

AUGMENTED KRNET FOR DENSITY ESTIMATION AND APPROXIMATION

XIAOLIANG WAN* AND KEJUN TANG†

Abstract. In this work, we have proposed augmented KRnets including both discrete and continuous models. One difficulty in flow-based generative modeling is to maintain the invertibility of the transport map, which is often a trade-off between effectiveness and robustness. The exact invertibility has been achieved in the real NVP using a specific pattern to exchange information between two separated groups of dimensions. KRnet has been developed to enhance the information exchange among data dimensions by incorporating the Knothe-Rosenblatt rearrangement into the structure of the transport map. Due to the maintenance of the exact invertibility, a full nonlinear update of all data dimensions needs three iterations in KRnet. To alleviate this issue, we will add augmented dimensions that act as a channel for the data dimensions to exchange information. In the augmented KRnet, a fully nonlinear update is achieved in two iterations. We also show that the augmented KRnet can be reformulated as the discretization of a neural ODE, where the exact invertibility is kept such that the adjoint method can be formulated with respect to the discretized ODE to obtain the exact gradient. Numerical experiments have been implemented to demonstrate the effectiveness of our models.

Key words. Deep learning, Density estimation, Optimal transport, Uncertainty quantification

1. Introduction. Density estimation and approximation play an important role in many fields such as variational Bayes, uncertainty quantification, unsupervised learning, etc. However, classical approaches or models such as kernel density estimator or the mixture of Gaussians are usually limited to low-dimensional cases due to the curse of dimensionality [25]. Recently deep generative modeling has received a lot of attention in deep learning, which is closely related to density estimation. The main motivation of deep generative modeling is to deal with the distribution of data that have a very large number of dimensions, e.g., high-resolution images. So the deep generative modeling needs to balance modeling capability and efficiency. For example, generative adversarial networks (GANs) [11, 1] are able to learn a mapping from a latent space to the data space without an explicit definition of likelihood. Due to such a flexibility, GANs have been successfully applied to many realistic applications; however, the lack of likelihood means that it is not suitable for density approximation, where a probability density function (PDF) is needed. Likelihood-based deep generative models have also been developed including the autoregressive models [12, 21, 22, 23], variational autoencoders (VAE) [16, 19], and flow-based generative models [6, 24, 7, 17, 31, 3]. The combination of different modeling strategies has also been actively explored. For instance, the flow-based model was coupled with GAN in [13] to obtain a likelihood; The VAE, flow-based model and GAN were coupled in [32] for more flexibility and efficiency; The flow-based model has been formulated as a discretized neural ordinary differential equations (ODE) [5, 8], where the velocity field of the ODE is modeled as a neural network.

We pay particular attention to flow-based generative models. The underlying idea of flow-based generative models is to construct a transport map from the data distribution to a prior distribution, e.g., the standard Gaussian. There are two ways to define such a transport map: continuous and discrete models. The continuous models refer to the dynamical evolution given by a neural ODE, which transforms

*Department of Mathematics, Center for Computation and Technology, Louisiana State University, Baton Rouge 70803 (xlwan@math.lsu.edu)

†Peng Cheng Laboratory, Shenzhen, China (tangkj@pcl.ac.cn).

the distribution of the initial data to another distribution within a certain amount of time. In discrete models, the transport maps are explicitly constructed by stacking a sequence of simple bijections modeled by shallow neural networks. Both continuous and discrete models need to maintain the invertibility of the transport map. The transport map and its inverse determine two important things. One mapping direction yields the PDF model of the data distribution, which can be written as a product of the PDF of the prior distribution and the determinant of the Jacobian matrix; and the other mapping direction yields sample generation, which maps the samples generated by the prior distribution to samples that are consistent with the data distribution. Simply speaking, flow-based generative models provide an PDF model, which can be easily sampled. This is similar to classical probabilistic model such as the mixture of Gaussians. However, deep generative models are usually much more complex and capable.

One interesting question is whether the flow-based generative model can serve as a generic PDF model for both density approximation and sample generation for problems in scientific computing. Note that density approximation and sample generation are usually addressed separately. To approximate a high-dimensional PDF, such as the posterior distribution in variational Bayes, the commonly used model is a Gaussian with diagonal or banded covariance matrix, which is often too simple although the statistics can be easily dealt with [2]. Given unlimited computational resources, sampling approaches such as Markov Chain Monte Carlo (MCMC) may eventually yield true samples for an arbitrary PDF under the assumption that the PDF is explicitly known up to a constant. It is challenging to compute the statistics of a high-dimensional random variable whose density is defined by a PDF equation, where both PDF approximation and sample generation may be expected simultaneously. We expect that the flow-based generative model can be capable enough to balance these two issues, e.g., we have applied the real NVP [7] to importance sampling for efficient probability estimation for a PDE subject to uncertainty [29], and KRnet [27] to approximate high-dimensional Fokker-Planck equations [28].

We have developed KRnet in [27] as a generalization of the real NVP [7] by incorporating the triangular structure of the Knothe-Rosenblatt rearrangement into the definition of the transport map. The real NVP separates all data dimensions into two groups. When updating the current data, one group of dimensions can receive nonlinear information of the other group but only linear information of itself, which is a compromise to maintain the exact invertibility of the transport map. Using such a method, a fully nonlinear update needs three iterations. The main idea of KRnet was to enhance the exchange of information among data dimensions through a more flexible partition of data dimensions, which, however, does not change the mechanism of information exchange in each iteration. So KRnet cannot deal with one-dimensional data since two groups of data dimensions are needed. To alleviate this issue, we introduce augmented dimensions in this paper, which serve as a channel for the data dimensions to send and receive nonlinear information. The augmented KRnet achieves a fully nonlinear update in two iterations and is able to deal with one-dimensional data. We then reformulate the augmented KRnet such that it can be regarded as a discretization of an ODE, where the exact invertibility is kept in the discretization. The advantage of a discretization with exact invertibility is that the adjoint method can be formulated in terms of the discrete model instead of the ODE such that the computation of the gradient is exact. The drawback is that the accuracy of such a discretization is only of first order.

The manuscript is organized as follows. In next section we briefly overview flow-based generative models and the KRnet. In section 3, we define augmented KRnet including both discrete and continuous models. Numerical experiments are implemented to demonstrate the effectiveness of the proposed strategies in section 4, followed by a summary section.

2. KRnet. KRnet is a discrete flow-based generative model. Generally speaking, the key component of a flow-based generative model is an invertible mapping $f(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^n$:

$$\begin{aligned} \mathbf{z} &= f(\mathbf{y}) = f_{[m]} \circ f_{[m-1]} \circ \dots \circ f_{[i]} \circ \dots \circ f_{[2]} \circ f_{[1]}(\mathbf{y}), \\ \mathbf{y} &= f^{-1}(\mathbf{z}) = f_{[1]}^{-1} \circ f_{[2]}^{-1} \circ \dots \circ f_{[i]}^{-1} \circ \dots \circ f_{[m-1]}^{-1} \circ f_{[m]}^{-1}(\mathbf{z}), \end{aligned}$$

which can be regarded as a composite mapping that consists of a sequence of intermediate bijections $f_{[i]}(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^n$. Let $p_{\mathbf{Z}}$ and $p_{\mathbf{Y}}$ be the probability density functions (PDF) of the random variables \mathbf{Z} and \mathbf{Y} respectively. The transformation $\mathbf{Z} = f(\mathbf{Y})$ yields the following relation

$$p_{\mathbf{Y}}(\mathbf{y}) = p_{\mathbf{Z}}(f(\mathbf{y})) |\det \nabla_{\mathbf{y}} f(\mathbf{y})|. \quad (2.1)$$

In other words, if we associate \mathbf{Z} with a certain prior distribution, e.g., the standard Gaussian distribution, the mapping $f(\cdot)$ induces explicitly a PDF model $p_{\mathbf{Y}}$, which can be used for either density estimation or approximation. Furthermore, the distribution $p_{\mathbf{Y}}$ can be easily sampled as $\mathbf{y}^{(i)} = f^{-1}(\mathbf{z}^{(i)})$ thanks to the invertibility of $f(\cdot)$, where $\mathbf{z}^{(i)}$ is sampled from the prior distribution $p_{\mathbf{Z}}$. Flow-based generative models share two main features: a large number of intermediate mappings $f_{[i]}(\cdot)$ and invertibility. Since $f(\cdot)$ intends to map the prior to an arbitrary distribution, a large number of intermediate mappings implies that the complexity of $f_{[i]}(\cdot)$ can be reduced. There are two ways to obtain $f_{[i]}(\cdot)$. One is explicit construction, e.g., NICE [6], real NVP [7], planar flow [24], inverse autoregressive flow [19], Sylvester flow [3], and KRnet [27]; the other one is through the discretization of a continuous model, e.g., neural ODE [5] and its variants subject to either augmentation [8] or regularization [30, 10]. Depending on the way to obtain $f_{[i]}(\cdot)$, we may classify the flow-based generative models as discrete or continuous models. Invertibility plays an important role because density estimation and sample generation use opposite directions of the same mapping. The invertibility of discrete models is usually exact and maintained locally by each $f_{[i]}(\cdot)$ while the invertibility of continuous models may only be kept at the continuous level and is not exact locally after the continuous model is discretized. For instance, for neural ODEs the two directions of the mapping $f(\cdot)$ correspond to forward and backward integration of the ODE. It is well known that although an ODE is theoretically invertible there does not exist a numerical scheme which is exactly invertible.

In terms of the optimal transport theory, the mapping $f(\cdot)$ corresponds to a transport map. Let $\mu_{\mathbf{Y}}$ and $\mu_{\mathbf{Z}}$ indicate the probability measures of \mathbf{Y} and \mathbf{Z} , respectively. The mapping $T : \mathbf{Z} \rightarrow \mathbf{Y}$ is called a transport map such that $T_{\#}\mu_{\mathbf{Z}} = \mu_{\mathbf{Y}}$, where $T_{\#}\mu_{\mathbf{Z}}$ is the push-forward of the law $\mu_{\mathbf{Z}}$ of \mathbf{Z} such that $\mu_{\mathbf{Y}}(B) = \mu_{\mathbf{Z}}(T^{-1}(B))$ for every Borel set B [9]. It is seen that T can be defined as $T(\mathbf{z}) = f^{-1}(\mathbf{z})$. The Knothe-Rosenblatt (K-R) rearrangement says that a transport map may have a lower-

triangular structure such that [26]

$$\mathbf{z} = T^{-1}(\mathbf{y}) = f(\mathbf{y}) = \begin{bmatrix} f_1(y_1) \\ f_2(y_1, y_2) \\ \vdots \\ f_n(y_1, y_2, \dots, y_n) \end{bmatrix}. \quad (2.2)$$

Such a mapping can be regarded as a limit of a sequence of optimal transport maps when the quadratic cost degenerates [4]. We have defined a flow-based generative model called KRnet in [27, 28] which generalizes the real NVP [7] by adapting the triangular structure of the K-R rearrangement into the model. For more flexibility, we consider a partition $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_K)$ in KRnet, where $\mathbf{y}_i = (y_{i,1}, \dots, y_{i,m})$, $1 \leq K \leq n$, $1 \leq m \leq n$, and $\sum_{i=1}^K \dim(\mathbf{y}_i) = n$. We employ a block-version of the K-R rearrangement

$$\mathbf{z} = f(\mathbf{y}) = \begin{bmatrix} f_1(\mathbf{y}_1) \\ f_2(\mathbf{y}_1, \mathbf{y}_2) \\ \vdots \\ f_K(\mathbf{y}_1, \dots, \mathbf{y}_K) \end{bmatrix}. \quad (2.3)$$

To integrate the K-R rearrangement into the data flow of $f(\cdot)$, we need to associate the mappings $f_i(\cdot)$, $i = 1, \dots, K$, with an order. We let the data flow from f_K to f_1 :

$$\mathbf{y} \xrightarrow{f_K} \mathbf{y}_{t_1} \xrightarrow{f_{K-1}} \mathbf{y}_{t_2} \xrightarrow{f_{K-2}} \dots \mathbf{y}_{t_{K-1}} \xrightarrow{f_1} \mathbf{y}_{t_K} = \mathbf{z}.$$

At step t_i , a certain group of dimensions will be deactivated. Thus KRnet has a lower triangular overall structure in the sense that the number of effective dimensions decreases similarly to the transition from $f_K(\cdot)$ to $f_1(\cdot)$ in equation (2.3).

2.1. An overview of the layers in KRnet. We now briefly overview some main building blocks $f_{[i]}(\cdot)$ used in KRnet. More details can be found in [27, 28]. We let $\mathbf{y}_{[i]}$ indicate an intermediate state of data after the mapping $f_{[i-1]}(\cdot)$, i.e., $\mathbf{y}_{[i]} = f_{[i]}(\mathbf{y}_{[i-1]})$ with $\mathbf{y}_{[0]} = \mathbf{y}$.

1. *Squeezing layer* deactivates a certain number of components using a mask

$$\mathbf{q} = (\underbrace{1, \dots, 1}_k, \underbrace{0, \dots, 0}_{n-k}).$$

The first k components given by $\mathbf{q} \odot \mathbf{y}_{[i]}$ will keep being updated and the rest $(n - k)$ components given by $(1 - \mathbf{q}) \odot \mathbf{y}_{[i]}$ will be deactivated from then on. Here \odot indicates the Hadamard product or component-wise product.

2. *Rotation layer* provides a simple and trainable strategy to determine the dimensions that will be deactivated first. The rotation layer defines a rotation of the coordinate system through an orthogonal matrix for the current active dimensions:

$$\mathbf{y}_{[i+1]} = \hat{\mathbf{W}} \mathbf{y}_{[i]} = \begin{bmatrix} \mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \mathbf{y}_{[i]} = \begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{U} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \mathbf{y}_{[i]},$$

where $\mathbf{W} \in \mathbb{R}^{k \times k}$ with k being the number of 1's in \mathbf{q} , and $\mathbf{I} \in \mathbb{R}^{(n-k) \times (n-k)}$ is an identity matrix, and $\mathbf{W} = \mathbf{L}\mathbf{U}$ is the LU decomposition of \mathbf{W} . The entries in the lower triangular part of \mathbf{L} and the upper triangular part of \mathbf{U} will be treated as trainable

parameters of the model except for the diagonal entries of \mathbf{L} which are equal to 1. Intuitively we expect the rotation may put the most important dimensions at the beginning, which need further modifications. We need to clarify one thing. Although the purpose of this layer can be understood through a rotation of the coordinate system, we simply train \mathbf{L} and \mathbf{U} in practice without enforcing the unity of $\hat{\mathbf{W}}$.

