

GLOBAL CONVERGENCE OF DEEP GALERKIN AND PINNS METHODS FOR SOLVING PARTIAL DIFFERENTIAL EQUATIONS

DEQING JIANG*, JUSTIN SIRIGNANO*, AND SAMUEL N. COHEN*

Abstract. Numerically solving high-dimensional partial differential equations (PDEs) is a major challenge. Conventional methods, such as finite difference methods, are unable to solve high-dimensional PDEs due to the curse-of-dimensionality. This is in particular a fundamental challenge for the solution of financial models, which are often inherently high-dimensional. Option pricing, hedging, mean-field financial models, order book models, and dynamic portfolio investment can all require the solution of high-dimensional PDEs. A variety of deep learning methods have been recently developed to try and solve high-dimensional PDEs by approximating the solution using a neural network. These deep learning methods have been widely applied to high-dimensional PDEs in financial engineering. In this paper, we prove global convergence for one of the commonly-used deep learning algorithms for solving PDEs, the Deep Galerkin Method (DGM). DGM trains a neural network approximator to solve the PDE using stochastic gradient descent. We prove that, as the number of hidden units in the single-layer network goes to infinity (i.e., in the “wide network limit”), the trained neural network converges to the solution of an infinite-dimensional linear ordinary differential equation (ODE). The PDE residual of the limiting approximator converges to zero as the training time $\rightarrow \infty$. Under mild assumptions, this convergence also implies that the neural network approximator converges to the solution of the PDE. A closely related class of deep learning methods for PDEs is Physics Informed Neural Networks (PINNs), which has been widely-used in a variety of fields (including financial mathematics but also physics and engineering). Using the same mathematical techniques, we can prove a similar global convergence result for the PINN neural network approximators. Both proofs require analyzing a kernel function in the limit ODE governing the evolution of the limit neural network approximator. A key technical challenge is that the kernel function, which is a composition of the PDE operator and the neural tangent kernel (NTK) operator, lacks a spectral gap, therefore requiring a careful analysis of its properties.

Key words. Neural network, partial differential equation, neural tangent kernel, machine learning, gradient flow

MSC codes. 65N75

1. Introduction. Deep learning methods have become widely used for solving high-dimensional PDEs in financial engineering. Similar deep learning methods have also been used to model physics data governed by PDEs. Although low-dimensional PDEs can be efficiently solved with existing numerical techniques, such as finite difference methods, high-dimensional PDEs are computationally intractable due to the curse of dimensionality. An alternative approach that has been widely employed is to approximate the PDE solution with a neural network and then train the neural network with stochastic gradient descent to satisfy the PDE and its boundary conditions – e.g., the deep Galerkin method (DGM) in [34]. The DGM algorithm – as well as many subsequent algorithms which have further developed the algorithm – has been applied to a wide range of high-dimensional PDEs in financial mathematics. A similar method – physics-informed neural networks (PINNs) in [29] – was developed to model physics data by training a neural network to both satisfy the corresponding governing PDE and match a sparse set of experimental observations. Subsequently, PINNs methods have also been used for a variety of applications in financial mathematics.

There are now hundreds of articles applying deep learning methods to PDE applications in mathematical finance. Examples include [3], [6], [4], [5], [30], [15], [16], [11], [2], [17], [13], [32], and [7]. There are many other important papers which have studied

*Mathematical Institute, University of Oxford (cohens@maths.ox.ac.uk, jiangd@maths.ox.ac.uk, sirignano@maths.ox.ac.uk).

48 the application of deep learning to solving PDEs in financial engineering, including
 49 mathematical analysis of neural networks in PDE applications. A critical gap in the
 50 mathematical theory for such methods is a proof that a neural network *trained with*
 51 *gradient descent* to minimize the PDE residual will converge to the solution of the
 52 PDE. In this paper, we develop a convergence analysis for DGM and PINNs methods,
 53 proving that they converge to the solution of a class of linear PDEs when trained with
 54 gradient descent. Our mathematical techniques also provide a foundation to study
 55 other classes of PDEs in financial mathematics.

56 Both the DGM and PINNs methods share a common feature of training a neural
 57 network with (stochastic) gradient descent to satisfy an objective function with a PDE.
 58 The neural network is trained (with gradient descent) to satisfy the PDE operator
 59 in the interior of the domain (i.e., to minimize the PDE residual) and to satisfy the
 60 boundary conditions. Numerous other articles in the literature have also explored
 61 solving PDEs with neural networks (see [1] for an overview). Solving PDEs with
 62 neural network approximators is a natural idea that has been considered in different
 63 forms for decades, for instance, [22], [21], [25] and [31]. These papers propose to use
 64 neural networks to solve differential equations by estimating neural network solutions
 65 on an *a priori* fixed mesh.

66 Due to the universal approximation properties of neural networks [18], it is clear
 67 that *there exists* a neural network which can approximate the solution to a given
 68 PDE (interpreting the PDE solution as a function in the Sobolev space \mathcal{H}^2). Universal
 69 approximation results have been proven for PINNs in [27] and [9]. Recently, important
 70 results were proven by [33] and [10], showing that the global minimizer of PINNs with
 71 Hölder and ridge regularization will converge to the solution of the PDE as the scale
 72 of the network goes to infinity. However, these results prove that a neural network
 73 *exists* which is an arbitrarily-accurate approximation to the PDE solution; they do
 74 not prove that a neural network *trained with gradient descent* will converge to the
 75 PDE solution. In particular, in the DGM algorithm, the PDE solution is *a priori*
 76 unknown. The neural network is therefore trained to satisfy the PDE operator in the
 77 interior of the domain and the boundary conditions on the domain boundary. This
 78 is a highly non-convex objective function; therefore, it is unclear whether the neural
 79 network (trained with gradient descent) will converge to the PDE solution.

80 Convergence analysis for the optimization of neural network approximators to
 81 PDEs must address several mathematical challenges. First, the neural network is non-
 82 convex in its parameters, which is further exacerbated by applying a PDE operator
 83 to the neural network in the objective function. Consequently, as the number of
 84 hidden units $\rightarrow \infty$, the standard neural tangent kernel (NTK) does not arise [19].
 85 Instead, the kernel function involves the PDE operator, requiring the development
 86 of new mathematical analysis. Finally, as is also true for the standard NTK setting,
 87 the kernel lacks a spectral gap, which makes analysis of infinite-dimensional systems
 88 (such as approximators to PDEs) challenging. Our proof leverages a careful analysis
 89 of the eigendecomposition of the limit ODE and its kernel function.

90 For both the DGM and PINN algorithms, we prove that, as the number of hidden
 91 units in the single-layer network goes to infinity (i.e., in the “wide network limit”),
 92 the trained neural network converges to the solution of an infinite-dimensional linear
 93 ODE. The PDE residual of the limiting approximator converges to zero as the training
 94 time $\rightarrow \infty$. Under mild assumptions, this convergence also implies that the neural
 95 network approximator converges to the solution of the PDE. [38] prove that, under
 96 some assumptions, as the number of neurons goes to infinity, the training process of
 97 PINNs will converge to a process characterized by a kernel matrix. However, [38] does

98 not prove global convergence of the neural network approximator.

99 Our paper provides a rigorous mathematical analysis of the DGM and PINN training
100 ing process for solving PDEs with neural networks. In summary, the key contributions
101 of this paper are:

- 102 1. We prove that as the number of hidden units in the neural network $\rightarrow \infty$ (i.e.,
103 in the “wide network limit), the training process of the neural approximator
104 trained to minimize the PDE residual converges to an infinite-dimensional
105 linear ODE characterized by a kernel function.
- 106 2. The kernel function is different than the standard NTK kernel and involves
107 the PDE operator.
- 108 3. We prove that even though the kernel is only positive semi-definite and there
109 is no spectral gap, the objective function (i.e., the PDE residual of the wide-
110 limit neural network) converges to zero as the training time $t \rightarrow \infty$. This
111 result establishes global convergence. Furthermore, under an additional mild
112 assumption, the wide-limit neural network converges to the PDE solution.

113 For our analysis of PINNs, we study a version of the PINNs algorithm where
114 spatial points are randomly sampled at each training iteration according to a fixed
115 probability measure μ . This leads to a term in the objective function which is an
116 expectation, with respect to the probability measure μ , of the PDE residual of the
117 neural network approximator over the entire domain.

118 Finally, we review some of the relevant literature on the mathematical analysis
119 of neural network methods for PDEs. In [24], a gradient descent training algorithm
120 for neural net approximations for PDEs is considered, under different assumptions
121 than in our paper and with different mathematical techniques. In particular, for
122 a particular choice of activation functions, assuming non-degeneracy of the Gram
123 matrix associated with the PDE kernel at finitely many (randomly chosen) points,
124 they argue that the gradient descent method will converge as training time $\rightarrow \infty$ to
125 a regularized approximation of the PDE. Under the assumption that the PDE source
126 term has bounded Rademacher complexity, the resulting sequence of approximations
127 will converge. However, it is not possible to interchange these two limiting steps (the
128 number of points and number of training steps) using their analysis. Our approach
129 instead proves that the trained neural network converges to an infinite-dimensional
130 evolution equation and then directly proves convergence of the limit network to the
131 PDE solution as the training time $\rightarrow \infty$, using eigenfunction decomposition analysis.
132 We do not assume that the eigenvalues of the corresponding NTK kernel are uniformly
133 lower bounded (which is not true as the number of hidden units in the network $\rightarrow \infty$)
134 and we do not make assumptions on the Rademacher complexity of the PDE.

135 A theoretical analysis for solving PDEs in Reproducing Kernel Hilbert Spaces
136 (RKHS) with applications to neural networks is given in [23]. Their analysis requires
137 the assumption that the PDE operator has the same eigenfunctions as the kernel,
138 which is not satisfied by the kernel induced by many common activation functions
139 used in deep learning. As an example, the eigenfunctions of the kernel for ReLU,
140 Sigmoid, and Tanh activation functions are not the eigenfunctions of a simple heat
141 equation (typically sine functions). Therefore, we cannot assume that the neural
142 network kernel eigenfunctions coincide with the PDE operator eigenfunctions. Our
143 proof does not require this assumption and therefore our convergence theorem applies
144 to a broad class of neural network models.

145 Our paper is organized as follows. In Section 2, we introduce the class of PDEs
146 that will be considered. Section 3 describes the neural network training algorithm for
147 solving PDEs and then proves that the neural network approximator converges to the

148 limit ODE as the number of hidden units $\rightarrow \infty$, and Section 4 presents the analysis
 149 on the kernel generated by the PDE and the activation function . Section 5 studies
 150 the properties of the kernel function that characterizes the limit ODE. Then, we prove
 151 that the PDE residual converges to zero as the training time $\rightarrow \infty$. Then, it is proven
 152 that – with an additional mild assumption on the PDE – the wide-limit neural network
 153 also converges to the PDE solution. In Section 6, we prove global convergence for the
 154 PINN algorithm. Lemmas, corollaries, and theorems are presented in the main part
 155 of the paper. All mathematical proofs are in the Appendix A.

156 **2. Mathematical Framework.** We will study the convergence of neural net-
 157 work algorithms – such as DGM and PINNs – for solving PDEs. In particular, we
 158 will analyze the convergence of such algorithms for the following class of second-order
 159 linear PDEs with Dirichlet boundary conditions:

$$160 \quad (2.1) \quad \begin{cases} \mathcal{A}v = h, & \text{in } \Omega \\ v = f, & \text{on } \partial\Omega, \end{cases}$$

161 where $\Omega \in \mathbb{R}^d$ is a compact set with a smooth boundary. We will study strong Sobolev
 162 solutions to the PDE (2.1); that is, we are interested in solutions $u \in \mathcal{H}^2$, where for
 163 a finite measure μ on Ω ,

$$164 \quad (2.2) \quad \mathcal{H}^p = \left\{ f \in L^2(\Omega, \mu) : \|f\|_{\mathcal{H}^p} := \left(\sum_{|\alpha| \leq p} \|D_\alpha f\|_{L^2} \right) < \infty \right\},$$

165 where Du is the weak derivative of u (see [12]), where (2.1) is taken to hold μ -almost
 166 everywhere on the interior of our domain, and in a trace sense on the boundary.
 167 We assume μ is equivalent to Lebesgue measure, and the logarithm of its Radon–
 168 Nikodym derivative is bounded (which ensures it generates the same \mathcal{H}^2 space as
 169 Lebesgue measure). For simplicity, in what follows we will simply say ‘solution’ for
 170 ‘strong solution’; we note that if a strong solution v is also known to be \mathcal{C}^2 , then it
 171 is also a classical solution. For notational convenience, we will write $\mathcal{H}_{(0)}^2 = \mathcal{H}^2 \cap$
 172 \mathcal{H}_0^1 , representing the \mathcal{H}^2 functions with zero value (in trace sense) on the boundary,
 173 equipped with the \mathcal{H}^2 norm.

