

# Modeling from Features: a Mean-field Framework for Over-parameterized Deep Neural Networks

**Cong Fang**

*Department of Machine Intelligence, Peking University*

FANGCONG@PKU.EDU.CN

**Jason D. Lee**

*Electrical and Computer Engineering, Princeton University*

JASONLEE@PRINCETON.EDU

**Pengkun Yang**

*Center for Statistical Science, Tsinghua University*

YANGPENGKUN@TSINGHUA.EDU.CN

**Tong Zhang**

*Department of Mathematics and Computer Science, HKUST*

TONGZHANG@TONGZHANG-ML.ORG

**Editors:** Mikhail Belkin and Samory Kpotufe

## Abstract

This paper proposes a new mean-field framework for over-parameterized deep neural networks (DNNs), which can be used to analyze neural network training. In this framework, a DNN is represented by probability measures and functions over its features (that is, the function values of the hidden units over the training data) in the continuous limit, instead of the neural network parameters as most existing studies have done. This new representation overcomes the degenerate situation where all the hidden units essentially have only one meaningful hidden unit in each middle layer, leading to a simpler representation of DNNs. Moreover, we construct a non-linear dynamics called *neural feature flow*, which captures the evolution of an over-parameterized DNN trained by Gradient Descent. We illustrate the framework via the Residual Network (Res-Net) architecture. It is shown that when the neural feature flow process converges, it reaches a *global* minimal solution under suitable conditions.

**Keywords:** deep residual network, mean-field theory, non-linear dynamics, global minimum.

## 1. Introduction

In recent years, deep neural networks (DNNs) have achieved great success empirically. However, the theoretical understanding of the practical success is still limited. One main conceptual difficulty is the non-convexity of DNN models. More recently, there has been remarkable progress in understanding the over-parameterized neural networks (NNs), which are NNs with massive hidden units. The over-parameterization is capable of circumventing the hurdles in analyzing non-convex functions under specific settings:

- (i) Under a specific scaling and initialization, it is sufficient to study the NN weights in a small region around the initial values given sufficiently many hidden units - the aptly named “lazy training” regime (Jacot et al., 2018; Li and Liang, 2018; Du et al., 2019a; Arora et al., 2019; Du et al., 2019b; Allen-Zhu et al., 2018; Allen-Zhu and Li, 2019; Zou et al., 2018; Chizat et al., 2019). The NN in this regime is nearly a linear model fitted with a random kernel in the tangent space, and provably achieves minimum training error. However, this regime does not explain why NNs can effectively learn representative features, and the expressive power of random kernels is limited (Yehudai and Shamir, 2019).

- (ii) Another line of research applies the mean-field analysis for NNs (Mei et al., 2018; Chizat and Bach, 2018; Sirignano and Spiliopoulos, 2019a; Rotskoff and Vanden-Eijnden, 2018; Mei et al., 2019; Dou and Liang, 2019; Wei et al., 2018; Sirignano and Spiliopoulos, 2019b; Fang et al., 2019; Araújo et al., 2019; Nguyen and Pham, 2020; Chen et al., 2020). Learning a two-layer over-parameterized NN can be approximately described as optimizing a functional over probability distributions of the NN weights. The evolution of NN weights trained by the (noisy) Gradient Descent algorithm corresponds to a Wasserstein gradient flow called “distributional dynamics”, solution to a non-linear partial differential equation (PDE) of McKean-Vlasov type (Sznitman, 1991). In the mean-field limit, the Wasserstein gradient flow converges to the globally optimal solution for two-layer NNs (Mei et al., 2018; Chizat and Bach, 2018; Fang et al., 2019). Compared with lazy training, the mean-field view can characterize the entire training process of NNs.

However, the mean-field analysis on DNNs is a challenging task. First of all, it is not easy to formulate the mean-field limit of DNNs. As we will discuss in Section 2.1, extending existing formulations to DNNs, hidden units in a middle layer essentially behave as a single unit along the training. This degenerate situation arguably cannot fully characterize the training process of actual DNNs. Moreover, understandings for the global convergence of Gradient Descent on DNNs are still required in the mean-field regime.

In this paper, we propose a new mean-field framework for over-parameterized DNNs to analyze NN training. In contrast to existing studies focusing on the NN weights, this framework represents a DNN in the continuous limit by probability measures and functions over its features, that is, the outputs of the hidden units over the training data. This new representation overcomes the degenerate situation in previous studies (Araújo et al., 2019; Nguyen and Pham, 2020).

We further describe a non-linear dynamic called *neural feature flow* that captures the evolution of a DNN trained by Gradient Descent. We illustrate the framework by Res-Nets (He et al., 2016). Neural feature flow involves the evolution of the features and does not require the boundedness of the weights. Under the standard initialization method of discrete Res-Nets (Glorot and Bengio, 2010; He et al., 2015), the NN weights scale to infinity with the growth of the number of hidden units. There are empirical studies, e.g. Zhang et al. (2019), which show that properly rescaling the standard initialization stabilizes training. We introduce a simple  $\ell_2$ -regression at initialization (see Algorithm 2). We prove that Gradient Descent from the regularized initialization with a suitable time scale on Res-Nets can be well-approximated by its limit, i.e., neural feature flow, when the number of hidden units is sufficiently large.

Finally, we consider the global convergence of neural feature flow for Res-Nets. Surprisingly, we show that when the neural feature flow process converges, it reaches a *globally* optimal solution under suitable conditions. We summarize the contributions of the paper below:

- (A) We propose a new mean-field framework of DNNs which characterizes DNNs via probability measures and functions over the features and introduce neural feature flow to capture the evolution of DNNs trained by the Gradient Descent algorithm.
- (B) We illustrate our framework by Res-Net model. We show that neural feature flow can find a global minimal solution of the learning task under certain conditions.

Our mean-field description can also be used to study the standard DNNs, which is discussed in Appendix E. However, it still remains open to achieve the global convergence of neural feature flow for standard DNNs.

