



# Adaptive deep density approximation for Fokker-Planck equations



Kejun Tang<sup>a,b</sup>, Xiaoliang Wan<sup>c</sup>, Qifeng Liao<sup>a,\*</sup>

<sup>a</sup> School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China

<sup>b</sup> Peng Cheng Laboratory, Shenzhen 518055, China

<sup>c</sup> Department of Mathematics and Center for Computation and Technology, Louisiana State University, Baton Rouge 70803, USA

## ARTICLE INFO

### Article history:

Received 19 March 2021

Received in revised form 14 December 2021

Accepted 12 February 2022

Available online 22 February 2022

### Keywords:

Density estimation

Flow-based generative models

Fokker-Planck equations

Deep learning

## ABSTRACT

In this paper we present an adaptive deep density approximation strategy based on KRnet (ADDA-KR) for solving the steady-state Fokker-Planck (F-P) equations. F-P equations are usually high-dimensional and defined on an unbounded domain, which limits the application of traditional grid based numerical methods. With the Knothe-Rosenblatt rearrangement, our newly proposed flow-based generative model, called KRnet, provides a family of probability density functions to serve as effective solution candidates for the Fokker-Planck equations, which has a weaker dependence on dimensionality than traditional computational approaches and can efficiently estimate general high-dimensional density functions. To obtain effective stochastic collocation points for the approximation of the F-P equation, we develop an adaptive sampling procedure, where samples are generated iteratively using the approximate density function at each iteration. We present a general framework of ADDA-KR, validate its accuracy and demonstrate its efficiency with numerical experiments.

© 2022 Elsevier Inc. All rights reserved.

## 1. Introduction

During the past few decades there has been a rapid development in numerical methods for Fokker-Planck equations. This explosion in interest has been driven by the need of assessing time evolution of probability density functions in randomly perturbed dynamical systems, which are widely used in physical and biological modeling [1–3]. It is known that there exist two main challenges for efficiently solving the Fokker-Planck equations: the spatial variable can be high-dimensional, which causes difficulties in applying grid based numerical methods, e.g. finite element methods [4,5]; the original spatial domain is typically unbounded, and it is challenging to derive a well-posed boundary condition for a bounded computational domain. To alleviate these difficulties, new numerical methods based on deep learning currently gain a lot of attention [6,7], and this paper is devoted to deep learning for the Fokker-Planck equations.

Deep learning methods for partial differential equations (PDEs) are under active development. In [8,9], a deep Ritz method is proposed based on variational methods. In [10–13], physics-informed neural networks are developed through infusing PDEs into networks as a constraint. A deep Galerkin method is proposed in [14]. Bayesian deep convolutional encoder-decoder networks for PDEs with high-dimensional random inputs are developed in [15,16]. Deep learning strategies are also introduced to discover physical laws [17,18]. In addition, efficient deep learning methods based on domain

\* Corresponding author.

E-mail addresses: tangkj@shanghaitech.edu.cn (K. Tang), xlwan@lsu.edu (X. Wan), liaoqf@shanghaitech.edu.cn (Q. Liao).

decomposition are studied in [19–24], and deep neural network methods for complex geometries and irregular domains are proposed in [25,26]. The main idea of deep learning methods for PDEs is to reformulate a PDE problem as an optimization problem and train deep neural networks through minimizing the corresponding loss functional. In these methods, stochastic collocation points are required to estimate the loss functional. We note that the stochastic collocation points herein are for the spatial variable, while stochastic collocation for PDEs with random inputs (especially for parameters) are discussed in detail in [27–33]. To result in an efficient deep learning strategy for PDEs, properly choosing the collocation points is crucial. Intuitively, the distribution of the collocation points should be consistent with the properties of the PDE solution in a certain sense. In our recent work [19], a hierarchical sampling procedure is proposed based on domain decomposition iterations, while it focuses on low-dimensional problems. As the spatial variable of the Fokker-Planck equation can be high-dimensional, it remains an open challenging problem to generate effective collocation points. We develop an effective adaptive sampling procedure to alleviate this issue in this work. Adaptivity is widely used in machine learning techniques to make the training process more effective by exploring the relation between the model and the data, e.g., active learning selects the most helpful samples to increase efficiency [34,35] and meta-learning tries to match learning algorithms with task properties [36]. In our problem, we will update the training set partially or completely according to the learned model, i.e., the approximate solution of the F-P equation, and the updated training set will yield a better approximate solution.

As the solution of the Fokker-Planck equation is a probability density function, solving this problem can also be considered as a density estimation problem. It is known that density estimation is a central topic in unsupervised learning, and it still remains an open challenge for high-dimensional density estimation [37]. Recently, two kinds of deep learning models have shown great promise for estimating high-dimensional probability density functions (PDFs), which include the flow-based generative model [38,39] and the neural ordinary differential equation model [40,41]. In this work, we focus on the flow-based generative model, which is to construct invertible mappings from a prescribed prior distribution to the empirical distribution given by data and build explicit probability density functions using the change of variables. The Knothe-Rosenblatt (KR) rearrangement [42] shows that such an invertible mapping can be achieved with a triangular structure. Incorporating with the KR rearrangement, we propose an invertible block-triangular mapping, called KRnet, which generalizes the flow-based generative model given by real NVP [38]. We note that there are a lot of generative models which can efficiently generate samples of the distributions under consideration but do not explicitly give the corresponding density functions, e.g., generative adversarial networks (GANs) [43] and the variational autoencoder (VAE) [44]. In addition, coupling flow-based generative models and reduced-order models into an importance sampling estimator is studied in [45].

In this work, we propose an adaptive deep density approximation method based on KRnet (ADDA-KR) for solving Fokker-Planck equations. We first provide additional details and results for KRnet that was outlined in the letter [46]. After that, we use KRnet to construct solutions of the Fokker-Planck equation. Since KRnet can induce a family of probability density functions, normality and vanishing boundary conditions are satisfied naturally. Like other deep learning algorithms for solving PDEs, our method is also meshfree. The PDE problem is converted into an optimization problem and it can be solved through stochastic gradient descent on a set of collocation points, while traditional grid-based numerical methods (e.g. finite element methods) rapidly become computationally infeasible since the number of grid points grows exponentially with the dimensionality. The choice of the collocation points plays a crucial role in a meshless method. The distribution of the collocation points should be consistent with the regularity of the solution for both accuracy and efficiency. Since the solution of the F-P equation is a probability density function, one way to achieve this is to use the samples of the solution PDF as the collocation points. Based on such an idea, we propose an adaptive approach ADDA-KR that has two main steps: training a KRnet to approximate the solution of the Fokker-Planck equation, and using the trained KRnet to generate collocation points for the next iteration. After each iteration, the distribution of the collocation points is more consistent with the solution PDF.

The rest of the paper is organized as follows. In the next section, the Fokker-Planck equations and the problem setting are introduced. Our KRnet is presented in section 3. In section 4, our novel adaptive deep density approximation approach for the Fokker-Planck equation is presented. In section 5, we demonstrate the efficiency of our adaptive sampling approach with numerical experiments. Finally section 6 concludes the paper.

## 2. Problem setup

Consider the state  $X_t$  modeled by the following stochastic differential equation

$$dX_t = \boldsymbol{\mu}(X_t, t)dt + \mathbf{G}(X_t, t)d\mathbf{w}_t, \quad (1)$$

where  $\boldsymbol{\mu} = [\mu_1, \dots, \mu_d]^\top$  is a vector field,  $\mathbf{G}(X_t, t) \in \mathbb{R}^{d \times d}$  is a matrix-valued function and  $\mathbf{w}_t$  is a  $d$ -dimensional standard Wiener process. The Fokker-Planck equation, which describes the probability density function of  $X_t$ , is

$$\begin{aligned} \frac{\partial p(\mathbf{x}, t)}{\partial t} &= \mathcal{L}p(\mathbf{x}, t) := \nabla \cdot [p(\mathbf{x}, t)\nabla V(\mathbf{x}, t)] + \nabla \cdot [\nabla \cdot (p(\mathbf{x}, t)\mathbf{D}(\mathbf{x}, t))], \quad \forall(\mathbf{x}, t) \in \mathbb{R}^d \times \mathbb{R}^+, \\ \int_{\mathbb{R}^d} p(\mathbf{x}, t)d\mathbf{x} &= 1, \quad p(\mathbf{x}, t) \geq 0, \quad \forall(\mathbf{x}, t) \in \mathbb{R}^d \times \mathbb{R}^+, \\ p(\mathbf{x}, 0) &= p_0(\mathbf{x}), \end{aligned} \quad (2)$$

where  $\mathbf{x} \in \mathbb{R}^d$  denotes a random vector,  $V(\mathbf{x}, t)$  is a potential function,  $\mathbf{D}(\mathbf{x}, t)$  is a diffusion matrix,  $p(\mathbf{x}, t)$  is the unknown probability density function (PDF) of  $\mathbf{x}$  with the initial PDF  $p_0(\mathbf{x})$ , and  $\mathcal{L}$  denotes the partial differential operator. Following [1], the potential function  $V(\mathbf{x}, t)$  and the diffusion matrix  $\mathbf{D}(\mathbf{x}, t)$  can be expressed as

$$\begin{aligned}\nabla V(\mathbf{x}, t) &= -\boldsymbol{\mu}(\mathbf{x}, t), \\ \mathbf{D}(\mathbf{x}, t) &= \frac{1}{2} \mathbf{G}(\mathbf{x}, t) \mathbf{G}(\mathbf{x}, t)^\top.\end{aligned}$$

In this work, we focus on the stationary solution of Eq. (2), i.e., the invariant measure independent of time,

$$\mathcal{L}p(\mathbf{x}) = \nabla \cdot [p(\mathbf{x}) \nabla V(\mathbf{x})] + \nabla \cdot [p(\mathbf{x}) \mathbf{D}(\mathbf{x})] = 0, \quad (3)$$

with the boundary condition

$$p(\mathbf{x}) \rightarrow 0 \quad \text{as} \quad \|\mathbf{x}\|_2 \rightarrow \infty, \quad (4)$$

and some extra constraints on  $p(\mathbf{x})$

$$\int_{\mathbb{R}^d} p(\mathbf{x}) d\mathbf{x} = 1, \quad \text{and} \quad p(\mathbf{x}) \geq 0, \quad (5)$$

where  $\|\mathbf{x}\|_2$  indicates the  $\ell_2$  norm of  $\mathbf{x}$ .

There are several difficulties for the approximation of equation (3). First, the boundary condition and the constraints of  $p(\mathbf{x})$  may not be easily satisfied when we employ the traditional approaches such as the finite element method. Since the support of  $p(\mathbf{x})$  is  $\mathbb{R}^d$ , the computation domain has to be truncated, implying that the boundary condition must be approximated, e.g., a homogeneous boundary condition. To preserve the nonnegativity of  $p(\mathbf{x})$ , a projection step is needed for the box constraint. Second, it requires a fine mesh to capture the whole information when the target density is multimodal, i.e., the potential function  $V(\mathbf{x})$  has many local minima [47], which is computationally infeasible when the dimension  $d$  is even moderately large. We also note that a homogeneous boundary condition usually requires a large computational domain, which makes a uniform refinement even more challenging, if no prior information can be used for certain adaptivity on mesh generation. To address these issues, we will propose an adaptive deep density approximation method to solve the Fokker-Planck equation (3) using a deep generative model for  $p(\mathbf{x})$ . The flow-based generative model not only provides an explicit density function that satisfies naturally all constraints on  $p(\mathbf{x})$ , but also suggests a simple but effective adaptive strategy for the approximation of equation (3) through sampling the current approximation of  $p(\mathbf{x})$ .

### 3. KRnet

KRnet is a flow-based generative model for density estimation or approximation. In this section we briefly overview KRnet that has been outlined in our recently published letter [46] and present more details that were not included in [46] due to the page limit. Let  $X \in \mathbb{R}^d$  be a random vector associated with a given data set, and its probability density function (PDF) is denoted by  $p_X(\mathbf{x})$ . The target is to estimate  $p_X(\mathbf{x})$  using available data. Let  $Z \in \mathbb{R}^d$  be a random vector associated with a PDF  $p_Z(\mathbf{z})$ , where  $p_Z(\mathbf{z})$  is a prior distribution (e.g., Gaussian distribution). The flow-based generative modeling is to seek an invertible mapping  $\mathbf{z} = f(\mathbf{x})$  where  $f(\cdot)$  is a bijection:  $f: \mathbf{x} \mapsto \mathbf{z}$  [38]. By the change of variables, we have the PDF of  $X = f^{-1}(Z)$  as

$$p_X(\mathbf{x}) = p_Z(f(\mathbf{x})) |\det \nabla_{\mathbf{x}} f|. \quad (6)$$

Once the prior distribution  $p_Z(\mathbf{z})$  is specified, equation (6) provides an explicit PDF of  $X$ . Given a set of training data, the invertible mapping  $f(\cdot)$  can be learned by maximizing the likelihood or minimizing the cross entropy. The inverse of  $f(\cdot)$  provides a convenient way to sample  $X$  as  $X = f^{-1}(Z)$ .

