
Supplementary material: Neural approximation of Wasserstein distance via a universal architecture for symmetric and factorwise group invariant functions

Samantha Chen
Department of Computer Science Engineering
University of California - San Diego
La Jolla, CA 92092
sac003@ucsd.edu

Yusu Wang
Halıcıoğlu Data Science Institute
University of California - San Diego
La Jolla, CA 92092
yusuwang@ucsd.edu

A Approximating Hausdorff distance

Suppose the underlying metric for \mathcal{X} (the space of weighted point sets) is Hausdorff distance, d_H . To clarify, given $S_1 = \{(x_i, w_i) : x_i \in \mathbb{R}^d, \sum_{i=1}^n w_i\}$ and $S_2 = \{(y_i, w'_i) : y_i \in \mathbb{R}^d, \sum_{i=1}^m w'_i\}$,

$$d_H(S_1, S_2) = d_H(\{x_1, \dots, x_n\}, \{y_1, \dots, y_m\}).$$

To get an encoding/decoding for \mathcal{X} equipped with Hausdorff distance using max-pooling, we will use the construction given in the proof of universality in [11], which we will briefly recap here. As in Lemma 3.4, we consider a δ -net for our original space X , $\{y_1, \dots, y_a\}$. For each y_i , we will again let $h_i(x) = e^{-d_X(x, B_\delta(y_i))}$ and define $h : X \rightarrow \mathbb{R}^a$ as $h(x) = [h_1(x), \dots, h_a(x)]$. However, instead of using sum-pooling to define $\mathbf{h} : \mathcal{X} \rightarrow \mathbb{R}^a$, we will use max-pooling as

$$\mathbf{h}(S) = \max_{x \in S} h(x)$$

where max is element-wise max. Then the decoding function $g : \mathbb{R}^a \rightarrow \mathcal{X}$ is defined as

$$g(v) = \{(y_i, \frac{v_i}{\|v\|_1}) : v_i \geq 1\}$$

Note that $v_i \geq 1$ when there is an $x_i \in S$ such that $x_i \in B_\delta(y_i)$. So for each $x_i \in S$ there is $y_i \in g \circ h(S)$ such that $d_X(x_i, y_i) < \delta$. Thus, $d_H(S, g \circ h(S)) < \delta$.

The final $\mathcal{N}_{\text{ProductNet}}$ for approximating uniformly continuous functions for Hausdorff distance with max-pooling is

$$\mathcal{N}_{\text{ProductNet}}^{\max}(A, B) = \rho_{\theta_3} \left(\phi_{\theta_2} \left(\max_{x \in A} (h_{\theta_1}(x)) \right) + \phi_{\theta_2} \left(\max_{y \in B} (h_{\theta_1}(y)) \right) \right).$$

Note that $\mathcal{N}_{\text{ProductNet}}^{\max}$ is able to approximate Hausdorff distance to ϵ -accuracy.

B Comparison with other neural network architectures for approximating Wasserstein distance

Through the lens of 1-Wasserstein approximation, we compare the power of $\mathcal{N}_{\text{ProductNet}}$ with other natural architectures, namely standard MLPs and Siamese networks. The standard MLP and Siamese networks are the natural points of comparison for $\mathcal{N}_{\text{ProductNet}}$ as (1) MLPs are the go-to choice in deep learning due to their versatility and efficacy across diverse problem domains and (2) Siamese networks have gained significant prominence in the context of metric learning, frequently employed

for tasks aimed at learning effective similarity metrics between pairs of inputs[5, 2, 7] However, we discuss the limitations of both the standard MLP and Siamese networks for the tasks of learning Wasserstein distances. For this setting, we will let (X, d_X) be some compact Euclidean space, i.e. $X \subseteq \mathbb{R}^d$ is compact and d_X is some ℓ_p norm. Then \mathcal{X} be the set of weighted point sets of cardinality at most N i.e. $S = \{(x_i, w_i) : w_i \in \mathbb{R}_+, \sum_i w_i = 1, x_i \in X\}$.

Comparison with multilayer perceptrons (MLPs). First, we consider a standard MLP model where the model takes a single tensor as input. Note that for each S , there is an $\tilde{S} \in \mathcal{X}$ such that

$$\tilde{S} = \{(x_i, w_i) : (x_i, w_i) \in S\} \cup \{(x_1, 0)_{\times(N-|S|)}\}$$

where $(x_1, 0)_{\times(N-|S|)}$ means that we repeat the element $(x_1, 0)$ $(N - |S|)$ times. Notice that $W_p(S, \tilde{S}) = 0$ and for any $S' \in \mathcal{X}$, $W_p(S, S') = W_p(\tilde{S}, S')$. If $|S| = N$, $\tilde{S} = S$. To use an MLP model to approximate Wasserstein distance, we want to represent any $(S_1, S_2) \in \mathcal{X} \times \mathcal{X}$ as an element of $\mathbb{R}^{(d+1) \times 2N}$ where $S_1 = \{(x_i, w_i) : w_i \in \mathbb{R}_+, \sum_i w_i = 1, x_i \in X\}$ and $S_2 = \{(x'_i, w'_i) : w'_i \in \mathbb{R}_+, \sum_i w'_i = 1, x'_i \in X\}$. Informally, given any "empty" element in $\mathbb{R}^{(d+1) \times 2N}$, we will map S_1 to the first N columns and map S_2 to the last N columns. Formally, we will define the mapping $\beta : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^{(d+1) \times 2N}$ given S_1 and S_2 as defined above as

$$\beta(S_1, S_2) = \beta(\tilde{S}_1, \tilde{S}_2) = \begin{bmatrix} x_1 & \cdots & x_N & x'_1 & \cdots & x'_N \\ w_1 & \cdots & w_N & w'_1 & \cdots & w'_N \end{bmatrix}$$

where we abuse notation slightly and use (x_i, w_i) and (x'_i, w'_i) to describe the elements in \tilde{S}_1 and \tilde{S}_2 respectively. Given $M \in \mathbb{R}^{(d+1) \times 2N}$, define $\beta^{-1} : \text{Image}(\beta) \rightarrow \mathcal{X} \times \mathcal{X}$ as

$$\beta^{-1} \left(\begin{bmatrix} x_1 & \cdots & x_N & x'_1 & \cdots & x'_N \\ w_1 & \cdots & w_N & w'_1 & \cdots & w'_N \end{bmatrix} \right) = (\{(x_i, w_i)\}, \{(x'_i, w'_i)\})$$