### 1.1. Notations

Let  $[m_1 : m_2] := \{m_1, m_1 + 1, \dots, m_2\}$  for  $m_1, m_2 \in \mathbb{N}$  with  $m_1 \leq m_2$  and  $[m_2] := [1 : m_2]$  for  $m_2 \geq 1$ . Let  $\mathcal{P}^n$  be the set of probability distributions over  $\mathbb{R}^n$ . For a matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$ , let  $\|\mathbf{A}\|_2$ ,  $\|\mathbf{A}\|_F$ , and  $\|\mathbf{A}\|_\infty$  denote its operator, Frobenius, max norms, respectively. If  $\mathbf{A}$  is symmetric, let  $\lambda_{\min}(\mathbf{A})$  be its smallest eigenvalue. Vectors are treated as columns. For a vector  $\mathbf{a} \in \mathbb{R}^n$ , let  $\|\mathbf{a}\|_2$  and  $\|\mathbf{a}\|_\infty$  denote its  $\ell_2$  and  $\ell_\infty$  norms, respectively. The  $i$ -th coordinate is denoted by  $\mathbf{a}(i)$ . For  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ , denote the entrywise product by  $\mathbf{a} \circ \mathbf{b}$  that  $[\mathbf{a} \circ \mathbf{b}](i) := \mathbf{a}(i) \circ \mathbf{b}(i)$  for  $i \in [n]$ . For  $c > 0$  and  $p \in [1, \infty]$ , let  $\mathcal{B}_p(\mathbf{a}, c)$  denote the  $\ell_p$ -ball centered at  $\mathbf{a}$  of radius  $c$ . For an unary function  $f : \mathbb{R} \rightarrow \mathbb{R}$ , define  $\hat{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  as the entrywise operation that  $\hat{f}(\mathbf{a})(i) = f(\mathbf{a}(i))$  for  $i \in [n]$  and  $\mathbf{a} \in \mathbb{R}^n$ . Denote  $n$ -dimensional identity matrix by  $\mathbf{I}^n$ . Denote  $m$ -by- $n$  zero matrix and  $n$ -dimensional zero vector by  $\mathbf{0}^{n \times m}$  and  $\mathbf{0}^n$ , respectively. We say a univariate distribution  $p$  is  $\sigma$ -sub-gaussian if  $\mathbb{E}_{x \sim p} \exp(x^2/\sigma^2) \leq e^1$ ; we say a  $d$ -dimensional distribution  $p$  is  $\sigma$ -sub-gaussian if the law of  $u^\top \mathbf{x}$  is  $\sigma$ -sub-gaussian for  $\mathbf{x} \sim p$  and any  $u \in \mathbb{S}^{d-1}$ . For two positive sequences  $\{p_n\}$  and  $\{q_n\}$ ,  $p_n = \mathcal{O}(q_n)$  if  $p_n \leq Cq_n$  for some positive constant  $C$ , and  $p_n = \Omega(q_n)$  if  $q_n = \mathcal{O}(p_n)$ .

## 2. Related Deep Learning Theory

In recent years, there have been a number of significant developments to obtain better theoretical understandings of NNs. One remarkable direction is to restrict the NN training in a small region. In this lazy training regime, the analysis cannot explain how NNs learn discriminative features. This is observed in real applications and argued to be one of the contributors to the success of deep learning. Beyond lazy training, one promising direction is to conduct mean-field analysis. However, in section 2.1, we show the challenges of analysis on DNNs. Specifically, if we still model from the weights, the standard initialization, e.g., (Glorot and Bengio, 2010; He et al., 2015) scales the weights to  $\sqrt{m}$ , which diverges in the mean-field limit, where  $m$  is the number of hidden units. On the other side, if we initialize the weights from a fixed distribution that is independent of  $m$  as existing mean-field works (Araújo et al., 2019; Nguyen and Pham, 2020) considered, DNNs would be stuck in a degenerate situation where the middle layers essentially only have one single feature. Both the issues motivate us to study DNN directly from the features. As a result, we propose a mean-field framework from tracking the distributions of features to analyze the DNN training.

### 2.1. Challenges on Mean-field Theory for DNNs

We discuss related mean-field studies and point out the challenges in modeling DNNs. For two-layer NNs, most of the existing works (Mei et al., 2018; Chizat and Bach, 2018; Sirignano and Spiliopoulos, 2019a; Rotskoff and Vanden-Eijnden, 2018) formulate the continuous limit as

$$f(\mathbf{x}; p) = \int w_2 h(\mathbf{w}_1^\top \mathbf{x}) dp(w_2, \mathbf{w}_1),$$

---

1. Here the value  $e$  can be replaced by any number greater than one. See Vershynin (2010, Remark 5.6).

where  $p$  is the probability distribution over the pair of weights  $(w_2, \mathbf{w}_1)$ . The weights of the second layer  $w_2$  can be viewed as functions of  $\mathbf{w}_1$ , which is a  $d$ -dimensional vector. However, following the approach, the higher-layer weights, say  $w_3$ , are functions over features of the hidden layer, with a diverging dimensionality in the mean-field limit. For 3-layer NNs,  $w_3$  as the last hidden layer is indexed by the connection to the output units in [Nguyen and Pham \(2020\)](#), which is not generalizable when middle layers present. An alternative approach is to model DNNs with nested measures (also known as multi-layer measures; see [Dawson et al. \(1982\)](#); [Dawson \(2018\)](#) and references therein), which however suffers the closure problem to establish a well-defined limit (see discussions in [Sirignano and Spiliopoulos \(2019b\)](#), Section 4.3).

