# Physics Informed Convex Artificial Neural Networks (PICANNs) for Optimal Transport based Density Estimation

Amanpreet Singh[1], Sarang Joshi[1], and Martin Bauer[2]*

[1] University of Utah, Salt Lake City, UT 84112, USA
[2] Florida State University, Tallahassee, FL 32306, USA

**Abstract.** Optimal Mass Transport (OMT) is a well studied problem with a variety of applications in a diverse set of fields ranging from Physics to Computer Vision and in particular Statistics and Data Science. Since the original formulation of Monge in 1781 significant theoretical progress been made on the existence, uniqueness and properties of the optimal transport maps. The actual numerical computation of the transport maps, particularly in high dimensions, remains a challenging problem. By Brenier's theorem, the continuous OMT problem can be reduced to that of solving a non-linear PDE of Monge-Ampere type whose solution is a convex function. In this paper, building on recent developments of input convex neural networks and physics informed neural networks for solving PDE's, we propose a Deep Learning approach to solve the continuous OMT problem. To demonstrate the versatility of our framework we focus on the ubiquitous density estimation and generative modeling tasks in statistics and machine learning. Finally as an example we show how our framework can be incorporated with an autoencoder to estimate an effective probabilistic generative model.

## 1 Introduction

Density estimation and random sampling are fundamental problems in machine learning and statistical inference. The density estimation problem is to estimate a smooth probability density based on a discrete finite set of observations.

In traditional parametric density estimation techniques, it is assumed that the data is drawn from a known parametric family of distributions and it only remains to best estimate these parameters. These methods require that one has a basis to believe that the data is indeed derived from a specific family of distributions and are consequently limited in their applicability to many modern tasks. One of the most ubiquitous parametric techniques is the Gaussian Mixture Modeling GMM [23].

Non-parametric techniques were first proposed by Fix and Hodges [12,32] to move away from such rigid distributional assumptions. The most used approach is the kernel density estimation, which dates back to Rosenblatt [29] and Parzen [26]. Despite decades of work in this field, there remain many challenges regarding the implementation and practical performance of kernel density estimators. This includes in particular

the bandwidth selection and the lack of local adaptivity resulting in a large sensitivity to outliers [18]. This is particularly exacerbated in high dimensions with the curse of dimensionality.

Recently, diffeomorphic transformation based algorithms have been proposed to tackle this problem [11,21,35,3]. The basic concept of transformation based algorithms is to find a diffeomorphic mapping between a reference probability distribution and the unknown target distribution, from which the observed data is drawn. Consequently, transformation based density estimation leads at the same time to an efficient generative model, as new samples from the estimated density can be generated at a low cost by sampling from the reference density and transforming the samples by the estimated transformation.

The fundamental problem in diffeomorphic transformation based approaches is how to estimate the transformation: from a theoretical point of view there exists an infinite set of transformations that map two given probability densities onto each other. In Real-NVP [11] the transformations are restricted to the class of diffeomorphisms with triangular Jacobians that are easy to invert. This is closely related to the Knothe-Rosenblatt rearrangement [16,30]. Optimal mass transport [34,33] on the other hand formulates the transport map selection as the minimizer of a cost function. The optimal transportation cost induces a metric structure, the Wasserstein metric, on the space of probability densities and is sometimes referred to as the Earth Mover's Distance. This theory dates back to 1781, where it was originally formulated by the French Mathematician Gaspard Monge [24]. The difficulty in applying this framework to the proposed density estimation problem lies in solving the corresponding optimization problem, which in dimension greater than one is highly non-trivial. Note, that the fully discrete OMT problem (optimal assignment problem) can be solved using linear programming and can be approximated by the Sinkhorn algorithm [9,25]. However, this does not lead to a continuous transformation map and thus it can't be used for the proposed diffeomorphic density estimation and generative modelling. Previous algorithmic solutions for the continuous OMT problem include fluid mechanics based approaches [4], finite element or finite difference based methods [6,5] and steepest descent based energy minimization approaches [2,8,19]. Recent work by Makkuva et.al. [20] proposed to approximate the OMT map as the solution of min-max optimization using input convex neural networks [1].

**Contributions:** In this paper we propose a different deep learning based framework to approximate the optimal transport maps. The approach we present relies on Brenier's celebrated theorem [7], thereby reducing the optimal transport problem to that of solving a partial differential equation: a Monge-Ampere type equation. We frame this PDE in the recently developed paradigm of Physical Informed Neural Networks (PINNs) [28]. Using this framework we directly inherit the dimensional scalability of neural networks [31], which traditional finite element or finite difference methods for solving PDEs do not posses. Brenier's theorem further states that the optimal transport map is given by the gradient of a convex function- the Brenier potential. To incorporate this information in our PINN approach, we parametrize the Brenier potential using an Input Convex Neural Network (ICNN) [1] architecture thereby guaranteeing it's convexity.

We quantify the accuracy of our OMT solver on numerous synthetic examples for which analytical solutions are known. We show that the deep learning based approach approximates the true solution well, even in high dimensions. As an explicit application of our solution of OMT we focus on the density estimation problem. In synthetic examples we show that one can estimate the true density based on a limited amount of samples. We also demonstrate the generative power of our framework by combining it with a traditional autoencoder and applying it to the MNIST data set.

In accordance with the best practices for reproducible research we are providing an open-source version of the code, is publicly available on github.