174 We make the following (standard) assumptions on our problem:

175 **ASSUMPTION 2.1** (Smoothness of the boundary Ω). *The boundary $\partial\Omega$ is $C^{3,\alpha}$ for*
 176 *some $\alpha \in (0, 1)$; i.e., three times continuously differentiable with α -Hölder continuous*
 177 *derivatives of order 3.*

178 **ASSUMPTION 2.2** (Auxiliary function η). *There exists a (known) function $\eta \in$*
 179 *$C_b^3(\mathbb{R}^n)$, which satisfies $\eta > 0$ in Ω , and $\eta = 0$ on $\partial\Omega$. Furthermore, its first order*
 180 *derivative does not vanish at the boundary (that is, for $x \in \partial\Omega$ and \mathbf{n}_x an outward*
 181 *unit normal vector at x , we have $\nabla\eta(x) \cdot \mathbf{n}_x \neq 0$).*

182 *Remark 2.3.* If the boundary of our domain is the unit sphere, a natural choice
 183 of η is given by $\eta(x) = 1 - \|x\|^2$. More generally, a mollification of the distance to
 184 the boundary of the domain can usually be used, if this is easy to compute. In other
 185 cases the specification of the domain can suggest choices of η , for example, if we are
 186 given a domain in the form $\Omega = \{\eta(x) > 0\}$ for some function η , then this provides a
 187 natural choice of η provided $\nabla\eta \neq 0$ when $\eta = 0$.

188 **ASSUMPTION 2.4** (Interpolation of the boundary condition function). *There exists*
 189 *a (known) function $\bar{f} \in \mathcal{H}^2$ such that $\bar{f}|_{\partial\Omega} = f$. In the rest of this paper, we identify*
 190 *f with its extension \bar{f} defined on $\bar{\Omega}$ for notational simplicity.*

191 We can reformulate the PDE as

$$192 \quad (2.3) \quad \begin{cases} \mathcal{A}u = g, & \text{in } \Omega \\ u = 0, & \text{on } \partial\Omega, \end{cases}$$

193 where $u := v - f$ and $g := h - \mathcal{A}f$. Finally, we assume that the PDE operator satisfies
194 a certain type of Lipschitz condition:

195 **ASSUMPTION 2.5** (Lipschitz condition). *There exists a constant $k > 0$ such that*
196 *for any $f_1, f_2 \in \mathcal{H}^2$ and any $x \in \Omega$, the linear operator \mathcal{A} satisfies*

$$197 \quad (2.4) \quad |\mathcal{A}f_1(x) - \mathcal{A}f_2(x)| \leq k \left[\sum_{0 \leq |\alpha| \leq 2} |D_\alpha f_1(x) - D_\alpha f_2(x)| \right].$$

198 **3. Deep Learning Algorithms for Solving PDEs.** Deep learning PDE algo-
199 rithms – such as DGM and PINNs – train a neural network approximator to satisfy the
200 PDE and its boundary conditions using either gradient descent or stochastic gradient
201 descent.

202 Consider the following single-layer neural network, with N hidden units equipped
203 with activation function σ , S^N :

$$204 \quad (3.1) \quad S^N(x; \theta^N) = \frac{1}{N^\beta} \sum_{i=1}^N c^i \sigma(w^i x + b^i),$$

205 where $\frac{1}{N^\beta}$ is a normalization factor and $\frac{1}{2} < \beta < 1$. A normalization factor in front
206 of the summation $\sum_{i=1}^N c^i \sigma(w^i x + b^i)$ is necessary to avoid the neural network output
207 diverging as $N \rightarrow \infty$. The choice of $\frac{1}{N^\beta}$ with $\beta \in (\frac{1}{2}, 1)$ as a normalization, combined
208 with the learning rate, will lead to an NTK-type limit for the neural network as
209 $N \rightarrow \infty$. We train a neural network Q^N to approximate the solution u to the PDE
210 where

$$211 \quad (3.2) \quad Q^N(x; \theta^N) := \eta(x) S^N(x; \theta^N) = \eta(x) \cdot \frac{1}{N^\beta} \sum_{i=1}^N c^i \sigma(w^i x + b^i).$$

212 Here $\eta(x)$ is a fixed function which vanishes on the boundary $x \in \partial\Omega$; therefore Q^N
213 automatically satisfies the boundary conditions of the PDE (2.3). This method, which
214 was first introduced by [26], simplifies the training of the neural network model. The
215 parameters $\theta^N = (c^i, w^i, b^i)_{i=1}^N$ must be trained using gradient descent to satisfy the
216 PDE in the interior of the domain. Specifically, we will minimize the PDE residual
217 error for the neural network by minimizing the following objective function:

$$218 \quad (3.3) \quad J(\theta^N) = \|\mathcal{A}Q^N - g\|_{L^2(\mu)}^2 := \int_{\Omega} [\mathcal{A}Q^N(x; \theta^N) - g(x)]^2 d\mu(x),$$

219 where μ is a sampling measure (satisfying the regularity assumptions stated after
220 (2.2) above). If the residual term $\mathcal{A}Q^N(x; \theta^N) - g(x)$ equals zero for all $x \in \Omega$,
221 then $Q^N = u$ is the solution of the PDE. We will minimize the objective function
222 (3.3) using continuous-time gradient descent with clipping. The parameter updates
223 in continuous-time gradient descent satisfy the system of ODEs in equation (3.5),
224 which can be numerically solved using ODE numerical methods (e.g., Runge-Kutta

225 or Euler). If an Euler scheme is used, it becomes identical to standard discrete-time
 226 gradient descent. The gradient of (3.3) is:

$$227 \quad (3.4) \quad \nabla_{\theta} J(\theta^N) = \int_{\Omega} [\mathcal{A}Q^N(x; \theta^N) - g(x)] \nabla_{\theta} \mathcal{A}Q^N(x; \theta^N) d\mu(x).$$

228 Gradient clipping is widely used in deep learning, see for instance [39], [28] and chap-
 229 ters 10 and 11 of [14]. The continuous-time gradient descent training with clipping is
 230 given by:

$$231 \quad (3.5) \quad \frac{d\theta_t^N}{dt} = -\alpha^N G^N(\theta_t^N),$$

232 where the learning rate is $\alpha^N = N^{2\beta-1}$ and

$$233 \quad (3.6) \quad G^N(\theta_t^N) = - \int_{\Omega} \psi^N(\mathcal{A}Q^N(x; \theta_t^N) - g(x)) \Phi^N(\nabla_{\theta} \mathcal{A}Q^N(x; \theta_t^N)) d\mu(x).$$

234 Here Φ^N is a vector function that applies elementwise clipping with the function ϕ^N .
 235 For the explicit definition of clipping functions being applied here, see Definition 3.3,
 236 and Assumption 3.4. For each entry of vector $\nabla_{\theta} \mathcal{A}Q^N$, we clip its value with the
 237 scalar function ϕ^N .

238 In practice, (3.5) can be approximated by discretizing in time and, at each time
 239 step, generating Monte Carlo samples from the measure μ to approximate the inte-
 240 gral, which is highly computationally efficient even for high-dimensional PDEs. This
 241 is also equivalent to the stochastic gradient descent version of (3.5). Although not
 242 investigated in this paper, standard methods can be used to prove that stochastic gra-
 243 dient – using the correct learning rate – will converge to the continuous-time gradient
 244 flow (3.5); for example, weak convergence analysis such as in [36] and [35] could be
 245 used.

246 Let $Q_t^N(x) = Q^N(x; \theta_t^N)$ be the neural network at training time t with N hidden
 247 units. We will analyze the trained neural network Q_t^N as the number of hidden units
 248 $N \rightarrow \infty$ and the training time $t \rightarrow \infty$. First, we prove that the trained neural network
 249 Q_t^N will converge to the solution of an infinite-dimensional ODE as $N \rightarrow \infty$. That
 250 is, in the “wide limit” where the number of hidden units $\rightarrow \infty$, Q_t^N converges to the
 251 solution of an ODE. Then, we prove that the wide-limit neural network converges
 252 to the global minimizer of the objective function (with zero PDE residual) as the
 253 training time $t \rightarrow \infty$. Under additional mild assumptions, this global minimizer is
 254 also a solution to the PDE. These convergence results can also be proven for the
 255 PINNs algorithm for solving PDEs; see Section 6.

256 Our convergence results will be proven under the following assumptions on the
 257 neural network architecture:

258 **ASSUMPTION 3.1** (Activation function). *The activation function $\sigma \in C_b^4(\mathbb{R})$ is*
 259 *non-constant.*

260 While Assumption 3.1 is potentially restrictive, it is satisfied by commonly used acti-
 261 vation functions such as Sigmoid and Tanh functions. In the remainder of this paper,
 262 we use both $\sigma(x; w, b)$ and $\sigma_{w,b}(x)$ to denote the activation function.

263 **ASSUMPTION 3.2** (Neural network initialization). *The initialization of the pa-*
 264 *rameters θ_0^N , for all $i \in \{1, 2, \dots, N\}$, satisfies:*

- 265 • *The parameters c_0^i, w_0^i, b_0^i are i.i.d. random variables.*

- 266 • The random variables c_0^i are bounded, $|c_0^i| < K_0$, and $\mathbb{E}[c_0^i] = 0$.
- 267 • The distribution of the random variables w_0^i, b_0^i has full support. That is, for
- 268 any open set $D \subseteq \mathbb{R}^{n+1}$, we have $\mathbb{P}((w_0^i, b_0^i) \in D) > 0$.
- 269 • The moments $\mathbb{E}[|(w_0^i)_k|^3]$ and $\mathbb{E}[|b_0^i|]$ are bounded where $(w_0^i)_k$ is the k -element
- 270 of w_0^i .

271 **DEFINITION 3.3** (Smooth clipping function). A function class $\{h^N\}_{N \in \mathbb{N}^+}$ forms

272 a family of smooth clipping functions with parameter $\gamma > 0$ if for any $N \in \mathbb{N}^+$

- 273 • $h^N \in \mathcal{C}_b^2(\mathbb{R})$ is increasing on \mathbb{R} .
- 274 • $|h^N|$ is bounded by $2N^\gamma$.
- 275 • $h^N(x) = x$ for $x \in [-N^\gamma, N^\gamma]$.
- 276 • $|(h^N)'| \leq 1$ on \mathbb{R} .

277 **ASSUMPTION 3.4.** Functions $\{\psi^N\}_{N \in \mathbb{N}^+}$ and $\{\phi^N\}_{N \in \mathbb{N}^+}$ are families of smooth

278 clipping functions with parameter δ and $\epsilon - \beta$ where $\epsilon > \delta > 0$, $\beta \in (\frac{1}{2}, 1)$ and

279 $\epsilon + \delta < \frac{1-\beta}{2}$.

280 **LEMMA 3.5.** There exists a constant k independent of N such that the change of

281 each component of c^i, w^i, b^i from its initial condition (e.g. $|c_t^i - c_0^i|$) is bounded by

282 $ktN^{2\beta-1+\delta+\epsilon-\beta} = ktN^{\beta+\delta+\epsilon-1}$.

283 *Proof.* Notice that in (3.5), three terms $\alpha^N, \psi^N(\mathcal{A}Q^N - g)$ and $\Phi^N(\nabla_\theta \mathcal{A}Q^N)$ are

284 bounded by $N^{2\beta-1}, N^\delta$ and $N^{\epsilon-\beta}$ respectively. □

285

286 **Remark 3.6.** Gradient clipping is a method for preventing gradient explosion (i.e.,

287 gradient updates becoming very large/unbounded) and is a standard technique in

288 deep learning. Assumption 3.4 and Lemma 3.5 are introduced to make sure that

289 when $N \rightarrow \infty$, the truncation due to the clipping function vanishes while also guar-

290 anteeing stability (appropriately bounded parameter updates) during training. From

291 a mathematical perspective, the gradient clipping allows us to prove certain uniform

292 bounds which are useful for the convergence proof that the neural network training

293 converges to the limit ODE as the number of hidden units $N \rightarrow \infty$. Since the clipping

294 function vanishes in the limit $N \rightarrow \infty$, it does not appear in the limit ODE and there-

295 fore it also does not appear in the convergence analysis for when the training time

296 $t \rightarrow \infty$. From a numerical perspective, gradient clipping is a standard method in deep

297 learning to reduce instability in training (due to, for example, exploding gradients)

298 and improve training convergence.

