
Neural tangent kernel analysis of PINN for advection-diffusion equation

M. H. Saadat, B. Gjorgiev, L. Das and G. Sansavini

Reliability and Risk Engineering Laboratory, Institute of Energy and Process Engineering,
Department of Mechanical and Process Engineering, ETH Zurich, 8092 Zurich, Switzerland

gsansavini@ethz.ch

Abstract

Physics-informed neural networks (PINNs) numerically approximate the solution of a partial differential equation (PDE) by incorporating the residual of the PDE along with its initial/boundary conditions into the loss function. In spite of their partial success, PINNs are known to struggle even in simple cases where the closed-form analytical solution is available. In order to better understand the learning mechanism of PINNs, this work focuses on a systematic analysis of PINNs for the linear advection-diffusion equation (LAD) using the Neural Tangent Kernel (NTK) theory. Thanks to the NTK analysis, the effects of the advection speed/diffusion parameter on the training dynamics of PINNs are studied and clarified. We show that the training difficulty of PINNs is a result of 1) the so-called 'spectral bias', which leads to difficulty in learning high-frequency behaviours; and 2) convergence rate disparity between different loss components that results in training failure. The latter occurs even in the cases where the solution of the underlying PDE does not exhibit high-frequency behaviour. Furthermore, we observe that this training difficulty manifests itself, to some extent, differently in advection-dominated and diffusion-dominated regimes. Different strategies to address these issues are also discussed. In particular, it is demonstrated that periodic activation functions can be used to partly resolve the spectral bias issue.

1 Introduction

Developing accurate and efficient numerical schemes for solving partial differential equations (PDEs) is an active research field with enormous applications in different scientific and engineering disciplines. Classical numerical methods, such as finite difference, finite volume or finite element schemes, rely on discretizing the governing equations over an underlying grid points and marching forward in time, and have been successfully applied to a wide range of applications from fluid dynamics and quantum mechanics to neuroscience and finance.

In addition to above mentioned established numerical methods, neural networks (NNs) have recently received attention as an alternative approach to conventional PDE solvers [1, 2]. In particular, the first realization of the physics-informed neural network (PINN) [3] showed promising results for solving various linear and non-linear PDEs. PINN approximates the solution of a PDE by turning it into an unsupervised optimization problem with the constraints being the residual of the PDE under consideration and its initial and boundary conditions. The major attractiveness of PINN, however, lies in the computation of derivatives which is done through the automatic differentiation (AD) technique during back-propagation. This results in a mesh-free algorithm that opens up the possibility of accurately solving complex and computationally demanding PDEs including high-dimensional PDEs [4] or Integral-differential equations [5]. The pioneering work of [3] has inspired various improved

PINN models which extend the operating domain of the original model [3] and are to some extent capable of solving PDEs in relevant practical applications [6, 7, 8, 9, 10, 11, 12, 13]. See [14] for a full review of PINNs development.

The original realization of PINN [3], however, struggles in dealing with PDEs which exhibit high-frequency or multi-scale behaviour, even if the underlying PDE has a closed-form analytical solution. It has been shown that in those scenarios the PINN fails to learn the correct solution, no matter how large the architecture is [15]. Different studies have attempted to identify the failure modes of PINN and to understand its learning mechanism [15, 16, 17, 18]. In particular, PINN failure has been attributed to the complexity of the loss function in challenging physical regimes [15] or to the unbalanced gradients of different loss terms during training, which results in numerical stiffness and unbalanced back-propagation via Gradient Descent (GD) [16]. Yet, the most rigorous mathematical explanation of the PINNs failure has been presented in [19] based on the so-called Neural Tangent Kernel (NTK) theory [20].

NTK is a recent theory which studies fully-connected NNs in the infinite-width limit, i.e, in the limit as the number of neurons in the hidden layers tend to infinity [20, 21]. It has been shown that in this limit, NNs behave like a linear model obtained from the first-order Taylor expansion of the NN around its initial parameters [20, 21, 22]. In other words, according to NTK theory, fully-connected infinite width NNs trained by the GD algorithm behave like a kernel regression model with a deterministic kernel called NTK [20, 21, 22]. The NTK theory has clarified different features of the NN training, including the success of over-parameterized networks. Another feature worth mentioning, is the so-called 'spectral-bias' or 'F-principle' [23, 24], which describe the tendency of NNs to be biased towards low-frequency solutions.

Wang et al. [19] derived the NTK for PINNs and proved its convergence to a deterministic kernel, provided that the network is wide and the learning rate is small enough. Based on that, they showed that, similar to NNs, PINNs also suffer from the 'spectral bias' [19, 25]. Furthermore, they demonstrated that the convergence rate discrepancy between different components in the loss function can lead to PINN failure [19, 25].

Regardless of limitations of the NTK theory, the framework presented by [19] provides an exceptional step forward towards analysing and understanding the behaviour of PINNs in dealing with different class of PDEs. Building on this work [19], this paper aims at (1) apply the NTK theory to PINNs for the linear advection-diffusion (LAD) equation, (2) clarify the exact mechanisms behind PINNs failure in advection-dominated and diffusion-dominated regimes, and (3) exploring the effectiveness of different training strategies in the view of NTK theory, and particularly the role of activation function on training dynamics. All theoretical results are supplemented by numerical experiments. We choose to study the advection-diffusion equation, as it describes two fundamental physical processes with enormous applications from fluid dynamics to biological systems and financial mathematics. In addition, there is a closed-form analytical solution, as well as a mature literature on different numerical methods for dealing with LAD equations. As a result, studying PINNs for this class of PDEs would provide valuable insights into application of PINNs for more challenging and complicated PDEs.

The remainder of the paper is organized as follows: for the sake of completeness, Section 2 briefly presents the main steps of deriving NTK for PINNs. Section 3, studies the PINN for the LAD and clarifies the impact of the PDE parameters on the training dynamics of PINN through NTK analysis backed by numerical experiments. Furthermore, the effectiveness of different improved learning strategies are investigated. Particularly, the role of activation function is studied and analyzed. Finally, Section 4 discusses our findings and provides conclusions.

2 Methodology

In this section, we briefly summarize the neural tangent kernel derivation of PINNs for the linear advection-diffusion equations following [19]. Interested readers are referred to [19, 25] for the full derivation and proofs of the NTK theory for PINNs.

Let us, for simplicity, consider the following scalar advection-diffusion equation,

$$\begin{aligned} u_t + au_x - \epsilon u_{xx} &= 0, x \in \Omega, \\ u(x, 0) &= f(x), x \in \Omega, \\ u(x, t) &= g(x), x \in \partial\Omega, \end{aligned} \quad (1)$$

where $u(x, t) : \Omega \times [0, T] \rightarrow \mathbb{R}$ denotes the unknown solution in one-dimensional physical space Ω , a is the advection speed, and ϵ represents the diffusion parameter. For the sake of notational simplicity, we can consider time t as an additional dimension in x , i.e. $\mathbf{x} = (x, t)$.

