

Mean Field Analysis of Deep Neural Networks

Justin Sirignano* and Konstantinos Spiliopoulos^{†‡}

October 20, 2020

Abstract

We analyze multi-layer neural networks in the asymptotic regime of simultaneously (A) large network sizes and (B) large numbers of stochastic gradient descent training iterations. We rigorously establish the limiting behavior of the multi-layer neural network output. The limit procedure is valid for any number of hidden layers and it naturally also describes the limiting behavior of the training loss. The ideas that we explore are to (a) take the limits of each hidden layer sequentially and (b) characterize the evolution of parameters in terms of their initialization. The limit satisfies a system of deterministic integro-differential equations. The proof uses methods from weak convergence and stochastic analysis. We show that, under suitable assumptions on the activation functions and the behavior for large times, the limit neural network recovers a global minimum (with zero loss for the objective function).

1 Introduction

Machine learning, and in particular deep learning, has achieved immense success, revolutionizing fields such as image, text, and speech recognition. It is also increasingly being used in engineering, medicine, and finance. However, despite their success in practice, there is currently limited mathematical understanding of deep neural networks. This has motivated recent mathematical research on multi-layer learning models such as [39], [40], [41], [20], [21], [42], [49], [50], [43], and [48].

Neural networks are nonlinear statistical models whose parameters are estimated from data using stochastic gradient descent (SGD) methods. Deep learning uses neural networks with many layers (i.e., “deep” neural networks), which produces a highly flexible, powerful and effective model in practice. Typically, a neural network with multiple layers between the input and the output layer is called a “deep” neural network, see for example [24]. We analyze multi-layer neural networks that have a fixed number of layers between the input and output layer, and where the number of hidden units in each layer becomes large.

Applications of deep learning include image recognition (see [35] and [24]), facial recognition [59], driverless cars [6], speech recognition (see [35], [4], [36], and [60]), and text recognition (see [62] and [57]). Neural networks have also been applied in engineering, robotics, medicine, and finance (see [37], [38], [58], [26], [47], [3], [51], [52], [53], and [54]).

In this paper we characterize multi-layer neural networks in the asymptotic regime of large network sizes and large numbers of stochastic gradient descent iterations. We rigorously prove the limit of the neural network output as the number of hidden units increases to infinity. The proof relies upon weak convergence analysis for stochastic processes. The result can be considered a “law of large numbers” for the neural network’s output when both the network size and the number of stochastic gradient descent steps grow to infinity. We show that the neural network output in the large hidden-units and large SGD-iterates limit depends on paths of representative weights that go from input to output layer. This result is then used to show that, under suitable assumptions, the limit neural network seeks to minimize the limit objective function and achieve zero loss.

Recently, law of large numbers and central limit theorems have been established for neural networks with a single hidden layer [10, 30, 43, 48, 49, 50]. For a single hidden layer, one can directly study the weak convergence of the empirical measure of the parameters. However, in a neural network with multiple

*Department of Industrial & Systems Engineering, University of Illinois at Urbana-Champaign, Urbana, E-mail: jasirign@illinois.edu

[†]Department of Mathematics and Statistics, Boston University, Boston, E-mail: kspiliop@math.bu.edu

[‡]K.S. was partially supported by the National Science Foundation (DMS 1550918) and Simons Foundation Award 672441

layers, there is a closure problem when studying the empirical measure of the parameters (which is explained in Section 4.3). Consequently, the law of large numbers for a multi-layer network is not a straightforward extension of the single-layer network result and the analysis involves unique challenges which require new approaches. In this paper we establish the limiting behavior of the output of the neural network.

To illustrate the idea, we consider a multi-layer neural network with two hidden layers:

$$g_{\theta}^{N_1, N_2}(x) = \frac{1}{N_2} \sum_{i=1}^{N_2} C^i \sigma \left(\frac{1}{N_1} \sum_{j=1}^{N_1} W^{2,i,j} \sigma(W^{1,j} \cdot x) \right). \quad (1.1)$$

As we will see in Section 4.2, the limit procedure can be extended to neural networks with three layers and subsequently to neural networks with any fixed number of hidden layers.

Notice now that (1.1) can be also written as

$$\begin{aligned} H^{1,j}(x) &= \sigma(W^{1,j} \cdot x), \quad j = 1, \dots, N_1, \\ Z^{2,i}(x) &= \frac{1}{N_1} \sum_{j=1}^{N_1} W^{2,i,j} H^{1,j}(x), \quad i = 1, \dots, N_2, \\ H^{2,i}(x) &= \sigma \left(Z^{2,i}(x) \right), \\ g_{\theta}^{N_1, N_2}(x) &= \frac{1}{N_2} \sum_{i=1}^{N_2} C^i H^{2,i}(x). \end{aligned} \quad (1.2)$$

where $C^i, W^{2,i,j} \in \mathbb{R}$ and $x, W^{1,j} \in \mathbb{R}^d$. The neural network model has parameters

$$\theta = (C^1, \dots, C^{N_2}, W^{2,1,1}, \dots, W^{2,N_1,N_2}, W^{1,1}, \dots, W^{1,N_1}),$$

which must be estimated from data. The number of hidden units in the first layer is N_1 and the number of hidden units in the second layer is N_2 . The multi-layer neural network (1.2) includes a normalization factor of $\frac{1}{N_1}$ in the first hidden layer and $\frac{1}{N_2}$ in the second hidden layer.

The loss function that we focus on in this paper is the mean squared error

$$L^{N_1, N_2}(\theta) = \frac{1}{2} \mathbb{E}_{Y, X} \left[(Y - g_{\theta}^{N_1, N_2}(X))^2 \right], \quad (1.3)$$

where the data $(X, Y) \sim \pi(dx, dy)$. The goal is to estimate a set of parameters θ which minimizes the objective function (1.3).

The literature frequently refers to minimizing the mean squared error loss as regression. Our results also hold for a more general class of error functions. In particular, we could have also considered the error function $\mathbb{E}_{Y, X} \Psi(Y - g_{\theta}^{N_1, N_2}(X))$ for a function Ψ that is smooth, convex and satisfies $\min_{x \in \mathbb{R}} \Psi(x) = \Psi(0) = 0$. However, for the purposes of simplicity, we will focus on the standard regression task (1.3).

The stochastic gradient descent (SGD) algorithm for estimating the parameters θ is, for $k \in \mathbb{N}$,

$$\begin{aligned}
C_{k+1}^i &= C_k^i + \frac{\alpha_C^{N_1, N_2}}{N_2} (y_k - g_{\theta_k}^{N_1, N_2}(x_k)) H_k^{2,i}(x_k), \\
W_{k+1}^{1,j} &= W_k^{1,j} + \frac{\alpha_{W,1}^{N_1, N_2}}{N_1} (y_k - g_{\theta_k}^{N_1, N_2}(x_k)) \left(\frac{1}{N_2} \sum_{i=1}^{N_2} C_k^i \sigma'(Z_k^{2,i}(x_k)) W_k^{2,i,j} \right) \sigma'(W_k^{1,j} \cdot x_k) x_k, \\
W_{k+1}^{2,i,j} &= W_k^{2,i,j} + \frac{\alpha_{W,2}^{N_1, N_2}}{N_1 N_2} (y_k - g_{\theta_k}^{N_1, N_2}(x_k)) C_k^i \sigma'(Z_k^{2,i}(x_k)) H_k^{1,j}(x_k), \\
H_k^{1,i}(x_k) &= \sigma(W_k^{1,i} \cdot x_k), \\
Z_k^{2,i}(x_k) &= \frac{1}{N_1} \sum_{j=1}^{N_1} W_k^{2,i,j} H_k^{1,j}(x_k), \\
H_k^{2,i}(x_k) &= \sigma(Z_k^{2,i}(x_k)), \\
g_{\theta_k}^{N_1, N_2}(x_k) &= \frac{1}{N_2} \sum_{i=1}^{N_2} C_k^i H_k^{2,i}(x_k).
\end{aligned} \tag{1.4}$$

where $\alpha_C^{N_1, N_2}$, $\alpha_{W,1}^{N_1, N_2}$, and $\alpha_{W,2}^{N_1, N_2}$ are the learning rates. The learning rates may depend upon N_1 and N_2 . The parameters at step k are $\theta_k = (C_k^1, \dots, C_k^{N_2}, W_k^{2,1,1}, \dots, W_k^{2,N_1, N_2}, W_k^{1,1}, \dots, W_k^{1, N_1})$. (x_k, y_k) are samples of the random variables (X, Y) .

The goal of this paper is to characterize the limit of an appropriate rescaling of the multi-layer neural network output $g_{\theta_k}^{N_1, N_2}(x)$ as both the number of hidden units (N_1, N_2) and the stochastic gradient descent iterates k become large. This is the topic of Theorem 2.3. The idea is to first take $N_1 \rightarrow \infty$ with N_2 fixed. In Lemma 2.2, we prove that the empirical measure of the parameters converges to a limit measure as $N_1 \rightarrow \infty$ (with N_2 fixed) which satisfies a measure evolution equation. This naturally implies a limit for the neural network output g^{N_1, N_2} as $N_1 \rightarrow \infty$. The next step is to take $N_2 \rightarrow \infty$. Theorem 2.3 proves that the limiting distribution can be represented via a system of ODEs.

As previously discussed, related limiting results for the single-layer neural network case have been investigated in [10, 30, 43, 48, 49, 50]. In those papers, it is proven that as the number of hidden units and stochastic gradient descent steps, in the appropriate scaling, diverge to infinity, the empirical distribution of the neural network parameters converges to the weak solution of a non-local PDE. This non-local PDE turns out to be a gradient flow for the limiting objective function in the space of probability measures endowed with the Wasserstein metric (this result is analogous to our Theorem 3.4). More results with non-asymptotic bounds, again for the single-layer neural network, can be found in [44].

Let us also mention that after our current paper appeared as a preprint on arXiv, [2] studied the limit of a multi-layer neural network, but under some important differences as compared to our paper. [45] also derives asymptotics for multi-layer neural networks. In [2], the weights in the first and last hidden layers are held fixed throughout training, and the number of hidden units in the rest of the layers is sent to infinity. In our paper, we train all parameters in all layers of the neural network, which introduces additional technical challenges. There is also the related papers [9, 31], whose authors look again at the limit as the units per layer go to infinity, but under a $\frac{1}{\sqrt{N}}$ normalization instead of our $\frac{1}{N}$ normalization. In the $\frac{1}{\sqrt{N}}$ normalization, a completely different limit equation will appear. The $\frac{1}{\sqrt{N}}$ case results in a perturbation around the network's randomized initialization, leading to a kernel-type regression limit.

We address the multi-layer neural network case in the mean field scaling of $\frac{1}{N}$. In particular, we study the behavior of the neural network's output (a) as the number of hidden units in each layer go to infinity one by one, and (b) as the number of stochastic gradient descent iterates (i.e., during training of the network) goes to infinity at the speed of the number of the hidden units of the first layer. In this asymptotic regime, we are able to obtain a well-defined limit in Theorem 2.3. Consequences of this result along with a simulation study are presented in Sections 3 and 4.

The rest of the paper is organized as follows. Our main result, which characterizes the asymptotic behavior of a neural network with two hidden layers when the number of hidden units becomes large, is presented in Section 2. The result can be easily extended to an arbitrary number of hidden layers. Section

3 is devoted to the global convergence arguments. In particular, we show that under the proper assumption the limiting problem derived in Section 2 seeks to minimize the limiting objective function, recovering the global minimum. Section 4 discusses further the theoretical results, includes a numerical study to showcase some of the theoretical implications, and, as an example, presents the limit for a three-layer neural network. The proof of the convergence theorem is in Section 5. The uniqueness of a solution to the limiting system is established in Section 6. The proof of the limit of the first layer (i.e., the proof of Lemma 2.2) and a few other results are provided in the Appendix. Section 7 has concluding remarks and possible directions for future work.

2 Main Results

Let us start by presenting our assumptions, which will hold throughout the paper. We shall work on a filtered probability space $(\Omega, \mathcal{F}, \mathbb{P})$ on which all the random variables are defined. The probability space is equipped with a filtration \mathfrak{F}_t that is right continuous and \mathfrak{F}_0 contains all \mathbb{P} -negligible sets.

Assumption 2.1. We assume the following conditions throughout the paper.

- $\sigma(\cdot) \in C_b^2$, i.e., it is twice continuously differentiable and bounded.
- The distribution $\pi(dx, dy)$ has compact support, i.e. the data (x_k, y_k) takes values in the compact set $\mathcal{X} \times \mathcal{Y}$.
- The random initialization of the parameters, i.e. $\{C_\circ^i\}_i, \{W_\circ^{2,i,j}\}_{i,j}, \{W_\circ^{1,j}\}_j$, are i.i.d. and take values in compact sets $\mathcal{C}, \mathcal{W}^1$, and \mathcal{W}^2 .
- The probability distributions of the initial parameters $(C_\circ^i, W_\circ^{2,i,j}, W_\circ^{1,j})_{i,j}$ admit continuous probability density functions.

We denote by $\mu_c(dc)$, $\mu_{W^2}(du)$, and $\mu_{W^1}(dw)$ the probability distributions of $\{C_\circ^i\}_i$, $\{W_\circ^{2,i,j}\}_{i,j}$, and $\{W_\circ^{1,j}\}_j$ respectively.

For reasons that will become clearer later on, we shall choose the learning rates to be

$$\alpha_C^{N_1, N_2} = \frac{N_2}{N_1}, \quad \alpha_{W^1}^{N_1, N_2} = 1 \text{ and } \alpha_{W^2}^{N_1, N_2} = N_2. \quad (2.1)$$

Note that the weights in the second layer are trained faster than the other parameters. This choice of learning rates is necessary for convergence to a non-trivial limit as $N_1, N_2 \rightarrow \infty$. If the parameters in all the layers are trained with the same learning rate, it can be mathematically shown that the network will not train as N_1, N_2 become large. We further explore this interesting fact in Section 4.1.

Define the empirical measure

$$\tilde{\gamma}_k^{N_1, N_2} := \frac{1}{N_1} \sum_{j=1}^{N_1} \delta_{W_k^{1,j}, W_k^{2,1,j}, \dots, W_k^{2, N_2, j}, C_k^1, \dots, C_k^{N_2}}. \quad (2.2)$$

If f is an appropriate test function on some space X and γ is a finite measure on X , then we denote the inner product $\langle f, \gamma \rangle = \int_X f(x) \gamma(dx)$. Using this inner product, the neural network's output can be re-written in terms of the empirical measure:

$$g_{\theta_k}^{N_1, N_2}(x) = \frac{1}{N_2} \sum_{i=1}^{N_2} \langle c_i, \tilde{\gamma}_k^{N_1, N_2} \rangle \sigma \left(\langle w^{2,i} \sigma(w^1 \cdot x), \tilde{\gamma}_k^{N_1, N_2} \rangle \right).$$

Let us next define the time-scaled empirical measure

$$\gamma_t^{N_1, N_2} := \tilde{\gamma}_{\lfloor N_1 t \rfloor}^{N_1, N_2},$$

and the corresponding time-scaled neural network output is

$$g_t^{N_1, N_2}(x) := g_{\theta_{\lfloor N_1 t \rfloor}}^{N_1, N_2}(x) = \frac{1}{N_2} \sum_{i=1}^{N_2} \langle c_i, \gamma_t^{N_1, N_2} \rangle \sigma \left(\langle w^{2,i} \sigma(w^1 \cdot x), \gamma_t^{N_1, N_2} \rangle \right).$$

At any time t , $\gamma_t^{N_1, N_2}$ is measure-valued. The scaled empirical measure $(\gamma_t^{N_1, N_2})_{0 \leq t \leq 1}$ is a random element of $D_E([0, 1])^1$ with $E = \mathcal{M}(\mathbb{R}^{d+2N_2})$.

We study convergence using iterated limits. We first let $N_1 \rightarrow \infty$ where the number of units in the first layer is N_1 and the number of stochastic gradient descent steps is $\lfloor N_1 \rfloor$. Then, we let the number of units in the second layer $N_2 \rightarrow \infty$.

We begin by letting the number of hidden units in the first layer $N_1 \rightarrow \infty$.

Lemma 2.2. The process $\gamma^{N_1, N_2} := (\gamma_t^{N_1, N_2})_{0 \leq t \leq 1}$ converges in distribution to the measure valued process γ^{N_2} that takes values in $D_E([0, 1])$ as $N_1 \rightarrow \infty$. For every $f \in C_b^2(\mathbb{R}^{d+2N_2})$, γ^{N_2} satisfies the measure evolution equation

$$\begin{aligned} \langle f, \gamma_t^{N_2} \rangle - \langle f, \gamma_0^{N_2} \rangle &= \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} (y - g_s^{N_2}(x)) \langle H_s^{N_2}(x) \cdot \nabla_c f, \gamma_s^{N_2} \rangle \pi(dx, dy) ds \\ &+ \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} (y - g_s^{N_2}(x)) \langle \sigma(w^1 \cdot x) (\sigma'(Z_s(x)) \odot c) \cdot \nabla_{w^2} f, \gamma_s^{N_2} \rangle \pi(dx, dy) ds \\ &+ \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} (y - g_s^{N_2}(x)) \frac{1}{N_2} \sum_{i=1}^{N_2} \langle c_i \sigma'(Z_s^{i, N_2}(x)) w^{2, i} \sigma'(w^1 \cdot x) x \cdot \nabla_{w^1} f, \gamma_s^{N_2} \rangle \pi(dx, dy) ds, \end{aligned} \quad (2.3)$$

where

$$Z_s^{i, N_2}(x) = \langle w^{2, i} \sigma(w^1 \cdot x), \gamma_s^{N_2} \rangle,$$

$$H_s^{i, N_2}(x) = \sigma(Z_s^{i, N_2}(x)),$$

$$g_s^{N_2}(x) = \frac{1}{N_2} \sum_{i=1}^{N_2} \langle c_i, \gamma_s^{N_2} \rangle H_s^{i, N_2}(x),$$

$$\gamma_0^{N_2}(dw^1, dw^2, dc) = \mu_{W^1}(dw^1) \times \mu_{W^2}(dw^{2,1}) \times \cdots \times \mu_{W^2}(dw^{2, N_2}) \times \delta_{C_0^1}(dc^1) \times \cdots \times \delta_{C_0^{N_2}}(dc^{N_2}). \quad (2.4)$$

Proof. The proof of this lemma is related to the limit in the first layer as the number of hidden units in the first layer grows with the number of hidden units in the second layer is held fixed. The proof is analogous to the proof in [49] and the details are presented for completeness in the Appendix A. \square

Lemma 2.2 studies the limit of the empirical measure $\gamma_t^{N_1, N_2}$ as $N_1 \rightarrow \infty$ with N_2 fixed. The limit is characterized by the stochastic evolution equation (2.3)-(2.4). Notice that Lemma 2.2 immediately implies that

$$\lim_{N_1 \rightarrow \infty} g_t^{N_1, N_2}(x) = g_t^{N_2}(x),$$

in probability, as $N_1 \rightarrow \infty$

The next step is to study the limit as $N_2 \rightarrow \infty$. To do so, we study the limit of the random ODE as $N_2 \rightarrow \infty$ whose law is characterized by (2.3)-(2.4). Our main goal is the characterization of the limit neural network output $g_t^{N_1, N_2}(x)$. The following convergence result characterizes the neural network output $g_t^{N_1, N_2}(x)$ for large N_1 and N_2 .