Now, in order to use the classical universal approximation theorem to approximate W_p to an ϵ -approximation error, we must show that there is a continuous function $f : \text{Image}(\beta) \rightarrow \mathbb{R}$ such that $W_p(S_1, S_2) = f(\beta(S_1, S_2))$.

Lemma B.1 *Given $\beta : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^{(d+1) \times 2N}$ as defined previously, there is a continuous function $f : \text{Image}(\beta) \rightarrow \mathbb{R}$ such that for any $\epsilon > 0$, and for any $(S_1, S_2) \in \mathcal{X} \times \mathcal{X}$, $|W_p(S_1, S_2) - f(\beta(S_1, S_2))| < \epsilon$.*

Proof: Let $\epsilon > 0$. From Corollary 3.5, there is a $\delta > 0$ such that there are continuous $h : X \rightarrow \mathbb{R}^{a(\delta)}$, $\phi : \mathbb{R}^{a(\delta)} \rightarrow \mathbb{R}^{a'}$, and $\rho : \mathbb{R}^{a'} \rightarrow \mathbb{R}$, such that for any $A, B \in \mathcal{X}$

$$\left| W_p(A, B) - \rho \left(\phi \left(\sum_{x \in A} w_x h(x) \right) + \phi \left(\sum_{x \in B} w_x h(x) \right) \right) \right| < \epsilon.$$

Let $(S_1, S_2) \in \mathcal{X} \times \mathcal{X}$ and let $M = \beta(S_1, S_2)$. Note that $\beta^{-1}(M) = (S_1, S_2)$. Let $M \in \text{Image}(\beta)$, let $M[:, i]$ denote the i th column vector of M , let $M[: d, i]$ denote vector of the first d elements in the i th column vector of M and let $M[i, j]$ denote the ij th entry of M . We know that $M[:, i] \in \mathbb{R}^{d+1}$. Since $M[:, i] \in \text{Image}(\beta)$, we know that $M[:, i] \in X$ and $M[d, i]$ corresponds to a weight. Furthermore, since h is a continuous mapping from $X \rightarrow \mathbb{R}$, $M[d, i]h(M[:, i])$ is continuous. Thus, since ϕ is also continuous, the function $f' : \text{Image}(\beta) \rightarrow \mathbb{R}^{a'}$

$$\phi \left(\sum_{i=1}^N M[d, i] h(M[:, i]) \right) + \phi \left(\sum_{i=N+1}^{2N} M[d, i] h(M[:, i]) \right)$$

is continuous. Since ρ is also continuous, we get that $\rho \circ f' : \text{Image}(\beta) \rightarrow \mathbb{R}$ is continuous. Thus, $f : \text{Image}(\beta) \rightarrow \mathbb{R}$ which is $f = \rho \circ f'$ is also continuous and $|f(\beta(S_1, S_2)) - W_p(S_1, S_2)| < \epsilon$ by the property from Corollary 3.5 \blacksquare

Since f in the above Lemma operates on fixed dimensional space, we can now use an MLP to approximate it. Then by the above Lemma and universal approximation of MLPs, there is an MLP

which can approximate $W_p : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ to an arbitrary ϵ additive error. However, this approach has the following issues: while an MLP can approximate W_p to an arbitrary ϵ additive error, the model complexity of the MLP depends on not only the approximation error, but also on the maximum size N of the input point set. In contrast, our neural network $\mathcal{N}_{\text{ProductNet}}$ introduced in Eqn (2) has model complexity which is independent of N ; see Corollary 3.6. Furthermore, the function represented by the MLP (that we use to approximate ϕ) however may not be symmetric to the two point sets, and may not guarantee permutation invariance with respect to $G \times G$. In practice, computationally expensive techniques such as data augmentation are often required in order for an unconstrained MLP to approximate a structured function such as our Wasserstein distance function, which is an SFGI-function.

Comparison with Siamese architectures. As mentioned previously, one common way of learning product functions is to use a Siamese network which embeds \mathcal{X} to Euclidean space and then use some ℓ_p norm to approximate the desired product function. Consider the learning of Wasserstein distances again. To approximate $W_p(A, B)$ for two point sets A and B . The Siamese network will first embed each point set to a fixed dimensional Euclidean space \mathbb{R}^D via a function ϕ_θ modeled by a neural network, and then compute $\|\phi_\theta(A) - \phi_\theta(B)\|_q$ for some $1 \leq q < \infty$. Intuitively, compared to our neural network $\mathcal{N}_{\text{ProductNet}}$ introduced in Eqn (2) (and recall Figure 1), the Siamese network will replace the outer function ρ by simply the L_q -norm of $\phi_\theta(A) - \phi_\theta(B)$. However, this approach intuitively requires that one can find a near-isometric embedding of the Wasserstein distance to the L_q distance in a Euclidean space. However, there exists lower bound results on the distortion incurred when embedding Wasserstein distance to L_p space. In the following section, we will review one such result on the lower bound of metric distortions when embedding the Wasserstein distance to L_1 space; implying that if we chose $q = 1$, then in the worst case the Siamese network will incur at least as much distortion as that lower bound.