Following the original formulation of PINN [3], we approximate the solution $u(x, t) = u(\mathbf{x})$ with a deep neural network $u(\mathbf{x}, \boldsymbol{\theta})$, where $\boldsymbol{\theta} \in \mathbb{R}^{n_{par}}$ is a collection of all network parameters including weights and biases, and by considering the following loss function,

$$\mathcal{L} = \mathcal{L}_r + \mathcal{L}_b, \quad (2)$$

where

$$\mathcal{L}_r = \frac{1}{2} \sum_{i=1}^{N_r} \|u_t(\mathbf{x}_r^i, \boldsymbol{\theta}) + au_x(\mathbf{x}_r^i, \boldsymbol{\theta}) - \epsilon u_{xx}(\mathbf{x}_r^i, \boldsymbol{\theta})\|^2 = \frac{1}{2} \sum_{i=1}^{N_r} \|\mathcal{N}u(\mathbf{x}_r^i, \boldsymbol{\theta})\|^2, \quad (3)$$

$$\mathcal{L}_b = \frac{1}{2} \sum_{i=1}^{N_b} \|u(\mathbf{x}_b^i, \boldsymbol{\theta}) - g(\mathbf{x}_b^i)\|^2. \quad (4)$$

Here, \mathcal{N} represents the differential operator on the LHS of 1, \mathcal{L}_r denotes the loss over residual of the PDE (1) at N_r collocation points sampled from inside of the domain $\{\mathbf{x}_r^i\}_{i=1}^{N_r}$, and \mathcal{L}_b is the loss associated with initial/boundary conditions evaluated at N_b initial/boundary points, $\{\mathbf{x}_b^i\}_{i=1}^{N_b}$.

Learning the network parameters $\boldsymbol{\theta}$, implies minimization of the loss function (2) through the gradient descent (GD) algorithm. Assuming an infinitesimal learning rate, GD can be written in terms of the so-called continuous time gradient flow differential equation [19],

$$\frac{d\boldsymbol{\theta}}{dt} = -\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}). \quad (5)$$

Applying the gradient leads to,

$$\frac{d\boldsymbol{\theta}}{dt} = - \left[\sum_{i=1}^{N_r} (\mathcal{N}u(\mathbf{x}_r^i, \boldsymbol{\theta}(t))) \frac{\partial \mathcal{N}u(\mathbf{x}_r^i, \boldsymbol{\theta}(t))}{\partial \boldsymbol{\theta}} + \sum_{i=1}^{N_b} (u(\mathbf{x}_b^i, \boldsymbol{\theta}(t)) - g(\mathbf{x}_b^i)) \frac{\partial u(\mathbf{x}_b^i, \boldsymbol{\theta}(t))}{\partial \boldsymbol{\theta}} \right]. \quad (6)$$

Using the chain rule, we can write for $0 \leq k \leq N_r$ and $0 \leq j \leq N_b$,

$$\frac{d\mathcal{N}u(\mathbf{x}_r^k, \boldsymbol{\theta}(t))}{dt} = \frac{d\mathcal{N}u(\mathbf{x}_r^k, \boldsymbol{\theta}(t))^T}{d\boldsymbol{\theta}} \cdot \frac{d\boldsymbol{\theta}}{dt}, \quad (7)$$

$$\frac{du(\mathbf{x}_b^j, \boldsymbol{\theta}(t))}{dt} = \frac{du(\mathbf{x}_b^j, \boldsymbol{\theta}(t))^T}{d\boldsymbol{\theta}} \cdot \frac{d\boldsymbol{\theta}}{dt}. \quad (8)$$

Substituting (6) into (7) results in,

$$\begin{aligned} \frac{d\mathcal{N}u(\mathbf{x}_r^k, \boldsymbol{\theta})}{dt} &= - \left[\sum_{i=1}^{N_r} (\mathcal{N}u(\mathbf{x}_r^i, \boldsymbol{\theta})) \left\langle \frac{d\mathcal{N}u(\mathbf{x}_r^k, \boldsymbol{\theta})}{d\boldsymbol{\theta}}, \frac{d\mathcal{N}u(\mathbf{x}_r^i, \boldsymbol{\theta})}{d\boldsymbol{\theta}} \right\rangle \right. \\ &\quad \left. + \sum_{i=1}^{N_b} (u(\mathbf{x}_b^i, \boldsymbol{\theta}) - g(\mathbf{x}_b^i)) \left\langle \frac{d\mathcal{N}u(\mathbf{x}_r^k, \boldsymbol{\theta})}{d\boldsymbol{\theta}}, \frac{du(\mathbf{x}_b^i, \boldsymbol{\theta})}{d\boldsymbol{\theta}} \right\rangle \right], \end{aligned} \quad (9)$$

and similarly for (8),

$$\begin{aligned} \frac{du(\mathbf{x}_b^j, \boldsymbol{\theta})}{dt} &= - \left[\sum_{i=1}^{N_r} (\mathcal{N}u(\mathbf{x}_r^i, \boldsymbol{\theta})) \left\langle \frac{du(\mathbf{x}_b^j, \boldsymbol{\theta})}{d\boldsymbol{\theta}}, \frac{d\mathcal{N}u(\mathbf{x}_r^i, \boldsymbol{\theta})}{d\boldsymbol{\theta}} \right\rangle \right. \\ &\quad \left. + \sum_{i=1}^{N_b} (u(\mathbf{x}_b^i, \boldsymbol{\theta}) - g(\mathbf{x}_b^i)) \left\langle \frac{du(\mathbf{x}_b^j, \boldsymbol{\theta})}{d\boldsymbol{\theta}}, \frac{du(\mathbf{x}_b^i, \boldsymbol{\theta})}{d\boldsymbol{\theta}} \right\rangle \right], \end{aligned} \quad (10)$$

where the part within $\langle \rangle$ indicates the inner product over all neural network parameters $\boldsymbol{\theta}$. By defining the following kernels,

$$(\mathbf{K}_u)_{ij}(t) = \left\langle \frac{du(\mathbf{x}_b^i, \boldsymbol{\theta}(t))}{d\boldsymbol{\theta}}, \frac{du(\mathbf{x}_b^j, \boldsymbol{\theta}(t))}{d\boldsymbol{\theta}} \right\rangle, \quad (11)$$

$$(\mathbf{K}_{ur})_{ij}(t) = \left\langle \frac{du(\mathbf{x}_b^i, \boldsymbol{\theta}(t))}{d\boldsymbol{\theta}}, \frac{d\mathcal{N}u(\mathbf{x}_r^j, \boldsymbol{\theta}(t))}{d\boldsymbol{\theta}} \right\rangle, \quad (12)$$

$$(\mathbf{K}_r)_{ij}(t) = \left\langle \frac{d\mathcal{N}u(\mathbf{x}_r^i, \boldsymbol{\theta}(t))}{d\boldsymbol{\theta}}, \frac{d\mathcal{N}u(\mathbf{x}_r^j, \boldsymbol{\theta}(t))}{d\boldsymbol{\theta}} \right\rangle, \quad (13)$$