#### 3.1. A new affine coupling layer

In flow-based generative models, the invertible mapping  $f(\cdot)$  is constructed by stacking a sequence of simple bijections, each of which is a shallow neural network, and thus the overall mapping is a deep net. The mapping  $f(\cdot)$  can be written in a composite form:

$$\mathbf{z} = f(\mathbf{x}) = f_{[L]} \circ \dots \circ f_{[1]}(\mathbf{x}) \quad \text{and} \quad \mathbf{x} = f^{-1}(\mathbf{z}) = f_{[1]}^{-1} \circ \dots \circ f_{[L]}^{-1}(\mathbf{z}), \quad (7)$$

where  $f_{[i]}$  is called an affine coupling layer at stage  $i$ . The Jacobian matrix can be obtained by the chain rule

$$|\det \nabla_{\mathbf{x}} f| = \prod_{i=1}^L |\det \nabla_{\mathbf{x}_{[i-1]}} f_{[i]}|, \quad (8)$$

where  $\mathbf{x}_{[i-1]}$  indicate the intermediate variables with  $\mathbf{x}_{[0]} = \mathbf{x}$  and  $\mathbf{x}_{[L]} = \mathbf{z}$ . Let  $\mathbf{x}_{[i]} = [\mathbf{x}_{[i],1}, \mathbf{x}_{[i],2}]^T$  be a partition of  $\mathbf{x}_{[i]}$  with  $\mathbf{x}_{[i],1} \in \mathbb{R}^m$  and  $\mathbf{x}_{[i],2} \in \mathbb{R}^{d-m}$  for  $i = 0, \dots, L-1$ . One technique to define the affine coupling layer is the real NVP [38]:

$$\begin{aligned}\mathbf{x}_{[i],1} &= \mathbf{x}_{[i-1],1} \\ \mathbf{x}_{[i],2} &= \mathbf{x}_{[i-1],2} \odot \exp(\log \mathbf{s}_i(\mathbf{x}_{[i-1],1})) + \mathbf{t}_i(\mathbf{x}_{[i-1],1}),\end{aligned}\quad (9)$$

where  $\mathbf{s}_i : \mathbb{R}^m \mapsto \mathbb{R}^{d-m}$  and  $\mathbf{t}_i : \mathbb{R}^m \mapsto \mathbb{R}^{d-m}$  are the scaling and the translation depending on  $\mathbf{x}_{[i-1],1}$ , and  $\odot$  is the Hadamard product or element-wise product. Note that  $\mathbf{x}_{[i-1],1}$  remains fixed and the modification of  $\mathbf{x}_{[i-1],2}$  is linear with respect to  $\mathbf{x}_{[i-1],2}$  and nonlinear in terms of  $\mathbf{x}_{[i-1],1}$ . This way, the Jacobian matrix  $\nabla_{\mathbf{x}_{[i-1]}} f_{[i]}$  is lower-triangular whose determinant can be evaluated efficiently. Furthermore,  $(\mathbf{s}_i, \mathbf{t}_i)$  is usually modeled by a neural network  $\text{NN}_{[i]}$

$$(\mathbf{s}_i, \mathbf{t}_i) = \text{NN}_{[i]}(\mathbf{x}_{[i-1],1}). \quad (10)$$

We proposed a new affine coupling layer  $f_{[i]}$  as follows [46]

$$\begin{aligned}\mathbf{x}_{[i],1} &= \mathbf{x}_{[i-1],1} \\ \mathbf{x}_{[i],2} &= \mathbf{x}_{[i-1],2} \odot (1 + \alpha \tanh(\mathbf{s}_i(\mathbf{x}_{[i-1],1}))) + e^{\beta_i} \odot \tanh(\mathbf{t}_i(\mathbf{x}_{[i-1],1})),\end{aligned}\quad (11)$$

where  $0 < \alpha < 1$  is a hyperparameter and the parameter  $\beta_i \in \mathbb{R}^{d-m}$  is trainable. Our affine coupling layer keeps the mechanism of the real NVP when updating the data, and it has the following advantages. First, the second equation in Eq. (11) adapts the trick of ResNet [48], where an identity mapping is added to improve the training process. Second, the constant  $\alpha \in (0, 1)$  is introduced to improve numerical stability. It is seen that the range of  $\det \nabla_{\mathbf{x}_{[i-1]}} f_{[i]}$  is  $[(1-\alpha)^{d-m}, (1+\alpha)^{d-m}]$  for our affine coupling layer and  $(0, +\infty)$  for the original real NVP. Our formulation can alleviate the illnesses when the determinant of the Jacobian in the original real NVP occasionally becomes too large or too small. Third, the trainable factor  $e^{\beta_i}$  depends on the whole training set, which helps avoid possible large oscillation in  $\mathbf{t}_i(\mathbf{x}_{[i-1],1})$  such that the number of outliers can be reduced for sample generation [46]. In our numerical experiments, we set  $\alpha = 0.6$  and it works well.

Since the affine coupling layer  $f_{[i]}$  only updates a part of  $\mathbf{x}_{[i-1]}$ , another affine coupling layer is needed for a complete update. In other words, the next affine coupling layer  $f_{[i+1]}$  can be defined as

$$\begin{aligned}\mathbf{x}_{[i+1],1} &= \mathbf{x}_{[i],1} \odot (1 + \alpha \tanh(\mathbf{s}_{i+1}(\mathbf{x}_{[i],2}))) + e^{\beta_{i+1}} \odot \tanh(\mathbf{t}_{i+1}(\mathbf{x}_{[i],2})) \\ \mathbf{x}_{[i+1],2} &= \mathbf{x}_{[i],2},\end{aligned}$$

where the components  $\mathbf{x}_{[i],1}$  are updated and  $\mathbf{x}_{[i],2}$  remains unchanged. From the dynamical point of view, a long chain of affine coupling layers may result in a highly nonlinear transformation of the input. To enhance the performance and efficiency of the mapping  $f(\mathbf{x})$ , we proposed KRnet to address the following questions: 1) How should we partition the vector? 2) How can we increase the modeling capability except for increasing the depth  $L$ ? 3) Can we provide a robust nonlinear bijection at least in a component-wise way?

### 3.2. The overall structure of KRnet

The basic idea of KRnet is to define the structure of  $f(\mathbf{x})$  in terms of the Knothe-Rosenblatt rearrangement. Let  $\mu_Z$  and  $\mu_X$  be the probability measures of two random variables  $X, Z \in \mathbb{R}^d$  respectively. A mapping  $\mathcal{T} : Z \mapsto X$  is called a transport map such that  $\mathcal{T}_\# \mu_Z = \mu_X$ , where  $\mathcal{T}_\# \mu_Z$  is the push-forward of  $\mu_Z$  such that  $\mu_X(B) = \mu_Z(\mathcal{T}^{-1}(B))$  for every Borel set  $B$  [42]. The Knothe-Rosenblatt rearrangement tells us that the transport map  $\mathcal{T}$  may have a lower-triangular structure

$$\mathbf{z} = \mathcal{T}^{-1}(\mathbf{x}) = \begin{bmatrix} \mathcal{T}_1(x_1) \\ \mathcal{T}_2(x_1, x_2) \\ \vdots \\ \mathcal{T}_d(x_1, \dots, x_d) \end{bmatrix}. \quad (12)$$

This mapping can be regarded as a limit of sequence of optimal transport maps when the quadratic cost degenerates [42]. Noticing that the invertible mapping  $f(\mathbf{x})$  also defines a transport map, we then incorporate the triangular structure of the Knothe-Rosenblatt rearrangement into the definition of  $f(\mathbf{x})$  which results in KRnet as a generalization of real NVP [38]. Let  $\mathbf{x} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}]^T$  be a partition of  $\mathbf{x}$ , where  $\mathbf{x}^{(i)} = [x_1^{(i)}, \dots, x_m^{(i)}]^T$  with  $1 \leq K \leq d, 1 \leq m \leq d$ , and  $\sum_{i=1}^K \dim(\mathbf{x}^{(i)}) = d$ . Our KRnet takes an overall form

$$\mathbf{z} = f_{\text{KR}}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}^{(1)}) \\ f_2(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \\ \vdots \\ f_K(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}) \end{bmatrix}, \quad (13)$$

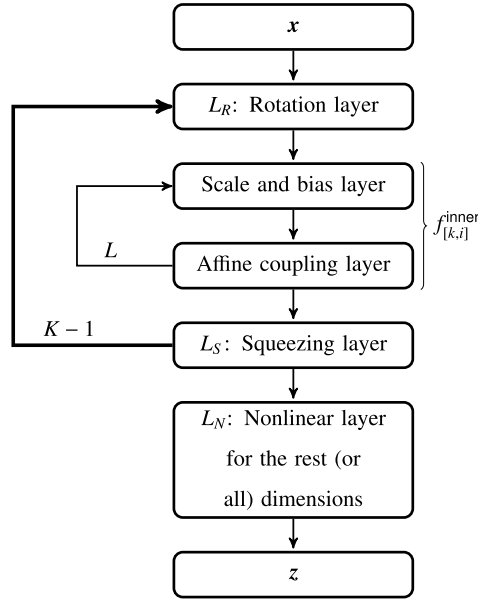


Fig. 1. The flow chart of KRnet.

where each  $f_i$  is an invertible mapping defined as in equation (7) for  $i = 2, \dots, K$ . Note that  $f_1$  is not included if  $K = d$  because we need to partition a vector to two parts to define the affine coupling layer. KRnet consists of one outer loop and  $K - 1$  inner loops. The outer loop has  $K - 1$  stages, corresponding to the  $K - 1$  mappings  $f_i$  in equation (13) with  $i = 2, \dots, K$ , and for each stage, an inner loop of  $L$  affine coupling layers is defined. More specifically, we have

$$\mathbf{z} = f_{KR}(\mathbf{x}) = L_N \circ f_{[K-1]}^{outer} \circ \dots \circ f_{[1]}^{outer}(\mathbf{x}), \tag{14}$$

where  $f_{[i]}^{outer}$  is defined as

$$f_{[k]}^{outer} = L_S \circ f_{[k,L]}^{inner} \circ \dots \circ f_{[k,1]}^{inner} \circ L_R. \tag{15}$$

Here  $f_{[k,i]}^{inner}$  indicates a combination of one affine coupling layer and one scale and bias layer, and  $L_N$ ,  $L_S$  and  $L_R$  indicate the nonlinear layer, the squeezing layer and the rotation layer, respectively, which will be briefly overviewed in the next section.

The flow chart of KRnet is illustrated in Fig. 1. Let us look at how the information flows in the KRnet. Each  $\mathbf{x}_{[k]} = [\mathbf{x}_{[k]}^{(1)}, \dots, \mathbf{x}_{[k]}^{(K)}]^T$  has the same partition with  $\mathbf{x}_{[k]} = f_{[k]}^{outer}(\mathbf{x}_{[k-1]})$  with  $\mathbf{x}_{[0]} = \mathbf{x}$ ,  $k = 1, \dots, K - 1$ . At the beginning, a sequence of affine coupling layers in  $f_{[1]}^{outer}$  is applied to the partition  $\mathbf{x}_{[0]} = [\mathbf{x}_{[0]}^{(1:K-1)}, \mathbf{x}_{[0]}^{(K)}]^T$ , where  $\mathbf{x}_{[0]}^{(1:K-1)}$  includes  $\mathbf{x}_{[0]}^{(i)}$ ,  $i = 1, \dots, K - 1$ . From then on, the last partition  $\mathbf{x}_{[k]}^{(K)}$  will remain fixed for  $k > 1$ . For the next iteration  $f_{[2]}^{outer}$ , the partition  $[\mathbf{x}_{[1]}^{(1:K-2)}, \mathbf{x}_{[1]}^{(K-1)}]^T$  will be used with  $\mathbf{x}_{[1]}^{(K)}$  being deactivated. In general, after the stage  $K - i + 1$  of the outer loop, the  $i$ -th partition of  $\mathbf{x}_{[k]}^{(i)}$  will become deactivated, in addition to the dimensions that are deactivated in the previous stages.

### 3.3. Other types of layer used in KRnet

Except for the affine coupling layers, several other types of layers are needed for the definition of KRnet. We briefly overview these layers in this section and provide some details excluded in the letter [46]. Since each  $\mathbf{x}_{[k]}$  has the same partition, we will drop the subscript for simplicity.

*Squeezing layer*  $L_S$  is used to deactivate some dimensions using a mask

$$\mathbf{q} = [\underbrace{1, \dots, 1}_n, \underbrace{0, \dots, 0}_{d-n}]^T, \tag{16}$$

where the components  $\mathbf{q} \odot \mathbf{x}$  will keep being updated and the rest components  $(1 - \mathbf{q}) \odot \mathbf{x}$  will be fixed from then on.

*Scale and bias layer* provides a simplification of the batch normalization [49], which is defined as

$$\hat{\mathbf{x}} = \mathbf{a} \odot \mathbf{x} + \mathbf{b}, \tag{17}$$

where  $\mathbf{a}$  and  $\mathbf{b}$  are trainable and initialized by the mean and standard deviation of the data. After the initialization,  $\mathbf{a}$  and  $\mathbf{b}$  will be treated as regular trainable parameters that are independent of the data. Numerical experiments show that the scale and bias layer is simple but effective, which provides a comparable performance to the batch normalization layer in our problem setting.