C Non-embeddability theorems for Wasserstein distances

Here we summarize results pertaining to the limitations of embedding Wasserstein distance to the L_q distance in a Euclidean space. Consider probability distributions over a grid of points in \mathbb{R}^2 , $\{0, 1, \dots, D\}^2$ equipped with the 1-Wasserstein distance, $(\mathcal{P}(\{0, 1, \dots, D\}^2), W_1)$. Let L_1 denote the space of Lebesgue measureable functions $f : [0, 1] \rightarrow \mathbb{R}$ such that

$$\|f\|_1 = \int_0^1 |f(t)| dt.$$

Given a mapping $F : (\mathcal{P}(\{0, 1, \dots, n\}^2), W_p) \rightarrow L_1$ such that for any $\mu, \rho \in (\mathcal{P}(\{0, 1, \dots, D\}^2), W_p)$,

$$W_p(\mu, \rho) \leq \|F(\mu) - F(\rho)\| \leq L \cdot W_p(\mu, \rho) \quad (1)$$

the *distortion* is the value of L .

Theorem C.1 ([9]) Any embedding $(\mathcal{P}(\{0, 1, \dots, n\}^2), W_1) \rightarrow L_1$ must incur distortion at least $\Omega(\sqrt{\log D})$

For $(\mathcal{P}(\{0, 1, \dots, D\}^d), W_1)$ where $d \geq 2$, $\mathcal{P}(\{0, 1, \dots, D\}^d)$ contains $\mathcal{P}(\{0, 1, \dots, D\}^2)$ so the lower bound $O(\sqrt{\log D})$ still applies for L_1 embeddings from $\mathcal{P}(\{0, 1, \dots, D\}^d)$, $d > 2$.

From the above results, we can see that for any Siamese architecture (even in the simple case of finite point sets in \mathbb{R}^2 and \mathbb{R}^3), we are unable to approximate the Wasserstein distances via any L_1 embedding.

Corollary C.2 Given a neural network $\mathcal{N}_{\text{Siamese}} : \mathcal{X} \rightarrow \mathbb{R}^d$ where \mathcal{X} are weighted point sets over $\{0, 1, \dots, D\}^2$, $\|\mathcal{N}_{\text{Siamese}}(\mu) - \mathcal{N}_{\text{Siamese}}(\nu)\|_1$ incurs distortion at least $\Omega(\sqrt{\log D})$.

Note that if we consider our input for a Siamese architecture to be finite point sets over $\{0, 1, \dots, D\}^2$, we allow multisets so the input set size is not bounded by D .

Table 1: Size of the output point set from the WPCE decoder for each dataset.

	noisy-sphere-2	noisy-sphere-6	uniform	ModelNet-small	ModelNet-large	RNaseq
Size	200	200	256	100	2048	100

D Experimental details

D.1 Baseline models for comparison with $\mathcal{N}_{\text{ProductNet}}$

Here, we will detail two baseline neural networks in our experiments.

Wasserstein point cloud embedding (WPCE) network First defined by [6], the WPCE network is an Siamese autoencoder architecture. It consists of an encoder network $\mathcal{N}_{\text{encoder}}$ and a decoder network $\mathcal{N}_{\text{decoder}}$. WPCE takes as input two point sets $P, Q \subseteq \mathbb{R}^d$. $\mathcal{N}_{\text{encoder}}$ is a permutation invariant neural network - which we chose to be DeepSets. In other words,

$$\mathcal{N}_{\text{encoder}}(P) = \phi_{\theta_2} \left(\sum_{x \in P} h_{\theta_1}(x) \right).$$

Note that one may also choose PointNet to be $\mathcal{N}_{\text{encoder}}$. However, in our experiments, we did not see a large difference in approximation quality between using PointNet and DeepSets (where the difference between the two amounts to using a sum vs. max aggregator). For sake of consistency with $\mathcal{N}_{\text{ProductNet}}$, we chose to use a sum aggregator for $\mathcal{N}_{\text{encoder}}$ (DeepSets). The decoder network, $\mathcal{N}_{\text{decoder}}$, is a fully connected neural network which outputs a fixed-size point set in \mathbb{R}^d . WPCE then uses a Sinkhorn reconstruction loss term to regularize the embedding produced by the encoder. Thus, given a set of paired input point sets and their Wasserstein distance $\{(P_i, Q_i, W_1(P_i, Q_i)) : i \in [N]\}$ the loss which we optimize over for WPCE is

$$\begin{aligned} \mathcal{L}(P, Q) = & \frac{1}{N} \sum \left(\|\mathcal{N}_{\text{encoder}}(P_i) - \mathcal{N}_{\text{decoder}}(Q_i)\|_2 - W_1(P_i, Q_i) \right)^2 \\ & + \lambda \left(\frac{1}{N} \sum S_\epsilon(\mathcal{N}_{\text{decoder}}(\mathcal{N}_{\text{encoder}}(P_i))) + \frac{1}{N} \sum S_\epsilon(\mathcal{N}_{\text{decoder}}(\mathcal{N}_{\text{encoder}}(Q_i))) \right) \end{aligned} \quad (2)$$

where λ is some constant in \mathbb{R} that controls the balance between the two loss terms and S_ϵ is the Sinkhorn divergence between P_i and Q_i . Note that ϵ in S_ϵ refers to the regularization parameter for the Sinkhorn divergence. For our experiments, we chose $\lambda = 0.1$ and $\epsilon = 0.1$. Furthermore, for each dataset used in our experiments, we set used a range of different sizes for the fixed-size output point set. This parameter is summarized per dataset in Table 1.