we can rewrite Eqs. (9) and (10) in a compact form as,

$$\begin{bmatrix} \frac{du(\mathbf{x}_b, \boldsymbol{\theta}(t))}{dt} \\ \frac{d\mathcal{N}u(\mathbf{x}_r, \boldsymbol{\theta}(t))}{dt} \end{bmatrix} = -\mathbf{K}(t) \cdot \begin{bmatrix} u(\mathbf{x}_b, \boldsymbol{\theta}(t)) - g(\mathbf{x}_b) \\ \mathcal{N}u(\mathbf{x}_r, \boldsymbol{\theta}(t)) \end{bmatrix}, \quad (14)$$

where

$$\mathbf{K}(t) = \begin{bmatrix} \mathbf{K}_u(t) & \mathbf{K}_{ur}(t) \\ \mathbf{K}_{ur}^T(t) & \mathbf{K}_r(t) \end{bmatrix}, \quad (15)$$

is called the neural tangent kernel (NTK) of our PINN. It can be verified that for an over-parameterized PINN, the NTK is a positive semi-definite (PSD) matrix [19].

According to the neural tangent kernel (NTK) theory, in the limit where the width of the deep neural network goes to infinity (and also under some other conditions [19]), the kernel $\mathbf{K}(t)$ converges to a static kernel \mathbf{K}^* and remains almost constant during the training. In most cases, the kernel can be well approximated with its value at the initialization, i.e.,

$$\mathbf{K}(t) \approx \mathbf{K}^* \approx \mathbf{K}(0). \quad (16)$$

Therefore, learning dynamics can be simplified into the following system of ordinary differential equations,

$$\begin{bmatrix} \frac{du(\mathbf{x}_b, \boldsymbol{\theta})}{dt} \\ \frac{d\mathcal{N}u(\mathbf{x}_r, \boldsymbol{\theta})}{dt} \end{bmatrix} \approx -\mathbf{K}(0) \cdot \begin{bmatrix} u(\mathbf{x}_b, \boldsymbol{\theta}) - g(\mathbf{x}_b) \\ \mathcal{N}u(\mathbf{x}_r, \boldsymbol{\theta}) \end{bmatrix}, \quad (17)$$

which implies that,

$$\begin{bmatrix} u(\mathbf{x}_b, \boldsymbol{\theta}(t)) - g(\mathbf{x}_b) \\ \mathcal{N}u(\mathbf{x}_r, \boldsymbol{\theta}(t)) \end{bmatrix} \approx e^{-\mathbf{K}(0)t} \cdot \begin{bmatrix} u(\mathbf{x}_b, \boldsymbol{\theta}(0)) - g(\mathbf{x}_b) \\ \mathcal{N}u(\mathbf{x}_r, \boldsymbol{\theta}(0)) \end{bmatrix}. \quad (18)$$

Spectral decomposition of $\mathbf{K}(0)$ as $\mathbf{K}(0) = P^T \Lambda P$, with P being an orthogonal matrix whose columns are eigenvectors of $\mathbf{K}(0)$ and Λ being a diagonal matrix with diagonal entries λ_i corresponding to eigenvalues of $\mathbf{K}(0)$, allows us to write the solution in the following form,

$$\begin{bmatrix} u(\mathbf{x}_b, \boldsymbol{\theta}(t)) - g(\mathbf{x}_b) \\ \mathcal{N}u(\mathbf{x}_r, \boldsymbol{\theta}(t)) \end{bmatrix} \approx P^T e^{-\Lambda t} P \cdot \begin{bmatrix} u(\mathbf{x}_b, \boldsymbol{\theta}(0)) - g(\mathbf{x}_b) \\ \mathcal{N}u(\mathbf{x}_r, \boldsymbol{\theta}(0)) \end{bmatrix}. \quad (19)$$

Some remarks about NTKs and the information they provide according to Eq.(19) are in order:

- Given that $\mathbf{K}(0)$ is a PSD matrix with non-negative eigenvalues $\lambda_i \geq 0$, the training loss decays monotonically with time t and converges to zero in the limit of $t \rightarrow \infty$.
- The i -th component of the training error decays at a rate proportional to $e^{-\lambda_i t}$. As a result, severe eigenvalue disparity between different loss terms would lead to domination of one term and, therefore, difficulty in training.
- It has been shown that larger eigenvalues correspond to low-frequency solutions. This would result in the 'spectral bias' of the PINNs, meaning that the network learns the low-frequency solutions faster.

In the next section, we shall illustrate the usefulness of NTKs in analyzing the behaviour of PINNs dealing with LAD equations.

3 Numerical experiments

In this section, through the use of the NTK, we study PINN for the linear advection-diffusion equation (1) and study the impact of the advection speed/diffusion parameter on the training dynamics of PINNs. Unless otherwise stated, we consider the computational domain as $(x, t) : [0, 1] \times [0, 1]$, the initial condition $u(x, 0) = \sin(2\pi x)$ and periodic boundary condition. The LAD has a closed form analytical solution given by,

$$u_{Exact} = \sin(2\pi(x - at)) e^{(-\epsilon a^2 t)}. \quad (20)$$

For simplicity, periodicity is imposed through Dirichlet boundary condition and according to the analytical solution (20). The network architecture used consists of 3 hidden layers with 200 neurons in each layer with \tanh as activation function. The network is initialized using the *Xavier* scheme, and *Adam* optimizer is used to minimize the loss function (2) with learning rate of 0.001. Furthermore, we generate $N_r = N_b = 300$ collocation points randomly sampled from the computational domain (x, t) .

3.1 Advection-dominated regime

We first investigate the pure linear advection equation ($\epsilon = 0$) at three different advection speeds, namely $a = 0.1, 1$ and 15 . We train the network for 80000 steps and monitor the eigenvalues of NTKs $\mathbf{K}_u(t)$ and $\mathbf{K}_r(t)$ over the training. Figure 1 illustrates the eigenvalues in descending order for the three different training steps. In all three cases, the figure shows that the eigenvalues decay rapidly to zero such that most of them, except the first few, are near zero. As mentioned in the introduction, the training error decays proportionally to the eigenvalue magnitude and larger eigenvalues correspond to low-frequency eigenvectors. This fact explains that, similar to conventional neural networks, PINNs also tend to learn low-frequency solutions faster and therefore, exhibit 'spectral bias' [19, 16]. In addition, Fig. 1 shows that all eigenvalues converge, which in turn means that the NTKs converge to deterministic kernels during the training. This is consistent with the NTK theory [16] discussed in Section 2.

Another key thing to note is that, for small values of the advection speed, there is a small disparity between eigenvalues and the predicted solution converges to the exact solution with good accuracy, as shown in Fig. 2. However, increasing the advection speed to $a = 15$ results in a severe disparity between \mathbf{K}_r and \mathbf{K}_u eigenvalues in a way that \mathbf{K}_r dominates \mathbf{K}_u (see Fig. 1). This results in faster decay of residual error \mathcal{L}_r , compared to initial/boundary condition error \mathcal{L}_b , as illustrated in Fig. 3. In spite of this fast convergence, the solution converges to a trivial non-physical solution after some time, as it can be clearly seen from Fig. 2. Domination of \mathbf{K}_r , indeed, makes the problem ill-posed, as it prevents the network to learn the initial/boundary conditions, leading to trivial constant solution.