Rotation layer  $L_R$  defines a linear mapping of the input  $\mathbf{x}$

$$\hat{\mathbf{x}} = \hat{\mathbf{W}}\mathbf{x},$$

through a trainable matrix

$$\hat{\mathbf{W}} = \begin{bmatrix} \mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{U} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix},$$

where  $\mathbf{W} \in \mathbb{R}^{n \times n}$ ,  $n$  is defined in the mask  $\mathbf{q}$ ,  $\mathbf{I} \in \mathbb{R}^{(d-n) \times (d-n)}$  is an identity matrix, and  $\mathbf{W} = \mathbf{L}\mathbf{U}$  is the LU factorization of  $\mathbf{W}$ . We expect  $\hat{\mathbf{W}}$  to provide a rotation such that the less important dimensions will be put at the end and deactivated by the next squeezing layer. Entries below the main diagonal of  $\mathbf{L}$  and entries in the upper triangle of  $\mathbf{U}$  are trainable. In practice, we simply optimize the trainable entries of  $\mathbf{L}$  and  $\mathbf{U}$  without enforcing the orthonormality of  $\hat{\mathbf{W}}$  and such a simplification works well.

Nonlinear layer  $L_N$  provides a component-wise nonlinear transformation. For simplicity, we only consider one component  $x$  of the data. We start with a nonlinear mapping  $F(s) : [0, 1] \mapsto [0, 1]$ :

$$F(s) = \int_0^s p(t)dt, \quad (18)$$

where  $p(s)$  is a probability density function. Let  $0 = s_0 < s_1 < \dots < s_{\hat{m}+1} = 1$  be a mesh of the interval  $[0, 1]$  with element size  $h_i = s_{i+1} - s_i$ . Define  $p(s)$  as a piece-wise linear polynomial

$$p(s) = \frac{w_{i+1} - w_i}{h_i}(s - s_i) + w_i, \quad \forall s \in [s_i, s_{i+1}], \quad (19)$$

where

$$p(s_i) = w_i.$$

Then  $F(s)$ , corresponding to a cumulative density function, is a quadratic function

$$F(s) = \frac{w_{i+1} - w_i}{2h_i}(s - s_i)^2 + w_i(s - s_i) + \sum_{k=0}^{i-1} \frac{w_k + w_{k+1}}{2}h_k, \quad \forall s \in [s_i, s_{i+1}], \quad (20)$$

whose inverse and derivative can be explicitly computed.

As the support of each dimension of  $\mathbf{x}$  is  $(-\infty, \infty)$ , a question is how to apply  $F(s)$  to the data. A straightforward strategy is to map  $(-\infty, \infty)$  to  $(0, 1)$  before  $F(s)$  is applied. However, when the inverse is considered, the singularity of mapping a finite interval to an infinite one may introduce issues on robustness. To alleviate this problem, we decompose  $(-\infty, \infty) = (-\infty, -a) \cup [-a, a] \cup (a, \infty)$  with  $a > 0$ , and define the following nonlinear mapping

$$\hat{F}(x) = \begin{cases} \beta_s(x + a) - a, & x \in (-\infty, -a) \\ 2aF\left(\frac{x+a}{2a}\right) - a, & x \in [-a, a] \\ \beta_s(x - a) + a, & x \in (a, \infty), \end{cases} \quad (21)$$

where  $\beta_s > 0$  is a scaling factor. It is seen that we only consider a nonlinear mapping for the data located in  $[-a, a]$  and  $\hat{F}(x)$  maps  $[-a, a]$  to itself. On  $(-\infty, -a) \cup (a, \infty)$ ,  $\hat{F}(x)$  is simply a linear mapping. The reasoning of such a strategy is that the range of data in the training set is always finite, and after being well scaled and shifted the data will be roughly centered at the origin, implying that a nonlinear mapping on  $[-a, a]$  is sufficient as long as  $a$  is large enough. To maintain the invertibility, we require some regularity at  $x = \pm a$ . More specifically,  $\hat{F}'(x)$  should exist at  $x = \pm a$ . Since  $\hat{F}'(x) = \beta_s$  on  $(-\infty, -a) \cup (a, \infty)$ , we have, on  $[-a, a]$ ,  $\hat{F}'(x)|_{x=\pm a} = F'(s)|_{s=0,1} = p(s)|_{s=0,1} = \beta_s$ . So the trainable parameters include  $p(s_i) = w_i$ ,  $i = 1, \dots, \hat{m}$ , subject to the constraint  $\int_0^1 p(s)ds = 1$ .

**Remark 1.** The nonlinear layer is only employed before the final output (see Fig. 1), which can be applied to all dimensions or simply the dimensions that have not been deactivated by the squeezing layer. In both cases, the nonlinear layer enlarges the prescribed prior distribution by a nonlinear component-wise transformation. The parameter  $\beta_s$  acts as an estimate of the density  $p(s)$  at  $s = 0, 1$ . If  $a$  is sufficiently large,  $\beta_s$  can be small accordingly. The prior distribution is often chosen as the standard Gaussian, which means that the density is larger around the origin when the data pass the nonlinear layer. This suggests we may consider an adaptive mesh for more effectiveness, in other words, the mesh is finer around  $s = 1/2$  and coarser around  $s = 0, 1$ .

### 3.4. The complexity of KRnet

We count the number of trainable parameters in KRnet. For simplicity, we assume that each  $f_{[k]}^{\text{outer}}$  has  $L$  general coupling layers  $f_{[k,i]}^{\text{inner}}$ . Let  $d_k$  be the number of effective dimensions for  $f_{[k]}^{\text{outer}}$  and  $N_{\text{NN},k}$  the number of model parameters for the neural network Eq. (10) used in  $f_{[k,i]}^{\text{inner}}$ . We note that the main characteristic of KRnet is that a portion of dimensions will be deactivated as  $k$  increases. As  $d_k$  decreases with  $k$ , we expect that the neural network Eq. (10) in  $f_{[k,i]}^{\text{inner}}$  should become simpler for a larger  $k$ . In other words,  $N_{\text{NN},k}$  may decrease as  $k$  increases. For simplicity, we let  $N_{\text{NN},k} = rN_{\text{NN},k-1}$ , where  $0 < r < 1$ , without worrying about the detailed configuration of the neural network. The number of trainable parameters is  $d_k^2$  for  $L_R$ , and  $\hat{m}d$  for  $L_N$ , and  $2d_k$  for the scale and bias layer. Assume that  $d = mK$ . We have  $d_k = d - (k-1)m$ ,  $k = 1, \dots, K$ . According to the flow chart in Fig. 1, we have the total number of model parameters as

$$N_{\text{dof}} = \hat{m}d + \sum_{k=1}^{K-1} (N_{\text{NN},1} r^{k-1} L + (K-k+1)^2 m^2 + 2(K-k+1)mL). \quad (22)$$

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

### 3.5. KRnet for density estimation

We study the performance of KRnet for density estimation in this part and provide more results on the comparison between the real NVP and the KRnet that were not included in [46]. Once the KRnet is constructed, we train the model  $p_X(\mathbf{x}; \Theta)$  by maximizing the likelihood of the data or minimizing the cross entropy between the data distribution and the density model, where  $\Theta$  includes all the trainable model parameters. Let  $\mathcal{S} = \{\mathbf{x}^{(i)}\}_{i=1}^{N_t}$  be the training set and  $p_{X,\text{data}}(\mathbf{x})$  the underlying data distribution. The Kullback-Leibler (KL) divergence between  $p_{X,\text{data}}(\mathbf{x})$  and  $p_X(\mathbf{x}; \Theta)$  is

$$\min_{\Theta} D_{\text{KL}}(p_{X,\text{data}}(\mathbf{x}) || p_X(\mathbf{x}; \Theta)) = \mathbb{E}_{\mathbf{x} \sim p_{X,\text{data}}(\mathbf{x})} \left[ \log \frac{p_{X,\text{data}}(\mathbf{x})}{p_X(\mathbf{x}; \Theta)} \right] = H(p_{X,\text{data}}(\mathbf{x}), p_X(\mathbf{x}; \Theta)) - H(p_{X,\text{data}}(\mathbf{x})) \quad (23)$$

where  $H(p_{X,\text{data}}(\mathbf{x}))$  is the entropy of  $p_{X,\text{data}}(\mathbf{x})$  and  $H(p_{X,\text{data}}(\mathbf{x}), p_X(\mathbf{x}; \Theta))$  is the cross entropy of  $p_{X,\text{data}}(\mathbf{x})$  and  $p_X(\mathbf{x}; \Theta)$ . Since  $p_{X,\text{data}}(\mathbf{x})$  is independent of  $\Theta$ , minimizing the KL divergence is equivalent to minimizing the cross entropy. Note that

$$H(p_{X,\text{data}}(\mathbf{x}), p_X(\mathbf{x}; \Theta)) \approx -\frac{1}{N_t} \sum_{i=1}^{N_t} \log p_X(\mathbf{x}^{(i)}; \Theta), \quad (24)$$

which corresponds to the negation of the log-likelihood.

To measure the quality of KRnet, we compute the KL divergence Eq. (23) on a validation set between a reference PDF and the trained density model. The training data sets  $\mathcal{S}$  is generated as follows. Assume that  $X$  has i.i.d. components and each component  $X_i \sim \text{Logistic}(0, s)$  has a PDF  $\rho(x_i; 0, s)$ . We generate a sample  $\mathbf{x}^{(i)}$  of  $X$ , and then check if it satisfies the following constraint:

$$\left\| \mathbf{R}_{\gamma, \theta_j} [x_j^{(i)}, x_{j+1}^{(i)}]^\top \right\|_2 \geq C, \quad j = 1, \dots, d-1, \quad (25)$$

where  $C$  is a specified constant, and

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

The sample  $\mathbf{x}^{(i)}$  will be accepted if the constraint Eq. (25) is satisfied and rejected otherwise. This way, an elliptic hole is generated for any two adjacent dimensions of data points. The reference PDF is then defined as

$$p_{X,\text{ref}}(\mathbf{x}) = \frac{I_B(\mathbf{x}) \prod_{i=1}^d \rho(x_i; 0, s)}{\mathbb{E}[I_B(X)]}, \quad (26)$$

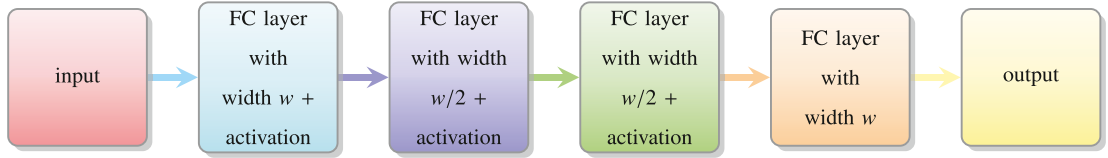
where  $B$  is the set defined by equation (25) and  $I_B(\cdot)$  is an indicator function with  $I_B(\mathbf{x}) = 1$  if  $\mathbf{x} \in B$ ; 0, otherwise. For this test problem, we set  $d = 8$ ,  $\gamma = 3$  and  $C = 7.6$ . This case has been studied in [46], where the rotation layers and nonlinear layers are turned off. In [46] an algebraic convergence has been observed numerically for both the real NVP and the KRnet, where the convergence rate of KRnet is about twice as large as that of the real NVP. We here only demonstrate the effectiveness of the rotation layer and the nonlinear layer.

We now compare the performance of KRnet and real NVP numerically. In KRnet, we deactivate the dimensions by one, i.e.,  $K = 7$ . We let  $N_{\text{NN},k} = 0.9N_{\text{NN},k-1}$  by adjusting the width of the neural network  $\text{NN}_{[i]}$ ,  $i = 1, 2, 3$ , which consists of two fully connected hidden layers of the same width. Other configurations of  $\text{NN}_{[i]}$  can also be considered. One example is given

**Table 1**

The effects of rotation and nonlinear layers in KRnet.  $\delta_I$ ,  $\delta_{II}$  and  $\delta_{III}$  are relative errors of KRnet, respectively, for the aforementioned three stages.  $\delta$  is the relative error of real NVP, whose depth is chosen to roughly match the DOFs of the KRnets from the same column. For the nonlinear layers, we use 32 nonuniform elements to decompose  $[-30, 30]$ , i.e.,  $a = 30$ . Note that the rotation layers and the nonlinear layers do not introduce a significant increase in the total number of DOFs. The percentages in parentheses indicate the degree of drop in terms of  $\delta_I$ .

