

---

# Spherical Sliced-Wasserstein

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many variants of the Wasserstein distance have been introduced to reduce its  
2 original computational burden. In particular the Sliced-Wasserstein distance (SW),  
3 which leverages one-dimensional projections for which a closed-form solution  
4 of the Wasserstein distance is available, has received a lot of interest. Yet, it is  
5 restricted to data living in Euclidean spaces, while the Wasserstein distance has  
6 been studied and used recently on manifolds. We focus more specifically on the  
7 sphere, for which we define a novel SW discrepancy, which we call spherical Sliced-  
8 Wasserstein, making a first step towards defining SW discrepancies on manifolds.  
9 Our construction is notably based on closed-form solutions of the Wasserstein  
10 distance on the circle, together with a new spherical Radon transform. Along  
11 with efficient algorithms and the corresponding implementations, we illustrate its  
12 properties in several machine learning use cases where spherical representations  
13 of data are at stake: density estimation on the sphere, variational inference or  
14 hyperspherical auto-encoders.

## 15 1 Introduction

16 Optimal transport (OT) [101] has received a lot of attention in machine learning in the past few years.  
17 As it allows to compare distributions with metrics, it has been used for different tasks such as domain  
18 adaptation [24] or generative models [8], to name a few. The most classical distance used in OT is  
19 the Wasserstein distance. However, calculating it can be computationally expensive. Hence, several  
20 variants were proposed to alleviate the computational burden, such as the entropic regularization  
21 [26, 97], minibatch OT [35] or the sliced-Wasserstein distance (SW) for distributions supported on  
22 Euclidean spaces [90].

23 Although embedded in larger dimensional Euclidean spaces, data generally lie in practice on manifolds  
24 [36]. A simple manifold, but with lots of practical applications, is the hypersphere  $S^{d-1}$ . Several  
25 types of data are by essence spherical: a good example is found in directional data [71, 87] for  
26 which dedicated machine learning solutions are being developed [98], but other applications concern  
27 for instance geophysical data [32], meteorology [11], cosmology [86] or extreme value theory  
28 for the estimation of spectral measures [44]. Remarkably, in a more abstract setting, considering  
29 hyperspherical latent representations of data is becoming more and more common (*e.g.* [28, 70, 110]).  
30 For example, in the context of variational autoencoders [58], using priors on the sphere has been  
31 demonstrated to be beneficial [28]. Also, in the context of self-supervised learning (SSL), where  
32 one wants to learn discriminative representations in an unsupervised way, the hypersphere is usually  
33 considered for the latent representation [20, 21, 43, 104, 108]. It is thus of primary importance to  
34 develop machine learning tools that accommodate well with this specific geometry.

35 The OT theory on manifolds is well developed [38, 73, 101] and several works started to use  
 36 it in practice, with a focus mainly on the approximation of OT maps. For example, Cohen et al.  
 37 [23], Rezende and Racanière [91] approximate the OT map to define normalizing flows on Riemannian  
 38 manifolds, Cui et al. [25], Hamfeldt and Turnquist [45, 46] derive algorithms to approximate the OT  
 39 map on the sphere, Alvarez-Melis et al. [5], Hoyos-Idrobo [51] learn the transport map on hyperbolic  
 40 spaces. However, the computational bottleneck to compute the Wasserstein distance on such spaces  
 41 remains, and, as underlined in the conclusion of [74], defining SW distances on manifolds would be  
 42 of much interest.

43 **Contributions.** Therefore, by leveraging properties of the Wasserstein distance on the circle [89],  
 44 we define the first, to the best of our knowledge, SW discrepancy on a non trivial manifold, namely the  
 45 sphere  $S^{d-1}$ , and hence we make a first step towards defining SW distances on Riemannian manifolds.  
 46 We make connections with a new spherical Radon transform and analyze some of its properties.  
 47 We discuss the underlying algorithmic procedure, and notably provide an efficient implementation  
 48 when computing the discrepancy against a uniform distribution. Then, we show that we can use this  
 49 discrepancy on different tasks such as density estimation, variational inference or generative modeling.

## 50 2 Background

51 The aim of this paper is to define a Sliced-Wasserstein discrepancy on the hypersphere  $S^{d-1} = \{x \in$   
 52  $\mathbb{R}^d, \|x\|_2 = 1\}$ . Therefore, in this section, we introduce the Wasserstein distance on manifolds and  
 53 the classical SW distance on  $\mathbb{R}^d$ .

### 54 2.1 Wasserstein distance

55 Since we are interested in defining a SW discrepancy on the sphere, we start by introducing the  
 56 Wasserstein distance on a Riemannian manifold  $M$  endowed with the Riemannian distance  $d$ . We  
 57 refer to [38, 101] for more details.

58 Let  $p \geq 1$  and  $\mu, \nu \in \mathcal{P}_p(M) = \{\mu \in \mathcal{P}(M), \int_M d^p(x, x_0) d\mu(x) < \infty \text{ for some } x_0 \in M\}$ . Then,  
 59 the  $p$ -Wasserstein distance between  $\mu$  and  $\nu$  is defined as

$$W_p^p(\mu, \nu) = \inf_{\gamma \in \Pi(\mu, \nu)} \int_{M \times M} d^p(x, y) d\gamma(x, y), \quad (1)$$

60 where  $\Pi(\mu, \nu) = \{\gamma \in \mathcal{P}(M \times M), \forall A \subset M, \gamma(M \times A) = \nu(A) \text{ and } \gamma(A \times M) = \mu(A)\}$   
 61 denotes the set of couplings.

62 For discrete probability measures, the Wasserstein distance can be computed using linear programs  
 63 [88]. However, these algorithms have a  $O(n^3 \log n)$  complexity *w.r.t.* the number of samples  $n$   
 64 which is computationally intensive. Therefore, a whole literature consists of defining alternative  
 65 discrepancies which are cheaper to compute. On Euclidean spaces, one of them is the Sliced-  
 66 Wasserstein distance.

### 67 2.2 Sliced-Wasserstein distance

68 On  $M = \mathbb{R}^d$  with  $d(x, y) = \|x - y\|_p^p$ , a more attractive distance is the Sliced-Wasserstein (SW)  
 69 distance. This distance relies on the appealing fact that for one dimensional measures  $\mu, \nu \in \mathcal{P}(\mathbb{R})$ ,  
 70 we have the following closed-form [88, Remark 2.30]

$$W_p^p(\mu, \nu) = \int_0^1 |F_\mu^{-1}(u) - F_\nu^{-1}(u)|^p du, \quad (2)$$

71 where  $F_\mu^{-1}$  (resp.  $F_\nu^{-1}$ ) is the quantile function of  $\mu$  (resp.  $\nu$ ). From this property, Bonnotte  
 72 [16], Rabin et al. [90] defined the SW distance as

$$\forall \mu, \nu \in \mathcal{P}_p(\mathbb{R}^d), SW_p^p(\mu, \nu) = \int_{S^{d-1}} W_p^p(P_{\#}^\theta \mu, P_{\#}^\theta \nu) d\lambda(\theta), \quad (3)$$

73 where  $P^\theta(x) = \langle x, \theta \rangle$ ,  $\lambda$  is the uniform distribution on  $S^{d-1}$  and for any Borel set  $A \in \mathcal{B}(\mathbb{R}^d)$ ,  
 74  $P^\theta_\# \mu(A) = \mu((P^\theta)^{-1}(A))$ .

75 This distance can be approximated efficiently by using a Monte-Carlo approximation [75], and  
 76 amounts to a complexity of  $O(Ln \log n)$  where  $L$  denotes the number of projections used for the  
 77 Monte-Carlo approximation and  $n$  the number of samples.

78 SW can also be written through the Radon transform [15]. Let  $f \in L^1(\mathbb{R}^d)$ , then the Radon transform  
 79  $R : L^1(\mathbb{R}^d) \rightarrow L^1(\mathbb{R} \times S^{d-1})$  is defined as [48]

$$\forall \theta \in S^{d-1}, \forall t \in \mathbb{R}, Rf(t, \theta) = \int_{\mathbb{R}^d} f(x) \mathbb{1}_{\{\langle x, \theta \rangle = t\}} dx. \quad (4)$$

80 Its dual  $R^* : C_0(\mathbb{R} \times S^{d-1}) \rightarrow C_0(\mathbb{R}^d)$  (also known as back-projection operator), where  $C_0$   
 81 denotes the set of continuous functions that vanish at infinity, satisfies for all  $f, g, \langle Rf, g \rangle_{\mathbb{R} \times S^{d-1}} =$   
 82  $\langle f, R^*g \rangle_{\mathbb{R}^d}$  and can be defined as [13, 15]

$$\forall g \in C_0(\mathbb{R} \times S^{d-1}), \forall x \in \mathbb{R}^d, R^*g(x) = \int_{S^{d-1}} g(\langle x, \theta \rangle, \theta) d\theta. \quad (5)$$

83 Therefore, by duality, we can define the Radon transform of a measure  $\mu \in \mathcal{M}(\mathbb{R}^d)$  as the measure  
 84  $R\mu \in \mathcal{M}(\mathbb{R} \times S^{d-1})$  such that for all  $g \in C_0(\mathbb{R} \times S^{d-1})$ ,  $\langle R\mu, g \rangle_{\mathbb{R} \times S^{d-1}} = \langle \mu, R^*g \rangle_{\mathbb{R}^d}$ . Since  $R\mu$   
 85 is a measure on the product space  $\mathbb{R} \times S^{d-1}$ , we can disintegrate it *w.r.t.*  $\lambda$ , the uniform measure  
 86 on  $S^{d-1}$  [6], as  $R\mu = \lambda \otimes K$  with  $K$  a probability kernel on  $S^{d-1} \times \mathcal{B}(\mathbb{R})$ , *i.e.* for all  $\theta \in S^{d-1}$ ,  
 87  $K(\theta, \cdot)$  is a probability on  $\mathbb{R}$ , for any Borel set  $A \in \mathcal{B}(\mathbb{R})$ ,  $K(\cdot, A)$  is measurable, and

$$\forall \phi \in C(\mathbb{R} \times S^{d-1}), \int_{\mathbb{R} \times S^{d-1}} \phi(t, \theta) d(R\mu)(t, \theta) = \int_{S^{d-1}} \int_{\mathbb{R}} \phi(t, \theta) K(\theta, dt) d\lambda(\theta). \quad (6)$$

88 By Proposition 6 in [15], we have that for  $\lambda$ -almost every  $\theta \in S^{d-1}$ ,  $(R\mu)^\theta = P^\theta_\# \mu$  where we denote  
 89  $K(\theta, \cdot) = (R\mu)^\theta$ . Therefore, we have

$$\forall \mu, \nu \in \mathcal{P}_p(\mathbb{R}^d), SW_p^p(\mu, \nu) = \int_{S^{d-1}} W_p^p((R\mu)^\theta, (R\nu)^\theta) d\lambda(\theta). \quad (7)$$

90 Variants of SW have been defined in recent works, either by integrating *w.r.t.* different distributions  
 91 [31, 80, 81], by projecting on  $\mathbb{R}$  using different projections [78, 79] or Radon transforms [22, 60], or  
 92 by projecting on subspaces of higher dimensions [52, 66, 67, 85].

### 93 3 A Sliced-Wasserstein discrepancy on the sphere

94 Our goal here is to define a sliced-Wasserstein distance on the sphere  $S^{d-1}$ . To that aim, we proceed  
 95 analogously to the classical Euclidean space. We first rely on the nice properties of the Wasserstein  
 96 distance on the circle [89] and then propose to project distributions lying on the sphere to great circles.  
 97 Hence, circles play the role of the real line for the hypersphere. In this section, we first describe the  
 98 OT problem on the circle, then we define a sliced-Wasserstein discrepancy on the sphere and discuss  
 99 some of its properties. Notably, we derive a new spherical Radon transform which is linked to our  
 100 newly defined spherical SW. We refer to Appendix A for the proofs.

#### 101 3.1 Optimal transport on the circle

102 On the circle  $S^1 = \mathbb{R}/\mathbb{Z}$  equipped with the geodesic distance  $d_{S^1}$ , an appealing formulation of the  
 103 Wasserstein distance is available [30]. First, let us parametrize  $S^1$  by  $[0, 1[$ , then the geodesic distance  
 104 can be written as [89], for all  $x, y \in [0, 1[$ ,  $d_{S^1}(x, y) = \min(|x - y|, 1 - |x - y|)$ . Then, for the cost  
 105 function  $c(x, y) = h(d_{S^1}(x, y))$  with  $h : \mathbb{R} \rightarrow \mathbb{R}^+$  an increasing convex function, the Wasserstein  
 106 distance between  $\mu \in \mathcal{P}(S^1)$  and  $\nu \in \mathcal{P}(S^1)$  can be written as

$$W_c(\mu, \nu) = \inf_{\alpha \in \mathbb{R}} \int_0^1 h(|F_\mu^{-1}(t) - (F_\nu - \alpha)^{-1}(t)|) dt, \quad (8)$$

107 where  $F_\mu : [0, 1[ \rightarrow [0, 1]$  denotes the cumulative distribution function (cdf) of  $\mu$ ,  $F_\mu^{-1}$  its quantile  
 108 function and  $\alpha$  is a shift parameter. The optimization problem over the shifted cdf  $F_\nu - \alpha$  can be seen  
 109 as looking for the best “cut” (or origin) of the circle into the real line because of the 1-periodicity.  
 110 Indeed, the proof of this result for discrete distributions in [89] consists in cutting the circle at  
 111 the optimal point and wrapping it around the real line, for which the optimal transport map is the  
 112 increasing rearrangement  $F_\nu^{-1} \circ F_\mu$  which can be obtained for discrete distributions by sorting the  
 113 points [88].

114 Rabin et al. [89] showed that the minimization problem is convex and coercive in the shift parameter  
 115 and Delon et al. [30] derived a binary search algorithm to find it. For the particular case of  $h = \text{Id}$ , it  
 116 can further be shown [19, 106] that

$$W_1(\mu, \nu) = \inf_{\alpha \in \mathbb{R}} \int_0^1 |F_\mu(t) - F_\nu(t) - \alpha| dt. \quad (9)$$

117 In this case, we know exactly the minimum which is attained at the level median [53]. For  
 118  $f : [0, 1[ \rightarrow \mathbb{R}$ ,

$$\text{LevMed}(f) = \min \left\{ \operatorname{argmin}_{\alpha \in \mathbb{R}} \int_0^1 |f(t) - \alpha| dt \right\} = \inf \left\{ t \in \mathbb{R}, \beta(\{x \in [0, 1[, f(x) \leq t\}) \geq \frac{1}{2} \right\}, \quad (10)$$

119 where  $\beta$  is the Lebesgue measure. Therefore, we also have

$$W_1(\mu, \nu) = \int_0^1 |F_\mu(t) - F_\nu(t) - \text{LevMed}(F_\mu - F_\nu)| dt. \quad (11)$$

120 Since we know the minimum, we do not need the binary search and we can approximate the integral  
 121 very efficiently as we only need to sort the samples to compute the level median and the cdfs.

122 Another interesting setting in practice is to compute  $W_2$ , *i.e.* with  $h(x) = x^2$ , *w.r.t.* a uniform  
 123 distribution  $\nu$  on the circle. We derive here the optimal shift  $\hat{\alpha}$  for the Wasserstein distance between  $\mu$   
 124 an arbitrary distribution on  $S^1$  and  $\nu$ . We also provide a closed-form when  $\mu$  is a discrete distribution.

125 **Proposition 1.** *Let  $\mu \in \mathcal{P}_2(S^1)$  and  $\nu = \text{Unif}(S^1)$ . Then,*

$$W_2^2(\mu, \nu) = \int_0^1 |F_\mu^{-1}(t) - t - \hat{\alpha}|^2 dt \quad \text{with} \quad \hat{\alpha} = \int x d\mu(x) - \frac{1}{2}. \quad (12)$$

126 *In particular, if  $x_1 < \dots < x_n$  and  $\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ , then*

$$W_2^2(\mu_n, \nu) = \frac{1}{n} \sum_{i=1}^n x_i^2 - \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 + \frac{1}{n^2} \sum_{i=1}^n (n+1-2i)x_i + \frac{1}{12}. \quad (13)$$

127 This proposition offers an intuitive interpretation: the optimal cut point between an empirical and a  
 128 uniform distributions is the antipodal point of the circular mean of the discrete samples. Moreover, a  
 129 very efficient algorithm can be derived from this property, as it solely requires a sorting operation to  
 130 compute the order statistics of the samples.

### 131 3.2 Definition of SW on the sphere

132 On the hypersphere, the counterpart of straight lines are the great circles, which correspond to the  
 133 geodesics. Moreover, we can compute the Wasserstein distance on the circle fairly efficiently. Hence,  
 134 to define a sliced-Wasserstein discrepancy on this manifold, we propose, analogously to the classical  
 135 SW distance, to project measures on great circles. The most natural way to project points from  $S^{d-1}$   
 136 to a great circle  $C$  is to use the geodesic projection [40, 55] defined as

$$\forall x \in S^{d-1}, P^C(x) = \operatorname{argmin}_{y \in C} d_{S^{d-1}}(x, y), \quad (14)$$

137 where  $d_{S^{d-1}}(x, y) = \arccos(\langle x, y \rangle)$  is the geodesic distance. See Figure 1 for an illustration of the  
 138 geodesic projection on a great circle. Note that the projection is unique for almost every  $x$  (see

139 [9, Proposition 4.2] and Appendix B.1) and hence the pushforward  $P_{\#}^C \mu$  of absolutely continuous  
 140 measures *w.r.t.* the Lebesgue measure  $\mu \in \mathcal{P}_{p,ac}(S^{d-1})$  is well defined.

141 Great circles can be obtained by intersecting  $S^{d-1}$  with a 2-dimensional plane [56]. Therefore,  
 142 to average over all great circles, we propose to integrate over the Grassmann manifold  $\mathcal{G}_{d,2} =$   
 143  $\{E \subset \mathbb{R}^d, \dim(E) = 2\}$  [2, 10] and then to project the distribution onto the intersection with the  
 144 hypersphere. Since the Grassmannian is not very practical, we consider the identification using the  
 145 set of rank 2 projectors:

$$\mathcal{G}_{d,2} = \{P \in \mathbb{R}^{d \times d}, P^T = P, P^2 = P, \text{Tr}(P) = 2\} = \{UU^T, U \in \mathbb{V}_{d,2}\}, \quad (15)$$

146 where  $\mathbb{V}_{d,2} = \{U \in \mathbb{R}^{d \times 2}, U^T U = I_2\}$  is the Stiefel manifold [10].