Theorem 2.3. For any $t \in [0, 1]$ and $x \in \mathcal{X}$,

$$\lim_{N_2 \rightarrow \infty} \lim_{N_1 \rightarrow \infty} g_t^{N_1, N_2}(x) = g_t(x),$$

in probability², where we have that

$$g_t(x) = \int_{\mathcal{C}} \tilde{C}_t^c \tilde{H}_t^{2, c}(x) \mu_c(dc), \quad (2.5)$$

¹ $D_S([0, 1])$ is the set of maps from $[0, 1]$ into S which are right-continuous and which have left-hand limits.

² $\lim_{N_2 \rightarrow \infty} \lim_{N_1 \rightarrow \infty} X^{N_1, N_2} = X$ in probability if, for all $\epsilon > 0$, $\lim_{N_2 \rightarrow \infty} \lim_{N_1 \rightarrow \infty} \mathbb{P} \left[\|X^{N_1, N_2} - X\| > \epsilon \right] = 0$.

with

$$\begin{aligned}
d\tilde{C}_t^c &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t(x)) \tilde{H}_t^{2,c}(x) \pi(dx, dy) dt, & \tilde{C}_0^c &= c, \\
d\tilde{W}_t^{1,w} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t(x)) V_t^w(x) \sigma'(\tilde{W}_t^{1,w} \cdot x) x \pi(dx, dy) dt, & \tilde{W}_0^{1,w} &= w, \\
d\tilde{W}_t^{2,c,w,u} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t(x)) \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \tilde{H}_t^{1,w}(x) \pi(dx, dy) dt, & \tilde{W}_0^{2,c,w,u} &= u, \\
\tilde{H}_t^{1,w}(x) &= \sigma(\tilde{W}_t^{1,w} \cdot x), \\
\tilde{Z}_t^c(x) &= \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}_t^{2,c,w,u} \tilde{H}_t^{1,w}(x) \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw), \\
\tilde{H}_t^{2,c}(x) &= \sigma(\tilde{Z}_t^c(x)), \\
V_t^w(x) &= \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \left(\int_{\mathcal{W}^2} \tilde{W}_t^{2,c,w,u} \mu_{\mathcal{W}^2}(du) \right) \mu_{\mathcal{C}}(dc).
\end{aligned} \tag{2.6}$$

The system in (2.6) has a unique solution. In addition, letting $g_t^{N_2}(x)$ defined through Lemma 2.2 we have the following rate of convergence

$$\sup_{x \in \mathcal{X}} \mathbb{E} \left[|g_t^{N_2}(x) - g_t(x)| \right] \leq K N_2^{-1/2},$$

for some constant $K < \infty$.

Notice that we can also write that $g_t(x)$ satisfies

$$g_t(x) = \int_{\mathcal{C}} \tilde{C}_t^c \sigma \left(\int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}_t^{2,c,w,u} \sigma(\tilde{W}_t^{1,w} \cdot x) \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) \right) \mu_{\mathcal{C}}(dc). \tag{2.7}$$

The proof of Theorem 2.3 is given in Section 5. Theorem 2.3 indicates that the neural network output in the large hidden units and large SGD-iterates limit depends on paths of representative weights that connect the input layer to the output layer. Even though, we restricted the statement of Theorem 2.3 in the interval $t \in [0, 1]$, its proof makes it clear that the statement is true for $t \in [0, T]$ for any $0 < T < \infty$. In addition, even though this does not mean that we can take $T = \infty$, we can still examine what happens to the limit problem as t grows. In particular, in Section 3 we show that under the proper assumptions, one does expect to recover the global minimum as $t \rightarrow \infty$.

Section 4 discusses some further consequences of Theorem 2.3 as well as challenges that come up in the study of the limiting behavior of multi-layer neural networks as the number of the hidden units grows.

3 On global convergence

The goal of this section is to demonstrate that, under appropriate assumptions, it can be expected that the global minimum is recovered as $t \rightarrow \infty$. Let the target data y be a function $f(x)$ that we seek to learn, i.e. $y = f(x)$.

Assumption 3.1. The activation function $\sigma(\cdot)$ is real analytic, bounded, and $\sigma'(\cdot) > 0$.

Notice that by [29], the fact that σ is bounded and non-constant by Assumption 3.1 implies that σ is also discriminatory in the sense of [11, 29]. Namely, if we have that

$$\int_{\mathcal{X}} h(x) \sigma(w \cdot x) \pi(dx) = 0 \text{ a.e. in } w \in \mathbb{R}^d,$$

then $h(x) = 0$ a.e. in \mathcal{X} .

Remark 3.2. An example of a real analytic, bounded activation function where $\sigma'(\cdot) > 0$ (and therefore discriminatory too), i.e. that satisfies Assumption 3.1, is the sigmoid function $\sigma(z) = \frac{e^z}{1+e^z}$.

Theorem 2 of [23] shows that for activation functions σ that are non-constant, bounded and monotone (in which case Assumption 3.1 holds), then two layer neural networks, such as $g_\theta^{N_1, N_2}(x)$, (and more generally multilayer neural networks) are dense in $\mathcal{C}(\mathcal{X})$. This result basically implies that for large enough N_1, N_2 functions of the form $g_\theta^{N_1, N_2}(x)$ approximate to arbitrary accuracy functions in $\mathcal{C}(\mathcal{X})$. Training the weights with stochastic gradient descent and sending N_1, N_2 to infinity, leads then, by Theorem 2.3, to function approximator $g_t(x)$, which is given through a system of integro-differential equations of the form

$$\begin{aligned}
d\tilde{C}_t^c &= \int_{\mathcal{X}} (f(x) - g_t(x)) \tilde{H}_t^{2,c}(x) \pi(dx) dt, & \tilde{C}_0^c &= c, \\
d\tilde{W}_t^{1,w} &= \int_{\mathcal{X}} (f(x) - g_t(x)) V_t^w \sigma'(\tilde{W}_t^{1,w} \cdot x) x \pi(dx) dt, & \tilde{W}_0^{1,w} &= w, \\
d\tilde{W}_t^{2,c,w,u} &= \int_{\mathcal{X}} (f(x) - g_t(x)) \tilde{C}_t^c \sigma'(\tilde{Z}_t^c) \tilde{H}_t^{1,w}(x) \pi(dx) dt, & \tilde{W}_0^{2,c,w,u} &= u, \\
\tilde{H}_t^{1,w}(x) &= \sigma(\tilde{W}_t^{1,w} \cdot x), \\
\tilde{Z}_t^c(x) &= \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}_t^{2,c,w,u} \tilde{H}_t^{1,w}(x) \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw), \\
\tilde{H}_t^{2,c}(x) &= \sigma(\tilde{Z}_t^c(x)), \\
V_t^w(x) &= \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \left(\int_{\mathcal{W}^2} \tilde{W}_t^{2,c,w,u} \mu_{\mathcal{W}^2}(du) \right) \mu_{\mathcal{C}}(dc),
\end{aligned} \tag{3.1}$$

where the neural network prediction is

$$g_t(x) = \int_{\mathcal{C}} \tilde{C}_t^c \tilde{H}_t^{2,c}(x) \mu_{\mathcal{C}}(dc), \tag{3.2}$$

We establish in this section that under reasonable assumptions $f(x)$ can indeed be recovered by $g_t(x)$ as $t \rightarrow \infty$.

Let us denote $\Theta_t(c, w, u) = (\tilde{C}_t^c, \tilde{W}_t^{1,w}, \tilde{W}_t^{2,c,w,u})$ for the components of the ODE in (3.1). For notational convenience, we shall often write $\theta = (c, w, u)$. Then, we obviously have that $g_t(x)$ depends on t only through $[\Theta_t] = \{\Theta_t(\theta)\}_{\theta \in \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2}$. In order to emphasize that we shall write $g_t(x) = g(x; [\Theta_t])$.

Analogously, we will denote $[\Theta_t(c, \cdot)] = \{\Theta_t(c, w, u)\}_{(w,u) \in \mathcal{W}^1 \times \mathcal{W}^2}$ and $[\Theta_t(\cdot, w, \cdot)] = \{\Theta_t(c, w, u)\}_{(c,u) \in \mathcal{C} \times \mathcal{W}^2}$, which then leads to the notation $\tilde{Z}_t^c(x) = \tilde{Z}(x; [\Theta_t(c, \cdot)])$ and $V_t^w(x) = V(x; [\Theta_t(\cdot, w, \cdot)])$. For notational convenience let us also denote $h(x; [\Theta_t]) = (f(x) - g_t(x))$.

With these definitions in place, and for $(c, w) \in \mathcal{C} \times \mathcal{W}^1$, let us also define

$$\begin{aligned}
R_1([\Theta_t], c) &\equiv \int_{\mathcal{X}} h(x; [\Theta_t]) \sigma(\tilde{Z}(x; [\Theta_t(c, \cdot)])) \pi(dx) \\
R_2([\Theta_t], \tilde{W}_t^{1,w}, w) &\equiv \int_{\mathcal{X}} h(x; [\Theta_t]) V(x; [\Theta_t(\cdot, w, \cdot)]) \sigma'(\tilde{W}_t^{1,w} \cdot x) x \pi(dx) \\
R_3([\Theta_t], \tilde{C}_t^c, \tilde{W}_t^{1,w}, c) &\equiv \int_{\mathcal{X}} h(x; [\Theta_t]) \tilde{C}_t^c \sigma'(\tilde{Z}(x; [\Theta_t(c, \cdot)])) \sigma(\tilde{W}_t^{1,w} \cdot x) \pi(dx)
\end{aligned} \tag{3.3}$$

and let us set

$$H([\Theta_t], \Theta_t(\theta), \theta) = (R_1([\Theta_t], c), R_2([\Theta_t], \tilde{W}_t^{1,w}, w), R_3([\Theta_t], \tilde{C}_t^c, \tilde{W}_t^{1,w}, c)).$$

The purpose of the notation above is to make it clear that the functions $H(\cdot)$ depends on $[\Theta_t]$, on $\Theta_t(\theta)$ and on θ separately. With these identifications we can write that the ODE system in (3.1) can be written in the form

$$\dot{\Theta}_t(\Theta_0) = H([\Theta_t], \Theta_t(\Theta_0), \Theta_0), \text{ such that } \Theta_0 = (c, w, u). \tag{3.4}$$

Next we investigate the limiting objective function

$$\begin{aligned}\bar{L}(\Theta_t) &= \lim_{N_2 \rightarrow \infty} \lim_{N_1 \rightarrow \infty} L^{N_1, N_2}(\theta_{[N_1 t]}) = \lim_{N_2 \rightarrow \infty} \lim_{N_1 \rightarrow \infty} \frac{1}{2} \mathbb{E}_X \left[(f(X) - g_{\theta_{[N_1 t]}}^{N_1, N_2}(X))^2 \right] \\ &= \frac{1}{2} \int_{\mathcal{X}} \left[(f(x) - g(x; [\Theta_t]))^2 \right] \pi(dx)\end{aligned}$$

where $g(x; [\Theta_t]) = g_t(x)$ is given by (3.2).

Notice that the function $\Theta \mapsto \bar{L}(\Theta)$ is always non-negative and becomes zero only at a global minimum. In particular, our goal is to demonstrate that as $t \rightarrow \infty$, there are sufficient conditions that guarantee that the global minimum is realized in the sense that

$$\lim_{t \rightarrow \infty} g(x; [\Theta_t]) = f(x), \text{ for almost all } x \in \mathcal{X}.$$

We notice that \bar{L} acts as a Lyapunov function for the dynamical system (3.4) in that

$$\frac{d}{dt} \bar{L}(\Theta_t) = \nabla_{\Theta} \bar{L}(\Theta_t) \cdot \dot{\Theta}_t = \nabla_{\Theta} \bar{L}(\Theta_t) \cdot H(\Theta_t) \leq 0$$

Indeed, we calculate

$$\begin{aligned}\frac{d}{dt} \bar{L}(\Theta_t) &= \int_{\mathcal{X}} (f(x) - g_t(x)) \dot{g}_t(x) \pi(dx) \\ &= \int_{\mathcal{X}} (f(x) - g_t(x)) \left[\int_{\mathcal{C}} \left(\dot{\tilde{C}}_t^c \tilde{H}_t^{2,c}(x) + \right. \right. \\ &\quad \left. \left. + \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \left(\int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \left(\dot{\tilde{W}}_t^{2,c,w,u} \tilde{H}_t^{1,w}(x) + \tilde{W}_t^{2,c,w,u} \sigma'(\tilde{W}_t^{1,w} \cdot x) \dot{\tilde{W}}_t^{1,w} \cdot x \right) \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) \right) \right) \mu_c(dc) \right] \pi(dx) \\ &= \int_{\mathcal{X}} (f(x) - g_t(x)) \int_{\mathcal{C}} \dot{\tilde{C}}_t^c \tilde{H}_t^{2,c}(x) \mu_c(dc) \pi(dx) \\ &\quad + \int_{\mathcal{X}} (f(x) - g_t(x)) \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}_t^{2,c,w,u} \sigma'(\tilde{W}_t^{1,w} \cdot x) \dot{\tilde{W}}_t^{1,w} \cdot x \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) \mu_c(dc) \pi(dx) \\ &\quad + \int_{\mathcal{X}} (f(x) - g_t(x)) \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \dot{\tilde{W}}_t^{2,c,w,u} \tilde{H}_t^{1,w}(x) \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) \mu_c(dc) \pi(dx) \\ &= \int_{\mathcal{C}} \dot{\tilde{C}}_t^c \left[\int_{\mathcal{X}} (f(x) - g_t(x)) \tilde{H}_t^{2,c}(x) \pi(dx) \right] \mu_c(dc) \\ &\quad + \int_{\mathcal{W}^1} \dot{\tilde{W}}_t^{1,w} \cdot \left[\int_{\mathcal{X}} (f(x) - g_t(x)) \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \int_{\mathcal{W}^2} \tilde{W}_t^{2,c,w,u} \sigma'(\tilde{W}_t^{1,w} \cdot x) x \mu_{\mathcal{W}^2}(du) \mu_c(dc) \pi(dx) \right] \mu_{\mathcal{W}^1}(dw) \\ &\quad + \int_{\mathcal{C}} \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \dot{\tilde{W}}_t^{2,c,w,u} \left[\int_{\mathcal{X}} (f(x) - g_t(x)) \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \tilde{H}_t^{1,w}(x) \pi(dx) \right] \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) \mu_c(dc) \\ &= \int_{\mathcal{C}} \dot{\tilde{C}}_t^c \left[\int_{\mathcal{X}} (f(x) - g_t(x)) \tilde{H}_t^{2,c}(x) \pi(dx) \right] \mu_c(dc) \\ &\quad + \int_{\mathcal{W}^1} \dot{\tilde{W}}_t^{1,w} \cdot \left[\int_{\mathcal{X}} (f(x) - g_t(x)) V_t^w \sigma'(\tilde{W}_t^{1,w} \cdot x) x \pi(dx) \right] \mu_{\mathcal{W}^1}(dw) \\ &\quad + \int_{\mathcal{C}} \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \dot{\tilde{W}}_t^{2,c,w,u} \left[\int_{\mathcal{X}} (f(x) - g_t(x)) \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \tilde{H}_t^{1,w}(x) \pi(dx) \right] \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) \mu_c(dc) \\ &= - \int_{\mathcal{C}} \left(\int_{\mathcal{X}} (f(x) - g_t(x)) \tilde{H}_t^{2,c}(x) \pi(dx) \right)^2 \mu_c(dc) \\ &\quad - \int_{\mathcal{W}^1} \left| \int_{\mathcal{X}} (f(x) - g_t(x)) V_t^w \sigma'(\tilde{W}_t^{1,w} \cdot x) x \pi(dx) \right|^2 \mu_{\mathcal{W}^1}(dw) \\ &\quad - \int_{\mathcal{C}} \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \left(\int_{\mathcal{X}} (f(x) - g_t(x)) \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \tilde{H}_t^{1,w}(x) \pi(dx) \right)^2 \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) \mu_c(dc)\end{aligned}$$

$$\begin{aligned}
&= - \left[\int_{\mathcal{C}} (R_1([\Theta_t], c))^2 \mu_c(dc) + \int_{\mathcal{W}^1} \left| R_2([\Theta_t], \tilde{W}_t^{1,w}, w) \right|^2 \mu_{\mathcal{W}^1}(dw) \right. \\
&\quad \left. + \int_{\mathcal{C}} \int_{\mathcal{W}^1} \left(R_3([\Theta_t], \tilde{C}_t^c, \tilde{W}_t^{1,w}, c) \right)^2 \mu_{\mathcal{W}^1}(dw) \mu_c(dc) \right] \\
&\leq 0.
\end{aligned} \tag{3.5}$$

The fact that $\frac{d}{dt} \bar{L}(\Theta_t) \leq 0$ from (3.5) means that $\bar{L}(\Theta_t)$ is decreasing in the gradient direction of the paths governing the limiting behavior of the weights.

Let us define η_0 to be the joint measure for the random initialization of the parameters, i.e. for $\{C_\circ^i\}_i$, $\{W_\circ^{1,j}\}_j$, $\{W_\circ^{2,i,j}\}_{i,j}$. By Assumption 2.1 we have that this is the product measure $\eta_0 = \mu_c \times \mu_{\mathcal{W}^1} \times \mu_{\mathcal{W}^2}$. Let us make a convenient assumption for the support of η_0 . In particular,

Assumption 3.3. We have that $\text{support}(\eta_0) = \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2$.

Notice now that we can write

$$\frac{d}{dt} \bar{L}(\Theta_t) = - \int_{\mathcal{C}} \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} |H([\Theta_t], \Theta_t(\theta), \theta)|^2 \eta_0(d\theta) \leq 0. \tag{3.6}$$

As a matter of fact, more is true. We investigate how the objective function changes with perturbation from Θ to $\hat{\Theta}$. In particular, we have for $\epsilon = (\epsilon_1, \epsilon_2, \epsilon_3)$

$$\left. \frac{\partial}{\partial \epsilon_i} \bar{L}(\Theta + \epsilon \odot \hat{\Theta}) \right|_{\epsilon=0} = \int_{\mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2} H_i([\Theta], \Theta(\theta), \theta) \hat{\Theta}_i(\theta) \eta_0(d\theta),$$

which, in combination with (3.6) and Assumption 3.3, then says that a minimizer for $\bar{L}(\Theta)$, say $\Theta^*(\theta)$, will satisfy for almost all $(c, w) \in \mathcal{C} \times \mathcal{W}^1$ the relations

$$R_1([\Theta^*], c) = R_2([\Theta^*], \tilde{W}^{1,*}, w) = R_3([\Theta^*], \tilde{C}^{*,c}, \tilde{W}^{1,*}, w) = 0. \tag{3.7}$$

By analogy, let us now define η_t to be the measure at time t of the random vector $\Theta_t = (\tilde{C}_t, \tilde{W}_t^1, \tilde{W}_t^2)$ as governed by the solution to the random ODE system (3.1). Then, η_t is the pushforward of η_0 under Θ_t given by (3.4) (see for example Chapter 8 in [1]), i.e.,

$$\eta_t = (\Theta_t)_\# \eta_0. \tag{3.8}$$

Theorem 3.4. Let assumptions 3.1 and 3.3 hold. If $\eta_t \rightarrow \eta^*$, weakly, where η^* is a non-degenerate measure that admits a density with finite first moments, then we have that η^* is a global minimum with zero loss.

Before proving Theorem 3.4 we discuss its assumptions in the remarks that follow.

Remark 3.5. It is interesting to note that Theorem 3.4 actually shows that any stationary point is a global minimum with zero loss. In addition, even though we do not show this here, one can see from the proof of Theorem 3.4 that the assumption on analyticity of the activation function can be replaced by an assumption on η^* having full support.