Siamese DeepSets. The baseline Siamese DeepSets model, $\mathcal{N}_{\text{SDepSets}}(\cdot, \cdot)$, consists of a single DeepSets model which maps two input point sets to some Euclidean space \mathbb{R}^d . The ℓ_2 -norm between the final embeddings is then used as the final estimate for Wasserstein distance. Formally, let $\mathcal{N}_{\text{SDepSets}}(P) = \phi_{\theta_2} \left(\sum_{x \in P} h_{\theta_1}(x) \right)$, where ϕ_{θ_2} and h_{θ_1} are both MLPs, be the DeepSets model. Then given two point sets P, Q , the final approximation of Wasserstein distance given by $\mathcal{N}_{\text{SDepSets}}(P, Q)$ is

$$\mathcal{N}_{\text{SDepSets}}(P, Q) = \|\mathcal{N}_{\text{SDepSets}}(P) - \mathcal{N}_{\text{SDepSets}}(Q)\|_2.$$

D.2 Training and implementation details

Datasets. We used several synthetic datasets as well as the ModelNet40 point cloud dataset. Furthermore, we used two different types of synthetic datasets. We construct the ‘noisy-sphere- d ’ dataset by sampling pairs of point clouds from four d -dimensional spheres centered at the origin with increasing radiuses of 0.25, 0.5, 0.75, and 1.0. For our experiments, we used ‘noisy-sphere-3’ and ‘noisy-sphere-6’. Finally, the ‘uniform’ dataset of points sets in \mathbb{R}^2 is constructed by sampling points sets from the uniform distribution on $[-4, 4] \times [-4, 4]$. The full details and names of each dataset are summarized in Table 2.

Table 2: Summary of details for each dataset. Note that min and max refer to the minimum and maximum number of points per input point set.

	dim	min	max	training pairs	val pairs
noisy-spheres-2	3	100	300	2000	200
noisy-spheres-6	6	100	300	3600	400
uniform	2	256	256	3000	300
ModelNet-small	3	20	200	3000	300
ModelNet-large	3	2048	2048	4000	400
RNAseq	2000	20	200	3000	300

Implementation. We implement all neural network architectures using the PyTorch[10] and GeomLoss [3] libraries. The ground truth 1-Wasserstein distances and Sinkhorn approximations were computed using Python Optimal Transport (POT) [4]. Note that for large point sets and for higher dimensional datasets, there is often a high degree of numerical instability in the POT implementation of the Sinkhorn algorithm. In these cases (ModelNet-large and RNAseq) we used our own implementation of the Sinkhorn algorithm. For each model, we used the Leaky-ReLU as the non-linearity. To train each model, we set the batch size for each dataset to be 128 and the learning rate to 0.001. All models were trained on an Nvidia RTX A6000 GPU. For both $\mathcal{N}_{\text{ProductNet}}$ and $\mathcal{N}_{\text{SDeepSets}}$, given two input point sets, we minimize on the mean squared error between the approximation produced by the network and the true Wasserstein distance. In other words, given two point sets P, Q , the loss for $\mathcal{N}_{\text{ProductNet}}$ is defined as $\mathcal{L}_{\mathcal{N}_{\text{ProductNet}}}(P, Q) = \left(\mathcal{N}_{\text{ProductNet}}(P, Q) - W_1(P, Q) \right)^2$

and the loss for Siamese DeepSets is $\mathcal{L}_{\mathcal{N}_{\text{SDeepSets}}}(P, Q) = \left(\mathcal{N}_{\text{SDeepSets}}(P, Q) - W_1(P, Q) \right)^2$. For WPCE, we train the network using the loss function defined in Eq. 2. Note that for $\mathcal{N}_{\text{ProductNet}}$, the hyperparameters are the width, depth, and output dimension for the MLPs which represent h_{θ_1} and ϕ_{θ_2} and the width and depth the MLP which represents ρ_{θ_3} . For WPCE, we set the decoder to a three layer neural network with width 100 and adjusted the width, depth, and output dimension for the MLPs which represent ϕ_{θ_2} and h_{θ_1} in $\mathcal{N}_{\text{encoder}}$. To find the best model for each architecture, we randomly sampled hyperparameter configurations and conducted a hyperparameter search over 85 models for $\mathcal{N}_{\text{ProductNet}}$ and 75 models for both WPCE and $\mathcal{N}_{\text{SDeepSets}}$.

D.3 Approximating 2-Wasserstein distance

To further show the use of our model, we additionally approximate the 2-Wasserstein distance; see Table 5 for results. The experimental set-up is the same as for 1-Wasserstein distance and we largely see the same trends as we see for 1-Wasserstein distance; that is, $\mathcal{N}_{\text{ProductNet}}$ outperforms all other neural network implementations. Note that Table 5 shows that the Sinkhorn approximation with $\epsilon = 0.01$ is more accurate than $\mathcal{N}_{\text{ProductNet}}$. However, as the ϵ parameter for the Sinkhorn approximation decreases, the computation time increases. In Table ??, we show that the Sinkhorn approximation is already much slower than $\mathcal{N}_{\text{ProductNet}}$ at $\epsilon = 0.1$ while also having a less accurate approximation. The Sinkhorn approximation with $\epsilon = 0.01$ (reported in Table 5) is slower the Sinkhorn approximation with $\epsilon = 0.1$ and additionally, is also much slower than $\mathcal{N}_{\text{ProductNet}}$ at inference time.