KRnet	$L = 2$	$L = 4$	$L = 6$	$L = 8$
$\delta_I$	7.54e-2	2.45e-2	1.44e-2	9.50e-3
$\delta_{II}$	6.53e-2 ( $\downarrow 13\%$ )	2.24e-2 ( $\downarrow 9\%$ )	1.39e-2 ( $\downarrow 3\%$ )	9.11e-3 ( $\downarrow 4\%$ )
$\delta_{III}$	4.93e-2 ( $\downarrow 35\%$ )	1.95e-2 ( $\downarrow 20\%$ )	1.26e-2 ( $\downarrow 13\%$ )	8.34e-3 ( $\downarrow 12\%$ )
Real NVP	$L = 10$	$L = 20$	$L = 32$	$L = 42$
$\delta$	2.17e-2	1.98e-2	2.11e-2	2.05e-2



**Fig. 2.** The architecture of  $NN_{[i]}$  for affine coupling layers, for  $i = 0, \dots, L - 1$  (FC layers refer to fully connected layers).

in Fig. 2, which is used in section 5. The neural network  $NN_{[i]}$  (for  $i = 0, \dots, L - 1$ ) consists of three hidden layers and one linear layer, where the first hidden layer and the linear layer have  $w$  neurons, and the middle two layers have  $w/2$  neurons. In this experiment, we combine the two middle hidden layers to one hidden layer with  $w$  neurons. We set  $w = 24$  and use the rectified linear unit function (ReLU) as the activation function [50]. The depth of the real NVP will be determined by  $N_{\text{dof}}$  of the KRnet, since we split the dimensions into two halves in real NVP. The KRnet will be implemented as follows. We train KRnet with three stages and record the errors of each stage. In the first stage, we switch off both the rotation layers and the nonlinear layers and train the model for 8000 epochs; in the second stage, we switch on the rotation layers and restart the training process for another 2000 epochs; finally, we switch on both the rotation layers and the nonlinear layers and continue the training process for another 2000 epochs. For the real NVP, we simply run 8000 epochs. For each epoch, we compute the relative error

$$\delta = \frac{D_{KL}(p_{X,\text{ref}}(\mathbf{x}) || p_X(\mathbf{x}; \Theta))}{H(p_{X,\text{ref}}(\mathbf{x}))} \quad (27)$$

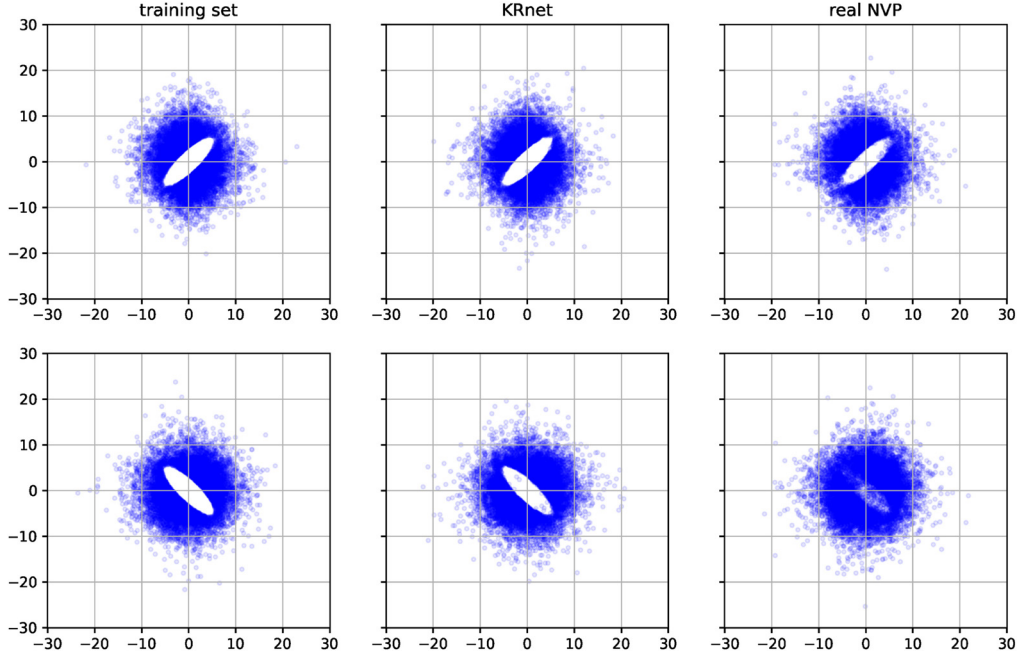
using the validation set, since the cross entropy should converge to the differential entropy of the reference PDF. We record the minimum relative error of all epochs. Furthermore, to reduce the bias of  $\delta$ , we will sample 10 independent training sets and repeat the training process ten times to obtain an averaged relative error  $\delta$ . The relative errors corresponding to the above three stages of training KRnet are denoted as  $\delta_I$ ,  $\delta_{II}$  and  $\delta_{III}$ . We will sample  $3.2 \times 10^5$  data points for both the training set and the validation set. We employ the Adam optimizer [51] with learning rate 0.001 and batch size 80000.

The results of numerical experiments have been summarized in Table 1. First of all, both  $\delta_{II}$  and  $\delta_{III}$  are smaller than  $\delta_I$ , indicating that the rotation layers and nonlinear layers are able to improve the model performance. Such an improvement is more noticeable for a smaller  $L$ . Second, for the specific setup of the numerical experiments, the errors  $\delta_i$ ,  $i = I, II, III$ , of the KRnet decay consistently as  $L$  increases while the errors of the real NVP do not show consistent decay. Since we compute the errors after 8000 epochs for all  $L$ , this shows that for a comparable model complexity the KRnet needs less epochs to obtain a substantial decrease in error than the real NVP. Third, as also shown in [46], the real NVP performance better than KRnet for a small  $L$ . The real NVP can be regarded as a KRnet with a half-half partition, i.e.,  $K = 2$  and  $m = \frac{d}{2} = 4$ . For a fixed complexity, the performance of KRnet depends on both  $K$  and  $L$ . In Fig. 3, we compare the approximated distributions given by the real NVP with  $L = 42$  and the KRnet with  $L = 8$ , where both the rotation layers and the nonlinear layers are switched on.

#### 4. Adaptive deep density approximation for the stationary Fokker-Planck equation

We intend to use KRnet as a PDF model to approximate the Fokker-Planck equation to alleviate the difficulties from the curse of dimensionality. In particular, we will develop an adaptive deep density approximation (ADDA) approach, which consists of two components: 1) solving the Fokker-Planck equation on a certain set of collocation points by a machine learning technique; 2) choosing a new set of collocation points to refine the current approximate solution. These two components are implemented alternately to achieve adaptivity such that both the accuracy and the efficiency will be improved.





**Fig. 3.** Training data, and data sampled from KRnet and real NVP. The first row shows the components  $x_1$  and  $x_2$ , and the second row shows the components  $x_4$  and  $x_5$ . We pick the two pairs of adjacent dimensions, where the real NVP performs the best and the worst, respectively.

#### 4.1. Stochastic gradient descent based on stochastic collocation points

Let  $p_X(\mathbf{x}; \Theta)$  be a probability density function associated with the random vector  $X$ , which is based on the KRnet. All the constraints in Eq. (4) and Eq. (5) are naturally satisfied since  $p_X(\mathbf{x}; \Theta)$  is a family of probability density functions, implying that the difficulties caused by the boundary conditions and the nonnegativity of PDF have disappeared. We seek to approximate the solution  $p(\mathbf{x})$  of the Fokker-Planck equation by  $p_X(\mathbf{x}; \Theta)$  to take advantage of the weaker dependence of deep neural networks on dimensionality than traditional computational approaches such as the finite element methods [52–54].

The main idea of a machine learning approach to solve PDEs is to consider an optimization problem defined on a set of collocation points where the equation is constrained. Let  $p_{\text{data}}(\mathbf{x})$  be a probability density function, based on which we define a loss functional

$$J(p_X(\mathbf{x}; \Theta)) = \mathbb{E}_{p_{\text{data}}(\mathbf{x})} (r^2(\mathbf{x}; \Theta)) = \mathbb{E}_{p_{\text{data}}(\mathbf{x})} (\mathcal{L}^2(p_X(\mathbf{x}; \Theta))) \quad (28)$$

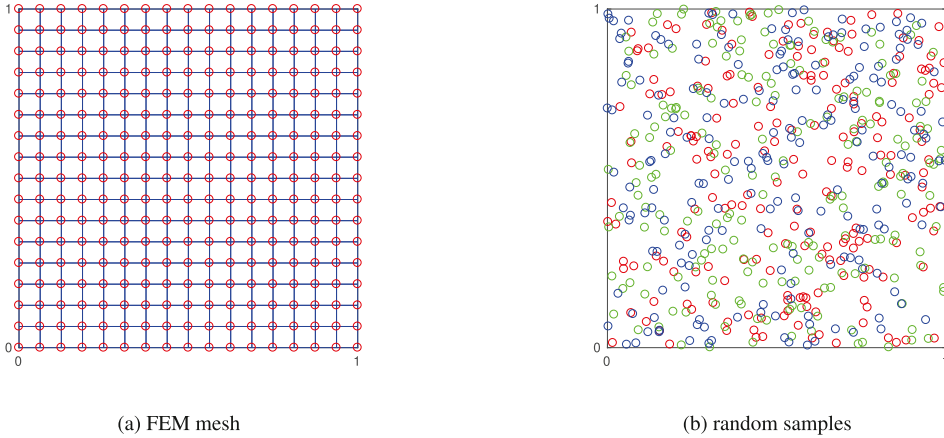
where  $\mathbb{E}_{p_{\text{data}}(\mathbf{x})}$  denotes the expectation with respect to the training set, and  $r$  is the residual loss. The solution  $p(\mathbf{x})$  of Eq. (3) can be approximated by  $p_X(\mathbf{x}; \Theta)$  through minimizing the loss functional  $J(p_X(\mathbf{x}; \Theta))$ . In reality, we usually do not have much prior understanding about the residual, and simply assign  $p_{\text{data}}(\mathbf{x})$  a simple distribution, e.g., a uniform distribution defined on a finite computational domain. We then use  $p_{\text{data}}(\mathbf{x})$  to sample a set  $\mathcal{C} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  of collocation points to approximate the loss functional, i.e.,

$$\hat{J}(p_X(\mathbf{x}; \Theta)) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}^2(p_X(\mathbf{x}^{(i)}; \Theta)) \approx J(p_X(\mathbf{x}; \Theta)), \quad (29)$$

based on which we choose the optimal parameter  $\Theta^*$ :

$$\Theta^* = \arg \min_{\Theta} \hat{J}(p_X(\mathbf{x}; \Theta)). \quad (30)$$

The optimization problem Eq. (30) will be solved by stochastic gradient-based optimization [55,51], which is summarized as follows. The set of collocation points can be divided into  $n_b$  mini-batches  $\{C_{i_b}\}_{i_b=1}^{n_b}$ , where every mini-batch  $C_{i_b}$  contains  $m$  samples such that  $N = mn_b$ . Denoting the parameters at  $i_b$ -th iteration of a certain epoch  $j$  as  $\Theta_{i_b}^{(j)}$ , for every mini-batch  $C_{i_b}$  and  $\mathbf{x}^{(l)} \in C_{i_b}$ ,  $l = 1, \dots, m$ , one can apply the mini-batch to estimate the expectation of the residual loss and the stochastic gradient, and then update the parameters  $\Theta$  based on the following scheme



**Fig. 4.** An example of linear finite element meshes and stochastic collocation points in  $[0, 1]^2$ .

$$\Theta_{i_b}^{(j)} = \Theta_{i_b-1}^{(j)} - \eta \nabla_{\Theta} \left[ \frac{1}{m} \sum_{l=1}^m \left( r(\mathbf{x}^{(l)}; \Theta_{i_b-1}^{(j)}) \right)^2 \right] \quad \text{for } i_b = 1, \dots, n_b, j = 1, 2, \dots \quad (31)$$

where  $\eta$  is a given learning rate. Compared with the gradient descent method, the stochastic gradient descent method only requires computing the gradient on the mini-batch  $C_{i_b}$ . In this work we employ the Adam optimizer, which is widely used to accelerate the training process for deep neural networks, as this method adopts adaptive learning rates for different components of parameters through estimates of first and second moments of the gradients [51].

#### 4.2. Adaptive sampling procedure

Compared with the standard finite element methods (FEM) [5], the deep learning approach does not require mesh generation to solve PDEs, which shares more similarities to meshless methods, and the approximation of Eq. (28) fits naturally with stochastic gradient-based optimization. Fig. 4 shows a linear finite element mesh in  $[0, 1]^2$  and the collocation points that are generated with a uniform distribution in  $[0, 1]^2$ .