## 2  OMT using Deep Learning

In this section we will present our framework for solving the Optimal Mass Transport (OMT) problem. Our approach will combine methods of deep learning with the celebrated theorem of Brenier, which reduces the solution of the OMT problem to solving a Monge Ampere type equation. To be more precise, we will tackle this problem by embedding the Monge Ampere equation into the broadly applicable concept of Physics Informed Neural Networks.

### 2.1  Mathematical Background of OMT

We start by summarizing the mathematical background of OMT, including a description of Brenier's theorem. For more information we refer to the vast literature on OMT, including [33,34].

Let $\Omega$ be a convex and bounded domain of $\mathbb{R}^n$ and let $dx$ denote the standard measure on $\mathbb{R}^n$. For simplicity we restrict our presentation to the set $\mathcal{P}(\Omega)$ of all absolutely continuous measures on $\Omega$, i.e., $\mathcal{P}(\Omega) \ni \mu = f dx$ with $f \in L^1(\Omega)$, such that $\int_\Omega f dx = 1$. From here on we will often identify the measure $\mu$ with it's density function $f$.

We aim to minimize the cost of of transporting a density $\mu$ to a density $\nu$ using a (transport) map $T$, which leads to the so-called Monge Optimal Transport Problem. To keep the presentation as simple as possible we will only consider the special case of a quadratic cost function.

**Definition 2.1 ($L^2$-Monge Optimal Transport Problem)** *Given $\mu, \nu \in \mathcal{P}(\Omega)$, minimize*

$$\mathbb{M}(T) = \int_\Omega \|x - T(x)\|^2 d\mu(x)$$

*over all $\mu$-measureable maps $T : \Omega \to \Omega$ subject to $\nu = T_*\mu$. We will call an optimal $T$ an optimal transport map.*

Here the constraint is formulated in terms of the push forward action of a measurable map $T : \Omega \to \Omega$, which is defined via

$$T_*\mu(B) = \mu(T^{-1}(B)), \tag{1}$$

for every measurable set $A \subset \Omega$. By a change of coordinates the constraint $T_*\mu = T_*(f dx) = \nu = g dx$ can be thus reduced to the equation

$$f(x) = g(T(x))|\det(DT(x))|. \tag{2}$$

The above equation can be also expressed via the pullback action as $\mu = T^*\nu$.

The existence of an optimal transport map is not always guaranteed. We will, however, see, that in our situation, i.e., for absolutely continuous measures the existence and uniqueness is indeed guaranteed. First, we will introduce a more general formulation of the Monge problem, the Kantorovich formulation of OMT.

Therefore we define the space of all transport plans $\Pi(\mu, \nu)$, i.e., of all measures on the product space $\Omega \times \Omega$, such that the first marginal is $\mu$ and the second marginal is $\nu$. Then the OMT problem in the Kantorovich formulation reads as:

**Definition 2.2 ($L^2$-Kantorovich's Optimal Transport Problem)** *Given $\mu, \nu \in \mathcal{P}(\Omega)$, minimize*

$$\mathbb{K}(\pi) = \int_{\Omega \times \Omega} \|x - y\|^2 d\pi(x, y)$$

*over all $\pi \in \Pi(\mu, \nu)$.*

Note, that the $L^2$-Wasserstein metric $W_2(\mu, \nu)$ between $\mu$ and $\nu$ is defined as the infimum of $\mathbb{K}$.

We will now formulate Brenier's theorem, which guarantees the existence of an optimal transport map and will be the central building block of our algorithm:

**Theorem 2.3 (Brenier's Theorem [7])** *Let $\mu, \nu \in \mathcal{P}(\Omega)$. Then there exists an unique optimal transport plan $\pi^* \in \Pi(f, g)$, which is given by $\pi^*(x, y) = (\mathrm{id} \times T)$ where $T = \nabla u$ is the gradient of a convex function $u$ that pushes $\mu$ forward to $\nu$, i.e. $(\nabla u)_*\mu = \nu$. The inverse $T^{-1}$ is also given by the gradient of a convex function that is the Legendre transform of the convex function $u$.*

Thus Brenier's Theorem guarantees the existence and the uniqueness of the optimal transport map of the OMT problem. Consequently one can determine this optimal transport map by solving for the function $u$ in the form of a Monge-Ampère equation:

$$\det(D^2(u)(x)) \cdot g(\nabla u(x)) = f(x) \tag{3}$$

where $D^2$ is the Hessian, $\mu = f dx$ and $\nu = g dx$. Note, that we obtain equation (3) directly from equation (2) using the constraint that $T = \nabla u$ as required by Brenier's theorem. We will also refer to this map as the Brenier map. Note that this map is a diffeomorphism as it's a gradient of a strictly convex function.