The continuous limit of DNNs is investigated by [Araújo et al. \(2019\)](#); [Nguyen and Pham \(2020\)](#) under the initialization that all weights are i.i.d. realizations of a *fixed* distribution independent of the number of hidden units. However, under that setting, all neurons in a middle layer essentially behave as a single neuron. Consider the output  $\hat{\beta}$  of a middle-layer neuron connecting to  $m$  hidden neurons in the previous layer:

$$\hat{\beta} = \frac{1}{m} \sum_{i=1}^m h(\hat{\beta}'_i) w_i, \quad (1)$$

where  $\hat{\beta}'_i$  is the output of  $i$ -th hidden neuron in the previous layer with bounded variance,  $w_i$  is the connecting weight, and  $h$  is the activation function. If  $w_i$  is initialized independently from  $\mathcal{N}(0, 1)$ , it is clear that  $\text{var}[\hat{\beta}] \rightarrow 0$  as  $m \rightarrow \infty$ , and thus the hidden neurons in middle layers are indistinguishable at the initialization. Moreover, the phenomenon sustains along the entire training process, as shown in [Proposition 1](#). This phenomenon serves as the basis of [Araújo et al. \(2019\)](#); [Nguyen and Pham \(2020\)](#) to characterize the mean-field limit using finite-dimensional probability distributions. This degenerate situation arguably does not fully characterize the actual DNN training. In fact, similar calculations to (1) are carried out by [Glorot and Bengio \(2010\)](#); [He et al. \(2015\)](#) and motivate the popular initialization strategy with  $\mathcal{N}(0, \mathcal{O}(m))$  such that the variance of  $\hat{\beta}$  is non-vanishing.

**Proposition 1** *Consider fully-connected  $L$ -layer DNNs with  $m$  units in each hidden layer trained by Gradient Descent. Suppose the activation and loss functions satisfy Assumption 1. Let the weights be initialized from a distribution with  $\mathcal{O}(1)$  variance. Let  $\hat{\beta}_{\ell,i}^k$  denote the output of  $i$ -th hidden neuron at  $\ell$ -th layer and  $k$ -th iteration, and define  $\Delta_{\ell,m} := \max_{i \neq i', k \in [K]} \|\hat{\beta}_{\ell,i}^k - \hat{\beta}_{\ell,i'}^k\|_{\infty}$ . Then, for every  $\ell \in [2 : L - 1]$ , almost surely,*

$$\lim_{m \rightarrow \infty} \Delta_{\ell,m} = 0.$$

To overcome this degenerate situation, we consider the popular initialization strategy with a simple  $\ell_2$  regression (See [Algorithm 2](#)). The regression preserves all initial features, thus the variance of the features is now non-vanishing. Moreover, to accurately characterize DNNs in the mean-field regime, we introduce a probability measure over the features instead of the weights, which leads to a new dynamic system called neural feature flow to describe the neural nets trained by Gradient Descent. We achieve to describe a more realistic learning process.

## 2.2. Comparisons of Dynamics with Other Mean-field Works

It is known that the evolution of a two-level NN trained by Gradient Descent can be described as a McKean-Vlasov process ([Mei et al., 2018](#); [Chizat and Bach, 2018](#)). There are lots of works that

studies the evolution of DNN trained by Gradient Descent. One important work from Araújo et al. (2019) shows that the evolution of DNNs can also be characterized by PDEs of McKean-Vlasov type when weights of DNNs in the first and the last layers are not updated. More recently, Nguyen and Pham (2020); Pham and Nguyen (2020) proposed another very interesting attempt by directly tracking the trajectories of the weights. In their description called neuronal embedding, the evolution of Gradient Descent is characterized by systems of ODEs, avoiding the presence of the conditional probabilities which is the main issue proposed by Araújo et al. (2019). Our description follows from the idea of Nguyen and Pham (2020). However, our dynamic further involves the evolution of the features and does not require the boundedness of the weights. Moreover, we introduce the conception of skip-connected paths to deal with the Res-Net architecture. Neural feature flow is more meaningful training dynamics as it is no longer restricted in the degenerate situation as Section 2.1 described.

For the analysis to achieve a global minimum for the dynamics, the work from Lu et al. (2020) considers the DNN as a relatively simple composition of multiple two-layer NNs. Their global convergence result requires a very restricted assumption that the limiting distribution has full support. The work from Nguyen and Pham (2020); Pham and Nguyen (2020) proved the global convergence for DNNs by a very novel topology argument under the degenerate initialization. Pham and Nguyen (2020) also mentioned the possibility of non-degenerate initialization leading to global convergence guarantees. Our proof idea is similar to Nguyen and Pham (2020), whereas, we take our concentration on the features.

### 2.3. Beyond Lazy Training

In the “lazy training” regime, e.g. Jacot et al. (2018); Du et al. (2019a); Allen-Zhu and Li (2019); Zou et al. (2018), the weights are restricted in an infinitely small region. The DNN in this regime essentially corresponds to a linear model on random features associated with a kernel termed neural tangent kernel. In the limit, the features are fixed. Encouragingly, we consider NNs beyond the lazy training regime and further allow the feature to move in a constant region. To arrive the goal, we study a special Res-Net architecture that bounds residuals by a bounded mapping  $h_2$  (see Section 3.1). Note that in our analysis, the bound is not needed to be small enough but can be arbitrarily large and less than infinity. Therefore, our setting allows the DNN to learn the targeted features. From the technical aspect, we require a different treatment to show the full rank of the feature matrices; this is achieved by Brouwer’s fixed-point theorem (see the Proof of Theorem 8 in Appendix B.2).

## 3. Formulation of Continuous Res-Nets

We consider the empirical minimization problem over  $N$  training samples  $\{\mathbf{x}^i, y^i\}_{i=1}^N$ , where  $\mathbf{x}^i \in \mathbb{R}^d$  and  $y^i \in \mathcal{Y}$ . For regression problems,  $\mathcal{Y}$  is typically  $\mathbb{R}$ ; for classification problems,  $\mathcal{Y}$  is often  $[K]$  for an integer  $K$ . We first present the formulation of  $L$ -layer Res-Nets.

### 3.1. Discrete Res-Nets

For discrete Res-Nets, let  $m_\ell$  denote the number of units at layer  $\ell$  for  $\ell = [0 : L + 1]$ . Suppose each hidden layer has  $m$  hidden units that  $m_\ell = m$  for  $\ell \in [L]$ . Let  $m_0 = d$  and node  $i$  outputs the value of  $i$ -th coordinate of the training data for  $i \in [d]$ . Let  $m_{L+1} = 1$  that is the unit of the final network output. For  $\ell \in [L + 1]$ , the output of node  $i$  for the  $N$  training samples in layer  $\ell$

is denoted by  $\hat{\beta}_{\ell,i} \in \mathbb{R}^N$ ; the weight that connects the node  $i$  at layer  $\ell - 1$  to node  $j$  at layer  $\ell$  is denoted by  $\hat{v}_{\ell,i,j} \in \mathbb{R}$ .