Remark 3.6. The assumption on convergence of η_t as $t \rightarrow \infty$ is a required assumption for convergence in the limit as $t \rightarrow \infty$ to occur. For completeness we mention here that verifying this assumption for the (commonly used in practice) unregularized case that we study here is still an open problem even for the corresponding unregularized case in the single layer case, see Appendix C.5 in [10] or [43]. Even though we do not show this here, one way to guarantee that such an assumption holds is to add an appropriate regularization term in the loss function, or to add an appropriate non-degenerate stochastic perturbation in the training algorithm (1.4) for the evolution of each one of $C_{k+1}^i, W_{k+1}^{1,j}, W_{k+1}^{2,i,j}$ (noisy SGD); see [30, 48] for the corresponding arguments in the single layer case. See also [8] for conditions under which the assumption on convergence of η_t as $t \rightarrow \infty$ holds for non-degenerate sparse problems that include an appropriate regularizing term in the loss function.

Remark 3.7. It is clear that convergence of $[\Theta_t]$ as $t \rightarrow \infty$ is linked to the convergence of the measure η_t as $t \rightarrow \infty$. At this point, we also mention that convergence of the trajectories $[\Theta_t]$ as $t \rightarrow \infty$ is linked to the

study of Łojasiewicz inequalities [7, 27]. To the best of our knowledge no general results in this direction are currently known neither for our case of interest, nor for problems corresponding to the single-layer case (see Appendix C.5 of [10]) (i.e. no general results currently exist for the non-geodesically convex case that we are dealing with here).

Proof of Theorem 3.4. We will use the notation $\tilde{\mu}_c(c)$, $\tilde{\mu}_{W^1}(w)$, and $\tilde{\mu}_{W^2}(w)$ for the densities of the measure $\mu_c(dc)$, $\mu_{W^1}(dw)$, and $\mu_{W^2}(dw)$ respectively. From equations (3.7), we have that any stationary solution must satisfy

$$\begin{aligned} 0 &= \int_{\mathcal{C}} \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \left(\int_{\mathcal{X}} (f(x) - g(x; [\Theta^*])) \tilde{C}^{*,c} \sigma'(\tilde{Z}^{*,c}(x)) \sigma(\tilde{W}^{1,*},w \cdot x) \pi(dx) \right)^2 \mu_{W^2}(du) \mu_{W^1}(dw) \mu_c(dc) \\ &= \int_{\mathcal{C}} \int_{\mathcal{W}^1} \left(\int_{\mathcal{X}} (f(x) - g(x; [\Theta^*])) \tilde{C}^{*,c} \sigma'(\tilde{Z}^{*,c}(x)) \sigma(\tilde{W}^{1,*},w \cdot x) \pi(dx) \right)^2 \tilde{\mu}_{W^1}(w) \tilde{\mu}_c(c) dw dc. \end{aligned} \quad (3.9)$$

Equation (3.9) and Assumption 3.3 imply that

$$\int_{\mathcal{X}} (f(x) - g(x; [\Theta^*])) \tilde{C}^{*,c} \sigma'(\tilde{Z}^{*,c}(x)) \sigma(\tilde{W}^{1,*},w \cdot x) \pi(dx) \tilde{\mu}_{W^1}(w) \tilde{\mu}_c(c) = 0, \quad (3.10)$$

for almost every $w \in \mathcal{W}^1$ and $c \in \mathcal{C}$. Recalling now that

$$Z^{*,c}(x) = \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}^{2,*},c,w,u \sigma(\tilde{W}^{1,*},w \cdot x) \mu_{W^2}(du) \mu_{W^1}(dw),$$

we obtain

$$\begin{aligned} \int_{\mathcal{C}} \sup_{x \in \mathcal{X}} |Z^{*,c}(x)| \mu_c(dc) &\leq \int_{\mathcal{C}} \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} |\tilde{W}^{2,*},c,w,u| \mu_{W^2}(du) \mu_{W^1}(dw) \mu_c(dc) \\ &\leq \int_{\mathbb{R}^d} |u| \tilde{\eta}_{W^2}^*(u) du < \infty, \end{aligned}$$

where $\tilde{\eta}_{W^2}^*(u)$ denotes the marginal density of the limiting measure η^* with respect to \tilde{W}^2 . Analogously we shall define the marginal densities $\tilde{\eta}_c^*(c)$ and $\tilde{\eta}_{W^1}^*(w)$ with respect to the variables \mathcal{C} and \tilde{W}^1 respectively. Similar to before, $\tilde{\eta}^*(c, w, u)$ is the density of the measure $\eta^*(dc, dw, du)$.

We must therefore have that $\sup_{x \in \mathcal{X}} |Z^{*,c}(x)| \tilde{\mu}_c(c) < \infty$ for almost every $c \in \mathcal{C}$. In addition, since η^* has a density, we have that there exists a set $B = \{c \in \mathbb{R} : |c - c_0| < \delta\}$ with $c_0 \neq 0$ and $0 < \delta < |c_0|$ such that

$$\begin{aligned} 0 &< \int_{\mathbb{R}} \mathbf{1}_{\{c \in B\}} \tilde{\eta}_c^*(c) dc \\ &= \int_{\mathcal{C}} \mathbf{1}_{\{C^{*,c} \in B\}} \tilde{\mu}_c(c) dc \\ &= \int_{\mathcal{C}} \mathbf{1}_{\{C^{*,c} \in B, \sup_{x \in \mathcal{X}} |Z^{*,c}(x)| \tilde{\mu}_c(c) < \infty\}} \tilde{\mu}_c(c) dc. \end{aligned}$$

Therefore, there must be a set $K \subset \mathbb{R}$ with Lebesgue measure $\lambda(K) > 0$ such that, for almost every $c \in K$,

$$\mathbf{1}_{\{C^{*,c} \in B, \sup_{x \in \mathcal{X}} |Z^{*,c}(x)| \tilde{\mu}_c(c) < \infty\}} \tilde{\mu}_c(c) > 0, \quad (3.11)$$

which in turn implies that $C^{*,c} \neq 0$ and $\sup_{x \in \mathcal{X}} |Z^{*,c}(x)| < \infty$ for almost every $c \in K$. Due to equations (3.11) and (3.10), for almost every $(c, w) \in K \times \mathcal{W}^1$,

$$\begin{aligned} C^{*,c} &\neq 0, \\ \sup_{x \in \mathcal{X}} |Z^{*,c}(x)| &< \infty, \\ 0 &= \int_{\mathcal{X}} (f(x) - g(x; [\Theta^*])) \tilde{C}^{*,c} \sigma'(\tilde{Z}^{*,c}(x)) \sigma(\tilde{W}^{1,*},w \cdot x) \pi(dx) \tilde{\mu}_{W^1}(w) \tilde{\mu}_c(c). \end{aligned} \quad (3.12)$$

Thus, there exists a $c_0 \in K$ such that

$$\int_{\mathcal{X}} (f(x) - g(x; [\Theta^*])) C^{*,c_0} \sigma'(\tilde{Z}^{*,c_0}(x)) \sigma(\tilde{W}^{1,*,w} \cdot x) \pi(dx) \tilde{\mu}_{W^1}(w) = 0,$$

for almost every $w \in \mathcal{W}^1$. This of course implies that, for almost every $w \in \mathcal{W}^1$,

$$\int_{\mathcal{X}} (f(x) - g(x; [\Theta^*])) \sigma'(\tilde{Z}^{*,c_0}(x)) \sigma(\tilde{W}^{1,*,w} \cdot x) \pi(dx) \tilde{\mu}_{W^1}(w) = 0.$$

Consequently, we have that

$$\int_{\mathcal{W}^1} \int_{\mathcal{X}} (f(x) - g(x; [\Theta^*])) \sigma'(\tilde{Z}^{*,c_0}(x)) \sigma(\tilde{W}^{1,*,w} \cdot x) \pi(dx) \tilde{\mu}_{W^1}(w) dw = 0.$$

Therefore, since $\eta_{W^1}^*$ has a density, there exists a compact set $A \subset \mathbb{R}^d$ such that, for every $w \in A$,

$$\Gamma(w) = \int_{\mathcal{X}} (f(x) - g(x; [\Theta^*])) \sigma'(\tilde{Z}^{*,c_0}(x)) \sigma(w \cdot x) \pi(dx) = 0. \quad (3.13)$$

Since $\sigma(\cdot)$ is analytic, $\Gamma(w)$ is analytic (see Lemma 3.8). Therefore, by the identity theorem for real-analytic functions, $\Gamma(w) = 0$ for every $w \in \mathbb{R}^d$. It then follows, since $\sigma(\cdot)$ is discriminatory,

$$(f(x) - g(x; [\Theta^*])) \sigma'(\tilde{Z}^{*,c_0}(x)) = 0,$$

for every $x \in \mathcal{X}$. Finally, since $\sigma'(\tilde{Z}^{*,c_0}(x)) > 0$ due to equation (3.12),

$$f(x) - g(x; [\Theta^*]) = 0,$$

for every $x \in \mathcal{X}$, concluding the proof of the theorem. \square

We conclude with the the result on analyticity of $\Gamma(w)$ defined in (3.13).

Lemma 3.8. Assume \mathcal{X} is a compact set and the activation unit $\sigma(\cdot)$ is real analytic and bounded. Then,

$$\Gamma(w) = \int_{\mathcal{X}} (f(x) - g(x; [\Theta^*])) \sigma'(\tilde{Z}^{*,c_0}(x)) \sigma(w \cdot x) \pi(dx)$$

defined in (3.13) is a real analytic function as well.

Proof of Lemma 3.8. The partial derivative of $\sigma(w \cdot x)$ is

$$\frac{\partial^{|\mu|}}{\partial w^\mu} \sigma(w \cdot x) = \sigma^{|\mu|}(w \cdot x) \prod_{i=1}^d x^{\mu_i}.$$

Before we prove analyticity of $\Gamma(w)$, we show that $g(x)$ is finite. Recall that

$$g(x; [\Theta^*]) = \int_{\mathcal{C}} C^{*,c} \sigma(Z^{*,c}(x)) \tilde{\mu}_c(c) dc.$$

Therefore, due to $\sigma(\cdot)$ being bounded and η^* having marginals with finite first moments,

$$\begin{aligned} |g(x)| &\leq K_0 \int_{\mathcal{C}} |C^{*,c}| \tilde{\mu}_c(c) dc \\ &\leq K_0 \int_{\mathcal{C}} |c| \tilde{\eta}_c^*(c) dc \\ &\leq K. \end{aligned}$$

Next,

$$\frac{\partial^{|\mu|} \Gamma}{\partial w^\mu}(w) = \int_{\mathcal{X}} (f(x) - g(x; [\Theta^*])) \sigma'(\tilde{Z}^{*,c_0}(x)) \sigma^{|\mu|}(w \cdot x) \prod_{i=1}^d x^{\mu_i} \pi(dx).$$

Due to the compactness of \mathcal{X} , $\sup_{x \in \mathcal{X}} |\sigma'(\tilde{Z}^{*,c_0}(x))| < \infty$, and the fact that σ is a real analytic function, Proposition 2.2.10 of [32] shows there exist constants $C_1, C_2 > 0$ such that

$$\begin{aligned} \left| \frac{\partial^{|\mu|} \Gamma}{\partial w^\mu}(w) \right| &\leq C_1 \frac{\mu!}{R_1^{|\mu|}} \prod_{i=1}^d C_2^{\mu_i} \\ &= C_1 \frac{\mu!}{R_1^{|\mu|}} C_2^{|\mu|} \\ &= C_1 \frac{\mu!}{R_1^{|\mu|}} (R_2^{-1})^{|\mu|} \\ &= C_1 \frac{\mu!}{(R_1 R_2)^{|\mu|}} \\ &= C_1 \frac{\mu!}{R^{|\mu|}}, \end{aligned}$$

where $R_2 = \frac{1}{C_2}$ and $R = R_1 R_2$. It follows then again from Proposition 2.2.10 of [32] that $\Gamma(w)$ is a real analytic function. This concludes the proof of the lemma. \square

4 Discussion on the limiting results and extensions to multi-layer networks with greater depth

In Subsection 4.1, we discuss some of the implications of our theoretical convergence results and presents related numerical results. In Section 4.2, we show that the procedure can be extended to treat deep neural networks with more than two hidden layers. General challenges in the study of multi-layer neural networks are explored Subsection 4.3.

4.1 Discussion on the limiting results

It is instructive to notice that the results of this paper recover the results of [49] (see also [43, 48]) if we restrict attention to the one-layer case. Indeed, let us set $N_2 = 1$, $N_1 = N$, $C^i = 1$ and $H^{2,i} = Z^{2,i}$ in (1.1)-(1.2), and we get the single-layer neural network

$$g_\theta^N(x) = \frac{1}{N} \sum_{j=1}^N W^{2,j} \sigma(W^{1,j} \cdot x),$$

with the corresponding empirical measure of the parameters becoming

$$\gamma_t^N = \tilde{\gamma}_{[Nt]}^N, \text{ where } \tilde{\gamma}_k^N = \frac{1}{N} \sum_{j=1}^N \delta_{W_k^{1,j}, W_k^{2,j}}. \quad (4.1)$$

In that case notice that we can simply write

$$g_{\theta_{[Nt]}^N}^N(x) = \langle w^2 \sigma(w^1 \cdot x), \gamma_t^N \rangle.$$

Then, it is relatively straightforward to notice that the result of Lemma 2.2 boils down to the one layer convergence results of [49], see also [43, 48]. Namely, if we write γ_t for the limit in probability of γ_t^N we get that

$$\lim_{N \rightarrow \infty} g_{\theta_{[Nt]}^N}^N(x) = \langle w^2 \sigma(w^1 \cdot x), \gamma_t \rangle. \quad (4.2)$$

It is useful to compare the limits of the neural network output in the one layer and two layer cases, (4.2) and (2.7) respectively.

Some general remarks of interest follow.

- It is clear that the two layer case (and more generally the multi-layer case) is more complex than the single-layer case, which provides some intuition for the increased complexity of deep neural networks when compared to shallow neural networks. For instance, in the single layer case the neural network output is given explicitly by (4.2). On the other hand, in the multi-layer case, the neural network’s output is the solution to an integral equation as given by (2.6)-(2.7).
- In contrast to the single-layer case, in the multi-layer case the asymptotic weight distribution is layer dependent and weights do not necessarily become independent in the large hidden units limit. In fact, as our result demonstrates, see (2.5)-(2.6), the neural network output in the large hidden units limit and large SGD-iterates does not depend on the individual weights, but on paths of weights that connect the input layer to the output layer. One could regard these paths, i.e. the solution to the random ODE’s in (2.6), as typical representative paths of the weights in the limit connecting the different layers of the “limiting” neural network.
- In the multi-layer case, units at a given layer receive an aggregate signal from units of the previous and/or of the next layer.
- The law of large numbers for a single-layer network indicates that the network will converge in probability to a deterministic limit. That is, after a certain point, adding more hidden units will not increase the accuracy. Our main result Theorem 2.3 suggests that the same conclusion is true for multi-layer neural networks as well.
- Theorem 3.4 combined with Theorem 2.3 characterizes the limiting behavior of the objective function $L^{N_1, N_2}(\theta)$ from (1.3). Section 3 shows that, under the proper assumptions, the limit objective function decreases in the gradient direction of the paths governing the limiting behavior of the weights. The limit ODEs seek to minimize the limit objective function and the global minimum is expected to be recovered.

One practical consequence of our analysis is that the parametrization of the learning rates, see (2.1), indicates that one should use larger learning rates for the weights that connect the different layers (e.g., $W^{2,i,j}$) as compared to the weights for the input or output layers (e.g., C^i and $W^{1,j}$). Notice that this is also the case for the three layer case outlined in Subsection 4.2 below.

As will be explained in Section 4.2, the law of large numbers can be extended to multi-layer neural networks with an arbitrary number of layers. The law of large numbers will only hold under a certain choice of the learning rates. The learning rates need to be scaled with the number of hidden units in each layer. For a multi-layer network with L layers, the learning rates are

$$\alpha_C = \frac{N_L}{N_1}, \quad \alpha_{W,1} = 1, \quad \alpha_{W,2} = N_2, \quad \alpha_{W,L} = \frac{N_{L-1}N_L}{N_1}, \quad \alpha_{W,\ell} = \frac{N_{\ell-1}N_\ell}{N_1}, \quad (4.3)$$

where N_ℓ is the number of hidden units in the ℓ -th layer.

If the learning rates are constant in the number of hidden units N_1, \dots, N_L , it turns out that the network will not train as $N_1, \dots, N_L \rightarrow \infty$ (i.e., in the limit, the network parameters will remain at their initial conditions).

The necessity of scaling the learning rates in the asymptotic regime of large numbers of hidden units (i.e., wide layers) is one of the interesting products of the mean-field limit analysis. A numerical example is presented in Figure 1 below. A deep neural network is trained to classify images for the CIFAR10 dataset [33]. The CIFAR10 dataset contains 60,000 color images in 10 classes (airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck). The dataset is divided into 50,000 training images and 10,000 test images. Each image has $32 \times 32 \times 3$ pixels. The goal is to train a neural network to correctly classify each image based solely on the image pixels as an input. The neural network we use has the mean-field normalizations $\frac{1}{N_\ell}$ in each layer ℓ . There are 8 convolution layers which each have 64 channels. This is followed by two fully-connected layers which each have 500 units. We first train the neural network using the scaled learning rates. Then, we also train the neural network with the standard stochastic gradient descent algorithm (no scaling of the learning rates). Using the scaled learning rates, we achieve a high test accuracy. However, without the scalings, the neural network does not train (i.e., it remains at a very low accuracy).

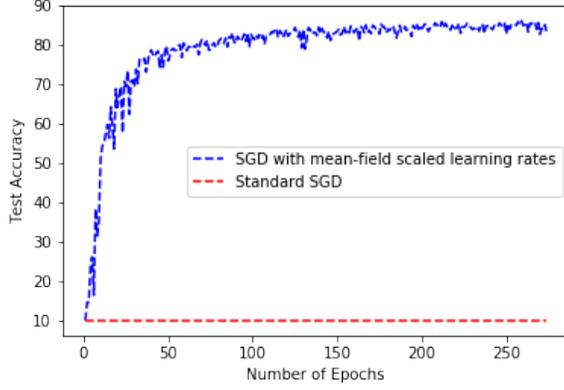


Figure 1: Performance of deep neural network on CIFAR10 dataset with and without scaled learning rates.

4.2 Extension to multi-layer neural networks with more layers

The procedure developed in this paper naturally extends to multi-layer neural networks with more layers than two layers. For brevity, let us present the result in the case of three layers. The situation for more layers is the same, albeit with more complicated algebra. A multi-layer neural network with three layers takes the form

$$g_{\theta}^{N_1, N_2, N_3}(x) = \frac{1}{N_3} \sum_{i=1}^{N_3} C^i \sigma \left(\frac{1}{N_2} \sum_{j=1}^{N_2} W^{3,i,j} \sigma \left(\frac{1}{N_1} \sum_{\nu=1}^{N_1} W^{2,j,\nu} \sigma (W^{1,\nu} \cdot x) \right) \right), \quad (4.4)$$

which can be also written as

$$\begin{aligned} H^{1,\nu}(x) &= \sigma(W^{1,\nu} \cdot x), \quad \nu = 1, \dots, N_1, \\ Z^{2,j}(x) &= \frac{1}{N_1} \sum_{\nu=1}^{N_1} W^{2,j,\nu} H^{1,\nu}(x), \quad j = 1, \dots, N_2 \\ H^{2,j}(x) &= \sigma(Z^{2,j}(x)), \\ Z^{3,i}(x) &= \frac{1}{N_2} \sum_{j=1}^{N_2} W^{3,i,j} H^{2,j}(x), \quad i = 1, \dots, N_3 \\ H^{3,i}(x) &= \sigma(Z^{3,i}(x)), \\ g_{\theta}^{N_1, N_2, N_3}(x) &= \frac{1}{N_3} \sum_{i=1}^{N_3} C^i H^{3,i}(x). \end{aligned} \quad (4.5)$$

where $C^i, W^{2,j,\nu}, W^{3,i,j} \in \mathbb{R}$ and $x, W^{1,\nu} \in \mathbb{R}^d$. The neural network model has parameters

$$\theta = (C^1, \dots, C^{N_3}, W^{2,1,1}, \dots, W^{2,N_2,N_1}, W^{3,1,1}, \dots, W^{3,N_3,N_2}, W^{1,1}, \dots, W^{1,N_1}),$$

which must be estimated from data. The number of hidden units in the first layer is N_1 , the number of hidden units in the second layer is N_2 , and the number of hidden units in the third layer is N_3 . Naturally, the loss function now becomes

$$L^{N_1, N_2, N_3}(\theta) = \frac{1}{2} \mathbb{E}_{Y, X} \left[(Y - g_{\theta}^{N_1, N_2, N_3}(X))^2 \right],$$

where the data $(X, Y) \sim \pi(dx, dy)$.