3. *Scale and bias layer* provides a simplification of batch normalization which is defined as [15, 17]

$$\mathbf{y}_{[i+1]} = \mathbf{a} \odot \mathbf{y}_{[i]} + \mathbf{b}, \quad (2.4)$$

where \mathbf{a} and \mathbf{b} are trainable, and initialized by the mean and standard deviation of data. After the initialization, \mathbf{a} and \mathbf{b} will be treated as regular trainable parameters that are independent of the data. The scale and bias layer helps to improve the conditioning of deep nets.

4. *Affine coupling layer* is the most important layer for evolving the data. Consider a partition $\mathbf{y}_{[i]} = (\mathbf{y}_{[i],1}, \mathbf{y}_{[i],2})$ with $\mathbf{y}_{[i],1} \in \mathbb{R}^m$ and $\mathbf{y}_{[i],2} \in \mathbb{R}^{n-m}$. The affine coupling layer is defined as [27, 7]

$$\begin{cases} \mathbf{z}_1 = \mathbf{y}_{[i],1}, \\ \mathbf{z}_2 = \mathbf{y}_{[i],2} \odot (1 + \alpha \tanh(\mathbf{s}(\mathbf{y}_{[i],1})) + e^{\beta} \odot \tanh(\mathbf{t}(\mathbf{y}_{[i],1}))), \end{cases} \quad (2.5)$$

where $\mathbf{s}, \mathbf{t} \in \mathbb{R}^{n-m}$ stand for scaling and translation functions depending only on $\mathbf{y}_{[i],1}$, $0 < \alpha < 1$ and $\beta \in \mathbb{R}$. Note that $\mathbf{y}_{[i],2}$ is updated linearly while the mappings $\mathbf{s}(\mathbf{y}_{[i],1})$ and $\mathbf{t}(\mathbf{y}_{[i],1})$ can be arbitrarily complicated, which are modeled as a neural network (NN),

$$(\mathbf{s}, \mathbf{t}) = \text{NN}(\mathbf{y}_{[i],1}). \quad (2.6)$$

Then the Jacobi matrix is lower-triangular, and an inverse can be easily computed. The two parts of $\mathbf{y}_{[i]}$ will be updated alternately by a sequence of affine coupling layers, e.g., at the next affine coupling layer, the first partition will be modified while the second partition remains fixed.

5. *Nonlinear invertible layer* defines a component-wise one-dimensional nonlinear mapping to map \mathbb{R} to itself. We decompose $\mathbb{R} = (-\infty, -a) \cup [-a, a] \cup (a, \infty)$ for $0 < a < \infty$, and define

$$z = \hat{F}(y) = \begin{cases} \beta(y - a) + a, & y \in (-\infty, -a) \\ \phi^{-1} \circ F \circ \phi(y), & y \in [-a, a] \\ \beta(y + a) - a, & y \in (a, \infty), \end{cases} \quad (2.7)$$

where $\phi : [-a, a] \rightarrow [0, 1]$ is an affine mapping, $\beta > 0$ is a scaling factor, and

$$F(x) = \int_0^x p(x) dx, \quad \forall x \in [0, 1]. \quad (2.8)$$

Here $p(x)$ can be regarded a PDF and $F(x)$ a cumulative distribution function. In particular, $p(x)$ will be defined as a piecewise linear function such that $F(x)$ is a quadratic function whose inverse can be computed explicitly.

2.2. Main structure of KRnet. The main structure of KRnet is illustrated in Figure 2.1. KRnet is mainly defined by two loops: outer loop $f_{[k]}^{\text{outer}}$ and inner loop $f_{[k,i]}^{\text{inner}}$, $k = 1, \dots, K - 1$; $i = 1, \dots, L$, where the outer loop has $K - 1$ stages induced by the K mappings f_i in equation (2.3), and the inner loop has L stages indicating the length of a chain that consists of general coupling layers.

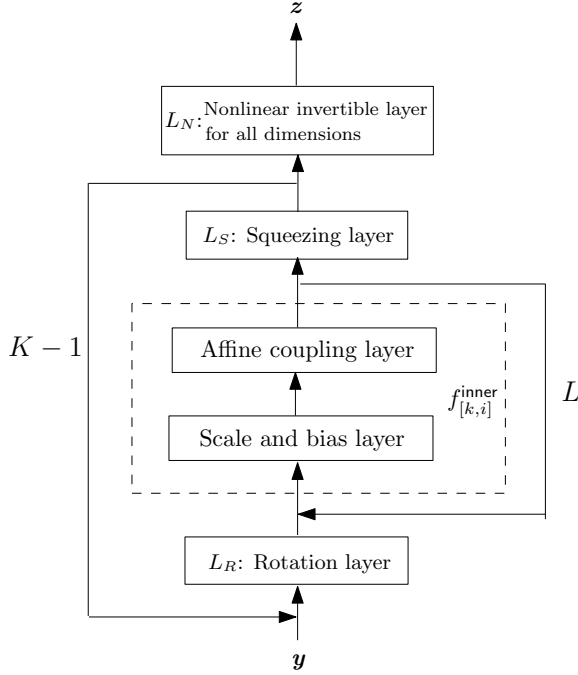


Fig. 2.1: Left: the flow chart of the block-triangular invertible mapping KRnet.

- Outer loop. The outer loop defines the main structure of KRnet that is consistent with the KR arrangement:

$$z = f(\mathbf{y}) = L_N \circ f_{[K-1]}^{\text{outer}} \circ \dots \circ f_{[1]}^{\text{outer}}(\mathbf{y}). \quad (2.9)$$

Let $\mathbf{y}_{[k]} = f_{[k]}^{\text{outer}}(\mathbf{y}_{[k-1]})$ with $\mathbf{y}_{[0]} = \mathbf{y}$, and $i = 1, \dots, K-1$. Each $\mathbf{y}_{[k]} = (\mathbf{y}_{[k],1}, \dots, \mathbf{y}_{[k],K})$ has the same partition. The i th partition will remain unchanged after $K-i+1$ stages. For example, $\mathbf{y}_{[k],K}$ will be updated only when $k=1$ and $\mathbf{y}_{[k],K-i}$ will be deactivated when $k > i+1$. This way, the number of effective dimensions decreases as k increases. Assuming that the prior distribution is chosen as a standard Gaussian, all dimensions of $\mathbf{y}_{[K-1]}$ are supposed to be independent of each other after the outer loop is completed. We then activate all the dimensions and apply the nonlinear invertible layer to $\mathbf{y}_{[K-1]}$ component-wisely before the final output. The nonlinear invertible layer generalizes the prior distribution through a component-wise nonlinear transformation.

- Inner loop. The inner loop mainly consists of a sequence of general coupling layers $f_{[k,i]}^{\text{inner}}$. Each $f_{[k,i]}^{\text{inner}}$ includes one scale and bias layer, and one affine coupling layer. $f_{[k]}^{\text{outer}}$ can be represented as:

$$f_{[k]}^{\text{outer}} = L_S \circ f_{[k,L]}^{\text{inner}} \circ \dots \circ f_{[k,1]}^{\text{inner}} \circ L_R, \quad (2.10)$$

where L_R is a rotation layer, and L_S is a squeezing layer. The affine coupling layers in $f_{[k,i]}^{\text{inner}}$ are defined between $\mathbf{y}_{[k],K+1-k}$ and the other active parts $\mathbf{y}_{[k],i}$, $i = 1, \dots, K-k$. Since we need at least two affine coupling layers for a full update of all data dimensions, we usually assume that L is even.

3. Augmented KRnet. In affine coupling layers (2.5) we need to update a certain part of the data using a mapping of the other part. The main motivation of such a strategy is to maintain the exact invertibility. The main drawback of such a strategy is that the change of a certain component y_i for each update is at most a linear function of y_i (see equation (2.5)). To alleviate such a limitation, we implement the affine coupling layer in a higher dimensional space such that the update of y_i may be in terms of all the components of \mathbf{y} through the augmented dimensions.

3.1. Introduce augmented dimensions. Suppose that $\{\mathbf{y}^{(i)}\}_{i=1}^N$ consists of samples from $\mathbf{Y} \in \mathbb{R}^n$ subject to a PDF $p_{\mathbf{Y}}(\mathbf{y})$. We augment \mathbf{Y} by another vector $\boldsymbol{\gamma} \in \mathbb{R}^m$ such that $\mathbf{Y}_{\boldsymbol{\gamma}} = (\boldsymbol{\gamma}, \mathbf{Y})$. Let $\mathbf{Z} \in \mathbb{R}^n$ have the prior distribution $p_{\mathbf{Z}}(\mathbf{z})$. The random variable \mathbf{Z} is augmented similarly by a vector $\boldsymbol{\xi}$ such that $\mathbf{Z}_{\boldsymbol{\xi}} = (\boldsymbol{\xi}, \mathbf{Z})$. Instead of considering an invertible mapping between \mathbf{Y} and \mathbf{Z} , we construct an invertible mapping between $\mathbf{Y}_{\boldsymbol{\gamma}} \in \mathbb{R}^{n+m}$ and $\mathbf{Z}_{\boldsymbol{\xi}} \in \mathbb{R}^{n+m}$ such that

$$z_{\boldsymbol{\xi}} = f_{\text{aug}}(\mathbf{y}_{\boldsymbol{\gamma}}). \quad (3.1)$$

Assuming that \mathbf{Y} and $\boldsymbol{\gamma}$ are independent. We have the PDF of $\mathbf{Y}_{\boldsymbol{\gamma}}$ as

$$p_{\mathbf{Y}}(\mathbf{y})p_{\boldsymbol{\gamma}}(\boldsymbol{\gamma}) = p_{\mathbf{Z}_{\boldsymbol{\xi}}}(f_{\text{aug}}(\mathbf{y}_{\boldsymbol{\gamma}}))|\nabla_{\mathbf{y}_{\boldsymbol{\gamma}}} f_{\text{aug}}(\mathbf{y}_{\boldsymbol{\gamma}})|, \quad \forall \mathbf{y}_{\boldsymbol{\gamma}}. \quad (3.2)$$

We now look at how the information flows through the augmented dimensions. Applying the affine coupling layers to the partition given by the data dimensions and the augmented dimensions, we have two adjacent affine coupling layers as

$$\boldsymbol{\gamma}_{[i+1]} = \boldsymbol{\gamma}_{[i]} \odot \mathbf{w}_{[i]}(\mathbf{y}_{[i]}) + \mathbf{b}_{[i]}(\mathbf{y}_{[i]}), \quad (3.3)$$

$$\mathbf{y}_{[i+1]} = \mathbf{y}_{[i]} \quad (3.4)$$

and

$$\boldsymbol{\gamma}_{[i+2]} = \boldsymbol{\gamma}_{[i+1]}, \quad (3.5)$$

$$\mathbf{y}_{[i+2]} = \mathbf{y}_{[i+1]} \odot \mathbf{w}_{[i+1]}(\boldsymbol{\gamma}_{[i+1]}) + \mathbf{b}_{[i+1]}(\boldsymbol{\gamma}_{[i+1]}), \quad (3.6)$$

where we let

$$\begin{aligned} \mathbf{w}_{[i]}(\mathbf{y}_{[i]}) &= 1 + \alpha \tanh(\mathbf{s}_{[i]}(\mathbf{y}_{[i]})), \\ \mathbf{b}_{[i]}(\mathbf{y}_{[i]}) &= e^{\boldsymbol{\beta}_{[i]}} \odot \tanh(\mathbf{t}_{[i]}(\mathbf{y}_{[i]})). \end{aligned}$$

We observe the following flow of information:

$$\mathbf{y}_{[i]} \rightarrow \boldsymbol{\gamma}_{[i+1]} \rightarrow \mathbf{y}_{[i+2]} \rightarrow \boldsymbol{\gamma}_{[i+3]} \rightarrow \dots$$

which implies that $\mathbf{y}_{[i+2]}$ may be affected by all the components of $\mathbf{y}_{[i]}$ although such a dependence is not explicit.

Similarly, the two adjacent steps in a regular KRnet can be written as

$$\mathbf{y}_{[i+1],1} = \mathbf{y}_{[i],1} \odot \mathbf{w}_{[i]}(\mathbf{y}_{[i],2}) + \mathbf{b}_{[i]}(\mathbf{y}_{[i],2}),$$

$$\mathbf{y}_{[i+1],2} = \mathbf{y}_{[i],2}$$

and

$$\mathbf{y}_{[i+2],1} = \mathbf{y}_{[i+1],1},$$

$$\mathbf{y}_{[i+2],2} = \mathbf{y}_{[i+1],2} \odot \mathbf{w}_{[i+1]}(\mathbf{y}_{[i+1],1}) + \mathbf{b}_{[i+1]}(\mathbf{y}_{[i+1],1}),$$

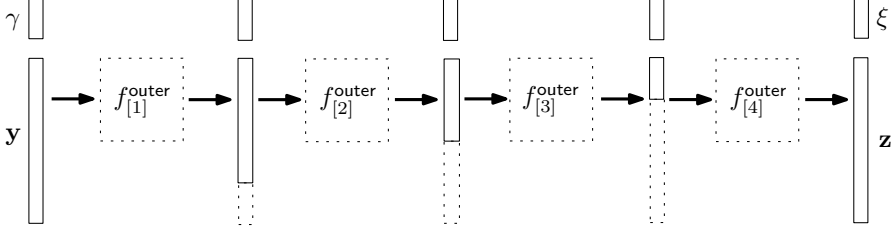


Fig. 3.1: The structure of the augmented KRnet, where the data will evolve from \mathbf{y} to \mathbf{z} and the augmented vector will evolve from $\boldsymbol{\gamma}$ to $\boldsymbol{\xi}$. In this example, the data vector has been partitioned into four parts. After each iteration $f_{[i]}^{\text{outer}}$, $i = 1, 2, 3$, an extra part of the data vector will be deactivated as illustrated by the dotted line.

where $\mathbf{y}_{[i]} = (\mathbf{y}_{[i],1}, \mathbf{y}_{[i],2})^\top$ has been partitioned to two parts. It is seen that after two steps, $\mathbf{y}_{[i+2]}$ will not depend on $\mathbf{y}_{[i]}$ in a fully nonlinear way, where $\mathbf{y}_{[i+2],1}$ depends on $\mathbf{y}_{[i],1}$ linearly and only $\mathbf{y}_{[i+2],2}$ depends on both $\mathbf{y}_{[i],1}$ and $\mathbf{y}_{[i],2}$ nonlinearly.