147 Finally, we can define the Spherical Sliced-Wasserstein distance (SSW) for  $p \geq 1$  between locally  
 148 absolutely continuous measures *w.r.t.* the Lebesgue measure [9]  $\mu, \nu \in \mathcal{P}_{p,ac}(S^{d-1})$  as

$$SSW_p^p(\mu, \nu) = \int_{\mathbb{V}_{d,2}} W_p^p(P_{\#}^U \mu, P_{\#}^U \nu) d\sigma(U), \quad (16)$$

149 where  $\sigma$  is the uniform distribution over the Stiefel manifold  $\mathbb{V}_{d,2}$ ,  $P^U$  is the geodesic projection on  
 150 the great circle generated by  $U$  and then projected on  $S^1$ , *i.e.*

$$\forall U \in \mathbb{V}_{d,2}, \forall x \in S^{d-1}, P^U(x) = U^T \underset{y \in \text{span}(UU^T) \cap S^{d-1}}{\text{argmin}} d_{S^{d-1}}(x, y) = \underset{z \in S^1}{\text{argmin}} d_{S^{d-1}}(x, Uz), \quad (17)$$

151 and the Wasserstein distance is defined with the geodesic distance  $d_{S^1}$ .

152 Moreover, we can derive a closed form expression which  
 153 will be very useful in practice:

154 **Lemma 1.** *Let  $U \in \mathbb{V}_{d,2}$  then for a.e.  $x \in S^{d-1}$ ,*

$$P^U(x) = \frac{U^T x}{\|U^T x\|_2}. \quad (18)$$

155 Hence, we notice from this expression of the projection  
 156 that we recover almost the same formula as Lin et al. [66]  
 157 but with an additional  $\ell^2$  normalization which projects  
 158 the data on the circle. As in [66], we could project on  
 159 a higher dimensional subsphere by integrating over  $\mathbb{V}_{d,k}$   
 160 with  $k \geq 2$ . However, we would lose the computational  
 161 efficiency provided by the properties of the Wasserstein  
 162 distance on the circle.

### 163 3.3 A Spherical Radon Transform

164 As for the classical SW distance, we can derive a second  
 165 formulation using a Radon transform. Let  $f \in L^1(S^{d-1})$ , we define a spherical Radon transform  
 166  $\tilde{R} : L^1(S^{d-1}) \rightarrow L^1(S^1 \times \mathbb{V}_{d,2})$  as

$$\forall z \in S^1, \forall U \in \mathbb{V}_{d,2}, \tilde{R}f(z, U) = \int_{S^{d-1}} f(x) \mathbb{1}_{\{z=P^U(x)\}} dx. \quad (19)$$

167 This is basically the same formulation as the classical Radon transform [48, 77] where we replaced  
 168 the real line coordinate  $t$  by the coordinate on the circle  $z$  and the projection is the geodesic one  
 169 which is well suited to the sphere. This transform is actually new since we integrate over different  
 170 sets compared to existing works on spherical Radon transforms.

171 Then, analogously to the classical Radon transform, we can define the back-projection operator  
 172  $\tilde{R}^* : C_0(S^1 \times \mathbb{V}_{d,2}) \rightarrow C_b(S^{d-1})$ ,  $C_b(S^{d-1})$  being the space of continuous bounded functions, for  
 173  $g \in C_0(S^1 \times \mathbb{V}_{d,2})$  as for a.e.  $x \in S^{d-1}$ ,

$$\tilde{R}^*g(x) = \int_{\mathbb{V}_{d,2}} g(P^U(x), U) d\sigma(U). \quad (20)$$

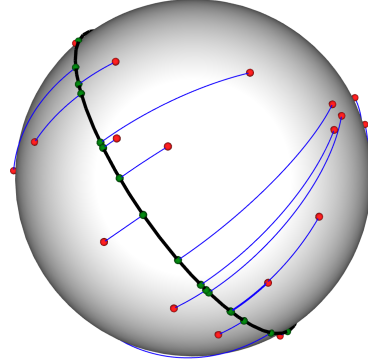


Figure 1: Illustration of the geodesic projections on a great circle (in black). In red, random points sampled on the sphere. In green the projections and in blue the trajectories.

174 **Proposition 2.**  $\tilde{R}^*$  is the dual operator of  $\tilde{R}$ , i.e. for all  $f \in L^1(S^{d-1})$ ,  $g \in C_0(S^1 \times \mathbb{V}_{d,2})$ ,

$$\langle \tilde{R}f, g \rangle_{S^1 \times \mathbb{V}_{d,2}} = \langle f, \tilde{R}^*g \rangle_{S^{d-1}}. \quad (21)$$

175 Now that we have a dual operator, we can also define the Radon transform of an absolutely continuous  
176 measure  $\mu \in \mathcal{M}_{ac}(S^{d-1})$  by duality [13, 15] as the measure  $\tilde{R}\mu$  satisfying

$$\forall g \in C_0(S^1 \times \mathbb{V}_{d,2}), \int_{S^1 \times \mathbb{V}_{d,2}} g(z, U) d(\tilde{R}\mu)(z, U) = \int_{S^{d-1}} \tilde{R}^*g(x) d\mu(x). \quad (22)$$

177 Since  $\tilde{R}\mu$  is a measure on the product space  $S^1 \times \mathbb{V}_{d,2}$ ,  $\tilde{R}\mu$  can be disintegrated [6, Theorem 5.3.1]  
178 w.r.t.  $\sigma$  as  $\tilde{R}\mu = \sigma \otimes K$  where  $K$  is a probability kernel on  $\mathbb{V}_{d,2} \times S^1$  with  $S^1$  the Borel  $\sigma$ -field of  
179  $S^1$ . We will denote for  $\sigma$ -almost every  $U \in \mathbb{V}_{d,2}$ ,  $(\tilde{R}\mu)^U = K(U, \cdot)$  the conditional probability.

180 **Proposition 3.** Let  $\mu \in \mathcal{M}_{ac}(S^{d-1})$ , then for  $\sigma$ -almost every  $U \in \mathbb{V}_{d,2}$ ,  $(\tilde{R}\mu)^U = P_{\#}^U \mu$ .

181 Finally, we can write SSW (16) using this Radon transform:

$$\forall \mu, \nu \in \mathcal{P}_{p,ac}(S^{d-1}), SSW_p^p(\mu, \nu) = \int_{\mathbb{V}_{d,2}} W_p^p((\tilde{R}\mu)^U, (\tilde{R}\nu)^U) d\sigma(U). \quad (23)$$

182 Note that a natural way to define SW distances can be through already known Radon transforms using  
183 the formulation (23). It is for example what was done in [60] using generalized Radon transforms  
184 [34, 50] to define generalized SW distances, or in [22] with the spatial Radon transform. However,  
185 for known spherical Radon transforms [1, 7] such as the Minkowski-Funk transform [27] or more  
186 generally the geodesic Radon transform [95], there is no natural way that we know of to integrate  
187 over some product space and allowing to define a SW distance using disintegration.

188 As observed by Kolouri et al. [60] for the generalized SW distances (GSW), studying the injectivity  
189 of the related Radon transforms allows to study the set on which SW is actually a distance. While  
190 the classical Radon transform integrates over hyperplanes of  $\mathbb{R}^d$ , the generalized Radon transform  
191 over hypersurfaces [60] and the Minkowski-Funk transform over “big circles”, i.e. the intersection  
192 between a hyperplane and  $S^{d-1}$  [96], the set of integration here is a half of a big circle. Hence,  $\tilde{R}$  is  
193 related to the hemispherical transform [94] on  $S^{d-2}$ . We refer to Appendix A.6 for more details on  
194 the links with the hemispherical transform. Using these connections, we can derive the kernel of  $\tilde{R}$  as  
195 the set of even measures which are null over all hyperplanes intersected with  $S^{d-1}$ .

196 **Proposition 4.**  $\ker(\tilde{R}) = \{\mu \in \mathcal{M}_{\text{even}}(S^{d-1}), \forall H \in \mathcal{G}_{d,d-1}, \mu(H \cap S^{d-1}) = 0\}$  where  $\mu \in$   
197  $\mathcal{M}_{\text{even}}$  if for all  $f \in C(S^{d-1})$ ,  $\langle \mu, f \rangle = \langle \mu, f_+ \rangle$  with  $f_+(x) = (f(x) + f(-x))/2$  for all  $x$ .

198 We leave for future works checking whether this set is null or not. Hence, we conclude here that SSW  
199 is a pseudo-distance, but a distance on the sets of injectivity of  $\tilde{R}$  [4].

200 **Proposition 5.** Let  $p \geq 1$ ,  $SSW_p$  is a pseudo-distance on  $\mathcal{P}_{p,ac}(S^{d-1})$ .

## 201 4 Implementation

202 In practice, we approximate the distributions with empirical approximations and, as for the classical  
203 SW distance, we rely on the Monte-Carlo approximation of the integral on  $\mathbb{V}_{d,2}$ . We first need to  
204 sample from the uniform distribution  $\sigma \in \mathcal{P}(\mathbb{V}_{d,2})$ . This can be done by first constructing  $Z \in \mathbb{R}^{d \times 2}$   
205 by drawing each of its component from the standard normal distribution  $\mathcal{N}(0, 1)$  and then applying  
206 the QR decomposition [67]. Once we have  $(U_\ell)_{\ell=1}^L \sim \sigma$ , we project the samples on the circle  $S^1$  by  
207 applying Lemma 1 and we compute the coordinates on the circle using the atan2 function. Finally,  
208 we can compute the Wasserstein distance on the circle by either applying the binary search algorithm  
209 of [30] or the level median formulation (11) for  $SSW_1$ . In the particular case in which we want to  
210 compute  $SSW_2$  between a measure  $\mu$  and the uniform measure on the sphere  $\nu = \text{Unif}(S^{d-1})$ , we  
211 can use the appealing fact that the projection of  $\nu$  on the circle is uniform, i.e.  $P_{\#}^U \nu = \text{Unif}(S^1)$   
212 (particular case of Theorem 3.1 in [55], see Appendix B.3). Hence, we can use the Proposition  
213 1 to compute  $W_2$ , which allows a very efficient implementation either by the closed-form (13) or  
214 approximation by rectangle method of (12). This will be of particular interest for applications in  
215 Section 5 such as autoencoders. We sum up the procedure in Algorithm 1.

---

**Algorithm 1** SSW
 

---

**Input:**  $(x_i)_{i=1}^n \sim \mu, (y_j)_{j=1}^m \sim \nu, L$  the number of projections,  $p$  the order  
**for**  $\ell = 1$  **to**  $L$  **do**  
 Draw a random matrix  $Z \in \mathbb{R}^{d \times 2}$  with for all  $i, j, Z_{i,j} \sim \mathcal{N}(0, 1)$   
 $U = \text{QR}(Z) \sim \sigma$   
 Project on  $S^1$  the points:  $\forall i, j, \hat{x}_i^\ell = \frac{U^T x_i}{\|U^T x_i\|_2}, \hat{y}_j^\ell = \frac{U^T y_j}{\|U^T y_j\|_2}$   
 Compute the coordinates on the circle  $S^1$ :  $\forall i, j, \tilde{x}_i^\ell = (\pi + \text{atan2}(-x_{i,2}, -x_{i,1})) / (2\pi), \tilde{y}_j^\ell = (\pi + \text{atan2}(-y_{j,2}, -y_{j,1})) / (2\pi)$   
 Compute  $W_p^p(\frac{1}{n} \sum_{i=1}^n \delta_{\tilde{x}_i^\ell}, \frac{1}{m} \sum_{j=1}^m \delta_{\tilde{y}_j^\ell})$  by binary search or (11) for  $p = 1$   
**end for**  
 Return  $SSW_p^p(\mu, \nu) \approx \frac{1}{L} \sum_{\ell=1}^L W_p^p(\frac{1}{n} \sum_{i=1}^n \delta_{\tilde{x}_i^\ell}, \frac{1}{m} \sum_{j=1}^m \delta_{\tilde{y}_j^\ell})$

---

216 **Complexity.** Let us note  $n$  (resp.  $m$ ) the number of samples of  $\mu$  (resp.  $\nu$ ), and  $L$  the number of  
 217 projections. First, we need to compute the QR factorization of  $L$  matrices of size  $d \times 2$ . This can be  
 218 done in  $O(Ld)$  by using *e.g.* Householder reflections [42, Chapter 5.2] or the Scharwz-Rutishauser  
 219 algorithm [41]. Projecting the points on  $S^1$  by Lemma 1 is in  $O((n+m)dL)$  since we need to  
 220 compute  $L(n+m)$  products between  $U_\ell^T \in \mathbb{R}^{2 \times d}$  and  $x \in \mathbb{R}^d$ . For the binary search or particular  
 221 case formula (11) and (13), we need first to sort the points. But the binary search also adds a cost  
 222 of  $O((n+m) \log(\frac{1}{\epsilon}))$  to approximate the solution with precision  $\epsilon$  [30] and the computation of  
 223 the level median requires to sort  $(n+m)$  points. Hence, for the general  $SSW_p$ , the complexity  
 224 is  $O(L(n+m)(d + \log(\frac{1}{\epsilon})) + Ln \log n + Lm \log m)$  versus  $O(L(n+m)(d + \log(n+m)))$  for  
 225  $SSW_1$  with the level median and  $O(Ln(d + \log n))$  for  $SSW_2$  against a uniform with the particular  
 226 advantage that we do not need uniform samples in this case.

227 **Runtime Comparison.** We perform here some runtime comparisons. Using Pytorch [83], we imple-  
 228 mented the binary search algorithm of [30] and used it with  $\epsilon = 10^{-6}$ . We also implemented  $SSW_1$   
 229 using the level median formula (11) and  $SSW_2$  against a uniform measure (12). All experiments are  
 230 conducted on GPU.

231 On Figure 2, we compare the runtime between  
 232 two distributions on  $S^2$  between SSW, the  
 233 Wasserstein distance and the entropic approx-  
 234 imation using the Sinkhorn algorithm [26]  
 235 with the geodesic distance as cost function.  
 236 The distributions were approximated using  
 237  $n \in \{10^2, 10^3, 10^4, 5 \cdot 10^4, 10^5\}$  samples of  
 238 each distribution and we report the mean over  
 239 20 computations. We use the Python Optimal  
 240 Transport (POT) library [39] to compute the  
 241 Wasserstein distance and the entropic approx-  
 242 imation. For large enough batches, we observe  
 243 that SSW is much faster than its Wasserstein  
 244 counterpart, and it also scales better in term of  
 245 memory because of the need to store the  $n \times n$   
 246 cost matrix. For small batches, the computation  
 247 of SSW actually takes longer because of the  
 248 computation of the QR factorizations and of the  
 249 projections. For bigger batches, it is bounded  
 250 by the sorting operation and we recover the quasi-linear slope. Furthermore, as expected, the fastest  
 251 algorithms are  $SSW_1$  with the level median and  $SSW_2$  against a uniform as they have a quasilinear  
 252 complexity. We report in Appendix C.2 other runtimes experiments *w.r.t.* to *e.g.* the number of  
 253 projections or the dimension.

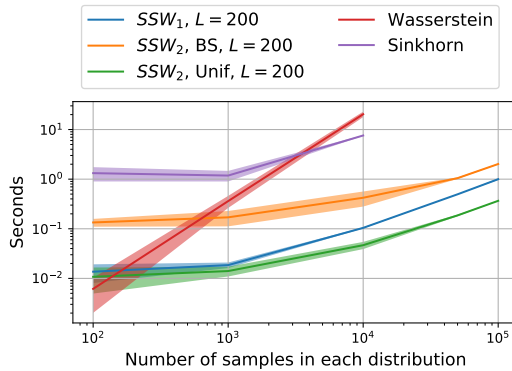


Figure 2: Runtime comparison in log-log scale between W, Sinkhorn with the geodesic distance,  $SSW_2$  with the binary search (BS) and uniform distribution (12) and  $SSW_1$  with formula (11) between two distributions on  $S^2$ . The time includes the calculation of the distance matrices.

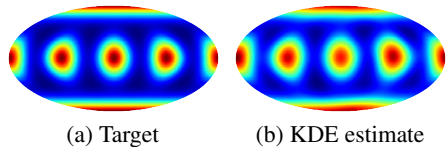


Figure 3: Minimization of SSW with respect to a mixture of vMF.

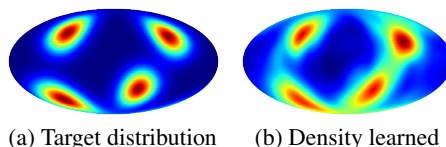


Figure 4: Amortized SSWVI with a normalizing flow *w.r.t.* a mixture of vMF.

## 254 5 Applications

255 In this section, we first illustrate the ability to approximate different distributions by minimizing SSW  
 256 *w.r.t.* some target distributions on  $S^2$ . We first use distributions from which we can draw samples.  
 257 Then, we use target distributions from which we know the density only up to a constant. Finally, we  
 258 apply SSW for generative modeling tasks using the framework of Sliced-Wasserstein autoencoder  
 259 and we show that we obtain competitive results with other Wasserstein autoencoder based methods  
 260 using a prior on the hypersphere. We also add in Appendix C.6 some experiments where we use SSW  
 261 in order to enforce uniformity in a contrastive self-supervised learning context.

### 262 5.1 SSW as a loss

263 We verify on the two first experiments that we can learn some target distribution  $\nu \in \mathcal{P}(S^{d-1})$  by  
 264 minimizing SSW, *i.e.* we consider the minimization problem  $\operatorname{argmin}_{\mu} SSW_p^p(\mu, \nu)$ .

265 **Gradient flow.** First, we suppose that we have access to the target distribution  $\nu$  through samples,  
 266 *i.e.* through  $\hat{\nu}_m = \frac{1}{m} \sum_{j=1}^m \delta_{y_j}$  where  $(y_j)_{j=1}^m$  are i.i.d samples of  $\nu$ . We choose as target distribution  
 267 a mixture of 6 well separated von Mises-Fisher distributions [72]. This is a fairly challenging distribu-  
 268 tion since there are 6 modes which are not connected. We show on Figure 3 the Mollweide projection  
 269 of the density approximated by a kernel density estimator for a distribution with 500 particles. To  
 270 optimize directly over particles, we can either perform a Riemannian gradient descent on the sphere  
 271 [3] or a projected gradient descent. We report in Appendix C.3 additional details and experiments.