(1) At the input layer, for  $i \in [d]$ , let  $\hat{\beta}_{0,i} := [\mathbf{x}^1(i), \mathbf{x}^2(i), \dots, \mathbf{x}^N(i)]^\top$ .

(2) At the first layer, for  $j \in [m]$ , let

$$\hat{\beta}_{1,j} = \frac{1}{m_0} \sum_{i=1}^{m_0} \hat{v}_{1,i,j} \hat{\beta}_{0,i}. \quad (2)$$

(3) We recursively define the upper layers for  $\ell \in [2 : L]$ . Let  $\hat{\alpha}_{\ell,j} \in \mathbb{R}^N$  be the residual term at node  $j$  at layer  $\ell$ :

$$\hat{\alpha}_{\ell,j} = \frac{1}{m} \sum_{i=1}^m \hat{v}_{\ell,i,j} \dot{h}_1 \left( \hat{\beta}_{\ell-1,i} \right), \quad j \in [m], \quad (3)$$

where  $h_1 : \mathbb{R} \rightarrow \mathbb{R}$  is the activation function and  $\dot{h}_1 : \mathbb{R}^N \rightarrow \mathbb{R}^N$  is the entrywise operation for  $h_1$ , which satisfies  $\dot{h}_1(\mathbf{a})(i) = h_1(\mathbf{a}(i))$  for  $i \in [N]$  and  $\mathbf{a} \in \mathbb{R}^N$ . Furthermore, we consider the following coupling between the residual and the previous feature:

$$\hat{\beta}_{\ell,j} = \dot{h}_2(\hat{\alpha}_{\ell,j}) + \hat{\beta}_{\ell-1,j}, \quad j \in [m]. \quad (4)$$

where  $h_2 : \mathbb{R} \rightarrow \mathbb{R}$ .

(4) For the output,

$$\hat{\beta}_{L+1,1} = \frac{1}{m} \sum_{i=1}^m \hat{v}_{L+1,i,1} \dot{h}_1 \left( \hat{\beta}_{L,i} \right). \quad (5)$$

We collect the weights from all layers into a single vector  $\hat{v} \in \mathbb{R}^{D_1}$  with  $D_1 := m^2(L - 1) + (d + 1)m$ . We also collect the residuals, and features from layers 2 to  $L$  into single vectors  $\hat{\alpha} \in \mathbb{R}^{D_2}$ , and  $\hat{\beta} \in \mathbb{R}^{D_2}$ , respectively, where  $D_2 := Nm(L - 1)$ . The learning problem for Res-Nets is given by

$$\min_{\hat{v}, \hat{\alpha}, \hat{\beta}} \hat{\mathcal{L}}(\hat{v}, \hat{\alpha}, \hat{\beta}) = \frac{1}{N} \sum_{n=1}^N \phi \left( \hat{\beta}_{L+1,1}(n), y^n \right), \quad (6)$$

where  $(\hat{v}, \hat{\alpha}, \hat{\beta})$  satisfies (2) – (5), and  $\phi : \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}$  denotes the loss function. One noteworthy feature in the architecture is (4), where we introduce a mapping  $h_2$  on the residual  $\hat{\alpha}_{\ell,j}$  before fusing it with  $\hat{\beta}_{\ell-1,j}$ . As have been mentioned, we assume that  $h_2$  is bounded by a constant  $L_1$ , and hence  $\|\hat{\beta}_{\ell,j} - \hat{\beta}_{\ell-1,j}\|_\infty \leq L_1$ . Therefore, the high-layer features can be regarded as perturbations of the low-layer ones. Similar ideas have also appeared in [Du et al. \(2019a\)](#); [Hardt and Ma \(2016\)](#), but realized in a different way. For example, in the lazy training regime, [Du et al. \(2019a\)](#) achieved it by scaling  $\hat{\alpha}_{\ell,j}$  with a vanishing  $\mathcal{O}(\frac{1}{\sqrt{m}})$  factor.

### 3.2. Continuous Res-Nets Formulation

We propose our formulation for the continuous Res-Nets. We consider the Res-Net with the architecture described in Section 3.1 that is initialized by a combination of a standard initialization and an additional regression procedure (See Algorithm 2). One can find that the regression procedure preserves all initial features but reduces the redundancy of the weights, making us introduce real-value functions to characterize the weights in the mean-field limit.

In fact, displayed in Appendix E.2, the continuous limit of a standard DNN can be described by introducing probability measures over features for individual layers. This idea is much clear. Here, we study Res-Net with the skip connections as in (4). In the continuous case, the discrete index  $j$  no longer makes sense and the skip connections need to be properly parametrized by an infinite set. To overcome the hurdle of infinite skip connections, we introduce  $\Theta = (v_1, \alpha_2, \dots, \alpha_L) \in \mathbb{R}^D$  for  $D = d + (N - 1)L$  to parametrize the skip connections that are described in (8) and (9) below. Each  $\Theta$  consists of  $v_1, \alpha_2, \dots, \alpha_L$  that can be regarded as an input-output path  $v_1 \rightarrow \alpha_2 \rightarrow \dots \rightarrow \alpha_L$  and is called a *skip-connected path*. Our main technique is to characterize the overall state of the continuous Res-Nets by the density  $p$  over skip-connected paths. Thus the joint distribution  $p$  can be regarded as a description of the overall topological structure about the skip connections. We represent the features  $\beta_\ell$  in the hidden layer  $\ell \in [L]$  as functions of  $\Theta$  that we introduce next:

(1) At the input layer, let  $\mathbf{X} = [\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N]^\top \in \mathbb{R}^{N \times d}$ .