Using methods of classical numerical analysis this has been done in [27]. In the following we will propose a new discretization to this problem, which will make use of recent advances in deep learning.
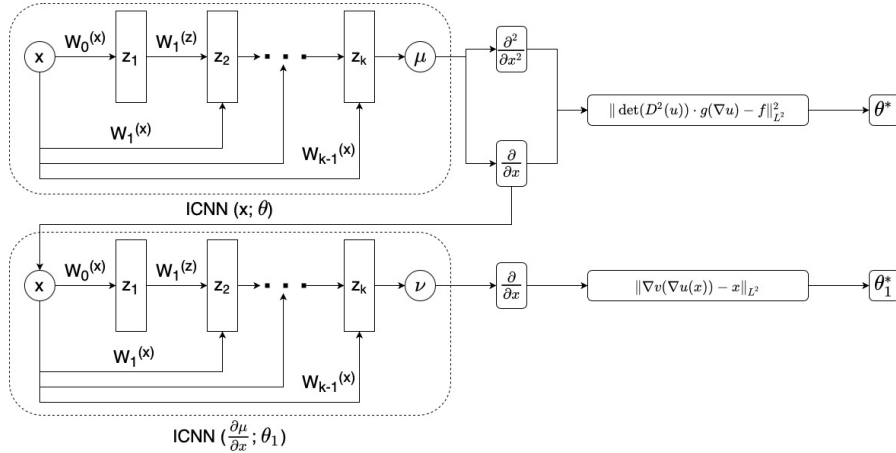
Fig. 1: PICANN Architecture. We present how a combination of two ICNN networks can be used to learn the forward and the inverse map between two distributions. Both these networks are trained independently with their respective loss functions. The inverse network uses the gradient of the output of the first network as it's input.

## 2.2 Solving OMT using PINNs

Physics Informed Neural Networks (PINNs) were proposed by Raissi et al [28] to solve general nonlinear partial differential equations (PDEs). The basic concept is to use the universal approximation property of Deep Neural Networks to represent the solution of a PDE via a network. Using the automatic differentiation capability of modern machine learning frameworks a loss function is formulated, such that it's minimizer solves the PDE in a weak sense. Such a loss function encodes the structured information which results in the amplification of the information content of the data the network sees [28]. This results in good generalization even when only few training examples are available.

PINNs have found widespread applications in a short period of time since their introduction. This includes a wide variety of PDEs, including the Navier-Stokes equation [15], nonlinear stochastic PDEs [36] or Allen Cahn PDEs [22].

In this work we propose to use the PINN approach to solve the Monge-Ampere equation (3) and hence implicitly the Optimal Mass Transport Problem. These equations have been extensively studied and the properties of its solutions are well established. By Theorem 2.3 we know that the solution is given by a convex function $u$.

Recently [1] proposed a new architecture of neural networks, Input Convex Neural Networks (ICNNs), that explicitly constraints the function approximated by the network to be convex. Consequently this architecture naturally lends itself to our proposed application, as it directly encodes Brenier's theorem.

In the ICNN architecture the activation function is a non-decreasing convex function and the internal weights $(W_n^{(x)})$ are constrained to be non-negative, see Figure 1 for a schematic description of this class of networks. This architecture is derived from

two simple facts: non-negative sums of convex functions are also convex and the composition of a convex and convex non-decreasing function is again convex.

In the following we assume that we are given $\mu = f dx$ and $\nu = g dx$. The loss function corresponding to equation (3) is then given by

$$\| \det(D^2(u)) \cdot g(\nabla u) - f \|_{L^2}^2 \tag{4}$$

where $u$ is expressed as the output of a ICNN of sufficient depth and width. Once we have estimated the optimal transport map, the $L^2$-Wasserstein metric between $\mu$ and $\nu$ is given by

$$\int \| x - \nabla u(x) \|^2 \, g(x) dx. \tag{5}$$

We call this combination of the PINN approach with the ICNN structure, Physics Informed Convex Artificial Neural Networks (PICANNs).

In several applications one is interested in computing the inverse transformation at the same time. By a duality argument we know that this map is also given by the gradient of a convex function. Thus we use a second ICNN to compute the inverse optimal transport map ($\nabla v$) by solving the minimisation problem:

$$\| \nabla v(\nabla u(x)) - x \|_{L^2}, \tag{6}$$

where $\nabla u$ is the optimal transport map solving $(\nabla u)_* \mu = \nu$.

### 2.3   Diffeomorphic Random Sampling

In many applications, such as Bayesian estimation one can evaluate the density rather easily but generating samples from a given density is not trivial. Traditional methods include Markov Chain Monte Carlo methods e.g. the Metropolis Hastings algorithm [13]. An alternative idea is to use diffeomorphic density matching between the given density $\nu$ and a standard density $\mu$ from which samples can be drawn easily. Once we have calculated the transport map, standard samples are transformed by the push forward diffeomorphism to generate samples from the target density $\nu$. This approach has been followed in several articles, where the optimal transport map selection was based on both the Fisher-Rao metric [3] or the Knothe–Rosenblatt rearrangement [21]. The efficient implementation of the present paper directly leads to an efficient random sampling algorithm in high dimensions.

### 2.4   Diffeomorphic Density Estimation

We now formulate the density estimation problem using the OMT framework. We are given samples $x_i$ drawn from an unknown density $\mu \in \mathcal{P}(\Omega)$ which we aim to estimate. The main idea of our algorithm is to represent the unknown density as the pull back via a (diffeomorphic) Brenier map $\nabla u$ of a given background density $\nu = g dx$, i.e., $(\nabla u)^* \nu = \mu$ or equivalently $(\nabla u)_* \mu = \nu$.