The stochastic gradient descent (SGD) algorithm for estimating the parameters θ is, for $k \in \mathbb{N}$, $\nu = 1, \dots, N_1$, $i = 1, \dots, N_3$ and $j = 1, \dots, N_2$ is

$$\begin{aligned}
C_{k+1}^i &= C_k^i + \frac{\alpha_C^{N_1, N_2, N_3}}{N_3} (y_k - g_{\theta_k}^{N_1, N_2, N_3}(x_k)) H_k^{3,i}(x_k), \\
W_{k+1}^{1,\nu} &= W_k^{1,\nu} + \frac{\alpha_{W,1}^{N_1, N_2, N_3}}{N_1} (y_k - g_{\theta_k}^{N_1, N_2, N_3}(x_k)) \left(\frac{1}{N_3} \sum_{i=1}^{N_3} C_k^i \sigma'(Z_k^{3,i}(x_k)) \left(\frac{1}{N_2} \sum_{j=1}^{N_2} W_k^{3,i,j} \sigma'(Z_k^{2,j}(x_k)) W_k^{2,j,\nu} \right) \right) \times \\
&\quad \times \sigma'(W_k^{1,\nu} \cdot x_k) x_k, \\
W_{k+1}^{3,i,j} &= W_k^{3,i,j} + \frac{\alpha_{W,3}^{N_1, N_2, N_3}}{N_2 N_3} (y_k - g_{\theta_k}^{N_1, N_2, N_3}(x_k)) C_k^i \sigma'(Z_k^{3,i}(x_k)) H_k^{2,j}(x_k), \\
W_{k+1}^{2,j,\nu} &= W_k^{2,j,\nu} + \frac{\alpha_{W,2}^{N_1, N_2, N_3}}{N_1 N_2} (y_k - g_{\theta_k}^{N_1, N_2, N_3}(x_k)) \frac{1}{N_3} \sum_{i=1}^{N_3} C_k^i \sigma'(Z_k^{3,i}(x_k)) W_k^{3,i,j} \sigma'(Z_k^{2,j}(x_k)) H_k^{1,\nu}(x_k),
\end{aligned} \tag{4.6}$$

where

$$\begin{aligned}
H_k^{1,\nu}(x_k) &= \sigma(W_k^{1,\nu} \cdot x_k), \quad \nu = 1, \dots, N_1, \\
Z_k^{2,j}(x_k) &= \frac{1}{N_1} \sum_{\nu=1}^{N_1} W_k^{2,j,\nu} H_k^{1,\nu}(x_k), \\
H_k^{2,j}(x_k) &= \sigma\left(Z_k^{2,j}(x_k)\right), \\
Z_k^{3,i}(x_k) &= \frac{1}{N_2} \sum_{j=1}^{N_2} W_k^{3,i,j} H_k^{2,j}(x_k), \\
H_k^{3,i}(x_k) &= \sigma\left(Z_k^{3,i}(x_k)\right), \\
g_{\theta_k}^{N_1, N_2, N_3}(x_k) &= \frac{1}{N_3} \sum_{i=1}^{N_3} C_k^i H_k^{3,i}(x_k).
\end{aligned} \tag{4.7}$$

where $\alpha_C^{N_1, N_2, N_3}$, $\alpha_{W,1}^{N_1, N_2, N_3}$, $\alpha_{W,2}^{N_1, N_2, N_3}$ and $\alpha_{W,3}^{N_1, N_2, N_3}$ are the learning rates. The parameters at step k are

$$\theta_k^{N_1, N_2, N_3} = (C_k^1, \dots, C_k^{N_3}, W_k^{2,1,1}, \dots, W_k^{2, N_2, N_1}, W_k^{3,1,1}, \dots, W_k^{3, N_3, N_2}, W_k^{1,1}, \dots, W_k^{1, N_1}).$$

(x_k, y_k) are samples of the random variables (X, Y) . We assume a condition analogous to Assumption 2.1.

Let us now choose the learning rates to be

$$\alpha_C^{N_1, N_2, N_3} = \frac{N_3}{N_1}, \quad \alpha_{W,1}^{N_1, N_2, N_3} = 1, \quad \alpha_{W,3}^{N_1, N_2, N_3} = \frac{N_2 N_3}{N_1}, \quad \alpha_{W,2}^{N_1, N_2, N_3} = N_2$$

Similar to before, define the empirical measure

$$\tilde{\gamma}_k^{N_1, N_2, N_3} = \frac{1}{N_1} \sum_{\nu=1}^{N_1} \delta_{W_k^{1,\nu}, W_k^{2,1,\nu}, \dots, W_k^{2, N_2, \nu}, W_k^{3,1,1}, \dots, W_k^{3, N_3, N_2}, C_k^1, \dots, C_k^{N_3}}.$$

The time-scaled empirical measure is

$$\gamma_t^{N_1, N_2, N_3} := \tilde{\gamma}_{\lfloor N_1 t \rfloor}^{N_1, N_2, N_3},$$

and the corresponding time-scaled neural network output is $g_t^{N_1, N_2, N_3}(x) = g_{\theta_{\lfloor N_1 t \rfloor}}^{N_1, N_2, N_3}(x)$.

Following the same procedure as for the two layer case one expects to get the following limit that describes $g_t(x)$ (the iterated limit of $g_t^{N_1, N_2, N_3}(x)$ as first $N_1 \rightarrow \infty$, then $N_2 \rightarrow \infty$, and then $N_3 \rightarrow \infty$):

$$\begin{aligned}
d\tilde{C}_t^c &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t(x)) \tilde{H}_t^{3,c}(x) \pi(dx, dy) dt, & \tilde{C}_0^c &= c, \\
d\tilde{W}_t^{1,w} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t(x)) V_t^w(x) \sigma'(W_t^{1,w} \cdot x) x \pi(dx, dy) dt, & \tilde{W}_0^{1,w} &= w, \\
d\tilde{W}_t^{3,c,v} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t(x)) \tilde{C}_t^c \sigma'(\tilde{Z}_t^{3,c}(x)) \tilde{H}_t^{2,v}(x) \pi(dx, dy) dt, & \tilde{W}_0^{3,c,v} &= v, \\
d\tilde{W}_t^{2,w,u,v} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t(x)) L_t^v(x) \sigma'(\tilde{Z}_t^{2,v}(x)) \tilde{H}_t^{1,w}(x) \pi(dx, dy) dt, & \tilde{W}_0^{2,w,u,v} &= u, \\
\tilde{H}_t^{1,w}(x) &= \sigma(\tilde{W}_t^{1,w} \cdot x), \\
\tilde{Z}_t^{2,v}(x) &= \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}_t^{2,w,u,v} \tilde{H}_t^{1,w}(x) \mu_{W^2}(du) \mu_{W^1}(dw), \\
\tilde{H}_t^{2,v}(x) &= \sigma(\tilde{Z}_t^{2,v}(x)), \\
\tilde{Z}_t^{3,c}(x) &= \int_{\mathcal{W}^3} \tilde{W}_t^{3,c,v} \tilde{H}_t^{2,v}(x) \mu_{W^3}(dv), \\
\tilde{H}_t^{3,c}(x) &= \sigma(\tilde{Z}_t^{3,c}(x)), \\
V_t^w(x) &= \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^{3,c}(x)) \left(\int_{\mathcal{W}^3} \tilde{W}_t^{3,c,v} \sigma'(\tilde{Z}_t^{2,v}(x)) \left(\int_{\mathcal{W}^2} \tilde{W}_t^{2,w,u,v} \mu_{W^2}(du) \right) \mu_{W^3}(dv) \right) \mu_{\mathcal{C}}(dc) \\
L_t^v(x) &= \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^{3,c}(x)) \tilde{W}_t^{3,c,v} \mu_{\mathcal{C}}(dc) \\
g_t(x) &= \int_{\mathcal{C}} \tilde{C}_t^c \tilde{H}_t^{3,c}(x) \mu_{\mathcal{C}}(dc), \tag{4.8}
\end{aligned}$$

In other words, one expects to be able to write that the neural network's output is

$$g_t(x) = \int_{\mathcal{C}} \tilde{C}_t^c \left(\sigma \left(\int_{\mathcal{W}^3} \tilde{W}_t^{3,c,v} \sigma \left(\int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}_t^{2,w,u,v} \tilde{\sigma}(\tilde{W}_t^{1,w} \cdot x) \mu_{W^2}(du) \mu_{W^1}(dw) \right) \mu_{W^3}(dv) \right) \right) \mu_{\mathcal{C}}(dc).$$

A computation analogous to the one in Section 3 shows that, as expected, the limit ODEs in (4.8) seek to minimize the corresponding limit objective function as $t \rightarrow \infty$. We leave the rigorous proof for the form of $g_t(x)$ in the three layer neural network case to the interested reader.

4.3 Challenges in the analysis of multi-layer neural networks

Challenges arise in the study of the asymptotics of multi-layer neural networks. [49], [43], and [48] leveraged traditional approaches from the mean-field and interacting particle system literature to characterize the asymptotics of single-layer neural networks. However, it turns out that the traditional mean-field approach cannot be used for multi-layer neural networks.

A standard approach for analyzing (1.4) as $N_1, N_2 \rightarrow \infty$ would be to construct an empirical measure $\rho_k^{N_1, N_2}$ of the parameters θ_k at training step k , as for example $\tilde{\gamma}_k^{N_1, N_2}$. Then, we could study the behavior of $\rho_k^{N_1, N_2}$ as $N_1, N_2 \rightarrow \infty$. This empirical measure $\rho_k^{N_1, N_2}$ needs to be designed such that the dynamics of $\rho_k^{N_1, N_2}$ can be written in terms of $\rho_k^{N_1, N_2}$ itself and the data (x_k, y_k) (plus martingale and remainder terms). That is, the dynamics of $\rho_k^{N_1, N_2}$ are *closed*.

This is straightforward for single-layer neural networks (see [49]), but it is challenging to do for multi-layer neural networks. In the case of single-layer networks, the empirical measure is simply given by (4.1) and its analysis has been successfully carried on in [49]. One is tempted to do the same thing for the multi-layer case, i.e. study the limit of the empirical measure defined in (2.2). The problem that one faces with this formulation is that, in contrast to the single-layer case, the dimension of the space on which the empirical measure takes values increases with N_2 . Therefore, the problem cannot be studied using such an empirical measure.

An alternative approach, which is also natural in this case, is to try and create “nested measures”, sometimes called multi-level measure valued processes in mathematical biology, see [12, 14, 15, 16, 17, 18, 19] and the review paper [13].

Let us explain how to construct a nested empirical measure for (1.4) and why it will not work. In order to simplify notation, let $N_1 = N_2 = N$ and $\rho_k^{N_1, N_2} = \rho_k^N$ in the following example. Considering $N_1 \neq N_2$ and taking subsequent limits does not alter the conclusions below.

- Let’s first examine the parameter $W^{1,j}$ in the first layer. The j -th unit in the first layer (i.e., $H^{1,j} = \sigma(W^{1,j} \cdot x)$) is connected to all of the hidden units in the second layer (i.e., $H^{2,i}$) via the weights $W^{2,i,j}$. Therefore, there is a measure associated with each $W^{1,j}$ which must track $\{W^{2,i,j}, Z^i, C^i\}_{i=1}^N$. $W^{2,i,j}$ and Z^i are required for calculating the SGD update for $W^{1,j}$ (see 1.4). This measure is $\nu_k^{N,j} = \frac{1}{N} \sum_{i=1}^N \delta_{W_k^{2,i,j}, Z_k^i, C_k^i}$.
- Let’s next examine the parameter C^i in the second layer. The i -th unit in the second layer is connected to all of the hidden units in the first layer via the weights $W^{2,i,j}$. Therefore, for each C^i , we must track $\{W^{2,i,j}, W^{1,j}\}_{j=1}^N$ in order to calculate the SGD update for C^i (see 1.4). Furthermore, updating C^i requires tracking $W^{1,j}$, and updating $W^{1,j}$ requires tracking ν^j . Therefore, updates to C^i require the empirical measure $\mu_k^{N,i} = \frac{1}{N} \sum_{j=1}^N \delta_{W_k^{2,i,j}, W_k^{1,j}, \nu_k^{N,j}} \in \mathcal{M}(\mathbb{R} \times \mathbb{R}^d \times \mathcal{M}(\mathbb{R}^3))$.
- Finally, the entire network at iteration k is specified by the empirical measure

$$\rho_k^N = \frac{1}{N} \sum_{i=1}^N \delta_{C_k^i, \mu_k^i} \in \mathcal{M}(\mathbb{R} \times \mathcal{M}(\mathbb{R} \times \mathbb{R}^d \times \mathcal{M}(\mathbb{R}^3))), \quad (4.9)$$

where $\mathcal{M}(E)$ is the space of measures on the metric space E . Notice that the process (4.9) involves *nested measures* (sometimes called “multi-level processes”). The process ρ_k takes values in a *space of nested measures* $\mathcal{M}(\mathcal{M}(\mathcal{M}(\dots)))$.

Careful inspection of ρ_k^N identifies a crucial problem: its dynamics are not closed. The evolution of $\nu_k^{N,j}$ (the innermost measure in the nested measures) cannot be written in terms of ρ_k^N . In particular, updating $\nu_k^{N,j}$ requires also updating $(W^{2,i,j}, Z_k^i, C_k^i)_{i=1}^N$. This would in fact require knowledge of $(W^{2,i,j}, Z_k^i, C_k^i, \mu_k^{N,i})_{i=1}^N$, i.e. we would have to re-define $\nu_k^{N,j}$ as $\frac{1}{N} \sum_{i=1}^N \delta_{W_k^{2,i,j}, Z_k^i, C_k^i, \eta_k^{N,i}}$ where $\eta_k^{N,i} = \frac{1}{N} \sum_{m=1}^N \delta_{W_k^{2,i,m}, W_k^m}$. This leads to re-defining ρ_k^N as taking values in $\mathcal{M}(\mathbb{R} \times \mathcal{M}(\mathbb{R} \times \mathbb{R}^d \times \mathcal{M}(\mathbb{R}^3 \times \mathcal{M}(\mathbb{R} \times \mathbb{R}^d))))$, i.e. the space of 4-nested measures (before it was 3-nested measures). However, the closure problem remains, since the evolution of $\eta_k^{N,i}$ cannot be written in terms ρ_k^N . In fact, there does not seem to exist any finite number of nested measures for which the empirical measure ρ_k^N ’s dynamics are closed and despite our best efforts we have not managed to find one. This closure problem then leads to non-trivial issues associated with establishing a well define limit.

For completeness, we also remark that a third alternative is to define the candidate empirical measure as an appropriately normalized double sum over both indices corresponding to the two hidden layers and then consider the limit of this measure as $N_1, N_2 \rightarrow \infty$. However, this approach also suffers a closure problem, similar to the situation described above.

The discussion in this section highlights some of the challenges that must be addressed in the analysis of multi-layer neural networks. Such problems are not present in the analysis of single-layer neural networks. Therefore, a different approach is required for the asymptotic analysis of multi-layer neural networks and this paper is a first step in this direction.

The problems that we describe above led us to the approach used in this paper. In particular, the approach in Section 2 first studies the limit of the empirical measure as the number of hidden units in the first layer grows to infinity. This is similar to [49]. The limit is a solution to an evolution equation and it is the law of a system of random ODEs. We then make the crucial observation that one can characterize the resulting system in terms of the initialization for the stochastic gradient descent iterates. This means that we can reformulate the limiting system of the first layer into an equivalent system of random ODEs and then consider the limit of the second layer. This allows us to obtain the limit of the output of the neural network as the number of hidden units of all layers grow to infinity by studying the limit of the random ODE in Theorem 2.3.

5 Proof of Theorem 2.3-Characterization of the limit

In preparation for the proof of Theorem 2.3, we first re-express the result from Lemma 2.2.

Corollary 5.1. Consider the particle system:

$$\begin{aligned}
dC_t^i &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) H_t^{2,i}(x) \pi(dx, dy) dt, \quad C_0^i = C_0^i, \quad i = 1, \dots, N_2, \\
dW_t^1 &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) \left(\frac{1}{N_2} \sum_{i=1}^{N_2} C_t^i \sigma'(Z_t^i(x)) W_t^{2,i} \right) \sigma'(W_t^1 \cdot x) \pi(dx, dy) dt, \quad W_0^1 \sim \mu_{W^1}(dw), \\
dW_t^{2,i} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) C_t^i \sigma'(Z_t^i(x)) H_t^1(x) \pi(dx, dy) dt, \quad W_0^{2,i} \sim \mu_{W^2}(du), \quad i = 1, \dots, N_2, \\
H_t^1(x) &= \sigma(W_t^1 \cdot x), \\
Z_t^i(x) &= \mathbb{E} \left[W_t^{2,i} H_t^1(x) \middle| C_0^1, \dots, C_0^{N_2} \right], \\
H_t^{2,i}(x) &= \sigma(Z_t^i(x)), \\
g_t^{N_2}(x) &= \frac{1}{N_2} \sum_{i=1}^{N_2} C_t^i H_t^{2,i}(x). \tag{5.1}
\end{aligned}$$

Let $\nu_{t, (c_1, \dots, c_{N_2})}$ be the conditional Law of $(W_t^1, W_t^{2,1}, \dots, W_t^{2,N_2}, C_t^1, \dots, C_t^{N_2})_{0 \leq t \leq T}$ given $(C_0^1, \dots, C_0^{N_2}) = (c_1, \dots, c_{N_2})$. Then, $\nu_{t, (C_0^1, \dots, C_0^{N_2})}$ is the unique solution to the evolution equation (2.3)-(2.4).