272 **Sliced-Wasserstein variational inference on the sphere.** Another setting of interest is when we  
 273 have access to some target distribution only up to a constant. For example in Bayesian inference,  
 274 we want to sample from a posterior distribution  $p(\cdot|x)$  for which the normalizing constant is costly  
 275 to compute, *i.e.* we can only evaluate some function  $\pi$  such that  $p(\cdot|x) \propto \pi$ . Popular methods to  
 276 solve these types of problems are MCMCs [93] or variational inference [12, 54].

277 Variational inference aims at approximating the target by a distribution  $q$  in some family of distribu-  
 278 tions  $\mathcal{Q}$ . The classical way of doing it is to minimize the Kullback-Leibler (KL) divergence. However,  
 279 the KL divergence suffers from some drawbacks such as under estimating the target distribution and  
 280 not being a distance. Recently, Yi and Liu [111] proposed to use the SW distance instead. The method  
 281 is called Sliced-Wasserstein Variation Inference (SWVI) and relies on running at each iteration few  
 282 MCMC steps and then performing gradient descent to learn the variational distribution. We refer to  
 283 Appendix C.4 and Algorithm 2 for further details on the method.

284 In the following, we replace SW by SSW in SWVI, which we denote SSWVI, and we perform  
 285 amortized variational inference on the sphere by using exponential map normalizing flows (see [92]  
 286 and Appendix B.4) to learn the distribution and the Geodesic Langevin algorithm [105] as MCMC  
 287 method. We use the same target as Rezende et al. [92] and we report on Figure 4 the Mollweide  
 288 projection of the learned density. Since we learn to sample from a noise distribution, here the uniform  
 289 distribution on the sphere, we do not have directly access to the density and we report a kernel density  
 290 estimate with a Gaussian kernel. On Figure 5, we plot the evolution of the effective samples size  
 291 (ESS) [33, 69] through the iterations. This indicates how well the flow matches the target. We observe  
 292 that using SSW gives slightly better results, or at least comparable, than SWVI with SW.



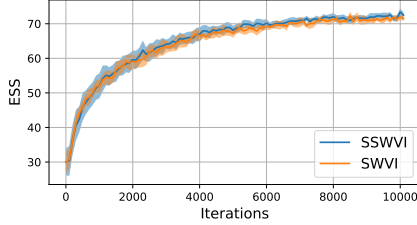


Figure 5: Comparison of the ESS between SWVI et SSWVI with the mixture target (mean and 95% confidence interval over 10 runs).

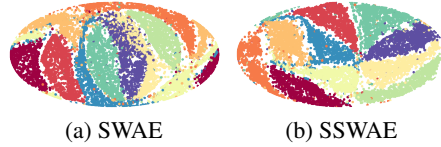


Figure 6: Latent space of SWAE and SSWAE for a uniform prior on  $S^2$ .

## 293 5.2 SSW autoencoders

294 In this section, we use SSW to learn the latent space of  
 295 autoencoders (AE). We rely on the SWAE framework in-  
 296 troduced by Kolouri et al. [59]. Let  $f$  be some encoder and  
 297  $g$  be some decoder, denote  $p_Z$  a prior distribution, then the  
 298 loss minimized in SWAE is

$$\mathcal{L}(f, g) = \int c(x, g(f(x))) d\mu(x) + \lambda SW_2^2(f_{\#}\mu, p_Z), \quad (24)$$

299 where  $\mu$  is the distribution of the data for which we have access to samples. One advantage of this  
 300 framework over more classical VAEs [58] is that no parametrization trick is needed here and therefore  
 301 the choice of the prior is more free.

302 In several concomitant works, it was shown that using a prior on the hypersphere can improve the  
 303 results [28, 110]. Hence, we propose in the same fashion as [59, 60, 84] to replace SW by SSW,  
 304 which we denote SSWAE, and to enforce a prior on the sphere. In the following, we use the MNIST  
 305 [64] and FashionMNIST [109] datasets, and we put an  $\ell^2$  normalization at the output of the encoder.  
 306 As a prior, we use the uniform distribution on  $S^{10}$  and we compare in Table 1 the Fréchet Inception  
 307 Distance (FID) [49], for 10000 samples and averaged over 5 trainings, obtained with the Wasserstein  
 308 Autoencoder (WAE) [99], the classical SWAE [59], the Sinkhorn Autoencoder (SAE) [84] and  
 309 circular GSWAE [60]. We observe that we obtain fairly competitive results. We add on Figure 6 the  
 310 latent space obtained with a uniform prior on  $S^2$ . We observe a better separation between classes for  
 311 SSWAE. We refer to appendix C.5 for more details and additional experiments.

## 312 6 Conclusion and discussion

313 In this work, we derive a new sliced-Wasserstein discrepancy on the hypersphere, that comes with  
 314 practical advantages when computing optimal transport distances on hyperspherical data. We notably  
 315 showed that it is competitive or even sometimes better than other metrics defined directly on  $\mathbb{R}^d$  on a  
 316 variety of machine learning tasks, including density estimation, variational inference or generative  
 317 models. Our work is, up to our knowledge, the first to adapt the sliced Wasserstein framework  
 318 to non-trivial manifolds. The three main ingredients are: *i*) a closed-form for Wasserstein on the  
 319 circle, *ii*) a closed-form solution to the projection onto great circles, and *iii*) a novel Radon transform  
 320 on the Sphere. An immediate extension of this work would be to consider sliced-Wasserstein  
 321 discrepancy in hyperbolic spaces, where geodesics are circular arcs as in the Poincaré disk. Beyond  
 322 the generalization to other, possibly well behaved, manifolds, statistical aspects need to be examined,  
 323 such as sample complexity or dependence to the hypersphere dimension. While we postulate that  
 324 results comparable to the Euclidean case might be reached, the fact that the manifold is closed might  
 325 bring interesting differences and justify further use of this type of discrepancies rather than their  
 326 Euclidean counterparts.

Table 1: FID (Lower is better).

Method / Dataset	MNIST	Fashion
SSWAE	<b>14.91 ± 0.32</b>	<b>43.94 ± 0.81</b>
SWAE	15.18 ± 0.32	44.78 ± 1.07
WAE-MMD IMQ	18.12 ± 0.62	68.51 ± 2.76
WAE-MMD RBF	20.09 ± 1.42	70.58 ± 1.75
SAE	19.39 ± 0.56	56.75 ± 1.7
Circular GSWAE	15.01 ± 0.26	44.65 ± 1.2

327 **References**

- 328 [1] Ahmed Abouelaz and Radouan Daher. Sur la transformation de radon de la sphère  $S^d$ . *Bulletin*  
329 *de la Société Mathématique de France*, 121(3):353–382, 1993. (Cited on p. 6)
- 330 [2] P-A Absil, Robert Mahony, and Rodolphe Sepulchre. Riemannian geometry of grassmann  
331 manifolds with a view on algorithmic computation. *Acta Applicandae Mathematica*, 80(2):  
332 199–220, 2004. (Cited on p. 5)
- 333 [3] P-A Absil, Robert Mahony, and Rodolphe Sepulchre. Optimization algorithms on matrix  
334 manifolds. In *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2009.  
335 (Cited on p. 8, 24, 27, 29, 30)
- 336 [4] Mark L Agranovskyt and Eric Todd Quintott. Injectivity of the spherical mean operator and  
337 related problems. *Complex analysis, harmonic analysis and applications*, 347:12, 1996. (Cited  
338 on p. 6)
- 339 [5] David Alvarez-Melis, Youssef Mroueh, and Tommi Jaakkola. Unsupervised hierarchy match-  
340 ing with optimal transport over hyperbolic spaces. In *International Conference on Artificial*  
341 *Intelligence and Statistics*, pages 1606–1617. PMLR, 2020. (Cited on p. 2)
- 342 [6] Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in*  
343 *the space of probability measures*. Springer Science & Business Media, 2005. (Cited on p. 3,  
344 6)
- 345 [7] Yuri A Antipov, Ricardo Estrada, and Boris Rubin. Inversion formulas for the spherical means  
346 in constant curvature spaces. *arXiv preprint arXiv:1107.5992*, 2011. (Cited on p. 6)
- 347 [8] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial  
348 networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.  
349 (Cited on p. 1)
- 350 [9] Eleonora Bardelli and Andrea Carlo Giuseppe Mennucci. Probability measures on infinite-  
351 dimensional stiefel manifolds. *Journal of Geometric Mechanics*, 9(3):291, 2017. (Cited on p.  
352 5, 18, 23)
- 353 [10] Thomas Bendokat, Ralf Zimmermann, and P-A Absil. A grassmann manifold handbook:  
354 Basic geometry and computational aspects. *arXiv preprint arXiv:2011.13699*, 2020. (Cited on  
355 p. 5)
- 356 [11] Camille Besombes, Olivier Pannekoucke, Corentin Lapeyre, Benjamin Sanderson, and Olivier  
357 Thuau. Producing realistic climate data with generative adversarial networks. *Nonlinear*  
358 *Processes in Geophysics*, 28(3):347–370, 2021. (Cited on p. 1)
- 359 [12] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for  
360 statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017. (Cited  
361 on p. 8, 28)
- 362 [13] Jan Boman and Filip Lindskog. Support theorems for the radon transform and cramér-wold  
363 theorems. *Journal of theoretical probability*, 22(3):683–710, 2009. (Cited on p. 3, 6)
- 364 [14] Silvere Bonnabel. Stochastic gradient descent on riemannian manifolds. *IEEE Transactions*  
365 *on Automatic Control*, 58(9):2217–2229, 2013. (Cited on p. 24)
- 366 [15] Nicolas Bonneel, Julien Rabin, Gabriel Peyré, and Hanspeter Pfister. Sliced and radon  
367 wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51(1):  
368 22–45, 2015. (Cited on p. 3, 6)
- 369 [16] Nicolas Bonnotte. *Unidimensional and evolution methods for optimal transportation*. PhD  
370 thesis, Paris 11, 2013. (Cited on p. 2)
- 371 [17] Nicolas Boumal. An introduction to optimization on smooth manifolds. To appear with  
372 Cambridge University Press, Apr 2022. URL <http://www.nicolasboumal.net/book>.  
373 (Cited on p. 24, 27)

- 374 [18] Simon Byrne and Mark Girolami. Geodesic monte carlo on embedded manifolds. *Scandinavian*  
375 *Journal of Statistics*, 40(4):825–845, 2013. (Cited on p. 29)
- 376 [19] Carlos A Cabrelli and Ursula M Molter. The kantorovich metric for probability measures on  
377 the circle. *Journal of Computational and Applied Mathematics*, 57(3):345–361, 1995. (Cited  
378 on p. 4)
- 379 [20] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Ar-  
380 mand Joulin. Unsupervised learning of visual features by contrasting cluster assign-  
381 ments. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Ad-*  
382 *vances in Neural Information Processing Systems*, volume 33, pages 9912–9924. Curran  
383 Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/file/  
384 70feb62b69f16e0238f741fab228fec2-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/70feb62b69f16e0238f741fab228fec2-Paper.pdf). (Cited on p. 1, 32)
- 385 [21] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
386 for contrastive learning of visual representations. In Hal Daumé III and Aarti Singh, editors,  
387 *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of  
388 *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 13–18 Jul 2020. URL  
389 <https://proceedings.mlr.press/v119/chen20j.html>. (Cited on p. 1, 32, 33, 34)
- 390 [22] Xiongjie Chen, Yongxin Yang, and Yunpeng Li. Augmented sliced wasserstein distances.  
391 *arXiv preprint arXiv:2006.08812*, 2020. (Cited on p. 3, 6)
- 392 [23] Samuel Cohen, Brandon Amos, and Yaron Lipman. Riemannian convex potential maps. In  
393 *International Conference on Machine Learning*, pages 2028–2038. PMLR, 2021. (Cited on p.  
394 2, 25)
- 395 [24] Nicolas Courty, Rémi Flamary, Devis Tuia, and Alain Rakotomamonjy. Optimal transport for  
396 domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):  
397 1853–1865, 2016. (Cited on p. 1)
- 398 [25] Li Cui, Xin Qi, Chengfeng Wen, Na Lei, Xinyuan Li, Min Zhang, and Xianfeng Gu. Spherical  
399 optimal transportation. *Computer-Aided Design*, 115:181–193, 2019. (Cited on p. 2)
- 400 [26] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in*  
401 *neural information processing systems*, 26, 2013. (Cited on p. 1, 7)
- 402 [27] Susanna Dann. On the minkowski-funk transform. *arXiv preprint arXiv:1003.5565*, 2010.  
403 (Cited on p. 6)
- 404 [28] Tim R. Davidson, Luca Falorsi, Nicola De Cao, Thomas Kipf, and Jakub M. Tomczak.  
405 Hyperspherical variational auto-encoders. In Amir Globerson and Ricardo Silva, editors,  
406 *Proceedings of the Thirty-Fourth Conference on Uncertainty in Artificial Intelligence, UAI*  
407 *2018, Monterey, California, USA, August 6-10, 2018*, pages 856–865. AUAI Press, 2018. URL  
408 <http://auai.org/uai2018/proceedings/papers/309.pdf>. (Cited on p. 1, 9, 25)
- 409 [29] Nicola De Cao and Wilker Aziz. The power spherical distribution. *arXiv preprint*  
410 *arXiv:2006.04437*, 2020. (Cited on p. 30)
- 411 [30] Julie Delon, Julien Salomon, and Andrei Sobolevski. Fast transport optimization for monge  
412 costs on the circle. *SIAM Journal on Applied Mathematics*, 70(7):2239–2258, 2010. (Cited on  
413 p. 3, 4, 6, 7)
- 414 [31] Ishan Deshpande, Yuan-Ting Hu, Ruoyu Sun, Ayis Pyrros, Nasir Siddiqui, Sanmi Koyejo,  
415 Zhizhen Zhao, David Forsyth, and Alexander G Schwing. Max-sliced wasserstein distance  
416 and its use for gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
417 *Pattern Recognition*, pages 10648–10656, 2019. (Cited on p. 3)
- 418 [32] Marco Di Marzio, Agnese Panzera, and Charles C Taylor. Nonparametric regression for  
419 spherical data. *Journal of the American Statistical Association*, 109(506):748–763, 2014.  
420 (Cited on p. 1)
- 421 [33] Arnaud Doucet, Nando De Freitas, Neil James Gordon, et al. *Sequential Monte Carlo methods*  
422 *in practice*, volume 1. Springer, 2001. (Cited on p. 8, 31)

- 423 [34] Leon Ehrenpreis. *The universality of the Radon transform*. OUP Oxford, 2003. (Cited on p. 6)
- 424 [35] Kilian Fatras, Younes Zine, Rémi Flamary, Remi Gribonval, and Nicolas Courty. Learning with  
425 minibatch wasserstein : asymptotic and gradient properties. In Silvia Chiappa and Roberto  
426 Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial  
427 Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages  
428 2131–2141. PMLR, 26–28 Aug 2020. URL [https://proceedings.mlr.press/v108/  
429 fatras20a.html](https://proceedings.mlr.press/v108/fatras20a.html). (Cited on p. 1)
- 430 [36] Charles Fefferman, Sanjoy Mitter, and Hariharan Narayanan. Testing the manifold hypothesis.  
431 *Journal of the American Mathematical Society*, 29(4):983–1049, 2016. (Cited on p. 1)
- 432 [37] Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trounev, and  
433 Gabriel Peyré. Interpolating between optimal transport and mmd using sinkhorn divergences.  
434 In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2681–  
435 2690, 2019. (Cited on p. 31)
- 436 [38] Alessio Figalli and Cédric Villani. Optimal transport and curvature. In *Nonlinear PDE’s and  
437 Applications*, pages 171–217. Springer, 2011. (Cited on p. 2)
- 438 [39] Rémi Flamary, Nicolas Courty, Alexandre Gramfort, Mokhtar Z. Alaya, Aurélie Boisbunon,  
439 Stanislas Chambon, Laetitia Chapel, Adrien Corenflos, Kilian Fatras, Nemo Fournier, Léo  
440 Gautheron, Nathalie T.H. Gayraud, Hicham Janati, Alain Rakotomamonjy, Ievgen Redko,  
441 Antoine Rolet, Antony Schutz, Vivien Seguy, Danica J. Sutherland, Romain Tavenard, Alexan-  
442 der Tong, and Titouan Vayer. Pot: Python optimal transport. *Journal of Machine Learning  
443 Research*, 22(78):1–8, 2021. URL <http://jmlr.org/papers/v22/20-451.html>. (Cited  
444 on p. 7, 31)
- 445 [40] P Thomas Fletcher, Conglin Lu, Stephen M Pizer, and Sarang Joshi. Principal geodesic  
446 analysis for the study of nonlinear statistics of shape. *IEEE transactions on medical imaging*,  
447 23(8):995–1005, 2004. (Cited on p. 4)
- 448 [41] Walter Gander. Algorithms for the qr decomposition. *Res. Rep*, 80(02):1251–1268, 1980.  
449 (Cited on p. 7)
- 450 [42] Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013. (Cited on p.  
451 7)
- 452 [43] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena  
453 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,  
454 Bilal Piot, koray kavukcuoglu, Remi Munos, and Michal Valko. Bootstrap your own latent -  
455 a new approach to self-supervised learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F.  
456 Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,  
457 pages 21271–21284. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.  
458 cc/paper/2020/file/f3ada80d5c4ee70142b17b8192b2958e-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/f3ada80d5c4ee70142b17b8192b2958e-Paper.pdf). (Cited on p.  
459 1)
- 460 [44] Armelle Guillou, Philippe Naveau, and Alexandre You. A folding methodology for multi-  
461 variate extremes: estimation of the spectral probability measure and actuarial applications.  
462 *Scandinavian Actuarial Journal*, 2015(7):549–572, 2015. (Cited on p. 1)
- 463 [45] Brittany Froese Hamfeldt and Axel GR Turnquist. A convergence framework for optimal  
464 transport on the sphere. *arXiv preprint arXiv:2103.05739*, 2021. (Cited on p. 2)
- 465 [46] Brittany Froese Hamfeldt and Axel GR Turnquist. A convergent finite difference method for  
466 optimal transport on the sphere. *arXiv preprint arXiv:2105.03500*, 2021. (Cited on p. 2)
- 467 [47] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
468 recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,  
469 pages 770–778, 2016. (Cited on p. 33)
- 470 [48] Sigurdur Helgason et al. *Integral geometry and Radon transforms*. Springer, 2011. (Cited on  
471 p. 3, 5)