As we do not have an explicit target density, but only a finite number of samples we need to find a replacement for the $L^2$-norm used in equation (5) to estimate the transport

map $\nabla u$. We do this by maximizing the log-likelihood of the data with respect to the density $(\nabla u)^* \nu$:

$$\frac{1}{N} \sum_i \log \left( \det(D^2(u(x_i))) \cdot g(\nabla u(x_i)) \right). \tag{7}$$

Using our PINNs framework we represent the convex function $u$ again via an ICNN, which serves as an implicit regularizer. Instead of maximizing the above log-likelihood we solve the minimization problem

$$-\frac{1}{N} \sum_i \log \left( \det(D^2(u(x_i))) \cdot g(\nabla u(x_i)) \right) \tag{8}$$

Note that this can be alternatively interpreted as minimizing the empirical Kullback-Leibler divergence between $\mu$ and the pull back of the background density $\nu$. To see this we recall that the Kullback-Leibler divergence between $\tilde{\mu} = f dx$ and $\tilde{\nu} = \tilde{g} dx$ is given by

$$\mathrm{KL}(f \| \tilde{g}) = \int f \log \left( \frac{f}{\tilde{g}} \right) dx \tag{9}$$

Applying this formula to

$$\tilde{\nu} = \det(\mathrm{D}^2(\mathrm{u})) \cdot g(\nabla u) dx$$

leads to

$$\mathrm{KL}(f \| \det(\mathrm{D}^2(\mathrm{u})) \cdot g(\nabla u))$$
$$= \int f \log \left( \frac{f}{\det(\mathrm{D}^2(\mathrm{u})) \cdot g(\nabla u)} \right) dx$$
$$= \int f \log f dx - \int f \log \left( \det(\mathrm{D}^2(\mathrm{u})) \cdot g(\nabla u) \right) dx.$$

The first term in the above equation is a constant and can be ignored for the minimization purposes. Thus we obtain

$$\int f \log \left( \det(D^2(u)) \cdot g(\nabla u) \right)$$
$$= \mathbb{E}_f \left( \log \left( \det(D^2(u)) \cdot g(\nabla u) \right) \right)$$
$$\approx \frac{1}{N} \sum_i \log \left( \det(D^2(u(x_i))) \cdot g(\nabla u(x_i)) \right)$$

where $x_i$ are the given samples from the unknown density $\mu = f dx$.

To generate new samples from the estimated density we use the inverse map to transform the samples from the background density $\nu$. Note, that we calculate the inverse map using a second neural network and explicit loss function given by equation (6).

(a) Ground truth          (b) Est. density (PICANN)          (c) Est. density (OT ICNN)



(d) Ground truth          (e) Est. map (PICANN)          (f) Est. map (OT ICNN)
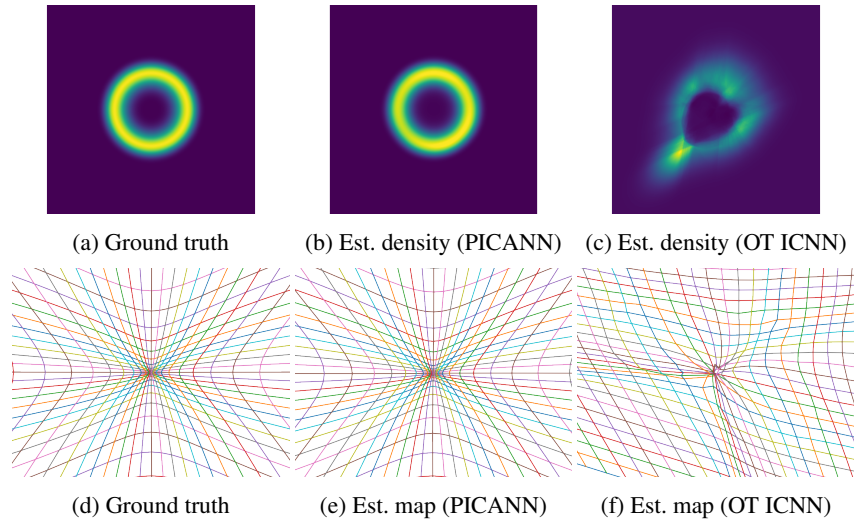
Fig. 2: Validation: Panel (a) shows the true annulus distribution, the estimated annulus distribution using the PICANN approach and OT ICNN approach are shown in Panels (b) and (c) respectively. Panel (d) shows the analytical optimal transport map between the unit Gaussian and the annulus distribution. The estimated optimal transport map using the PICANN approach is presented in Panel (e) while the map estimated using the OT ICNN method is shown in Panel (f).

## 3  Experimental Results

In this section we will detail our implementation and present selected examples for both random sampling and density estimation.

### 3.1  Network details

As explained in Section 2.2 we use an ICNN architecture for both the forward and the backwards map in all of our experiments, c.f. Figure 1. As with every deep learning approach we need to tune the hyper-parameters including width/depth of the network, activation functions and the batch size. The width of the network needs to increase with the dimension of the ambient space of the data to ensure sufficient flexibility. For our experiments in 2d we used a network with 5 hidden layers with 128 neurons in each layer, whereas for experiments in higher dimensions we used 256 neurons in each layer. To initialize the network, we first train the networks to learn the identity transformation, i.e. $\nabla u = I$, which we use as the initial starting point for all of our experiments.

To guarantee convexity of the output function, the activation functions need to be convex and non-decreasing. The simple ReLU's are not strictly convex and the second derivative is zero almost everywhere. Thus we experimented with the family of Rectified Power Units (RePUs) and the log exponential family or the 'Softplus' function.