(2) At the first layer, let the features be

$$\beta_1(\Theta) = \frac{1}{d}(\mathbf{X}v_1). \quad (7)$$

(3) At layer  $\ell \in [2 : L]$ , let  $v_\ell : \text{supp}(p) \times \text{supp}(p) \rightarrow \mathbb{R}$  denote the weights on the connections from layer  $\ell - 1$  to  $\ell$ , then for all  $\Theta = (v_1, \alpha_1, \alpha_2, \dots, \alpha_L) \in \mathbb{R}^D$ , we have the forward-propagation constraint for  $v_\ell$  and  $p$ :

$$\alpha_\ell = \int v_\ell(\Theta, \bar{\Theta}) \dot{h}_1(\beta_{\ell-1}(\bar{\Theta})) dp(\bar{\Theta}), \quad (8)$$

$$\beta_\ell(\Theta) = \dot{h}_2(\alpha_\ell) + \beta_{\ell-1}(\Theta). \quad (9)$$

Here,  $\Theta$  takes on values in  $\mathbb{R}^D$  and for each  $\Theta$ ,  $\alpha_\ell$  is one part of  $\Theta$  and  $\beta_\ell$  is a function of  $\Theta$ .  $\alpha_\ell$  represents the residual at layer  $\ell$  on the skip connected path described by  $\Theta$ . And  $\beta_\ell(\Theta)$  represents the corresponding feature.

(4) At the output layer, let  $v_{L+1} : \text{supp}(p) \rightarrow \mathbb{R}$  be the weights in the layer  $L + 1$ , and we have

$$\beta_{L+1} = \int v_{L+1}(\Theta) \dot{h}_1(\beta_L(\Theta)) dp(\Theta). \quad (10)$$

In our continuous formulation, a static Res-Net is characterized by  $p$  and  $v_2 \dots v_{L+1}$ . We will show in the next section that the continuous Res-Net  $(\{v_\ell\}_{\ell=2}^{L+1}, p)$  that satisfies the (7) – (10) will serve as a feasible initialization for neural feature flow.

---

**Algorithm 1** Scaled Gradient Descent for Training a Res-Net.

---

- 1: Input the data  $\{\mathbf{x}^i, y^i\}_{i=1}^N$ , step size  $\eta$ , and initial weights  $\hat{\mathbf{v}}^0$ .
- 2: **for**  $k = 0, 1, \dots, K - 1$  **do**
- 3:   Perform forward-propagation (2) – (5) to compute  $\hat{\beta}_{L+1,1}^k$ .
- 4:   Perform backward-propagation to compute the gradient  $\hat{\mathcal{G}}_{\ell,i,j}^k = \frac{\partial \hat{\mathcal{L}}}{\partial \hat{v}_{\ell,i,j}^k}$ .
- 5:   Perform Scaled Gradient Descent:

$$\hat{v}_{\ell,i,j}^{k+1} = \hat{v}_{\ell,i,j}^k - [\eta m_{\ell-1} m_{\ell}] \hat{\mathcal{G}}_{\ell,i,j}^k, \quad \ell \in [L+1], i \in [m_{\ell-1}], j \in [m_{\ell}].$$

- 6: **end for**
  - 7: Output the weights  $\hat{\mathbf{v}}^K$ .
- 

#### 4. Scaled Gradient Descent and Neural Feature Flow for Res-Nets

We focus on the dynamic of the Res-Net trained by Gradient Descent. We consider the scaled Gradient Descent algorithm<sup>2</sup>. Given initial weights  $\hat{\mathbf{v}}^0$ , the meta-algorithm of scaled Gradient Descent is presented in Algorithm 1. Note that Algorithm 1 differs from the standard Gradient Descent only on the step sizes (time scales). Such scaling is also adopted in existing works (Araújo et al., 2019; Nguyen and Pham, 2020).

Now we describe the continuous limit of the Res-Net trained by Algorithm 1 by the continuous trajectories of the Res-Nets. This idea follows from Nguyen and Pham (2020) for analyzing there-layer DNNs. A trajectory is denoted by  $\Phi$  that maps the initial Res-Nets at  $t = 0$  to a Res-Net process over  $[0, T]$ . Specifically, it consists of the following parts:

- $\Phi_{\ell}^{\beta} : \text{supp}(p) \rightarrow C([0, T], \mathbb{R}^N)$  is the trajectory of  $\beta_{\ell}$  for  $\ell \in [L]$ ;
- $\Phi_{\ell}^{\alpha} : \text{supp}(p) \rightarrow C([0, T], \mathbb{R}^N)$  is the trajectory of  $\alpha_{\ell}$  for  $\ell \in [2 : L]$ ;
- $\Phi_1^v : \text{supp}(p) \rightarrow C([0, T], \mathbb{R}^d)$  and  $\Phi_{L+1}^v : \text{supp}(p) \rightarrow C([0, T], \mathbb{R})$  are the trajectories of  $v_1$  and  $v_{L+1}$ , respectively;
- $\Phi_{\ell}^v : \text{supp}(p) \times \text{supp}(p) \rightarrow C([0, T], \mathbb{R})$  is the trajectory of  $v_{\ell}$  for  $\ell \in [2 : L]$ .

The continuous gradient for the weight can be obtained from the backward-propagation algorithm. For a given trajectory  $\Phi$ , the gradients of weights at time  $t \in [0, T]$  can be obtained from the backward-propagation algorithm. Similar to the usual backward-propagation, we first define gradients with respect to the features and residuals. Specifically, for all  $\Theta = (v_1, \alpha_2, \dots, \alpha_L) \in$

---

2. In practice, one often use stochastic gradient instead of the full counterpart for training. Under mild conditions, the dynamic of scaled Stochastic Gradient Descent will also converge to the neural feature flow in the continuous limit.