299 **3.1. Convergence of the trained neural network as the number of hid-**

300 **den units $N \rightarrow \infty$.** In this section, we analyze the evolution of the neural network

301 Q_t^N as it is trained to minimize the PDE residual in the objective function $J(\theta^N)$. We

302 can prove that, as the number of hidden units $N \rightarrow \infty$, the neural network $Q_t^N(y)$ will

303 converges to the solution $Q_t(y)$ of an infinite-dimensional linear ODE. By the chain

304 rule, the dynamics of Q^N satisfy

$$(3.7) \quad \begin{aligned} \frac{dQ_t^N}{dt}(y) &= \nabla_\theta Q_t^N(y) \cdot \frac{d\theta^N}{dt} \\ &= - \int_{\Omega} \psi^N(\mathcal{A}Q_t^N(x) - g(x)) [\alpha^N \Phi^N(\nabla_\theta \mathcal{A}Q_t^N(x))] \cdot \nabla_\theta Q_t^N(y) d\mu(x). \end{aligned}$$

306 We are interested in studying the limit of the dynamics of the neural network Q_t^N as

307 $N \rightarrow \infty$. Specifically, we will prove that Q_t^N will converge to Q_t as $N \rightarrow \infty$ where Q_t

308 satisfies

$$309 \quad (3.8) \quad \frac{dQ_t}{dt}(y) = - \int_{\Omega} [\mathcal{A}Q_t(x) - g(x)]U(x, y)d\mu(x), \quad Q_0 = 0,$$

310 where the function U is

$$311 \quad (3.9) \quad U(x, y) := \mathbb{E}_{c,w,b} \left[\nabla_{c,w,b} \mathcal{A}[\eta(x)c\sigma(x; w, b)] \cdot \nabla_{c,w,b} [\eta(y)c\sigma(y; w, b)] \right],$$

312 where the random variable (c, w, b) has the same distribution as (c_0^i, w_0^i, b_0^i) .

313 The ODE (3.8) is an infinite-dimensional linear ODE governing the evolution
 314 of the wide-limit neural network (i.e., a neural network with an “infinite” number
 315 of hidden units) during training. The right-hand side (RHS) of the ODE involves
 316 integral over the PDE residual $[\mathcal{A}Q_t(x) - g(x)]$ weighted by a kernel $U(x, y)$. It is
 317 important to notice that the kernel $U(x, y)$ is not the standard NTK kernel: it involves
 318 the PDE operator \mathcal{A} , which significantly complicates its analysis.

319 One of the consequences of the presence of the PDE operator \mathcal{A} in the kernel
 320 $U(x, y)$ is that $U(x, y)$ is asymmetric. This is a key difference from the standard NTK
 321 kernel, which is symmetric.

322 Define the integral operator $\mathcal{U} : L^2 \rightarrow \mathcal{H}_{(0)}^2 \subset L^2$ by

$$323 \quad (3.10) \quad \mathcal{U}f := \int_{\Omega} f(x)U(x, y)d\mu(x).$$

324 Note that it is straightforward to check that $U(x, y) = 0$ for $y \in \partial\Omega$ and that U is
 325 C_b^2 with respect to y , which ensures that $\mathcal{U}f$ takes values in $\mathcal{H}_{(0)}^2$. (This is the main
 326 motivation for Assumption 3.1.) Using this notation, the limit ODE (3.8) can be
 327 rewritten as a linear equation in $\mathcal{H}_{(0)}^2$:

$$328 \quad (3.11) \quad \frac{dQ_t}{dt} = -\mathcal{U}[\mathcal{A}Q_t - g], \quad Q_0 = 0.$$

329

330 LEMMA 3.7. *The ODE (3.11) admits a unique solution in $\mathcal{H}_{(0)}^2$.*

331 *Proof.* By Assumptions 2.5, 3.1 and 3.2, the operator \mathcal{U} is Lipschitz in \mathcal{H}^2 norm.
 332 Therefore it admits a unique solution in \mathcal{H}^2 . It is easy to verify that, as its initial
 333 condition is in $\mathcal{H}_{(0)}^2$, and \mathcal{U} has codomain $\mathcal{H}_{(0)}^2$, the solution lives in the Hilbert
 334 subspace $\mathcal{H}_{(0)}^2$. \square

335 Now we present one of this paper’s main results. The trajectory of Q_t^N during training,
 336 in the limit $N \rightarrow \infty$, can be characterized by the wide-limit network Q_t which satisfies
 337 the infinite-dimensional ODE (3.11).

338 THEOREM 3.8. *For any $t \geq 0$, the neural network Q_t^N converges to Q_t in \mathcal{H}^2 :*

$$339 \quad (3.12) \quad \lim_{N \rightarrow \infty} \mathbb{E}[\|Q_t^N - Q_t\|_{\mathcal{H}^2}] = 0.$$

340 *The proof of this theorem is presented in the Appendix.*

341 **4. Analysis of the kernel function.** In order to prove global convergence as
 342 $t \rightarrow \infty$ for the limit ODE (3.11), we first must prove some key properties for the
 343 integral operator. Specifically, we will study the properties of the operator $\mathcal{S} = \mathcal{A}\mathcal{U}$.

344 DEFINITION 4.1 (Operator \mathcal{S}). The operator $\mathcal{S} : L^2 \rightarrow L^2$ is defined by

$$345 \quad (4.1) \quad \mathcal{S}f := \mathcal{A}Uf = \mathcal{A}\left(\int_{\Omega} f(x)U(x, \cdot)d\mu(x)\right)$$

346 DEFINITION 4.2 (Kernel S). The kernel S is defined by

$$347 \quad (4.2) \quad S(x, y) := \mathbb{E}_{c,w,b} \left[\nabla_{c,w,b} \mathcal{A}[\eta(x)c\sigma(x; w, b)] \cdot \nabla_{c,w,b} \mathcal{A}[\eta(y)c\sigma(y; w, b)] \right].$$

348 By symmetry of second derivatives (Clairaut–Schwarz–Young theorem) we know that
349 $S(x, \cdot) := \mathcal{A}U(x, \cdot)$ and hence $\mathcal{S}f = \int_{\Omega} f(x)S(x, \cdot)d\mu(x)$.

350 While $U(x, y)$ is asymmetric, $S(x, y)$ is symmetric. The symmetric kernel S depends
351 upon the interaction of PDE operator \mathcal{A} applied to the activation function σ . It will
352 next be proven (Lemma 4.6) that the operator \mathcal{S} is discriminatory in the image set of
353 \mathcal{A} . This is an important property that will later be leveraged in the global convergence
354 proof.

355 LEMMA 4.3. The kernel S is uniformly bounded. That is, there exists a constant
356 $k > 0$ such that for any $(x, y) \in \Omega^2$ we know $|S(x, y)| \leq k$.

357 *Proof.* By definition

$$358 \quad (4.3) \quad S(x, y) = \mathbb{E} \left[\mathcal{A}[\eta(x)\sigma_{w,b}(x)]\mathcal{A}[\eta(y)\sigma_{w,b}(y)] + \sum_{i=1}^d \mathcal{A}[c\eta(x)x_i\sigma'_{w,b}(x)]\mathcal{A}[c\eta(y)y_i\sigma'_{w,b}(y)] \right. \\ \left. + \mathcal{A}[c\eta(x)\sigma'_{w,b}(x)]\mathcal{A}[c\eta(y)\sigma'_{w,b}(y)] \right].$$

359 By the fact that \mathcal{A} is Lipschitz, and that η , $\sigma_{w,b}$ and their partial derivatives are all
360 bounded, there exists constant $k_1 > 0$, $k_2 > 0$ such that

$$361 \quad (4.4) \quad |\mathcal{A}[\eta(x)\sigma_{w,b}(x)]| \leq k \sum_{\alpha} |D_{\alpha}\eta(x)\sigma_{w,b}(x)| \leq k_1 \sum_{1 \leq i, j \leq d} |w_i| + |w_i w_j| + 1.$$

362 Similarly, since c is bounded and Ω is bounded

$$363 \quad (4.5) \quad |\mathcal{A}[c\eta(x)\sigma'_{w,b}(x)]| \leq k_1 \sum_{1 \leq i, j \leq d} |w_i| + |w_i w_j| + 1 \\ |\mathcal{A}[c\eta(x)x_i\sigma'_{w,b}(x)]| \leq k_1 \sum_{1 \leq i, j \leq d} |w_i| + |w_i w_j| + 1.$$

364 Therefore,

$$365 \quad (4.6) \quad |S(x, y)| \leq \mathbb{E}_{c,w,b} [(d+2)k_1^2 \left(\sum_{1 \leq i, j \leq d} |w_i| + |w_i w_j| + 1 \right)^2] \leq k_2. \quad \square$$

366 LEMMA 4.4. The integral operator \mathcal{S} is Hilbert–Schmidt. In addition, \mathcal{S} is self-
367 adjoint and positive semi-definite.

368 *Proof.* Since $S(x, y)$ is uniformly bounded, we have $\iint_{\Omega^2} |S(x, y)|^2 d\mu(x)d\mu(y) <$
369 ∞ . Therefore \mathcal{S} is Hilbert–Schmidt. And since $S(x, y)$ is symmetric, the operator \mathcal{S}

370 is self-adjoint with respect to the L^2 inner product. Now it remains to prove that \mathcal{S}
371 is positive semi-definite. For $f \in L^2$

(4.7)

$$\begin{aligned} \langle f, \mathcal{S}f \rangle &= \iint_{\Omega^2} f(x)S(x, y)f(y)d\mu(x)d\mu(y) \\ &= \iint f(x)\mathbb{E}_{c,w,b}[\nabla_{c,w,b}\mathcal{A}[\eta(x)c\sigma(x; w, b)] \cdot \nabla_{c,w,b}\mathcal{A}[\eta(y)c\sigma(y; w, b)]]f(y)d\mu(x)d\mu(y) \end{aligned}$$

373 By Tonelli's theorem, swapping the order of expectation and the integral gives

(4.8)

$$\begin{aligned} \langle f, \mathcal{S}f \rangle &= \mathbb{E} \left[\iint_{\Omega^2} f(x)\nabla_{c,w,b}\mathcal{A}[\eta(x)c\sigma(x; w, b)] \cdot \nabla_{c,w,b}\mathcal{A}[\eta(y)c\sigma(y; w, b)]f(y)d\mu(x)d\mu(y) \right] \\ &\geq \mathbb{E}_{c,w,b} \left[\iint_{\Omega^2} f(x)\nabla_c\mathcal{A}[\eta(x)c\sigma(x; w, b)] \cdot \nabla_c\mathcal{A}[\eta(y)c\sigma(y; w, b)]f(y)d\mu(x)d\mu(y) \right] \\ &= \mathbb{E}_{c,w,b} \left[\iint_{\Omega^2} f(x)\mathcal{A}[\eta(x)\sigma(x; w, b)]\mathcal{A}[\eta(y)\sigma(y; w, b)]f(y)d\mu(x)d\mu(y) \right] \\ &= \mathbb{E}_{c,w,b} \left[\left(\int_{\Omega} f(x)\mathcal{A}[\eta(x)\sigma(x; w, b)]d\mu(x) \right)^2 \right] \geq 0, \end{aligned}$$

375 which concludes the proof. \square

376 LEMMA 4.5 (Spectral decomposition). *The integral operator \mathcal{S} is compact, in*
377 *particular, there exists an orthogonal basis of L^2 , $\{\varepsilon_i\}_{i \in \mathbb{N}^+} \cup \{\nu_i\}_{i \in \mathbb{N}^+}$ such that*

$$378 \quad (4.9) \quad \mathcal{S}\varepsilon_i = \lambda_i\varepsilon_i, \quad \mathcal{S}\nu_i = 0,$$

379 where $\lambda_1 \geq \lambda_2 \geq \dots > 0$.

380 *Proof.* As \mathcal{S} is a Hilbert–Schmidt integral operator, \mathcal{S} is compact. Since \mathcal{S} is
381 self-adjoint, the spectral theorem applies. From Lemma 4.4, we see that \mathcal{S} is positive
382 semi-definite. Thus, its eigenvalues are real, non-negative, and concentrate only at
383 zero. We use $\{\nu_i\}$ to represent the eigenfunctions of the zero eigenvalue, and $\{\varepsilon_i\}$ for
384 eigenfunctions of positive eigenvalues. \square

385 LEMMA 4.6 (Projection on $\ker(\mathcal{S})$). *For $h \in L^2$, if $\mathcal{S}h = 0$ then $\langle h, \mathcal{A}f \rangle = 0$ for*
386 *any $f \in \mathcal{H}_{(0)}^2$.*

387 We recall without proof the following technical result.

388 LEMMA 4.7 (Lemma 5 in [8]). *Given Assumptions 2.1 and 2.2:*

- 389 1. *The set of functions $C^3(\overline{\Omega}) \cap C_0(\overline{\Omega})$ is dense in $\mathcal{H}_{(0)}^2 = \mathcal{H}^2 \cap \mathcal{H}_0^1$ (under the*
390 *\mathcal{H}^2 topology).*
- 391 2. *For any function $u \in C^3(\overline{\Omega}) \cap C_0(\overline{\Omega})$, the function $\tilde{u} = u/\eta$ is in $C_b^2(\Omega) \subset \mathcal{H}^2$.*

392 We now proceed to the proof of Lemma 4.6.

393 *Proof of Lemma 4.6.* By (4.8), we have

$$394 \quad (4.10) \quad \langle h, \mathcal{S}h \rangle = \mathbb{E}_{c,w,b} \left[\left(\int_{\Omega} h(x)\mathcal{A}[\eta(x)\sigma(x; w, b)]d\mu(x) \right)^2 \right] = 0.$$

395 Therefore

$$396 \quad (4.11) \quad \int_{\Omega} h(x) \mathcal{A}[\eta(x) \sigma(x; w, b)] d\mu(x) = 0$$

397 for any (c, w, b) by the continuity of the objective with respect to parameters c, w, b .