Proof. See Appendix A. □

Due to exchangeability properties and without loss of generality, we can transform (5.1) into the equivalent particle system:

$$\begin{aligned}
dC_t^{C_0^i} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) H_t^{2,C_0^i}(x) \pi(dx, dy) dt, \quad C_0^{C_0^i} = C_0^i \quad i = 1, \dots, N_2, \\
dW_t^{1,W_0} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) V_t^{N_2, W_0} \sigma'(W_t^{1,W_0} \cdot x) \pi(dx, dy) dt, \quad W_0^{1,W_0} = W_0 \sim \mu_{W^1}(dw) \\
dW_t^{2,C_0^i, W_0, W_0^{2,i}} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) C_t^{C_0^i} \sigma'(Z_t^{C_0^i}(x)) H_t^{1,W_0}(x) \pi(dx, dy) dt, \quad W_0^{2,C_0^i, W_0, W_0^{2,i}} = W_0^{2,i} \sim \mu_{W^2}(du), \\
H_t^{1,W_0}(x) &= \sigma(W_t^{1,W_0} \cdot x), \\
Z_t^{C_0^i}(x) &= \mathbb{E} \left[W_t^{2,C_0^i, W_0, W_0^{2,i}} H_t^{1,W_0}(x) \middle| C_0^1, \dots, C_0^{N_2} \right], \quad i = 1, \dots, N_2, \\
H_t^{2,C_0^i}(x) &= \sigma(Z_t^{C_0^i}(x)), \quad i = 1, \dots, N_2, \\
g_t^{N_2}(x) &= \frac{1}{N_2} \sum_{i=1}^{N_2} C_t^{C_0^i} H_t^{2,C_0^i}(x), \\
V_t^{N_2, W_0}(x) &= \frac{1}{N_2} \sum_{i=1}^{N_2} C_t^{C_0^i} \sigma'(Z_t^{C_0^i}(x)) W_t^{2,C_0^i, W_0, W_0^{2,i}}. \tag{5.2}
\end{aligned}$$

Since for all $i = 1, \dots, N_2$, C_0^i have probability density function as described in Assumption 2.1, we have that

$$\mathbb{P}[\{C_0^i \neq C_0^j\}_{i \neq j, (i,j)=1,2,\dots,N_2}] = 1.$$

This allows us to substitute the variable names

$$(\hat{C}_t^i, \hat{W}_t^1, \hat{W}_t^{2,i}, \hat{Z}_t^i, \hat{H}_t^1, \hat{H}_t^{2,i}) \text{ for } (C_t^{C_0^i}, W_t^{1,W_0}, W_t^{2,C_0^i, W_0, W_0^{2,i}}, Z_t^{C_0^i}, H_t^{1,W_0}, H_t^{2,C_0^i}),$$

for $i = 1, \dots, N_2$.

This produces the system:

$$\begin{aligned}
d\hat{C}_t^i &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) \hat{H}_t^{2,i}(x) \pi(dx, dy) dt, & \hat{C}_0^i &= C_\circ^i, & i &= 1, \dots, N_2, \\
d\hat{W}_t^1 &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) V_t^{N_2, W_0} \sigma'(\hat{W}_t^1 \cdot x) x \pi(dx, dy) dt, & \hat{W}_0^1 &= W_0 \sim \mu_{W^1}(dw), \\
d\hat{W}_t^{2,i} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) \hat{C}_t^i \sigma'(\hat{Z}_t^i(x)) \hat{H}_t^1(x) \pi(dx, dy) dt, & \hat{W}_0^{2,i} &= W_0^{2,i} \sim \mu_{W^2}(du), \\
\hat{H}_t^1(x) &= \sigma(\hat{W}_t^1 \cdot x), \\
\hat{Z}_t^i(x) &= \mathbb{E} \left[\hat{W}_t^{2,i} \hat{H}_t^1(x) \middle| C_\circ^1, \dots, C_\circ^{N_2} \right], & i &= 1, \dots, N_2, \\
\hat{H}_t^{2,i}(x) &= \sigma(\hat{Z}_t^i(x)), & i &= 1, \dots, N_2, \\
g_t^{N_2}(x) &= \frac{1}{N_2} \sum_{i=1}^{N_2} \hat{C}_t^i \hat{H}_t^{2,i}(x), \\
V_t^{N_2, W_0}(x) &= \frac{1}{N_2} \sum_{i=1}^{N_2} \hat{C}_t^i \sigma'(\hat{Z}_t^i(x)) \hat{W}_t^{2,i}.
\end{aligned} \tag{5.3}$$

The system (5.3) is exactly the same system as (5.1). Notice also that \hat{C}_t^i in (5.3) depends also on $\{C_\circ^i\}_{i=1}^{N_2}$ in a symmetric way via $g_t^{N_2}(x)$. Similarly \hat{W}_t^1 and $\hat{W}_t^{2,i}$ depend also on $\{C_\circ^i\}_{i=1}^{N_2}$ and on $\{W_0^{2,i}\}_{i=1}^{N_2}$ symmetrically via $g_t^{N_2}(x)$ and $V_t^{N_2, W_0}(x)$. This is also the situation for (5.1). Then independence and identical distribution of the initial conditions together with the aforementioned exchangeability property imply that (5.1) and (5.2) are equivalent.

5.1 Limiting System

The goal is to prove that for any $t \in [0, 1]$ and $x \in \mathcal{X}$,

$$\lim_{N_2 \rightarrow \infty} g_t^{N_2}(x) = g_t(x),$$

in L^1 , where $g_t^{N_2}(x)$ is from (5.2) and the limit $g_t(x)$ is given by

$$g_t(x) = \int_{\mathcal{C}} \tilde{C}_t^c \tilde{H}_t^{2,c}(x) \mu_c(dc),$$

where,

$$\begin{aligned}
d\tilde{C}_t^c &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t(x)) \tilde{H}_t^{2,c}(x) \pi(dx, dy) dt, & \tilde{C}_0^c &= c, \\
d\tilde{W}_t^{1,w} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t(x)) V_t^w \sigma'(\tilde{W}_t^{1,w} \cdot x) x \pi(dx, dy) dt, & \tilde{W}_0^{1,w} &= w, \\
d\tilde{W}_t^{2,c,w,u} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t(x)) \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \tilde{H}_t^{1,w}(x) \pi(dx, dy) dt, & \tilde{W}_0^{2,c,w,u} &= u, \\
\tilde{H}_t^{1,w}(x) &= \sigma(\tilde{W}_t^{1,w} \cdot x), \\
\tilde{Z}_t^c(x) &= \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}_t^{2,c,w,u} \tilde{H}_t^{1,w}(x) \mu_{W^2}(du) \mu_{W^1}(dw), \\
\tilde{H}_t^{2,c}(x) &= \sigma(\tilde{Z}_t^c(x)), \\
V_t^w(x) &= \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \left(\int_{\mathcal{W}^2} \tilde{W}_t^{2,c,w,u} \mu_{W^2}(du) \right) \mu_c(dc).
\end{aligned} \tag{5.4}$$

Before presenting the proof of this result, let us define a quantity that will be of central interest in the

sequel. In particular, for $(c, w, u) \in \{(C_0^i, W_0^1, W_0^{2,i}), i = 1, \dots, N_2\}$ and $x \in \mathcal{X}$, let's define the error function

$$E_t^{N_2}(c, w, u, x) := (C_t^c - \tilde{C}_t^c)^2 + \left\| W_t^{1,w} - \tilde{W}_t^{1,w} \right\|^2 + (W_t^{2,c,w,u} - \tilde{W}_t^{2,c,w,u})^2 \\ + (H_t^{2,c}(x) - \tilde{H}_t^{2,c}(x))^2 + (Z_t^c(x) - \tilde{Z}_t^c(x))^2.$$

Note that we certainly have,

$$|C_0^c - \tilde{C}_0^c|^2 + \left\| W_0^{1,w} - \tilde{W}_0^{1,w} \right\|^2 + |W_0^{2,c,w,u} - \tilde{W}_0^{2,c,w,u}|^2 + |Z_0^c(x) - \tilde{Z}_0^c(x)|^2 = 0.$$

We will show below that $\mathbb{E} \left[E_t^{N_2}(C_0^i, W_0, W_0^{2,i}, x) \right]$ appropriately converges to zero as $N_2 \rightarrow \infty$ which will then imply that $g_t^{N_2}(x)$ converges to $g_t(x)$ as indicated.

5.2 A Priori Bounds

Let us first establish uniform bounds on the processes $C_t^c, W_t^{1,w}, W_t^{2,c,w,u}$, and $g_t^{N_2}(x)$ for the system (5.2). For any $t \in [0, T]$ and any $N_2 \in \mathbb{N}$,

$$\frac{1}{N_2} \sum_{i=1}^{N_2} |C_t^{C_0^i}| \leq \frac{1}{N_2} \sum_{i=1}^{N_2} |C_0^{C_0^i}| + \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} |y - g_s^{N_2}(x)| \frac{1}{N_2} \sum_{i=1}^{N_2} |H_s^{2,C_0^i}(x)| \pi(dx, dy) ds.$$

$|H_s^{2,C_0^i}(x)| < K$ since $\sigma(\cdot) \in C_b^2$. Then, using the fact that $\mathcal{X} \times \mathcal{Y}$ is compact,

$$\frac{1}{N_2} \sum_{i=1}^{N_2} |C_t^{C_0^i}| \leq \frac{1}{N_2} \sum_{i=1}^{N_2} |C_0^{C_0^i}| + K_1 t + K_2 \int_0^t \frac{1}{N_2} \sum_{i=1}^{N_2} |C_s^{C_0^i}| ds.$$

By Gronwall's inequality, we have that for any $N_2 \in \mathbb{N}$ and for any $t \in [0, T]$,

$$\frac{1}{N_2} \sum_{i=1}^{N_2} |C_t^{C_0^i}| \leq K.$$

Using the same approach, we can establish uniform bounds on the other processes $W_t^{1,w}$ and $W_t^{2,c,w,u}$. Therefore, for any $N_2 \in \mathbb{N}$, $t \in [0, 1]$, $x \in \mathcal{X}$, and $(c, w, u) \in \{(C_0^i, W_0^1, W_0^{2,i}), i = 1, \dots, N_2\}$,

$$\max\{|g_t^{N_2}(x)|, |C_t^c|, |W_t^{2,c,w,u}|, |V_t^{N_2,w}(x)|, |W_t^{1,w}|\} \leq K, \quad (5.5)$$

5.3 Bound for $\mathbb{E}|V_t^{N_2, W_0}(x) - V_t^{W_0}(x)|^2$.

We have:

$$|V_t^{N_2, W_0}(x) - V_t^{W_0}(x)| \leq \left| \frac{1}{N_2} \sum_{i=1}^{N_2} C_t^{C_0^i} \sigma'(Z_t^{C_0^i}(x)) W_t^{2,C_0^i, W_0, W_0^{2,i}} - \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \left(\int_{\mathcal{W}^2} \tilde{W}_t^{2,c, W_0, u} \mu_{\mathcal{W}^2}(du) \right) \mu_{\mathcal{C}}(dc) \right| \\ = \left| \frac{1}{N_2} \sum_{i=1}^{N_2} C_t^{C_0^i} \sigma'(Z_t^{C_0^i}(x)) W_t^{2,C_0^i, W_0, W_0^{2,i}} - \frac{1}{N_2} \sum_{i=1}^{N_2} \tilde{C}_t^{C_0^i} \sigma'(\tilde{Z}_t^{C_0^i}(x)) \tilde{W}_t^{2,C_0^i, W_0, W_0^{2,i}} \right| \\ + \left| \frac{1}{N_2} \sum_{i=1}^{N_2} \tilde{C}_t^{C_0^i} \sigma'(\tilde{Z}_t^{C_0^i}(x)) \tilde{W}_t^{2,C_0^i, W_0, W_0^{2,i}} - \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \left(\int_{\mathcal{W}^2} \tilde{W}_t^{2,c, W_0, u} \mu_{\mathcal{W}^2}(du) \right) \mu_{\mathcal{C}}(dc) \right| \\ := |\Gamma_t^{V,1, W_0}| + |\Gamma_t^{V,2, W_0}|. \quad (5.6)$$

The first term in (5.6) can be studied using a decomposition:

$$\begin{aligned}
\left| \Gamma_t^{V,1,W_0} \right| &= \left| \frac{1}{N_2} \sum_{i=1}^{N_2} C_t^{C_\circ^i} \sigma'(Z_t^{C_\circ^i}(x)) W_t^{2,C_\circ^i,W_0,W_0^{2,i}} - \frac{1}{N_2} \sum_{i=1}^{N_2} \tilde{C}_t^{C_\circ^i} \sigma'(\tilde{Z}_t^{C_\circ^i}(x)) \tilde{W}_t^{2,C_\circ^i,W_0,W_0^{2,i}} \right| \\
&= \left| \frac{1}{N_2} \sum_{i=1}^{N_2} \left[\left(C_t^{C_\circ^i} - \tilde{C}_t^{C_\circ^i} \right) \sigma'(Z_t^{C_\circ^i}(x)) W_t^{2,C_\circ^i,W_0,W_0^{2,i}} \right. \right. \\
&\quad \left. \left. + \tilde{C}_t^{C_\circ^i} \left(\sigma'(Z_t^{C_\circ^i}(x)) W_t^{2,C_\circ^i,W_0,W_0^{2,i}} - \sigma'(\tilde{Z}_t^{C_\circ^i}(x)) \tilde{W}_t^{2,C_\circ^i,W_0,W_0^{2,i}} \right) \right] \right| \\
&= \left| \frac{1}{N_2} \sum_{i=1}^{N_2} \left[\left(C_t^{C_\circ^i} - \tilde{C}_t^{C_\circ^i} \right) \sigma'(Z_t^{C_\circ^i}(x)) W_t^{2,C_\circ^i,W_0,W_0^{2,i}} + \tilde{C}_t^{C_\circ^i} \left(\sigma'(Z_t^{C_\circ^i}(x)) - \sigma'(\tilde{Z}_t^{C_\circ^i}(x)) \right) W_t^{2,C_\circ^i,W_0,W_0^{2,i}} \right. \right. \\
&\quad \left. \left. + \tilde{C}_t^{C_\circ^i} \sigma'(\tilde{Z}_t^{C_\circ^i}(x)) \left(W_t^{2,C_\circ^i,W_0,W_0^{2,i}} - \tilde{W}_t^{2,C_\circ^i,W_0,W_0^{2,i}} \right) \right] \right| \\
&\leq \frac{K}{N_2} \sum_{i=1}^{N_2} \left[\left| C_t^{C_\circ^i} - \tilde{C}_t^{C_\circ^i} \right| + \left| Z_t^{C_\circ^i}(x) - \tilde{Z}_t^{C_\circ^i}(x) \right| + \left| W_t^{2,C_\circ^i,W_0,W_0^{2,i}} - \tilde{W}_t^{2,C_\circ^i,W_0,W_0^{2,i}} \right| \right] \\
&\leq \frac{K}{N_2} \sum_{i=1}^{N_2} \sqrt{E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x)}.
\end{aligned}$$

where the uniform bounds from (5.5) were used. In addition, we also have for some constant $K < \infty$

$$\begin{aligned}
\mathbb{E} \left[\left(\Gamma_t^{V,2,W_0} \right)^2 \right] &= \mathbb{E} \left[\left(\frac{1}{N_2} \sum_{i=1}^{N_2} \tilde{C}_t^{C_\circ^i} \sigma'(\tilde{Z}_t^{C_\circ^i}(x)) \tilde{W}_t^{2,C_\circ^i,W_0,W_0^{2,i}} - \int_{\mathcal{C}} \tilde{C}_t^c \sigma'(\tilde{Z}_t^c(x)) \left(\int_{W_2} \tilde{W}_t^{2,c,W_0,u} \mu_{W_2}(du) \right) \mu_c(dc) \right)^2 \right] \\
&= \frac{1}{N_2^2} \sum_{i=1}^{N_2} \text{Var} \left[\tilde{C}_t^{C_\circ^i} \sigma'(\tilde{Z}_t^{C_\circ^i}(x)) \tilde{W}_t^{2,C_\circ^i,W_0,W_0^{2,i}} \right] \\
&\leq \frac{K}{N_2}.
\end{aligned}$$

where we used the assumed independence of C_\circ^i , W_0 , and $W_0^{2,i}$ as well as the a-priori bound from (5.5).

Hence, we obtain that for some unimportant constant $K < \infty$

$$\mathbb{E} |V_t^{N_2,W_0}(x) - V_t^{W_0}(x)|^2 \leq \frac{K}{N_2} \sum_{i=1}^{N_2} \mathbb{E} \left[E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] + \frac{K}{N_2}. \quad (5.7)$$

5.4 Bound for $|g_t^{N_2}(x) - g_t(x)|$.

We can write

$$\begin{aligned}
|g_t^{N_2}(x) - g_t(x)| &= \left| \frac{1}{N_2} \sum_{i=1}^{N_2} C_t^{C_\circ^i} H_t^{2,C_\circ^i}(x) - \int_{\mathcal{C}} \tilde{C}_t^c \tilde{H}_t^{2,c}(x) \mu_c(dc) \right| \\
&\leq \left| \frac{1}{N_2} \sum_{i=1}^{N_2} C_t^{C_\circ^i} H_t^{2,C_\circ^i}(x) - \frac{1}{N_2} \sum_{i=1}^{N_2} \tilde{C}_t^{C_\circ^i} \tilde{H}_t^{2,C_\circ^i}(x) \right| \\
&\quad + \left| \frac{1}{N_2} \sum_{i=1}^{N_2} \tilde{C}_t^{C_\circ^i} \tilde{H}_t^{2,C_\circ^i}(x) - \int_{\mathcal{C}} \tilde{C}_t^c \tilde{H}_t^{2,c}(x) \mu_c(dc) \right| \\
&:= \Gamma_t^{g,1}(x) + \Gamma_t^{g,2}(x). \quad (5.8)
\end{aligned}$$

Let's analyze the first term in (5.8). Using the uniform bounds from (5.5) we have for some unimportant

constant $K < \infty$

$$\begin{aligned}
\left| \Gamma_t^{g,1}(x) \right| &= \left| \frac{1}{N_2} \sum_{i=1}^{N_2} C_t^{C_\circ^i} H_t^{2,C_\circ^i}(x) - \frac{1}{N_2} \sum_{i=1}^{N_2} \tilde{C}_t^{C_\circ^i} \tilde{H}_t^{2,C_\circ^i}(x) \right| \\
&= \left| \frac{1}{N_2} \sum_{i=1}^{N_2} (C_t^{C_\circ^i} - \tilde{C}_t^{C_\circ^i}) H_t^{2,C_\circ^i}(x) + \frac{1}{N_2} \sum_{i=1}^{N_2} \tilde{C}_t^{C_\circ^i} (H_t^{2,C_\circ^i}(x) - \tilde{H}_t^{2,C_\circ^i}(x)) \right| \\
&\leq \frac{K}{N_2} \sum_{i=1}^{N_2} \left[|C_t^{C_\circ^i} - \tilde{C}_t^{C_\circ^i}| + |H_t^{2,C_\circ^i}(x) - \tilde{H}_t^{2,C_\circ^i}(x)| \right] \\
&\leq \frac{K}{N_2} \sum_{i=1}^{N_2} \sqrt{E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x)}.
\end{aligned}$$

The second term in (5.8) is bounded, as follows,

$$\mathbb{E} \left[\left(\Gamma_t^{g,2}(x) \right)^2 \right] = \mathbb{E} \left[\left(\frac{1}{N_2} \sum_{i=1}^{N_2} (\tilde{C}_t^{C_\circ^i} \tilde{H}_t^{2,C_\circ^i}(x) - \int_c \tilde{C}_t^c \tilde{H}_t^{2,c}(x) \mu_c(dc)) \right)^2 \right] = \frac{1}{N_2^2} \sum_{i=1}^{N_2} \text{Var} [\tilde{C}_t^{C_\circ^i} \tilde{H}_t^{2,C_\circ^i}(x)] \leq \frac{K}{N_2}, \quad (5.9)$$

where the independence of C_\circ^i was used. Hence, putting things together we get

$$\mathbb{E} |g_t^{N_2}(x) - g_t(x)| \leq \mathbb{E} \left[\frac{K}{N_2} \sum_{i=1}^{N_2} \sqrt{E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x)} \right] + \frac{K}{\sqrt{N_2}}.$$

$$\mathbf{5.5} \quad \text{Bound for } \mathbb{E} \left[(C_t^{C_\circ^i} - \tilde{C}_t^{C_\circ^i})^2 + (W_t^{1,W_0} - \tilde{W}_t^{1,W_0})^2 + (W_t^{2,C_\circ^i,W_0,W_0^{2,i}} - \tilde{W}_t^{2,C_\circ^i,W_0,W_0^{2,i}})^2 \right].$$