$\text{supp}(p)$ ,  $t \in [0, T]$ , and  $\ell \in [2 : L]$ , let

$$\beta_{L+1}(\Phi, t) := \int \Phi_{L+1}^v(\Theta)(t) \dot{h}_1(\Phi_L^\beta(\Theta)(t)) dp(\Theta), \quad (11a)$$

$$\begin{aligned} \underline{\mathcal{D}}_{L+1}(\Phi, t) &:= [\phi_1'(\beta_{L+1}(\Phi, t)(1), y^1), \dots, \phi_1'(\beta_{L+1}(\Phi, t)(N), y^N)]^\top, \\ \underline{\mathcal{D}}_L^\beta(\Theta; \Phi, t) &:= [\Phi_{L+1}^v(\Theta)(t) \underline{\mathcal{D}}_{L+1}(\Phi, t)] \circ \dot{h}_1'(\Phi_L^\beta(\Theta)(t)), \end{aligned} \quad (11b)$$

$$\underline{\mathcal{D}}_\ell^\alpha(\Theta; \Phi, t) := \underline{\mathcal{D}}_\ell^\beta(\Theta; \Phi, t) \circ \dot{h}_2'(\Phi_\ell^\alpha(\Theta)(t)), \quad (11c)$$

$$\underline{\mathcal{D}}_{\ell-1}^\beta(\Theta; \Phi, t) := \underline{\mathcal{D}}_\ell^\beta(\Theta; \Phi, t) + \left[ \int \Phi_\ell^v(\Theta, \bar{\Theta})(t) \underline{\mathcal{D}}_\ell^\alpha(\bar{\Theta}; \Phi, t) dp(\bar{\Theta}) \right] \circ \dot{h}_1'(\Phi_{\ell-1}^\beta(\Theta)(t)).$$

For all  $\Theta, \bar{\Theta} \in \text{supp}(p)$ , the drift term for the weights is given by

$$\underline{\mathcal{G}}_{L+1}^v(\Theta; \Phi, t) := \frac{1}{N} [\underline{\mathcal{D}}_{L+1}(\Phi, t)]^\top \dot{h}_1(\Phi_L^\beta(\Theta)(t)), \quad (12a)$$

$$\underline{\mathcal{G}}_\ell^v(\Theta, \bar{\Theta}; \Phi, t) := \frac{1}{N} [\underline{\mathcal{D}}_\ell^\alpha(\bar{\Theta}; \Phi, t)]^\top \dot{h}_1(\Phi_{\ell-1}^\beta(\Theta)(t)), \quad \ell \in [2 : L], \quad (12b)$$

$$\underline{\mathcal{G}}_1^v(\Theta; \Phi, t) := \frac{1}{N} \mathbf{X} \underline{\mathcal{D}}_1^\beta(\Theta; \Phi, t). \quad (12c)$$

Moreover, the changes of the weights will induce a change of the residuals and features. By the chain rule, we can obtain the drift term for the residuals and features: for  $\ell \in [L-1]$  and  $\Theta \in \text{supp}(p)$ ,

$$\underline{\mathcal{G}}_1^\beta(\Theta; \Phi, t) := \frac{1}{d} [\mathbf{X} \underline{\mathcal{G}}_1^v(\Theta; \Phi, t)], \quad (13a)$$

$$\begin{aligned} \underline{\mathcal{G}}_{\ell+1}^\alpha(\Theta; \Phi, t) &:= \int \Phi_{\ell+1}^v(\bar{\Theta}, \Theta)(t) \left[ \dot{h}_1'(\Phi_\ell^\beta(\bar{\Theta})(t)) \circ \underline{\mathcal{G}}_\ell^\beta(\bar{\Theta}; \Phi, t) \right] dp(\bar{\Theta}) + \\ &\quad + \int \dot{h}_1(\Phi_\ell^\beta(\bar{\Theta})(t)) \circ \underline{\mathcal{G}}_{\ell+1}^v(\bar{\Theta}, \Theta; \Phi, t) dp(\bar{\Theta}), \end{aligned} \quad (13b)$$

$$\underline{\mathcal{G}}_{\ell+1}^\beta(\Theta; \Phi, t) := \underline{\mathcal{G}}_\ell^\beta(\Theta; \Phi, t) + \underline{\mathcal{G}}_{\ell+1}^\alpha(\Theta; \Phi, t) \circ \dot{h}_2'(\Phi_{\ell+1}^\alpha(\Theta)(t)). \quad (13c)$$

In all, the process of a continuous Res-Net trained by scaled Gradient Descent can be defined below.

**Definition 2 (Neural Feature Flow for Res-Net)** *Given an initial continuous Res-Net  $(\{v_\ell\}_{\ell=2}^{L+1}, p)$  that satisfies the (7) – (10) and  $T < \infty$ , we say a trajectory  $\Phi_*$  is a neural feature flow if for all  $\Theta = (v_1, \alpha_2, \dots, \alpha_L) \in \text{supp}(p)$ ,  $\bar{\Theta} \in \text{supp}(p)$ , and  $t \in [0, T]$ ,*

$$\Phi_{*,\ell}^\beta(\Theta)(t) = \left[ \frac{1}{d} \mathbf{X} v_1 + \sum_{i=2}^{\ell} \dot{h}_2(\alpha_i) \right] - \int_0^t \underline{\mathcal{G}}_\ell^\beta(\Theta; \Phi_*, s) ds, \quad \ell \in [L],$$

$$\Phi_{*,\ell}^\alpha(\Theta)(t) = \alpha_\ell - \int_0^t \underline{\mathcal{G}}_\ell^\alpha(\Theta; \Phi_*, s) ds, \quad \ell \in [2 : L],$$

$$\Phi_{*,1}^v(\Theta)(t) = v_1 - \int_0^t \underline{\mathcal{G}}_1^v(\Theta; \Phi_*, s) ds,$$

$$\Phi_{*,\ell}^v(\Theta, \bar{\Theta})(t) = v_\ell(\Theta, \bar{\Theta}) - \int_0^t \underline{\mathcal{G}}_\ell^v(\Theta, \bar{\Theta}; \Phi_*, s) ds, \quad \ell \in [2 : L],$$

$$\Phi_{*,L+1}^v(\Theta)(t) = v_{L+1}(\Theta) - \int_0^t \underline{\mathcal{G}}_{L+1}^v(\Theta; \Phi_*, s) ds.$$

We call the process as neural feature flow because it characterizes the evolution of both weights and features.

## 5. Main Results

### 5.1. Assumptions

We first present the assumptions that are needed in our analysis.