D.4 Generalization to large point sets

In addition to the results reported in Tables 1 and 5, which record the average relative error for point sets of unseen sizes during training, we have also included several plots in Figures 2 and 1 which demonstrate how the error for each approximation method changes as the input point set size increases for both 1-Wasserstein distance and 2-Wasserstein distance. Observe that for ModelNet-small, noisy-sphere-3, and noisy-sphere-6, the error for $\mathcal{N}_{\text{ProductNet}}$ exhibits a significantly slower increase as compared to WPCE. The rate at which the error increases is not as evident for the uniform and ModelNet-large datasets. It is worth mentioning that we trained the model on fixed-size input for both of these datasets. It is possible that training with fixed-size input leads a rapid deterioration in approximation quality for WPCE and $\mathcal{N}_{\text{SDeepSets}}$ when dealing with point sets of sizes not seen at training time. Furthermore, consider that WPCE may be especially sensitive to differences in input

Table 3: Time for 500 1-Wasserstein distance computations in seconds. Note that we chose input point set size to be the maximum possible point set size that we trained on. Additionally, the Sinkhorn distance reported uses $\epsilon = 0.1$ as the regularization parameter. Note that as ϵ decreases, the error incurred by the Sinkhorn approximation will decrease but the computation time will also increase. Here, when $\epsilon = 0.1$, the Sinkhorn approximation is already much slower than the neural network approximations while being less accurate.

Dataset	Input size	Models			Sinkhorn	Ground truth
		$\mathcal{N}_{\text{ProductNet}}$	WPCE	$\mathcal{N}_{\text{SDepSets}}$		
noisy-sphere-3	300	1.050	0.676	0.4904	2.591	2.813
noisy-sphere-6	300	0.752	0.491	0.503	1.986	6.6770
uniform	256	0.155	0.184	0.137	15.113	1.018
ModelNet-small	200	0.330	0.330	0.191	2.074	1.615
ModelNet-large	2048	1.174	1.571	0.612	239.448	254.947
RNAseq	200	1.128	0.856	0.792	92.153	105.908

Table 4: Comparison of the mean relative error versus overall computation time for 500 approximations of 1-Wasserstein distance for $\mathcal{N}_{\text{ProductNet}}$ and the Sinkhorn distance. Note that the input point set size is the same as in Table 3 for each dataset. The parameter ϵ controls the accuracy of the Sinkhorn approximation with lower ϵ corresponding to a more accurate approximation once the Sinkhorn algorithm converges. However, notice that in some cases, the Sinkhorn algorithm with $\epsilon = 0.01$ has a higher relative error than the Sinkhorn algorithm with $\epsilon = 0.1$ as the algorithm fails to converge within a reasonable number of iterations (1000).

Dataset		$\mathcal{N}_{\text{ProductNet}}$	Sinkhorn ($\epsilon = 0.10$)	Sinkhorn ($\epsilon = 0.01$)
ModelNet-small	Error	0.084 ± 0.077	0.187 ± 0.232	0.011 ± 0.003
	Time (s)	0.330	2.074	104.712
ModelNet-large	Error	0.140 ± 0.206	0.148 ± 0.048	0.026 ± 0.008
	Time (s)	1.174	239.448	1930.852
Uniform	Error	0.097 ± 0.073	0.073 ± 0.009	0.023 ± 0.098
	Time (s)	0.155	15.113	63.028
noisy-sphere-3	Error	0.046 ± 0.043	0.187 ± 0.232	0.162 ± 0.132
	Time (s)	1.050	2.591	214.185
noisy-sphere-6	Error	0.015 ± 0.014	0.137 ± 0.122	0.326 ± 0.135
	Time (s)	0.752	1.986	101.763
RNAseq	Error	0.012 ± 0.010	0.040 ± 0.009	0.035 ± 0.013
	Time (s)	1.128	92.153	91.573

sizes at testing time as training WPCE depends on minimizing the Wasserstein difference between the input point set and a fixed-size decoder point set which may cause the model to be overly specialized to point sets of a fixed input size. This observation could provide an explanation for the observed plots in both the ModelNet-large and uniform cases. Finally, as predicted in our theoretical analysis, the performance of the model degrades for higher dimensional datasets i.e. the RNAseq dataset.

E Extra proofs

E.1 Proof of DeepSets Universality.

Here we will provide extra details on the proof of Theorem 2.1 using multisymmetric polynomials. Note that multisymmetric polynomials were previously used in [12] and [8] to show universality for equivariant set networks and arbitrary G -invariant neural networks for any permutation group G .

Proof: To begin, we will first define multisymmetric polynomials and power sum multisymmetric polynomials. Let $A[y_1, \dots, y_n]$ be the ring of polynomials in n variables with coefficients in a ring A .