Let us write for notational convenience $c = C_\circ^i$. We have

$$d(C_t^c - \tilde{C}_t^c)^2 = 2(C_t^c - \tilde{C}_t^c) \int_{\mathcal{X} \times \mathcal{Y}} \left[(y - g_t^{N_2}(x)) (H_t^{2,c}(x) - \tilde{H}_t^{2,c}(x)) + (g_t(x) - g_t^{N_2}(x)) \tilde{H}_t^{2,c}(x) \right] \pi(dx, dy) dt.$$

Using Young's inequality and the uniform bounds from (5.5),

$$\begin{aligned}
(C_t^c - \tilde{C}_t^c)^2 &\leq (C_0^c - \tilde{C}_0^c)^2 + K \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} \pi(dx, dy) \left[(C_s^c - \tilde{C}_s^c)^2 + (H_s^{2,c}(x) - \tilde{H}_s^{2,c}(x))^2 \right. \\
&\quad \left. + |C_s^c - \tilde{C}_s^c| \cdot |g_s(x) - g_s^{N_2}(x)| \right] ds \\
&= K \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} \pi(dx, dy) \left[(C_s^c - \tilde{C}_s^c)^2 + (H_s^{2,c}(x) - \tilde{H}_s^{2,c}(x))^2 + |C_s^c - \tilde{C}_s^c| (|\Gamma_s^{g,1}(x) + \Gamma_s^{g,2}(x)|) \right] ds \\
&= K \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} \pi(dx, dy) \left[(C_s^c - \tilde{C}_s^c)^2 + (H_s^{2,c}(x) - \tilde{H}_s^{2,c}(x))^2 + |C_s^c - \tilde{C}_s^c| (|\Gamma_s^{g,1}(x) + \Gamma_s^{g,2}(x)|)^2 \right] ds. \quad (5.10)
\end{aligned}$$

For the term $|C_s^c - \tilde{C}_s^c| \Gamma_s^{g,1}(x)$, where we recall $c = C_\circ^i$, we have

$$\begin{aligned}
|C_s^{C_\circ^i} - \tilde{C}_s^{C_\circ^i}| \Gamma_s^{g,1}(x) &\leq |C_s^{C_\circ^i} - \tilde{C}_s^{C_\circ^i}| \frac{K}{N_2} \sum_{j=1}^{N_2} \sqrt{E_s^{N_2}(C_\circ^j, W_0, W_0^{2,j}, x)} \\
&= \frac{K}{N_2} \sum_{j=1}^{N_2} |C_s^{C_\circ^i} - \tilde{C}_s^{C_\circ^i}| \sqrt{E_s^{N_2}(C_\circ^j, W_0, W_0^{2,j}, x)} \\
&\leq K E_s^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) + \frac{K}{N_2} \sum_{j=1}^{N_2} E_s^{N_2}(C_\circ^j, W_0, W_0^{2,j}, x)
\end{aligned}$$

Therefore, using (5.9) and (5.10) we obtain

$$\begin{aligned} \mathbb{E} \left[(C_t^{C_\circ^i} - \tilde{C}_t^{C_\circ^i})^2 \right] &\leq \frac{K_1}{N_2} + K_2 \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] ds \\ &\quad + \frac{K_3}{N_2} \sum_{j=1}^{N_2} \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^j, W_0, W_0^{2,j}, x) \right] ds \end{aligned}$$

Using similar arguments, we can also show, using (5.7), that

$$\begin{aligned} \mathbb{E} \left[(W_t^{1, W_0} - \tilde{W}_t^{1, W_0})^2 \right] &\leq \frac{K_1}{N_2} + K_2 \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] ds \\ &\quad + \frac{K_3}{N_2} \sum_{j=1}^{N_2} \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^j, W_0, W_0^{2,j}, x) \right] ds \\ \mathbb{E} \left[(W_t^{2, C_\circ^i, W_0, W_0^{2,i}} - \tilde{W}_t^{2, C_\circ^i, W_0, W_0^{2,i}})^2 \right] &\leq \frac{K_1}{N_2} + K_2 \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] ds \\ &\quad + \frac{K_3}{N_2} \sum_{j=1}^{N_2} \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^j, W_0, W_0^{2,j}, x) \right] ds \quad (5.11) \end{aligned}$$

Therefore, we overall get that

$$\begin{aligned} &\mathbb{E} \left[(C_t^{C_\circ^i} - \tilde{C}_t^{C_\circ^i})^2 + (W_t^{1, W_0} - \tilde{W}_t^{1, W_0})^2 + (W_t^{2, C_\circ^i, W_0, W_0^{2,i}} - \tilde{W}_t^{2, C_\circ^i, W_0, W_0^{2,i}})^2 \right] \leq \\ &\leq \frac{K_1}{N_2} + K_2 \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] ds + \frac{K_3}{N_2} \sum_{j=1}^{N_2} \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^j, W_0, W_0^{2,j}, x) \right] ds \quad (5.12) \end{aligned}$$

5.6 Bound for $\mathbb{E} \left[(Z_t^{C_\circ^i}(x) - \tilde{Z}_t^{C_\circ^i}(x))^2 \right]$.

We next consider $(Z_t^{C_\circ^i}(x) - \tilde{Z}_t^{C_\circ^i}(x))^2$. For the following calculations, we define $\mathcal{F}_C^{N_2} = (C_\circ^1, C_\circ^2, \dots, C_\circ^{N_2})$. For $i = 1, 2, \dots, N_2$, we have

$$\begin{aligned} (Z_t^{C_\circ^i}(x) - \tilde{Z}_t^{C_\circ^i}(x))^2 &= \left(\mathbb{E} \left[W_t^{2, C_\circ^i, W_0, W_0^{2,i}} H_t^{1, W_0}(x) \middle| \mathcal{F}_C^{N_2} \right] - \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}_t^{2, C_\circ^i, w, u} \tilde{H}_t^{1, w}(x) \mu_{W^2}(du) \mu_{W^1}(dw) \right)^2 \\ &= \left(\mathbb{E} \left[(W_t^{2, C_\circ^i, W_0, W_0^{2,i}} - \tilde{W}_t^{2, C_\circ^i, W_0, W_0^{2,i}}) H_t^{1, W_0}(x) \middle| \mathcal{F}_C^{N_2} \right] \right. \\ &\quad \left. + \mathbb{E} \left[\tilde{W}_t^{2, C_\circ^i, W_0, W_0^{2,i}} H_t^{1, W_0}(x) \middle| \mathcal{F}_C^{N_2} \right] - \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}_t^{2, C_\circ^i, w, u} \tilde{H}_t^{1, w}(x) \mu_{W^2}(du) \mu_{W^1}(dw) \right)^2 \\ &\leq \mathbb{E} \left[(W_t^{2, C_\circ^i, W_0, W_0^{2,i}} - \tilde{W}_t^{2, C_\circ^i, W_0, W_0^{2,i}})^2 \middle| \mathcal{F}_C^{N_2} \right] \\ &\quad + \left(\mathbb{E} \left[[\tilde{W}_t^{2, C_\circ^i, W_0, W_0^{2,i}} H_t^{1, W_0}(x) - \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \tilde{W}_t^{2, C_\circ^i, w, u} \tilde{H}_t^{1, w}(x) \mu_{W^2}(du) \mu_{W^1}(dw)] \middle| \mathcal{F}_C^{N_2} \right] \right)^2 \\ &\leq \mathbb{E} \left[(W_t^{2, C_\circ^i, W_0, W_0^{2,i}} - \tilde{W}_t^{2, C_\circ^i, W_0, W_0^{2,i}})^2 \middle| \mathcal{F}_C^{N_2} \right] \\ &\quad + \left(\mathbb{E} \left[[\tilde{W}_t^{2, C_\circ^i, W_0, W_0^{2,i}} H_t^{1, W_0}(x) - \tilde{W}_t^{2, C_\circ^i, W_0, W_0^{2,i}} \tilde{H}_t^{1, W_0}(x)] \middle| \mathcal{F}_C^{N_2} \right] \right)^2 \\ &\leq \mathbb{E} \left[(W_t^{2, C_\circ^i, W_0, W_0^{2,i}} - \tilde{W}_t^{2, C_\circ^i, W_0, W_0^{2,i}})^2 \middle| \mathcal{F}_C^{N_2} \right] + K_0 \mathbb{E} \left[(W_t^{1, W_0} - \tilde{W}_t^{1, W_0})^2 \middle| \mathcal{F}_C^{N_2} \right], \end{aligned}$$

where the uniform bounds from (5.5) were used together with the compactness of the state space assumption from Assumption 2.1.

Therefore, using iterated expectations, we have that

$$\begin{aligned}
\mathbb{E} \left[(Z_t^{C_\circ^i}(x) - \tilde{Z}_t^{C_\circ^i}(x))^2 \right] &\leq \mathbb{E} \left[(W_t^{2,C_\circ^i, W_0, W_0^{2,i}} - \tilde{W}_t^{2,C_\circ^i, W_0, W_0^{2,i}})^2 \right] + K_0 \mathbb{E} \left[(W_t^{1, W_0} - \tilde{W}_t^{1, W_0})^2 \right] \\
&\leq \frac{K_1}{N_2} + K_2 \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] ds \\
&\quad + \frac{K_3}{N_2} \sum_{j=1}^{N_2} \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^j, W_0, W_0^{2,j}, x) \right] ds
\end{aligned} \tag{5.13}$$

where (5.11) was used.

Finally, using the assumption that $\sigma(\cdot)$ is globally Lipschitz,

$$\begin{aligned}
\mathbb{E} \left[(H_t^{2,C_\circ^i}(x) - \tilde{H}_t^{2,C_\circ^i}(x))^2 \right] &\leq K_0 \mathbb{E} \left[(Z_t^{C_\circ^i} - \tilde{Z}_t^{C_\circ^i})^2 \right] \\
&\leq \frac{K_1}{N_2} + K_2 \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] ds \\
&\quad + \frac{K_3}{N_2} \sum_{j=1}^{N_2} \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^j, W_0, W_0^{2,j}, x) \right] ds,
\end{aligned} \tag{5.14}$$

for possibly different, but finite constants $K_1, K_2, K_3 < \infty$.

5.7 Bound for $\mathbb{E} \left[E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right]$

Collecting our results from (5.12), (5.13) and (5.14), together with the definition of the error function $E_t^{N_2}$, we have, for $i = 1, \dots, N_2$, the bound

$$\begin{aligned}
\sup_{x \in \mathcal{X}} \mathbb{E} \left[E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] &\leq \frac{K_1}{N_2} + K_2 \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] ds, \\
&\quad + \frac{K_3}{N_2} \sum_{j=1}^{N_2} \int_0^t \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_s^{N_2}(C_\circ^j, W_0, W_0^{2,j}, x) \right] ds
\end{aligned}$$

Therefore, by averaging over all $i = 1, \dots, N_2$ and then using Gronwall's inequality, we get for any $0 \leq t \leq T$,

$$\frac{1}{N_2} \sum_{i=1}^{N_2} \sup_{x \in \mathcal{X}} \mathbb{E} \left[E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] \leq \frac{K}{N_2}.$$

for an appropriate constant $K < \infty$. Combining the last two displays we naturally get, again using Gronwall's inequality, that for $i = 1, \dots, N_2$

$$\sup_{x \in \mathcal{X}} \mathbb{E} \left[E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right] \leq \frac{K}{N_2}. \tag{5.15}$$

where the constant $K < \infty$ does not depend on i or N_2 .

5.8 Convergence of neural network prediction

The bound (5.15) of course proves the (uniform) convergence in probability of the neural network output $g_t^{N_2}(x) \rightarrow g_t(x)$. Recall, by (5.8), that

$$\mathbb{E} |g_t^{N_2}(x) - g_t(x)| \leq \mathbb{E} \left| \Gamma_t^{g,1}(x) \right| + \frac{K}{\sqrt{N_2}}.$$

where

$$\left| \Gamma_t^{g,1}(x) \right| \leq \frac{K}{N_2} \sum_{i=1}^{N_2} \sqrt{E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x)}.$$

Then, using the Cauchy-Schwartz inequality and (5.15),

$$\begin{aligned}
\mathbb{E} \left| \Gamma_t^{g,1}(x) \right| &\leq \frac{K}{N_2} \sum_{i=1}^{N_2} \mathbb{E} \left[\sqrt{E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x)} \right] \\
&\leq \frac{K}{N_2} \sum_{i=1}^{N_2} \sqrt{\mathbb{E} \left[E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right]} \\
&\leq \frac{K}{N_2} \sum_{i=1}^{N_2} \sqrt{\sup_{x \in \mathcal{X}} \mathbb{E} \left[E_t^{N_2}(C_\circ^i, W_0, W_0^{2,i}, x) \right]} \\
&\leq \frac{K}{\sqrt{N_2}}.
\end{aligned}$$

Therefore, for $0 \leq t \leq T$, and for some unimportant constant $K < \infty$

$$\sup_{x \in \mathcal{X}} \mathbb{E} \left[|g_t^{N_2}(x) - g_t(x)| \right] \leq \frac{K}{\sqrt{N_2}},$$

concluding the identification of the limit in Theorem 2.3.

6 Proof of Theorem 2.3-Uniqueness of the limit

Lemma 6.1. The solution to the limiting system (2.6) is unique in $C([0, 1], \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X})$.

Proof. Suppose there are two solutions to (2.6). Let's denote the first solution as $(\hat{W}_t^{1,w}, \hat{W}_t^{2,c,w,u}, \hat{C}_t^c, \hat{Z}_t^c(x), \hat{V}_t^w(x), \hat{g}_t(x))$ and the second solution as $(\bar{W}_t^{1,w}, \bar{W}_t^{2,c,w,u}, \bar{C}_t^c, \bar{Z}_t^c(x), \bar{V}_t^w(x), \bar{g}_t(x))$. Define the function

$$\begin{aligned}
Q_t = \sup_{(c,w,u,x) \in \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X}} &\left\{ (\hat{W}_t^{2,c,w,u} - \bar{W}_t^{2,c,w,u})^2 + \left\| \hat{W}_t^{1,w} - \bar{W}_t^{1,w} \right\|^2 + (\hat{V}_t^w(x) - \bar{V}_t^w(x))^2 + (\hat{Z}_t^c(x) - \bar{Z}_t^c(x))^2 \right. \\
&\left. + (\hat{C}_t^c - \bar{C}_t^c)^2 + (\hat{g}_t(x) - \bar{g}_t(x))^2 \right\}
\end{aligned}$$

Note that

$$Q_0 = 0.$$

We next study the evolution of Q_t for $t > 0$.

Using the same approach as in Section 5.2, we can show that any $\tilde{W}_t^{2,c,w,u}, \tilde{C}_t^c, \tilde{W}_t^{1,w}, \tilde{H}_t^{1,w}(x), \tilde{Z}_t^c, \tilde{V}_t^w$, and $g_t(x)$ which solve (2.6) are uniformly bounded on $\mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X} \times [0, 1]$.

We can then prove the inequality

$$\begin{aligned}
(\hat{W}_t^{2,c,w,u} - \bar{W}_t^{2,c,w,u})^2 &= 2 \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} (\hat{W}_s^{2,c,w,u} - \bar{W}_s^{2,c,w,u}) \left((y - \hat{g}_s(x)) \hat{C}_s^c \sigma'(\hat{Z}_s^c(x)) \hat{H}_s^{1,w}(x) \right. \\
&\quad \left. - (y - \bar{g}_s(x)) \bar{C}_s^c \sigma'(\bar{Z}_s^c(x)) \bar{H}_s^{1,w}(x) \right) \pi(dx, dy) ds \\
&\leq K \int_0^t Q_s ds.
\end{aligned}$$

where we have also used Young's inequality and the fact that $\mathcal{X} \times \mathcal{Y}$ is compact. Therefore,

$$\sup_{(c,w,u,x) \in \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X}} (\hat{W}_t^{2,c,w,u} - \bar{W}_t^{2,c,w,u})^2 \leq K \int_0^t Q_s ds.$$

Similarly,

$$\begin{aligned}
\sup_{(c,w,u,x) \in \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X}} \left\| \hat{W}_t^{1,w} - \bar{W}_t^{1,w} \right\|^2 &\leq K \int_0^t Q_s ds. \\
\sup_{(c,w,u,x) \in \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X}} (\hat{C}_t^c - \bar{C}_t^c)^2 &\leq K \int_0^t Q_s ds.
\end{aligned}$$

Next, using the Cauchy-Schwarz inequality and Young's inequality,

$$\begin{aligned}
(\hat{Z}_t^c(x) - \bar{Z}_t^c(x))^2 &= \left(\int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \hat{W}_t^{2,c,w,u} \hat{H}_t^{1,w}(x) \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) - \right. \\
&\quad \left. - \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \bar{W}_t^{2,c,w,u} \bar{H}_t^{1,w}(x) \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) \right)^2 \\
&= \left(\int_{\mathcal{W}^1} \int_{\mathcal{W}^2} (\hat{W}_t^{2,c,w,u} \hat{H}_t^{1,w}(x) - \bar{W}_t^{2,c,w,u} \bar{H}_t^{1,w}(x)) \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) \right)^2 \\
&\leq K \int_{\mathcal{W}^1} \int_{\mathcal{W}^2} \left[(\hat{W}_t^{2,c,w,u} - \bar{W}_t^{2,c,w,u}) + (\hat{W}_t^{1,w} - \bar{W}_t^{1,w})^2 \right] \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{W}^1}(dw) \\
&\leq K \sup_{(c,w,u,x) \in \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X}} \left[(\hat{W}_t^{2,c,w,u} - \bar{W}_t^{2,c,w,u}) + (\hat{W}_t^{1,w} - \bar{W}_t^{1,w})^2 \right] \\
&\leq K \int_0^t Q_s ds.
\end{aligned}$$

Using a similar approach,

$$\begin{aligned}
(\hat{V}_t^w(x) - \bar{V}_t^w(x))^2 &= \left(\int_{\mathcal{C}} \int_{\mathcal{W}^2} \hat{C}_t^c \sigma'(\hat{Z}_t^c(x)) \hat{W}_t^{2,c,w,u} \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{C}}(dc) \right. \\
&\quad \left. - \int_{\mathcal{C}} \int_{\mathcal{W}^2} \bar{C}_t^c \sigma'(\bar{Z}_t^c(x)) \bar{W}_t^{2,c,w,u} \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{C}}(dc) \right)^2 \\
&\leq \int_{\mathcal{C}} \int_{\mathcal{W}^2} \left(\hat{C}_t^c \sigma'(\hat{Z}_t^c(x)) \hat{W}_t^{2,c,w,u} - \bar{C}_t^c \sigma'(\bar{Z}_t^c(x)) \bar{W}_t^{2,c,w,u} \right)^2 \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{C}}(dc) \\
&\leq K \int_{\mathcal{C}} \int_{\mathcal{W}^2} \left[(\hat{C}_t^c - \bar{C}_t^c)^2 + (\hat{Z}_t^c(x) - \bar{Z}_t^c(x))^2 + (\hat{W}_t^{2,c,w,u} - \bar{W}_t^{2,c,w,u})^2 \right] \mu_{\mathcal{W}^2}(du) \mu_{\mathcal{C}}(dc) \\
&\leq K \sup_{(c,w,u,x) \in \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X}} \left[(\hat{C}_t^c - \bar{C}_t^c)^2 + (\hat{Z}_t^c(x) - \bar{Z}_t^c(x))^2 + (\hat{W}_t^{2,c,w,u} - \bar{W}_t^{2,c,w,u})^2 \right] \\
&\leq K \int_0^t Q_s ds.
\end{aligned}$$