**Assumption 1 (Activation Functions and Loss Function)** *For the activation functions, we assume that there exist constants  $L_1, L_2, L_3 > 0$  such that, for all  $x \in \mathbb{R}$ ,*

$$|h_1(x)| \leq L_1, \quad |h_2(x)| \leq L_1, \quad |h'_1(x)| \leq L_2, \quad |h'_2(x)| \leq L_2.$$

Moreover, for all  $x, y \in \mathbb{R}$ ,

$$|h'_1(x) - h'_1(y)| \leq L_3|x - y|, \quad |h'_2(x) - h'_2(y)| \leq L_3|x - y|.$$

For the loss function, we assume that there exist constants  $L_4, L_5 > 0$  such that, for all  $y \in \mathcal{Y}$ ,  $x_1 \in \mathbb{R}$ , and  $x_2 \in \mathbb{R}$ ,

$$|\phi'_1(x_1, y)| \leq L_4, \quad |\phi'_1(x_1, y) - \phi'_1(x_2, y)| \leq L_5|x_1 - x_2|.$$

Assumption 1 is easy to be satisfied. It only requires some boundedness, continuity, and smoothness for the activation and loss functions. It is adopted in most mean-field analysis, such as [Mei et al. \(2018\)](#); [Araújo et al. \(2019\)](#).

**Assumption 2 (Strong Universal Approximation Property)** *Assume that for any function  $f_2 : \mathbb{R}^d \rightarrow \mathbb{R}^N$  that is bounded by  $C_B$ , i.e., for all  $\mathbf{v}_1 \in \mathbb{R}^d$ ,  $\|f_2(\mathbf{v}_1)\|_\infty \leq C_B$ , we have*

$$\lambda_{\min} \left[ \int \left[ \dot{h}_1 \left( \frac{1}{d} \mathbf{X} \mathbf{v}_1 + f_2(\mathbf{v}_1) \right) \right] \left[ \dot{h}_1 \left( \frac{1}{d} \mathbf{X} \mathbf{v}_1 + f_2(\mathbf{v}_1) \right) \right]^\top dp_1(\mathbf{v}_1) \right] \geq \bar{\lambda} > 0. \quad (14)$$

where  $\bar{\lambda}$  only depends on  $\mathbf{X}$ ,  $C_B$ , and  $h_1$ , and  $p_1 = \mathcal{N}(\mathbf{0}^d, \mathbf{I}^d)$ .

Assumption 2 is a technical assumption that we conjecture to hold under fairly general conditions. Notably when  $C_B = 0$ , it is shown in [Du et al. \(2019a\)](#), Lemma F.1) that the assumption holds for all analytic non-polynomial  $h_1$ . Lemma 3 affords many examples that satisfy the assumption for constant  $C_B$ .

**Lemma 3** *Suppose that the data is non-parallel, i.e.,  $\mathbf{x}_i \notin \text{Span}(\mathbf{x}_j)$  for all  $i \neq j$ .*

- (i) *If  $g : \mathbb{R} \rightarrow \mathbb{R}$  is a non-polynomial function that is bounded and has Lipschitz continuous gradient, then  $h_1(x) := g(cx)$  satisfies Assumption 2 when  $c > 0$  is sufficiently small.*
- (ii) *The Relu-type function  $h_1(x) = (x)_+^\alpha$  for  $\alpha > 0$  satisfies Assumption 2.*
- (iii) *If  $h_1(x) = c|x|^{-\alpha}$  or  $h_1(x) = c(x)_+^{-\alpha}$  for  $|x| > c'$ , where  $c, c', \alpha > 0$ , then  $h_1$  satisfies Assumption 2.*

The condition in Lemma 3 is standard and widely used in the analysis of lazy training (see Du et al. (2019a, Lemma F.1)). It only requires the data are not parallel: for every  $j \neq i$ ,  $\mathbf{x}_j \neq c\mathbf{x}_i$  for any scalar  $c$ . In the following, we propose the conditions for the initial continuous Res-Net  $(\{v_\ell\}_{\ell=2}^{L+1}, p)$ . In the next section, we will show concrete examples that realize these assumptions.

**Assumption 3 (Initialization for Continuous Res-Net)** *We first assume that the initial continuous Res-Net  $(\{v_\ell\}_{\ell=2}^{L+1}, p)$  is a feasible continuous Res-Net that satisfies the forward propagation constraints, i.e., (7) – (10). Moreover,  $p$  is  $\sigma$ -sub-gaussian distribution and has a full support<sup>3</sup>. For all  $\ell \in [2 : L]$ ,  $v_\ell(\cdot, \cdot)$  has sublinear growth on the second argument, that is, there is a constant  $C_1$  such that*

$$|v_\ell(\Theta, \bar{\Theta})| \leq C_1 (1 + \|\bar{\Theta}\|_\infty), \quad \text{for all } \Theta, \bar{\Theta} \in \text{supp}(p), \ell \in [2 : L]. \quad (15)$$

Besides,  $v_\ell(\cdot, \cdot)$  are locally Lipschitz continuous where the Lipschitz constant has sub-linear growth on the second argument. In detail, there is a constant  $C_2$ , such that for  $\Theta_1 \in \text{supp}(p)$ ,  $\tilde{\Theta}_1 \in \text{supp}(p) \cap \mathcal{B}_\infty(\Theta_1, 1)$ ,  $\Theta_2 \in \text{supp}(p)$ , and  $\tilde{\Theta}_2 \in \text{supp}(p) \cap \mathcal{B}_\infty(\Theta_2, 1)$ , we have

$$|v_\ell(\Theta_1, \Theta_2) - v_\ell(\tilde{\Theta}_1, \tilde{\Theta}_2)| \leq C_2 (1 + \|\Theta_2\|_\infty) \left( \|\Theta_1 - \tilde{\Theta}_1\|_\infty + \|\Theta_2 - \tilde{\Theta}_2\|_\infty \right).$$

Finally, for the last layer, there exist constants  $C_3$  and  $C_4$ , such that for all  $\Theta, \bar{\Theta} \in \text{supp}(p)$ , we have

$$|v_{L+1}(\Theta)| \leq C_3 \quad \text{and} \quad |v_{L+1}(\Theta) - v_{L+1}(\bar{\Theta})| \leq C_4 \|\Theta - \bar{\Theta}\|_\infty.$$

## 5.2. Property and Approximation of Neural Feature Flow

We analyze the neural feature flow for the continuous Res-Net under the initial conditions in Assumption 3. The following theorem guarantees the existence and uniqueness.