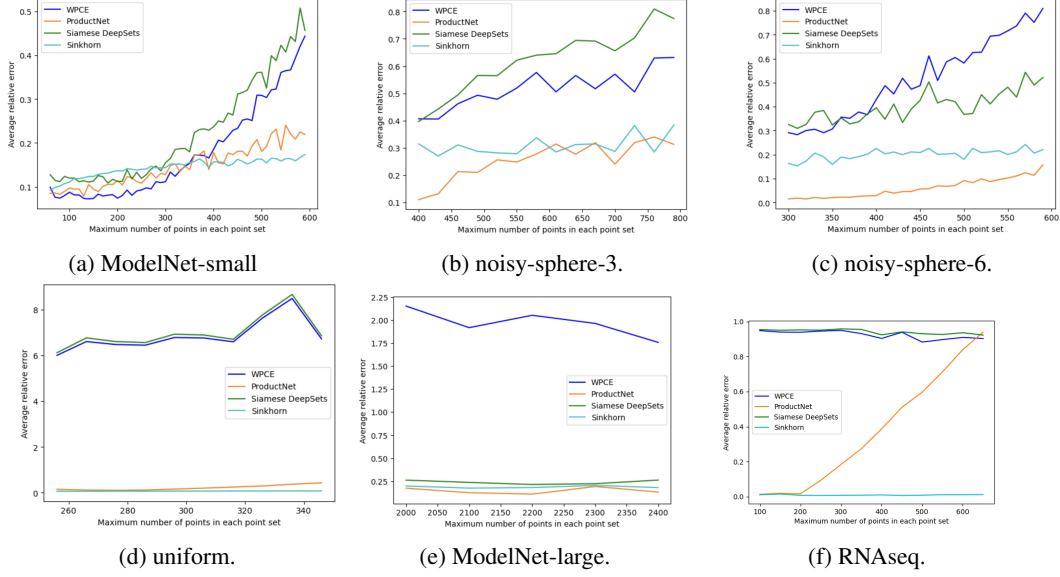


Figure 1: Average error for 1-Wasserstein approximation for each model as the maximum number of points increases. Note that the Sinkhorn approximation is the $\epsilon = 0.1$. These graphs include point set sizes not seen at training time to display how each approximation performs on unseen examples. Note that especially for ModelNet-small, noisy-sphere-3, and noisy-sphere-6, we can see that the error for $\mathcal{N}_{\text{ProductNet}}$ increases at a slower rate than WPCE.

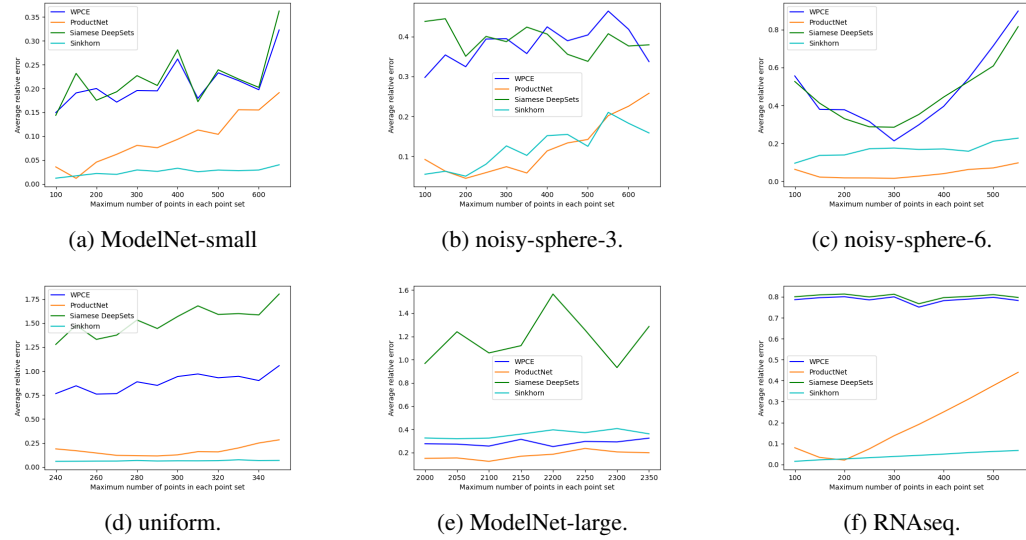


Figure 2: Average 2-Wasserstein error for each model as the maximum number of points increases. Note that for all datasets, the Sinkhorn error is with $\epsilon = 0.10$.

Table 5: Mean relative error between approximations and 2-Wasserstein distance between point sets. The top row for each dataset shows the error for point sets with input sizes that were seen at training time; while the bottom row shows the error for point sets with input sizes that were not seen at training time. Note that the Sinkhorn approximation is computed with the regularization parameter set to $\epsilon = 0.01$.

Dataset	Input size	$\mathcal{N}_{\text{ProductNet}}$	WPCE	$\mathcal{N}_{\text{SDepSets}}$	Sinkhorn
noisy-sphere-3	[100, 300]	0.054 \pm 0.071	0.291 \pm 0.201	0.400 \pm 0.336	0.078 \pm 0.186
	[400, 600]	0.188 \pm 0.197	0.387 \pm 0.386	0.427 \pm 0.375	0.161 \pm 0.311
noisy-sphere-6	[100, 300]	0.024 \pm 0.010	0.331 \pm 0.237	0.358 \pm 0.231	0.019 \pm 0.057
	[400, 600]	0.092 \pm 0.074	0.434 \pm 0.598	0.623 \pm 0.596	0.050 \pm 0.039
uniform	256	0.112 \pm 0.082	0.221 \pm 0.162	0.241 \pm 0.171	0.182 \pm 0.044
	[200, 300]	0.175 \pm 0.123	2.431 \pm 2.162	4.058 \pm 3.324	0.055 \pm 0.053
ModelNet-small	[20, 200]	0.078 \pm 0.095	0.178 \pm 0.148	0.183 \pm 0.148	0.023 \pm 0.059
	[400, 600]	0.163 \pm 0.151	0.216 \pm 0.252	0.227 \pm 0.179	0.034 \pm 0.031
ModelNet-large	2048	0.187 \pm 0.335	0.281 \pm 0.203	0.538 \pm 0.298	0.172 \pm 0.065
	[1800, 2000]	0.185 \pm 0.302	0.523 \pm 0.526	33.086 \pm 28.481	0.046 \pm 0.039
RNAseq	[20, 200]	0.049 \pm 0.029	0.508 \pm 0.291	0.490 \pm 0.271	0.024 \pm 0.009
	[400, 600]	0.281 \pm 0.057	0.533 \pm 0.300	0.568 \pm 0.317	0.987 \pm 0.0002

Table 6: Comparison of the mean relative error versus overall computation time for 300 approximations of 2-Wasserstein distance for $\mathcal{N}_{\text{ProductNet}}$ and the Sinkhorn distance. The parameter ϵ controls the accuracy of the Sinkhorn approximation with lower ϵ corresponding to a more accurate approximation.