398 Since μ is a finite measure, from Theorem 4 in [18] we have that the linear
 399 span of $\{\sigma(w \cdot x + b)\}_{w, b \in \mathbb{R}}$ is dense in \mathcal{H}^2 . For a general $f \in C^3(\Omega) \cap C_0(\Omega)$, by
 400 Lemma 4.7(2), we can approximate the function f/η within the linear span of $\{\sigma(w \cdot$
 401 $x + b)\}_{w, b \in \mathbb{R}}$. Multiplying by η , by Lemma 4.7(1) it follows that the function class
 402 $\{\eta(x) \sigma(x; w, b)\}$ is dense in the function space $\mathcal{H}_{(0)}^2(\Omega)$, which consists of \mathcal{H}^2 function
 403 with boundary value zero. For any $f \in \mathcal{H}_{(0)}^2$, there exists a function sequence $\{F_N :$
 404 $\eta \sum_{i=1}^N c_i \sigma_{w, b}\}_{N \geq 1}$ such that

$$405 \quad (4.12) \quad \lim_{N \rightarrow \infty} \|F_N - f\|_{\mathcal{H}^2} = 0.$$

406 Therefore,

$$407 \quad (4.13) \quad \langle h, \mathcal{A}f \rangle = \int_{\Omega} h(x) \mathcal{A}f d\mu(x) = \lim_{N \rightarrow \infty} \langle h, \mathcal{A}F^N \rangle = 0. \quad \square$$

408 *Remark 4.8.* If u is a solution to the PDE (2.3), Lemma 4.6 implies that for any
 409 $t \geq 0$, the residual of our approximator $\mathcal{A}Q_t - g = \mathcal{A}[Q_t - u]$ has zero projection on
 410 the eigenfunction family $\{\nu_i\}$ since $\mathcal{S}\nu_i = 0$.

411 **COROLLARY 4.9.** *Assuming the PDE (2.3) admits a solution $u \in \mathcal{H}_{(0)}^2$, any sta-*
 412 *tionary point $Q^* \in \mathcal{H}_{(0)}^2$ of the limit ODE (3.8) is a solution of the PDE (2.3). That*
 413 *is, any stationary point of the limit training algorithm (3.8) is a global minimizer and*
 414 *a solution of (2.3).*

415 *Proof.* Suppose that Q^* is a stationary point. Then we have

$$416 \quad (4.14) \quad \mathcal{U}[\mathcal{A}Q^* - g] = 0.$$

417 It follows that

$$418 \quad (4.15) \quad \mathcal{S}\mathcal{A}[Q^* - u] = \mathcal{S}[\mathcal{A}Q^* - g] = \mathcal{A}\mathcal{U}[\mathcal{A}Q^* - g] = 0.$$

419 By Lemma 4.6, the inner product term $\langle \mathcal{A}[Q^* - u], \mathcal{A}f \rangle = 0$ for any $f \in \mathcal{H}_{(0)}^2$.
 420 Therefore, by taking $f = Q^* - u$, we have $\|\mathcal{A}Q^* - g\|_2^2 = \langle \mathcal{A}[Q^* - u], \mathcal{A}[Q^* - u] \rangle = 0$.
 421 At the same time, since $Q^* \in \mathcal{H}_{(0)}^2$ satisfies the zero boundary condition, we conclude
 422 that Q^* is a solution of the PDE. \square

423 *Remark 4.10.* The existence of solutions is needed in this result. For example,
 424 consider the trivial case where $\mathcal{A}u \equiv 0$. Then the kernel satisfies $\mathcal{U} \equiv 0$, so $\frac{d}{dt}Q_t = 0$,
 425 and we must have a fixed point. However, if $g \neq 0$, then $Q_t = Q_0$ does not solve the
 426 PDE. Nevertheless, no uniqueness of solutions is needed, and consequently we cannot
 427 expect that $Q^* = u$.

428 **5. Global convergence of the limit ODE as $t \rightarrow \infty$.** By analyzing the PDE
 429 residual term's projection on the eigenfunctions of \mathcal{S} , we can prove that the PDE
 430 residual (which is the objective function that is being minimized) converges to zero
 431 as the training time $t \rightarrow \infty$.

432 Applying \mathcal{A} with respect to y on both sides of the limit ODE (3.11) yields

$$433 \quad (5.1) \quad \frac{d\mathcal{A}Q_t}{dt} = \frac{d[\mathcal{A}Q_t - g]}{dt} = -\mathcal{S}[\mathcal{A}Q_t - g].$$

434

435 **THEOREM 5.1** (Convergence of PDE residual under DGM). *Assuming the PDE*
 436 *(2.3) admits a solution in $\mathcal{H}_{(0)}^2$, the PDE residual for the wide-limit neural network*
 437 *Q_t converges to zero:*

$$438 \quad (5.2) \quad \lim_{t \rightarrow \infty} \|\mathcal{A}Q_t - g\|_{L^2(\mu)} = 0.$$

439 *Proof.* In (5.1), multiplying $\mathcal{A}Q_t - g$ on both sides and integrating with respect
 440 to $\mu(dy)$:

$$441 \quad (5.3) \quad \frac{d\|\mathcal{A}[Q_t - u]\|_2^2}{dt} = -2\langle \mathcal{A}[Q_t - u], \mathcal{S}\mathcal{A}[Q_t - u] \rangle \leq 0$$

442 with strict inequality unless $\|\mathcal{A}Q_t - g\|_2 = 0$, which corresponds to the PDE solution.

443 Consider $\mathcal{A}Q_t - g = \mathcal{A}[Q_t - u]$ projected on $\{\varepsilon_i\}_{i \in \mathbb{N}^+} \cup \{\nu_i\}_{i \in \mathbb{N}^+}$. By Lemma 4.6,
 444 we have

$$445 \quad (5.4) \quad \mathcal{A}Q_t - g = \sum_i h_t^i \varepsilon_i + \sum_j 0\nu_j = \sum_i h_t^i \varepsilon_i.$$

446 Therefore $\|\mathcal{A}Q_t - g\|_2^2 = \sum_i h_t^{i2}$. Now consider its projection on each ε_i

$$447 \quad (5.5) \quad \begin{aligned} \frac{d}{dt} \langle \mathcal{A}Q_t - g, \varepsilon_i \rangle &= \left\langle \frac{d\mathcal{A}Q_t - g}{dt}, \varepsilon_i \right\rangle = \langle -\mathcal{S}[\mathcal{A}Q_t - g], \varepsilon_i \rangle = \langle \mathcal{A}Q_t - g, -\mathcal{S}\varepsilon_i \rangle \\ &= -\lambda_i \langle \mathcal{A}Q_t - g, \varepsilon_i \rangle. \end{aligned}$$

448 Therefore

$$449 \quad (5.6) \quad \frac{d}{dt} \langle \mathcal{A}Q_t - g, \varepsilon_i \rangle = \frac{d}{dt} [h_t^i] = -\lambda_i h_t^i.$$

450 Consequently $h_t^i = h_0^i e^{-\lambda_i t}$ whose absolute value decays exponentially, and $|h_t^i| \leq |h_0^i|$
 451 for any $t \geq 0$ and any i .

452 Now, by the dominated convergence theorem

$$453 \quad (5.7) \quad \lim_{t \rightarrow \infty} \|\mathcal{A}Q_t - g\|_2^2 = \lim_{t \rightarrow \infty} \sum_i |h_t^i|^2 = \sum_i \lim_{t \rightarrow \infty} |h_t^i|^2 = 0. \quad \square$$

454 This result establishes global convergence of the (wide-limit) training algorithm since
 455 the objective function is the PDE residual.

456 *Remark 5.2.* Due to the lack of a spectral gap for the eigenvalues of the kernel, it
 457 is difficult to establish a convergence rate since some eigenfunction projections may be
 458 converging very slowly (if their corresponding eigenvalue is close to zero). As shown in
 459 our analysis, for an eigenfunction ε_i , we can establish an exponential convergence rate
 460 for eigenfunction projections of the PDE residual $\langle \varepsilon_i, \mathcal{A}Q_t - g \rangle$ where Q_t is the neural
 461 network and $\mathcal{A}u = g$ is the PDE we are trying to solve. In particular, $|\langle \varepsilon_i, \mathcal{A}Q_t - g \rangle| \leq$
 462 $C_i \exp(-\lambda_i t)$ where λ_i is the eigenvalue corresponding to the eigenfunction ε_i .

463 However, the convergence of the PDE residual to zero does not necessarily guarantee
 464 that $\|Q_t - u\|_{L^2(\mu)}$ as $t \rightarrow \infty$ where u is the solution of the PDE (2.3). We show that,
 465 under a mild additional assumption on the PDE operator \mathcal{A} , we can guarantee that
 466 Q_t converges to the solution of the PDE u .

467 **ASSUMPTION 5.3** (Bounded inverse of \mathcal{A}). *The inverse of \mathcal{A} is a bounded operator*
 468 *on $L^2 \rightarrow L^2$. That is, there exists a constant $k > 0$, such that for any $g \in L^2$, the*
 469 *PDE (2.3) has a unique solution $u \in \mathcal{H}_{(0)}^2$ satisfying*

$$470 \quad \|u\|_2 \leq k \|g\|_{L^2(\mu)}.$$

471 *Remark 5.4.* Second-order uniformly elliptic PDEs naturally have this property.
 472 For reference, see Theorem 6, Chapter 6 of [12].

473 **THEOREM 5.5** (Convergence of Q_t under DGM). *If \mathcal{A} has a bounded inverse, Q_t*
 474 *converges to the solution u of the PDE:*

$$475 \quad \lim_{t \rightarrow \infty} \|Q_t - u\|_{L^2(\mu)} = 0.$$

476 *Proof.* Writing \mathcal{A}^{-1} for the inverse operator of \mathcal{A} , we have

$$477 \quad (5.8) \quad \|Q_t - u\|_2 = \|\mathcal{A}^{-1}[\mathcal{A}[Q_t - u]]\|_2 \leq k \|\mathcal{A}[Q_t - u]\|_2 = k \|\mathcal{A}Q_t - g\|_2 \rightarrow 0. \quad \square$$

478 **6. Global Convergence of the PINN Algorithm.** In many real-world engi-
 479 neering applications, the PDE solution can be observed at a sparse set of points in
 480 the interior of the domain Ω . We consider a PINN algorithm which trains a neural
 481 network model to predict the solution $u(x)$ for all $x \in \Omega$ by minimizing (via gradient
 482 descent) both the PDE residual sampled by μ on a random set of points in each epoch
 483 as well as the distance between the neural network and the observations at the finite
 484 set of sparse points. It is assumed that the boundary condition is known.

485 Let us denote these extra observation points in Ω as $\mathbf{x} := \{x_i\}_{i=1,2,\dots,M}$ where the
 486 observations are $u(x_i) = u_i$ for $1 \leq i \leq M$. Define the measure $\mu_{\mathbf{x}} = \frac{1}{M} \sum_{i=1}^M \delta_{x_i}$ and
 487 the kernel function

$$488 \quad (6.1) \quad B(x, y) := \mathbb{E}_{c,w,b} \left[\nabla_{c,w,b}[c\eta(x)\sigma_{w,b}(x)] \cdot \nabla_{c,w,b}[c\eta(y)\sigma_{w,b}(y)] \right],$$

489 and the corresponding ‘integral’ operator $\mathcal{B} : L^2(\mu_{\mathbf{x}}) \equiv \mathbb{R}^M \rightarrow L^2(\mu)$ by

$$490 \quad (6.2) \quad \mathcal{B}\mathbf{v}(y) := \frac{1}{M} \sum_{i=1}^M v_i B(x_i, y).$$

491 For any operator \mathcal{C} , we define $\bar{\mathcal{C}}$ to be the evaluation of \mathcal{C} on each of our training
 492 points,

$$493 \quad \bar{\mathcal{C}}\mathbf{v} = [\mathcal{C}\mathbf{v}(x_1), \dots, \mathcal{C}\mathbf{v}(x_M)]^\top \in \mathbb{R}^M,$$

494 in particular

$$495 \quad (6.3) \quad \bar{\mathcal{B}}\mathbf{v} := \left[\frac{1}{M} \sum_{i=1}^M v_i B(x_i, x_1), \dots, \frac{1}{M} \sum_{i=1}^M v_i B(x_i, x_M) \right]^\top,$$

496 and for generic $g \in L^2(\mu)$ and $\mathbf{v} \in \mathbb{R}^M$, direct calculation shows that

$$\begin{aligned}
 (6.4) \quad \langle \bar{\mathcal{U}}g, \mathbf{v} \rangle_{L^2(\mu_{\mathbf{x}})} &= \frac{1}{M} \sum_i (\mathcal{U}g)(x_i) v_i \\
 497 \quad &= \frac{1}{M} \sum_i \int \mathbb{E} \left[v_i g(x) \nabla_{c,w,b} \mathcal{A} [c\eta(x) \sigma_{w,b}(x)] \cdot \nabla_{c,w,b} [c\eta(x_i) \sigma_{w,b}(x_i)] \right] d\mu \\
 &= \langle g, \mathcal{A}\mathcal{B}\mathbf{v} \rangle_{L^2(\mu)},
 \end{aligned}$$

498 in other words, $\bar{\mathcal{U}} : L^2(\mu) \rightarrow \mathbb{R}^M$ is the adjoint of $\mathcal{A}\mathcal{B} : \mathbb{R}^M \rightarrow L^2(\mu)$, so we write
 499 $\mathcal{A}\mathcal{B} = \bar{\mathcal{U}}^*$.