Although more nonlinear dependence of $\mathbf{y}_{[i+2]}$ on $\mathbf{y}_{[i]}$ has been introduced through the augmented dimensions, we need to deal with $(m+n)$ -dimensional mapping, implying a higher requirement on the complexity of the model due to the curse of dimensionality. Both issues are related to the choice of m . We note that the dependence of $\mathbf{y}_{[i+2]}$ on $\mathbf{y}_{[i]}$ prefers a larger m while a smaller m is preferred from the viewpoint of model complexity. The KRnet provides an effective way to balance these two issues.

Suppose that KRnet uses a uniform partition of $\mathbf{y} \in \mathbb{R}^n$ as $\mathbf{y} = (\mathbf{y}_1^\top, \dots, \mathbf{y}_K^\top)^\top$ with $K = n/m$. We then consider an augmented vector $\mathbf{y}_\gamma = (\boldsymbol{\gamma}^\top, \mathbf{y}^\top)^\top$ with $\boldsymbol{\gamma} \in \mathbb{R}^m$, i.e., we let \mathbf{y}_i and $\boldsymbol{\gamma}$ have the same number of dimensions. The overall structure of the augmented KRnet is illustrated in figure 3.1, which is similar to the regular KRnet. The main difference is that the augmented part will never be deactivated since it is used as a buffer zone for communicating information. Due to the refined partition of data in KRnet, only a small number of augmented dimensions is needed. For example, if we deactivate dimensions one by one in KRnet, we only need one extra dimension, i.e., $m = 1$.

3.2. Loss and the marginal PDF. Now let us look at the loss. Let p_{data} correspond to the PDF for the data. By the definition of the augmented KRnet, we need to minimize the KL divergence

$$\begin{aligned} D_{\text{KL}}(p_{\text{data}}(\mathbf{y})p_\gamma(\boldsymbol{\gamma}) \| p_{\mathbf{Y}_\gamma, \boldsymbol{\theta}}(\mathbf{y}, \boldsymbol{\gamma})) &= \int p_{\text{data}}p_\gamma \ln \frac{p_{\text{data}}p_\gamma}{p_{\mathbf{Y}_\gamma, \boldsymbol{\theta}}} d\mathbf{y}d\boldsymbol{\gamma} \\ &= \int p_{\text{data}} \ln p_{\text{data}} d\mathbf{y} + \int p_{\text{data}}p_\gamma \ln \frac{p_\gamma}{p_{\mathbf{Y}_\gamma, \boldsymbol{\theta}}} d\mathbf{y}d\boldsymbol{\gamma}, \end{aligned}$$

where we use $\boldsymbol{\theta}$ to indicate the model parameters. Since the first term on the right-hand side is only related to the data, it is equivalent to minimize the second term, which defines the loss

$$L(p_{\mathbf{Y}_\gamma, \boldsymbol{\theta}}) = \frac{1}{N} \sum_{i=1}^N \ln \frac{p_\gamma(\boldsymbol{\gamma}^{(i)})}{p_{\mathbf{Y}_\gamma, \boldsymbol{\theta}}(\boldsymbol{\gamma}^{(i)}, \mathbf{y}^{(i)})} \approx \mathbb{E}_{p_{\text{data}}p_\gamma} \left[\ln \frac{p_\gamma}{p_{\mathbf{Y}_\gamma, \boldsymbol{\theta}}} \right], \quad (3.7)$$

where for each $\mathbf{y}^{(i)}$ we associate one $\boldsymbol{\gamma}^{(i)}$ sampled independently from p_γ .

Finally, we look at the approximation of the marginal distribution $p_{\mathbf{Y}}(\mathbf{y})$. By the construction of the augmented KRnet, we have

$$p_{\mathbf{Y}}(\mathbf{y})p_{\gamma}(\gamma) \approx p_{\mathbf{Y}_{\gamma},\theta}(\mathbf{y}_{\gamma}), \quad \forall \mathbf{y}_{\gamma}. \quad (3.8)$$

To get rid of γ , we have at least two choices:

1. Integrating out γ , we have

$$p_{\mathbf{Y}}(\mathbf{y}) \approx \mathbb{E}_{p_{\gamma}} \left[\frac{p_{\mathbf{Y}_{\gamma},\theta}}{p_{\gamma}} \right] \approx \frac{1}{N} \sum_{i=1}^N \frac{p_{\mathbf{Y}_{\gamma},\theta}(\mathbf{y}, \gamma^{(i)})}{p_{\gamma}(\gamma^{(i)})}, \quad \forall \mathbf{y}, \quad (3.9)$$

where $\{\gamma^{(i)}\}_{i=1}^N$ are sampled from p_{γ} . If equation (3.8) is well approximated, the variance of the integrand should be very small, meaning that a small sample size N is sufficient.

2. Picking a certain γ^* , such that

$$p_{\mathbf{Y}}(\mathbf{y}) \approx p_{\mathbf{Y}_{\gamma},\theta}(\mathbf{y}, \gamma^*)p_{\gamma}^{-1}(\gamma^*), \quad \forall \mathbf{y}. \quad (3.10)$$

In particular, we may choose

$$\gamma^* = \arg \max p_{\gamma}(\gamma), \quad (3.11)$$

such that

$$p_{\mathbf{Y}}(\mathbf{y}) = p_{\gamma}^{-1}(\gamma^*)p_{\mathbf{Z}_{\xi}}(f_{\text{aug},\theta}(\mathbf{y}_{\gamma=\gamma^*}))|\nabla_{\mathbf{y}_{\gamma=\gamma^*}} f_{\text{aug},\theta}(\mathbf{y}_{\gamma=\gamma^*})|. \quad (3.12)$$

If we let p_{γ} and $p_{\mathbf{Z}_{\xi}}$ be the standard Gaussian, i.e.,

$$p_{\gamma} = \mathcal{N}(0, \mathbf{I}), \quad p_{\mathbf{Z}_{\xi}} = \mathcal{N}(0, \mathbf{I}), \quad (3.13)$$

we have $\gamma^* = \mathbf{0}$ and

$$p_{\mathbf{Y}}(\mathbf{y}) \approx (2\pi)^{m/2} p_{\mathbf{Z}_{\xi}}(f_{\text{aug},\theta}(\mathbf{y}_{\gamma=\mathbf{0}}))|\nabla_{\mathbf{y}_{\gamma=\mathbf{0}}} f_{\text{aug},\theta}(\mathbf{y}_{\gamma=\mathbf{0}})|. \quad (3.14)$$

3.3. The complexity of the augmented KRnet. We count the number of trainable parameters in KRnet. Let us first exclude the rotation layers and the non-linear layers, and assume that each $f_{[k]}^{\text{outer}}$ has L general coupling layers $f_{[k,i]}^{\text{inner}}$. Let n_k be the number of effective dimensions for $f_{[k]}^{\text{outer}}$ and $\text{NN}_{k,i}$ the neural network (see equation (2.6)) used in $f_{[k,i]}^{\text{inner}}$. Assume all $\text{NN}_{k,i}$ are a plain neural network with two fully connected hidden layers of width m_k . Let $\text{NN}_{k,i}$ define a mapping from $\mathbb{R}^{n_{k,1}}$ to $\mathbb{R}^{2n_{k,2}}$ with $n_k = n_{k,1} + n_{k,2}$. The number of model parameters in $\text{NN}_{k,i}$ is:

$$m_k^2 + 2m_k + m_k n_k + (m_k + 3)n_{k,2}.$$

By definition, $\text{NN}_{k,i+1}$ defines a mapping from $\mathbb{R}^{n_{k,2}}$ to $\mathbb{R}^{2n_{k,1}}$ with the number of model parameters as

$$m_k^2 + 2m_k + m_k n_k + (m_k + 3)n_{k,1}.$$

If we combine the two adjacent affine coupling layers, we obtain

$$2m_k^2 + 4m_k + 3(m_k + 1)n_k,$$

which only depends on m_k and n_k , and is independent of the partition introduced by the affine coupling layer. If L is even, we can simply regard that $\text{NN}_{k,i}$ have the same number of model parameters as

$$N_{\text{NN}_k} = m_k^2 + 2m_k + 3(m_k + 1)n_k/2.$$

We note that the main characteristic of KRnet is that a portion of dimensions will be deactivated as k increases. As n_k decreases with k , we expect that the neural network $\text{NN}_{k,i}$ may become simpler for a larger k . In other words, $N_{\text{NN},k}$ decreases as k increases. A simple choice to achieve this is to decrease the width of $\text{NN}_{k,i}$ in terms of k . We let $m_{k+1} = \lceil rm_k \rceil$ with $0 < r < 1$. The number of trainable parameters is $2n_k$ for the scale and bias layer. Assume that $n = mK$. We have $n_k = (n + m) - (k - 1)m$, $k = 1, \dots, K$. According to figures 3.1 and 2.1, we have the total number of model parameters of the augmented KRnet as

$$N_{\text{dof}} = \sum_{k=1}^K (N_{\text{NN}_k} L + 2(K - k + 2)mL) = N_{\text{aKR}} L, \quad (3.15)$$

with

$$N_{\text{aKR}} = \sum_{k=1}^K (N_{\text{NN}_k} + 2(K - k + 2)m). \quad (3.16)$$

The model complexity is mainly determined by the depth L and the number K for the partition of data.

We now look at the rotation and nonlinear invertible layers. The total number of parameters from rotation layers is

$$\sum_{i=2}^K n_k^2 = \sum_{i=2}^K (im)^2 = \frac{mn(K+1)(2K+1) - 6m^2}{6}, \quad (3.17)$$

where we assume that $n = mK$. The summation is based on two constraints: (1) only the data dimensions are rotated, and (2) at stage K , no rotation is needed for deactivation since it is right before the final output. The total number of parameter from nonlinear invertible layers is nn_p , where n_p is the number of grid points for the partition of the interval $[-a, a]$, see equation (2.7). Since both rotation layers and nonlinear invertible layers do not depend on the inner loop $f_{[k,i]}^{\text{inner}}$, the portion of the DOFs from these two types of layers is in general small.

3.4. An augmented neural ODE. Recently the connection between Resnet and the discretization of ODE has been observed and exploited to construct deep nets subject to a certain type of recursive structure [14, 20, 5]. Using an ODE to describe the evolution of \mathbf{x} in terms of t , i.e.,

$$\frac{d\mathbf{x}}{dt} = \mathbf{v}(\mathbf{x}; \boldsymbol{\theta}), \quad \forall \mathbf{x} \in [0, T], \quad (3.18)$$

neural ODE models the velocity field with a neural network and treats the learning process as a parameter estimation problem for the ODE model. In terms of our problem, we associate $\mathbf{x}(0)$ with the data distribution, and expect the distribution at $\mathbf{x}(T)$ is consistent with the prior distribution. The transformation from $\mathbf{x}(0)$ to

$\mathbf{x}(T)$, or from $\mathbf{x}(T)$ to $\mathbf{x}(0)$, will be achieved by a forward or backward numerical discretization of the ODE (3.18) respectively. We note that no numerical schemes can maintain exactly the invertibility between the forward and backward integration. This implies that if we want to maintain the exact invertibility for a continuous model, we should not assume that the velocity field $\mathbf{v}(\mathbf{x})$ is simply modeled by a general neural network. We intend to incorporate the structure of the augmented KRnet into the definition of the velocity field of an ODE and maintain the exact invertibility in the discretization as well.

3.4.1. Neural ODE from an exactly invertible mapping. We first reformulate two consecutive affine coupling layers (2.5) as

$$\begin{cases} \mathbf{z}_1 = \mathbf{y}_{[i],1} + [\mathbf{y}_{[i],1} \odot \mathbf{w}_1(\mathbf{y}_{[i],2}) + \mathbf{b}_1(\mathbf{y}_{[i],2})] \Delta t, \\ \mathbf{z}_2 = \mathbf{y}_{[i],2}, \end{cases} \quad (3.19)$$

and

$$\begin{cases} \mathbf{y}_{[i+1],1} = \mathbf{z}_1, \\ \mathbf{y}_{[i+1],2} = \mathbf{z}_2 + [\mathbf{z}_2 \odot \mathbf{w}_2(\mathbf{z}_1) + \mathbf{b}_2(\mathbf{z}_1)] \Delta t, \end{cases} \quad (3.20)$$

where $\mathbf{w}_i(\cdot)$ and $\mathbf{b}_i(\cdot)$ take the following form

$$\mathbf{w}_i(\mathbf{x}) = e^\alpha \odot \tanh(\mathbf{s}_i(\mathbf{x})), \quad (3.21)$$

$$\mathbf{b}_i(\mathbf{x}) = e^\beta \odot \tanh(\mathbf{t}_i(\mathbf{x})), \quad (3.22)$$

with $(\mathbf{s}_i, \mathbf{t}_i) = \text{NN}_i(\mathbf{x})$ is an neural network with input \mathbf{x} . Compared to equations (2.5), we replace the constant α with a trainable scaling factor e^α and a constant Δt .