Therefore,

$$\begin{aligned}
\sup_{(c,w,u,x) \in \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X}} (\hat{Z}_t^c(x) - \bar{Z}_t^c(x))^2 &\leq K \int_0^t Q_s ds, \\
\sup_{(c,w,u,x) \in \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X}} (\hat{V}_t^w(x) - \bar{V}_t^w(x))^2 &\leq K \int_0^t Q_s ds.
\end{aligned}$$

Finally,

$$\begin{aligned}
(\hat{g}_t(x) - \bar{g}_t(x))^2 &= \left(\int_{\mathcal{C}} \hat{C}_t^c \hat{H}_t^{2,c}(x) \mu_{\mathcal{C}}(dc) - \int_{\mathcal{C}} \bar{C}_t^c \bar{H}_t^{2,c}(x) \mu_{\mathcal{C}}(dc) \right)^2 \\
&\leq \int_{\mathcal{C}} \left[\hat{C}_t^c \hat{H}_t^{2,c}(x) \mu_{\mathcal{C}}(dc) - \bar{C}_t^c \bar{H}_t^{2,c}(x) \right]^2 \mu_{\mathcal{C}}(dc) \\
&\leq K \int_{\mathcal{C}} \left[(\hat{C}_t^c - \bar{C}_t^c)^2 + (\hat{Z}_t^c(x) - \bar{Z}_t^c(x))^2 \right] \mu_{\mathcal{C}}(dc) \\
&\leq K \int_0^t Q_s ds.
\end{aligned}$$

Consequently,

$$\sup_{(c,w,u,x) \in \mathcal{C} \times \mathcal{W}^1 \times \mathcal{W}^2 \times \mathcal{X}} (\hat{g}_t(x) - \bar{g}_t(x))^2 \leq K \int_0^t Q_s ds.$$

Collecting our results,

$$\begin{aligned} Q_t &\leq K \int_0^t Q_s ds, \\ Q_0 &= 0. \end{aligned}$$

Therefore, by Gronwall's inequality,

$$Q_t = 0,$$

for all $t \in [0, 1]$, completing the proof of the lemma. □

7 Conclusions and future work

In this paper, we have developed an approach that allows us to obtain the limiting behavior of multi-layer neural networks in the mean-field scaling as the number of units per layer and stochastic gradient steps grow to infinity. We have demonstrated that the limit is characterized by paths of weights that connecting the input layer to the output layer. The limit procedure shows that the main hyperparameters (the learning rates) must be chosen using a specific scaling (as a function of the number of hidden units per layer) in order to get a well-defined limit. A numerical study demonstrates that the mean-field scaling of the learning rates has a significant effect on how well the network can be trained. In addition, we have shown that the limit neural network seeks to minimize the limit objective function.

There are of course a number of open questions. The discussion in Section 4.3 shows that a genuine mean field formulation in terms of convergence of empirical distributions of trained parameters has yet to be understood. Secondly, our methodology works for convex error functions like the regression task. It is of interest to characterize the limit in other cases such as Wassenstein distances or regularized entropies. Thirdly, the universal approximation theorems say that a neural network approximator in the infinite unit per layer limit can approximate any continuous function [5, 23, 28, 29]. It would be of interest to investigate whether the limit procedure developed in this paper could lead to more specific information on convergence rates. Fourthly, the, commonly used in practice, unregularized algorithm and its finer properties (such as more detailed behavior as time grows to infinity) need to be better understood and we hope that the results in this paper build towards this direction.

A Limit of First Layer, Lemma 2.2

In this appendix we prove Lemma 2.2, which is about the limit of the first layer. The proof is analogous to [49]. Hence, instead of repeating all the arguments we will only present the general outline emphasizing the differences and refer the interested reader to [49] when the arguments are the same.

We let $N_1 \rightarrow \infty$ (with N_2 fixed). We want to prove that for each $N_2 \in \mathbb{N}$, $\tilde{\gamma}^{N_1, N_2} \xrightarrow{d} \gamma^{N_2}$ in $D_E([0, 1])$. $\gamma_t^{N_2}$ is a random measure-valued process which, for every $f \in C_b^2(\mathbb{R}^{d+2N_2})$, γ^{N_2} , satisfies the evolution equation (2.3)-(2.4).

A.1 Relative Compactness

Let's first establish a bound on C_k^i . Recall that $\sigma(\cdot)$ is bounded and $\alpha_C^{N_1, N_2} = \frac{N_2}{N_1}$. In the following calculations, the unimportant constants, K, K_0, K_1 , and K_2 may change from line to line.

For $k = 0, 1, \dots, \lfloor N_1 \rfloor$,

$$C_{k+1}^i = C_k^i + \frac{\alpha_C^N}{N_2} \left(y_k - \frac{1}{N_2} \sum_{m=1}^{N_2} C_k^m H_k^{2,m}(x_k) \right) H_k^{2,i}(x_k),$$

Since $\sigma(\cdot)$ is bounded, $|H_k^{2,i}| < K$. Therefore,

$$|C_{k+1}^i| \leq |C_k^i| + \frac{1}{N_1} \left(K_1 + \frac{1}{N_2} \sum_{m=1}^{N_2} |C_k^m| \right).$$

This yields

$$\frac{1}{N_2} \sum_{i=1}^{N_2} |C_{k+1}^i| \leq \frac{1}{N_2} \sum_{i=1}^{N_2} |C_k^i| + \frac{1}{N_1} \left(K_1 + \frac{1}{N_2} \sum_{i=1}^{N_2} |C_k^i| \right).$$

This implies that

$$\begin{aligned} \frac{1}{N_2} \sum_{i=1}^{N_2} |C_{k+1}^i| &\leq \frac{1}{N_2} \sum_{i=1}^{N_2} |C_0^i| + K_1 \frac{k}{N_1} + \frac{K_2}{N_1} \sum_{j=1}^k \frac{1}{N_2} \sum_{i=1}^{N_2} |C_k^i| \\ &\leq K_0 + K_1 \frac{k}{N_1} + \frac{K_2}{N_1} \sum_{j=1}^k \frac{1}{N_2} \sum_{i=1}^{N_2} |C_k^i| \end{aligned}$$

By Gronwall's inequality, for $k \leq \lfloor N_1 \rfloor$,

$$\frac{1}{N_2} \sum_{i=1}^{N_2} |C_k^i| \leq K \exp(K).$$

Then, we also have that

$$\begin{aligned} |C_{k+1}^i| &\leq |C_k^i| + \frac{K}{N_1}. \\ &\leq K \frac{k}{N_1}. \end{aligned}$$

This immediately yields that for $k \leq \lfloor N_1 \rfloor$

$$|C_k^i| < K.$$

Let's now address the parameters $W_k^{2,i,j}$. Recall that $\alpha_{W,2}^{N_1,N_2} = N_2$. Then,

$$\begin{aligned} |W_{k+1}^{2,i,j}| &\leq |W_k^{2,i,j}| + \frac{\alpha_{W,2}^{N_1,N_2}}{N_1 N_2} \left| (y_k - g_{\theta_k}^N(x_k)) C_k^i \sigma'(Z_k^i) H_k^{1,i} \right| \\ &\leq |W_k^{2,i,j}| + \frac{K}{N_1}. \end{aligned}$$

Therefore, for $k \leq \lfloor N_1 \rfloor$,

$$|W_{k+1}^{2,i,j}| \leq K.$$

Similarly, we have

$$\begin{aligned} |W_{k+1}^{1,j}| &\leq |W_k^{1,j}| + \frac{1}{N_1} \left| (y_k - g_{\theta_k}^N(x_k)) \left(\frac{1}{N_2} \sum_{i=1}^{N_2} C_k^i \sigma'(Z_k^i) W_k^{2,i,j} \right) \sigma'(W_k \cdot x) x_j \right| \\ &\leq |W_k^{1,j}| + \frac{K}{N_1}. \end{aligned}$$

Therefore, we obtain

$$|W_k^{1,j}| \leq K.$$

Collecting our results, for all $k/N_1 \leq 1$ and for all $j = 1, \dots, N_2$, we have the uniform bound

$$|C_k^i| + \left\| W_k^{1,j} \right\| + |W_k^{2,i,j}| < K. \quad (\text{A.1})$$

We now prove relative compactness of the family $\{\gamma^{N_1, N_2}\}_{N_1 \in \mathbb{N}}$ in $D_E([0, 1])$ where $E = \mathcal{M}(\mathbb{R}^{d+2N_2})$. It is sufficient to show compact containment and regularity of the γ^{N_1, N_2} 's (see for example Chapter 3 of [22]).

Lemma A.1. For each $\eta > 0$ and $t \geq 0$, there is a compact subset \mathcal{K} of E such that

$$\sup_{N_1 \in \mathbb{N}, 0 \leq t \leq T} \mathbb{P}[\gamma_t^{N_1, N_2} \notin \mathcal{K}] < \eta.$$

Proof. This uniform bound (A.1) actually implies the stronger statement of compact support. In particular, notice that the set $[-K, K]^{d+2N_2}$ is compact, and define

$$\mathcal{K} = \{\omega \in M(\mathbb{R}^{d+2N_2}) : \omega([-K, K]^{d+2N_2}) = 1\}.$$

Then $\mathcal{K} \subset M(\mathbb{R}^{d+2N_2})$, and \mathbb{P} -a.s. $\gamma_t^{N_1, N_2} \in \mathcal{K}$ for all $N_1 \in \mathbb{N}$ and $t \in [0, 1]$. This concludes the proof. \square

We now establish regularity of the γ^{N_1, N_2} 's. Define the function $q(z_1, z_2) = \min\{|z_1 - z_2|, 1\}$ where $z_1, z_2 \in \mathbb{R}$.

Lemma A.2. There is a constant $K < \infty$ such that for $0 \leq u \leq \delta$, $0 \leq v \leq \delta \wedge t$, $t \in [0, 1]$,

$$\mathbb{E} \left[q(\langle f, \gamma_{t+u}^{N_1, N_2} \rangle, \langle f, \gamma_t^{N_1, N_2} \rangle) q(\langle f, \gamma_t^{N_1, N_2} \rangle, \langle f, \gamma_{t-v}^{N_1, N_2} \rangle) | \mathcal{F}_t^N \right] \leq C\delta + \frac{K}{N_1}.$$

Proof. Let $0 \leq s < t \leq T$ and $\delta < 1$, such that $t - s < \delta < 1$. We then have

$$\begin{aligned} |C_{[N_1 t]}^i - C_{[N_1 s]}^i| &= \left| \sum_{k=[N_1 s]}^{[N_1 t]-1} (C_{k+1}^i - C_k^i) \right| \\ &\leq \sum_{k=[N_1 s]}^{[N_1 t]-1} \frac{1}{N_1} |(y_k - g_{\theta_k}^{N_1, N_2}(x_k)) H_k^{2,i}(x_k)| \\ &\leq \frac{1}{N_1} \sum_{k=[N_1 s]}^{[N_1 t]-1} K \\ &\leq K_1 \delta + \frac{K_2}{N_1}. \end{aligned}$$

Using the same approach, we can establish similar bounds for the other parameters:

$$\begin{aligned} |W_{[N_1 t]}^{1,j} - W_{[N_1 s]}^{1,j}| &\leq K_1 \delta + \frac{K_2}{N_1}, \\ |W_{[N_1 t]}^{2,i,j} - W_{[N_1 s]}^{2,i,j}| &\leq K_1 \delta + \frac{K_2}{N_1}. \end{aligned}$$

The desired result then follows. \square

We conclude this section now with the required relative compactness of the sequence $\{\gamma^{N_1, N_2}\}_{N_1 \in \mathbb{N}}$. This implies that for each fixed N_2 , every subsequence γ^{N_1, N_2} 's has a convergent sub-subsequence as $N_1 \rightarrow \infty$.

Lemma A.3. The sequence of probability measures $\{\gamma^{N_1, N_2}\}_{N_1 \in \mathbb{N}}$ is relatively compact in $D_E([0, 1])$.

Proof. Given Lemmas A.1 and A.2, Theorem 8.6 and Remark 8.7 B of Chapter 3 of [22], gives the statement of the lemma. \square

A.2 Identification of the Limit

We consider the evolution of the empirical measure $\gamma_t^{N_1, N_2}$ via test functions $f \in C_b^2(\mathbb{R}^{d+2N_2})$. Using a Taylor expansion based argument similarly to [49], we can show that the scaled empirical measure satisfies

$$\begin{aligned}
\langle f, \gamma_t^{N_1, N_2} \rangle - \langle f, \gamma_0^{N_1, N_2} \rangle &= \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} (y - g_s^{N_1, N_2}(x)) \langle H_s^{N_1, N_2}(x) \cdot \nabla_c f, \gamma_s^{N_1, N_2} \rangle \pi(dx, dy) ds \\
&+ \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} (y - g_s^{N_1, N_2}(x)) \langle \sigma(w^1 \cdot x) (\sigma'(Z_s^{N_1, N_2}(x)) \odot c) \cdot \nabla_{w^2} f, \gamma_s^{N_1, N_2} \rangle \pi(dx, dy) ds \\
&+ \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} (y - g_s^{N_1, N_2}(x)) \frac{1}{N_2} \sum_{i=1}^{N_2} \langle c_i \sigma'(Z_s^{N_1, N_2}(x)) w^{2,i} \sigma'(w^1 \cdot x) x \cdot \nabla_{w^1} f, \gamma_s^{N_1, N_2} \rangle \pi(dx, dy) ds \\
&+ M^{N_1, N_2}(t) + O(N_1^{-1}),
\end{aligned} \tag{A.2}$$

where $M^{N_1, N_2}(t)$ is a martingale term,

$$\begin{aligned}
Z_s^{i, N_1, N_2}(x) &= \langle w^{2,i} \sigma(w^1 \cdot x), \gamma_s^{N_1, N_2} \rangle, \\
H_s^{i, N_1, N_2}(x) &= \sigma(Z_s^{i, N_1, N_2}(x)), \\
g_s^{N_1, N_2}(x) &= \frac{1}{N_2} \sum_{i=1}^{N_2} H_s^{i, N_1, N_2}(x) \langle c_i, \gamma_s^{N_1, N_2} \rangle.
\end{aligned}$$

and we have defined $H_t^{N_1, N_2}$ as the vector $(H_t^{1, N_1, N_2}, \dots, H_t^{N_2, N_1, N_2})$, c as the vector (c_1, \dots, c_{N_2}) , and $Z_t^{N_1, N_2}$ as the vector $(Z_t^{1, N_1, N_2}, \dots, Z_t^{N_2, N_1, N_2})$. The martingale term $M^{N_1, N_2}(t)$ converges to 0 in L^2 as $N_1 \rightarrow \infty$. That is,

$$\lim_{N_1 \rightarrow \infty} \mathbb{E} \left[\left(M^{N_1, N_2}(t) \right)^2 \right] = 0. \tag{A.3}$$

The proof for (A.3) follows as in Lemma 3.1 of [49] and thus it is omitted. The limit point of γ^{N_1, N_2} , as $N_1 \rightarrow \infty$ and for a fixed N_2 , will satisfy the evolution equation

$$\begin{aligned}
\langle f, \gamma_t^{N_2} \rangle - \langle f, \gamma_0^{N_2} \rangle &= \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} \alpha(y - g_s^{N_2}(x)) \langle H_s^{N_2}(x) \cdot \nabla_c f, \gamma_s^{N_2} \rangle \pi(dx, dy) ds \\
&+ \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} \alpha(y - g_s^{N_2}(x)) \langle \sigma(w^1 \cdot x) (\sigma'(Z_s) \odot c) \cdot \nabla_{w^2} f, \gamma_s^{N_2} \rangle \pi(dx, dy) ds \\
&+ \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} \alpha(y - g_s^{N_2}(x)) \frac{1}{N_2} \sum_{i=1}^{N_2} \langle c_i \sigma'(Z_s^{i, N_2}(x)) w^{2,i} \sigma'(w^1 \cdot x) x \cdot \nabla_{w^1} f, \gamma_s^{N_2} \rangle \pi(dx, dy) ds,
\end{aligned} \tag{A.4}$$

where

$$\begin{aligned}
Z_s^{i, N_2}(x) &= \langle w^{2,i} \sigma(w^1 \cdot x), \gamma_s^{N_2} \rangle, \\
H_s^{i, N_2}(x) &= \sigma(Z_s^{i, N_2}(x)), \\
g_s^{N_2}(x) &= \frac{1}{N_2} \sum_{i=1}^{N_2} H_s^{i, N_2}(x) \langle c_i, \gamma_s^{N_2} \rangle,
\end{aligned}$$

and

$$\gamma_0^{N_2}(dw^1, dw^2, dc) = \mu_{W^1}(dw^1) \times \mu_{W^2}(dw^{2,1}) \times \dots \times \mu_{W^2}(dw^{2, N_2}) \times \delta_{C_1^1}(dc^1) \times \dots \times \delta_{C_1^{N_2}}(dc^{N_2}).$$

Let π^{N_1, N_2} be the probability measure of $(\gamma_t^{N_1, N_2})_{0 \leq t \leq 1}$. Each π^{N_1, N_2} takes values in the set of probability measures $\mathcal{M}(D_E([0, 1]))$. Relative compactness, proven in Section A.1, implies that there is a subsequence $\pi^{N_{1_k}, N_2}$ which weakly converges. We must prove that any limit point π^{N_2} of a convergent subsequence $\pi^{N_{1_k}, N_2}$ will satisfy the evolution equation (A.4).

Lemma A.4. Let π^{N_1, N_2} be a convergent subsequence with a limit point π^{N_2} . Then, π^{N_2} is a Dirac measure concentrated on $\gamma^{N_2} \in D_E([0, 1])$ and γ^{N_2} satisfies the measure evolution equation (A.4).

Proof. We define a map $F(\gamma) : D_E([0, 1]) \rightarrow \mathbb{R}_+$ for each $t \in [0, T]$, $f \in C_b^2(\mathbb{R}^{d+2N_2})$, $g_1, \dots, g_p \in C_b(\mathbb{R}^{d+2N_2})$ and $0 \leq s_1 < \dots < s_p \leq t$.

$$\begin{aligned}
F(\gamma) &= \left| \left(\langle f, \gamma_t \rangle - \langle f, \gamma_0 \rangle - \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} \alpha(y - g_s^{N_2}(x)) \langle H_t^{N_2}(x) \cdot \nabla_c f, \gamma_s^{N_2} \rangle \pi(dx, dy) ds \right. \right. \\
&\quad - \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} \alpha(y - g_s^{N_2}(x)) \langle \sigma(w^1 \cdot x) (\sigma'(Z_s) \odot c) \cdot \nabla_{w^2} f, \gamma_s^{N_2} \rangle \pi(dx, dy) ds \\
&\quad - \left. \int_0^t \int_{\mathcal{X} \times \mathcal{Y}} \alpha(y - g_s^{N_2}(x)) \frac{1}{N_2} \sum_{i=1}^{N_2} \langle c_i \sigma'(Z_s^{i, N_2}(x)) w^{2, i} \sigma'(w^1 \cdot x) x \cdot \nabla_{w^1} f, \gamma_s^{N_2} \rangle \pi(dx, dy) ds \right) \\
&\quad \times \langle g_1, \gamma_{s_1} \rangle \times \dots \times \langle g_p, \gamma_{s_p} \rangle \Big|. \tag{A.5}
\end{aligned}$$

Then,

$$\begin{aligned}
\mathbb{E}_{\pi^{N_1, N_2}}[F(\gamma)] &= \mathbb{E}[F(\gamma^{N_1, N_2})] \\
&= \mathbb{E} \left| \left(M^{N_1, N_2}(t) + O(N_1^{-1}) \prod_{i=1}^p \langle g_i, \gamma_{s_i}^{N_1, N_2} \rangle \right) \right| \\
&\leq \mathbb{E}[|M^{N_1, N_2}(t)|] + O(N_1^{-1}) \\
&\leq \mathbb{E}[(M^{N_1, N_2}(t))^2]^{1/2} + O(N_1^{-1}) \\
&\leq K \left(\frac{1}{\sqrt{N}} + \frac{1}{N} \right).
\end{aligned}$$

Therefore,

$$\lim_{N_1 \rightarrow \infty} \mathbb{E}_{\pi^{N_1, N_2}}[F(\gamma)] = 0.$$

Since $F(\cdot)$ is continuous and $F(\gamma^{N_1, N_2})$ is uniformly bounded (due to the uniform boundedness results of Section A.1),

$$\mathbb{E}_{\pi^{N_2}}[F(\gamma)] = 0.$$

Since this holds for each $t \in [0, T]$, $f \in C_b^2(\mathbb{R}^{d+2N_2})$ and $g_1, \dots, g_p \in C_b(\mathbb{R}^{d+2N_2})$, γ^{N_2} satisfies the evolution equation (A.4). \square

It remains to prove that the evolution equation (A.4) has a unique solution. This is the content of Section A.3.