To study the importance of properly learning the initial/boundary conditions and the effect of eigenvalue disparity on that, we repeat the training for the case with $a = 15$ using the adaptive weights training proposed by Wang et al. [19]. The idea of adaptive weights strategy is to assign proper weights to different loss components in order to reduce the eigenvalue disparity [19, 25]. The results illustrated in Figure 4 demonstrate that while this strategy is, to some extent effective, the predicted solution still converges to trivial solution, albeit at a later time. The failure of PINN, in this case after resolving the eigenvalue disparity issue, can be explained through the lens of the 'spectral bias', since by increasing the advection speed the solution exhibits multi-scale high-frequency behaviour in the space-time domain (see the exact solution in Fig. 2). To further investigate this issue, we repeat the training with same setup but on a smaller sub-domain of $(x, t) : [0, 1] \times [0, 0.5]$. This strategy effectively reduces the complexity of the underlying solution and results in robust training of PINN with accurate results. This confirms that spectral bias is an important source of training difficulty. Noteworthy to say that, this is the idea behind the sequence-to-sequence learning [15], or time adaptive strategies [26, 27], which splits the entire space-time domain into multiple subdomains, trains the PINN in sub-domains while marching forward in time. It is important to emphasise that all of these improved learning strategies come at the price of increasing computational and algorithmic complexity.

3.1.1 Effect of activation function

Most of the state-of-the-arts neural network architectures in the literature utilize the rectified linear unit (ReLU) activation function, which have shown success in various deep learning applications.

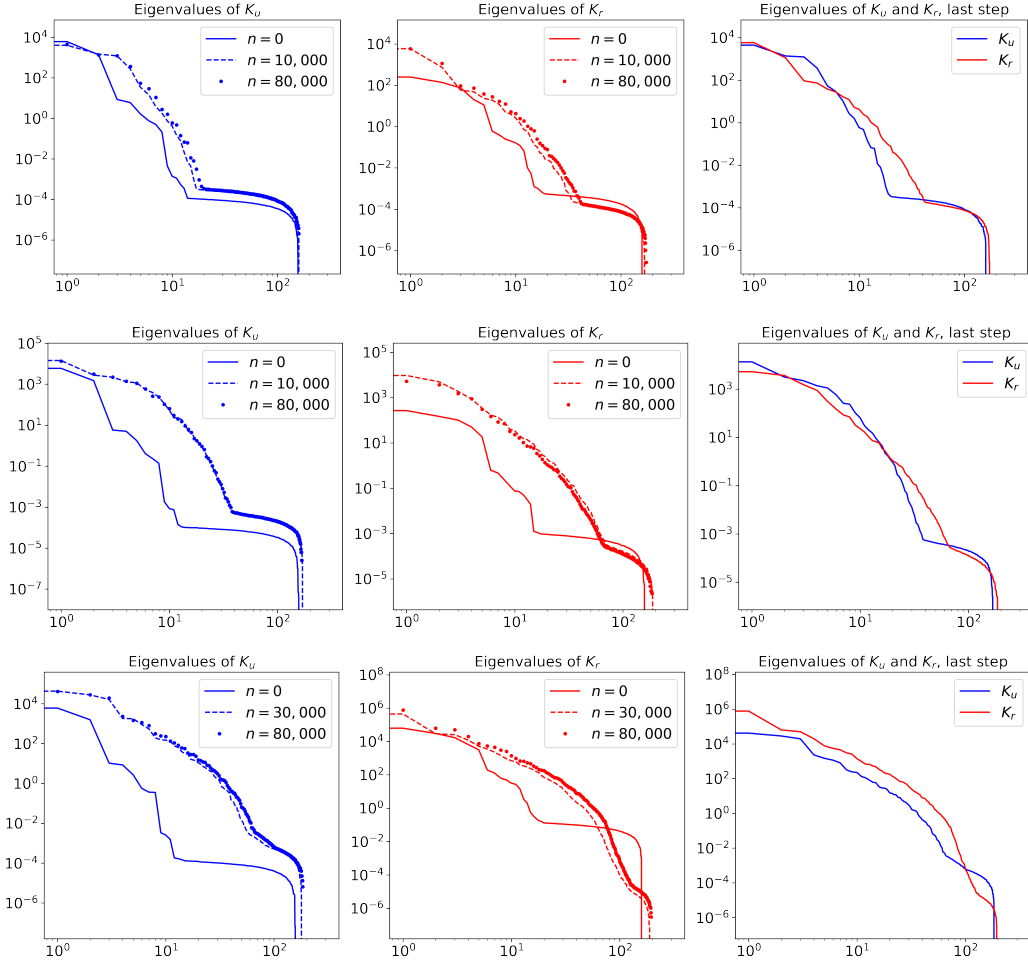


Figure 1: Evolution of NTK eigenvalues for the advection with $a = 0.1$ (top), $a = 1$ (middle) and $a = 15$ (bottom).

This is, however, not possible in PINNs, as they require sufficiently smooth activation functions, while ReLU is a piecewise linear function incapable of representing second-order derivative and beyond. As such, tanh is the most widely used activation function for PINN models, despite the fact that it suffers from the vanishing gradient problem and, therefore, is not suitable for deep network architectures.

Recently, there has been a growing interest in using periodic functions as an alternative to traditional activation functions [28, 29, 30, 31]. While the first use of periodic activation goes back to three decades ago [32], it was not until recently that neural networks based on periodic functions were used to solve challenging problems and shown to outperform traditional networks in some applications [29, 30, 31]. Noteworthy, is the so-called *SIREN* network which leverages sin activation function and has been successful in representing complex physical signals [30].