- 472 [49] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
473 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances*  
474 *in neural information processing systems*, 30, 2017. (Cited on p. 9)
- 475 [50] Andrew Homan and Hanming Zhou. Injectivity and stability for a generic class of generalized  
476 radon transforms. *The Journal of Geometric Analysis*, 27(2):1515–1529, 2017. (Cited on p. 6)
- 477 [51] Andrés Hoyos-Idrobo. Aligning hyperbolic representations: an optimal transport-based  
478 approach. *arXiv preprint arXiv:2012.01089*, 2020. (Cited on p. 2)
- 479 [52] Minhui Huang, Shiqian Ma, and Lifeng Lai. A riemannian block coordinate descent method  
480 for computing the projection robust wasserstein distance. In *International Conference on*  
481 *Machine Learning*, pages 4446–4455. PMLR, 2021. (Cited on p. 3)
- 482 [53] Shayan Hundrieser, Marcel Klatt, and Axel Munk. The statistics of circular optimal transport.  
483 *arXiv preprint arXiv:2103.15426*, 2021. (Cited on p. 4)
- 484 [54] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An  
485 introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233,  
486 1999. (Cited on p. 8, 28)
- 487 [55] Sungkyu Jung. Geodesic projection of the von mises–fisher distribution for projection pursuit  
488 of directional data. *Electronic Journal of Statistics*, 15(1):984–1033, 2021. (Cited on p. 4, 6,  
489 24)
- 490 [56] Sungkyu Jung, Ian L Dryden, and James Stephen Marron. Analysis of principal nested spheres.  
491 *Biometrika*, 99(3):551–568, 2012. (Cited on p. 5, 20)
- 492 [57] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv*  
493 *preprint arXiv:1412.6980*, 2014. (Cited on p. 27, 30, 31, 33)
- 494 [58] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
495 *arXiv:1312.6114*, 2013. (Cited on p. 1, 9)
- 496 [59] Soheil Kolouri, Phillip E Pope, Charles E Martin, and Gustavo K Rohde. Sliced wasserstein  
497 auto-encoders. In *International Conference on Learning Representations*, 2018. (Cited on p.  
498 9, 31)
- 499 [60] Soheil Kolouri, Kimia Nadjahi, Umut Simsekli, Roland Badeau, and Gustavo Rohde. General-  
500 ized sliced wasserstein distances. *Advances in Neural Information Processing Systems*, 32,  
501 2019. (Cited on p. 3, 6, 9, 19, 31)
- 502 [61] Alex Krizhevsky. Learning multiple layers of features from tiny images, 2009. (Cited on p.  
503 33)
- 504 [62] Gerhard Kurz and Uwe D Hanebeck. Stochastic sampling of the hyperspherical von mises–  
505 fisher distribution without rejection methods. In *2015 Sensor Data Fusion: Trends, Solutions,*  
506 *Applications (SDF)*, pages 1–6. IEEE, 2015. (Cited on p. 24)
- 507 [63] Shiwei Lan, Bo Zhou, and Babak Shahbaba. Spherical hamiltonian monte carlo for constrained  
508 target distributions. In *International Conference on Machine Learning*, pages 629–637. PMLR,  
509 2014. (Cited on p. 29)
- 510 [64] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>. (Cited on p. 9)
- 511  
512 [65] Mufan Bill Li and Murat A Erdogdu. Riemannian langevin algorithm for solving semidefinite  
513 programs. *arXiv preprint arXiv:2010.11176*, 2020. (Cited on p. 29)
- 514 [66] Tianyi Lin, Chenyou Fan, Nhat Ho, Marco Cuturi, and Michael Jordan. Projection robust  
515 wasserstein distance and riemannian optimization. *Advances in neural information processing*  
516 *systems*, 33:9383–9397, 2020. (Cited on p. 3, 5)
- 517 [67] Tianyi Lin, Zeyu Zheng, Elynn Chen, Marco Cuturi, and Michael I Jordan. On projection  
518 robust optimal transport: Sample complexity and model misspecification. In *International*  
519 *Conference on Artificial Intelligence and Statistics*, pages 262–270. PMLR, 2021. (Cited on p.  
520 3, 6)

- 521 [68] Chang Liu, Jun Zhu, and Yang Song. Stochastic gradient geodesic mcmc methods. *Advances*  
522 *in neural information processing systems*, 29, 2016. (Cited on p. 29, 30)
- 523 [69] Jun S Liu and Rong Chen. Blind deconvolution via sequential imputations. *Journal of the*  
524 *american statistical association*, 90(430):567–576, 1995. (Cited on p. 8, 31)
- 525 [70] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Spheraface:  
526 Deep hypersphere embedding for face recognition. In *The IEEE Conference on Computer*  
527 *Vision and Pattern Recognition (CVPR)*, 2017. (Cited on p. 1)
- 528 [71] Kanti V Mardia, Peter E Jupp, and KV Mardia. *Directional statistics*, volume 2. Wiley Online  
529 Library, 2000. (Cited on p. 1, 24)
- 530 [72] Kantilal Varichand Mardia. Statistics of directional data. *Journal of the Royal Statistical*  
531 *Society: Series B (Methodological)*, 37(3):349–371, 1975. (Cited on p. 8)
- 532 [73] Robert J McCann. Polar factorization of maps on riemannian manifolds. *Geometric &*  
533 *Functional Analysis GAFA*, 11(3):589–608, 2001. (Cited on p. 2)
- 534 [74] Kimia Nadjahi. *Sliced-Wasserstein distance for large-scale machine learning: theory, method-*  
535 *ology and extensions*. PhD thesis, Institut polytechnique de Paris, 2021. (Cited on p. 2)
- 536 [75] Kimia Nadjahi, Alain Durmus, Umut Simsekli, and Roland Badeau. Asymptotic guarantees  
537 for learning generative models with the sliced-wasserstein distance. *Advances in Neural*  
538 *Information Processing Systems*, 32, 2019. (Cited on p. 3)
- 539 [76] Kimia Nadjahi, Alain Durmus, Lénaïc Chizat, Soheil Kolouri, Shahin Shahrampour, and Umut  
540 Simsekli. Statistical and topological properties of sliced probability divergences. *Advances in*  
541 *Neural Information Processing Systems*, 33:20802–20812, 2020. (Cited on p. 26)
- 542 [77] Frank Natterer. *The mathematics of computerized tomography*. SIAM, 2001. (Cited on p. 5)
- 543 [78] Khai Nguyen and Nhat Ho. Amortized projection optimization for sliced wasserstein generative  
544 models. *arXiv preprint arXiv:2203.13417*, 2022. (Cited on p. 3)
- 545 [79] Khai Nguyen and Nhat Ho. Revisiting sliced wasserstein on images: From vectorization to  
546 convolution. *arXiv preprint arXiv:2204.01188*, 2022. (Cited on p. 3)
- 547 [80] Khai Nguyen, Son Nguyen, Nhat Ho, Tung Pham, and Hung Bui. Improving relational  
548 regularized autoencoders with spherical sliced fused gromov wasserstein. *arXiv preprint*  
549 *arXiv:2010.01787*, 2020. (Cited on p. 3)
- 550 [81] Khai Nguyen, Nhat Ho, Tung Pham, and Hung Bui. Distributional sliced-wasserstein and  
551 applications to generative modeling. In *9th International Conference on Learning Representations,*  
552 *ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL  
553 <https://openreview.net/forum?id=QYj07OACDK>. (Cited on p. 3)
- 554 [82] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji  
555 Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of*  
556 *Machine Learning Research*, 22(57):1–64, 2021. (Cited on p. 25)
- 557 [83] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
558 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas  
559 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,  
560 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style,  
561 high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer,  
562 F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing*  
563 *Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. (Cited on p. 7)
- 564 [84] Giorgio Patrini, Rianne van den Berg, Patrick Forre, Marcello Carioni, Samarth Bhargav,  
565 Max Welling, Tim Genewein, and Frank Nielsen. Sinkhorn autoencoders. In *Uncertainty in*  
566 *Artificial Intelligence*, pages 733–743. PMLR, 2020. (Cited on p. 9, 31)
- 567 [85] François-Pierre Paty and Marco Cuturi. Subspace robust wasserstein distances. In *International*  
568 *conference on machine learning*, pages 5072–5081. PMLR, 2019. (Cited on p. 3)

- 569 [86] Nathanaël Perraudin, Michaël Defferrard, Tomasz Kacprzak, and Raphael Sgier. DeepSphere:  
570 Efficient spherical convolutional neural network with healpix sampling for cosmological  
571 applications. *Astronomy and Computing*, 27:130–146, 2019. (Cited on p. 1)
- 572 [87] Arthur Pewsey and Eduardo García-Portugués. Recent advances in directional statistics. *Test*,  
573 30(1):1–58, 2021. (Cited on p. 1)
- 574 [88] Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data  
575 science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019. (Cited on p.  
576 2, 4)
- 577 [89] Julien Rabin, Julie Delon, and Yann Gousseau. Transportation distances on the circle. *Journal*  
578 *of Mathematical Imaging and Vision*, 41(1):147–167, 2011. (Cited on p. 2, 3, 4, 17)
- 579 [90] Julien Rabin, Gabriel Peyré, Julie Delon, and Marc Bernot. Wasserstein barycenter and its  
580 application to texture mixing. In *International Conference on Scale Space and Variational*  
581 *Methods in Computer Vision*, pages 435–446. Springer, 2011. (Cited on p. 1, 2)
- 582 [91] Danilo J Rezende and Sébastien Racanière. Implicit riemannian concave potential maps. *arXiv*  
583 *preprint arXiv:2110.01288*, 2021. (Cited on p. 2, 25)
- 584 [92] Danilo Jimenez Rezende, George Papamakarios, Sébastien Racanière, Michael Albergo, Gurtej  
585 Kanwar, Phiala Shanahan, and Kyle Cranmer. Normalizing flows on tori and spheres. In  
586 *International Conference on Machine Learning*, pages 8083–8092. PMLR, 2020. (Cited on p.  
587 8, 25, 30, 31)
- 588 [93] Christian P Robert, George Casella, and George Casella. *Monte Carlo statistical methods*,  
589 volume 2. Springer, 1999. (Cited on p. 8)
- 590 [94] Boris Rubin. Inversion and characterization of the hemispherical transform. *Journal d'Analyse*  
591 *Mathématique*, 77(1):105–128, 1999. (Cited on p. 6, 21)
- 592 [95] Boris Rubin. Inversion formulas for the spherical radon transform and the generalized cosine  
593 transform. *Advances in Applied Mathematics*, 29(3):471–497, 2002. (Cited on p. 6)
- 594 [96] Boris Rubin. Notes on radon transforms in integral geometry. *Fractional Calculus and Applied*  
595 *Analysis*, 6(1):25–72, 2003. (Cited on p. 6, 19, 20)
- 596 [97] Meyer Scetbon, Marco Cuturi, and Gabriel Peyré. Low-rank sinkhorn factorization. In  
597 *International Conference on Machine Learning*, pages 9344–9354. PMLR, 2021. (Cited on p.  
598 1)
- 599 [98] Suvrit Sra. Directional statistics in machine learning: a brief review. *Applied Directional*  
600 *Statistics: Modern Methods and Case Studies*, 225:6, 2018. (Cited on p. 1)
- 601 [99] Ilya O. Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schölkopf. Wasserstein auto-  
602 encoders. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver,*  
603 *BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.  
604 URL <https://openreview.net/forum?id=HkL7n1-0b>. (Cited on p. 9, 31)
- 605 [100] Gary Ulrich. Computer generation of distributions on the m-sphere. *Journal of the Royal*  
606 *Statistical Society: Series C (Applied Statistics)*, 33(2):158–163, 1984. (Cited on p. 24)
- 607 [101] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer, 2009. (Cited on p. 1,  
608 2, 23)
- 609 [102] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David  
610 Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J.  
611 van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew  
612 R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W.  
613 Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A.  
614 Quintero, Charles R. Harris, Anne M. Archibald, António H. Ribeiro, Fabian Pedregosa, Paul  
615 van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific  
616 Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.  
617 (Cited on p. 30)

- 618 [103] Dilin Wang and Qiang Liu. Learning to draw samples: With application to amortized mle for  
619 generative adversarial learning. *arXiv preprint arXiv:1611.01722*, 2016. (Cited on p. 29)
- 620 [104] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through  
621 alignment and uniformity on the hypersphere. In *International Conference on Machine*  
622 *Learning*, pages 9929–9939. PMLR, 2020. (Cited on p. 1, 32, 33, 34)
- 623 [105] Xiao Wang, Qi Lei, and Ioannis Panageas. Fast convergence of langevin dynamics on manifold:  
624 Geodesics meet log-sobolev. *Advances in Neural Information Processing Systems*, 33:18894–  
625 18904, 2020. (Cited on p. 8, 29)
- 626 [106] Michael Werman, Shmuel Peleg, and Azriel Rosenfeld. A distance metric for multidimensional  
627 histograms. *Computer Vision, Graphics, and Image Processing*, 32(3):328–336, 1985. (Cited  
628 on p. 4)
- 629 [107] Andrew TA Wood. Simulation of the von mises fisher distribution. *Communications in*  
630 *statistics-simulation and computation*, 23(1):157–164, 1994. (Cited on p. 24)
- 631 [108] Zhirong Wu, Yuanjun Xiong, Stella X. Yu, and Dahua Lin. Unsupervised feature learning via  
632 non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer*  
633 *Vision and Pattern Recognition (CVPR)*, June 2018. (Cited on p. 1, 32)
- 634 [109] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for  
635 benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017. (Cited  
636 on p. 9)
- 637 [110] Jiacheng Xu and Greg Durrett. Spherical latent spaces for stable variational autoencoders. In  
638 Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii, editors, *Proceedings of the*  
639 *2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium,*  
640 *October 31 - November 4, 2018*, pages 4503–4513. Association for Computational Linguistics,  
641 2018. (Cited on p. 1, 9, 25)
- 642 [111] Mingxuan Yi and Song Liu. Sliced wasserstein variational inference. In *Fourth Symposium on*  
643 *Advances in Approximate Bayesian Inference*, 2021. (Cited on p. 8, 29)



644 **A Proofs**

645 **A.1 Proof of Proposition 1**

646 **Optimal  $\alpha$ .** Let  $\mu \in \mathcal{P}_2(S^1)$ ,  $\nu = \text{Unif}(S^1)$ . Since  $\nu$  is the uniform distribution on  $S^1$ , its cdf is  
 647 the identity on  $[0, 1]$  (where we identified  $S^1$  and  $[0, 1]$ ). We can extend the cdf  $F$  on the real line as  
 648 in [89] with the convention  $F(y + 1) = F(y) + 1$ . Therefore,  $F_\nu = \text{Id}$  on  $\mathbb{R}$ . Moreover, we know  
 649 that for all  $x \in S^1$ ,  $(F_\nu - \alpha)^{-1}(x) = F_\nu^{-1}(x + \alpha) = x + \alpha$  and

$$W_2^2(\mu, \nu) = \inf_{\alpha \in \mathbb{R}} \int_0^1 |F_\mu^{-1}(t) - (F_\nu - \alpha)^{-1}(t)|^2 dt. \quad (25)$$

650 For all  $\alpha \in \mathbb{R}$ , let  $f(\alpha) = \int_0^1 (F_\mu^{-1}(t) - (F_\nu - \alpha)^{-1}(t))^2 dt$ . Then, we have:

$$\begin{aligned} \forall \alpha \in \mathbb{R}, f(\alpha) &= \int_0^1 (F_\mu^{-1}(t) - t - \alpha)^2 dt \\ &= \int_0^1 (F_\mu^{-1}(t) - t)^2 dt + \alpha^2 - 2\alpha \int_0^1 (F_\mu^{-1}(t) - t) dt \\ &= \int_0^1 (F_\mu^{-1}(t) - t)^2 dt + \alpha^2 - 2\alpha \left( \int_0^1 x d\mu(x) - \frac{1}{2} \right), \end{aligned} \quad (26)$$

651 where we used that  $(F_\mu^{-1})_\# \text{Unif}([0, 1]) = \mu$ .

652 Hence,  $f'(\alpha) = 0 \iff \alpha = \int_0^1 x d\mu(x) - \frac{1}{2}$ .

653 **Closed-form for empirical distributions.** Let  $(x_i)_{i=1}^n \in [0, 1]^n$  such that  $x_1 < \dots < x_n$  and let  
 654  $\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  a discrete distribution.

655 To compute the closed-form of  $W_2$  between  $\mu_n$  and  $\nu = \text{Unif}(S^1)$ , we first have that the optimal  $\alpha$   
 656 is  $\alpha_n = \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{2}$ . Moreover, we also have:

$$\begin{aligned} W_2^2(\mu_n, \nu) &= \int_0^1 (F_{\mu_n}^{-1}(t) - (t + \hat{\alpha}_n))^2 dt \\ &= \int_0^1 F_{\mu_n}^{-1}(t)^2 dt - 2 \int_0^1 t F_{\mu_n}^{-1}(t) dt - 2\hat{\alpha}_n \int_0^1 F_{\mu_n}^{-1}(t) dt + \frac{1}{3} + \hat{\alpha}_n + \hat{\alpha}_n^2. \end{aligned} \quad (27)$$

657 Then, by noticing that  $F_{\mu_n}^{-1}(t) = x_i$  for all  $t \in [F(x_i), F(x_{i+1})]$ , we have

$$\int_0^1 t F_{\mu_n}^{-1}(t) dt = \sum_{i=1}^n \int_{\frac{i-1}{n}}^{\frac{i}{n}} t x_i dt = \frac{1}{2n^2} \sum_{i=1}^n x_i (2i - 1), \quad (28)$$

658

$$\int_0^1 F_{\mu_n}^{-1}(t)^2 dt = \frac{1}{n} \sum_{i=1}^n x_i^2, \quad \int_0^1 F_{\mu_n}^{-1}(t) dt = \frac{1}{n} \sum_{i=1}^n x_i, \quad (29)$$

659 and we also have:

$$\hat{\alpha}_n + \hat{\alpha}_n^2 = \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{2} + \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 + \frac{1}{4} - \frac{1}{n} \sum_{i=1}^n x_i = \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 - \frac{1}{4}. \quad (30)$$

660 Then, by plugging these results into (27), we obtain

$$\begin{aligned} W_2^2(\mu_n, \nu) &= \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{1}{n^2} \sum_{i=1}^n (2i - 1)x_i - 2 \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 + \frac{1}{n} \sum_{i=1}^n x_i + \frac{1}{3} + \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 - \frac{1}{4} \\ &= \frac{1}{n} \sum_{i=1}^n x_i^2 - \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 + \frac{1}{n^2} \sum_{i=1}^n (n + 1 - 2i)x_i + \frac{1}{12}. \end{aligned} \quad (31)$$

661 **A.2 Proof of Equation (17)**

662 Let  $U \in \mathbb{V}_{d,2}$ . Then the great circle generated by  $U \in \mathbb{V}_{d,2}$  is defined as the intersection between  
 663  $\text{span}(UU^T)$  and  $S^{d-1}$ . And we have the following characterization:

$$\begin{aligned} x \in \text{span}(UU^T) \cap S^{d-1} &\iff \exists y \in \mathbb{R}^d, x = UU^T y \text{ and } \|x\|_2^2 = 1 \\ &\iff \exists y \in \mathbb{R}^d, x = UU^T y \text{ and } \|UU^T y\|_2^2 = y^T UU^T y = \|U^T y\|_2^2 = 1 \\ &\iff \exists z \in S^1, x = Uz. \end{aligned}$$

664 And we deduce that

$$\forall U \in \mathbb{V}_{d,2}, x \in S^{d-1}, P^U(x) = \underset{z \in S^1}{\text{argmin}} d_{S^{d-1}}(x, Uz). \quad (32)$$

665 **A.3 Proof of Lemma 1**

666 Let  $U \in \mathbb{V}_{d,2}$  and  $x \in S^{d-1}$  such that  $U^T x \neq 0$ . Denote  $U = (u_1 \ u_2)$ , i.e. the 2-plane  $E$   
 667 is  $E = \text{span}(UU^T) = \text{span}(u_1, u_2)$  and  $(u_1, u_2)$  is an orthonormal basis of  $E$ . Then, for all  
 668  $x \in S^{d-1}$ , the projection on  $E$  is  $p^E(x) = \langle u_1, x \rangle u_1 + \langle u_2, x \rangle u_2 = UU^T x$ .