A.3 Uniqueness

Lemma A.5. There exists a unique solution to the evolution equation (A.4).

Proof. We only sketch the proof as the argument is similar to the uniqueness argument of [49]. Consider the

particle system:

$$\begin{aligned}
dC_t^i &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) H_t^{2,i}(x) \pi(dx, dy) dt, \quad C_0^i = C_\circ^i, \quad i = 1, \dots, N_2, \\
dW_t^1 &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) \left(\frac{1}{N_2} \sum_{i=1}^{N_2} C_t^i \sigma'(Z_t^i(x)) W_t^{2,i} \right) \sigma'(W_t^1 \cdot x) x \pi(dx, dy) dt, \quad W_0^1 \sim \mu_{W^1}(dw), \\
dW_t^{2,i} &= \int_{\mathcal{X} \times \mathcal{Y}} (y - g_t^{N_2}(x)) C_t^i \sigma'(Z_t^i) H_t^1(x) \pi(dx, dy) dt, \quad W_0^{2,i} \sim \mu_{W^2}(dw^2), \quad i = 1, \dots, N_2, \\
H_t^1(x) &= \sigma(W_t^1 \cdot x), \\
Z_t^i(x) &= \left\langle w^{2,i} H_t^1(x), \gamma_t^{N_2} \right\rangle, \\
H_t^{2,i}(x) &= \sigma(Z_t^i(x)), \\
g_t^{N_2}(x) &= \frac{1}{N_2} \sum_{i=1}^{N_2} C_t^i H_t^{2,i}(x). \tag{A.6}
\end{aligned}$$

Let $\nu_{t,(c_1, \dots, c_{N_2})}$ be the conditional law of $(W_t^1, W_t^{2,1}, \dots, W_t^{2,N_2}, C_t^1, \dots, C_t^{N_2})_{0 \leq t \leq T}$ given $(C_\circ^1, \dots, C_\circ^{N_2}) = (c_1, \dots, c_{N_2})$. Similarly to the general results of [34], $\nu_{t,(C_\circ^1, \dots, C_\circ^{N_2})} = \gamma_t^{N_2}$ is a solution to the evolution equation (A.4) and conversely, any solution $\nu_{t,(C_\circ^1, \dots, C_\circ^{N_2})}$ to the evolution equation (A.4) must also be the law of the solution to (A.6).

Let us next prove that we can write $Z_t^i(x) = \mathbb{E} \left[W_t^{2,i} H_t^1(x) \middle| C_\circ^1, \dots, C_\circ^{N_2} \right]$. Recall that $\mathcal{F}_C^{N_2} = (C_\circ^1, \dots, C_\circ^{N_2})$ and let us similarly define $\mathcal{F}_{W^2}^{N_2} = (W_0^{2,1}, \dots, W_0^{2,N_2})$. Inspecting (A.6) it becomes clear that the random variables $(W_t^1, W_t^{2,1}, \dots, W_t^{2,N_2}, C_t^1, \dots, C_t^{N_2})_{0 \leq t \leq T}$ depend, in addition to their own initial conditions, also on $\mathcal{F}_C^{N_2}$ and $\mathcal{F}_{W^2}^{N_2}$ in a symmetric way through the terms $g_t^{N_2}(x)$ and $\frac{1}{N_2} \sum_{i=1}^{N_2} C_t^i \sigma'(Z_t^i(x)) W_t^{2,i}$. In order to make this dependence specific in the notation, we write in particular that

$$W_t^1 = W_t^{1, W_0; \mathcal{F}_C^{N_2}, \mathcal{F}_{W^2}^{N_2}}, \quad W_t^{2,i} = W_t^{2, C_\circ^i, W_0, W_0^{2,i}; \mathcal{F}_C^{N_2}, \mathcal{F}_{W^2}^{N_2}} \quad \text{and} \quad C_t^i = C_t^{C_\circ^i; \mathcal{F}_C^{N_2}}.$$

This, then leads to the following calculations

$$\begin{aligned}
Z_t^i(x) &= \left\langle w^{2,i} H_t^1(x), \gamma_t^{N_2} \right\rangle \\
&= \mathbb{E} \left[W_t^{2,i} \sigma(W_t^1 \cdot x) \right] \\
&= \mathbb{E} \left[W_t^{2, C_\circ^i, W_0, W_0^{2,i}; \mathcal{F}_C^{N_2}, \mathcal{F}_{W^2}^{N_2}} \sigma(W_t^{1, W_0; \mathcal{F}_C^{N_2}, \mathcal{F}_{W^2}^{N_2}} \cdot x) \right] \\
&= \mathbb{E} \left[\mathbb{E} \left[W_t^{2, C_\circ^i, W_0, W_0^{2,i}; \mathcal{F}_C^{N_2}, \mathcal{F}_{W^2}^{N_2}} \sigma(W_t^{1, W_0; \mathcal{F}_C^{N_2}, \mathcal{F}_{W^2}^{N_2}} \cdot x) \middle| W_0, \mathcal{F}_{W^2}^{N_2}, \mathcal{F}_C^{N_2} \right] \right] \\
&= \left\langle W_t^{2, C_\circ^i, W_0, W_0^{2,i}; \mathcal{F}_C^{N_2}, \mathcal{F}_{W^2}^{N_2}} \sigma(W_t^{1, W_0; \mathcal{F}_C^{N_2}, \mathcal{F}_{W^2}^{N_2}} \cdot x), \gamma_0^{N_2} \right\rangle \\
&= \mathbb{E} \left[W_t^{2, C_\circ^i, W_0, W_0^{2,i}; \mathcal{F}_C^{N_2}, \mathcal{F}_{W^2}^{N_2}} \sigma(W_t^{1, W_0; \mathcal{F}_C^{N_2}, \mathcal{F}_{W^2}^{N_2}} \cdot x) \middle| \mathcal{F}_C^{N_2} \right] \\
&= \mathbb{E} \left[W_t^{2,i} \sigma(W_t^1 \cdot x) \middle| C_\circ^1, \dots, C_\circ^{N_2} \right]
\end{aligned}$$

The second to the last equality is due to the assumed independence of the initial conditions via Assumption 2.1 together with the fact that the marginal of $\gamma_0^{N_2}$ with respect to $C_\circ^1, \dots, C_\circ^{N_2}$ is a product of delta Dirac distributions.

It remains to show that the solution to (A.6) is unique. We do this via a fixed point argument, completely analogous to Section 4 of [49]. The details are omitted due to the similarity of the argument. \square

A.4 Convergence Result for First Layer

Let π^{N_1, N_2} be the probability measure corresponding to γ^{N_1, N_2} . Each π^{N_1, N_2} takes values in the set of probability measures $\mathcal{M}(D_E([0, 1]))$. Relative compactness, proven in Section A.1, implies that every subsequence $\pi^{N_{1_k}, N_2}$ has a further sub-sequence $\pi^{N_{1_{k_m}}, N_2}$ which weakly converges. Section A.2 proves that any limit point π^{N_2} of $\pi^{N_{1_{k_m}}, N_2}$ will satisfy the evolution equation (A.4). Section A.3 proves that the solution of the evolution equation (A.4) is unique. Therefore, by Prokhorov's Theorem, π^{N_1, N_2} weakly converges to π^{N_2} , where π^{N_2} is the distribution of γ^{N_2} , the unique solution of (A.4). That is, γ^{N_1, N_2} converges in distribution to γ^{N_2} .

Bibliography

- [1] L. Ambrosio, N. Gigli, and G. Savaré. Gradient flows: in metric spaces and in the space of probability measures. Springer Science & Business Media, 2008.
- [2] D. Araújo, R. I. Oliveira and D. Yukimura. A mean-field limit for certain deep neural networks. 2019, arXiv: 1906.00193.
- [3] B. Alipanahi, A. Delong, M. Weirauch, and B. Frey. Predicting the sequence specificities of DNA-and RNA-binding proteins by deep learning. *Nature Biotechnology*, 33(8): 831, 2015.
- [4] S. Arik, M. Chrzanowski, A. Coates, G. Diamos, A. Gibiansky, Y. Kang, X. Li, J. Miller, A. Ng, J. Raiman, S. Sengputa. Deep voice: Real-time neural text-to-speech. arXiv:1702.07825., 2017.
- [5] A. Barron. Approximation and estimation bounds for artificial neural networks. *Machine Learning*, 14(1), 115-133, 1994.
- [6] M. Bojarski, D. Del Test, D. Dworakowski, B. Firnier, B. Flepp, P. Goyal, L. Jackel, M. Monfort, U. Muller, J. Zhang, and X. Zhang. End to end learning for self-driving cars, arXiv:1604.07316, 2016.
- [7] A. Blanchet and J. Bolte. A family of functional inequalities, Łojasiewicz inequalities and displacement convex functions. (2016), arXiv:1612.02619.
- [8] L. Chizat. Sparse Optimization on Measures with Over-parameterized Gradient Descent. (2019), arXiv preprint: 1907.10300.
- [9] L. Chizat, and F. Bach. A note on lazy training in supervised differentiable programming. (2018), arXiv:1812.07956.
- [10] L. Chizat, and F. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in Neural Information Processing Systems (NeurIPS)*. (2018), pp. 3040-3050.
- [11] G. Cybenko. Approximation by superposition of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2, (1989), pp. 303–314.
- [12] D.A. Dawson, Hierarchical and mean-eld stepping stone models, in *Progress in Population Genetics and Human Evolution*, eds. P. Donnelly and S. Tavare, Springer, IMA Volume in Mathematics and its Applications vol. 87. (1997).
- [13] D.A. Dawson, Multilevel mutation-selection systems and set-valued duals, *Journal of Mathematical Biology*, Vol. 76, Issue 1–2, (2018), pp 295–378
- [14] D.A. Dawson and K.J. Hochberg Wandering random measures in the Fleming-Viot model, *Annals of Probability* 10 (1982), pp. 554-580.
- [15] D.A. Dawson and K.J. Hochberg, A multilevel branching model, *Advances in Applied Probability* 23, (1991), pp. 701-715.
- [16] D.A. Dawson, K.J. Hochberg and V. Vinogradov, High-density limits of hierarchically structured branching-diuising populations, *Stochastic Processes and their Applications*, 62, (1996), pp. 191-222.

- [17] D.A. Dawson and J. Gärtner, Analytic aspects of multilevel large deviations, in *Asymptotic Methods in Probability and Statistics* (ed. B. Szyszkowicz), Elsevier, Amsterdam, (1998), pp. 401-440.
- [18] D.A. Dawson and A. Greven, Hierarchically interacting Fleming-Viot processes with selection and mutation: multiple space-time scale analysis and quasi equilibria, *Electronic J. of Prob.*, vol. 4, paper 4, (1993), pp. 1-81.
- [19] D.A. Dawson and Y. Wu, Multilevel multitype models of an information system, *I.M.A. Volume 84, Classical and Modern Branching Processes*, eds. K.B. Athreya and P. Jagers, Springer-Verlag, (1996), pp. 57-72.
- [20] S. Du, J. Lee, H. Li, L. Wang, and X. Zhai. Gradient Descent Finds Global Minima of Deep Neural Networks. arXiv: 1811.03804, 2019.
- [21] S. Du, X. Zhai, B. Póczos, and A. Singh. Gradient Descent Provably Optimizes Over-parameterized Neural Networks. arXiv: 1810.02054, 2019.
- [22] S. Ethier and T. Kurtz. *Markov Processes: Characterization and Convergence*. 1986, Wiley, New York, MR0838085.
- [23] Ken-Ichi Funahashi On the approximate realization of continuous mappings by neural networks. *Neural Networks*, Vol. 2, Issue 3, (1989), pp. 183-192.
- [24] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. Cambridge: MIT Press, 2016.
- [25] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pp. 249-256, 2010.
- [26] S. Gu, E. Holly, T. Lillicrap, and S. Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. *IEEE Conference on Robotics and Automation*, 3389-3396, 2017.
- [27] D. Hauer and J. Mazón. Kurdyka-Łojasiewicz-Simon inequality for gradient flows in metric spaces. *Trans. Amer. Math. Soc.*, Vol. 372, (2019), pp. 4917–4976
- [28] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5), 359-366, 1989.
- [29] K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2), 251-257, 1991.
- [30] K. Hu, Z. Ren, D. Šiška, and L. Szpruch. Mean-Field Langevin Dynamics and Energy Landscape of Neural Networks arXiv preprint: 1905.07769.
- [31] A. Jacot, F. Gabriel, and C. Hongler. Neural Tangent Kernel: Convergence and Generalization in Neural Networks. *32nd Conference on Neural Information Processing Systems (NeurIPS 2018)*, Montreal, Canada.
- [32] S.G. Krantz and H.R. Parks. *A Primer of Real Analytic Functions*, second edition. Springer Science and Business Media, New York , 2002
- [33] A. Krizhevsky Learning Multiple Layers of Features from Tiny Images, Technical Report, 2009.
- [34] V.N. Kolokoltsov. *Nonlinear Markov processes and kinetic equations* Vol. 182, Cambridge University Press, 2010.
- [35] Y. LeCun, Y. Bengio, and G. Hinton. Deep Learning. *Nature*, 521(7553), 436, 2015.
- [36] Y. Leviathan and Y. Matias. Google Duplex: An AI System for Accomplishing Real-World Tasks Over the Phone. Google, 2018.

- [37] J. Ling, A. Kurzawski, and J. Templeton. Reynolds averaged turbulence modelling using deep neural networks with embedded invariance. *Journal of Fluid Mechanics*, 807, 155-166, 2016.
- [38] J. Ling, R. Jones, and J. Templeton. Machine learning strategies for systems with invariance properties. *Journal of Computational Physics*, 318, 22-35, 2016.
- [39] D. Li, T. Ding, and R. Sun. Over-Parameterized Deep Neural Networks Have No Strict Local Minima For Any Continuous Activations. arXiv: 1812.11039, 2018.
- [40] S. Liang, R. Sun, J. Lee, and R. Srikant. Adding One Neuron Can Eliminate All Bad Local Minima. NIPS, 2018.
- [41] S. Liang, R. Sun, Y. Li, and R. Srikant. Understanding the Loss Surface of Neural Networks for Binary Classification. ICML, 2018.
- [42] S. Mallat. Understanding deep convolutional neural networks. *Philosophical Transactions of the Royal Society A*. 374.2065, 20150203, 2016.
- [43] S. Mei, A. Montanari, and P. Nguyen A mean field view of the landscape of two-layer neural networks *Proceedings of the National Academy of Sciences*, Vol. 115, Issue 33, (2018), pp. E7665-E767.
- [44] S. Mei, T. Misiakiewicz and A. Montanari A mean field theory of two-layers neural networks: dimension free bounds and kernel limit. 2019, arXiv:1902.06015.
- [45] P.-M. Nguyen. Mean Field Limit of the Learning Dynamics of Multilayer Neural Networks. 2019, arXiv:1902.02880.
- [46] R. Nallapati, B. Zhou, C. Gulcehre, and B. Xiang. Abstractive text summarization using sequence-to-sequence RNNs and beyond. arXiv:1602.06023, 2016.
- [47] H. Pierson and M. Gashler. Deep learning in robotics: a review of recent research. *Advanced Robotics*, 31(16): 821-835, 2017.
- [48] G. M. Rotskoff and E. Vanden-Eijnden Neural Networks as Interacting Particle Systems: Asymptotic Convexity of the Loss Landscape and Universal Scaling of the Approximation Error. 2018, arXiv:1805.00915.
- [49] J. Sirignano and K. Spiliopoulos. Mean Field Analysis of Neural Networks: A Law of Large Numbers *SIAM Journal on Applied Mathematics*, Vol. 80, Issue 2, (2020), pp. 725–752.
- [50] J. Sirignano and K. Spiliopoulos. Mean Field Analysis of Neural Networks: A central limit theorem *Stochastic Processes and their applications*, Vol. 130, Issue 3, (2020), pp. 1820–1852.
- [51] J. Sirignano, A. Sadhwani, and K. Giesecke. Deep Learning for Mortgage Risk. arXiv:1607.02470, 2016.
- [52] J. Sirignano and R. Cont. Universal features of price formation in financial markets: perspectives from Deep Learning. arXiv:1803.06917, 2018.
- [53] J. Sirignano and K. Spiliopoulos, Stochastic gradient descent in continuous time, *SIAM Journal on Financial Mathematics*, Vol. 8, Issue 1, (2017), pp. 933–961.
- [54] J. Sirignano and K. Spiliopoulos, DGM: A deep learning algorithm for solving partial differential equations, *Journal of Computational Physics*, (2018), Vol. 375, pp. 1339-1364.
- [55] A-S. Sznitman. Topics in propagation of chaos. in *Ecole d’Eté de Probabilités de Saint-Flour XIX - 1989*. series, Lecture Notes in Mathematics, P.-L. Hennequin, Ed. Springer, Berlin Heidelberg. 1464, 165-251, 1991.
- [56] A-S. Sznitman. Nonlinear reflecting diffusion processes, and the propagation of chaos and fluctuations associated. *Journal of Functional Analysis*. 56, 311-336, 1984.

- [57] I. Sutskever, O. Vinyals, and Q. Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, 3104-3112, 2014.
- [58] N. Sunderhauf, O. Brock, W. Cheirer, R. Hadsell, D. Fox, J. Leitner, B. Upcroft, P. Abbeel, W. Burgard, M. Milford, and P. Corke. The limits and potentials of deep learning for robotics. *The International Journal of Robotics Research*, 37(4): 405-420, 2018.
- [59] Y. Taigman, M. Yang, M. Ranzato, L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1701-1708, 2014.
- [60] Y. Zhang, W. Chan, and N. Jaitly. Very deep convolutional networks for end-to-end speech recognition. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*. 4845-4849, 2017.
- [61] C. Wang, J. Mattingly, and Y. Lu. Scaling limit: Exact and tractable analysis of online learning algorithms with applications to regularized regression and PCA. 2017, arXiv:1712.04332.
- [62] Y. Wu, M. Schuster, Z. Chen, Q. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, and J. Klingner. Google's neural machine translation system: Bridging the gap between human and machine translation. arXiv:1609.08144, 2016.