The Softplus function to the power of $\alpha$, which defined via

$$\text{Softplus}^{\alpha}(x) = (\log{(1 + \exp{x})})^{\alpha}$$

turned out to be best suited for our applications, where we choose $\alpha = 1.1$; in particular our experiments suggested that networks with this activation function were able to generalize well to regions where no or only limited training data was available.

| Dimensions | Number of Samples | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | 5k | | 10k | | 20k | |
| | Avg % error (Std dev) | | Avg % error (Std dev) | | Avg % error (Std dev) | |
| | PICANN | OT ICNN | PICANN | OT ICNN | PICANN | OT ICNN |
| 2 | **2.34** (2.16) | 15.46 (10.01) | **2.82** (1.89) | 18.88 (6.59) | **2.04** (1.53) | 7.10 (7.75) |
| 3 | **1.36** (1.26) | 27.75 (7.07) | **1.95** (1.38) | 14.61 (5.53) | **1.69** (1.60) | 3.68 (2.77) |
| 5 | **0.69** (0.44) | 19.37 (3.80) | **0.88** (0.61) | 6.31 (1.48) | **0.54** (0.39) | 0.87(0.59) |
| 8 | **0.28** (0.19) | 15.74 (1.77) | **0.33** (0.31) | 4.03 (0.70) | 0.39 (0.27) | **0.32** (0.19) |
| 15 | **0.51** (0.68) | 6.84 (0.53) | **0.19** (0.14) | 1.54 (0.22) | **0.16** (0.10) | 0.24 (0.18) |
| 30 | **0.16** (0.15) | 2.32 (2.42) | **0.16** (0.11) | 0.67 (0.13) | **0.19** (0.12) | **0.19** (0.12) |

Table 1: Validation: We present a comparison between our PICANN approach and the ICNN min-max approach of [20]. In all experiments we compare the average errors over 20 realizations between the true Wasserstein metric and the approximated Wasserstein metric using these algorithms between a unit Gaussian and a randomly sampled Gaussian. The best results in each experiment are boldfaced.

## 3.2  Validation and accuracy of PICANNs in computing optimal transport maps

We validate the accuracy of the optimal transport maps as computed using our PICANN approach. Towards this aim we consider special situations, where one can analytically calculate the OMT distance. We then compare this ground truth to the approximations as calculated using our approach. To further demonstrate the quality of our results we compare them to results obtained with our main competitor, the previous deep learning ICNNs approach of [20]. We do not present comparisons of our approach to the Sinkhorn algorithm [10] or the linear programming approaches [27], as these frameworks, although they approximate the OMT distances, do not compute the continuous optimal transport map, which is essential for the proposed density estimation and generative modeling. While finite element or finite difference based Monge-Ampere solvers, see e.g. [5,14,4], calculate the continuous OMT map, they are not suitable in dimensions greater than two or three.

Analytic solutions for the OMT-problem are in general not available. For the special case when both densities are from a family of Gaussians distributions the OMT map is

simply given by an affine transform and the OMT distance is given in closed form:

$$\mathrm{w}_2(\mathcal{N}(m_1, \Sigma_1); \mathcal{N}(m_2, \Sigma_2))^2 =$$

$$\|m_1 - m_2\|^2 + \mathrm{Tr}(\Sigma_1 + \Sigma_2 - 2(\Sigma_1^{\frac{1}{2}} \Sigma_2 \Sigma_1^{\frac{1}{2}})^{\frac{1}{2}}) \quad (10)$$

here $\mathcal{N}(m_1, \Sigma_1)$ and $\mathcal{N}(m_2, \Sigma_2)$ are two Gaussians with mean $m_i$ and covariance $\Sigma_i$.

To statistically quantify the accuracy we generate Gaussian distributions with random mean and covariances, where the the means are sampled from a uniform distribution on $[-1, 1]$. To construct a random covariance matrix, we recall that we need to enforce the matrix to be positive definite and symmetric. Therefore we generate a random matrix $A$ of dimension $d \times 3d$, where $d$ is the dimension of the space and where the entries are i.i.d. chosen from a uniform distribution on $[0, 0.75]$. Then a random covariance matrix can be constructed by letting $\Sigma = AA^T$, does almost surely guaranteeing positive definiteness and symmetricity.

In table 1 we reports the mean errors and standard deviations in different dimensions for both using our method (PICANN) and the method of [20] (OT ICNN). In all these experiments we fix $\mathcal{N}(m_1, \Sigma_1)$ to the unit Gaussian i.e. $\mathcal{N}(0, I)$ and we consider the second Gaussian using the randomly sampled mean $m_2$ and covariance matrix $\Sigma_2$. We also present the scalability of our algorithm by performing the same experiment in dimensions 5, 8, 15 and 30. We also show the effect of using different sizes of data samples during the training process. Each experiment was repeated 20 times to compute the statistics. In all dimensions one can clearly see the superior performance of our method for 5k and 10k data samples; for 20k samples both methods perform comparably.

To further validate our algorithm we choose a more challenging problem: an annulus density for which we know the transport map in closed form. The annulus distribution is given by a push forward of the Gaussian distribution by a gradient of a radially symmetric convex function. This distribution is given by:

$$f = g((x^2 + y^2) \cdot (x, y)) \cdot 3(x^2 + y^2)^2 \quad (11)$$

where $g$ is the unit Gaussian. Figure 2 shows the true annulus distribution, the distribution predicted using our PICANN approach and the error in those two. It should be noted that the transport map $X \mapsto (X^T X)X$ is the gradient of the convex function $\frac{1}{4}(X^T X)^2$. Thus, this is not only the transport map but the optimal transport map for the Gaussian distribution to the annulus distribution. The inverse map of the optimal transport map is given by $X \mapsto X(X^T X)^{-\frac{1}{3}}$. As we have the transport map in closed form, also shown in the same figure, we can compute the Wasserstein metric between the annulus distribution and the unit Gaussian using this map. Comparing this theoretical value with the Wasserstein metric computed using our approach we have again a quantitative measure about the accuracy of our algorithm. In our experiments we found that the theoretical Wasserstein metric between the annulus distribution and the unit Gaussian is $1.0097$ while the approximated metric using our method is $1.0174$. This is a difference of $0.76\%$ between the two. Using the OT ICNN code we were not able to get a satisfactory approximation of the transport map and even with 20k samples we obtained an error of $58.16\%$ in the estimated Wasserstein distance. See also Fig. 2. As this example in 2d, we were also able to compare it to classical methods such as [14], but also this did not give any satisfactory results.