Adaptivity plays an important role in classical numerical methods for the approximation of PDEs. Considering a finite element method subject to a certain mesh of the computation domain, we expect that the element-wise approximation errors are distributed in a nearly uniform way. This means that the most effective mesh should be non-uniform since the regularity of the solution varies in the computation domain. In our problem, the distribution  $p_{\text{data}}(\mathbf{x})$  of the collocation points will affect the approximation of  $J(p_X(\mathbf{x}; \Theta))$  and the optimal parameter  $\Theta^*$  as well. Apparently a uniform distribution is not an optimal choice for  $p_{\text{data}}(\mathbf{x})$  especially for high-dimensional problems. For a certain amount of collocation points, the curse of dimensionality will weaken the contribution of each collocation point to our learning problem, which will be worsen for the approximation of PDF if the exact solution  $p(\mathbf{x})$  is far away from being uniform. We then expect to use samples from a nonuniform distribution  $p_{\text{data}}(\mathbf{x})$  for the approximation  $J(p_X(\mathbf{x}; \Theta))$ , where a simple criterion is that  $p_{\text{data}}(\mathbf{x})$  should be consistent with the true solution  $p(\mathbf{x})$  to some extent. This will result in adaptive deep density approximation (ADDA) for the approximation of the Fokker-Planck equation.

It is, in general, difficult to generate samples that are adaptive to the true solution  $p(\mathbf{x})$ . Fortunately, flow-based deep generative models provide an opportunity for us to do this thanks to the invertible mapping. Our strategy is as follows. Starting with an initial set of collocation points  $C_0 = \{\mathbf{x}_{(0)}^{(i)}\}_{i=1}^N$  drawn from a uniform distribution, we train and obtain the KRnet  $Z = f_{\text{KR},(0)}(X; \Theta^{*,(0)})$ , which corresponds to the PDF  $p_X^{(0)}(\mathbf{x}; \Theta^{*,(0)})$ . We then generate a new set  $C_1 = \{\mathbf{x}_{(1)}^{(i)}\}_{i=1}^N$  of collocation points by  $X = f_{\text{KR},(0)}^{-1}(Z)$  using  $N$  samples from the prior distribution of  $Z$ . Then  $C_1$  is a set of samples from  $p_X^{(0)}(\mathbf{x}; \Theta^{*,(0)})$ . We continue to update the KRnet using  $\Theta^{*,(0)}$  as the initial parameters and  $C_1$  as the training set, which yields  $f_{\text{KR},(1)}(X; \Theta^{*,(1)})$ . Then another iteration starts. In general, we sample the current optimal PDF model  $p_X^{(k)}(\mathbf{x}; \Theta^{*,(k)})$  to generate a new training set  $C_{k+1} = \{\mathbf{x}_{(k+1)}^{(i)}\}_{i=1}^N$  and update the KRnet to  $f_{\text{KR},(k+1)}(\mathbf{x}; \Theta^{*,(k+1)})$ . This way, the samples for the training process become more and more consistent with the true solution, if  $p_X^{(k)}(\mathbf{x}; \Theta^{*,(k)})$  approaches  $p(\mathbf{x})$  as  $k$  increases. In other words, more collocation points will be chosen in the region of high density while less collocation points in the region of low density. Our adaptive training process has been summarized in Algorithm 1, where  $N_{\text{adaptive}} \in \mathbb{N}$  is a given number of maximum adaptivity iterations, and this strategy is called the adaptive deep density approximation based on KRnet (ADDA-KR) from now on. The final KRnet-induced PDF is the ADDA-KR approximation for the steady state Fokker-Planck problem (3)–(5).

We note that the adaptivity in Algorithm 1 can be further tuned. One possible strategy is to update the training set gradually for each training stage, e.g., up to a certain percentage. In this work, we replace the whole training set from the previous stage just for simplicity.

### 4.3. Implementation issues

When minimizing the loss functional Eq. (29), numerical underflow issues can be encountered, especially when  $\mathbf{x}$  is relatively high-dimensional. That is, the loss functional can be too small to provide an effective gradient descent direction. To alleviate this issue, we develop the following scaling strategy in our implementation. Multiplying both sides of equation (3) by a constant  $C_s > 0$  gives

$$\mathcal{L}(C_s p(\mathbf{x})) = \nabla \cdot [C_s p(\mathbf{x}) \nabla V(\mathbf{x})] + \nabla \cdot [\nabla \cdot (C_s p(\mathbf{x}) \mathbf{D}(\mathbf{x}))] = 0. \quad (32)$$

The solution of the above equation is the same as the solution of the original stationary Fokker-Planck equation (3). However, if  $C_s$  is large enough, Eq. (32) is numerically more stable than Eq. (3), and the loss functional Eq. (29) associated with Eq. (32) can typically provide effective gradient descent directions to optimize the parameters  $\Theta$ . In our practical implementation, we usually set  $C_s = 100$ .

---

#### Algorithm 1 Adaptive deep density approximation based on KRnet (ADDA-KR) for the Fokker-Planck equation.

---

**Input:** Initial KRnet  $p_X^{(0)}(\mathbf{x}; \Theta_0^{(0)})$ , maximum epoch number  $N_e$ , maximum iteration number  $N_{\text{adaptive}}$ , learning rate  $\eta$ , batch size  $m$ , and initial training set

```

 $C_0 = \left\{ \mathbf{x}_{(0)}^{(i)} \right\}_{i=1}^N$ .
1: Divide  $C_0 = \left\{ \mathbf{x}_{(0)}^{(i)} \right\}_{i=1}^N$  into  $n_b$  mini-batch  $\{C_{i_b}\}_{i_b=1}^{n_b}$ .
2: for  $k = 1 : N_{\text{adaptive}}$  do
3:   for  $j = 0 : N_e - 1$  do
4:     for  $i_b = 1 : n_b$  do
5:       Compute the values of the residual loss  $r(\mathbf{x}_{(k-1)}^{(l)}; \Theta_{i_b-1}^{(j)})$  (see Eq. (29)) for  $l = 1, \dots, m$ , on the mini-batch  $C_{i_b}$ .
6:       Update the parameters  $\Theta_{i_b}^{(j)}$  using the Adam optimizer with learning rate  $\eta$ .
7:     end for
8:     if  $j = N_e - 1$  then
9:       Let  $\Theta^{*,(k)} := \Theta_{n_b}^{(N_e-1)}$ .
10:    else
11:      Let  $\Theta_0^{(j+1)} := \Theta_{n_b}^{(j)}$ .
12:    end if
13:    Shuffle the set of collocation points  $C_{k-1} = \left\{ \mathbf{x}_{(k-1)}^{(i)} \right\}_{i=1}^N$ .
14:    Divide  $C_{k-1} = \left\{ \mathbf{x}_{(k-1)}^{(i)} \right\}_{i=1}^N$  into  $n_b$  mini-batch  $\{C_{i_b}\}_{i_b=1}^{n_b}$ .
15:  end for
16:  if  $k = N_{\text{adaptive}}$  then
17:    Let  $\Theta := \Theta^{*,(k)}$ .
18:  else
19:    Generate  $C_{k+1} = \left\{ \mathbf{x}_{(k+1)}^{(i)} \right\}_{i=1}^N$  by  $p_X^{(k)}(\mathbf{x}; \Theta^{*,(k)})$ .
20:    Let  $\Theta_0^{(0)} := \Theta^{*,(k)}$ .
21:  end if
22: end for
23: Obtain the ADDA-KR solution  $p_X(\mathbf{x}; \Theta) := p_X^{(N_{\text{adaptive}})}(\mathbf{x}; \Theta)$ .

```

**Output:** The ADDA-KR solution  $p_X(\mathbf{x}; \Theta)$ .

---

## 5. Numerical study

In this section, numerical experiments are conducted to illustrate the effectiveness of our ADDA-KR (adaptive deep density approximation based on KRnet) approach presented in Algorithm 1. Five test problems for the Fokker-Planck equation are studied—one one-dimensional test problem, two two-dimensional test problems (one is a single modal distribution, and the other is a bimodal distribution), one four-dimensional test problem, and one eight-dimensional test problem. The activation function of  $\text{NN}_{|ij|}$  (see Eq. (10)) is set to the hyperbolic tangent function for all test problems. For comparison, we also test the performance of a direct adaptive version of classic real NVP, and as the real NVP utilizes a half-half partition (see section 3.5), we refer to it as ADDA-HH. The implementation of ADDA-HH is to replace the KRnet in ADDA-KR (Algorithm 1) by the classical real NVP, and we set the same input parameters for both ADDA-KR and ADDA-HH in all our test problems. In addition, results of non-adaptive versions of KRnet and real NVP are included for high-dimensional test problems (the four-dimensional and the eight-dimensional test problems), which are referred to as Uniform-KR and Uniform-HH. In Uniform-KR and Uniform-HH, collocation points are generated through uniform distributions, and other settings of KRnet and real NVP are the same as the settings for ADDA-KR in these test problems.

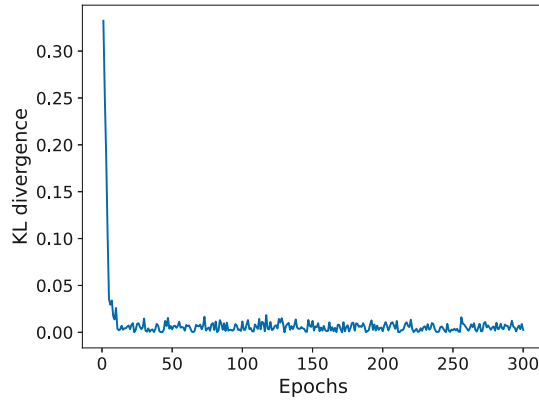


Fig. 5. The KL divergence with respect to epochs, one-dimensional test problem.

### 5.1. A one-dimensional test problem

We start with this one-dimensional case, where the governing equation is

$$\begin{aligned} \frac{\partial(xp(x))}{\partial x} + \frac{1}{2} \frac{\partial^2(p(x))}{\partial x^2} &= 0, \\ \int_{\mathbb{R}} p(x) dx &= 1, \quad p(x) \geq 0, \end{aligned} \quad (33)$$

and the exact solution is

$$p(x) = \frac{\exp(-x^2)}{\sqrt{\pi}}. \quad (34)$$

For this one-dimensional problem, KRnet is the same as the classical real NVP, meaning that only the affine coupling layers are needed. As the affine coupling layers (see section 3.1) need at least two-dimensions, we use  $[x, x]$  as an input in our implementation of KRnet. We generate the initial parameters  $\Theta_0^{(0)}$  for the inputs of Algorithm 1, using Glorot Gaussian initialization [56], and then construct the initial KRnet  $p_X^{(0)}(\mathbf{x}; \Theta_0^{(0)})$ . The number of epochs is set to  $N_e = 300$ , and only one adaptivity iteration is conducted for this one-dimensional problem, i.e.,  $N_{\text{adaptive}} = 1$ . The learning rate for Adam optimizer is set to  $\eta = 0.0002$ , and the batch size is set to  $m = 500$ . The initial training set  $\mathcal{C}_0$  is generated through the uniform distribution with range  $[-5, 5]$ , and the sample size is set to  $|\mathcal{C}_k| = 3000$  for each iteration step  $k$  for  $k = 0, \dots, N_{\text{adaptive}}$ . In addition, we take  $L = 8$  affine coupling layers, and two fully connected layers with  $w = 48$  neurons for  $\text{NN}_{[ij]}$  (see Eq. (10)).

To assess the accuracy of our ADDA-KR approach (Algorithm 1), we compute the KL divergence between the exact solution  $p(x)$  and our ADDA-KR solution  $p_X(x; \Theta)$ :

$$\begin{aligned} D_{KL}(p(x) || p_X(x; \Theta)) &= \int_{-\infty}^{\infty} p(x) \log p(x) dx - \int_{-\infty}^{\infty} p(x) \log p_X(x; \Theta) dx \\ &= -\frac{1}{2}(1 + \log \pi) - \int_{-\infty}^{\infty} p(x) \log p_X(x; \Theta) dx \end{aligned}$$

where the last term of the above equation is approximated by Monte Carlo integration with  $10^4$  samples. Fig. 5 shows the KL divergence decreases to zero quickly. Fig. 6 shows the exact solution  $p(x)$  and our ADDA-KR solution  $p_X(x; \Theta)$ , where it can be seen that they are visually indistinguishable.

### 5.2. Two-dimensional test problems

In this part, two-dimensional Fokker-Planck equations are considered, where the solution of the first one is a single modal distribution and the solution of the second one is a bimodal distribution.

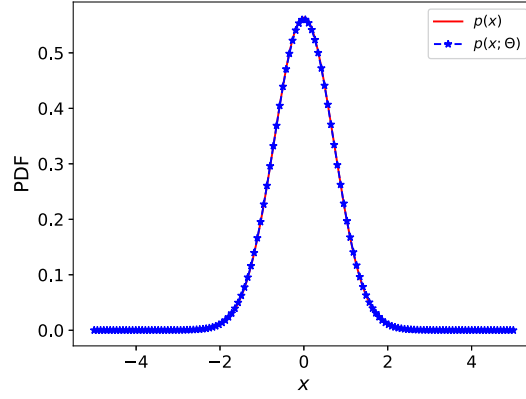


Fig. 6. The exact solution and the ADDA-KR solution, one-dimensional test problem.