Dataset		$\mathcal{N}_{\text{ProductNet}}$	Sinkhorn ($\epsilon = 0.10$)	Sinkhorn ($\epsilon = 0.01$)
ModelNet-small	Error	0.078 \pm 0.097	0.232 \pm 0.132	0.019 \pm 0.057
	Time (s)	0.208	1.165	9.857
ModelNet-large	Error	0.187 \pm 0.335	0.363 \pm 0.255	0.172 \pm 0.089
	Time (s)	2.841	6.079	36.265
uniform	Error	0.112 \pm 0.082	0.0303 \pm 0.022	0.182 \pm 0.044
	Time (s)	0.712	16.515	29.312
noisy-sphere-3	Error	0.054 \pm 0.071	0.225 \pm 0.093	0.078 \pm 0.186
	Time (s)	0.677	1.591	12.760
noisy-sphere-6	Error	0.024 \pm 0.010	0.324 \pm 0.316	0.023 \pm 0.059
	Time (s)	0.428	0.877	9.831
RNAseq	Error	0.049 \pm 0.029	0.031 \pm 0.014	0.024 \pm 0.009
	Time (s)	0.716	47.701	81.016

Definition E.1 (Multisymmetric polynomials) *The multisymmetric polynomials on the n families of k variables $\mathbf{x}_1, \dots, \mathbf{x}_n$ where $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,k})$ are those polynomials that remain unchanged under every permutation of the n families, $\mathbf{x}_1, \dots, \mathbf{x}_n$.*

Let A be a ring. The algebra of multisymmetric polynomials in n families of k variables with coefficients in A is denoted $\mathfrak{J}_n^k(A)$.

Furthermore, we define the multisymmetric power-sum polynomials:

Definition E.2 (Multisymmetric power-sum polynomials) *Let $\alpha = (\alpha_1, \dots, \alpha_k) \in \mathbb{N}^k$. Given $\mathbf{x} = (x_1, \dots, x_k)$, let*

$$\mathbf{x}^\alpha = x_1^{\alpha_1} \cdots x_k^{\alpha_k}$$

The multisymmetric power sum with multi-degree α is

$$p_\alpha = \sum_i^n \mathbf{x}_i^\alpha$$

Among them, we will consider the set of elementary multisymmetric power sums to be those with $|\alpha| \leq n$.

Notice that there $t = \binom{n+k}{k}$ multisymmetric power sums. Let $\alpha_1, \dots, \alpha_t$ be the list of all $\alpha \in \mathbb{N}^k$ such that $|\alpha| \leq n$.

Theorem E.3 ([1]) *Let A be a ring in which $n!$ is invertible. Then the multisymmetric power sums generate $\mathfrak{J}_n^r(A)$ as an A -algebra.*

If we take $A = \mathbb{R}$, we get that the multisymmetric power sum polynomials generate \mathfrak{J}_N^k . Now given a continuous function $f : \mathbb{R}^{k \times N} \rightarrow \mathbb{R}$ which is invariant to permutations of the columns, we know that f can be approximated by a polynomial p which is invariant to permutations of columns (see [8] for a detailed argument). Such a polynomial p is a multisymmetric polynomial in N families of k variables with coefficients in \mathbb{R} i.e., $p \in \mathfrak{J}_N^k$. Given $\mathbf{x} \in \mathbb{R}^k$,

$$\phi(\mathbf{x}) = [x^{\alpha_1}, \dots, x^{\alpha_t}]$$

Then

$$\sum_{i=1}^N \phi(\mathbf{x}_i) = \begin{bmatrix} \sum_{i=1}^N x_i^{\alpha_1} \\ \sum_{i=1}^N x_i^{\alpha_2} \\ \vdots \\ \sum_{i=1}^N x_i^{\alpha_t} \end{bmatrix} = \begin{bmatrix} p_{\alpha_1} \\ p_{\alpha_2} \\ \vdots \\ p_{\alpha_t} \end{bmatrix}$$

By Theorem E.3, we have that $p_{\alpha_1}, \dots, p_{\alpha_t}$ will generate any polynomial in \mathfrak{J}_N^k . Then we have some polynomial $q \in \mathbb{R}[y_1, \dots, y_t]$ such that $p = q(p_{\alpha_1}, \dots, p_{\alpha_t})$. ■

E.2 Proofs from Section 3

E.2.1 Proof of Lemma 3.2

Proof: Let $\epsilon > 0$. Since f is uniformly continuous, $\exists \delta$ such that for all $(A_1, \dots, A_m), (A'_1, \dots, A'_m) \in \mathcal{X}_1 \times \dots \times \mathcal{X}_m$ where $d_{\mathcal{X}_1 \times \dots \times \mathcal{X}_m}((A_1, \dots, A_m), (A'_1, \dots, A'_m)) < \delta$, we have $|f(A_1, \dots, A_m) - f(A'_1, \dots, A'_m)| < \epsilon$. Since for any $\delta > 0$ and any $i \in [m]$, $(\mathcal{X}_i, d_{\mathcal{X}_i})$ has a (δ, a, G) -sketch, we know that there is an $a_i \in \mathbb{N}^+$ where there are continuous $h_i : \mathcal{X}_i \rightarrow \mathbb{R}^{a_i}$ and $g_i : \mathbb{R}^{a_i} \rightarrow \mathcal{X}_i$ where $d_{\mathcal{X}_i}(g_i \circ h_i(A), A) < \delta/m$ for each $A \in \mathcal{X}_i$.