Motivated by this, here we explore the performance of sin activation function in comparison with the tanh activation used before by repeating the advection-dominated case with $a = 15$. Figure 5 shows that eigenvalues corresponding to sin activation decay slower compared to the ones corresponding to tanh. This means that sin networks might be able to learn higher frequency behaviours than tanh networks. This is confirmed in Fig. 6, where the predicted solution with sin network agrees well with the exact solution, while tanh network fails to learn the correct solution (see Fig. 4 top). Consequently, we can conclude that employing periodic activation functions in PINNs can be an

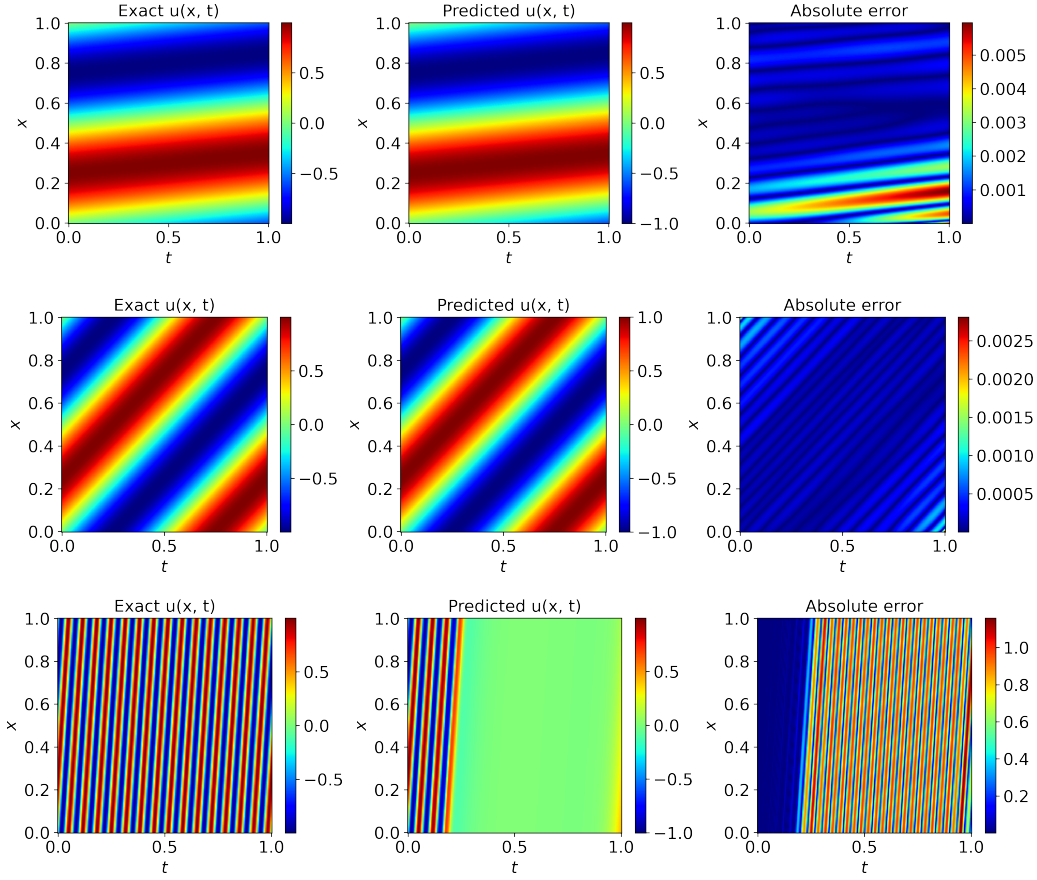


Figure 2: Snapshots of the predicted solution versus exact solution for the advection with $a = 0.1$ (top), $a = 1$ (middle) and $a = 15$ (bottom).

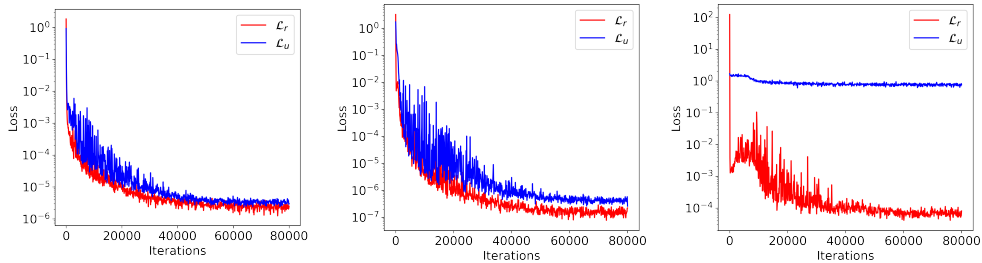


Figure 3: Time history of the different terms in the loss function for the advection with $a = 0.1$ (left), $a = 1$ (middle) and $a = 15$ (right).

effective approach to address the spectral bias issue without any additional computation overhead or major change in the algorithm (contrary to the previous strategies).

3.2 Diffusion-dominated regime

The present section focuses on clarifying the effect of diffusion parameter ϵ on the training dynamics of PINN. Similar to the previous section, the network is trained for 80000 steps and for two different parameter values, namely $\epsilon = 0.1$ and 10.

Figure 7 depicts the NTKs for three cases. Qualitatively speaking, NTKs follow similar trends observed in the previous section. As diffusion parameter increases to $\epsilon = 10$, the training becomes

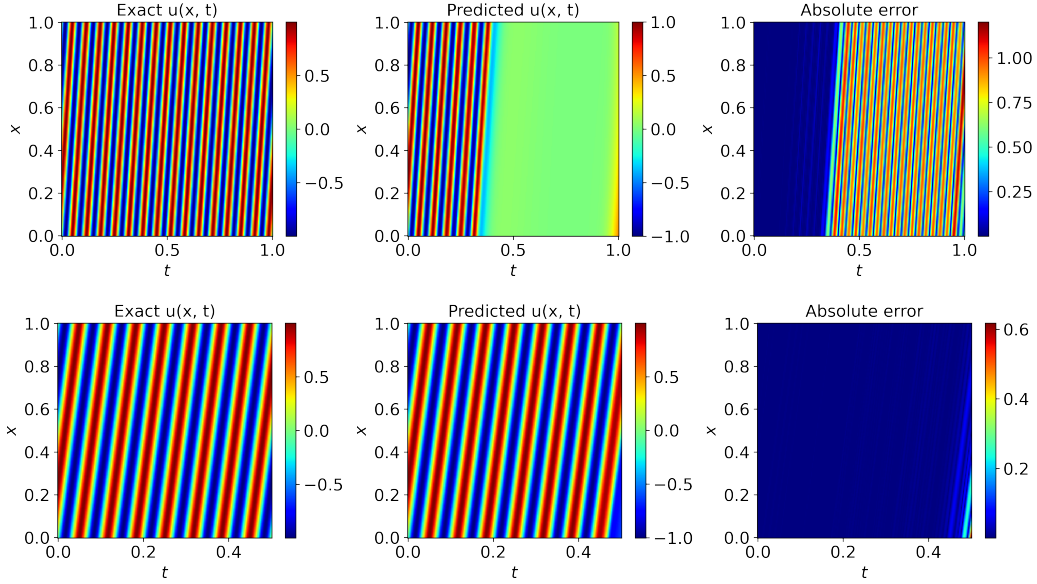


Figure 4: Snapshots of the predicted solution versus exact solution for the advection with $a = 15$ using adaptive weights strategy (top) and sub-domain learning on $t = [0, 0.5]$ (bottom).