### 3.3    Random Sampling and Generative Autoencoders



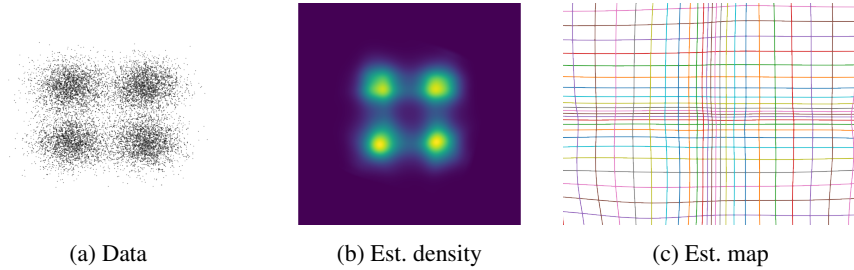(a) Data                    (b) Est. density                    (c) Est. map

Fig. 3: Density Estimation 1: In this figure we show an example for density estimation using a simple Gaussian mixture. Panel (a) shows the given data; the approximated density and the inverse map as found using our PICANN approach are shown in Panels (b) and (c).



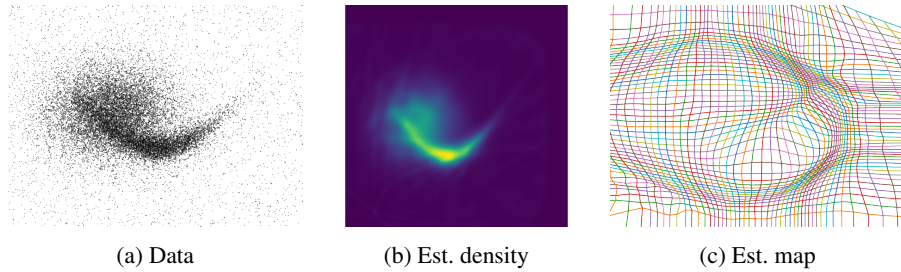(a) Data                    (b) Est. density                    (c) Est. map

Fig. 4: Density Estimation 2: In this figure we show a second example for density estimation. Panel (a) shows the given data; the approximated density and the inverse map as found using our PICANN approach are shown in Panels (b) and (c).

In this section we present the results for using the PICANN framework for density estimation. We consider the problem of estimating a continuous density from discrete finite samples as described in Section 2.4. Shown in Figure 3 are 10k random samples generated from a known Gaussian mixture model of 4 Gaussians. We use the standard normal distribution as the reference distribution and estimate the Optimal Transport Map between the data and the reference using our PICANN approach. The push forward of the reference distribution by the estimated transport map and the estimated transport map are both shown in Figure3. One can see that the estimated density matches the original Gaussian mixture.

Next we consider a more challenging example, see Figure 4 which shows 20k random samples from an non-symmetric distribution, that has been constructed in [3]. We

again present the estimated density and transport map. Similar as in the first example we obtain a good match with the original distribution, where one needs a highly non-linear transport map.

### 3.4  Generative Modeling with PICANNs

Finally we present an application of the PICANN approach to develop a generative model. For this task we operate within the autoencoder setting, which is essentially a dimensionality reduction technique [17]. The autoencoder maps a high dimensional data to a latent space of lower dimensions, which then can be transformed back to the original space using the "decoder" part of the network. Fig 5 summarizes how the autoencoder can be naturally included in our PICANN approach to obtain an efficient generative model for high-dimensional data.
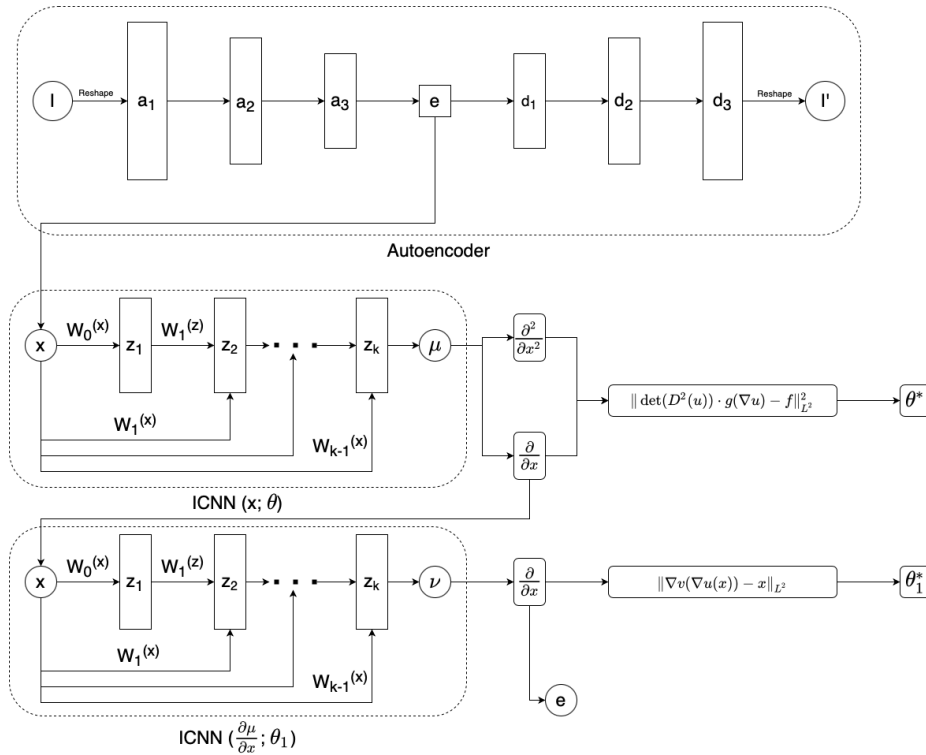