669 Now, let us compute the geodesic distance between  $x \in S^{d-1}$  and  $\frac{p^E(x)}{\|p^E(x)\|_2} \in E \cap S^{d-1}$ :

$$d_{S^{d-1}} \left( x, \frac{p^E(x)}{\|p^E(x)\|_2} \right) = \arccos \left( \left\langle x, \frac{p^E(x)}{\|p^E(x)\|_2} \right\rangle \right) = \arccos(\|p^E(x)\|_2), \quad (33)$$

670 using that  $x = p^E(x) + p^{E^\perp}(x)$ .

671 Let  $y \in E \cap S^{d-1}$  another point on the great circle. By the Cauchy-Schwarz inequality, we have

$$\langle x, y \rangle = \langle p^E(x), y \rangle \leq \|p^E(x)\|_2 \|y\|_2 = \|p^E(x)\|_2. \quad (34)$$

672 Therefore, using that  $\arccos$  is decreasing on  $(-1, 1)$ ,

$$d_{S^{d-1}}(x, y) = \arccos(\langle x, y \rangle) \geq \arccos(\|p^E(x)\|_2) = d_{S^{d-1}} \left( x, \frac{p^E(x)}{\|p^E(x)\|_2} \right). \quad (35)$$

673 Moreover, we have equality if and only if  $y = \lambda p^E(x)$ . And since  $y \in S^{d-1}$ ,  $|\lambda| = \frac{1}{\|p^E(x)\|_2}$ . Using  
 674 again that  $\arccos$  is decreasing, we deduce that the minimum is well attained in  $y = \frac{p^E(x)}{\|p^E(x)\|_2} =$   
 675  $\frac{UU^T x}{\|UU^T x\|_2}$ .

676 Finally, using that  $\|UU^T x\|_2 = x^T UU^T UU^T x = x^T UU^T x = \|U^T x\|_2$ , we deduce that

$$P^U(x) = \frac{U^T x}{\|U^T x\|_2}. \quad (36)$$

677 Finally, by noticing that the projection is unique if and only if  $U^T x = 0$ , and using [9, Proposition  
 678 4.2] which states that there is a unique projection for a.e.  $x$ , we deduce that  $\{x \in S^{d-1}, U^T x = 0\}$   
 679 is of measure null and hence, for a.e.  $x \in S^{d-1}$ , we have the result.

680 **A.4 Proof of Proposition 2**

681 Let  $f \in L^1(S^{d-1})$ ,  $g \in C_0(S^1 \times \mathbb{V}_{d,2})$ , then by Fubini's theorem,

$$\begin{aligned}
\langle \tilde{R}f, g \rangle_{S^1 \times \mathbb{V}_{d,2}} &= \int_{\mathbb{V}_{d,2}} \int_{S^1} \tilde{R}f(z, U) g(z, U) \, dz d\sigma(U) \\
&= \int_{\mathbb{V}_{d,2}} \int_{S^1} \int_{S^{d-1}} f(x) \mathbb{1}_{\{z=PU(x)\}} g(z, U) \, dx dz d\sigma(U) \\
&= \int_{S^{d-1}} f(x) \int_{\mathbb{V}_{d,2}} \int_{S^1} g(z, U) \mathbb{1}_{\{z=PU(x)\}} \, dz d\sigma(U) dx \\
&= \int_{S^{d-1}} f(x) \int_{\mathbb{V}_{d,2}} g(P^U(x), U) \, d\sigma(U) dx \\
&= \int_{S^{d-1}} f(x) \tilde{R}^* g(x) \, dx \\
&= \langle f, \tilde{R}^* g \rangle_{S^{d-1}}.
\end{aligned} \tag{37}$$

682 **A.5 Proof of Proposition 3**

683 Let  $g \in C_0(S^1 \times \mathbb{V}_{d,2})$ ,

$$\begin{aligned}
\int_{\mathbb{V}_{d,2}} \int_{S^1} g(z, U) (\tilde{R}\mu)^U(dz) \, d\sigma(U) &= \int_{S^1 \times \mathbb{V}_{d,2}} g(z, U) \, d(\tilde{R}\mu)(z, U) \\
&= \int_{S^{d-1}} \tilde{R}^* g(x) \, d\mu(x) \\
&= \int_{S^{d-1}} \int_{\mathbb{V}_{d,2}} g(P^U(x), U) \, d\sigma(U) d\mu(x) \\
&= \int_{\mathbb{V}_{d,2}} \int_{S^{d-1}} g(P^U(x), U) \, d\mu(x) d\sigma(U) \\
&= \int_{\mathbb{V}_{d,2}} \int_{S^1} g(z, U) \, d(P_{\#}^U \mu)(z) d\sigma(U).
\end{aligned} \tag{38}$$

684 Hence, for  $\sigma$ -almost every  $U \in \mathbb{V}_{d,2}$ ,  $(\tilde{R}\mu)^U = P_{\#}^U \mu$ .

685 **A.6 Study of the Spherical Radon transform  $\tilde{R}$**

686 In this Section, we first discuss the set of integration of the spherical Radon transform  $\tilde{R}$  (19). We  
687 further show that it is related to the hemispherical Radon transform and we derive its kernel.

688 **Set of integration.** While the classical Radon transform integrates over hyperplanes of  $\mathbb{R}^d$  and the  
689 generalized Radon transform integrates over hypersurfaces [60], the set of integration of the spherical  
690 Radon transform (19) is a half of a ‘‘big circle’’, *i.e.* half of the intersection between a hyperplane and  
691  $S^{d-1}$  [96]. We illustrate this on  $S^2$  in Figure 7. On  $S^2$ , the intersection between a hyperplane and  $S^2$   
692 is a great circle.

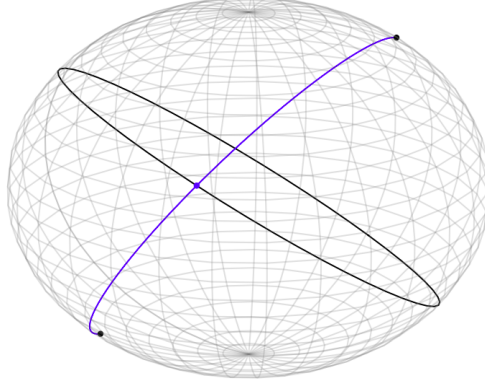


Figure 7: Set of integration of the spherical Radon transform (19). The great circle is in black and the set of integration in blue. The point  $Uz \in \text{span}(UU^T) \cap S^{d-1}$  is in blue.

693 **Proposition 6.** Let  $U \in \mathbb{V}_{d,2}$ ,  $z \in S^1$ . The set of integration of (19) is

$$\{x \in S^{d-1}, P^U(x) = z\} = \{x \in F \cap S^{d-1}, \langle x, Uz \rangle > 0\}, \quad (39)$$

694 where  $F = \text{span}(UU^T)^\perp \oplus \text{span}(Uz)$ .

695 *Proof.* Let  $U \in \mathbb{V}_{d,2}$ ,  $z \in S^1$ . Denote  $E = \text{span}(UU^T)$  the 2-plane generating the great circle,  
 696 and  $E^\perp$  its orthogonal complementary. Hence,  $E \oplus E^\perp = \mathbb{R}^d$  and  $\dim(E^\perp) = d - 2$ . Now, let  
 697  $F = E^\perp \oplus \text{span}(Uz)$ . Since  $Uz = UU^T Uz \in E$ , we have that  $\dim(F) = d - 1$ . Hence,  $F$  is a  
 698 hyperplane and  $F \cap S^{d-1}$  is a “big circle” [96], i.e. a  $(d - 2)$ -dimensional subsphere of  $S^{d-1}$ .

699 Now, for the first inclusion, let  $x \in \{x \in S^{d-1}, P^U(x) = z\}$ . First, we show that  $x \in F \cap S^{d-1}$ . By  
 700 Lemma 1 and hypothesis, we know that  $P^U(x) = \frac{U^T x}{\|U^T x\|_2} = z$ . By denoting by  $p^E$  the projection on  
 701  $E$ , we have:

$$p^E(x) = UU^T x = U(\|U^T x\|_2 z) = \|U^T x\|_2 Uz \in \text{span}(Uz). \quad (40)$$

702 Hence,  $x = p^E(x) + x_{E^\perp} = \|U^T x\|_2 Uz + x_{E^\perp} \in F$ . Moreover, as

$$\langle x, Uz \rangle = \|U^T x\|_2 \langle Uz, Uz \rangle = \|U^T x\|_2 > 0, \quad (41)$$

703 we deduce that  $x \in \{F \cap S^{d-1}, \langle x, Uz \rangle > 0\}$ .

704 For the other inclusion, let  $x \in \{F \cap S^{d-1}, \langle x, Uz \rangle > 0\}$ . Since  $x \in F$ , we have  $x = x_{E^\perp} + \lambda Uz$ ,  
 705  $\lambda \in \mathbb{R}$ . Hence, using Lemma 1,

$$P^U(x) = \frac{U^T x}{\|U^T x\|_2} = \frac{\lambda}{|\lambda|} \frac{z}{\|z\|_2} = \text{sign}(\lambda)z. \quad (42)$$

706 But, we also have  $\langle x, Uz \rangle = \lambda \|Uz\|_2^2 = \lambda > 0$ . Therefore,  $\text{sign}(\lambda) = 1$  and  $P^U(x) = z$ .

707 Finally, we conclude that  $\{x \in S^{d-1}, P^U(x) = z\} = \{x \in F \cap S^{d-1}, \langle x, Uz \rangle > 0\}$ .  $\square$

708 **Link with Hemispherical transform.** Since the intersection between a hyperplane and  $S^{d-1}$  is  
 709 isometric to  $S^{d-2}$  [56], we can relate  $\tilde{R}$  to the hemispherical transform  $\mathcal{H}$  [96] on  $S^{d-2}$ . First, the  
 710 hemispherical transform of a function  $f \in L^1(S^{d-1})$  is defined as

$$\forall x \in S^{d-1}, \mathcal{H}f(x) = \int_{S^{d-1}} f(y) \mathbb{1}_{\{\langle x, y \rangle > 0\}} dy. \quad (43)$$

711 From Proposition 6, we can write the spherical Radon transform (19) as a hemispherical transform  
 712 on  $S^{d-2}$ .

713 **Proposition 7.** Let  $f \in L^1(S^{d-1})$ ,  $U \in \mathbb{V}_{d,2}$  and  $z \in S^1$ , then

$$\tilde{R}f(z, U) = \int_{S^{d-2}} \tilde{f}(x) \mathbb{1}_{\{\langle x, \tilde{U}z \rangle > 0\}} dx = \mathcal{H}\tilde{f}(\tilde{U}z), \quad (44)$$

714 where for all  $x \in S^{d-2}$ ,  $\tilde{f}(x) = f(O^T Jx)$  with  $O$  the rotation matrix such that for all  $x \in F$ ,  
 715  $Ox \in \text{span}(e_1, \dots, e_{d-1})$  where  $(e_1, \dots, e_d)$  denotes the canonical basis, and  $J = \begin{pmatrix} I_{d-1} \\ 0_{1,d-1} \end{pmatrix}$ , and  
 716  $\tilde{U} = J^T O U \in \mathbb{R}^{(d-1) \times 2}$ .

717 *Proof.* Let  $f \in L^1(S^{d-1})$ ,  $z \in S^1$ ,  $U \in \mathbb{V}_{d,2}$ , then by Proposition 6,

$$\tilde{R}f(z, U) = \int_{S^{d-1} \cap F} f(x) \mathbb{1}_{\{\langle x, Uz \rangle > 0\}} dx. \quad (45)$$

718  $F$  is a hyperplane. Let  $O \in \mathbb{R}^{d \times d}$  be the rotation such that for all  $x \in F$ ,  $Ox \in \text{span}(e_1, \dots, e_{d-1}) =$   
 719  $\tilde{F}$  where  $(e_1, \dots, e_d)$  is the canonical basis. By applying the change of variable  $Ox = y$ , and since  
 720  $O^{-1} = O^T$ ,  $\det O = 1$ , we obtain

$$\tilde{R}f(z, U) = \int_{O(F \cap S^{d-1})} f(O^T y) \mathbb{1}_{\{\langle O^T y, Uz \rangle > 0\}} dy = \int_{\tilde{F} \cap S^{d-1}} f(O^T y) \mathbb{1}_{\{\langle y, OUz \rangle > 0\}} dy. \quad (46)$$

721 Now, we have that  $OU \in \mathbb{V}_{d,2}$  since  $(OU)^T(OU) = I_2$ , and since  $Uz \in F$ ,  $OUz \in \tilde{F}$ . For all  
 722  $y \in \tilde{F}$ , we have  $\langle y, e_d \rangle = y_d = 0$ . Let  $J = \begin{pmatrix} I_{d-1} \\ 0_{1,d-1} \end{pmatrix} \in \mathbb{R}^{d \times (d-1)}$ , then for all  $y \in \tilde{F} \cap S^{d-1}$ ,  
 723  $y = J\tilde{y}$  where  $\tilde{y} \in S^{d-2}$  is composed of the  $d-1$  first coordinates of  $y$ .

724 Let's define, for all  $\tilde{y} \in S^{d-2}$ ,  $\tilde{f}(\tilde{y}) = f(O^T J\tilde{y})$ ,  $\tilde{U} = J^T O U$ .

725 Then, since  $\tilde{F} \cap S^{d-1} \cong S^{d-2}$ , we can write:

$$\tilde{R}f(z, U) = \int_{S^{d-2}} \tilde{f}(\tilde{y}) \mathbb{1}_{\{\langle \tilde{y}, \tilde{U}z \rangle > 0\}} d\tilde{y} = \mathcal{H}\tilde{f}(\tilde{U}z). \quad (47)$$

726 □

727 **Kernel of  $\tilde{R}$ .** By exploiting the expression using the hemispherical transform in Proposition 7, we  
 728 can derive its kernel in Appendix A.7.

## 729 A.7 Proof of Proposition 4

730 First, we recall Lemma 2.3 of [94] on  $S^{d-2}$ .

731 **Lemma 2** (Lemma 2.3 [94]).  $\ker(\mathcal{H}) = \{\mu \in \mathcal{M}_{\text{even}}(S^{d-2}), \mu(S^{d-2}) = 0\}$  where  $\mathcal{M}_{\text{even}}$  is  
 732 the set of even measures, i.e. measures such that for all  $f \in C(S^{d-2})$ ,  $\langle \mu, f \rangle = \langle \mu, f^- \rangle$  where  
 733  $f^-(x) = f(-x)$  for all  $x \in S^{d-2}$ .

734 Let  $\mu \in \mathcal{M}_{\text{ac}}(S^{d-1})$ . First, we notice that the density of  $\tilde{R}\mu$  w.r.t.  $\lambda \otimes \sigma$  is, for all  $z \in S^1$ ,  $U \in \mathbb{V}_{d,2}$ ,  
 735

$$(\tilde{R}\mu)(z, U) = \int_{S^{d-1}} \mathbb{1}_{\{P^U(x)=z\}} d\mu(x) = \int_{F \cap S^{d-1}} \mathbb{1}_{\{\langle x, Uz \rangle > 0\}} d\mu(x). \quad (48)$$

736 Indeed, using Proposition 2, and Proposition 6, we have for all  $g \in C_0(S^1 \times \mathbb{V}_{d,2})$ ,

$$\begin{aligned} \langle \tilde{R}\mu, g \rangle_{S^1 \times \mathbb{V}_{d,2}} &= \langle \mu, \tilde{R}^*g \rangle_{S^{d-1}} = \int_{S^{d-1}} R^*g(x) d\mu(x) \\ &= \int_{S^{d-1}} \int_{\mathbb{V}_{d,2}} \int_{S^1} g(z, U) \mathbb{1}_{\{z=P^U(x)\}} dz d\sigma(U) d\mu(x) \\ &= \int_{\mathbb{V}_{d,2} \times S^1} g(z, U) \int_{S^{d-1}} \mathbb{1}_{\{z=P^U(x)\}} d\mu(x) dz d\sigma(U) \\ &= \int_{\mathbb{V}_{d,2} \times S^1} g(z, U) \int_{F \cap S^{d-1}} \mathbb{1}_{\{\langle x, Uz \rangle > 0\}} d\mu(x) dz d\sigma(U). \end{aligned} \quad (49)$$

737 Hence, using Proposition 7, we can write  $(\tilde{R}\mu)(z, U) = (\mathcal{H}\tilde{\mu})(\tilde{U}z)$  where  $\tilde{\mu} = J_{\#}^T O_{\#}\mu$ .

738 Now, let  $\mu \in \ker(\tilde{R})$ , then for all  $z \in S^1$ ,  $U \in \mathbb{V}_{d,2}$ ,  $\tilde{R}\mu(z, U) = \mathcal{H}\tilde{\mu}(\tilde{U}z) = 0$  and hence  
 739  $\tilde{\mu} \in \ker(\mathcal{H}) = \{\tilde{\mu} \in \mathcal{M}_{\text{even}}(S^{d-2}), \tilde{\mu}(S^{d-2}) = 0\}$ .