### 5.2.1. Two-dimensional single modal distribution

The stationary Fokker-Planck equation for this test problem is

$$\begin{aligned} \nabla \cdot [p(\mathbf{x})\mathbf{A}\mathbf{x}] + \nabla \cdot [\nabla \cdot (p(\mathbf{x})\mathbf{D})] &= 0, \\ \int_{\mathbb{R}^d} p(\mathbf{x})d\mathbf{x} &= 1, \quad p(\mathbf{x}) \geq 0, \end{aligned} \quad (35)$$

where  $\mathbf{A}$  and  $\mathbf{D}$  are two constant matrices. This equation is corresponding to the following Ornstein-Uhlenbeck process

$$dX_t = -\mathbf{A}X_t dt + \mathbf{G}d\mathbf{w}_t, \quad (36)$$

where  $\mathbf{D} = \mathbf{G}\mathbf{G}^T/2$ .

The solution of Eq. (35) exists if the real parts of the eigenvalues of  $\mathbf{A}$  are larger than zero [1], and it can be written as

$$p(\mathbf{x}) = (2\pi)^{-1} (\det \boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}\right), \quad (37)$$

where the covariance matrix  $\boldsymbol{\Sigma}$  is determined by the following Lyapunov equation

$$\mathbf{A}\boldsymbol{\Sigma} + \boldsymbol{\Sigma}\mathbf{A}^T = 2\mathbf{D}. \quad (38)$$

The above Lyapunov equation has a unique solution if and only if the eigenvalues  $\lambda_i$  of  $\mathbf{A}$  satisfy  $\lambda_i \neq -\lambda_j$  for all  $i, j = 1, 2$ . In this test problem, the constant matrix  $\mathbf{A}$  for the drift term and the diffusion matrix  $\mathbf{D}$  are set to

$$\mathbf{A} = \begin{bmatrix} 1.37096037 & -0.48306187 \\ -0.48306187 & 1.62903963 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 22.52429192 & -6.55821381 \\ -6.55821381 & 12.68972 \end{bmatrix},$$

which implies that the covariance matrix  $\boldsymbol{\Sigma}$  is

$$\boldsymbol{\Sigma} = \begin{bmatrix} 8.12186142 & -0.26372569 \\ -0.26372569 & 3.81664391 \end{bmatrix}.$$

We generate the initial parameters  $\Theta_0^{(0)}$  with Glorot Gaussian initialization [56], and then construct the initial KRnet  $p_X^{(0)}(\mathbf{x}; \Theta_0^{(0)})$  for Algorithm 1. The number of epochs is set to  $N_e = 300$ , and two adaptivity iterations are conducted for this problem, i.e.,  $N_{\text{adaptive}} = 2$ . The learning rate for Adam optimizer is set to  $\eta = 0.0002$ , and the batch size is set to  $m = 1000$ . The initial training set  $\mathcal{C}_0$  is generated through the uniform distribution with range  $[-6, 6]^2$ , and the sample size is set to  $|\mathcal{C}_k| = 6 \times 10^4$  for each iteration step  $k$  for  $k = 0, \dots, N_{\text{adaptive}}$ . In addition, we take  $L = 8$  affine coupling layers, and two fully connected layers with  $w = 48$  neurons for  $\text{NN}_{[ij]}$  (see Eq. (10)).

Fig. 7 shows the exact solution  $p(\mathbf{x})$  and our ADDA-KR solution  $p_X(\mathbf{x}; \Theta)$ , where it can be seen that they are visually indistinguishable. For this test problem, there is no significant difference between the ADDA-KR solution and the ADDA-HH solution, and we then only show the exact solution and our ADDA-KR solution. Fig. 8 shows samples drawn from the exact solution of Eq. (35) and our ADDA-KR solution, which confirms that the corresponding distributions ( $p(x)$  and  $p_X(\mathbf{x}; \Theta)$ ) are very close.

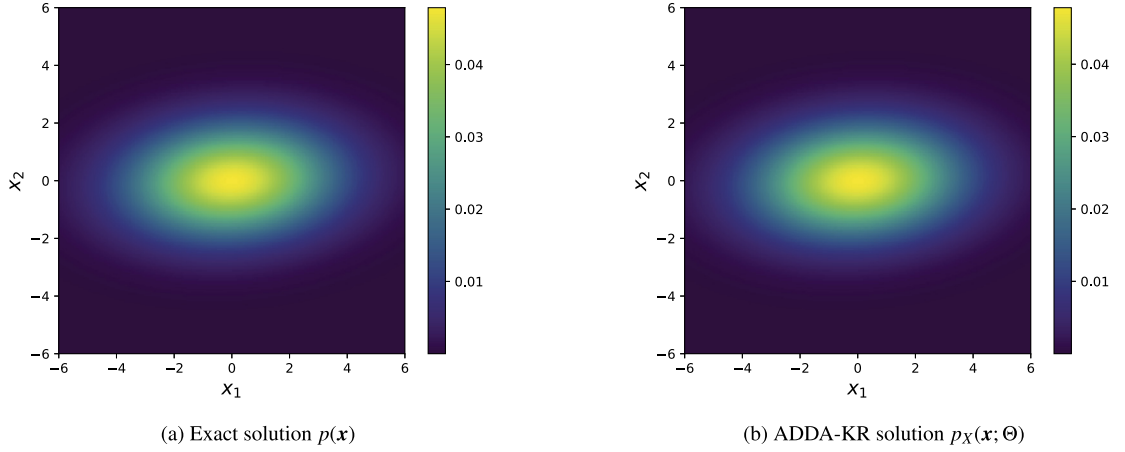


Fig. 7. Solutions, two-dimensional single modal test problem.

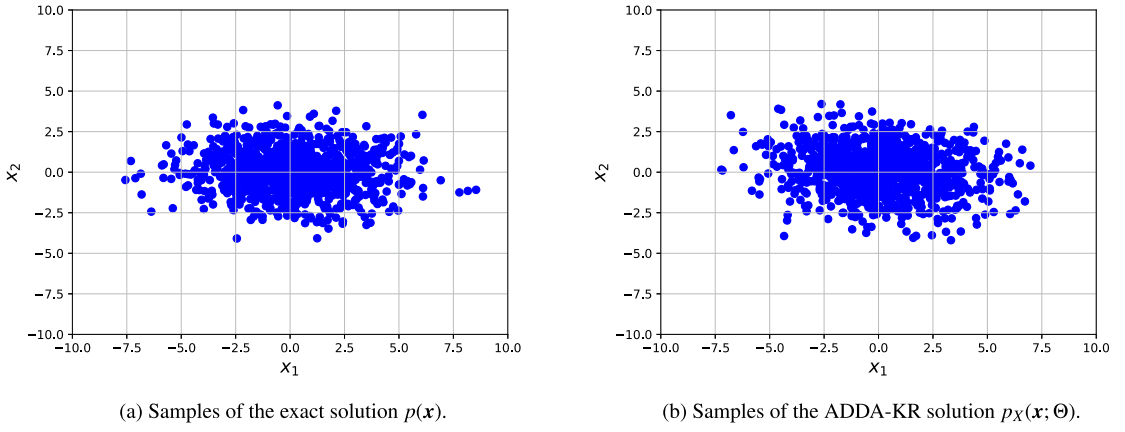


Fig. 8. Samples, two-dimensional single modal test problem.

### 5.2.2. Two-dimensional bimodal distribution

In this test problem, the Fokker-Planck equation considered is

$$\begin{aligned}
 -\nabla \cdot [p(\mathbf{x}) \nabla \log(\beta_1 p_1(\mathbf{x}) + \beta_2 p_2(\mathbf{x}))] + \nabla^2 p(\mathbf{x}) &= 0, \\
 \int_{\mathbb{R}^d} p(\mathbf{x}) d\mathbf{x} &= 1, \quad p(\mathbf{x}) \geq 0,
 \end{aligned} \tag{39}$$

where for  $k = 1, 2$ , each  $p_k(\mathbf{x})$  is the probability density function of the normal distribution with mean  $\mu_k$  and covariance  $\Sigma_k$ , and  $\beta_1 + \beta_2 = 1$ . The solution of Eq. (39) is the following Gaussian mixture distribution [47, p. 123],

$$p(\mathbf{x}) = \beta_1 p_1(\mathbf{x}) + \beta_2 p_2(\mathbf{x}). \tag{40}$$

Here, we set  $\mu_k$ ,  $\Sigma_k$  and  $\beta_k$  for  $k = 1, 2$  as

$$\begin{aligned}
 \beta_1 = 0.55, \quad \beta_2 = 0.45, \quad \mu_1 = [-1, -1]^T, \quad \mu_2 = [2, 2]^T \\
 \Sigma_1 = \begin{bmatrix} 6.12186142 & -0.26372569 \\ -0.26372569 & 1.81664391 \end{bmatrix}, \quad \Sigma_2 = \begin{bmatrix} 2.8828528 & -0.70234742 \\ -0.70234742 & 2.69199911 \end{bmatrix}.
 \end{aligned} \tag{41}$$

The matrices  $\Sigma_1$  and  $\Sigma_2$  are positive definite, and their entries are randomly constructed.

We again generate the initial parameters  $\Theta_0^{(0)}$  with Glorot Gaussian initialization, and then construct the initial KRnet  $p_X^{(0)}(\mathbf{x}; \Theta_0^{(0)})$ . The number of epochs is set to  $N_e = 200$ , and the maximum number of adaptivity iterations conducted for this problem is set to  $N_{\text{adaptive}} = 5$ . The learning rate for Adam optimizer is set to  $\eta = 0.0001$ , and the batch size is set to  $m = 1000$ . The initial training set  $C_0$  is generated through the uniform distribution with range  $[-5, 5]^2$ , and the sample size

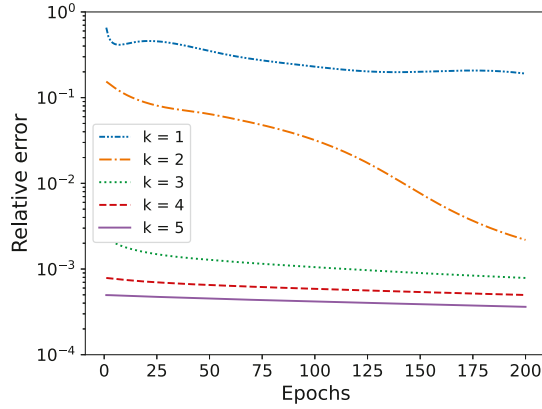


Fig. 9. The relative error for ADDA-KR, two-dimensional bimodal test problem.

is set to  $|\mathcal{C}_k| = 6 \times 10^4$  for each iteration step  $k$  for  $k = 0, \dots, N_{\text{adaptive}}$ . In addition, we take  $L = 8$  affine coupling layers for both KRnet and real NVP, and two fully connected layers with  $w = 48$  neurons for  $\text{NN}_{[i]}$  (see Eq. (10)). For KRnet, we set  $K = 2$  to focus on the effectiveness of the rotation layer and the nonlinear layer for this test problem. To assess the effectiveness of our ADDA-KR approach, we generate a validation data set  $\mathcal{C}_v = \{\mathbf{x}^{(i)}\}_{i=1}^{N_v}$ , and compute the relative error defined by Eq. (27). The KL divergence is approximated by Monte Carlo integration

$$D_{KL}(p(\mathbf{x})||p_X(\mathbf{x}; \Theta)) \approx \frac{1}{N_v} \sum_{i=1}^{N_v} \left( \log p(\mathbf{x}^{(i)}) - \log p(\mathbf{x}^{(i)}; \Theta) \right), \quad (42)$$

where  $\mathbf{x}^{(i)}$  is drawn from the exact solution  $p(\mathbf{x})$ , and the size of the validation data set is set to  $3.2 \times 10^5$  such that the KL-divergence can be approximated well.

Fig. 9 shows the relative error between the exact solution  $p(\mathbf{x})$  and our ADDA-KR solution  $p_X(\mathbf{x}; \Theta)$  at each adaptivity iteration step  $k$ . It is clear that, as the adaptivity iteration step increases, the relative error decreases quickly. In addition, it can be seen that as the number of epochs increases, the relative error decreases. Fig. 10 shows the comparison between our ADDA-KR and ADDA-HH. From Fig. 10(a), it can be seen that the relative error of ADDA-KR is smaller than that of ADDA-HH at each adaptivity iteration step. Fig. 10(b), Fig. 10(c) and Fig. 10(d) show the relative error decreases as the number of epochs increases, at adaptivity iteration steps  $k = 1, 3, 5$  respectively. It can be seen that the relative error of ADDA-KR is clearly smaller than that of ADDA-HH for each value of epochs, except for the situations that the epoch number is smaller than 125 at the first adaptivity iteration in 10(b). Fig. 11 shows the exact solution  $p(\mathbf{x})$  and the ADDA-KR solution  $p_X(\mathbf{x}; \Theta)$ , where it can be seen that this bimodal distribution is well approximated by our ADDA-KR solution.