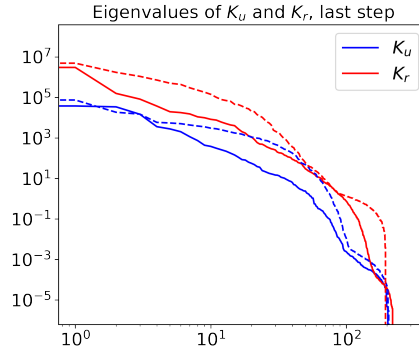


Figure 5: Comparison of NTK eigenvalues for the advection with $a = 15$. Solid lines: tanh activation; dashed lines: sin activation.

more difficult, which eventually leads to PINN collapse. However, Fig. 8 shows the underlying exact solution in the case of high diffusivity does not exhibit high-frequency behaviour, as the initial condition decays rapidly to trivial solution zero. This counter-intuitive failure of PINN in learning this relatively simple solution of the diffusion equation, could be explained by severe eigenvalue disparity that is observed for this case and is shown in Fig. 7. Similar to the previous case, the domination of the residual error term prevents the network from learning the initial/boundary conditions. Based on this observation and given the fact that the solution is more or less smooth and trivial in this case, we expect adaptive weights to be a more effective strategy than sub-domain learning or using sin. This is verified in Fig. 10, where the predicted solution with adaptive weights correctly recovers the exact solution with reasonable accuracy, while the predicted solution even on a very small time domain $t = [0, 0.1]$ does not show any significant improvement compared to the baseline PINN results shown in Fig. 8. This is in contrast to the advection regime, where sub-domain learning was an effective approach to resolve the training difficulty of PINN.

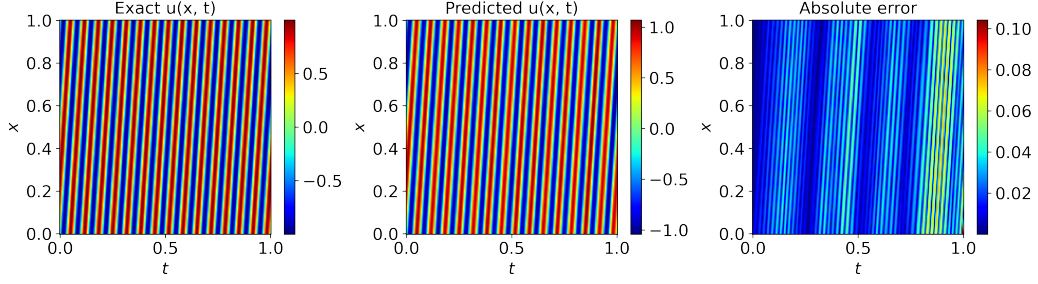


Figure 6: Snapshot of the predicted solution versus exact solution for the advection with $a = 15$ using sin activation function.

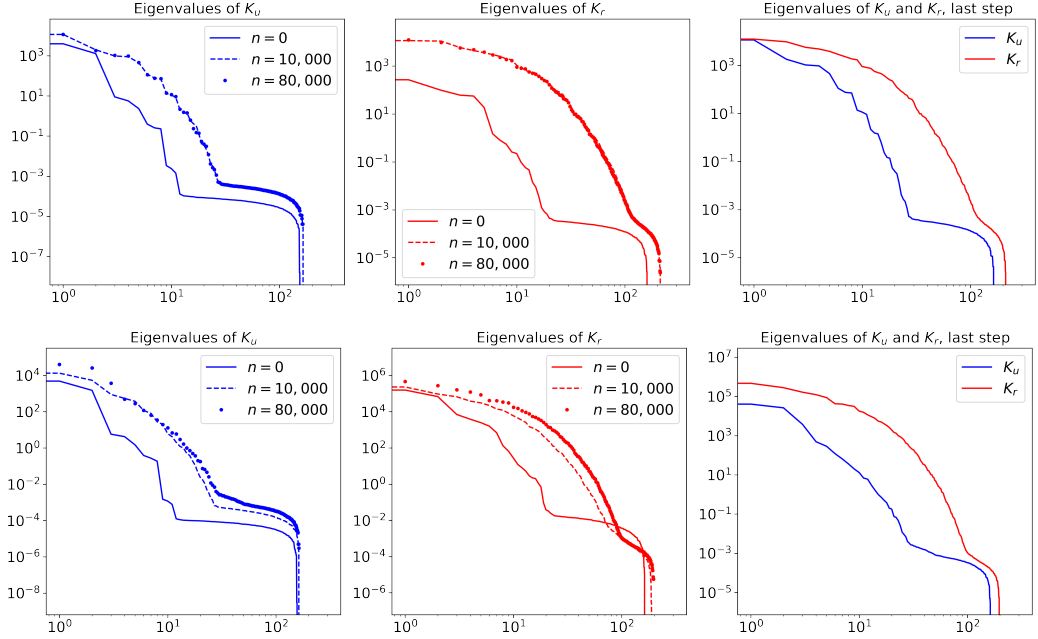


Figure 7: Evolution of NTK eigenvalues for the diffusion with $\epsilon = 0.1$ (top), $\epsilon = 10$ (bottom).

3.3 Advection-diffusion regime: a numerical example

Building on the previous sections' findings, we can argue that training failure mechanism of PINNs is different in the advection-dominated and diffusion-dominated regimes. While in both regimes PINNs suffer from the eigenvalue discrepancy, the training difficulty is more prominent in advection-dominated regime due to the underlying high-frequency behaviour of the solution in this regime and the spectral bias of PINNs. A successful PINN model in dealing with advection-diffusion equation (and possibly other class of PDEs) should, therefore, be able to address both eigenvalue disparity and spectral bias, at the same time.

As such, we illustrate the effectiveness of using sin activation function combined with adaptive weights strategy [16] for the LAD in a difficult regime with $a = 15$ and $\epsilon = 0.1$. Figure 11 illustrates the good accuracy of the proposed PINN for this case, with the note that the vanilla PINN with tanh activation [3] is unable to learn the solution. We shall discuss all our findings in the next Sec. 4

4 Discussion and Conclusions

In this work, we performed neural tangent kernel analysis of physics-informed neural networks for the linear advection-diffusion equation in advection-dominated and diffusion-dominated regimes.

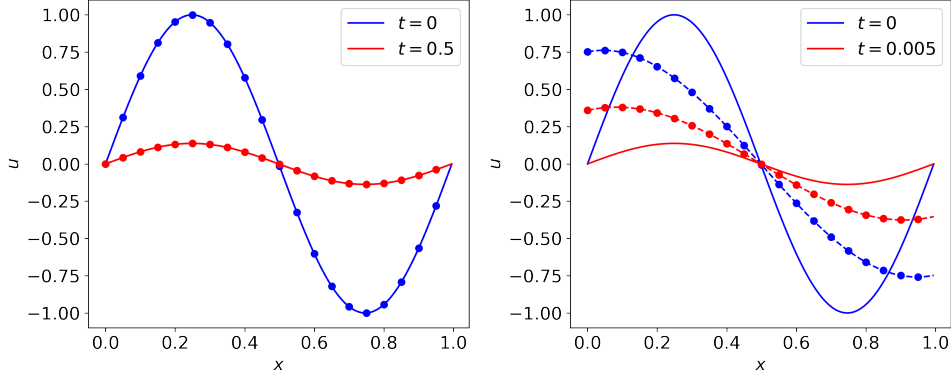


Figure 8: Predicted solution versus exact solution for the diffusion with $\epsilon = 0.1$ (left), $\epsilon = 10$ (right). Solid lines: exact solution; dashed lines: predicted solution.