740 First, let's show that  $\mu \in \mathcal{M}_{\text{even}}(S^{d-1})$ . Let  $f \in C(S^{d-1})$  and  $U \in \mathbb{V}_{d,2}$ , then, by using the same  
 741 notation as in Propositions 6 and 7, we have

$$\begin{aligned}
 \langle \mu, f \rangle_{S^{d-1}} &= \int_{S^{d-1}} f(x) d\mu(x) = \int_{S^{d-1}} \int_{S^1} f(x) \mathbb{1}_{\{z=P^U(x)\}} dz d\mu(x) \\
 &= \int_{S^1} \int_{S^{d-1}} f(x) \mathbb{1}_{\{z=P^U(x)\}} d\mu(x) dz \\
 &= \int_{S^1} \int_{F \cap S^{d-1}} f(x) \mathbb{1}_{\{(x,Uz) > 0\}} d\mu(x) dz \quad \text{by Prop. 6} \\
 &= \int_{S^1} \int_{S^{d-2}} \tilde{f}(y) \mathbb{1}_{\{(y,\tilde{U}z) > 0\}} d\tilde{\mu}(y) dz \\
 &= \int_{S^1} \langle \mathcal{H}\tilde{\mu}, \tilde{f} \rangle_{S^{d-2}} dz \\
 &= \int_{S^1} \langle \tilde{\mu}, \mathcal{H}\tilde{f} \rangle_{S^{d-2}} dz \\
 &= \int_{S^1} \langle \tilde{\mu}, (\mathcal{H}\tilde{f})^- \rangle_{S^{d-2}} dz \quad \text{since } \tilde{\mu} \in \mathcal{M}_{\text{even}} \\
 &= \int_{S^{d-1}} f^-(x) d\mu(x) = \langle \mu, f^- \rangle_{S^{d-1}},
 \end{aligned} \tag{50}$$

742 using for the last line all the opposite transformations. Therefore,  $\mu \in \mathcal{M}_{\text{even}}(S^{d-1})$ .

743 Now, we need to find on which set the measure is null. We have

$$\begin{aligned}
 \forall z \in S^1, U \in \mathbb{V}_{d,2}, \tilde{\mu}(S^{d-2}) &= 0 \\
 \iff \forall z \in S^1, U \in \mathbb{V}_{d,2}, \mu(O^{-1}((J^T)^{-1}(S^{d-2}))) &= \mu(F \cap S^{d-1}) = 0.
 \end{aligned} \tag{51}$$

744 Hence, we deduce that

$$\begin{aligned}
 \ker(\tilde{R}) &= \{\mu \in \mathcal{M}_{\text{even}}(S^{d-1}), \forall U \in \mathbb{V}_{d,2}, \forall z \in S^1, F = \text{span}(UU^T)^\perp \cap \text{span}(Uz), \\
 &\quad \mu(F \cap S^{d-1}) = 0\}.
 \end{aligned} \tag{52}$$

745 Moreover, we have that  $\cup_{U,z} F_{U,z} \cap S^{d-1} = \{H \cap S^{d-1} \subset \mathbb{R}^d, \dim(H) = d-1\}$ .

746 Indeed, on the one hand, let  $H$  an hyperplane,  $x \in H \cap S^{d-1}$ ,  $U \in \mathbb{V}_{d,2}$ , and note  $z = P^U(x)$ . Then,  
 747  $x \in F \cap S^{d-1}$  by Proposition 6 and  $H \cap S^{d-1} \subset \cup_{U,z} F_{U,z}$ .

748 On the other hand, let  $U \in \mathbb{V}_{d,2}$ ,  $z \in S^1$ ,  $F$  is a hyperplane since  $\dim(F) = d-1$  and therefore  
 749  $F \cap S^{d-1} \subset \{H, \dim(H) = d-1\}$ .

750 Finally, we deduce that

$$\ker(\tilde{R}) = \{\mu \in \mathcal{M}_{\text{even}}(S^{d-1}), \forall H \in \mathcal{G}_{d,d-1}, \mu(H \cap S^{d-1}) = 0\}. \tag{53}$$

## 751 A.8 Proof of Proposition 5

752 Let  $p \geq 1$ . First, it is straightforward to see that for all  $\mu, \nu \in \mathcal{P}_p(S^{d-1})$ ,  $SSW_p(\mu, \nu) \geq 0$ ,  
 753  $SSW_p(\mu, \nu) = SSW_p(\nu, \mu)$ ,  $\mu = \nu \implies SSW_p(\mu, \nu) = 0$  and that we have the triangular

754 inequality since

$$\begin{aligned}
\forall \mu, \nu, \alpha \in \mathcal{P}_p(S^{d-1}), \quad SSW_p(\mu, \nu) &= \left( \int_{\mathbb{V}_{d,2}} W_p^p(P_{\#}^U \mu, P_{\#}^U \nu) \, d\sigma(U) \right)^{\frac{1}{p}} \\
&\leq \left( \int_{\mathbb{V}_{d,2}} (W_p(P_{\#}^U \mu, P_{\#}^U \alpha) + W_p(P_{\#}^U \alpha, P_{\#}^U \nu))^p \, d\sigma(U) \right)^{\frac{1}{p}} \\
&\leq \left( \int_{\mathbb{V}_{d,2}} W_p^p(P_{\#}^U \mu, P_{\#}^U \alpha) \, d\sigma(U) \right)^{\frac{1}{p}} \\
&\quad + \left( \int_{\mathbb{V}_{d,2}} W_p^p(P_{\#}^U \alpha, P_{\#}^U \nu) \, d\sigma(U) \right)^{\frac{1}{p}} \\
&= SSW_p(\mu, \alpha) + SSW_p(\alpha, \nu),
\end{aligned} \tag{54}$$

755 using the triangular inequality for  $W_p$  and the Minkowski inequality. Therefore, it is at least a  
756 pseudo-distance.

757 To be a distance, we also need  $SSW_p(\mu, \nu) = 0 \implies \mu = \nu$ . Suppose that  $SSW_p(\mu, \nu) = 0$ .  
758 Since, for all  $U \in \mathbb{V}_{d,2}$ ,  $W_p^p(P_{\#}^U \mu, P_{\#}^U \nu) \geq 0$ ,  $SSW_p(\mu, \nu) = 0$  implies that for  $\sigma$ -ae  $U \in \mathbb{V}_{d,2}$ ,  
759  $W_p^p(P_{\#}^U \mu, P_{\#}^U \nu) = 0$  and hence  $P_{\#}^U \mu = P_{\#}^U \nu$  or  $(\tilde{R}\mu)^U = (\tilde{R}\nu)^U$  for  $\sigma$ -ae  $U \in \mathbb{V}_{d,2}$  since  $W_p$  is a  
760 distance on the circle. Therefore, it is a distance on the sets of injectivity of  $\tilde{R}$ .

## 761 A.9 Convergence Properties

762 **Proposition 8.** *Let  $(\mu_k), \mu \in \mathcal{P}_p(S^{d-1})$  such that  $\mu_k \xrightarrow[k \rightarrow \infty]{} \mu$ , then*

$$SSW_p(\mu_k, \mu) \xrightarrow[k \rightarrow \infty]{} 0. \tag{55}$$

763 *Proof.* Since the Wasserstein distance metrizes the weak convergence (Corollary 6.11 [101]), we have  
764  $P_{\#}^U \mu_k \xrightarrow[k \rightarrow \infty]{} P_{\#}^U \mu$  (by continuity)  $\iff W_p^p(P_{\#}^U \mu_k, P_{\#}^U \mu) \xrightarrow[k \rightarrow \infty]{} 0$  and hence by the dominated  
765 convergence theorem,  $SSW_p(\mu_k, \mu) \xrightarrow[k \rightarrow \infty]{} 0$ .  $\square$

## 766 B Background on the Sphere

### 767 B.1 Uniqueness of the Projection

768 Here, we discuss the uniqueness of the projection  $P^U$  for almost every  $x$ . For that, we recall some  
769 results of [9].

770 Let  $M$  be a closed subset of a complete finite-dimensional Riemannian manifold  $N$ . Let  $d$  be the  
771 Riemannian distance on  $N$ . Then, the distance from the set  $M$  is defined as

$$d_M(x) = \inf_{y \in M} d(x, y). \tag{56}$$

772 The infimum is a minimum since  $M$  is closed and  $N$  locally compact, but the minimum might  
773 not be unique. When it is unique, let's denote the point which attains the minimum as  $\pi(x)$ , i.e.  
774  $d(x, \pi(x)) = d_M(x)$ .

775 **Proposition 9** (Proposition 4.2 in [9]). *Let  $M$  be a closed set in a complete  $m$ -dimensional Riemannian  
776 manifold  $N$ . Then, for almost every  $x$ , there exists a unique point  $\pi(x) \in M$  that realizes the  
777 minimum of the distance from  $x$ .*

778 From this Proposition, they further deduce that the measure  $\pi_{\#}\gamma$  is well defined on  $M$  with  $\gamma$  a  
779 locally absolutely continuous measure w.r.t. the Lebesgue measure.

780 In our setting, for all  $U \in \mathbb{V}_{d,2}$ , we want to project a measure  $\mu \in \mathcal{P}(S^{d-1})$  on the great circle  
781  $\text{span}(UU^T) \cap S^{-1}$ . Hence, we have  $N = S^{d-1}$  which is a complete finite-dimensional Riemannian  
782 manifold and  $M = \text{span}(UU^T) \cap S^{d-1}$  a closed set in  $N$ . Therefore, we can apply Proposition 9  
783 and the push-forward measures are well defined for absolutely continuous measures.

784 **B.2 Optimization on the Sphere**

785 Let  $F : S^{d-1} \rightarrow \mathbb{R}$  be some functional on the sphere. Then, we can perform a gradient descent on a  
 786 Riemannian manifold by following the geodesics, which are the counterpart of straight lines in  $\mathbb{R}^d$ .  
 787 Hence, the gradient descent algorithm [3, 14] reads as

$$\forall k \geq 0, x_{k+1} = \exp_{x_k}(-\gamma \text{grad} f(x)), \quad (57)$$

788 where for all  $x \in S^{d-1}$ ,  $\exp_x : T_x S^{d-1} \rightarrow S^{d-1}$  is a map from the tangent space  $T_x S^{d-1} = \{v \in$   
 789  $\mathbb{R}^d, \langle x, v \rangle = 0\}$  to  $S^{d-1}$  such that for all  $v \in T_x S^{d-1}$ ,  $\exp_x(v) = \gamma_v(1)$  with  $\gamma_v$  the unique geodesic  
 790 starting from  $x$  with speed  $v$ , i.e.  $\gamma(0) = x$  and  $\gamma'(0) = v$ .

791 For  $S^{d-1}$ , the exponential map is known and is

$$\forall x \in S^{d-1}, \forall v \in T_x S^{d-1}, \exp_x(v) = \cos(\|v\|_2)x + \sin(\|v\|_2) \frac{v}{\|v\|_2}. \quad (58)$$

792 Moreover, the Riemannian gradient on  $S^{d-1}$  is known as [3, Eq. 3.37]

$$\text{grad} f(x) = \text{Proj}_x(\nabla f(x)) = \nabla f(x) - \langle \nabla f(x), x \rangle x, \quad (59)$$

793  $\text{Proj}_x$  denoting the orthogonal projection on  $T_x S^{d-1}$ .

794 For more details, we refer to [3, 17].

795 **B.3 Von Mises-Fisher Distribution**

796 The von Mises-Fisher (vMF) distribution is a distribution on  $S^{d-1}$  characterized by a concentration  
 797 parameter  $\kappa > 0$  and a location parameter  $\mu \in S^{d-1}$  through the density

$$\forall \theta \in S^{d-1}, f_{\text{vMF}}(\theta; \mu, \kappa) = \frac{\kappa^{d/2-1}}{(2\pi)^{d/2} I_{d/2-1}(\kappa)} \exp(\kappa \mu^T \theta), \quad (60)$$

798 where  $I_\nu(\kappa) = \frac{1}{2\pi} \int_0^\pi \exp(\kappa \cos(\theta)) \cos(\nu\theta) d\theta$  is the modified Bessel function of the first kind.

799 Several algorithms allow to sample from it, see e.g. [100, 107] for algorithms using rejection sampling  
 800 or [62] without rejection sampling.

801 For  $d = 1$ , the vMF coincides with the von Mises (vM) distribution, which has for density

$$\forall \theta \in [-\pi, \pi[, f_{\text{vM}}(\theta; \mu, \kappa) = \frac{1}{I_0(\kappa)} \exp(\kappa \cos(\theta - \mu)), \quad (61)$$

802 with  $\mu \in [0, 2\pi[$  the mean direction and  $\kappa > 0$  its concentration parameter. We refer to [71, Section  
 803 3.5 and Chapter 9] for more details on these distributions.

804 In particular, for  $\kappa = 0$ , the vMF (resp. vM) distribution coincides with the uniform distribution on  
 805 the sphere (resp. the circle).

806 Jung [55] studied the law of the projection of a vMF on a great circle. In particular, they showed that,  
 807 while the vMF plays the role of the normal distributions for directional data, the projection actually  
 808 does not follow a von Mises distribution. More precisely, they showed the following theorem:

809 **Theorem 1** (Theorem 3.1 in [55]). *Let  $d \geq 3$ ,  $X \sim \text{vMF}(\mu, \kappa) \in S^{d-1}$ ,  $U \in \mathbb{V}_{d,2}$  and  $T = P^U(X)$   
 810 the projection on the great circle generated by  $U$ . Then, the density function of  $T$  is*

$$\forall t \in [-\pi, \pi[, f(t) = \int_0^1 f_R(r) f_{\text{vM}}(t; 0, \kappa \cos(\delta)r) dr, \quad (62)$$

811 where  $\delta$  is the deviation of the great circle (geodesic) from  $\mu$  and the mixing density is

$$\forall r \in ]0, 1[, f_R(r) = \frac{2}{I_\nu^*(\kappa)} I_0(\kappa \cos(\delta)r) r (1 - r^2)^{\nu-1} I_{\nu-1}^*(\kappa \sin(\delta) \sqrt{1 - r^2}), \quad (63)$$

812 with  $\nu = (d - 2)/2$  and  $I_\nu^*(z) = (\frac{z}{2})^{-\nu} I_\nu(z)$  for  $z > 0$ ,  $I_\nu^*(0) = 1/\Gamma(\nu + 1)$ .

813 Hence, as noticed by Jung [55], in the particular case  $\kappa = 0$ , i.e.  $X \sim \text{Unif}(S^{d-1})$ , then

$$f(t) = \int_0^1 f_R(r) f_{\text{vM}}(t; 0, 0) dr = f_{\text{vM}}(t; 0, 0) \int_0^1 f_R(r) dr = f_{\text{vM}}(t; 0, 0), \quad (64)$$

814 and hence  $T \sim \text{Unif}(S^1)$ .



## 815 B.4 Normalizing Flows on the Sphere

816 Normalizing flows [82] are invertible transformations. There has been a recent interest in defining  
817 such transformations on manifolds, and in particular on the sphere [23, 91, 92].

818 Here, we implemented the Exponential map normalizing flows introduced in [92]. The transformation  
819  $T$  is

$$\forall x \in S^{d-1}, z = T(x) = \exp_x(\text{Proj}_x(\nabla\phi(x))), \quad (65)$$

820 where  $\phi(x) = \sum_{i=1}^K \frac{\alpha_i}{\beta_i} e^{\beta_i(x^T \mu_i - 1)}$ ,  $\alpha_i \geq 0$ ,  $\sum_i \alpha_i \leq 1$ ,  $\mu_i \in S^{d-1}$  and  $\beta_i > 0$  for all  $i$ .  $(\alpha_i)_i$ ,  
821  $(\beta_i)_i$  and  $(\mu_i)_i$  are the learnable parameters.

822 The density of  $z$  can be obtained as

$$p_Z(z) = p_X(x) \det(E(x)^T J_T(x)^T J_T(x) E(x))^{-\frac{1}{2}}, \quad (66)$$

823 where  $J_f$  is the Jacobian in the embedded space and  $E(x)$  it the matrix whose columns form an  
824 orthonormal basis of  $T_x S^{d-1}$ .

825 The common way of training normalizing flows is to use either the reverse or forward KL divergence.  
826 Here, we use them with a different loss, namely SSW.

## 827 C Additional Experiments

### 828 C.1 Evolution of SSW between von Mises-Fisher distributions

829 The KL divergence between the von Mises-Fisher distribution and the uniform distribution has been  
830 derived analytically in [28, 110] as

$$\begin{aligned} \text{KL}(\text{vMF}(\mu, \kappa) \parallel \text{vMF}(\cdot, 0)) &= \kappa \frac{I_{d/2}(\kappa)}{I_{d/2-1}(\kappa)} + \left(\frac{d}{2} - 1\right) \log \kappa - \frac{d}{2} \log(2\pi) - \log I_{d/2-1}(\kappa) \\ &+ \frac{d}{2} \log \pi + \log 2 - \log \Gamma\left(\frac{d}{2}\right). \end{aligned} \quad (67)$$

831 We plot on Figure 8 the evolution of KL and SSW *w.r.t.*  $\kappa$  for different dimensions. We observe a  
832 different trend. SSW seems to get lower with the dimension contrary to KL.

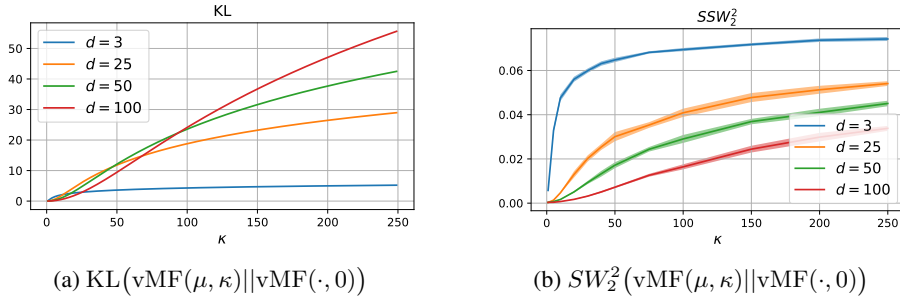


Figure 8: Evolution *w.r.t.*  $\kappa$  between  $\text{vMF}(\mu, \kappa)$  and  $\text{vMF}(\cdot, 0)$ . For SW, we used 100 projections (for memory reasons for  $d = 100$ ), and computed it for  $\kappa \in \{1, 5, 10, 20, 30, 40, 50, 75, 100, 150, 200, 250\}$ , 10 times by dimension and  $\kappa$ , and with 500 samples of both distributions.