500 The PINN objective function is

$$\begin{aligned}
 (6.5) \quad J(\theta^N) &= \|\mathcal{A}Q_t^N - g\|_{L^2(\mu)}^2 + \|Q_t^N - u\|_{L^2(\mu_{\mathbf{x}})}^2 \\
 &= \int_{\Omega} (\mathcal{A}Q_t^N - g)^2 d\mu + \int_{\Omega} (Q_t^N - u)^2 d\mu_{\mathbf{x}}.
 \end{aligned}$$

502 Informally, the wide-limit ODE (as $N \rightarrow \infty$) for the neural network trained with
 503 continuous-time gradient descent is:

$$(6.6) \quad \frac{dQ_t}{dt} = -\mathcal{U}[\mathcal{A}Q_t - g] - \mathcal{B}[Q_t - u].$$

505 Applying the operator \mathcal{A} and recalling $\mathcal{A}\mathcal{B}$ is the adjoint of $\bar{\mathcal{U}} : L^2(\mu) \rightarrow \mathbb{R}^M$, we
 506 derive that

$$(6.7) \quad \frac{d}{dt} \begin{bmatrix} \mathcal{A}Q_t - g \\ \mathbf{Q}_t - \mathbf{u} \end{bmatrix} = - \begin{bmatrix} \mathcal{S} & \bar{\mathcal{U}}^* \\ \bar{\mathcal{U}} & \bar{\mathcal{B}} \end{bmatrix} \begin{bmatrix} \mathcal{A}Q_t - g \\ \mathbf{Q}_t - \mathbf{u} \end{bmatrix},$$

508 where, on the left-hand side of the equation, the first term $\mathcal{A}Q_t - g \in L^2(\mu)$ and the
 509 second term $\mathbf{Q}_t - \mathbf{u} \in \mathbb{R}^M$ is

$$(6.8) \quad \mathbf{Q}_t - \mathbf{u} = \begin{bmatrix} Q_t(x_1) - u(x_1) \\ Q_t(x_2) - u(x_2) \\ \vdots \\ Q_t(x_M) - u(x_M) \end{bmatrix}.$$

511 In light of this, we define the operator $\mathcal{V} : L^2(\mu) \times \mathbb{R}^M \rightarrow L^2(\mu) \times \mathbb{R}^M$

$$(6.9) \quad \mathcal{V}f := \begin{bmatrix} \mathcal{S} & \bar{\mathcal{U}}^* \\ \bar{\mathcal{U}} & \bar{\mathcal{B}} \end{bmatrix} f.$$

513 Given the properties of \mathcal{S} (Lemma 4.4), it is easy to verify that \mathcal{V} is Hilbert–Schmidt,
 514 self-adjoint and positive semi-definite. It will not generally be strictly positive definite.
 515 However, using the representation properties of neural networks, we can prove that
 516 $(\mathcal{A}Q_t - g, \mathbf{Q}_t - \mathbf{u})^\top$ in (6.7) has zero projection on the kernel set of \mathcal{V} .

517 **LEMMA 6.1** (Projection on $\ker(\mathcal{V})$). *For $h \in L^2(\mu) \times \mathbb{R}^M$, if $\mathcal{V}h = 0$ then*
 518 *$\langle h, \mathcal{V}m \rangle = 0$ for any $m = [\mathcal{A}f - g, f(x_1) - u(x_1), \dots, f(x_M) - u(x_M)]^\top$ where $f \in \mathcal{H}_{(0)}^2$.*

519 *Proof.* Let us begin with the following remark.

520 *Remark 6.2.* The setting we consider has a small modification as compared to
 521 the original PINNs formulation: we consider the limit ODE for an algorithm which
 522 samples points at random using the measure μ instead of using an *a priori* fixed set
 523 of points. We study this setting because if the PDE residual is minimized only on
 524 a fixed sparse set of points, then there is no guarantee that the approximator will
 525 converge to the solution – given any function, it is possible to perturb the function
 526 locally around each measurement point to construct functions which will have zero
 527 residual at these points; such functions will generally have no relation to the solution
 528 of the PDE. For similar reasons, we assume that the boundary value of the PDE is
 529 perfectly known.

530 Denote by $\varrho = \begin{bmatrix} \varrho^d \\ \varrho^p \end{bmatrix}$ an eigenfunction of \mathcal{V} which has an eigenvalue zero:

$$531 \quad (6.10) \quad \mathcal{V} \begin{bmatrix} \varrho^d \\ \varrho^p \end{bmatrix} = \mathcal{V}\varrho = 0\varrho = \begin{bmatrix} \mathcal{S} & \bar{\mathcal{U}}^* \\ \bar{\mathcal{U}} & \bar{\mathcal{B}} \end{bmatrix} \begin{bmatrix} \varrho^d \\ \varrho^p \end{bmatrix}.$$

532 We now show that

$$533 \quad (6.11) \quad \left\langle \varrho, \begin{bmatrix} \mathcal{A}f - g \\ \mathbf{f} - \mathbf{u} \end{bmatrix} \right\rangle = 0,$$

534 for any $f \in \mathcal{H}_{(0)}^2$. This is because, similar to (4.8), we have

$$535 \quad (6.12) \quad \langle \varrho, \mathcal{V}\varrho \rangle \geq \mathbb{E} \left[\left(\int_{\Omega} \varrho^d(x) \mathcal{A}[\eta(x)\sigma_{w,b}(x)] d\mu(x) + \int_{\Omega} \varrho^p(x) \eta(x) \sigma_{w,b}(x) d\mu_{\mathbf{x}}(x) \right)^2 \right] \geq 0.$$

536 Therefore, $\mathcal{V}\varrho = 0$ implies that for any (w, b) pairs

$$537 \quad (6.13) \quad \int_{\Omega} \varrho^d(x) \mathcal{A}[\eta(x)\sigma(x; w, b)] d\mu(x) + \int_{\Omega} \varrho^p(x) \eta(x) \sigma(x; w, b) d\mu_{\mathbf{x}}(x) = 0.$$

538 Now, for any $f - u \in \mathcal{H}_{(0)}^2$, there exists a function sequence $\{F_N : \eta \sum_{i=1}^N c_i \sigma_{w,b}\}_{N \geq 1}$
 539 such that

$$540 \quad (6.14) \quad \lim_{N \rightarrow \infty} \|F_N - (f - u)\|_{\mathcal{H}^2}^2 = 0,$$

541 and simultaneously

$$542 \quad (6.15) \quad \lim_{N \rightarrow \infty} \|F_N - (f - u)\|_{L^2(\mu_{\mathbf{x}})}^2 = \frac{1}{M} \sum_{i=1}^M [F_N(x_i) - (f(x_i) - u(x_i))]^2 = 0,$$

543 we have

$$544 \quad (6.16) \quad \int_{\Omega} \varrho^d(x) \mathcal{A}F^N(x) d\mu(x) + \int_{\Omega} \varrho^p(x) F^N(x) d\mu_{\mathbf{x}}(x) = 0.$$

545 This implies that

$$546 \quad (6.17) \quad \int_{\Omega} \varrho^d(x) \mathcal{A}[f - u](x) d\mu(x) + \int_{\Omega} \varrho^p(x) [f - u](x) d\mu_{\mathbf{x}}(x) = 0,$$

547 and hence (6.11) is proven. \square

548 **THEOREM 6.3** (Global convergence of PINNs objective). *Assume that the PDE*
 549 *(2.3) admits a continuous \mathcal{H}^2 solution u . The objective function $\|\mathcal{A}Q_t - g\|_{L^2(\mu)}^2 +$
 550 $\|Q_t - u\|_{L^2(\mu_x)}^2$ then converges to zero:*

$$551 \quad (6.18) \quad \lim_{t \rightarrow \infty} \left(\|\mathcal{A}Q_t - g\|_{L^2(\mu)}^2 + \|Q_t - u\|_{L^2(\mu_x)}^2 \right) = 0.$$

552 *Proof.* From (6.7) we see that

$$553 \quad (6.19) \quad \frac{d}{dt} \left[\|\mathcal{A}Q_t - g\|_{L^2(\mu)}^2 + \|Q_t - u\|_{L^2(\mu_x)}^2 \right] = - \begin{bmatrix} \mathcal{A}Q_t - g \\ \mathbf{Q}_t - \mathbf{u} \end{bmatrix}^\top \mathcal{V} \begin{bmatrix} \mathcal{A}Q_t - g \\ \mathbf{Q}_t - \mathbf{u} \end{bmatrix} \leq 0,$$

554 and that the equality holds iff $Q_t = u$. Therefore, the optimization objective is
 555 decreasing.

556 Consider $\tilde{\mathbf{Q}}_t = [\mathcal{A}Q_t - g, \mathbf{Q}_t - \mathbf{u}]^\top$ projected on $\{\vartheta\}_{i \in \mathbb{N}^+} \cup \{\varrho_i\}_{i \in \mathbb{N}^+}$. By Lemma
 557 6.1, we have

$$558 \quad (6.20) \quad \tilde{\mathbf{Q}}_t = \sum_i h_t^i \vartheta_i + \sum_j 0 \varrho_j = \sum_i h_t^i \vartheta_i.$$

559 Therefore $\|\mathcal{A}Q_t - g\|_{L^2(\mu)}^2 + \|Q_t - u\|_{L^2(\mu_x)}^2 = \sum_i h_t^{i2}$. Now consider its projection on
 560 each ϑ_i

$$561 \quad (6.21) \quad \frac{d}{dt} \langle \tilde{\mathbf{Q}}_t, \vartheta_i \rangle = \left\langle \frac{d\tilde{\mathbf{Q}}_t}{dt}, \vartheta_i \right\rangle = \langle -\mathcal{V}\tilde{\mathbf{Q}}_t, \vartheta_i \rangle = \langle \tilde{\mathbf{Q}}_t, -\mathcal{V}\vartheta_i \rangle = -\lambda_i \langle \tilde{\mathbf{Q}}_t, \vartheta_i \rangle.$$

562 Therefore

$$563 \quad (6.22) \quad \frac{d}{dt} \langle \tilde{\mathbf{Q}}_t, \vartheta_i \rangle = \frac{d}{dt} [h_t^i] = -\lambda_i h_t^i.$$

564 Consequently $h_t^i = h_0^i e^{-\lambda_i t}$ whose absolute value decays exponentially, and $|h_t^i| \leq |h_0^i|$
 565 for any $t \geq 0$ and any i .

566 Now, by the dominated convergence theorem

$$567 \quad (6.23) \quad \lim_{t \rightarrow \infty} \|\mathcal{A}Q_t - g\|_{L^2(\mu)}^2 + \|Q_t - u\|_{L^2(\mu_x)}^2 = \lim_{t \rightarrow \infty} \sum_i |h_t^i|^2 = \sum_i \lim_{t \rightarrow \infty} |h_t^i|^2 = 0. \quad \square$$

568 *Remark 6.4.* The above theorem makes the additional assumption that the solu-
 569 tion to our PDE is continuous. This is needed to ensure that the evaluation at the
 570 points x_i is meaningful, as strong solutions are generally only defined almost every-
 571 where. It is also worth noting that, if the solution to the PDE is not unique, then
 572 Theorem 6.3 allows the user to choose a solution which achieves specific values $u(x_i)$,
 573 which can be specified in advance (assuming such a solution exists).