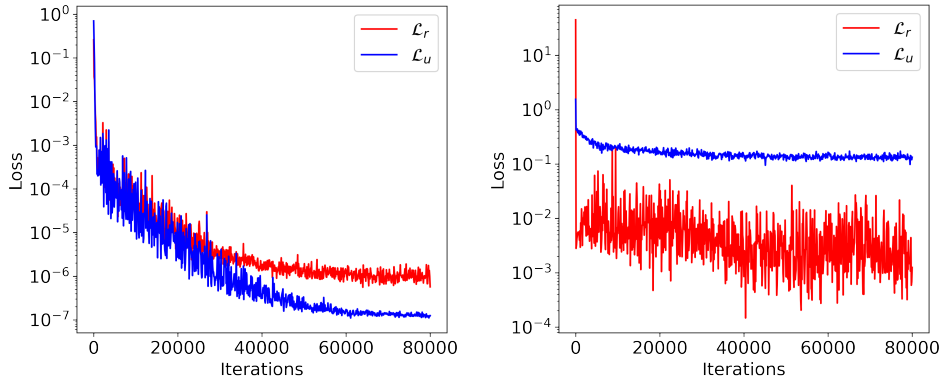


Figure 9: Time history of the different terms in the loss function for the diffusion with $\epsilon = 0.1$ (left), $\epsilon = 10$ (right).

Through numerical experiments and eigenvalue analysis of NTKs, we showed that in both regimes the training of PINNs become problematic at large PDE parameters, with possibility of total failure in extreme cases.

Going into more detail, we observed that increasing the advection speed results in: 1) severe eigenvalue disparity between the NTKs of the residual term \mathbf{K}_r and initial/boundary conditions term \mathbf{K}_u ; and 2) the solution in space-time domain exhibits high-frequency behaviour. The former results in domination of the residual term, and thus, it has faster convergence rate through the gradient descent training. This prevents the PINN from learning the correct initial/boundary conditions, which makes the problem ill-posed. The latter, on the other hand, poses challenge to PINNs due to the so-called 'spectral bias', a general feature of fully connected neural networks (like the ones used in PINNs) that tend to be biased towards low-frequency solutions. One remedy to the first issue is to reduce the eigenvalue disparity by manipulating the weights of different loss terms as, for example, it is done in the adaptive weights algorithm [19]. While using this algorithm improves the results to some extent, other strategies are needed to address the issue stemming from the 'spectral bias'. We showed that splitting the time domain, doing the training on a smaller sub-domain and then marching consecutively in time, as done in the sequence-to-sequence learning [15] or time adaptive re-sampling learning [27], is an effective approach to resolve this issue. These improved learning strategies, however, increase the computational cost as well as complexity of the algorithm. It was then demonstrated that a simple, yet effective strategy to address the spectral bias issue without any additional computational cost and major change in the algorithm is to use periodic activation functions and in particular sin activation. NTK analysis supplemented with numerical experiments showed that PINN models with sin activation are capable of representing higher-frequency solutions compared to tanh models. Last important thing to note is that, the behaviour of PINNs in advection-dominated

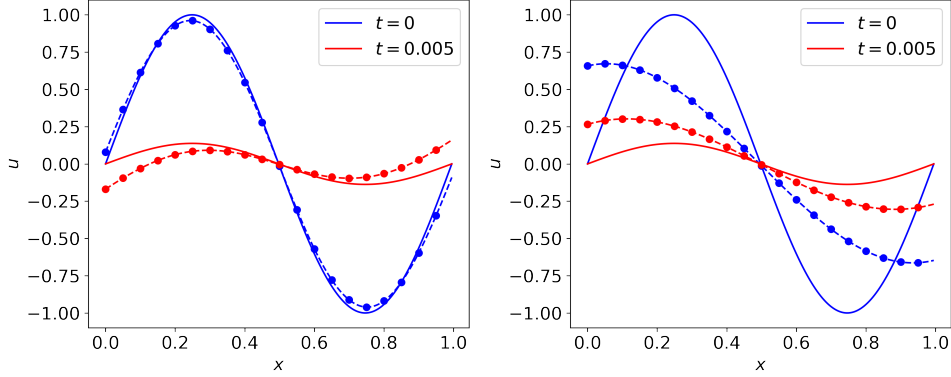


Figure 10: Predicted solution versus exact solution for the diffusion with $\epsilon = 10$ using adaptive weights strategy (left) and sub-domain learning on $t = [0, 0.1]$ (right). Solid lines: exact solution; dashed lines: predicted solution.

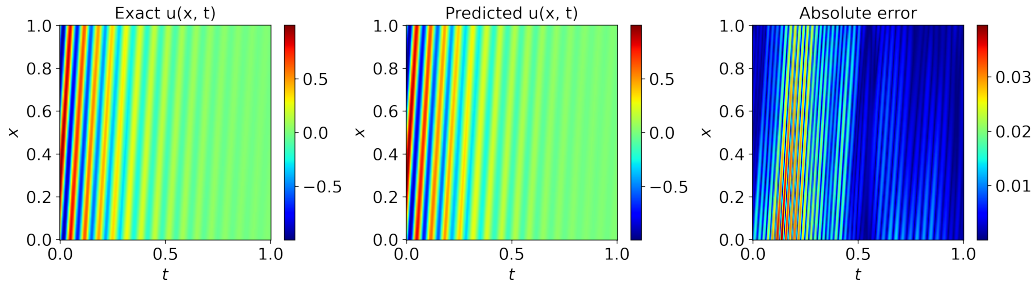


Figure 11: Snapshot of the predicted solution versus exact solution for the advection-diffusion with $a = 15$ and $\epsilon = 0.1$.

regimes is consistent with the behaviour of conventional numerical methods in this regime. It is well known that classical numerical methods are more prone to numerical instabilities in advection-dominated regimes [33] and there is a long tradition of designing numerical schemes suitable for advection dominated regime [34].

In the diffusion-dominated regime, the eigenvalue analysis of NTKs showed the same trend as in the advection-dominated regime, i.e. increasing the diffusion parameter leads to eigenvalue discrepancy, domination of the residual term in the training process, inability of PINNs in learning the correct initial/boundary conditions and eventually, PINNs failure. However, contrary to the advection-dominated regime, diffusion in general smooths out high-frequency behaviours in the solution. Yet, we showed that for the pure diffusion, where the solution dissipates quickly to a constant steady-state solution and seems to be trivial to learn, the PINN fails and cannot learn the exact solution. This surprising behaviour, however, could be explained by the eigenvalue disparity, which was resolved by employing the adaptive weights learning strategy [19]. It is also interesting to note that, when it comes to diffusion-dominated regimes, the described behaviour of PINNs is in sharp contrast to the behaviour of conventional numerical methods. From numerical point of view, it is in general desirable to have diffusion in our scheme as it dissipates spurious numerical oscillations leading to enhanced numerical stability [33]. Consequently, one can find an extensive literature on how to optimally add artificial diffusion in order to maintain the stability of a numerical scheme [34].