### 5.3. High-dimensional bimodal distributions (four-dimensional and eight-dimensional test problems)

In this part, we again consider the Fokker-Planck equation with two peaks Eq. (39) and set  $\beta_1 = 0.55, \beta_2 = 0.45$ . However, the dimensionality of the problem considered in this part is different from section 5.2. We here consider a four-dimensional ( $d = 4$ ) problem and an eight-dimensional ( $d = 8$ ) problem. The exact solution of Eq. (39) is a Gaussian mixture distribution Eq. (40). For  $d = 4$ , we set

$$\begin{aligned} \boldsymbol{\mu}_1 &= [-1, -1, -0.3, -0.3]^T, \quad \boldsymbol{\mu}_2 = [2, 2, 0.6, 0.6]^T \\ \boldsymbol{\Sigma}'_1 &= \begin{bmatrix} \boldsymbol{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & 0.6\boldsymbol{\Sigma}_1 \end{bmatrix}, \quad \boldsymbol{\Sigma}'_2 = \begin{bmatrix} \boldsymbol{\Sigma}_2 & \mathbf{0} \\ \mathbf{0} & 0.6\boldsymbol{\Sigma}_2 \end{bmatrix}, \end{aligned} \quad (43)$$

where  $\boldsymbol{\Sigma}_1$  and  $\boldsymbol{\Sigma}_2$  are given in Eq. (41), and  $\boldsymbol{\Sigma}'_1$  and  $\boldsymbol{\Sigma}'_2$  are the covariance matrices of  $p_1$  and  $p_2$  for this test problem.

Similarly to the previous settings, we generate the initial parameters  $\Theta_0^{(0)}$  with Glorot Gaussian initialization, and then construct the initial KRnet  $p_X^{(0)}(\mathbf{x}; \Theta_0^{(0)})$ . The number of epochs is set to  $N_e = 1$ , and the number of adaptivity iterations conducted for this problem is set to  $N_{\text{adaptive}} = 16$ . Here, KRnet is trained and sampled in an interleaved manner. That is for both the four-dimensional and the eight-dimensional test problems, samples  $\mathcal{C}_k$  drawn at the  $k$ -th adaptivity iteration are immediately used for training KRnet at the  $(k + 1)$ -th iteration, while  $p_X^{(k+1)}(\mathbf{x}; \Theta)$  is immediately used for sampling. The learning rate for Adam optimizer is set to  $\eta = 0.0001$ , and the batch size is set to  $m = 500$ . The initial training set  $\mathcal{C}_0$  is generated through the uniform distribution with range  $[-6, 6]^d$ , and two cases of the collocation sample size are considered: one is  $10^5$  and the other is  $2 \times 10^5$ . In addition, we take  $L = 8$  affine coupling layers for KRnet, and  $L = 16$  for real NVP. The architecture of  $\text{NN}_{[i]}$  is the same as that shown in Fig. 2 with  $w = 120$ . For KRnet, we set  $K = 3$ . The

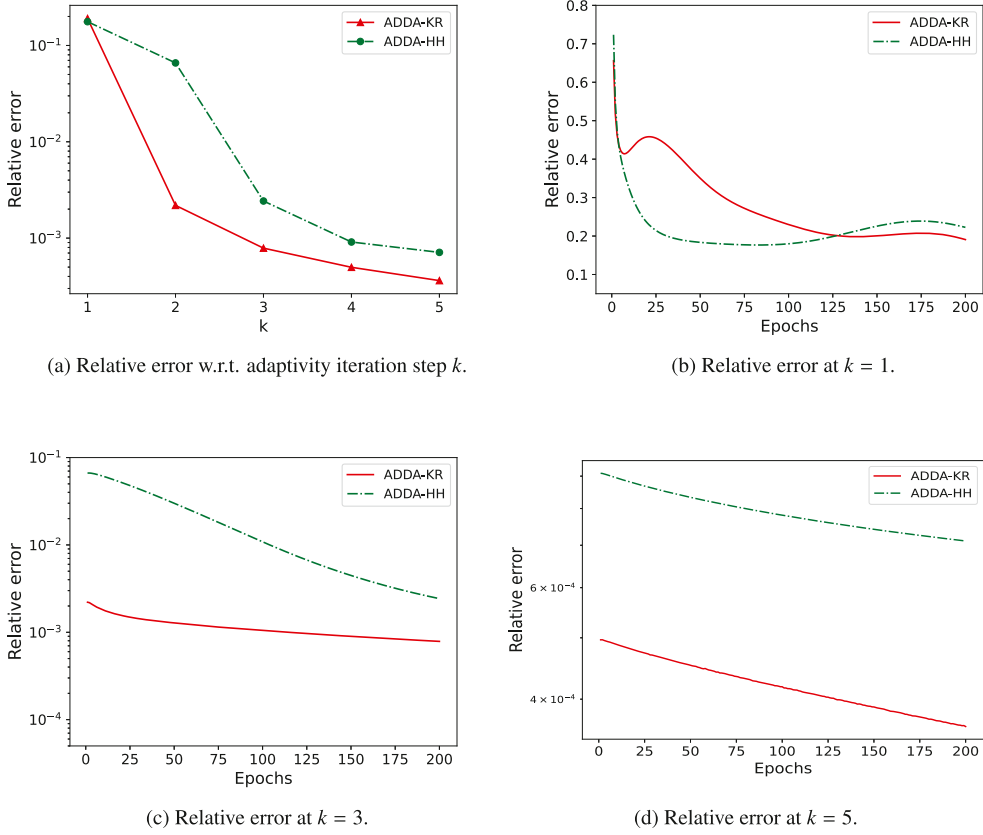


Fig. 10. The relative error for ADDA-KR and ADDA-HH, two-dimensional bimodal test problem.

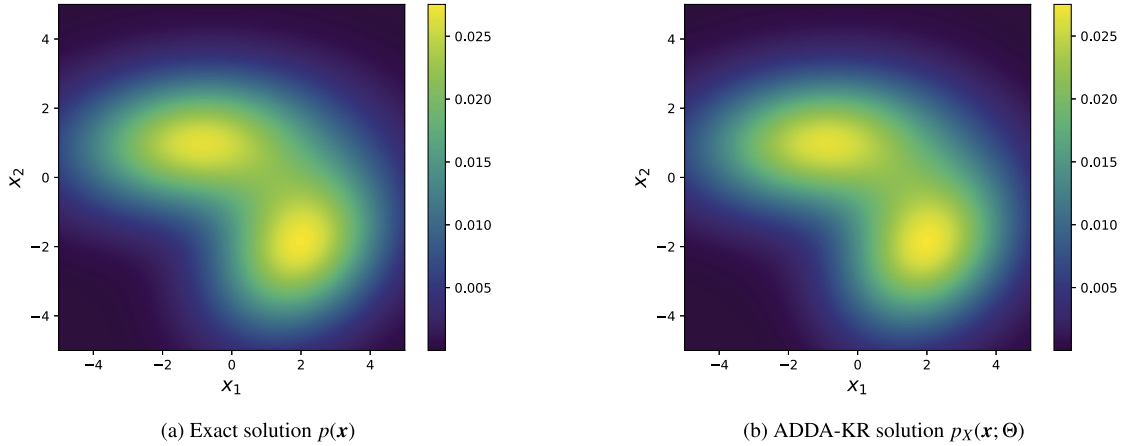


Fig. 11. Solutions, two-dimensional bimodal test problem.

rotation layer and the nonlinear layer are turned on. To assess the accuracy of ADDA-KR, we again compute the relative error Eq. (27) between  $p(\mathbf{x})$  and  $p_X(\mathbf{x}; \Theta)$  using  $3.2 \times 10^5$  validation samples drawn from the exact solution.

Fig. 12 shows the relative error between  $p(\mathbf{x})$  and  $p_X(\mathbf{x}; \Theta)$  for ADDA-KR and ADDA-HH, where different numbers of collocation points are considered. From Fig. 12(a), it can be seen that the relative error of ADDA-KR is smaller than that of ADDA-HH. From Fig. 12(b) and Fig. 12(c), as the number of epochs increases, the relative errors of ADDA-KR and ADDA-HH decrease quickly, while the relative errors of the uniform sampling strategies (Uniform-KR and Uniform-HH) decrease slowly. In addition, it can be seen that the relative error decreases as the number of training points increases from  $10^5$  to  $2 \times 10^5$  for ADDA-KR, ADDA-HH, Uniform-KR and Uniform-HH.



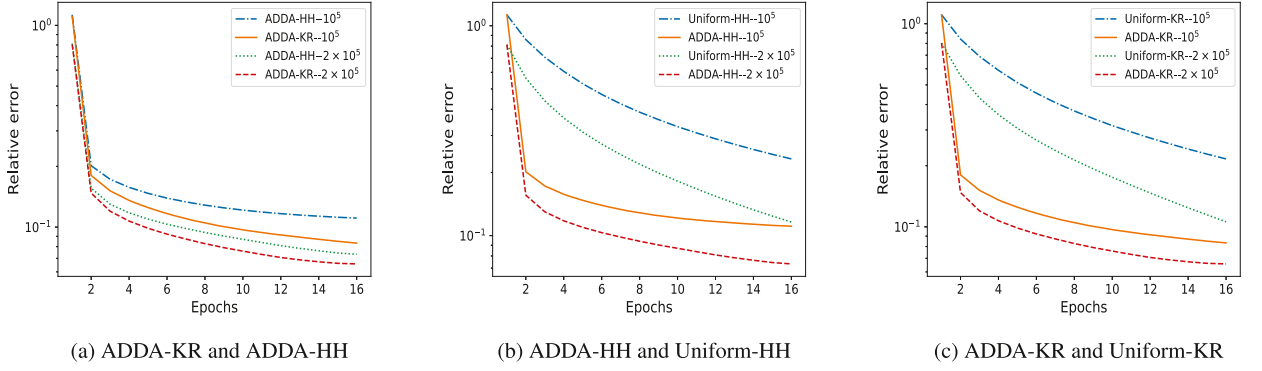


Fig. 12. Relative errors, four-dimensional test problem.

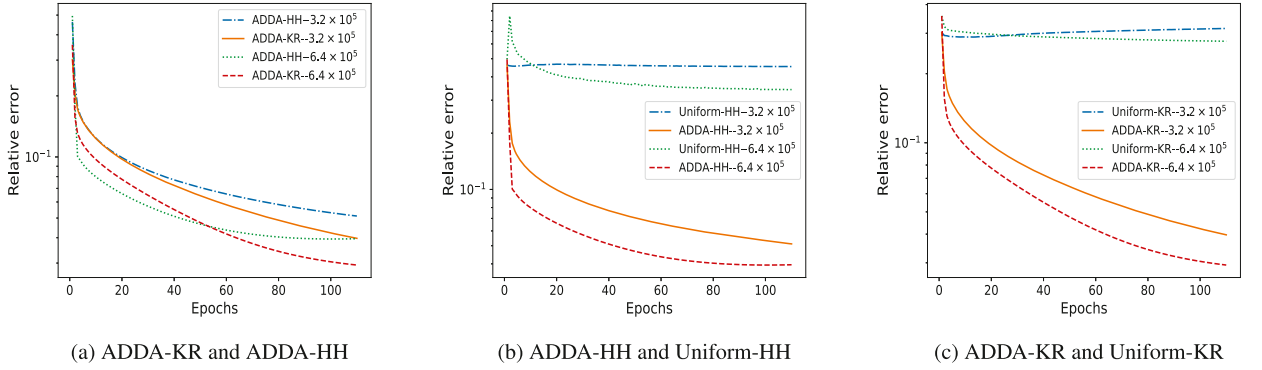


Fig. 13. Relative errors, eight-dimensional test problem.

Finally, we consider an eight-dimensional bimodal distribution. For this problem, we set

$$\mu_1 = [-1, -1, -0.3, -0.3, -0.4, -0.4, -1.6, -1.6]^T, \mu_2 = [2, 2, 0.6, 0.6, 0.8, 0.8, 2.3, 2.3]^T$$

$$\tilde{\Sigma}_1 = \begin{bmatrix} \Sigma_1 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 0.6\Sigma_1 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & 0.8\Sigma_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & 1.2\Sigma_1 \end{bmatrix}, \tilde{\Sigma}_2 = \begin{bmatrix} \Sigma_2 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 0.6\Sigma_2 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & 0.8\Sigma_2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & 1.2\Sigma_2 \end{bmatrix}, \quad (44)$$

where  $\Sigma_1$  and  $\Sigma_2$  are given in Eq. (41), and  $\tilde{\Sigma}_1$  and  $\tilde{\Sigma}_2$  are the covariance matrices of  $p_1$  and  $p_2$  for this test problem.

Again, we generate the initial parameters  $\Theta_0^{(0)}$  with Glorot Gaussian initialization, and then construct the initial KRnet  $p_X^{(0)}(\mathbf{x}; \Theta_0^{(0)})$ . The number of epochs is set to  $N_e = 1$ , and the maximum number of adaptivity iterations conducted for this problem is set to  $N_{\text{adaptive}} = 120$ . The learning rate for Adam optimizer is set to  $\eta = 0.0001$ , and the batch size is set to  $m = 4000$ . The initial training set  $C_0$  is generated through the uniform distribution with range  $[-6, 6]^8$ , and two cases of the collocation sample size are considered: one is  $3.2 \times 10^5$  and the other is  $6.4 \times 10^5$ . In addition, we take  $L = 10$  affine coupling layers for KR, and  $L = 20$  for real NVP. The architecture of  $\text{NN}_{[ij]}$  is the same as that shown in Fig. 2 with  $w = 160$ . For KRnet, we set  $K = 3$ . The rotation layer and the nonlinear layer are turned on. We again compute the relative error Eq. (27) using  $3.2 \times 10^5$  validation samples drawn from the exact solution.