Combining the two consecutive affine layers as one layer such that the whole vector gets updated, we have

$$\begin{cases} \mathbf{y}_{[i+1],1} = \mathbf{y}_{[i],1} + \mathbf{g}_1(\mathbf{y}_{[i],1}, \mathbf{y}_{[i],2}) \Delta t, \\ \mathbf{y}_{[i+1],2} = \mathbf{y}_{[i],2} + \mathbf{g}_2(\mathbf{y}_{[i+1],1}, \mathbf{y}_{[i],2}) \Delta t, \end{cases} \quad (3.23)$$

where

$$\begin{aligned} \mathbf{g}_1(\mathbf{y}_{[i],1}, \mathbf{y}_{[i],2}) &= \mathbf{y}_{[i],1} \odot \mathbf{w}_1(\mathbf{y}_{[i],2}) + \mathbf{b}_1(\mathbf{y}_{[i],2}) \\ \mathbf{g}_2(\mathbf{y}_{[i+1],1}, \mathbf{y}_{[i],2}) &= \mathbf{y}_{[i],2} \odot \mathbf{w}_2(\mathbf{y}_{[i+1],1}) + \mathbf{b}_2(\mathbf{y}_{[i+1],1}) \end{aligned}$$

Note that

$$\begin{aligned} \lim_{\Delta t \rightarrow 0} \frac{\mathbf{y}_{[i+1],1} - \mathbf{y}_{[i],1}}{\Delta t} &= \mathbf{g}_1(\mathbf{y}_{[i],1}, \mathbf{y}_{[i],2}), \\ \lim_{\Delta t \rightarrow 0} \frac{\mathbf{y}_{[i+1],2} - \mathbf{y}_{[i],2}}{\Delta t} &= \lim_{\Delta t \rightarrow 0} \mathbf{g}_2(\mathbf{y}_{[i+1],1}, \mathbf{y}_{[i],2}) = \mathbf{g}_2(\mathbf{y}_{[i],1}, \mathbf{y}_{[i],2}). \end{aligned}$$

So equation (3.23) can be regarded as an explicit one-step numerical method of a dynamical system

$$\begin{cases} \frac{d\mathbf{y}_1}{dt} = \mathbf{g}_1(\mathbf{y}_1, \mathbf{y}_2) = \mathbf{y}_1 \odot \mathbf{w}_1(\mathbf{y}_2) + \mathbf{b}_1(\mathbf{y}_2), \\ \frac{d\mathbf{y}_2}{dt} = \mathbf{g}_2(\mathbf{y}_1, \mathbf{y}_2) = \mathbf{y}_2 \odot \mathbf{w}_2(\mathbf{y}_1) + \mathbf{b}_2(\mathbf{y}_1), \end{cases} \quad (3.24)$$

where the only difference than a regular explicit Euler scheme is an updated \mathbf{y}_1 is used in the discretization of the second equation. Let $f_{\text{af},1}(\cdot)$ and $f_{\text{af},2}(\cdot)$ indicate the affine coupling layers given by equations (3.19) and (3.20) respectively. We define

$$f_{\text{af},1,2}^i(\cdot) = \underbrace{(f_{\text{af},1} \circ f_{\text{af},2}) \circ \dots \circ (f_{\text{af},1} \circ f_{\text{af},2})}_i(\cdot)$$

We see that by equations (3.19) and (3.20) the mapping

$$\mathbf{y}_{[i+1]} = \mathbf{y}_{[i]} + (f_{\text{af},1,2} - \text{Id})(\mathbf{y}_{[i]}) \quad (3.25)$$

is invertible, where $\mathbf{y}_{[i]} = (\mathbf{y}_{[i],1}, \mathbf{y}_{[i],2})$ and Id is an identity operator. In particular, the limit

$$\lim_{\Delta t \rightarrow 0} \frac{\mathbf{y}_{[i+1]} - \mathbf{y}_{[i]}}{\Delta t}$$

exists, which defines the dynamical system (3.24). This result can be generalized as the following lemma.

LEMMA 3.1. *Let $\mathbf{y}_{[i+1]} = f_{\text{af},1,2}^m(\mathbf{y}_{[i]})$ with $m \in \mathbb{N}_+$. There exists $\mathbf{g}_{[m]}(\mathbf{y}_{[i]}) \in \mathbb{R}^n$ such that*

$$\lim_{\Delta t \rightarrow 0} \frac{\mathbf{y}_{[i+1]} - \mathbf{y}_{[i]}}{\Delta t} = \mathbf{g}_{[m]}(\mathbf{y}_{[i]}),$$

i.e., $\mathbf{y}_{[i]}$ can be regarded as an approximation of the ODE

$$\dot{\mathbf{y}} = \mathbf{g}_{[m]}(\mathbf{y}),$$

if $\mathbf{g}_{[m]}(\mathbf{y})$ is sufficiently smooth.

Proof. We argue by induction. It is seen that it is true when $m = 1$. Assume the conclusion holds for $m \leq k$. We have

$$\mathbf{y}_{[i+1]} = f_{\text{af},1,2}^{k+1}(\mathbf{y}_{[i]}) = f_{\text{af},1,2} \circ f_{\text{af},1,2}^k(\mathbf{y}_{[i]}).$$

Let $\mathbf{z} = f_{\text{af},1,2}^k(\mathbf{y}_{[i]})$. We have

$$\lim_{\Delta t \rightarrow 0} \frac{\mathbf{y}_{[i+1]} - \mathbf{y}_{[i]}}{\Delta t} = \lim_{\Delta t \rightarrow 0} \frac{\mathbf{y}_{[i+1]} - \mathbf{z} + \mathbf{z} - \mathbf{y}_{[i]}}{\Delta t}.$$

According to the assumption, there exist $\mathbf{g}_{[1]}(\cdot)$ and $\mathbf{g}_{[k]}(\cdot)$ such that

$$\begin{aligned} \lim_{\Delta t \rightarrow 0} \frac{\mathbf{y}_{[i+1]} - \mathbf{z}}{\Delta t} &= \mathbf{g}_{[1]}(\mathbf{z}), \\ \lim_{\Delta t \rightarrow 0} \frac{\mathbf{z} - \mathbf{y}_{[i]}}{\Delta t} &= \mathbf{g}_{[k]}(\mathbf{y}_{[i]}). \end{aligned}$$

We can then let

$$\mathbf{g}_{[k+1]}(\mathbf{y}_{[i]}) = \mathbf{g}_{[1]}(\mathbf{y}_{[i]}) + \mathbf{g}_{[k]}(\mathbf{y}_{[i]}), \quad (3.26)$$

since

$$\lim_{\Delta t \rightarrow 0} \mathbf{z} = \lim_{\Delta t \rightarrow 0} f_{\text{af},1,2}^k(\mathbf{y}_{[i]}) = \mathbf{y}_{[i]}$$

by definition. \square

REMARK 3.2. *It is seen from equation (3.26) that every time $f_{\text{af},1,2}(\cdot)$ is introduced, a vector function $\mathbf{g}_{[1]}(\mathbf{y}_{[i]}; \boldsymbol{\theta}_j)$ is added to the velocity field such that*

$$\mathbf{g}_{[k]}(\mathbf{y}_{[i]}) = \sum_{j=1}^k \mathbf{g}_{[1]}(\mathbf{y}_{[i]}; \boldsymbol{\theta}_j), \quad (3.27)$$

where we include the model parameters $\boldsymbol{\theta}_j$ to differentiate the k functions $g_{[1]}(\mathbf{y}_{[i]}; \boldsymbol{\theta}_j)$, $j = 1, \dots, k$.

REMARK 3.3. The mapping $\mathbf{y}_{[i+1]} = f_{\text{af},1,2}^m(\mathbf{y}_{[i]})$ can be regarded as a multi-stage process that is defined on a time interval $[0, \Delta t]$. Let

$$\mathbf{y}_{[i+\frac{j}{m}]} = f_{\text{af},1,2} \left(\mathbf{y}_{[i+\frac{j-1}{m}]}, \boldsymbol{\theta}_j, \Delta t \right), \quad j = 1, \dots, m, \quad (3.28)$$

where the transform from $\mathbf{y}_{[i+\frac{j-1}{m}]}$ to $\mathbf{y}_{[i+\frac{j}{m}]}$ is achieved at stage j . We can then decompose $\mathbf{y}_{[i+1]} = f_{\text{af}}^{2m}(\mathbf{y}_{[i]})$ as

$$\begin{cases} \mathbf{y}_{[i+\frac{1}{m}]} &= f_{\text{af},1,2} \left(\mathbf{y}_{[i]}, \boldsymbol{\theta}_1, \Delta t \right) \\ \mathbf{y}_{[i+\frac{2}{m}]} &= f_{\text{af},1,2} \left(\mathbf{y}_{[i+\frac{1}{m}]}, \boldsymbol{\theta}_2, \Delta t \right) \\ &\dots \\ \mathbf{y}_{[i+1]} &= f_{\text{af},1,2} \left(\mathbf{y}_{[i+\frac{m-1}{m}]}, \boldsymbol{\theta}_m, \Delta t \right) \end{cases} \quad (3.29)$$

Note that the following two limits exist

$$\lim_{\Delta t \rightarrow 0} \frac{\mathbf{y}_{[i+\frac{j}{m}]} - \mathbf{y}_{[i+\frac{j-1}{m}]}}{\Delta t} = \mathbf{g}_{[1]}(\mathbf{y}_{[i+\frac{j-1}{m}]}, \boldsymbol{\theta}_j), \quad \lim_{\Delta t \rightarrow 0} \mathbf{y}_{[i+\frac{j}{m}]} = \mathbf{y}_{[i]}.$$

We then have

$$\lim_{\Delta t \rightarrow 0} \frac{\mathbf{y}_{[i+1]} - \mathbf{y}_{[i]}}{\Delta t} = \lim_{\Delta t \rightarrow 0} \sum_{j=1}^m \frac{\mathbf{y}_{[i+\frac{j}{m}]} - \mathbf{y}_{[i+\frac{j-1}{m}]}}{\Delta t} = \sum_{j=1}^m \mathbf{g}_{[1]}(\mathbf{y}_{[i]}, \boldsymbol{\theta}_j).$$

Compared to the multi-stage numerical schemes such as the Runge-Kutta method for the numerical approximation of ODE, we use multiple stages to achieve the exact invertibility rather than a better accuracy.

3.4.2. Generalize the model. We generalize the model in lemma 3.1 by integrating an augmented KRnet into the definition of the velocity field. We simply consider the recursive formula defined by an augmented KRnet:

$$\mathbf{y}_{\gamma,[i+1]} = f_{\text{KRnet}}(\mathbf{y}_{\gamma,[i]}) = (L_S \circ f_{[\text{af},1,2],K}^{m_K}) \circ \dots \circ (L_S \circ f_{[\text{af},1,2],1}^{m_1})(\mathbf{y}_{\gamma,[i]}), \quad (3.30)$$

where $f_{[\text{af},1,2],k}$ indicates two affine coupling layers given by equations (3.19) and (3.20) at stage k , and the active dimensions for $f_{[\text{af},1,2],k}^{m_k}$ are defined with respect to figure 3.1. Following remark 3.3, $f_{\text{KRnet}}(\cdot)$ can be understood as a one-step method, where multiple stages are used to maintain the exact invertibility. In other words, the following limit exists

$$\lim_{\Delta t \rightarrow 0} \frac{\mathbf{y}_{\gamma,[i+1]} - \mathbf{y}_{\gamma,[i]}}{\Delta t} = \lim_{\Delta t \rightarrow 0} \frac{f_{\text{KRnet}}(\mathbf{y}_{\gamma,[i]}) - \mathbf{y}_{\gamma,[i]}}{\Delta t} = \mathbf{g}_{\text{KRnet}}(\mathbf{y}_{\gamma,[i]}),$$

which suggests a dynamical system

$$\frac{d\mathbf{y}_{\gamma}}{dt} = \mathbf{g}_{\text{KRnet}}(\mathbf{y}_{\gamma}). \quad (3.31)$$

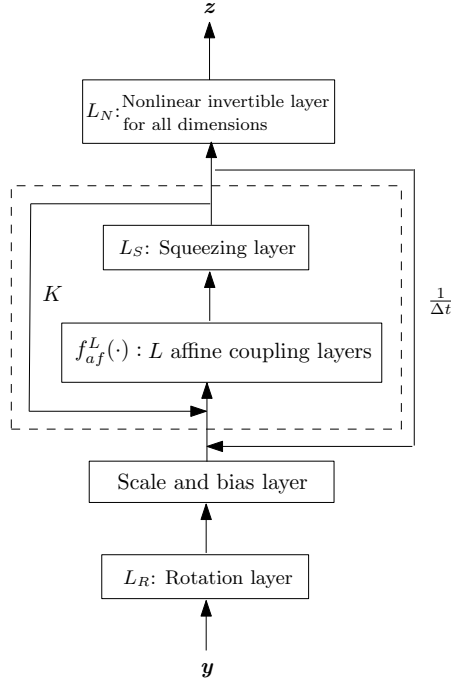


Fig. 3.2: The flow chart of an neural ODE based on an invertible recursive mapping given by KRnet. The dashed rectangle indicates the operation for each time step, where the time interval $[0, 1]$ is uniformly discretized with step size Δt .

3.4.3. The adjoint method for an invertible mapping. One difficulty of neural ODE is that as the time step size decreases the size of the computation graph for automatic differentiation may explode and exhaust the computer memory quickly. We then need to consider the adjoint method to compute the gradient for the optimizer. The adjoint method is defined with respect to an ODE, where the system is invertible. For a certain path from $\mathbf{x}(0)$ to $\mathbf{x}(T)$, the adjoint method needs to integrate the ODE (3.18) backwardly from $\mathbf{x}(T)$ to $\mathbf{x}(0)$. Once the ODE is discretized, the exact invertibility will be lost at the discrete level, implying that the adjoint method in general cannot yield the gradient up to the machine accuracy.