A unified training strategy that can address both spectral bias and eigenvalue disparity issues is, therefore, in great need. Recent works [26, 27] in the literature have shown promising results in that respect. Nevertheless, we showed that using periodic activation like \sin along with adaptive weights [16] is one possibility, which needs to be further investigated for more general and challenging PDEs with non-linearity. Analysing the performance of other periodic activation functions [29] in PINNs could be another line of research.

Acknowledgments and Disclosure of Funding

This work was supported by the Swiss Federal Office of Energy: “IMAGE - Intelligent Maintenance of Transmission Grid Assets” (Project Nr. SI/502073-01).

References

- [1] Bing Yu et al. The deep ritz method: a deep learning-based numerical algorithm for solving variational problems. *Communications in Mathematics and Statistics*, 6(1):1–12, 2018.
- [2] Justin Sirignano and Konstantinos Spiliopoulos. Dgm: A deep learning algorithm for solving partial differential equations. *Journal of computational physics*, 375:1339–1364, 2018.
- [3] Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics*, 378:686–707, 2019.
- [4] Jiequn Han, Arnulf Jentzen, and E Weinan. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [5] Lei Yuan, Yi-Qing Ni, Xiang-Yun Deng, and Shuo Hao. A-pinn: Auxiliary physics informed neural networks for forward and inverse problems of nonlinear integro-differential equations. *Journal of Computational Physics*, 462:111260, 2022.
- [6] Qin Lou, Xuhui Meng, and George Em Karniadakis. Physics-informed neural networks for solving forward and inverse flow problems via the boltzmann-bgk formulation. *Journal of Computational Physics*, 447:110676, 2021.
- [7] Ameya D Jagtap, Ehsan Kharazmi, and George Em Karniadakis. Conservative physics-informed neural networks on discrete domains for conservation laws: Applications to forward and inverse problems. *Computer Methods in Applied Mechanics and Engineering*, 365:113028, 2020.
- [8] Xiaowei Jin, Shengze Cai, Hui Li, and George Em Karniadakis. Nsfnets (navier-stokes flow nets): Physics-informed neural networks for the incompressible navier-stokes equations. *Journal of Computational Physics*, 426:109951, 2021.
- [9] Han Gao, Matthew J Zahr, and Jian-Xun Wang. Physics-informed graph neural galerkin networks: A unified framework for solving pde-governed forward and inverse problems. *Computer Methods in Applied Mechanics and Engineering*, 390:114502, 2022.
- [10] Ali Kashefi and Tapan Mukerji. Physics-informed pointnet: A deep learning solver for steady-state incompressible flows and thermal fields on multiple sets of irregular geometries. *arXiv preprint arXiv:2202.05476*, 2022.
- [11] Xuhui Meng, Zhen Li, Dongkun Zhang, and George Em Karniadakis. Ppinn: Parareal physics-informed neural network for time-dependent pdes. *Computer Methods in Applied Mechanics and Engineering*, 370:113250, 2020.
- [12] Maziar Raissi, Alireza Yazdani, and George Em Karniadakis. Hidden fluid mechanics: Learning velocity and pressure fields from flow visualizations. *Science*, 367(6481):1026–1030, 2020.
- [13] Zhiping Mao, Ameya D Jagtap, and George Em Karniadakis. Physics-informed neural networks for high-speed flows. *Computer Methods in Applied Mechanics and Engineering*, 360:112789, 2020.
- [14] Salvatore Cuomo, Vincenzo Schiano Di Cola, Fabio Giampaolo, Gianluigi Rozza, Maizar Raissi, and Francesco Piccialli. Scientific machine learning through physics-informed neural networks: Where we are and what’s next. *arXiv preprint arXiv:2201.05624*, 2022.
- [15] Aditi Krishnapriyan, Amir Gholami, Shandian Zhe, Robert Kirby, and Michael W Mahoney. Characterizing possible failure modes in physics-informed neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- [16] Sifan Wang, Yujun Teng, and Paris Perdikaris. Understanding and mitigating gradient flow pathologies in physics-informed neural networks. *SIAM Journal on Scientific Computing*, 43(5):A3055–A3081, 2021.

- [17] Siddhartha Mishra and Roberto Molinaro. Estimates on the generalization error of physics-informed neural networks for approximating a class of inverse problems for pdes. *IMA Journal of Numerical Analysis*, 42(2):981–1022, 2022.
- [18] Tim De Ryck and Siddhartha Mishra. Error analysis for physics informed neural networks (pinns) approximating kolmogorov pdes. *arXiv preprint arXiv:2106.14473*, 2021.
- [19] Sifan Wang, Xinling Yu, and Paris Perdikaris. When and why pinns fail to train: A neural tangent kernel perspective. *Journal of Computational Physics*, 449:110768, 2022.
- [20] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- [21] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *Advances in Neural Information Processing Systems*, 32, 2019.
- [22] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *Advances in neural information processing systems*, 32, 2019.
- [23] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Conference on Machine Learning*, pages 5301–5310. PMLR, 2019.
- [24] Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency principle: Fourier analysis sheds light on deep neural networks. *arXiv preprint arXiv:1901.06523*, 2019.
- [25] Sifan Wang, Hanwen Wang, and Paris Perdikaris. On the eigenvector bias of fourier feature networks: From regression to solving multi-scale pdes with physics-informed neural networks. *Computer Methods in Applied Mechanics and Engineering*, 384:113938, 2021.
- [26] Sifan Wang, Shyam Sankaran, and Paris Perdikaris. Respecting causality is all you need for training physics-informed neural networks. *arXiv preprint arXiv:2203.07404*, 2022.
- [27] Colby L Wight and Jia Zhao. Solving allen-cahn and cahn-hilliard equations using the adaptive physics informed neural networks. *arXiv preprint arXiv:2007.04542*, 2020.
- [28] Giambattista Parascandolo, Heikki Huttunen, and Tuomas Virtanen. Taming the waves: sine as activation function in deep neural networks. 2016.
- [29] Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2017.
- [30] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33:7462–7473, 2020.
- [31] Jian Cheng Wong, Chinchun Ooi, Abhishek Gupta, and Yew-Soon Ong. Learning in sinusoidal spaces with physics-informed neural networks. *IEEE Transactions on Artificial Intelligence*, 2022.
- [32] Josep M Sopena, Enrique Romero, and Rene Alquezar. Neural networks with periodic and monotonic activation functions: a comparative study in classification problems. 1999.
- [33] VK Suman, Tapan K Sengupta, C Jyothi Durga Prasad, K Surya Mohan, and Deepanshu Sanwalia. Spectral analysis of finite difference schemes for convection diffusion equation. *Computers & Fluids*, 150:95–114, 2017.
- [34] Sergio Pirozzoli. Numerical methods for high-speed flows. *Annual review of fluid mechanics*, 43:163–194, 2011.