574 **THEOREM 6.5** (Global convergence of PINNs). *If \mathcal{A} has an L^2 bounded inverse*
 575 *and the PDE (2.3) admits a continuous \mathcal{H}^2 solution, then Q_t converges to the solution*
 576 *of the PDE*

$$577 \quad \lim_{t \rightarrow \infty} \left(\|Q_t - u\|_{L^2(\mu)}^2 + \|Q_t - u\|_{L^2(\mu_x)}^2 \right) = 0.$$

578 *Proof.* As the existence of an inverse guarantees that the PDE (2.3) admits a
 579 (unique) solution, it is clear from Theorem 6.3 that $\|Q_t - u\|_{L^2(\mu_x)} \rightarrow 0$ and $\|\mathcal{A}Q_t -$
 580 $g\|_{L^2(\mu)}^2 \rightarrow 0$. As in the proof of Theorem 5.5, the convergence of the residual implies
 581 the convergence of Q_t to u , given \mathcal{A}^{-1} is a bounded operator. \square

582 **7. Conclusion.** In this paper, we develop a convergence theory for neural net-
583 work approximators of PDEs trained with gradient descent (e.g. DGM and PINNs).
584 It is proven that a neural network trained by corresponding algorithms to minimize
585 the PDE residual will converge to an infinite-dimensional linear ODE as the number
586 of hidden units $\rightarrow \infty$. The limit ODE's dynamics are characterized by a novel kernel
587 function involving the PDE operator and the neural network activation function. The
588 kernel lacks a spectral gap, making the analysis of the limit ODE challenging. Using
589 an eigendecomposition approach, we are able to prove that the PDE residual of the
590 limit neural network converges to zero. Furthermore, under mild assumptions, the
591 limit neural network converges to the solution of the PDE.

592 In this paper, we do not cover free-boundary problems nor non-linear PDEs,
593 which are both important types of problems in financial mathematics contexts. From
594 a numerical usage perspective, there already exist applications of PINNs/DGM type
595 algorithms on free-boundary problems or non-linear PDEs, see for instance, [37, 20].
596 However, in terms of proving a mathematical convergence for the algorithm, there
597 are challenges that must be addressed. For free boundary problems, it is unlikely to
598 match the boundary condition exactly when the algorithm starts, which is a crucial
599 point in our analysis. For non-linear problems, the kernel will become very different
600 from the kernel for linear PDEs and there will be unique mathematical challenges
601 which must be addressed.

602 **Acknowledgments.** SC acknowledges the support of the UKRI Prosperity Part-
603 nership Scheme (FAIR) under EPSRC Grant EP/V056883/1 and the "Mathemat-
604 ical Foundations of Intelligence: An 'Erlangen Programme' for AI" under grant
605 EP/Y028872/1. DJ's research was supported by the EPSRC Centre for Doctoral
606 Training in Industrially Focused Mathematical Modelling (EP/L015803/1).

607 **Appendix A. Proofs.**

608 **A.1. Proof of Theorem 3.8.** The main idea of the proof is the split the dif-
 609 ference into multiple residual terms and provide a bound for each of them. Then we
 610 apply Grönwall's inequality to provide an estimate of the difference between Q^N and
 611 Q . In this proof, constant C may vary from line to line, but it remains invariant with
 612 the number of neurons N and training time t . Expectations are taken with respect
 613 to random initialization of neural network parameters $\{c^i, w^i, b^i\}_{i=1, \dots, N}$.

614 We start with the following auxiliary lemmas which provide bounds for residual
 615 terms in our estimate. To simplify our notations, we may denote $\mathcal{L}Q := \mathcal{A}Q - g$.

616 **LEMMA A.1.** *We define residual term $M_1 :=$*

$$(A.1) \quad \int_0^T N^{2\beta+\delta-1} \left(\int_{\Omega^2} |\Phi^N(\nabla_\theta \mathcal{A}Q_s^N(x)) \cdot \sum_\alpha D_\alpha \nabla_\theta (Q_s^N - Q_0^N)(y)|^2 d\mu(x)d\mu(y) \right)^{\frac{1}{2}} ds,$$

618 *It satisfies*

$$(A.2) \quad \lim_{N \rightarrow \infty} \mathbb{E}M_1 = 0.$$

620 *Proof.* By Jensen's inequality, it suffices to prove that for any indices α

$$(A.3) \quad \int_0^T N^{2\beta+\delta-1} \left(\int \mathbb{E} \left[|\Phi^N(\nabla_\theta \mathcal{A}Q_s^N(x)) \cdot D_\alpha \nabla_\theta (Q_s^N - Q_0^N)(y)|^2 \right] d\mu(x)d\mu(y) \right)^{\frac{1}{2}} ds \rightarrow 0.$$

622 Notice that $|\Phi^N|$ is bounded elementwise by $N^{\epsilon-\beta}$. We also notice that by the mean
 623 value theorem and that η, σ , and their derivatives are up to polynomial growth

$$(A.4) \quad \|D_\alpha^y(\nabla_\theta Q_s^N - \nabla_\theta Q_0^N)(y)\|_1 \leq \frac{C}{N^\beta} \sum_{i=1}^N \left((\|w_s^i - w_0^i\| + \|b_s^i - b_0^i\| + \|c_s^i - c_0^i\|) f(w_0^i, c_0^i, b_0^i, y) \right)$$

625 where f is a polynomial up to 3rd order. We notice that $\|w_u^i - w_0^i\| + \|b_u^i - b_0^i\| +$
 626 $\|c_u^i - c_0^i\| \leq uCN^{\epsilon+\delta+\beta-1}$. Therefore

$$(A.5) \quad \mathbb{E} \left[|\Phi^N(\nabla_\theta \mathcal{A}Q_s^N(x)) \cdot D_\alpha(\nabla_\theta Q_s^N - \nabla_\theta Q_0^N)(y)|^2 \right] \leq CN^{2(2\epsilon+\delta-\beta)}.$$

628 Therefore, the left-hand side of (A.3) is smaller than $CN^{2\delta+2\epsilon+\beta-1}$, which goes to 0
 629 as $N \rightarrow \infty$. \square

630 **LEMMA A.2.** *We define residual term*

$$(A.6) \quad M_2 := \int_0^T N^{2\beta+\delta-1} \left(\int |\Phi^N(\nabla_\theta \mathcal{A}Q_s^N(x)) - \Phi^N(\nabla_\theta \mathcal{A}Q_0^N(x))| \cdot \sum_\alpha D_\alpha \nabla_\theta Q_0^N(y) \right)^2 d\mu(x)d\mu(y) \Big)^{\frac{1}{2}} ds$$

632 *It satisfies*

$$(A.7) \quad \lim_{N \rightarrow \infty} \mathbb{E}M_2 = 0.$$

634 *Proof.* Similarly, it suffices to show that for any α that

$$(A.8) \quad \int_0^T N^{2\beta+\delta-1} \left(\int_{\Omega^2} \mathbb{E} \left[|(\Phi^N(\nabla_\theta \mathcal{A}Q_s^N(x)) - \Phi^N(\nabla_\theta \mathcal{A}Q_0^N(x))) \cdot D_\alpha \nabla_\theta Q_0^N(y)|^2 \right] \right. \\ \left. d\mu(x)d\mu(y) \right)^{\frac{1}{2}} ds \rightarrow 0.$$

636 As, elementwise, $|\phi^N(x) - \phi^N(y)| \leq |x - y|$, we have

$$(A.9) \quad |(\Phi^N(\nabla_\theta \mathcal{A}Q_s^N(x)) - \Phi^N(\nabla_\theta \mathcal{A}Q_0^N(x))) \cdot D_\alpha \nabla_\theta Q_0^N(y)| \leq |\nabla_\theta \mathcal{A}Q_s^N(x) - \nabla_\theta \mathcal{A}Q_0^N(x)| \\ \cdot |D_\alpha \nabla_\theta Q_0^N(y)|$$

638 Then, since

$$(A.10) \quad |(\nabla_\theta \mathcal{A}Q_s^N(x) - \nabla_\theta \mathcal{A}Q_0^N(x)) \cdot D_\alpha \nabla_\theta Q_0^N(y)| \leq \frac{C}{N^{2\beta}} \sum_{i=1}^N \left((\|w_s^i - w_0^i\| + \|b_s^i - b_0^i\| \right. \\ \left. + \|c_s^i - c_0^i\|) g(w_0^i, c_0^i, b_0^i, x, y) \right)$$

640 where g is a polynomial up to 3rd order. Since $\|w_u^i - w_0^i\| + \|b_u^i - b_0^i\| + \|c_u^i - c_0^i\| \leq$
641 $uCN^{\epsilon+\delta+\beta-1}$, we have

$$(A.11) \quad \mathbb{E} \left[|(\Phi^N(\nabla_\theta \mathcal{A}Q_s^N(x)) - \Phi^N(\nabla_\theta \mathcal{A}Q_0^N(x))) \cdot D_\alpha \nabla_\theta Q_0^N(y)|^2 \right] \leq CN^{2(\epsilon+\delta-\beta)}.$$

643 The left-hand side of (A.8) is bounded by $CN^{\epsilon+\delta+\beta-1}$, which goes to 0 as N goes to
644 infinity. \square

645 **LEMMA A.3.** *We define residual term*

$$(A.12) \quad M_3 := \int_0^T N^{2\beta+\delta-1} \left(\int_{\Omega^2} |(\Phi^N(\nabla_\theta \mathcal{A}Q_0^N(x)) - \nabla_\theta \mathcal{A}Q_0^N(x)) \cdot \sum_\alpha D_\alpha \nabla_\theta Q_0^N(y)|^2 \right. \\ \left. d\mu(x)d\mu(y) \right)^{\frac{1}{2}} ds.$$

647 *We have*

$$(A.13) \quad \lim_{N \rightarrow \infty} \mathbb{E} M_3 = 0.$$

649 *Proof.* By Jensen's inequality, it suffices to show that

$$(A.14) \quad \int_0^T N^{2\beta+\delta-1} \left(\int_{\Omega^2} \mathbb{E} \left[|(\Phi^N(\nabla_\theta \mathcal{A}Q_0^N(x)) - \nabla_\theta \mathcal{A}Q_0^N(x)) \cdot D_\alpha \nabla_\theta Q_0^N(y)|^2 \right] \right. \\ \left. d\mu(x)d\mu(y) \right)^{\frac{1}{2}} ds \rightarrow 0$$

651 for any indices α where $|\alpha| \leq 2$. Let us denote the k -th unit of the initial approximator
 652 Q_0^N as q^k , where $q^k = c_0^k N^{-\beta} \eta \sigma_{w_0^k, b_0^k}$. Then $\mathbb{E}[|(\Phi^N(\nabla_\theta \mathcal{A} Q_0^N(x)) - \nabla_\theta \mathcal{A} Q_0^N(x)) \cdot$
 653 $D_\alpha \nabla_\theta Q_0^N(y)|^2]$

(A.15)

$$\begin{aligned} &= \mathbb{E} \left[\left(\sum_{k=1}^N [\Phi^N(\nabla_\theta \mathcal{A} q^k(x)) - \nabla_\theta \mathcal{A} q^k(x)] \cdot \nabla_\theta D_\alpha q^k(y) \right)^2 \right] \\ 654 &= N \mathbb{E} \left[(\Phi^N(\nabla_\theta \mathcal{A} q(x)) - \nabla_\theta \mathcal{A} q(x)) \cdot \nabla_\theta D_\alpha q(y) \right]^2 + N(N-1) \mathbb{E} \left[\Phi^N(\nabla_\theta \mathcal{A} q(x)) \right. \\ &\quad \left. - \nabla_\theta \mathcal{A} q(x) \right] \cdot \nabla_\theta D_\alpha q(y) \right]^2. \end{aligned}$$

655 Then, by definition

(A.16)

$$\begin{aligned} &\Phi^N(\nabla_\theta \mathcal{A} q(x)) - \nabla_\theta \mathcal{A} q(x) \cdot \nabla_\theta D_\alpha q(y) \\ &= [\phi^N(\mathcal{A}[N^{-\beta} \eta \sigma_{w,b}](x)) - \mathcal{A}[N^{-\beta} \eta \sigma_{w,b}](x)] \cdot D_\alpha [N^{-\beta} \eta \sigma_{w,b}](y) \\ 656 &+ \sum_{i=1}^d [\phi^N(\mathcal{A}[N^{-\beta} c x_i \eta \sigma'_{w,b}](x)) - \mathcal{A}[N^{-\beta} c x_i \eta \sigma'_{w,b}](x)] \cdot D_\alpha [N^{-\beta} c y_i \eta \sigma'_{w,b}](y) \\ &+ [\phi^N(\mathcal{A}[N^{-\beta} c \eta \sigma'_{w,b}](x)) - \mathcal{A}[N^{-\beta} c \eta \sigma'_{w,b}](x)] \cdot D_\alpha [N^{-\beta} c \eta \sigma'_{w,b}](y). \end{aligned}$$

657 We introduce a uniform bound for $|D_\alpha[\eta \sigma_{w,b}(x)]|$, $|D_\alpha[c x_i \eta \sigma_{w,b}(x)]|$ and also the term
 658 $|D_\alpha[c \eta \sigma_{w,b}(x)]|$ for any indices α and any x . By the fact that \mathcal{A} is Lipschitz,

$$659 \quad (\text{A.17}) \quad |\mathcal{A}[\eta \sigma_{w,b}](x)| \leq k_0 \sum_{0 \leq |\alpha| \leq 2} |D_\alpha[\eta \sigma_{w,b}](x)|.$$