In our model, we do not need to formulate the adjoint method in terms of the ODE since the exact invertibility is kept by definition. We first consider the following optimization problem

$$\min_{\boldsymbol{\theta}} L = -\mathbb{E}_{p_{\text{data}}} \log p_{\mathbf{Y}}(\mathbf{y}; \boldsymbol{\theta}) = -\frac{1}{N} \sum_{j=1}^N \log p_{\mathbf{Y}}(\mathbf{y}^{(j)}; \boldsymbol{\theta}) \quad (3.32)$$

subject to the following constraints:

$$\mathbf{y}_{[i+1]} = F(\mathbf{y}_{[i]}, \boldsymbol{\theta}), \quad i = 0, 1, \dots, n-1, \quad (3.33)$$

where F can be regarded as an invertible mapping defined by equation (3.23), and the subscript $*_i$ indicates the temporal discretization. We let $\mathbf{y}_{[0]} = \mathbf{y}$ and assume

$\mathbf{y}_{[n]} = \mathbf{z}$ has a standard Gaussian distribution. From equation (2.1), we have

$$\begin{aligned}\log p_{\mathbf{Y}}(\mathbf{y}) &= \log p_{\mathbf{Z}}(\mathbf{z}) + \sum_{i=0}^{n-1} \log |\det \nabla_{\mathbf{y}_{[i]}} \mathbf{y}_{[i+1]}| \\ &= \log p_{\mathbf{Z}}(\mathbf{z}) + \sum_{i=0}^{n-1} g_{[i]}(\mathbf{y}_{[i]}, \boldsymbol{\theta}),\end{aligned}\quad (3.34)$$

where $g_{[i]}(\mathbf{y}_{[i]}, \boldsymbol{\theta})$ can be explicitly computed by the definition of the affine coupling layer. For simplicity, we only consider one data point and ignore the superscript such that

$$L = -\log p_{\mathbf{Y}}(\mathbf{y}; \boldsymbol{\theta}). \quad (3.35)$$

To compute $\nabla_{\boldsymbol{\theta}} L$, we consider the following Lagrangian:

$$\mathcal{L} = -\log p_{\mathbf{Z}}(\mathbf{z}) - \sum_{i=0}^{n-1} g_{[i]}(\mathbf{y}_{[i]}, \boldsymbol{\theta}) - \sum_{i=0}^{n-1} \boldsymbol{\lambda}_{[i]}^{\top} (\mathbf{y}_{[i+1]} - F_{[i]}(\mathbf{y}_{[i]}, \boldsymbol{\theta})). \quad (3.36)$$

The key idea of the adjoint method is to choose appropriate Lagrange multipliers $\boldsymbol{\lambda}_{[i]}$ such that the computation of the gradient is convenient. We have

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} \mathcal{L} &= (\nabla_{\boldsymbol{\theta}} \mathbf{z})^{\top} (-\nabla_{\mathbf{z}} \log p_{\mathbf{Z}}(\mathbf{z})) - \sum_{i=1}^{n-1} ((\nabla_{\boldsymbol{\theta}} \mathbf{y}_{[i]})^{\top} \nabla_{\mathbf{y}_{[i]}} g_{[i]} + \nabla_{\boldsymbol{\theta}} g_{[i]}) - \nabla_{\boldsymbol{\theta}} g_{[0]} \\ &\quad - \sum_{i=1}^{n-1} (\nabla_{\boldsymbol{\theta}} \mathbf{y}_{[i+1]} - \nabla_{\mathbf{y}_{[i]}} F_{[i]} \nabla_{\boldsymbol{\theta}} \mathbf{y}_{[i]} - \nabla_{\boldsymbol{\theta}} F_{[i]})^{\top} \boldsymbol{\lambda}_{[i]} \\ &\quad - (\nabla_{\boldsymbol{\theta}} \mathbf{y}_{[1]} - \nabla_{\boldsymbol{\theta}} F_{[0]})^{\top} \boldsymbol{\lambda}_{[0]} \\ &= (\nabla_{\boldsymbol{\theta}} \mathbf{z})^{\top} (-\nabla_{\mathbf{z}} \log p_{\mathbf{Z}}(\mathbf{z}) - \boldsymbol{\lambda}_{[n-1]}) - \sum_{i=0}^{n-1} \nabla_{\boldsymbol{\theta}} g_{[i]} + \sum_{i=0}^{n-1} (\nabla_{\boldsymbol{\theta}} F_{[i]})^{\top} \boldsymbol{\lambda}_{[i]} \\ &\quad - \sum_{i=1}^{n-1} (\nabla_{\boldsymbol{\theta}} \mathbf{y}_{[i]})^{\top} (\boldsymbol{\lambda}_{[i-1]} - (\nabla_{\mathbf{y}_{[i]}} F_{[i]})^{\top} \boldsymbol{\lambda}_{[i]} + \nabla_{\mathbf{y}_{[i]}} g_{[i]}).\end{aligned}$$

We let

$$\begin{cases} \boldsymbol{\lambda}_{[i-1]} &= (\nabla_{\mathbf{y}_{[i]}} F_{[i]})^{\top} \boldsymbol{\lambda}_{[i]} - \nabla_{\mathbf{y}_{[i]}} g_{[i]}, & i = n-1, \dots, 1 \\ \boldsymbol{\lambda}_{[n-1]} &= -\nabla_{\mathbf{z}} \log p_{\mathbf{Z}}(\mathbf{z}) \end{cases} \quad (3.37)$$

and obtain that

$$\nabla_{\boldsymbol{\theta}} \mathcal{L} = \nabla_{\boldsymbol{\theta}} L = \sum_{i=0}^{n-1} ((\nabla_{\boldsymbol{\theta}} F_{[i]})^{\top} \boldsymbol{\lambda}_{[i]} - \nabla_{\boldsymbol{\theta}} g_{[i]}). \quad (3.38)$$

Except for the neural ODE layer, we may add other types of layers into the model. The following lemma provides a more general result for the adjoint method in terms of an invertible mapping:

LEMMA 3.4. *Consider a general invertible mapping:*

$$\mathbf{y}_{[i+1]} = F_{[i]}(\mathbf{y}_{[i]}; \boldsymbol{\theta}_{[i]}), \quad \mathbf{y}_{[0]} = \mathbf{y}, \quad \mathbf{y}_{[n]} = \mathbf{z}, \quad i = 0, \dots, n-1, \quad (3.39)$$

where we let

$$g_{[i]}(\mathbf{y}_{[i]}, \boldsymbol{\theta}_{[i]}) = \log |\det \nabla_{\mathbf{y}_{[i]}} \mathbf{y}_{[i+1]}|.$$

Assume that the loss L is given by equation (3.35). The following two sequences $\boldsymbol{\lambda}_{[i]}$ and $\mathbf{y}_{[i]}$ can be computed backwardly:

$$\begin{cases} \boldsymbol{\lambda}_{[i-1]} = (\nabla_{\mathbf{y}_{[i]}} F_{[i]})^\top \boldsymbol{\lambda}_{[i]} - \nabla_{\mathbf{y}_{[i]}} g_{[i]}, & i = n-1, \dots, 1 \\ \boldsymbol{\lambda}_{n-1} = -\nabla_{\mathbf{z}} \log p_{\mathbf{Z}}(\mathbf{z}), \\ \mathbf{y}_{[i]} = F_{[i]}^{-1}(\mathbf{y}_{[i+1]}, \boldsymbol{\theta}_{[i]}), & i = n-1, \dots, 0. \end{cases} \quad (3.40)$$

We have

$$\nabla_{\boldsymbol{\theta}_{[i]}} L = (\nabla_{\boldsymbol{\theta}_{[i]}} F_{[i]})^\top \boldsymbol{\lambda}_{[i]} - \nabla_{\boldsymbol{\theta}_{[i]}} g_{[i]}, \quad i = 0, 1, \dots, n-1. \quad (3.41)$$

If $\boldsymbol{\theta}_{[i]} = \boldsymbol{\theta}$ for $i \in \mathcal{I} \subset \{0, 1, \dots, n-1\}$, we have

$$\nabla_{\boldsymbol{\theta}} L = \sum_{i \in \mathcal{I}} ((\nabla_{\boldsymbol{\theta}} F_{[i]})^\top \boldsymbol{\lambda}_{[i]} - \nabla_{\boldsymbol{\theta}} g_{[i]}). \quad (3.42)$$

Proof. We consider the Lagrangian

$$\mathcal{L} = -\log p_{\mathbf{Z}}(\mathbf{z}) - \sum_{i=0}^{n-1} g_{[i]}(\mathbf{y}_{[i]}, \boldsymbol{\theta}_i) - \sum_{i=0}^{n-1} \boldsymbol{\lambda}_{[i]}^\top (\mathbf{y}_{[i+1]} - F_{[i]}(\mathbf{y}_{[i]}, \boldsymbol{\theta}_{[i]})). \quad (3.43)$$

For $\boldsymbol{\theta}_{[k]}$, $0 \leq k < n-1$, we have

$$\begin{aligned} \nabla_{\boldsymbol{\theta}_{[k]}} \mathcal{L} &= (\nabla_{\boldsymbol{\theta}_{[k]}} \mathbf{z})^\top (-\nabla_{\mathbf{z}} \log p_{\mathbf{Z}}(\mathbf{z})) - \sum_{i=k+1}^{n-1} (\nabla_{\boldsymbol{\theta}_{[k]}} \mathbf{y}_{[i]})^\top \nabla_{\mathbf{y}_{[i]}} g_{[i]} - \nabla_{\boldsymbol{\theta}_{[k]}} g_{[k]} \\ &\quad - \sum_{i=k+1}^{n-1} (\nabla_{\boldsymbol{\theta}_{[k]}} \mathbf{y}_{[i+1]} - \nabla_{\mathbf{y}_{[i]}} F_{[i]}(\mathbf{y}_{[i]}, \boldsymbol{\theta}_{[i]}))^\top \nabla_{\boldsymbol{\theta}_{[k]}} \mathbf{y}_{[i]} \\ &\quad - (\nabla_{\boldsymbol{\theta}_{[k]}} \mathbf{y}_{[k+1]} - \nabla_{\boldsymbol{\theta}_{[k]}} F_{[k]}(\mathbf{y}_{[k]}, \boldsymbol{\theta}_{[k]}))^\top \boldsymbol{\lambda}_{[k]} \\ &= (\nabla_{\boldsymbol{\theta}_{[k]}} \mathbf{z})^\top (-\nabla_{\mathbf{z}} \log p_{\mathbf{Z}}(\mathbf{z}) - \boldsymbol{\lambda}_{[n-1]}) \\ &\quad - \sum_{i=k+1}^{n-1} (\nabla_{\boldsymbol{\theta}_{[k]}} \mathbf{y}_{[i]})^\top (\nabla_{\mathbf{y}_{[i]}} g_{[i]} - (\nabla_{\mathbf{y}_{[i]}} F_{[i]}(\mathbf{y}_{[i]}, \boldsymbol{\theta}_{[i]}))^\top \boldsymbol{\lambda}_{[i]} + \boldsymbol{\lambda}_{[i-1]}) \\ &\quad + (\nabla_{\boldsymbol{\theta}_{[k]}} F_{[k]}(\mathbf{y}_{[k]}, \boldsymbol{\theta}_{[k]}))^\top \boldsymbol{\lambda}_{[k]} - \nabla_{\boldsymbol{\theta}_{[k]}} g_{[k]}. \end{aligned}$$

If $k = n-1$, we have

$$\begin{aligned} \nabla_{\boldsymbol{\theta}_{[k]}} \mathcal{L} &= (\nabla_{\boldsymbol{\theta}_{[n-1]}} \mathbf{z})^\top (-\nabla_{\mathbf{z}} \log p_{\mathbf{Z}}(\mathbf{z}) - \boldsymbol{\lambda}_{[n-1]}) \\ &\quad + (\nabla_{\boldsymbol{\theta}_{[n-1]}} F_{[n-1]}(\mathbf{y}_{[n-1]}, \boldsymbol{\theta}_{[n-1]}))^\top \boldsymbol{\lambda}_{[n-1]} - \nabla_{\boldsymbol{\theta}_{[n-1]}} g_{[n-1]}. \end{aligned}$$

Letting

$$\begin{cases} \boldsymbol{\lambda}_{[i-1]} &= (\nabla_{\mathbf{y}_{[i]}} F_{[i]})^\top \boldsymbol{\lambda}_{[i]} - \nabla_{\mathbf{y}_{[i]}} g_{[i]}, & i = n-1, \dots, k+1, \\ \boldsymbol{\lambda}_{[n-1]} &= -\nabla_{\mathbf{z}} \log p_{\mathbf{Z}}(\mathbf{z}), \end{cases} \quad (3.44)$$

we have

$$\nabla_{\boldsymbol{\theta}_{[k]}} \mathcal{L} = (\nabla_{\boldsymbol{\theta}_{[k]}} F_{[k]}(\mathbf{y}_{[k]}, \boldsymbol{\theta}_{[k]}))^{\top} \boldsymbol{\lambda}_{[k]} - \nabla_{\boldsymbol{\theta}_{[k]}} g_{[k]}. \quad (3.45)$$

Since the recursive formula (3.44) holds for any $0 \leq k \leq n-1$, we obtain the conclusion. \square

REMARK 3.5. For the gradient of the cross entropy (3.32) defined by N data points, we need to collect the contributions from all data points to the gradient using Lemma 3.4. We have

$$\nabla_{\boldsymbol{\theta}} L = \frac{1}{N} \sum_{i=1}^N \sum_{j \in \mathcal{I}} \left((\nabla_{\boldsymbol{\theta}} F_{[j]}(\mathbf{y}_{[j]}^{(i)}, \boldsymbol{\theta}_{[j]}))^{\top} \boldsymbol{\lambda}_{[j]}^{(i)} - \nabla_{\boldsymbol{\theta}} g_{[j]}(\mathbf{y}_{[j]}^{(i)}, \boldsymbol{\theta}_{[j]}) \right), \quad (3.46)$$

where $\boldsymbol{\theta}_{[j]} = \boldsymbol{\theta}$ for $j \in \mathcal{I}$, and the superscript $^{(i)}$ indicates each data point.