Fig. 5: This figure shows how a combination of an autoencoder with our PICANN approach can be used to develop a generative model. Note how the latent space of the autoencoder becomes the input to the PICANN network. In such a setting the PICANN estimates the latent space density and samples from the estimated distribution. Using random samples from this distribution one can pass them through the decoder to generate new samples.

To demonstrate the effectiveness of the generative algorithm we consider the MNIST dataset encoded using a simple fully connected autoencoder to a latent space of dimension 2. Figure 6 shows the encoded data points with each class assigned a unique color. We train a PICANN network of 5 hidden layers with 128 neurons in each layer, to learn the forward and the inverse mapping from a Gaussian with mean and covariance of the encoded data to the unknown "latent MNIST distribution". Figure 6 panel(c) presents the samples as generated from the Gaussian and then transported points using the approximated inverse map. Panel (d) displays the result of passing a random subset of these generated samples through the decoder.



(a) Encoded MNIST
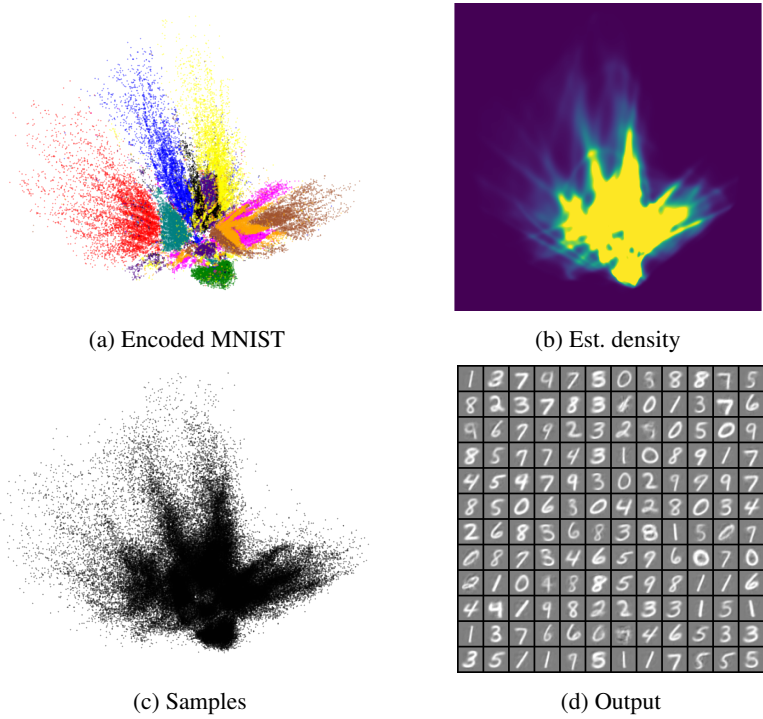
(b) Est. density

(c) Samples

(d) Output

Fig. 6: Generative model: this figure details the application of our framework to the MNIST dataset. Panel (a) shows the encoded samples from the MNIST database. A PI-CANN was trained to estimate this distribution and learn the forward and inverse transport map between the encoded 'MNIST Distribution' and a Gaussian. The estimated density can be seen in Panel (b) and Panel (c) shows 100k samples generated from this density using our obtained optimal transport map. The first 144 samples passed through the decoder part of the autoencoder are shown in Panel (d).

## 4    Conclusion

In this paper we use the $L^2$-Wasserstein metric and optimal mass transport (OMT) theory to formulate a density estimation estimation and generative modeling framework. Towards this aim we develop a deep learning based solver for the continuous OMT problem, which is rooted in Brenier's celebrated theorem. This theorem allows us to formulate the density estimation problem as a solution to a nonlinear PDE – a Monge-Ampere equation. Recent developments in deep learning for PDEs, namely PINNS and ICNNs, allow us to develop an efficient solver. We demonstrate the accuracy of our framework, by comparing our results to analytic Wasserstein distances. Finally, we present an example of how the diffeomorphic density estimation framework can be used as a generative model.