660 Notice that

$$\begin{aligned} 661 \quad (\text{A.18}) \quad |D_\alpha[\eta \sigma_{w,b}](x)| &= \left| \sum_{\alpha_1 + \alpha_2 = \alpha} D_{\alpha_1} \eta \cdot D_{\alpha_2} \sigma_{w,b}(x) \right| \leq k_\eta \left| \sum_{\alpha_1 + \alpha_2 = \alpha} D_{\alpha_2} \sigma_{w,b}(x) \right| \\ &\leq k_\eta \left| \sum_{0 \leq |\alpha_2| \leq 2} D_{\alpha_2} \sigma_{w,b}(x) \right| \\ &\leq k_\eta k_\sigma k_\Omega \left(1 + \sum_{i=1}^d |(w)_i| + \sum_{i,j=1}^d |(w)_i (w)_j| \right). \end{aligned}$$

662 Therefore

(A.19)

$$663 \quad |\mathcal{A}[\eta \sigma_{w,b}](x)| \leq k_0 \sum_{0 \leq |\alpha| \leq 2} |D_\alpha[\eta \sigma_{w,b}](x)| \leq k \left(1 + \sum_{i=1}^d |(w)_i| + \sum_{i,j=1}^d |(w)_i (w)_j| \right).$$

664 Similar results hold for $|D_\alpha[c x_i \eta \sigma_{w,b}](x)|$ and $|D_\alpha[c \eta \sigma_{w,b}](x)|$. We define $f(w) := k(1 +$
 665 $\sum_{i=1}^d |(w)_i| + \sum_{i,j=1}^d |(w)_i (w)_j|)$. Now $|\phi^N(\mathcal{A}[N^{-\beta} \eta \sigma_{w,b}](x)) - \mathcal{A}[N^{-\beta} \eta \sigma_{w,b}](x)|$

$$\begin{aligned} 666 \quad (\text{A.20}) \quad &\leq N^{-\beta} \left[|\mathcal{A}[\eta \sigma_{w,b}](x)| - N^\epsilon \right] \mathbf{1}_{\{|\mathcal{A}[\eta \sigma_{w,b}](x)| \geq N^\epsilon\}}(x) \\ &\leq N^{-\beta} [f(w) - N^\epsilon] \mathbf{1}_{\{f(w) \geq N^\epsilon\}}(x). \end{aligned}$$

667 Subsequently,

(A.21)

$$668 \quad |\Phi^N(\nabla_\theta \mathcal{A}q(x)) - \nabla_\theta \mathcal{A}q(x)] \cdot \nabla_\theta D_\alpha q(y)| \leq N^{-2\beta}(d+2)[f(w) - N^\epsilon] \mathbf{1}_{\{f(w) \geq N^\epsilon\}}(x).$$

669 Therefore

(A.22)

$$670 \quad N \mathbb{E} \left[\left(\Phi^N(\nabla_\theta \mathcal{A}q(x)) - \nabla_\theta \mathcal{A}q(x) \right) \cdot \nabla_\theta D_\alpha q(y) \right]^2 \leq N^{1-4\beta}(d+2) \mathbb{E}[f(w)^2] \leq kN^{1-4\beta}.$$

671 Meanwhile $\mathbb{E}[\Phi^N(\nabla_\theta \mathcal{A}q(x)) - \nabla_\theta \mathcal{A}q(x)] \cdot \nabla_\theta D_\alpha q(y)$

(A.23)

$$672 \quad \leq \mathbb{E}[N^{-2\beta}(d+2)[f(w) - N^\epsilon] \mathbf{1}_{\{f(w) \geq N^\epsilon\}}(x)] \leq N^{-2\beta}(d+2) \mathbb{E}[f(w) \mathbf{1}_{\{f(w) \geq N^\epsilon\}}(x)] \\ \leq k_0 N^{-2\beta-\epsilon} \mathbb{E}[f(w)^2] \leq kN^{-2\beta-\epsilon}.$$

673 Hence

$$674 \quad (A.24) \quad N(N-1) \mathbb{E} \left[\Phi^N(\nabla_\theta \mathcal{A}q(x)) - \nabla_\theta \mathcal{A}q(x) \right]^2 \leq kN^{2-4\beta-2\epsilon}.$$

675 Combining (A.22) and (A.24), we conclude that

$$676 \quad (A.25) \quad M_3 \leq k(N^{\delta-1/2} + N^{\delta-\epsilon}) \rightarrow 0$$

677 as $N \rightarrow \infty$. □

678 LEMMA A.4. *We define residual term*

(A.26)

$$679 \quad M_4 := \mathbb{E} \left[\int_0^T N^\delta \left(\int_{\Omega^2} |N^{2\beta-1} \nabla_\theta \mathcal{A}Q_0^N(x) \cdot \sum_\alpha D_\alpha \nabla_\theta Q_0^N(y) - \sum_\alpha D_\alpha^y U(x, y)|^2 \right. \right. \\ \left. \left. d\mu(x) d\mu(y) \right)^{\frac{1}{2}} ds \right].$$

680 *It satisfies*

$$681 \quad (A.27) \quad \lim_{N \rightarrow \infty} \mathbb{E} M_4 = 0.$$

682 *Proof.* It suffices to prove that for any indices α

(A.28)

$$683 \quad \int_0^T N^\delta \left(\int \mathbb{E} \left[|N^{2\beta-1} \nabla_\theta \mathcal{A}Q_0^N(x) \cdot D_\alpha \nabla_\theta Q_0^N(y) - D_\alpha^y U(x, y)|^2 \right] d\mu(x) d\mu(y) \right)^{\frac{1}{2}} ds$$

684 converges to 0 as $N \rightarrow \infty$. By the strong law of large numbers:

$$685 \quad (A.29) \quad \lim_{N \rightarrow \infty} N^{2\beta-1} \nabla_\theta \mathcal{A}Q_0^N(x) \cdot D_\alpha \nabla_\theta Q_0^N(y) = D_\alpha^y U(x, y).$$

686 Therefore

$$687 \quad (A.30) \quad \mathbb{E} \left[|N^{2\beta-1} \nabla_\theta \mathcal{A}Q_0^N(x) \cdot D_\alpha \nabla_\theta Q_0^N(y) - D_\alpha^y U(x, y)|^2 \right] \\ = \text{Var} [N^{2\beta-1} \nabla_\theta \mathcal{A}Q_0^N(x) \cdot D_\alpha \nabla_\theta Q_0^N(y)] = \frac{1}{N} \text{Var} [D_\alpha^y U_{c,w,b}(x, y)].$$

688 Consequently, the left-hand side of (A.28) is bounded by $CN^{\delta-1}$ which vanishes as
689 N goes to infinity. □

690 LEMMA A.5. *We have*

$$691 \quad (\text{A.31}) \quad \lim_{N \rightarrow \infty} \int_0^T \int_{\Omega} |\psi^N(\mathcal{L}Q_s(x)) - \mathcal{L}Q_s(x)| d\mu(x) ds = 0.$$

692 *Proof.* Notice that

$$693 \quad (\text{A.32}) \quad |\psi^N(\mathcal{L}Q_s(x)) - \mathcal{L}Q_s(x)| \leq |\mathcal{L}Q_s(x)| \mathbf{1}_{\{|\mathcal{L}Q_s| > N^\delta\}}(x).$$

694 By the dominated convergence theorem, we conclude our proof. \square

695 LEMMA A.6. *We have*

$$696 \quad (\text{A.33}) \quad \lim_{N \rightarrow \infty} \mathbb{E} \left[\left(\int_{\Omega} \left[\sum_{\alpha} |D_{\alpha}(Q_0^N - Q_0)(y)| \right]^2 d\mu(y) \right)^{\frac{1}{2}} \right] = 0.$$

697 *Proof.* By Jensen's inequality, it suffices to show that for any α

$$698 \quad (\text{A.34}) \quad \int_{\Omega} \mathbb{E}[D_{\alpha}(Q_0^N - Q_0)(y)^2] d\mu(y) \rightarrow 0.$$

699 Since $\mathbb{E}[Q_0^N] = Q_0 = 0$ holds for all y we have

$$700 \quad (\text{A.35}) \quad \mathbb{E}[D_{\alpha}(Q_0^N - Q_0)(y)^2] = \text{Var}[D_{\alpha}Q_0^N(y)] = N^{1-2\beta} \text{Var}[D_{\alpha}c\sigma(wy + b)] \leq CN^{1-2\beta}.$$

701 Integrating with respect to $\mu(y)$, and letting N go to infinity finish the proof. \square

702 Now we present our main proof.

703 *Proof of Theorem 3.8.* In the integral form, we have

$$704 \quad (\text{A.36}) \quad Q_t^N(y) = Q_0^N(y) - \int_0^t \int_{\Omega} \psi^N(\mathcal{L}Q_s^N(x)) \alpha^N \Phi^N(\nabla_{\theta} \mathcal{A}Q_s^N(x)) \cdot \nabla_{\theta} Q_s^N(y) d\mu(x) ds.$$

705 Similarly, at the same time

$$706 \quad (\text{A.37}) \quad Q_t(y) = Q_0(y) - \int_0^t \int_{\Omega} \mathcal{L}Q_s(x) U(x, y) d\mu(x) ds.$$

707 Subtracting (A.37) from (A.36) and taking partial derivative with respect to y with
708 indices α give

$$709 \quad (\text{A.38}) \quad \begin{aligned} & |D_{\alpha}(Q_t^N - Q_t)(y)| \\ & \leq \int_0^t \int_{\Omega} \left| \psi^N(\mathcal{L}Q_s^N(x)) \alpha^N \Phi^N(\nabla_{\theta} \mathcal{A}Q_s^N(x)) \cdot D_{\alpha} \nabla_{\theta} Q_s^N(y) - \mathcal{L}Q_t^N(x) D_{\alpha}^y U(x, y) \right| \\ & \quad d\mu(x) ds + |D_{\alpha}(Q_0^N - Q_0)(y)|, \end{aligned}$$

710 where the right-hand side can be further bounded by

(A.39)

$$\begin{aligned} &\leq \int_0^t \int_{\Omega} \left| \psi^N(\mathcal{L}Q_s^N(x)) \alpha^N \Phi^N(\nabla_{\theta} \mathcal{A}Q_s^N(x)) \cdot D_{\alpha}(\nabla_{\theta} Q_s^N - \nabla_{\theta} Q_0^N)(y) \right| d\mu(x) ds \\ &\quad + \int_0^t \int_{\Omega} \left| \psi^N(\mathcal{L}Q_s^N(x)) \alpha^N (\Phi^N(\nabla_{\theta} \mathcal{A}Q_s^N(x)) - \Phi^N(\nabla_{\theta} \mathcal{A}Q_0^N(x))) \cdot \right. \\ &\quad \left. D_{\alpha} \nabla_{\theta} Q_0^N(y) \right| d\mu(x) ds \end{aligned}$$

$$\begin{aligned} 711 \quad &+ \int_0^t \int_{\Omega} \left| \psi^N(\mathcal{L}Q_s^N(x)) \alpha^N (\Phi^N(\nabla_{\theta} \mathcal{A}Q_0^N(x)) - \nabla_{\theta} \mathcal{A}Q_0^N(x)) \cdot D_{\alpha} \nabla_{\theta} Q_0^N(y) \right| d\mu(x) ds \\ &+ \int_0^t \int_{\Omega} \left| \psi^N(\mathcal{L}Q_s^N(x)) [\alpha^N \nabla_{\theta} \mathcal{A}Q_0^N(x) \cdot D_{\alpha} \nabla_{\theta} Q_0^N(y) - D_{\alpha}^y U(x, y)] \right| d\mu(x) ds \\ &+ \int_0^t \int_{\Omega} \left| [\psi^N(\mathcal{L}Q_s^N(x)) - \psi^N(\mathcal{L}Q_s(x))] D_{\alpha}^y U(x, y) \right| d\mu(x) ds \\ &+ \int_0^t \int_{\Omega} \left| [\psi^N(\mathcal{L}Q_s(x)) - \mathcal{L}Q_s(x)] D_{\alpha}^y U(x, y) \right| d\mu(x) ds + |D_{\alpha}(Q_0^N - Q_0)(y)|. \end{aligned}$$