Fig. 13 shows the relative error between  $p(\mathbf{x})$  and  $p_X(\mathbf{x}; \Theta)$  for ADDA-KR and ADDA-HH. From Fig. 13(a), it can be seen that the relative error of ADDA-KR is smaller than that of ADDA-HH, when the number of epochs is larger than 60. From Fig. 13(b) and Fig. 13(c), as the number of epochs increases, the relative errors of ADDA-KR and ADDA-HH decrease quickly, while the relative errors of the uniform sampling strategies (Uniform-KR and Uniform-HH) decrease slowly. In addition, it can be seen that the relative error decreases as the number of collocation points increases from  $3.2 \times 10^5$  to  $6.4 \times 10^5$  for ADDA-KR, ADDA-HH, Uniform-KR and Uniform-HH.

## 6. Conclusions

Conducting adaptivity is of fundamental importance for the efficient approximation of high-dimensional Fokker-Planck equations. With a focus on deep learning methods, we have developed an adaptive deep density approximation strategy

based on KRnet (ADDA-KR) in this work. Our KRnet, which is built on a block-triangular structure inspired by the Knothe-Rosenblatt rearrangement, gives an explicit family of probability density functions, which can serve as solution candidates of the Fokker-Planck equation. We also showed that KRnet is effective for estimating high-dimensional density functions in general. The fact that KRnet can efficiently generate samples integrates the two main steps in our ADDA-KR strategy to achieve efficient iterations: train KRnet for the Fokker-Planck equation with current collocation points, and generate new collocation points using the KRnet for the next iteration. Compared to real NVP, which is a widely used generative model, numerical results show that our ADDA-KR gives much more accurate numerical solutions for the Fokker-Planck equation. ADDA-KR in general works very well for Fokker-Planck equations with dimension of  $O(10)$ . For higher-dimensional cases, the sparsity of high-dimensional data will induce more severe difficulties, where we may need to consider dimension reduction to adapt more problem properties into the algorithm.

### CRediT authorship contribution statement

**Kejun Tang:** Programming, Methodology, Writing-Original draft preparation. **Xiaoliang Wan:** Conceptualization, Methodology, Programming, Writing. **Qifeng Liao:** Conceptualization, Methodology, Writing-Reviewing and editing.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Acknowledgements

K. Tang and Q. Liao are supported by the National Natural Science Foundation of China (No. 12071291) and the Science and Technology Commission of Shanghai Municipality (No. 20JC1414300), and X. Wan's work was supported by the National Science Foundation under grant DMS-1913163.

### References

- [1] H. Risken, Fokker-Planck-Kolmogorov Equation, Springer, 1984.
- [2] S. Jin, B. Yan, A class of asymptotic-preserving schemes for the Fokker-Planck-Landau equation, *J. Comput. Phys.* 230 (2011) 6420–6437.
- [3] Y. Li, A data-driven method for the steady state of randomly perturbed dynamics, *Commun. Math. Sci.* 17 (2019) 1045–1059.
- [4] B. Spencer, L. Bergman, On the numerical solution of the Fokker-Planck equation for nonlinear stochastic systems, *Nonlinear Dyn.* 4 (4) (1993) 357–372.
- [5] H.C. Elman, D.J. Silvester, A.J. Wathen, *Finite Elements and Fast Iterative Solvers: With Applications in Incompressible Fluid Dynamics*, Oxford University Press, USA, 2014.
- [6] M. Dobson, Y. Li, J. Zhai, An efficient data-driven solver for Fokker-Planck equations: algorithm and analysis, arXiv:1906.02600, 2019.
- [7] X. Chen, L. Yang, J. Duan, G.E. Karniadakis, Solving inverse stochastic problems from discrete particle observations using the Fokker-Planck equation and physics-informed neural networks, *SIAM J. Sci. Comput.* 43 (3) (2021) B811–B830.
- [8] W. E, A proposal on machine learning via dynamical systems, *Commun. Math. Stat.* 5 (1) (2017) 1–11.
- [9] W. E, B. Yu, The deep Ritz method: a deep learning-based numerical algorithm for solving variational problems, *Commun. Math. Stat.* 6 (1) (2018) 1–12.
- [10] M. Raissi, P. Perdikaris, G.E. Karniadakis, Physics informed deep learning (part I): data-driven solutions of nonlinear partial differential equations, arXiv:1711.10561, 2017.
- [11] M. Raissi, P. Perdikaris, G.E. Karniadakis, Physics informed deep learning (part II): data-driven discovery of nonlinear partial differential equations, arXiv:1711.10566, 2017.
- [12] M. Raissi, P. Perdikaris, G.E. Karniadakis, Physics-informed neural networks: a deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, *J. Comput. Phys.* 378 (2019) 686–707.
- [13] G. Pang, L. Lu, G.E. Karniadakis, fPINNs: fractional physics-informed neural networks, *SIAM J. Sci. Comput.* 41 (4) (2019) A2603–A2626.
- [14] J. Sirignano, K. Spiliopoulos, DGM: a deep learning algorithm for solving partial differential equations, *J. Comput. Phys.* 375 (2018) 1339–1364.
- [15] Y. Zhu, N. Zabaras, Bayesian deep convolutional encoder-decoder networks for surrogate modeling and uncertainty quantification, *J. Comput. Phys.* 366 (2018) 415–447.
- [16] Y. Zhu, N. Zabaras, P.-S. Koutsourelakis, P. Perdikaris, Physics-constrained deep learning for high-dimensional surrogate modeling and uncertainty quantification without labeled data, *J. Comput. Phys.* 394 (2019) 56–81.
- [17] K. Wu, D. Xiu, Numerical aspects for approximating governing equations using data, *J. Comput. Phys.* 384 (2019) 200–221.
- [18] K. Wu, T. Qin, D. Xiu, Structure-preserving method for reconstructing unknown Hamiltonian systems from trajectory data, *SIAM J. Sci. Comput.* 42 (2020) A3704–A3729.
- [19] K. Li, K. Tang, T. Wu, Q. Liao, D3M: a deep domain decomposition method for partial differential equations, *IEEE Access* 8 (2020) 5283–5294.
- [20] A.D. Jagtap, E. Kharazmi, G.E. Karniadakis, Conservative physics-informed neural networks on discrete domains for conservation laws: applications to forward and inverse problems, *Comput. Methods Appl. Mech. Eng.* 365 (2020) 113028.
- [21] S. Dong, Z. Li, Local extreme learning machines and domain decomposition for solving linear and nonlinear partial differential equations, *Comput. Methods Appl. Mech. Eng.* 387 (2021) 114129.
- [22] W. Li, X. Xiang, Y. Xu, Deep domain decomposition method: elliptic problems, in: J. Lu, R. Ward (Eds.), *Proceedings of the First Mathematical and Scientific Machine Learning Conference*, in: *Proceedings of Machine Learning Research*, PMLR, vol. 107, Princeton University, Princeton, NJ, USA, 2020, pp. 269–286.
- [23] A. Heinlein, A. Klawonn, M. Lanser, J. Weber, Combining machine learning and domain decomposition methods—a review, Technical report, Universität zu Köln, October 2020.
- [24] E. Kharazmi, Z. Zhang, G.E. Karniadakis, hp-VPINNs: variational physics-informed neural networks with domain decomposition, *Comput. Methods Appl. Mech. Eng.* 374 (2021) 113547.

- [25] H. Sheng, C. Yang, PFNN: a penalty-free neural network method for solving a class of second-order boundary-value problems on complex geometries, *J. Comput. Phys.* (2020) 110085.
- [26] H. Gao, L. Sun, J.-X. Wang, Phygeonet: physics-informed geometry-adaptive convolutional neural networks for solving parameterized steady-state PDEs on irregular domain, *J. Comput. Phys.* 428 (2021) 110079.
- [27] D. Xiu, *Numerical Methods for Stochastic Computations: A Spectral Method Approach*, Princeton University Press, 2010.
- [28] D. Xiu, J.S. Hesthaven, High-order collocation methods for differential equations with random inputs, *SIAM J. Sci. Comput.* 27 (3) (2005) 1118–1139.
- [29] I. Babuška, F. Nobile, R. Tempone, A stochastic collocation method for elliptic partial differential equations with random input data, *SIAM J. Numer. Anal.* 45 (3) (2007) 1005–1034.
- [30] J. Foo, X. Wan, G.E. Karniadakis, The multi-element probabilistic collocation method (ME-PCM): error analysis and applications, *J. Comput. Phys.* 227 (22) (2008) 9572–9595.
- [31] X. Ma, N. Zabarar, An adaptive hierarchical sparse grid collocation algorithm for the solution of stochastic differential equations, *J. Comput. Phys.* 228 (8) (2009) 3084–3113.
- [32] A. Narayan, D. Xiu, Stochastic collocation methods on unstructured grids in high dimensions via interpolation, *SIAM J. Sci. Comput.* 34 (3) (2012) A1729–A1752.
- [33] H. Lei, X. Yang, B. Zheng, G. Lin, N.A. Baker, Constructing surrogate models of complex systems with enhanced sparsity: quantifying the influence of conformational uncertainty in biomolecular solvation, *Multiscale Model. Simul.* 13 (4) (2015) 1327–1353.
- [34] P. Ren, Y. Xiao, X. Cheang, P.-Y. Huang, Z. Li, X. Chen, X. Wang, A survey of deep active learning, arXiv:2009.00236v1, 2020.
- [35] R. Cang, H. Yao, Y. Ren, One-shot generation of near-optimal topology through theory-driven machine learning, *Comput. Aided Des.* 109 (2019) 12–21.
- [36] R. Vilalta, Y. Drissi, A perspective view and survey of meta-learning, *Artif. Intell. Rev.* 18 (2001) 77–95.
- [37] D.W. Scott, *Multivariate Density Estimation: Theory, Practice, and Visualization*, John Wiley & Sons, 2015.
- [38] L. Dinh, J. Sohl-Dickstein, S. Bengio, Density estimation using real NVP, arXiv:1605.08803, 2016.
- [39] D.P. Kingma, P. Dhariwal, Glow, Generative flow with invertible  $1 \times 1$  convolutions, in: *Advances in Neural Information Processing Systems*, 2018, pp. 10215–10224.
- [40] L. Zhang, W. E, L. Wang, Monge-Ampère flow for generative modeling, arXiv:1809.10188, 2018.
- [41] T.Q. Chen, Y. Rubanova, J. Bettencourt, D.K. Duvenaud, Neural ordinary differential equations, in: *Advances in Neural Information Processing Systems*, 2018, pp. 6571–6583.
- [42] G. Carlier, A. Galichon, F. Santambrogio, From Knothe’s transport to Brenier’s map and a continuation method for optimal transport, *SIAM J. Math. Anal.* 41 (6) (2010) 2554–2576.
- [43] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: *Advances in Neural Information Processing Systems*, 2014, pp. 2672–2680.
- [44] D.P. Kingma, M. Welling, Auto-encoding variational Bayes, *Stat* 1050 (2014) 1.
- [45] X. Wan, S. Wei, Coupling the reduced-order model and the generative model for an importance sampling estimator, *J. Comput. Phys.* 408 (2020) 109281.
- [46] K. Tang, X. Wan, Q. Liao, Deep density estimation via invertible block-triangular mapping, *Theor. Appl. Mech. Lett.* 10 (2020) 143.
- [47] G.A. Pavliotis, *Stochastic Processes and Applications: Diffusion Processes, the Fokker-Planck and Langevin Equations*, vol. 60, Springer, 2014.
- [48] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [49] S. Ioffe, C. Szegedy, Batch normalization: accelerating deep network training by reducing internal covariate shift, arXiv:1502.03167, 2015.
- [50] X. Glorot, A. Bordes, Y. Bengio, Deep sparse rectifier neural networks, in: *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 2011, pp. 315–323.
- [51] D.P. Kingma, J. Ba, Adam, A method for stochastic optimization, arXiv:1412.6980, 2014.
- [52] G. Cybenko, Approximation by superpositions of a sigmoidal function, *Math. Control Signals Syst.* 2 (4) (1989) 303–314.
- [53] M. Leshno, V.Y. Lin, A. Pinkus, S. Schocken, Multilayer feedforward networks with a nonpolynomial activation function can approximate any function, *Neural Netw.* 6 (6) (1993) 861–867.
- [54] Z. Lu, H. Pu, F. Wang, Z. Hu, L. Wang, The expressive power of neural networks: a view from the width, in: *Advances in Neural Information Processing Systems*, 2017, pp. 6231–6239.
- [55] L. Bottou, F.E. Curtis, J. Nocedal, Optimization methods for large-scale machine learning, *SIAM Rev.* 60 (2) (2018) 223–311.
- [56] X. Glorot, Y. Bengio, Understanding the difficulty of training deep feedforward neural networks, in: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 2010, pp. 249–256.