3.5. A summary of the main features of KRnet. To this end, we we have developed various techniques that either improve the performance of KRnet as a discrete model or reformulate it as a continuous model. We summarize some useful features of KRnet as follows:

1. *The Knothe-Rosenblatt rearrangement* defines the main structure of the KRnet for both the discrete and continuous models.
2. *The rotation layer* provides a mechanism, which is similar to the principle component analysis, to pick a certain set of dimensions to deactivate.
3. *The nonlinear layer* provides a much larger family of prior distributions than the commonly used standard Gaussian distributions through a component-wise nonlinear transformation.
4. *The augmented dimensions* provide a buffer zone for the data dimensions to exchange nonlinear information more effectively.
5. *The KRnet-ODE* integrates KRnet into a continuous model as neural ODE while the exact invertibility is maintained. The adjoint method can be formulated with respect to the discrete model instead of the continuous one such that the gradient of the loss can be computed exactly.

The features 1-4 can be coupled to improve the performance of a discrete model; The features 1, 4 and 5 can be coupled to improve the performance of a continuous model.

3.6. Density estimation and approximation via KRnet. The developed KRnets may be used to construct a PDF model for both density estimation and approximation. For density estimation, we assume that the empirical distribution $p_{\text{data}}(\mathbf{y})$ is given, and for density estimation, we assume that the unnormalized PDF $\hat{p}_{\mathbf{Y}}(\mathbf{y}) = C p_{\text{ref}, \mathbf{Y}}(\mathbf{y})$ is given, where $p_{\text{ref}, \mathbf{Y}}$ is the true PDF and C is an unknown constant. For both density estimation and approximation, we can use the Kullback-Leibler (KL) divergence to minimize the difference between the given distribution and the PDF model $p_{\text{KRnet}}(\mathbf{y})$ based on KRnet.

For density estimation, we consider the KL divergence

$$D_{\text{KL}}(p_{\text{data}} \| p_{\text{KRnet}, \mathbf{Y}}) = h(p_{\text{data}}, p_{\text{KRnet}, \mathbf{Y}}) - h(p_{\text{data}}), \quad (3.47)$$

where the first term on the right-hand side is the differential cross entropy of $p_{\text{KRnet}, \mathbf{Y}}$ relative to p_{data} , and the second term is the differential entropy of p_{data} . Since $h(p_{\text{data}})$ is independent of $p_{\text{KRnet}, \mathbf{Y}}$, minimizing the KL divergence $D_{\text{KL}}(p_{\text{data}} \| p_{\text{KRnet}, \mathbf{Y}})$ is equivalent to minimizing the differential cross entropy $h(p_{\text{data}}, p_{\text{KRnet}, \mathbf{Y}})$, which is also equivalent to maximizing the likelihood.

For density approximation, we consider the KL divergence

$$\begin{aligned} D_{\text{KL}}(p_{\text{KRnet}, \mathbf{Y}} \| p_{\text{ref}, \mathbf{Y}}) &= \mathbb{E}_{p_{\text{KRnet}, \mathbf{Y}}} \left[\ln \frac{p_{\text{KRnet}, \mathbf{Y}}}{\hat{p}_{\mathbf{Y}}} d\mathbf{y} \right] + \ln C, \\ &\approx \frac{1}{N} \sum_{i=1}^N \ln \frac{p_{\text{KRnet}, \mathbf{Y}}(\mathbf{y}^{(i)})}{\hat{p}_{\mathbf{Y}}(\mathbf{y}^{(i)})}, \end{aligned} \quad (3.48)$$

where $\{\mathbf{y}^{(i)}\}_{i=1}^N$ are samples from p_{KRnet} . It is seen that we only need to minimize the first term on the right-hand side and the unknown constant C does not affect the optimization. In contrast to the density estimation, we use the relative entropy from $p_{\text{ref}, \mathbf{Y}}$ to $p_{\text{KRnet}, \mathbf{Y}}$ to avoid the integration in terms of $p_{\text{ref}, \mathbf{Y}}$, where the integration with respect to p_{KRnet} can be easily approximated by the Monte Carlo method thanks to the generative model. For the augmented KRnet, we may consider the following KL divergence:

$$D_{\text{KL}}(p_{\text{KRnet_aug}, \mathbf{Y}, \gamma}(\mathbf{Y}, \gamma) \| p_{\text{ref}, \mathbf{Y}} p_{\gamma}), \quad (3.49)$$

where $p_{\text{KRnet_aug}, \mathbf{Y}, \gamma}(\mathbf{Y}, \gamma)$ is the joint PDF of \mathbf{Y} and γ induced by the augmented KRnet.

4. Numerical examples. In this section we present some numerical experiments including one-, two-, four- and eight-dimensional problems, where PDFs with different types of support are considered. All the models have been trained with ADAM method subject to a fixed learning rate 0.001 [18]. If no additional clarification is given, the neural networks (2.6) for the affine coupling layer always have two fully-connected hidden layers of 24 neurons. When the nonlinear invertible layers (see equation (2.7)) are needed, the interval $[-20, 20]$, i.e., $a = 20$, is discretized to 32 elements, and $\beta = 10^{-10}$. The elements are nonuniform, where the element size increases from the middle to both sides with a ratio 1.15.

4.1. The augmented ODE model for the approximation of 1d PDFs.

The simplest case of equation (3.31) includes one data dimension and one augmented dimension, which takes the following form:

$$\begin{cases} \dot{\gamma} &= v_1(\gamma, y) = \gamma w_1(y) + b_1(y), \\ \dot{y} &= v_2(\gamma, y) = y w_2(\gamma) + b_2(\gamma), \end{cases} \quad (4.1)$$

subject to the constraint

$$\rho_{t=0}(\gamma, y) = p(\gamma)f(y), \quad \rho_{t=1}(\gamma, y) = p(\gamma)p(y), \quad (4.2)$$

where $p(\cdot)$ is a standard Gaussian PDF and $f(\cdot)$ an arbitrary PDF. We know that ρ satisfies the Liouville equation:

$$\partial_t \rho + \nabla \cdot (\rho \mathbf{v}) = 0, \quad (4.3)$$

from which we have

$$\partial_t \ln \rho = \frac{1}{\rho} \partial_t \rho = -\frac{1}{\rho} (\rho \nabla \cdot \mathbf{v} + \mathbf{v} \cdot \nabla \rho) = -\nabla \cdot \mathbf{v} - \mathbf{v} \cdot \nabla \ln \rho, \quad (4.4)$$

i.e.,

$$\frac{d \ln \rho}{dt} = -\nabla \cdot \mathbf{v} = -w_1(y(t)) - w_2(\gamma(t)), \quad (4.5)$$

subject to the boundary conditions (4.2). Due to the exact invertibility, the right-hand side of equation (4.5) is given by two functions in terms of y and γ respectively. However, according to equation (4.1), $y(t)$ and $\gamma(t)$ depend on each other. In terms of equation (4.5) and the boundary conditions (4.2), $w_1(y)$, $w_2(\gamma)$, $b_1(y)$ and $b_2(\gamma)$ need to be chosen such that

$$\ln p(\gamma(1))p(y(1)) = \ln p(\gamma(0))f(y(0)) - \int_0^1 (w_1(y(t)) + w_2(\gamma(t)))dt. \quad (4.6)$$

Let us look at a simple case, where $w_1(\cdot) = w_2(\cdot) = 0$. In other words, the dynamics given by (4.7) preserves volume. We have the ODE as

$$\begin{cases} \dot{\gamma} &= b_1(y), \\ \dot{y} &= b_2(\gamma). \end{cases} \quad (4.7)$$

Let $\hat{b}'_1(y) = b_1(y)$ and $\hat{b}'_2(\gamma) = b_2(\gamma)$. We then have

$$\hat{b}_2(\gamma(t)) = \hat{b}_1(y(t)) + C, \quad (4.8)$$

where $C = \hat{b}_2(\gamma(0)) - \hat{b}_1(y(0))$ is determined by the initial condition. Due to the first integral (4.8), we expect that both $b_1(y)$ and $b_2(\gamma)$ are complex enough for a good approximation. For example, if we simply let $b_2(\gamma) = 1$, we have

$$y(t) = y(0) + t, \quad \gamma(t) = \hat{b}_1(y(t)) - \hat{b}_1(y(0)) - \gamma(0).$$

We then model $b_1(y)$ such that

$$p(y(0) + 1)p(\gamma(0) + \hat{b}_1(y(0) + 1) - \hat{b}_1(y(0))) \approx p(\gamma(0))f(y(0)).$$

It is easy to see that no matter how complex $b_1(y)$ is the above approximation may not be good enough since $y(t)$ and $\gamma(t)$ cannot be independent for the case that $b_2(\gamma) = 1$. However, we should note $y(t)$ and $\gamma(t)$ may be independent of each other if they both depend on $y(0)$ and $\gamma(0)$ in a certain way. One example is the Box–Muller transform, which maps two independent uniform random variables to two independent Gaussian random variables through an invertible mapping. So both $b_1(y)$ and $b_2(\gamma)$ need to be complex enough. Furthermore, when $w_1(y)$ and $w_2(\gamma)$ are included into the model, the modeling capability will be improved further.

Since γ corresponds to an augmented dimension, equation (4.1) can be regarded as a neural ODE for the approximation of an arbitrary PDF $f(y)$. We now look at how well model (4.1) can evolve a standard Gaussian distribution $p(y)$ to an arbitrary distribution $f(y)$. We will consider four cases, where the support of $f(y)$ is $(-\infty, \infty)$, $(0, \infty)$, $[-1, 1]$ and $[-1.5, -0.5] \cup [0.5, 1.5]$, respectively. The training set has 3.2×10^5 samples. The Adams method is subject to 4 minibatches. Let us refer to model (4.1) as augmented KRnet_ODE. We will compare its performance to the augmented KRnet. For the neural ODE, we consider a uniform temporal mesh with $\Delta t = 0.1$. Both the augmented KRnet and KRnet_ODE are defined by a sequence $f_{\text{af}}^L(\cdot)$ of affine coupling layers between y and γ , where L is the number of affine coupling layers. In the augmented KRnet, $f_{\text{af}}^L(\cdot)$ will achieve the whole transformation from data distribution to the prior distribution while in the augmented KRnet_ODE, $f_{\text{af}}^L(\cdot)$ only implements the transformation for one time step. Note that the definition of $f_{\text{af}}(\cdot)$ for the KRnet is slightly different than that for the KRnet_ODE (see equations (2.5) and (3.20)).

The prior distribution is always the standard Gaussian no matter that the target distribution has a compact support or not. When the model $p_{\mathbf{Y}_\gamma}$ converges to $f(y)p(\gamma)$, the loss function is

$$\mathbb{E}_{f(y)p(\gamma)} \left[\ln \frac{p(\gamma)}{p_{\mathbf{Y}_\gamma}} \right] \rightarrow -\mathbb{E}_{f(y)p(\gamma)} [\ln f(y)] = -\mathbb{E}_{f(y)} \ln f(y),$$

which is the differential entropy $h(f)$ of $f(y)$. We then define a relative error

$$\delta = \frac{|L - h(f)|}{h(f)} \quad (4.9)$$

to measure the quality of the corresponding PDF model. We consider the following cases:

Case (i): $f(y)$ is Logistic distribution on $(-\infty, \infty)$. Consider the logistic distribution with the location parameter $\mu = 0$ and the scale parameter $s = 2$. The differential entropy is $h(f) = 2.0 + \ln(2.0)$. The relative errors for this case are plotted in the left plot of figure 4.1. It is seen that $L = 2$ works well for both KRnet and KRnet_ODE. The high oscillations are due to the uncertainty from data since the loss function is an approximation of the differential entropy given by the Monte Carlo method.

Case (ii): $f(y)$ is Lognormal distribution on $(0, \infty)$. The lognormal distribution is given by the exponential function of a standard normal random variable. The differential entropy is $\ln(2\pi)/2 + 1/2$. The relative errors for this case are plotted in the right plot of figure 4.1. Since the positive densities on $(-\infty, \infty)$ needs to be mapped to $(0, \infty)$, the transformation is more demanding than the previous case. When $L = 2$, the KRnet_ODE has a slightly smaller error than the KRnet. When $L = 4$, both models have an error that is comparable to the uncertainty from data.

Case (iii): $f(y)$ is uniform on $[-1, 1]$. The differential entropy for the uniform distribution is $\ln(2)$. For this case, the positive densities on $(-\infty, \infty)$ needs to be mapped to $[-1, 1]$. As L increases, the performance of both KRnet and KRnet_ODE improves. It appears that the KRnet is more effective to reduce the loss while the KRnet_ODE is more robust. It is seen that the error given by KRnet with $L = 4$ is comparable to the error given by KRnet_ODE with $K = 8$. When $L = 2$, it takes KRnet a long time to find a good local minimizer.

Case (iv): $f(y)$ is uniform on $[-1.5, 0.5] \cup [0.5, 1.5]$. Compared to the previous uniform distribution, similar behavior is observed for both KRnet and KRnet_ODE except that the error is larger for the same configuration due to the more demanding requirements on the transformation. We plot some approximate PDFs in figure 4.3 for this case and the lognormal distribution in case (ii). It is seen that the KRnet handles discontinuities slightly better than the continuous flow defined by an ODE.