712 Now by the fact that $|\psi^N| < N^{\delta}$, $|D_{\alpha}^y U(x, y)| < C$ we have

(A.40)

$$\begin{aligned} &|D_{\alpha}(Q_t^N - Q_t)(y)| \\ &\leq \int_0^t \int_{\Omega} N^{2\beta+\delta-1} \left| \Phi^N(\nabla_{\theta} \mathcal{A}Q_s^N(x)) \cdot D_{\alpha}(\nabla_{\theta} Q_s^N - \nabla_{\theta} Q_0^N)(y) \right| d\mu(x) ds \\ &\quad + \int_0^t \int_{\Omega} N^{2\beta+\delta-1} \left| (\Phi^N(\nabla_{\theta} \mathcal{A}Q_s^N(x)) - \Phi^N(\nabla_{\theta} \mathcal{A}Q_0^N(x))) \cdot D_{\alpha} \nabla_{\theta} Q_0^N(y) \right| d\mu(x) ds \\ 713 \quad &+ \int_0^t \int_{\Omega} N^{2\beta+\delta-1} \left| (\Phi^N(\nabla_{\theta} \mathcal{A}Q_0^N(x)) - \nabla_{\theta} \mathcal{A}Q_0^N(x)) \cdot D_{\alpha} \nabla_{\theta} Q_0^N(y) \right| d\mu(x) ds \\ &+ \int_0^t \int_{\Omega} N^{\delta} \left| N^{2\beta-1} \nabla_{\theta} \mathcal{A}Q_0^N(x) \cdot D_{\alpha} \nabla_{\theta} Q_0^N(y) - D_{\alpha}^y U(x, y) \right| d\mu(x) ds \\ &+ C \int_0^t \int_{\Omega} \left| \psi^N(\mathcal{L}Q_s^N(x)) - \psi^N(\mathcal{L}Q_s(x)) \right| d\mu(x) ds \\ &+ C \int_0^t \int_{\Omega} \left| \psi^N(\mathcal{L}Q_s(x)) - \mathcal{L}Q_s(x) \right| d\mu(x) ds + |D_{\alpha}(Q_0^N - Q_0)(y)|. \end{aligned}$$

714 Let us denote $V_t^N(y) := \sum_{|\alpha| \leq 2} |D_{\alpha}(Q_t^N - Q_t)(y)|$. Then by summing inequality
715 (A.40) with respect to all indices and integrating with respect to $\mu(y)$

(A.41)

$$716 \quad \int_{\Omega} V_t^N(y)^2 d\mu(y) \leq \left(\int_{\Omega} V_t^N(y)^2 d\mu(y) \right)^{\frac{1}{2}} \left[C \int_0^t \left(\int_{\Omega} V_s^N(x)^2 d\mu(x) \right)^{\frac{1}{2}} ds + M \right]$$

717 where $M := M_1 + M_2 + M_3 + M_4 + M_5 + M_6$ denotes the sum of residual terms where

718 M_1 to M_4 have been defined in previous lemmas and

$$719 \quad (A.42) \quad \begin{aligned} M_5 &:= \int_0^T \int_{\Omega} |\psi^N(\mathcal{L}Q_s(x)) - \mathcal{L}Q_s(x)| d\mu(x) ds, \\ M_6 &:= \left(\int_{\Omega} \left[\sum_{\alpha} |D_{\alpha}(Q_0^N - Q_0)(y)| \right]^2 d\mu(y) \right)^{\frac{1}{2}}. \end{aligned}$$

720 Then by Grönwall's inequality, since from (A.41) we have

$$721 \quad (A.43) \quad \left(\int_{\Omega} V_t^N(y)^2 d\mu(y) \right)^{\frac{1}{2}} \leq C \int_0^t \left(\int_{\Omega} V_s^N(x)^2 d\mu(x) \right)^{\frac{1}{2}} ds + M,$$

722 we derive that

$$723 \quad (A.44) \quad \int_0^t \left(\int_{\Omega} V_s^N(x)^2 d\mu(x) \right)^{\frac{1}{2}} ds \leq \frac{M}{C} e^{Ct}.$$

724 Taking expectation on both sides, by Lemmas A.1, A.2, A.3, A.4, A.5, A.6 we can
725 control each of the terms in M . Hence

$$726 \quad (A.45) \quad \lim_{N \rightarrow \infty} \mathbb{E} \left[\int_0^t \left(\int_{\Omega} V_s^N(x)^2 d\mu(x) \right)^{\frac{1}{2}} ds \right] = 0.$$

727 By (A.43) we have for any $0 \leq s \leq T$

$$728 \quad (A.46) \quad \lim_{N \rightarrow \infty} \mathbb{E} \left[\left(\int_{\Omega} V_s^N(x)^2 d\mu(x) \right)^{\frac{1}{2}} \right] = 0.$$

729 Notice that $\int_{\Omega} V_s^N(x)^2 d\mu(x) \geq \|Q_s^N - Q_s\|_{\mathcal{H}^2}^2$, we conclude that for any $0 \leq s \leq T$

$$730 \quad (A.47) \quad \lim_{N \rightarrow \infty} \mathbb{E} \left[\|Q_s^N - Q_s\|_{\mathcal{H}^2} \right] = 0. \quad \square$$

731

REFERENCES

- 732 [1] C. BECK, M. HUTZENTHALER, A. JENTZEN, AND B. KUCKUCK, *An overview on deep*
733 *learning-based approximation methods for partial differential equations*, arXiv preprint
734 arXiv:2012.12348, (2020).
735 [2] J. BERNER, P. GROHS, AND A. JENTZEN, *Analysis of the generalization error: Empirical risk*
736 *minimization over deep artificial neural networks overcomes the curse of dimensionality in*
737 *the numerical approximation of black-scholes partial differential equations*, SIAM Journal
738 on Mathematics of Data Science, 2 (2020).
739 [3] M. BICHUCH AND K. CHEN, *A deep learning scheme for solving fully nonlinear partial differ-*
740 *ential equation*, Peter Carr Gedenkschrift: Research Advances in Mathematical Finance,
741 (2023), pp. 101–140.
742 [4] R. CARMONA AND M. LAURIÈRE, *Convergence analysis of machine learning algorithms for the*
743 *numerical solution of mean field control and games i: The ergodic case*, SIAM Journal on
744 Numerical Analysis, 59 (2021).
745 [5] ———, *Deep learning for mean field games and mean field control with applications to finance*,
746 arXiv:2107.04568, (2021).
747 [6] ———, *Convergence analysis of machine learning algorithms for the numerical solution of mean*
748 *field control and games: II—the finite horizon case*, The Annals of Applied Probability,
749 32 (2022), pp. 4065–4105.

- 750 [7] C. BECK, S. BECKER, P. CHERIDITO, A. JENTZEN, AND A. NEUFELD, *Deep splitting method for*
751 *parabolic PDEs*, SIAM Journal on Scientific Computing, 43 (2021).
- 752 [8] S. N. COHEN, D. JIANG, AND J. SIRIGNANO, *Neural Q-learning for solving elliptic PDEs*, arXiv
753 preprint arXiv:2203.17128, (2022).
- 754 [9] T. DE RYCK, A. D. JAGTAP, AND S. MISHRA, *Error estimates for physics-informed neural*
755 *networks approximating the navier–stokes equations*, IMA Journal of Numerical Analysis,
756 44 (2024), pp. 83–119.
- 757 [10] N. DOUMÈCHE, G. BLAU, AND C. BOYER, *Convergence and error analysis of PINNs*, arXiv
758 preprint arXiv:2305.01240, (2023).
- 759 [11] D. ELBRÄCHTER, P. GROHS, A. JENTZEN, AND C. SCHWAB, *DNN expression rate analysis of*
760 *high-dimensional PDEs: application to option pricing*, Constructive Approximation, 55
761 (2022).
- 762 [12] L. C. EVANS, *Partial differential equations*, American Mathematical Society, second ed., 2010.
- 763 [13] C. GAO, S. GAO, AND Z. ZHU, *Convergence of the backward deep BSDE method with applica-*
764 *tions to optimal stopping problems*, SIAM Journal on Financial Mathematics, 14 (2023).
- 765 [14] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep learning*, MIT press, 2016.
- 766 [15] P. GROHS, F. HORNING, A. JENTZEN, AND P. V. WURSTEMBERGER, *A proof that artificial*
767 *neural networks overcome the curse of dimensionality in the numerical approximation of*
768 *black-scholes partial differential equations*, Memoirs of the American Mathematical Society,
769 284 (2023).
- 770 [16] P. GROHS, A. JENTZEN, AND D. SALIMOVA, *Deep neural network approximations for solutions*
771 *of PDEs based on monte carlo algorithms*, Partial Differential Equations and Applications,
772 3 (2022).
- 773 [17] J. HAN, A. JENTZEN, AND W. E, *Solving high-dimensional partial differential equations using*
774 *deep learning*, Proceedings of the National Academy of Sciences, 115 (2018).
- 775 [18] K. HORNIK, *Approximation capabilities of multilayer feedforward networks*, Neural networks,
776 4 (1991), pp. 251–257.
- 777 [19] A. JACOT, F. GABRIEL, AND C. HONGLER, *Neural tangent kernel: Convergence and general-*
778 *ization in neural networks*, Advances in neural information processing systems, 31 (2018).
- 779 [20] X. JIANG, D. WANG, Q. FAN, M. ZHANG, C. LU, AND A. P. T. LAU, *Physics-informed neural*
780 *network for nonlinear dynamics in fiber optics*, Laser & Photonics Reviews, 16 (2022),
781 p. 2100483.
- 782 [21] I. E. LAGARIS, A. LIKAS, AND D. I. FOTIADIS, *Artificial neural networks for solving ordinary*
783 *and partial differential equations*, IEEE transactions on neural networks, 9 (1998), pp. 987–
784 1000.
- 785 [22] H. LEE AND I. S. KANG, *Neural algorithm for solving differential equations*, Journal of Com-
786 putational Physics, 91 (1990), pp. 110–131.
- 787 [23] Y. LU, J. BLANCHET, AND L. YING, *Sobolev acceleration and statistical optimality for learning*
788 *elliptic equations via gradient descent*, Advances in Neural Information Processing Systems,
789 35 (2022), pp. 33233–33247.
- 790 [24] T. LUO AND H. YANG, *Two-layer neural networks for partial differential equations: Optimiza-*
791 *tion and generalization theory*, arXiv preprint arXiv:2006.15733, (2020).
- 792 [25] A. MALEK AND R. S. BEIDOKHTI, *Numerical solution for high order differential equations using*
793 *a hybrid neural network—optimization method*, Applied Mathematics and Computation,
794 183 (2006), pp. 260–271.
- 795 [26] K. S. MCFALL AND J. R. MAHAN, *Artificial neural network method for solution of boundary*
796 *value problems with exact satisfaction of arbitrary boundary conditions*, IEEE Transactions
797 on Neural Networks, 20 (2009), pp. 1221–1233.
- 798 [27] S. MISHRA AND R. MOLINARO, *Estimates on the generalization error of physics-informed neural*
799 *networks for approximating PDEs.*, IMA Journal of Numerical Analysis, 43 (2023).
- 800 [28] R. PASCANU, T. MIKOLOV, AND Y. BENGIO, *On the difficulty of training recurrent neural*
801 *networks*, in International conference on machine learning, Pmlr, 2013, pp. 1310–1318.
- 802 [29] M. RAISSI, P. PERDIKARIS, AND G. E. KARNIADAKIS, *Physics-informed neural networks: A deep*
803 *learning framework for solving forward and inverse problems involving nonlinear partial*
804 *differential equations*, Journal of Computational physics, 378 (2019), pp. 686–707.
- 805 [30] A. REPPEN, H. SONER, AND V. TISSOT-DAGUETTE, *Deep stochastic optimization in finance*,
806 Digital Finance, 5 (2022).
- 807 [31] K. RUDD, *Solving partial differential equations using artificial neural networks*, PhD thesis,
808 Duke University, 2013.
- 809 [32] Y. SAPORITO AND Z. ZHANG, *Path-dependent deep Galerkin method: A neural network ap-*
810 *proach to solve path-dependent partial differential equations*, SIAM Journal on Financial
811 Mathematics, 12 (2021).

- 812 [33] Y. SHIN, J. DARBON, AND G. E. KARNIADAKIS, *On the convergence of physics informed neu-*
813 *ral networks for linear second-order elliptic and parabolic type PDEs*, arXiv preprint
814 arXiv:2004.01806, (2020).
- 815 [34] J. SIRIGNANO AND K. SPILIOPOULOS, *DGM: A deep learning algorithm for solving partial dif-*
816 *ferential equations*, Journal of computational physics, 375 (2018), pp. 1339–1364.
- 817 [35] ———, *Mean field analysis of neural networks: A law of large numbers*, SIAM Journal on
818 Applied Mathematics, 80 (2020).
- 819 [36] ———, *Asymptotics of reinforcement learning with neural networks*, Stochastic Systems, 12
820 (2022).
- 821 [37] S. WANG AND P. PERDIKARIS, *Deep learning of free boundary and stefan problems*, Journal of
822 Computational Physics, 428 (2021), p. 109914.
- 823 [38] S. WANG, X. YU, AND P. PERDIKARIS, *When and why PINNs fail to train: A neural tangent*
824 *kernel perspective*, Journal of Computational Physics, 449 (2022), p. 110768.
- 825 [39] J. ZHANG, T. HE, S. SRA, AND A. JADBABAIE, *Why gradient clipping accelerates training: A*
826 *theoretical justification for adaptivity*, arXiv preprint arXiv:1905.11881, (2019).