## References

1. Amos, B., Xu, L., Kolter, J.Z.: Input convex neural networks. In: International Conference on Machine Learning. pp. 146–155. PMLR (2017)
2. Angenent, S., Haker, S., Tannenbaum, A.: Minimizing flows for the monge–kantorovich problem. SIAM journal on mathematical analysis **35**(1), 61–97 (2003)
3. Bauer, M., Joshi, S., Modin, K.: Diffeomorphic random sampling using optimal information transport. In: International Conference on Geometric Science of Information. pp. 135–142. Springer (2017)
4. Benamou, J.D., Brenier, Y.: A computational fluid mechanics solution to the monge-kantorovich mass transfer problem. Numerische Mathematik **84**(3), 375–393 (2000)
5. Benamou, J.D., Duval, V.: Minimal convex extensions and finite difference discretisation of the quadratic monge–kantorovich problem. European Journal of Applied Mathematics **30**(6), 1041–1078 (2019)
6. Benamou, J.D., Froese, B.D., Oberman, A.M.: Two numerical methods for the elliptic monge-ampere equation. ESAIM: Mathematical Modelling and Numerical Analysis **44**(4), 737–758 (2010)
7. Brenier, Y.: Polar factorization and monotone rearrangement of vector-valued functions. Communications on pure and applied mathematics **44**(4), 375–417 (1991)
8. Carlier, G., Galichon, A., Santambrogio, F.: From knothe's transport to brenier's map and a continuation method for optimal transport. SIAM Journal on Mathematical Analysis **41**(6), 2554–2576 (2010)
9. Cuturi, M.: Sinkhorn distances: lightspeed computation of optimal transport. In: NIPS. vol. 2, p. 4 (2013)
10. Cuturi, M.: Sinkhorn distances: Lightspeed computation of optimal transport. Advances in neural information processing systems **26**, 2292–2300 (2013)
11. Dinh, L., Sohl-Dickstein, J., Bengio, S.: Density estimation using real NVP. In: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net (2017), `https://openreview.net/forum?id=HkpbnH9lx`
12. Fix, E., Hodges, J.L.: Nonparametric discrimination: Consistency properties. Randolph Field, Texas, Project pp. 21–49 (1951)
13. Hastings, W.K.: Monte carlo sampling methods using markov chains and their applications (1970)
14. Jacobs, M., Léger, F.: A fast approach to optimal transport: The back-and-forth method. Numerische Mathematik **146**(3), 513–544 (2020)

15. Jin, X., Cai, S., Li, H., Karniadakis, G.E.: Nsfnets (navier-stokes flow nets): Physics-informed neural networks for the incompressible navier-stokes equations. Journal of Computational Physics **426**, 109951 (2021)
16. Knothe, H.: Contributions to the theory of convex bodies. Michigan Mathematical Journal **4**(1), 39–52 (1957)
17. Kramer, M.A.: Nonlinear principal component analysis using autoassociative neural networks. AIChE journal **37**(2), 233–243 (1991)
18. Loader, C.R., et al.: Bandwidth selection: classical or plug-in? The Annals of Statistics **27**(2), 415–438 (1999)
19. Loeper, G., Rapetti, F.: Numerical solution of the monge–ampère equation by a newton's algorithm. Comptes Rendus Mathematique **340**(4), 319–324 (2005)
20. Makkuva, A., Taghvaei, A., Oh, S., Lee, J.: Optimal transport mapping via input convex neural networks. In: International Conference on Machine Learning. pp. 6672–6681. PMLR (2020)
21. Marzouk, Y., Moselhy, T., Parno, M., Spantini, A.: Sampling via measure transport: An introduction. Handbook of uncertainty quantification pp. 1–41 (2016)
22. McClenny, L., Braga-Neto, U.: Self-adaptive physics-informed neural networks using a soft attention mechanism. arXiv preprint arXiv:2009.04544 (2020)
23. McLachlan, G.J., Basford, K.E.: Mixture models: Inference and applications to clustering, vol. 38. M. Dekker New York (1988)
24. Monge, G.: Mémoire sur la théorie des déblais et des remblais. Histoire de l'Académie Royale des Sciences de Paris (1781)
25. Papadakis, N.: Optimal transport for image processing. Ph.D. thesis, Université de Bordeaux; Habilitation thesis (2015)
26. Parzen, E.: On estimation of a probability density function and mode. The annals of mathematical statistics **33**(3), 1065–1076 (1962)
27. Peyré, G., Cuturi, M., et al.: Computational optimal transport: With applications to data science. Foundations and Trends® in Machine Learning **11**(5-6), 355–607 (2019)
28. Raissi, M., Perdikaris, P., Karniadakis, G.E.: Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations. arXiv preprint arXiv:1711.10561 (2017)
29. Rosenblatt, M.: Remarks on a multivariate transformation. The annals of mathematical statistics **23**(3), 470–472 (1952)
30. Rosenblatt, M.: Remarks on a multivariate transformation. The annals of mathematical statistics **23**(3), 470–472 (1952)
31. Shin, Y., Darbon, J., Karniadakis, G.E.: On the convergence and generalization of physics informed neural networks. arXiv preprint arXiv:2004.01806 (2020)
32. Silverman, B.W., Jones, M.C.: E. fix and jl hodges (1951): An important contribution to nonparametric discriminant analysis and density estimation: Commentary on fix and hodges (1951). International Statistical Review/Revue Internationale de Statistique pp. 233–238 (1989)
33. Villani, C.: Topics in optimal transportation. No. 58, American Mathematical Soc. (2003)
34. Villani, C.: Optimal transport: old and new, vol. 338. Springer Science & Business Media (2008)
35. Younes, L.: Diffeomorphic learning. Journal of Machine Learning Research **21**(220), 1–28 (2020), http://jmlr.org/papers/v21/18-415.html
36. Zhang, D., Guo, L., Karniadakis, G.E.: Learning in modal space: Solving time-dependent stochastic pdes using physics-informed neural networks. SIAM Journal on Scientific Computing **42**(2), A639–A665 (2020)