Note that for all four cases, we map the prior Gaussian distribution defined on $(-\infty, \infty)$ to the data distribution whether the data are subject to a compact support or not. Both augmented KRnet and augmented KRnet_ODE demonstrate effectiveness and flexibility for the density estimation. Of course, we can integrate other techniques such as regularization and data preprocessing whenever necessary. For example, if the data are defined on a compact support, say $[\delta, 1 - \delta]$ with $\delta > 0$, we may use the Logistic transformation

$$y = \frac{s}{2} \log \frac{x}{1-x}, \quad x = \frac{1}{2} (\tanh(x/s) + 1) \quad (4.10)$$

to map $x \in (0, 1)$ to $y \in (-\infty, \infty)$ such that the data distribution and the prior distribution have the same support. The results of such a strategy are plotted in

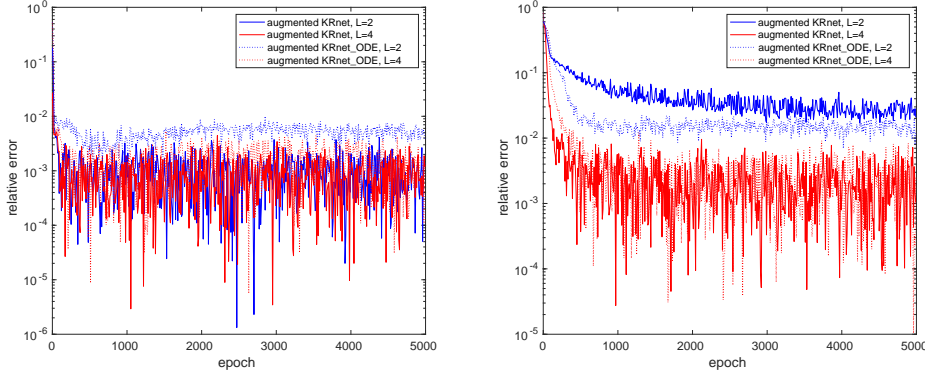


Fig. 4.1: Compare the convergence behavior of KRnet_aug and KRnet_ODE.

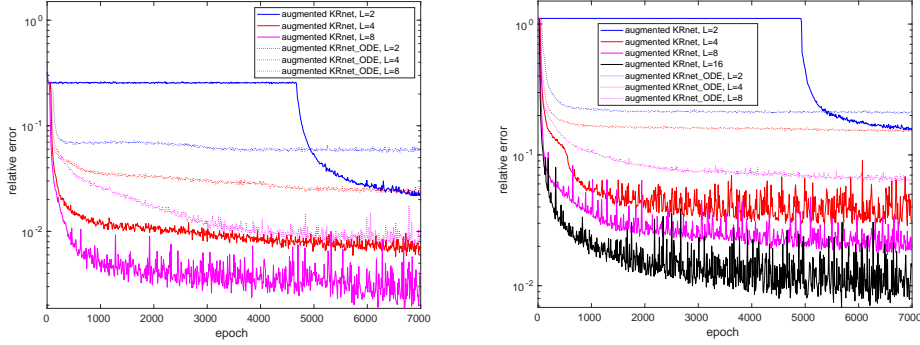


Fig. 4.2: Compare the convergence behavior of KRnet_aug and KRnet_ODE.

figure 4.4. It is seen that the transition of KRnet at discontinuities is much sharper than that of KRnet_ODE.

4.2. Two-dimensional mixture of Gaussians. We consider a mixture of Gaussians

$$p_{\mathbf{Y}}(\mathbf{y}) = \frac{1}{6} \sum_{i=1}^6 \mathcal{N}(\mathbf{y}_i, \mathbf{I}), \quad (4.11)$$

where $\mathbf{y}_i = (5 \cos \frac{i\pi}{3}, 5 \sin \frac{i\pi}{3})$. We have six standard Gaussians uniformly located on a circle of radius 5. We examine and compare the following modeling techniques:

- KRnet: This KRnet only keeps the triangular structure inspired by the K-R rearrangement. For two-dimensional problems, KRnet is consistent with the real NVP .
- KRnet_aug: One augmented dimension is added to KRnet.
- KRnet_R&N: The rotation layers and the nonlinear invertible layer are switched on for KRnet.
- KRnet_aug_R&N: The rotation layers and the nonlinear invertible layer are

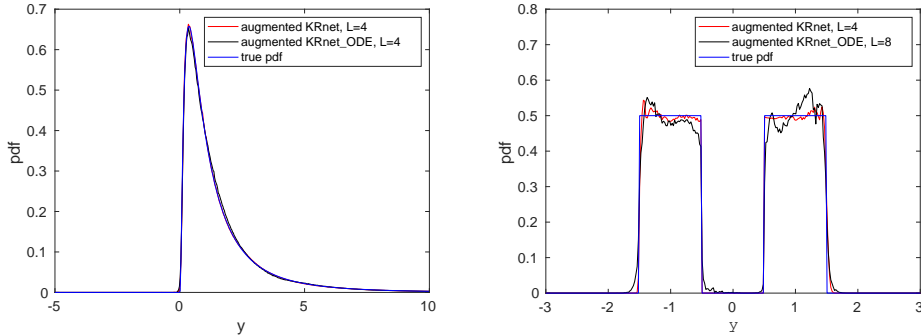


Fig. 4.3: The approximated PDFs for the lognormal distribution and the uniform distribution with a hole.

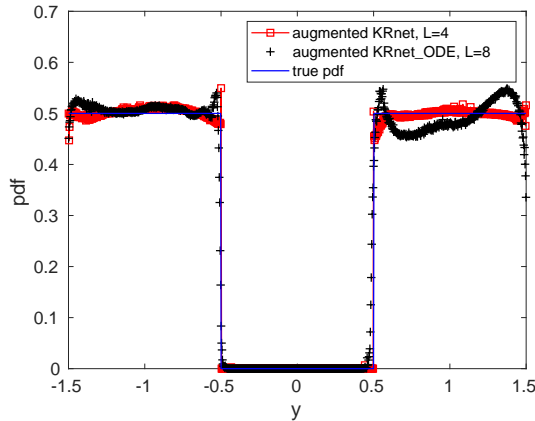


Fig. 4.4: The approximated PDFs for the lognormal distribution and the uniform distribution with a hole.

switched on for `KRnet_aug`, where the rotation only acts on the data dimensions and does not affect the augmented dimension.

- `KRnet_ODE`: This is the neural ODE model based on the `KRnet`.

For the numerical experiments, we obtain 6.4×10^5 samples from the mixture of Gaussians for the training set. We minimize the cross entropy between the empirical distribution and the model using 8 minibatches. The error is defined as the relative difference between the cross entropy and the differential entropy of the mixture of Gaussians, see equation (4.9), which can be regarded as the KL divergence between the model and the data distribution since the sample size is relatively large. For the `KRnet_ODE`, the ODE is discretized on $[0, 1]$ with a step size 0.05.

All models have been trained using the same training set. For each model, we implement the training process ten times and define the mean of the ten errors as the final error. This way the bias from random initialization is reduced. For each training process, we run up to 8000 epochs. The results have been summarized in

Table 4.1: Errors of some KRnet-based models for the density estimation of samples from the mixture of Gaussians (4.11). In affine coupling layers, the neural network (2.6) has two dense hidden layers, each of which has 24 neurons, and this number decays at a ratio $r = 0.9$ in terms of the index k of the outer loop of KRnet.

	KRnet	KRnet_aug	KRnet_aug_R&N	KRnet_R&N	KRnet_ODE
$L = 2$:	6.96e-2	1.02e-1	4.52e-2	1.50e-2	2.93e-2
$L = 4$:	1.74e-2	8.47e-3	1.29e-3	2.56e-3	1.67e-2
$L = 6$:	5.46e-3	1.53e-3	6.79e-4	1.56e-3	1.02e-2

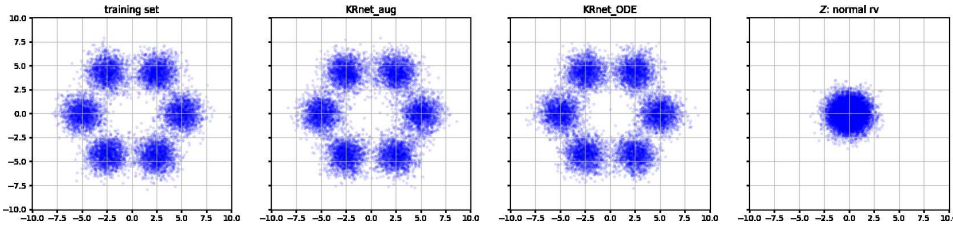


Fig. 4.5: Data distributions given by the training set, the augmented KRnet, and the KRnet-based neural ODE, where the prior Gaussian distribution has been mapped to the mixture of Gaussians (4.11). $L = 6$.

table 4.1. First of all, for each model the error decays as the number of affine coupling layers increases. Second, the KRnet_ODE demonstrates a better performance than KRnet when L is small, and is outperformed by KRnet when L is large. However, KRnet_ODE is significantly slower than KRnet. Third, the model KRnet_aug_R&N yields the best performance, implying that the dimensional augmentation, the rotation layer and the nonlinear invertible layer are effective. When $L = 2$, KRnet_aug performs the worst. This is reasonable since the number of dimensions is increased by one. However, the KRnet_aug has a fast decay in error.

In figure 4.5, we compare the data distributions given by the training set, KRnet_aug, and KRnet_ODE for the case $L = 6$ in Table 4.1. Both KRnet_aug and KRnet_ODE produce a distribution that is visually the same as the data distribution given by the training set.

We next consider the density approximation. We use KRnet_aug_R&N to approximate the PDF (4.11) by minimizing the KL divergence (3.49). For this case, there does not exist a training set. The samples for the approximation of the KL divergence are from the model KRnet_aug_R&N. Since every minibatch can be independently sampled from the model, the optimization solver can be regarded as a minibatch stochastic gradient method with a training set of infinitely many data. For KRnet_aug_R&N, we use $L = 6$, and the rest of the configuration is the same as before. The size of minibatch is 10^5 . In figure 4.6, we plot the convergence behavior of KRnet_aug_R&N, and in figure 4.7, we compare the samples from the true PDF and the approximated PDF. It is seen that the augmented KRnet is also effective for density approximation.

4.3. Logistic distribution with holes. The training data sets $\mathcal{S} = \{\mathbf{y}^{(i)}\}_{i=1}^{N_t}$ for density estimation are generated as follows. Assume that \mathbf{Y} has i.i.d. components

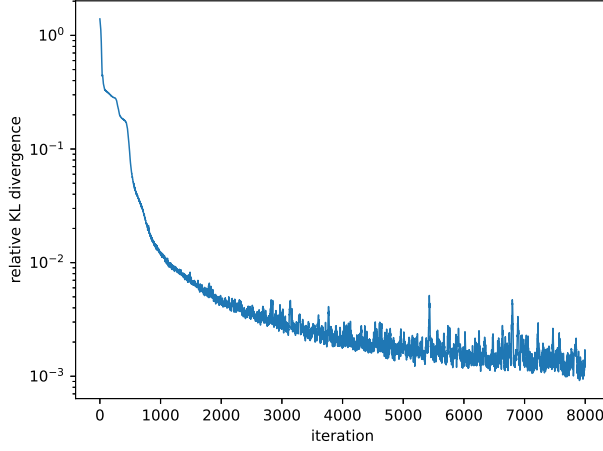


Fig. 4.6: The convergence behavior of KRnet_aug_R&N for the approximation of the 2d mixture of Gaussians (4.11).

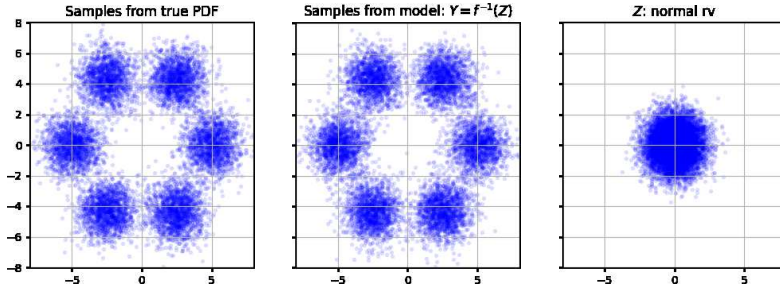


Fig. 4.7: Compare the data distributions from the 2d mixture of Gaussians (4.11) and the approximated PDF given by KRnet_aug_R&N. The sample size is $N = 10000$.

and each component $Y_i \sim \text{Logistic}(0, s)$ with PDF $\rho(y_i; 0, s)$. We propose the following constraint

$$\|R_{\gamma, \theta_j} [y_j^{(i)}, y_{j+1}^{(i)}]^\top\|_2 \geq C, \quad j = 1, \dots, d-1, \quad (4.12)$$

where C is a specified constant, and

$$R_{\gamma, \theta_j} = \begin{bmatrix} \gamma & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \cos \theta_j & -\sin \theta_j \\ \sin \theta_j & \cos \theta_j \end{bmatrix}, \quad \theta_j = \frac{\pi}{4}, \text{ if } j \text{ is even}; \frac{3\pi}{4}, \text{ otherwise.}$$

We then generate samples $\mathbf{y}^{(i)}$ of \mathbf{Y} , out of which we only accept those that satisfy the constraint (4.12). This way, an elliptic hole is generated for two adjacent dimensions. The reference PDF takes the form

$$p_{\mathbf{Y}, \text{ref}}(\mathbf{y}) = \frac{I_B(\mathbf{y}) \prod_{i=1}^n \rho(y_i; 0, s)}{\mathbb{E}[I_B(\mathbf{Y})]}, \quad (4.13)$$

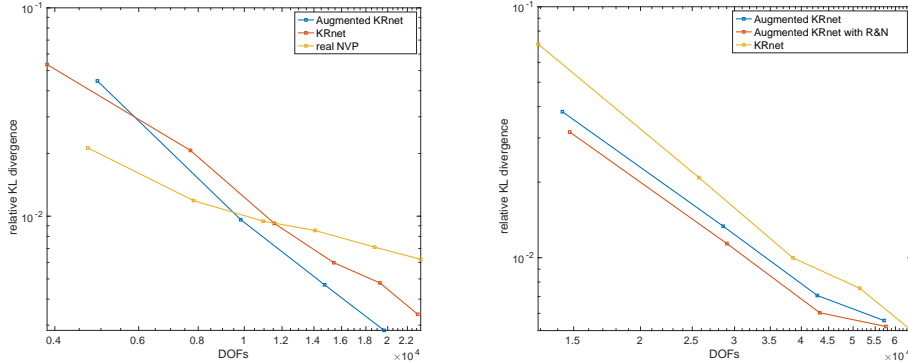


Fig. 4.8: Compare the convergence behavior of augmented KRnet, regular KRnet and real NVP. Left: $n = 4$; Right: $n = 8$.

where B is the set defined by equation (4.12) and $I_B(\cdot)$ is an indicator function with $I_B(\mathbf{y}) = 1$ if $\mathbf{y} \in B$; 0, otherwise.