833 As a sanity check, we compare on Figure 9 the evolution of SSW between vMF distributions  
834 where we fix  $\text{vMF}(\mu_0, 10)$  and we rotate the first vMF along a great circle. More precisely, we  
835 plot  $SW_2^2(\text{vMF}((1, 0, 0, \dots), 10), \text{vMF}((\cos(\theta), \sin(\theta), 0, \dots), 10))$  for  $\theta \in \{\frac{k\pi}{6}\}_{k \in \{0, \dots, 12\}}$ . As  
836 expected, we obtain a bell shape which is maximal when the second vMF distribution has for location

837 parameter  $-\mu_0$ . We observe a similar behavior between  $SSW_2$ ,  $SSW_1$  and  $SW_2$  with different  
 838 scales.

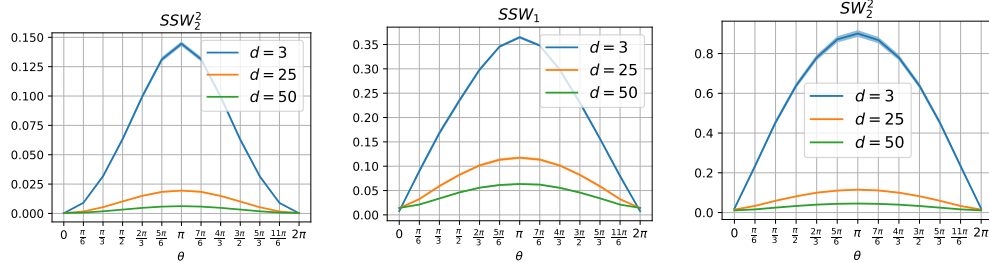


Figure 9: Evolution of  $SW$  between vMF samples in  $S^{d-1}$  (mean over 100 batch).

839 On Figure 10, we plot the evolution of SSW *w.r.t.* the number of projections for different dimensions.  
 840 We observe that for around 100 projections, the variance seems to be low enough.

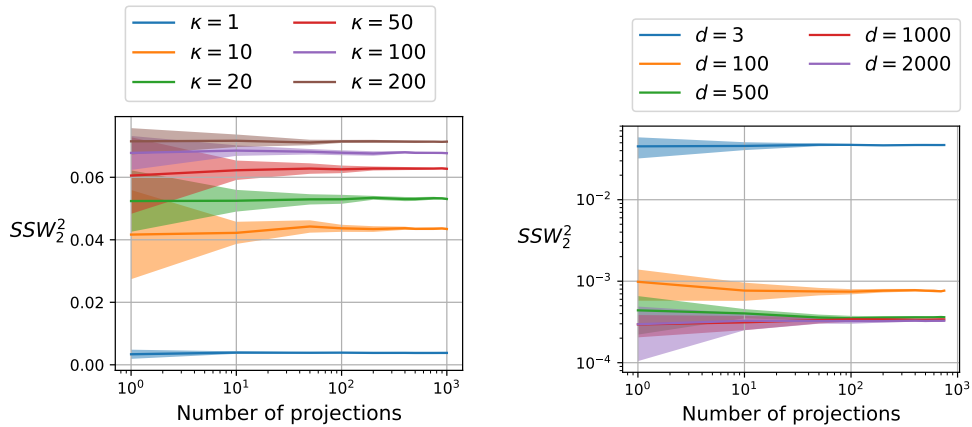


Figure 10: Influence of the number of projections. We compute  $SSW_2^2(\text{vMF}(\mu, \kappa) || \text{vMF}(\cdot, 0))$  20 times, for  $n = 500$  samples in dimension  $d = 3$ .

841 **Nadjahi et al. [76]** proved that, contrary to the Wasserstein distance, the classical sliced-Wasserstein  
 842 distance has a sample complexity independent of the dimension  $d$ . We show empirically on Figure 11  
 843 that we expect to have similar results for SSW by plotting SSW and the Wasserstein distance (with  
 844 geodesic distance) between samples of the uniform distribution on the sphere *w.r.t.* the number of  
 845 samples. We observe indeed that the convergence rate of SSW is independent of the dimension.

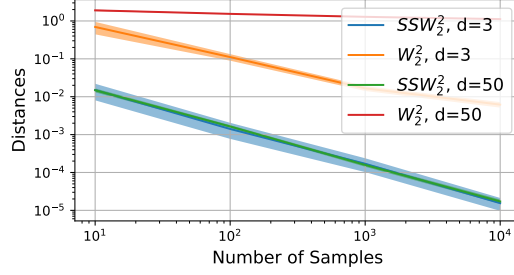


Figure 11: Spherical Sliced-Wasserstein and Wasserstein distance (with geodesic distance) between samples of the uniform distribution on the sphere. Results are averaged over 20 runs and the shaded are corresponds to the standard deviation.

## 846 C.2 Runtime Comparisons

847 We study here the evolution of the runtime *w.r.t.* different parameters. On Figure 12, we plot for  
 848 several dimensions the runtime to compute  $SSW_2$  *w.r.t.* the number of projections and the number of  
 849 samples. We observe the linearity *w.r.t.* the number of projections and the quasi-linearity *w.r.t.* the  
 850 number of samples.

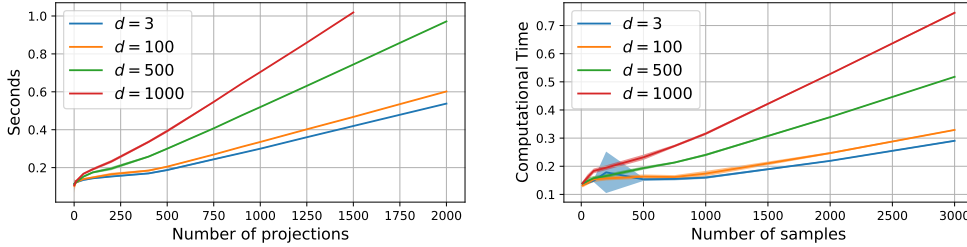


Figure 12: Computation time *w.r.t.* the number of projections or samples, taken for  $\kappa = 10$  and  $n = 500$  samples for the left figure, and  $\kappa = 10$  and 200 projections for the right figure, and for 20 times.

## 851 C.3 Gradient Flows

852 **Mixture of vMF distributions.** For the experiment in Section 5.1, we use as target distribution of  
 853 mixture of 6 vMF distributions from which we have access to samples. We refer to Appendix B.3 for  
 854 background on vMF distributions.

855 The 6 vMF distributions have weights  $1/6$ , concentration parameter  $\kappa = 10$  and location parameters  
 856  $\mu_1 = (1, 0, 0)$ ,  $\mu_2 = (0, 1, 0)$ ,  $\mu_3 = (0, 0, 1)$ ,  $\mu_4 = (-1, 0, 0)$ ,  $\mu_5 = (0, -1, 0)$  and  $\mu_6 = (0, 0, -1)$ .

857 We use two different approximation of the distribution. First, we approximate it using the empirical  
 858 distribution, *i.e.*  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  and we optimize over the particles  $(x_i)_{i=1}^n$ . To optimize over  
 859 particles, we can either use a projected gradient descent:

$$\begin{cases} x^{(k+1)} = x^{(k)} - \gamma \nabla_{x^{(k)}} SSW_2^2(\hat{\mu}_k, \nu) \\ x^{(k+1)} = \frac{x^{(k+1)}}{\|x^{(k+1)}\|_2}, \end{cases} \quad (68)$$

860 or a Riemannian gradient descent on the sphere [3] (see Appendix B.2 for more details). Note that  
 861 the projected gradient descent is a Riemannian gradient descent with retraction [17].

862 We can also use neural networks such as a multilayer perceptron (MLP). We used a MLP composed  
 863 of 5 layers of 100 units with leaky relu activation functions. The output of the MLP is normalized on  
 864 the sphere using a  $\ell^2$  normalization. We perform a gradient descent using Adam [57] as the optimizer

865 with a learning rate of  $10^{-4}$  for 2000 epochs. We approximate SSW with  $L = 1000$  projections and  
 866 a batch size of 500. The base distribution is choose as the uniform distribution on the sphere.

867 We report on Figure 13 a comparison of the 2 approximations where the density is estimated with a  
 868 Gaussian kernel density estimator.

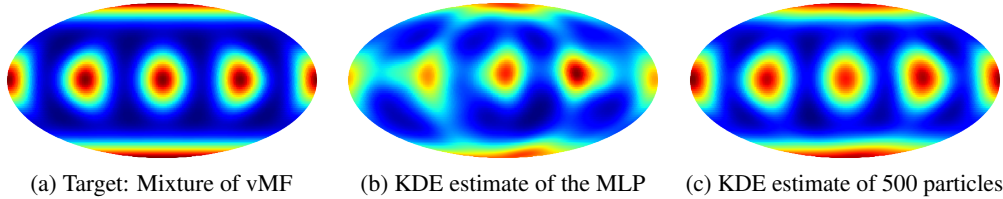


Figure 13: Minimization of SSW with respect to a mixture of vMF.

869 **vMF distribution.** A a simpler experiment, we choose a simple vMF distribution with  $\kappa = 10$ . We  
 870 report on Figure 14 the evolution of the density approximated using a KDE, and on Figure 15 the  
 871 evolution of particles.

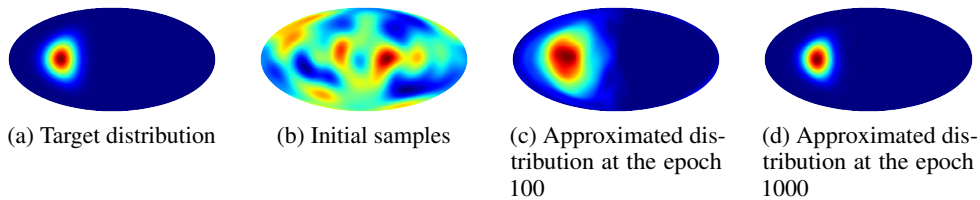


Figure 14: Gradient Flows on SW with a vMF target and Mollweide projections. The distributions are approximated using KDE.

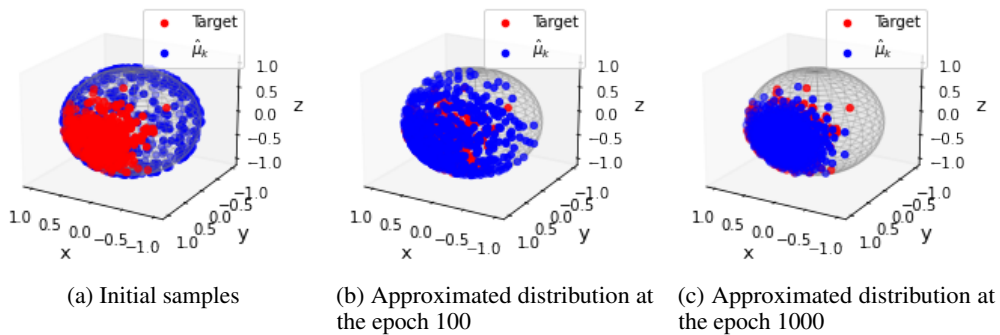


Figure 15: Gradient Flows on SW with a vMF target and Mollweide projections.

## 872 C.4 Sliced-Wasserstein Variational Inference

### 873 C.4.1 Variational Inference

874 In variational inference (VI) [12, 54], we have some observed data  $(x_i)_{i=1}^n$  and some latent data  
 875  $(z_i)_{i=1}^n$ . The goal of variational inference is to approximate the posterior distribution  $p(\cdot|x)$  by some  
 876 distribution  $q \in \mathcal{Q}$  where  $\mathcal{Q}$  is a family of probabilities. The usual way of doing that is to minimize

---

**Algorithm 2** SWVI [111]

---

**Input:**  $V$  a potential,  $K$  the number of iterations of SWVI,  $N$  the batch size,  $\ell$  the number of MCMC steps  
**Initialization:** Choose  $q_\theta$  a sampler  
**for**  $k = 1$  **to**  $K$  **do**  
  Sample  $(z_i^0)_{i=1}^N \sim q_\theta$   
  Run  $\ell$  MCMC steps starting from  $(z_i^0)_{i=1}^N$  to get  $(z_j^\ell)_{j=1}^N$   
  // Denote  $\hat{\mu}_0 = \frac{1}{N} \sum_{j=1}^N \delta_{z_j^0}$  and  $\hat{\mu}_\ell = \frac{1}{N} \sum_{j=1}^N \delta_{z_j^\ell}$   
  Compute  $J = SW_2^2(\hat{\mu}_0, \hat{\mu}_\ell)$   
  Backpropagate through  $J$  w.r.t.  $\theta$   
  Perform a gradient step  
**end for**

---

877 the Kullback-Leibler divergence among this family, *i.e.*

$$\min_{q \in \mathcal{Q}} \text{KL}(q||p(\cdot|x)) = \mathbb{E}_q[\log \left( \frac{q(Z)}{p(Z|x)} \right)]. \quad (69)$$

878 But the KL divergence suffers from some drawbacks, as it is only a divergence (*i.e.* it does not satisfy  
879 the triangular inequality, and it is non symmetric), but it also suffers from under estimating the target  
880 distribution (or over estimating it for the reverse KL).

881 Yi and Liu [111] propose to use an optimal transport distance instead, namely the SW distance  
882 which gives the sliced-Wasserstein variational inference method. Basically, given some unnormalized  
883 probability  $p(\cdot|x)$  that we want to approximate with some variational distribution  $q_\phi$ , we can first  
884 apply a MCMC algorithm and then learn  $q_\phi$  using a gradient descent on SW with the target being  
885 the empirical distributions of the samples given by the MCMC. But running long MCMC chain is  
886 time consuming and it might be difficult to diagnose burn-in period. Therefore, they propose to only  
887 run at each iteration some number of steps  $t$  of MCMC chain, and then learn by gradient descent the  
888 variational distribution. Therefore, the variational distribution is guided at each step by the MCMC  
889 samples toward the stationary distribution which is the target. This is called an amortized sampler  
890 (see Problem 1 in [103]). We sum up the procedure in Algorithm 2.

891 We propose here to substitute  $SW$  by  $SSW$  in order to perform SSWVI on the sphere. To do that,  
892 we first need a MCMC method on the sphere.

### 893 C.4.2 MCMC on the Sphere

894 Several MCMC methods on the sphere have been proposed. For example, Hamiltonian Monte-Carlo  
895 (HMC) methods were proposed in [18, 63, 68], and Riemannian Langevin algorithms were proposed  
896 in [65, 105].

897 In our experiments, we use the Geodesic Langevin algorithm (GLA) introduced by Wang et al.  
898 [105]. This algorithm is a natural generalization of the Unadjusted Langevin Algorithm (ULA) and it  
899 consists at simply following the geodesics of the regular ULA step, *i.e.*

$$\forall k > 0, x_{k+1} = \exp_{x_k} \left( \text{Proj}_{x_k}(-\gamma \nabla V(x_k) + \sqrt{2\gamma}Z) \right), Z \sim \mathcal{N}(0, I), \quad (70)$$

900 where for the sphere,

$$\forall x \in S^{d-1}, \forall v \in T_x S^{d-1}, \exp_x(v) = x \cos(\|v\|) + \frac{v}{\|v\|} \sin(\|v\|), \quad (71)$$

901  $\text{Proj}_x$  is the projection on the tangent space  $T_x S^{d-1} = \{v \in \mathbb{R}^d, \langle x, v \rangle = 0\}$  (which is the  
902 orthogonal space) and is defined as

$$\text{Proj}_x(v) = v - \langle x, v \rangle x. \quad (72)$$

903 For more details, we refer to [3].

904 We use GLA here for simplicity and as a proof of concept. But note that GLA, as ULA, is biased  
 905 and therefore the distribution learned will not be the exact true stationary distribution. However, a  
 906 Metropolis-Hastings step at each iteration could be used to enforce the reversibility *w.r.t.* the target  
 907 distribution or we could use other MCMC with more appealing convergence properties (see *e.g.* [68]).

### 908 C.4.3 Applications

909 **Target: Power spherical distribution.** First, as a simple example on  $S^2$ , we use the power spherical  
 910 distribution introduced by De Cao and Aziz [29]. This distribution has the advantage over the vMF  
 911 distribution to allow for the direct use of the reparameterization trick since it does not require rejection  
 912 sampling. The pdf is obtained as,

$$\forall x \in S^{d-1}, p_X(x; \mu, \kappa) \propto (1 + \mu^T x)^\kappa \quad (73)$$

913 with  $\mu \in S^{d-1}$  and  $\kappa > 0$ . We can sample from drawing first  $Z \sim \text{Beta}(\frac{d-1}{2} + \kappa, \frac{d-1}{2})$ ,  
 914  $v \sim \text{Unif}(S^{d-2})$ , then constructing  $T = 2Z - 1$  and  $Y = [T, v^T \sqrt{1 - T^2}]^T$ . Finally, apply a  
 915 Householder reflection about  $\mu$  to  $Y$ . All the operations are well differentiable and allow to apply the  
 916 reparametrization trick. For the algorithm, see Algorithm 1 in [29]. Hence, in this case, if we denote  
 917  $g_\theta$  the map which takes samples from a uniform distribution on  $S^{d-2}$  and from a Beta distribution as  
 918 input and outputs samples of power spherical distribution with parameters  $\theta = (\kappa, \mu)$ , we can use it  
 919 as the sampler. We test the algorithm with a target being a power spherical distribution of parameter  
 920  $\mu = (0, 1, 0)$  and  $\kappa = 10$ , starting from  $\mu = (1, 1, 1)$  and  $\kappa = 0.1$ . Performing 2000 optimization  
 921 steps with a gradient descent (Riemannian gradient descent on  $\mu$  to stay on the sphere), and 20 steps  
 922 of the GLA algorithm, we are getting close enough to the true distribution as we can see on Figure 16.

923 For the hyperparameters, we used a step size of  $10^{-3}$  for GLA, 1000 projections to approximate SSW,  
 924 a Riemannian gradient descent on the sphere [3] to learn the location parameter  $\mu$  with a learning rate  
 925 of 2, and a learning of 200 for  $\kappa$ . We performed  $K = 2000$  steps and used  $N = 500$  particles.

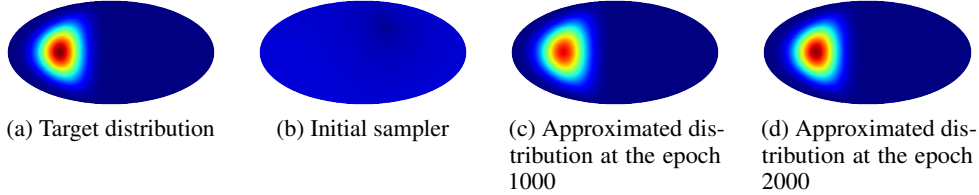


Figure 16: SWVI on Power Spherical Distributions with Mollweide projections.

