelebA, CIFAR10 with batch size = 512. On LSUN we use batch size = 4096. We use Adam optimizer for all models with learning rate=0.0005 and betas=(0.5, 0.999), $p = 2$. The range for hidden layer size of the MLP defining function of Max-GSW-NN is (32,100,784,1000). We tune λ_C of DSW and DGSW by grid searching in (1, 10, 100, 1000) in every experiment. The number of epochs for MNIST is 200, CelebA is 50, CIFAR10 is 100, and LSUN is 20.

In evaluation, we use empirical distribution with 10000 samples from two target distribution to compute discrete Wasserstein distance via linear programming..

Generator architecture was used for MNIST dataset:

$$z \in \mathbb{R}^{32} \rightarrow FC_{100} \rightarrow ReLU \rightarrow FC_{200} \rightarrow ReLU \rightarrow FC_{400} \rightarrow ReLU \rightarrow FC_{784} \rightarrow ReLU$$

Generator architecture was used for CelebA, CIFAR10 and LSUN dataset $z \in \mathbb{R}^{100} \rightarrow TransposeConv_{512} \rightarrow BatchNorm \rightarrow ReLU \rightarrow TransposeConv_{256} \rightarrow BatchNorm \rightarrow ReLU \rightarrow TransposeConv_{128} \rightarrow BatchNorm \rightarrow ReLU \rightarrow TransposeConv_{64} \rightarrow BatchNorm \rightarrow ReLU \rightarrow TransposeConv_1 \rightarrow Tanh$

Discriminator architecture was used for CelebA, CIFAR10 and LSUN dataset:

First part: $x \in \mathbb{R}^{64 \times 64 \times 3} \rightarrow Conv_{64} \rightarrow LeakyReLU_{0.2} \rightarrow Conv_{128} \rightarrow BatchNorm \rightarrow LeakyReLU_{0.2} \rightarrow Conv_{256} \rightarrow BatchNorm \rightarrow LeakyReLU_{0.2} \rightarrow Conv_{512} \rightarrow BatchNorm \rightarrow Tanh$

Second part: $Conv_1 \rightarrow Sigmoid$

Joint Contrastive inference encoder architecture on MNIST:

$$x \in \mathbb{R}^{28 \times 28} \rightarrow FC_{400} \rightarrow LeakyReLU_{0.2} \rightarrow FC_{200} \rightarrow LeakyReLU_{0.2} \rightarrow FC_{100} \rightarrow LeakyReLU_{0.2} \rightarrow FC_{32}$$

Joint Contrastive inference decoder architecture on MNIST:

$$z \in \mathbb{R}^{32} \rightarrow FC_{100} \rightarrow ReLU \rightarrow FC_{200} \rightarrow ReLU \rightarrow FC_{400} \rightarrow ReLU \rightarrow FC_{784} \rightarrow ReLU$$