Let $g' : \mathbb{R}^{a_1} \times \dots \times \mathbb{R}^{a_m} \rightarrow \mathcal{X}_1 \times \dots \times \mathcal{X}_m$ be defined as $(u_1, \dots, u_m) \mapsto (g_1(u_1), \dots, g_m(u_m))$. Since $d_{\mathcal{X}_i}(g_i \circ h_i(A_i), A_i) < \delta/k$ for $A_i \in \mathcal{X}_i$ and $i \in [m]$,

$$d_{\mathcal{X}_1 \times \dots \times \mathcal{X}_m}((g_1 \circ h_1(A_1), \dots, g_m \circ h_m(A_m)), (A_1, \dots, A_m)) < \delta.$$

Let $\rho = f \circ g'$ and $\phi_i = h_i$. Then

$$\begin{aligned} |f(A_1, \dots, A_m) - \rho(\phi_1(A_1), \dots, \phi_m(A_m))| &= |f(A_1, \dots, A_m) - f \circ g'(h_1(A_1), \dots, h_m(A_m))| \\ &= |f(A_1, \dots, A_m) - f(g \circ h(A_1), \dots, g \circ h(A_m))| < \epsilon. \end{aligned}$$

Note that if $\mathcal{X}_1 = \mathcal{X}_2 = \dots = \mathcal{X}_m$ and $G_1 = G_2 = \dots = G_m$, we can use the same encoding and decoding function - h_i and g_i , respectively - for all \mathcal{X}_i . Thus, in this case, $\phi_1 = \phi_2 = \dots = \phi_m$. ■

E.2.2 Proof of Lemma 3.3

Proof: Using the same argument as Theorem 3.2, we know that for $\epsilon/2$, there is a continuous $h : \mathcal{X} \rightarrow \mathbb{R}^a$ and $g : \mathbb{R}^a \rightarrow \mathbb{R}$ such that

$$|f(A_1, \dots, A_m) - f(g \circ h(A_1), \dots, g \circ h(A_m))| < \frac{\epsilon}{2}.$$

As before, let $g' : \mathbb{R}^{a \times m} \rightarrow \mathbb{R}$ be $g'(u_1, \dots, u_m) = (g(u_1), \dots, g(u_m))$. Take $F : \mathbb{R}^{a \times m} \rightarrow \mathbb{R}$ as $F(u_1, \dots, u_m) = f \circ g'(u_1, \dots, u_m) = f(g(u_1), \dots, g(u_m))$ where u_i represents the i th column of an element in $\mathbb{R}^{a \times m}$. Note that F is continuous and invariant to permutations of the columns. Let $t = \binom{a+m}{m}$. Therefore, by Theorem 2.1, there is a $\gamma : \mathbb{R}^a \rightarrow \mathbb{R}^t$ and ρ such that

$$\left| f \circ g'(v_1, \dots, v_m) - \rho\left(\sum_{i=1}^m \gamma(v_i)\right) \right| < \frac{\epsilon}{2}$$

Now set $\phi = \gamma \circ h$ which is a function from $\mathbb{R}^a \rightarrow \mathbb{R}^t$. We thus have that

$$\begin{aligned}
& \left| f(A_1, \dots, A_m) - \rho\left(\sum_{i=1}^m \phi(A_i)\right) \right| \\
&= \left| f(A_1, \dots, A_m) - f \circ g'(h(A_1), \dots, h(A_m)) + f \circ g'(h(A_1), \dots, h(A_m)) - \rho\left(\sum_{i=1}^m \phi(A_i)\right) \right| \\
&< \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon
\end{aligned}$$

This completes the proof. ■

References

- [1] E. Briand. When is the algebra of multisymmetric polynomials generated by the elementary multisymmetric polynomials? *Beiträge zur Algebra und Geometrie: Contributions to Algebra and Geometry*, 45 (2), 353-368., 2004.
- [2] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah. Signature verification using a "siamese" time delay neural network. *Advances in neural information processing systems*, 6, 1993.
- [3] J. Feydy, T. Séjourné, F.-X. Vialard, S.-i. Amari, A. Trounev, and G. Peyré. Interpolating between optimal transport and mmd using sinkhorn divergences. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2681–2690, 2019.
- [4] R. Flamary, N. Courty, A. Gramfort, M. Z. Alaya, A. Boissunon, S. Chambon, L. Chapel, A. Corenflos, K. Fatras, N. Fournier, L. Gautheron, N. T. Gayraud, H. Janati, A. Rakotomamonjy, I. Redko, A. Rolet, A. Schutz, V. Seguy, D. J. Sutherland, R. Tavenard, A. Tong, and T. Vayer. Pot: Python optimal transport. *Journal of Machine Learning Research*, 22(78):1–8, 2021.
- [5] R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006.
- [6] K. Kawano, S. Koide, and T. Kutsuna. Learning wasserstein isometric embedding for point clouds. In *2020 International Conference on 3D Vision (3DV)*, pages 473–482. IEEE, 2020.
- [7] G. Koch, R. Zemel, R. Salakhutdinov, et al. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, volume 2. Lille, 2015.
- [8] H. Maron, E. Fetaya, N. Segol, and Y. Lipman. On the universality of invariant networks. In *International conference on machine learning*, pages 4363–4371. PMLR, 2019.
- [9] A. Naor and G. Schechtman. Planar earthmover is not in l_1 . *SIAM Journal on Computing*, 37(3):804–826, 2007.
- [10] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [11] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017.
- [12] N. Segol and Y. Lipman. On universal equivariant set networks. *arXiv preprint arXiv:1910.02421*, 2019.