For this test problem, we set $\gamma = 3$ and $C = 7.6$. This case was studied in [27] and we use the same setup here. The size of the training set is $N_t = 9.6 \times 10^5$ and the errors are computed in terms of a validation set of size 3.2×10^5 . For each model configuration, we train the model 10 times respectively in terms of 10 independently sampled training sets. We then use the averaged error to reduce the bias. The neural network for the affine coupling layers has two dense hidden layers of 24 neurons for $n = 4$ and of 32 neurons for $n = 8$. The comparison of the augmented KRnet, the regular KRnet and the real NVP is summarized in figure 4.8 in terms of DOFs, where the relative Kullback-Leibler (KL) divergence is defined as

$$\frac{D_{\text{KL}}(p_{\mathbf{Y},\text{ref}} \| p_{\mathbf{Y}})}{h(p_{\mathbf{Y},\text{ref}})},$$

where $D_{\text{KL}}(p_{\mathbf{Y},\text{ref}} \| p_{\mathbf{Y}})$ is approximated by the validation set. It is seen that both the augmented KRnet and the regular KRnet yield a much better trend in terms of the convergence rate than the real NVP. The augmented KRnet and the regular KRnet have similar convergence behavior while the augmented KRnet is more effective than the regular KRnet for the same number of DOFs, which is verified by the simulation results for both $d = 4$ and $d = 8$. On the right plot, we also include the results for the augmented KRnet with rotation and nonlinear layers. With a slightly larger number of DOFs, the rotation and nonlinear layers further improve the performance of the augmented KRnet. It is seen that rotation and nonlinear layers do not improve the augmented KRnet for the last case. It is because a constant error has been reached since both the loss and the generalization error have been approximated by the Monte Carlo method.

In figure 4.9, we have compared the samples generated by some generative models to the training set. The data distribution to be learned is highly irregular. On any face given by two adjacent dimensions, a cylinder hole exists. On the boundary of this cylinder, the density can be large. The existence of sharp discontinuities in density implies that classical PDF models such as the mixture of Gaussians are not effective. However, the deep generative models can deal with this high-dimensional density

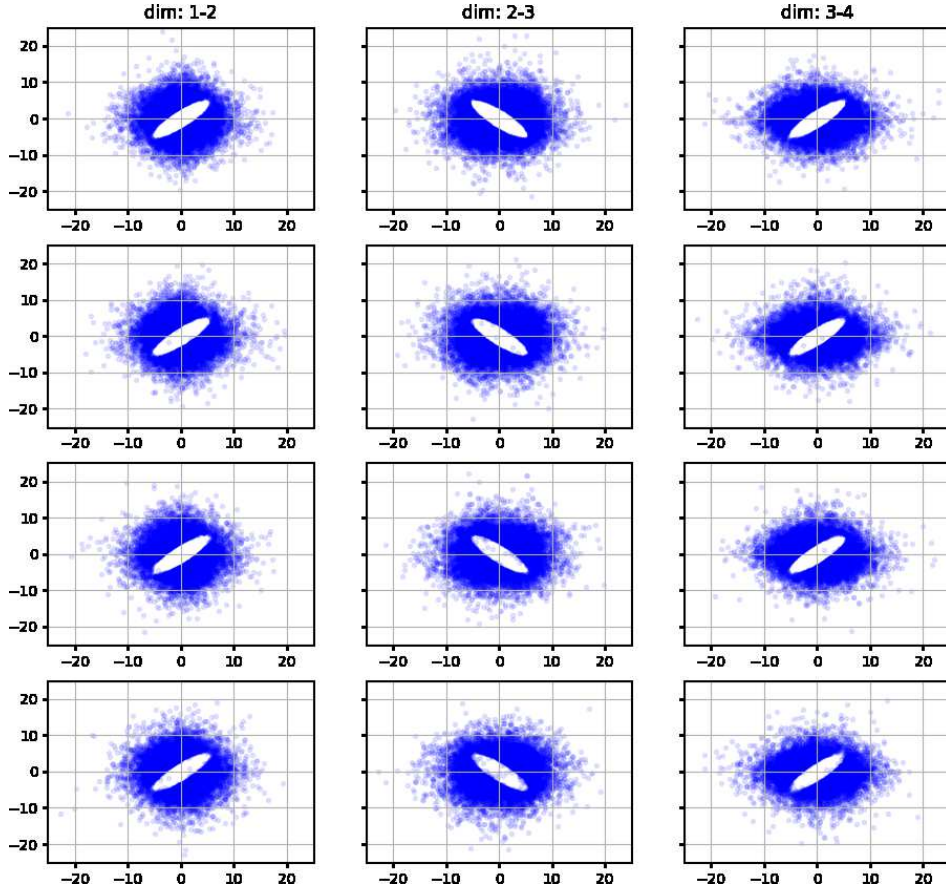


Fig. 4.9: Compare the samples from training set, the augmented KRnet, the regular KRnet and the real NVP for about the same number of DOFs. The three models correspond to the three cases in the left plot of figure 4.8 with DOFs about $2e4$. From top to bottom, each row shows three groups of adjacent dimensions $((y_1, y_2), (y_2, y_3), (y_3, y_4))$ for the data from the training set, the augmented KRnet, the regular KRnet and the real NVP, respectively. Each set has 10000 samples.

estimation problem quite well. Roughly speaking, we may tell the improvement from the real NVP to the augmented KRnet by the decreasing number of outliers in the hole, where the density is supposed to be zero. Out of 10000 samples only a few show up in the holes meaning that boundaries of the holes have been well captured.

In figure 4.10 we plot the distribution of samples generated by the augmented KRnet for the 8-dimensional Logistic distribution with elliptic holes. It is seen that for such a high-dimensional irregular distribution the sharp discontinuities in density can also be well resolved.

5. Summary and discussions. In this work we have developed augmented KRnet for both discrete and continuous models. The main idea is to introduce augmented dimensions to enhance the exchange of information between data dimensions

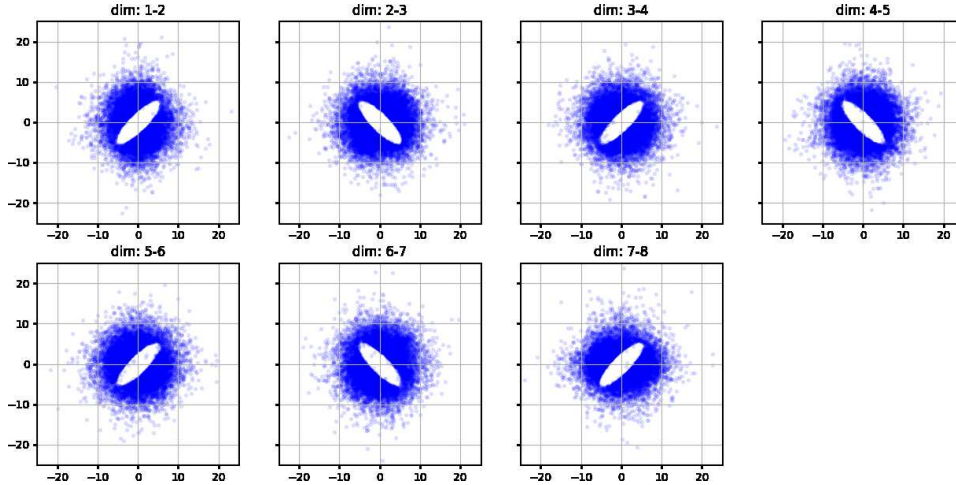


Fig. 4.10: Samples generated by augmented KRnet for the 8-dimensional Logistic distributions with elliptic holes. The sample size is $N = 10000$. The model corresponds to the third case of KRnet_aug_R&N in the right plot of figure 4.8, where the number of DOFs is about $4.3e4$.

such that the flow-based generative model induced by KRnet may further increase its modeling capability while maintaining the exact invertibility of the transport map. We have also formulated the augmented KRnet as a discretization of a neural ODE by a one-step method of first-order accuracy, where the exact invertibility has been kept locally. Although we are not able to discretize the neural ODE with a high-order numerical scheme, a dynamical model with a first-order invertible discretization is still of particular interest for the modeling of dynamical data since the gradient can be exactly computed. A number of numerical experiments have been implemented. Both discrete and continuous models based on the augmented KRnet are effective for both density estimation and approximation, where the algebraic convergence is observed as the number of DOFs increases. In particular, the augmented KRnets are able to deal with high-dimensional distributions that have sharp discontinuous boundaries. Based on these observations, we think that the augmented KRnet may serve as a generic PDF model for many applications. At this moment, our numerical experiments show that the discrete models are in general more effective and much faster than the continuous models. Further research is needed to improve the efficiency of KRnet_ODE.

Acknowledgment. This work was supported by NSF grant DMS-1913163.

REFERENCES

- [1] M. Arjovsky, S. Chintala, and L. Bottou, *Wasserstein GAN*, (2017), arXiv:1701.07875v3.
- [2] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe, *Variational inference: A review for statisticians*, (2018), arXiv:1601.00670v9.
- [3] R. van den Berg, L. Hasenclever, J. M. Tomczak and M. Welling, *Sylvester normalizing flows for variational inference*, (2019), arXiv:1601.00670v9.
- [4] G. Carlier, A. Galichon, and F. Santambrogio, *From Knothe's transport to Brenier's map*

- and a continuation method for optimal transport, *SIAM J. Math. Anal.*, 41(6) (2010), pp. 2554–2576.
- [5] R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, *Neural ordinary differential equations*, (2019), arXiv:1806.07366v5.
- [6] L. Dinh, D. Krueger, and S. Bengio, *Nice: non-linear independent components estimation*, (2014), arXiv:1410.8516.
- [7] L. Dinh, J. Sohl-Dickstein, and S. Bengio, *Density estimation using real NVP*, (2017), arXiv:1605.08803v3.
- [8] E. Dupont, A. Doucet, and Y. W. Teh, *Augmented neural ODEs*, (2019), arXiv:1904.01681v3.
- [9] F. Santambrogio, *Optimal Transport for Applied Mathematicians*, Birkhäuser, 2010.
- [10] C. Finlay, J.-H. Jacobsen, L. Nurbekyan, and A. M. Oberman *How to train your neural ODE: the world of Jacobian and kinetic regularization*, (2020), arXiv:2002.02798v3.
- [11] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, *Generative adversarial nets*, *Advances in Neural Information Processing Systems*, (2014), 2672–2680.
- [12] A. Graves, *Generating sequences with recurrent neural networks*, (2013), arXiv:1308.0850.
- [13] A. Grover, M. Dhar, and S. Ermon, *Flow-GAN: Combining maximum likelihood and adversarial learning in generative models*, (2018), arXiv:1705.08868v2.
- [14] K. He, X. Zhang, S. Ren, and J. Sun, *Deep residual learning for image recognition*, (2015), arXiv:1512.03385v1.
- [15] S. Ioffe, and C. Szegedy, *Batch normalization: Accelerating deep network training by reducing internal covariance shift*, (2015), arXiv:1502.03167v3.
- [16] D. P. Kingma, and M. Welling, *Auto-encoding variational Bayes*, (2014), arXiv:1312.6114v10.
- [17] D. P. Kingma, and P. Dhariwal, *Glow: Generative flow with invertible 1x1 convolutions*, (2018), arXiv:1807.03039v2.
- [18] D. P. Kingma, and J. L. Ba, *ADAM: A method for stochastic optimization*, (2017), arXiv:1412.6980v9.
- [19] D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, and M. Welling, *Improving variational inference with inverse autoregressive flow*, *Advances in Neural Information Processing Systems*, (2016), pp. 4743–4751.
- [20] Y. Lu, A. Zhang, Q. Li, and B. Dong, *Beyond finite layer neural networks: bridging deep architectures and numerical differential equations*, (2020), arXiv:1710.10121v3.
- [21] A. van den Oord, N. Kalchbrenner, and K. Kavukcuoglu, *Pixel recurrent neural networks*, (2016), arXiv:1601.06759.
- [22] A. van den Oord, N. Kalchbrenner, O. Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu, *Conditional image generation with PixelCNN decoders*, (2016), arXiv:1606.05328.
- [23] G. Papamakarios, T. Pavlakou, and I. Murray, *Masked autoregressive flow for density estimation*, (2018), arXiv:1705.07057v4.
- [24] D. Rezende, and S. Mohamed, *Variational inference with normalizing flows*, *ICML*, (2015), 1530–1538.
- [25] D. Scott, *Multivariate Density Estimation: Theory, Practice, and Visualization*, 2nd Edition, John Wiley & Sons, Inc., 2015.
- [26] A. Spatini, D. Bigoni, and Y. Marzouk, *Inference via low-dimensional couplings*, (2017), arXiv:1703.06131v4.
- [27] K. Tang, X. Wan, and Q. Liao, *Deep density estimation via invertible block-triangular mapping*, *Theoretical & Applied Mechanics Letters*, **10**, 2020, 000-5.
- [28] K. Tang, X. Wan, and Q. Liao, *Adaptive deep density approximation for Fokker-Planck equations*, (2021), arXiv:2013.11181v1.
- [29] X. Wan, and S. Wei, *Coupling the reduced-order model and the generative model for an importance sampling estimator*, *J. Compt. Phys.*, in press.
- [30] L. Yang, and G. E. Karniadakis, *Potential flow generator with L_2 optimal transport regularity for generative models*, (2019), arXiv:1908.11462v1.
- [31] L. Zhang, W. E, and L. Wang, *Monge-Ampère flow for generative modeling*, (2018), arXiv:1809.10188v1.
- [32] J. Zhu, D. Zhao, and B. Zhang, *LIA: Latently Invertible Autoencoder with Adversarial Learning*, (2019), arXiv:1906.08090v1.