926 **Target: mixture of vMFs.** In Section 5.1, we perform amortized variational inference with a  
 927 mixture of vMF distributions as target. For this, we train exponential map normalizing flows (see  
 928 [92] and Appendix B.4). Moreover, we use the same target as Rezende et al. [92], *i.e.* the target  
 929  $\nu$  has a density  $p(x) \propto \sum_{k=1}^4 e^{10x^T T_{s \rightarrow e}(\mu_k)}$  with  $\mu_1 = (0.7, 1.5)$ ,  $\mu_2 = (-1, 1)$ ,  $\mu_3 = (0.6, 0.5)$   
 930 and  $\mu_4 = (-0.7, 4)$ . These are spherical coordinates which are converted to euclidean using  
 931  $T_{s \rightarrow e}(\theta, \phi) = (\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi)$ .

932 The exponential map normalizing flow is composed of  $N = 6$  blocks with  $K = 5$  components. We  
 933 run the algorithm for 10000 iterations, with at each iteration 20 steps of GLA with  $\gamma = 10^{-1}$  as  
 934 learning rate, and one step of backpropagation through SSW using the Adam [57] optimizer with a  
 935 learning rate of  $10^{-3}$ .

936 We report on Figure 4 the Mollweide projection of the learned density. Since we learn to samples from  
 937 a noise distribution, here the uniform distribution on the sphere, we do not have directly access to the  
 938 density and we report a kernel density estimate with a Gaussian kernel using the implementation of  
 939 Scipy [102].

940 We also report in Figure 5 the effective sample size (ESS) [33, 69] over the iterations. The ESS is  
 941 estimated by [92]

$$\text{ESS} = \frac{\text{Var}_{\text{Unif}}(e^{-\beta u(X)})}{\text{Var}_q\left(\frac{e^{-\beta u(X)}}{q_\eta(X)}\right)} \approx \frac{\left(\sum_{s=1}^S w_s\right)^2}{\sum_{s=1}^S w_s^2}, \quad (74)$$

942 where  $w_s = e^{-\beta u(x_s)}/q_\eta(x_s)$ . The ESS is reported as a percentage of the sample size. Higher ESS  
 943 indicates that the flow matches the target better [92].

### 944 C.5 Sliced-Wasserstein Autoencoder

945 We recall that in the WAE framework, we want to minimize

$$\mathcal{L}(f, g) = \int c(x, g(f(x))) d\mu(x) + \lambda D(f_{\#}\mu, p_Z), \quad (75)$$

946 where  $f$  is an encoder,  $g$  a decoder,  $p_Z$  a prior distribution,  $c$  some cost function and  $D$  is a divergence  
 947 in the latent space. Several  $D$  were proposed. For example, Tolstikhin et al. [99] proposed to use  
 948 the MMD, Kolouri et al. [59] used the SW distance, Patrini et al. [84] used the Sinkhorn divergence,  
 949 Kolouri et al. [60] used the generalized SW distance. Here, we use  $D = \text{SSW}_2^2$ .

950 **Architecture and procedure.** For the encoder  $f$  and the decoder  $g$ , we use the same architecture  
 951 as Kolouri et al. [59].

952 For both the encoder and the decoder architecture, we use fully convolutional architectures with 3x3  
 953 convolutional filters. More precisely, the architecture of the encoder is

$$\begin{aligned} x \in \mathbb{R}^{28 \times 28} &\rightarrow \text{Conv2d}_{16} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv2d}_{16} \rightarrow \text{LeakyReLU}_{0.2} \rightarrow \text{AvgPool}_2 \\ &\rightarrow \text{Conv2d}_{32} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv2d}_{32} \rightarrow \text{LeakyReLU}_{0.2} \rightarrow \text{AvgPool}_2 \\ &\rightarrow \text{Conv2d}_{64} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv2d}_{64} \rightarrow \text{LeakyReLU}_{0.2} \rightarrow \text{AvgPool}_2 \\ &\rightarrow \text{Flatten} \rightarrow \text{FC}_{128} \rightarrow \text{ReLU} \\ &\rightarrow \text{FC}_{d_Z} \rightarrow \ell^2 \text{ normalization} \end{aligned}$$

954 where  $d_Z$  is the dimension of the latent space (either 11 for  $S^{10}$  or 3 for  $S^2$ ).

955 The architecture of the decoder is

$$\begin{aligned} z \in \mathbb{R}^{d_Z} &\rightarrow \text{FC}_{128} \rightarrow \text{FC}_{1024} \rightarrow \text{ReLU} \\ &\rightarrow \text{Reshape}(64 \times 4 \times 4) \rightarrow \text{Upsample}_2 \rightarrow \text{Conv}_{64} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv}_{64} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Upsample}_2 \rightarrow \text{Conv}_{64} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv}_{32} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Upsample}_2 \rightarrow \text{Conv}_{32} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv}_1 \rightarrow \text{Sigmoid} \end{aligned}$$

956 To compare the different autoencoders, we used as the reconstruction loss the binary cross entropy,  
 957  $\lambda = 10$ , Adam [57] as optimizer with a learning rate of  $10^{-3}$  and Pytorch’s default momentum  
 958 parameters for 800 epochs with batch of size  $n = 500$ . Moreover, when using SW type of distance,  
 959 we approximated it with  $L = 1000$  projections.

960 We report in Table 1 the FID obtained using 10000 samples and we report the mean over 5 trainings.

961 For SSW, we used the formulation using the uniform distribution (12). To compute SW, we used the  
 962 POT library [39]. To compute the Sinkhorn divergence, we used the GeomLoss package [37].

963 **Additional experiments.** We report on Figure 17 samples obtained with SSW for a uniform prior  
 964 on  $S^{10}$ .

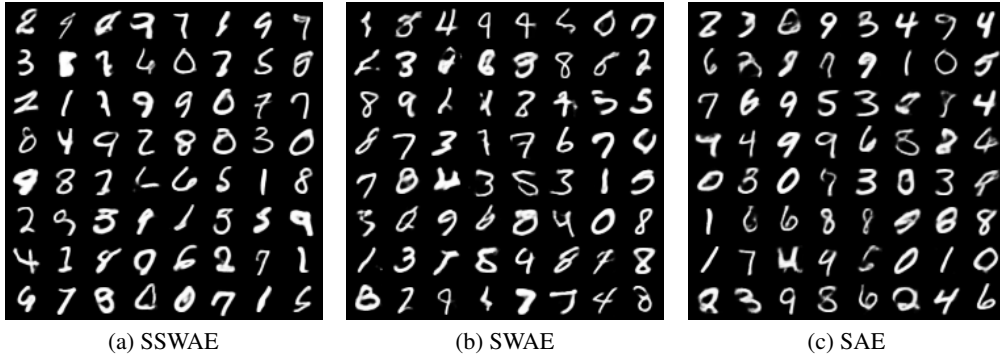


Figure 17: Samples generated with Sliced-Wasserstein Autoencoders with a uniform prior on  $S^{10}$ .

965 On Figure 18, we add the evolution over epochs of the Wasserstein distance between generated  
 966 images and samples from the test set.

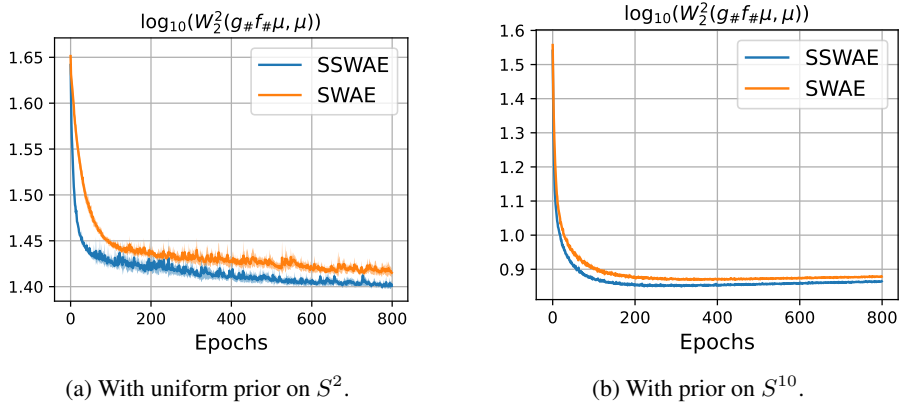


Figure 18: Comparison of the evolution of the Wasserstein distance over epochs between SWAE and SSWAE on MNIST (averaged over 5 trainings).

967 **C.6 Self-supervised learning**

968 We conduct experiments using SSW to  
 969 prevent collapsing representations in contrastive self-supervised learning (SSL)  
 970 models. Such contrastive losses on the hypersphere have exhibited great representat-  
 971 ive capacity [20, 21, 108] on unlabelled datasets by learning robust image represen-  
 972 tations invariantly to augmentations. As proposed in [104], the contrastive objec-  
 973 tive can be decomposed into an alignment loss which forces positive representations  
 974 coming from the same image to be similar and a uniformity loss which preserves maximal information of the feature distribution and hence  
 975 avoids collapsing representations. Without the uniformity loss, the representations tend to converge  
 976  
 977  
 978  
 979  
 980  
 981

Table 2: Linear evaluation on CIFAR10. The features are taken either on the encoder output or directly on the sphere  $S^2$ .

Method	Encoder output	$S^2$
Supervised	82.26	81.43
Chen et al. [21]	66.55	59.09
Wang and Isola [104]	60.53	55.86
SW-SSL, $\lambda = 1, L = 10$	62.65	57.77
SW-SSL, $\lambda = 1, L = 3$	62.46	57.64
SSW-SSL, $\lambda = 20, L = 10$	64.89	58.91
SSW-SSL, $\lambda = 20, L = 3$	63.75	59.75



982 towards a constant representation which yields the best alignment loss possible but also contains  
 983 no information about original images. Wang and Isola [104] propose to enforce uniformity by  
 984 leveraging the Gaussian potential kernel which is bound to the uniform distribution on the sphere.  
 985 This formulation is also related to the denominator of the contrastive loss as specified in Chen et al.  
 986 [21]. We propose to replace the Gaussian kernel uniformity loss with SSW for which the complexity  
 987 is more linear *w.r.t.* the number of batch samples. A simple choice of the alignment loss is to  
 988 minimize the mean squared euclidean distance between pairs of different augmented versions of the  
 989 same image. A self-supervised learning network is pre-trained using this alignment loss added with  
 990 an uniformity term. Our overall self-supervised loss can be defined as:

$$\mathcal{L}_{\text{SSW-SSL}} = \underbrace{\frac{1}{n} \sum_{i=1}^n \|z_i^A - z_i^B\|_2^2}_{\text{Alignment loss}} + \frac{\lambda}{2} \underbrace{\left( \text{SSW}_2^2(z^A, \nu) + \text{SSW}_2^2(z^B, \nu) \right)}_{\text{Uniformity loss}}, \quad (76)$$

991 where  $z^A, z^B \in \mathbb{R}^{n \times d}$  are the representations from the network projected on the hypersphere of  
 992 two augmented versions of the same images,  $\nu = \text{Unif}(S^{d-1})$  is the uniform distribution on the  
 993 hypersphere and  $\lambda > 0$  is used to balance the two terms.

994 We pretrain a ResNet18 [47] model on the CIFAR10 [61] data with projections projected onto the  
 995 sphere  $S^2$ . This feature dimension allow us to visualize the entire validation set of CIFAR10 and  
 996 its distribution on the sphere. The visualization of the projections on  $S^2$  are visible on Figure 19.  
 997 We then evaluate the performance of each contrastive objective by fitting a linear classifier on top  
 998 of the output of the layer before the projection on the sphere on the training dataset as is common  
 999 for SSL methods. For comparison, we also report the results when the features are taken directly on  
 1000 the sphere. As a baseline, we also train a predictive supervised encoder by training jointly the linear  
 1001 classifier and the image encoder in a supervised manner using cross entropy.

1002 We use a ResNet18 [47] encoder which outputs 1024 features that are then projected onto the sphere  
 1003  $S^2$  using a last fully connected layer followed by a  $\ell^2$  normalization. We pretrain the model for 200  
 1004 epochs using minibatch stochastic gradient descent (SGD) with a momentum of 0.9, a weight decay  
 1005 of 0.001 and an initial learning rate of 0.05. We use a batch size of 512 samples. The images are  
 1006 augmented using a standard set of random augmentations for SSL: random crops, horizontal flipping,  
 1007 color jittering and gray scale transformation as done in Wang and Isola [104]. For the trade-off  
 1008 parameter  $\lambda$ , we  $\lambda = 20$  for SSW and  $\lambda = 1$  for SW.

1009 To evaluate the performance of representations, we use the common linear evaluation protocol where  
 1010 a linear classifier is fitted on top of the pre-trained representations and the best validation accuracy  
 1011 is reported. The linear classifiers are trained for 100 epochs using the Adam [57] optimizer with a  
 1012 learning rate of 0.001 with a decay of 0.2 at epoch 60 and 80. We compare our methods with two  
 1013 other contrastive objectives, Chen et al. [21] with the normalized temperature-scaled cross-entropy  
 1014 (NT-Xent) loss and Wang and Isola [104] which proposes to decompose the objective in two distinct  
 1015 terms  $\mathcal{L}_{\text{align}}$  and  $\mathcal{L}_{\text{uniform}}$ . We recall the respective uniformity loss of each method in Table 3. As  
 1016 one can see in Table 2, our method achieves here comparable performances to two state-of-the-art  
 1017 approaches, yet slightly under-performing compared to [21]. We suspect that a finer validation of  
 1018 the balancing parameter  $\lambda$  is needed. Especially since the representations on Figure 19b are not  
 1019 completely uniformly distributed around the sphere after pre-training compared to other contrastive  
 1020 methods. Nevertheless, these preliminary results show that SSW-SSL is a promising contrastive  
 1021 learning approach without explicit distances between negative samples, especially compared to SW  
 1022 on the sphere. To this end, further works should be devoted to finding a good balance between the  
 1023 alignment and uniformity objectives.

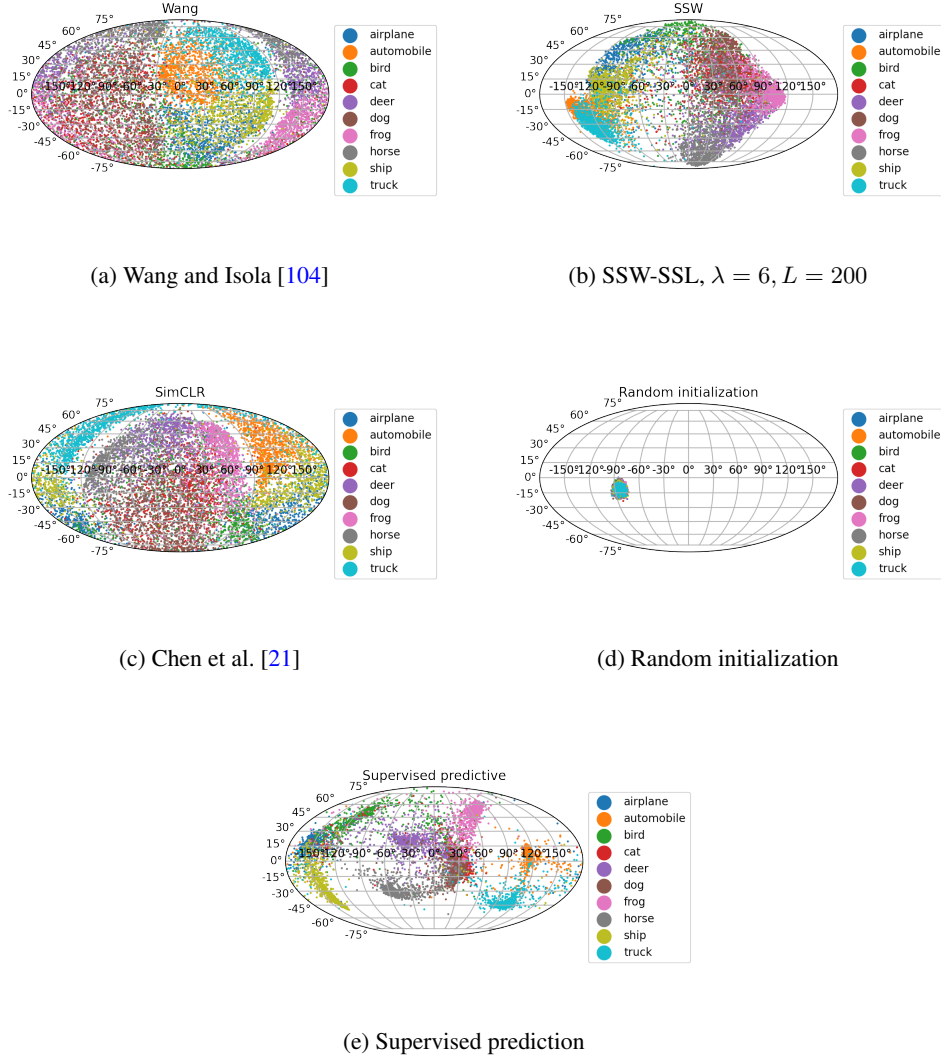


Figure 19: The CIFAR10 validation set on  $S^2$  after pre-training.

Table 3: Comparison of contrastive methods and their respective uniformity objective where  $z^A, z^B \in \mathbb{R}^{n \times d}$  are representations from two augmented versions of the same set of images and  $\nu = \text{Unif}(S^{d-1})$  is the uniform distribution on the hypersphere.

Method	$\mathcal{L}_{\text{uniform}}(z^A) + \mathcal{L}_{\text{uniform}}(z^B)$	Complexity
Chen et al. [21]	$\frac{1}{2n} \sum_{i=1}^n \log \sum_{j \neq i} \exp\left(\frac{\langle \hat{z}_i, \hat{z}_j \rangle}{\tau}\right), \hat{z} = \text{cat}(z^A, z^B)$	$O(n^2 d)$
Wang and Isola [104]	$\sum_{z \in \{z^A, z^B\}} \log \frac{2}{n(n-1)} \sum_{i>j} \exp(-t \ z_i - z_j\ _2^2)$	$O(n^2 d)$
SSW-SSL (Ours)	$\frac{1}{2} (SSW_2^2(z^A, \nu) + SSW_2^2(z^B, \nu))$	$O(Ln(d + \log n))$