**Theorem 4 (Existence and Uniqueness of Neural Feature Flow on Res-Net)** *Under Assumptions 1 and 3, for any  $T < \infty$ , there exists a unique neural feature flow  $\Phi_*$ .*

In fact, we also have that  $\Phi_*$  is a continuous mapping on  $\Theta$  given time  $t$  (see Theorem 13 in Appendix C.1). A similar continuity argument has also been observed by Nguyen and Pham (2020). The proofs of Theorems 4 and 13 can be obtained by the technique of Picard iterations (see, e.g., Hartman (1964)) with a special consideration on the search space to deal with the unboundedness of parameters. The latter differs from the former by introducing a more restrictive space in which all the candidates satisfy the desired property. We defer the proofs of this paper to Appendix.

Now we consider the approximation between a discrete DNN trained by scaled Gradient Descent and a continuous one evolving as neural feature flow. Based on the initial condition for the continuous Res-Net, we introduce the initial condition for discrete Res-Net.

**Definition 5 ( $\varepsilon_1$ -independent Initial Discrete Res-Net)** *We say an initial discrete Res-Net  $(\hat{\mathbf{v}}, \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}})$  is  $\varepsilon_1$ -independent if there exists a continuous initial Res-Net  $(\{v_\ell\}_{\ell=2}^{L+1}, p)$  satisfying Assumption 3 and  $(\bar{\mathbf{v}}, \bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}})$  such that*

3. The assumption that  $p$  has a full support will only be used in Theorem 8. It can be replaced by a slightly weaker assumption that there exists a continuous function  $f_1 : \mathbb{R}^d \rightarrow \mathbb{R}^{D-d}$  such that  $\text{supp}(p) \supseteq \{(\mathbf{v}_1, f_1(\mathbf{v}_1)) : \mathbf{v}_1 \in \mathbb{R}^d\}$ . When  $p_1$  has a full support,  $f_1$  can be simply chosen as  $f_1(\mathbf{v}_1) \equiv \mathbf{0}^{D-d}$ .

$$(1) \bar{\Theta}_i = (\bar{v}_{1,i}, \bar{\alpha}_{2,i}, \dots, \bar{\alpha}_{L,i}) \stackrel{i.i.d.}{\sim} p;$$

(2) For  $\bar{\beta}$  and  $\bar{v}$ ,

- $\bar{\beta}_{\ell,i} = \frac{1}{d} (\mathbf{X} \bar{v}_{1,i}) + \sum_{\ell_1=2}^{\ell} \dot{h}_2(\bar{\alpha}_{\ell_1,i})$  for  $\ell \in [L]$  and  $i \in [m]$ ;
- $\bar{v}_{\ell,i,j} = v_{\ell}(\bar{\Theta}_i, \bar{\Theta}_j)$  for  $\ell \in [2:L]$ ,  $i, j \in [m]$ ;
- $\bar{v}_{L+1,i,1} = v_{L+1}(\bar{\Theta}_i)$  for  $i \in [m]$ ;

(3)  $\varepsilon_1$ -closeness:

- $\|\bar{v}_{1,i} - \hat{v}_{1,i}\|_{\infty} \leq (1 + \|\bar{\Theta}_i\|_{\infty}) \varepsilon_1$  for  $i \in [m]$ ;
- $|\bar{v}_{\ell+1,i,j} - \hat{v}_{\ell+1,i,j}| \leq (1 + \|\bar{\Theta}_i\|_{\infty} + \|\bar{\Theta}_j\|_{\infty}) \varepsilon_1$  for  $\ell \in [L-1]$ ,  $i, j \in [m]$ ;
- $|\bar{v}_{L+1,i,1} - \hat{v}_{L+1,i,1}| \leq (1 + \|\bar{\Theta}_i\|_{\infty}) \varepsilon_1$  for  $i \in [m]$ .

As an  $\varepsilon_1$ -independent initialization relates to a continuous Res-Net satisfying Assumption 3, which yields an unique neural feature flow  $\Phi_*$  by Theorem 4, we show that scaled Gradient Descent from an  $\varepsilon_1$ -independent initialization is well-approximated by the corresponding neural feature flow when the number of hidden units is  $\tilde{\Omega}(\varepsilon_1^{-2})$ , where  $\tilde{\Omega}$  hides poly-logarithmic factors. This resembles a ‘‘propagation of chaos’’ argument from a Kac’s chaotic initial system (Sznitman, 1991). We compare the scaled Gradient Descent with an ideal discrete process determined by  $\Phi_*$  as specified below:

- Actual process  $(\hat{v}^{[0:K]}, \hat{\alpha}^{[0:K]}, \hat{\beta}^{[0:K]})$  by executing Algorithm 1 in  $K = \frac{T}{\eta}$  steps on the discrete Res-Net from  $(\hat{v}, \hat{\alpha}, \hat{\beta})$ ;
- Ideal process  $(\bar{v}^{[0,T]}, \bar{\alpha}^{[0,T]}, \bar{\beta}^{[0,T]})$  that evolves as neural feature flow:

$$\begin{aligned} \bar{\beta}_{\ell,i}^t &= \Phi_{*,\ell}^{\beta}(\bar{\Theta}_i)(t), \quad \ell \in [L], i \in [m], t \in [0, T], \\ \bar{\alpha}_{\ell,i}^t &= \Phi_{*,\ell}^{\alpha}(\bar{\Theta}_i)(t), \quad \ell \in [2:L], i \in [m], t \in [0, T], \\ \bar{v}_{1,i}^t &= \Phi_{*,1}^v(\bar{\Theta}_i)(t), \quad i \in [m], t \in [0, T], \\ \bar{v}_{\ell,i,j}^t &= \Phi_{*,\ell}^v(\bar{\Theta}_i, \bar{\Theta}_j)(t), \quad \ell \in [2:L], i \in [m], j \in [m], t \in [0, T], \\ \bar{v}_{L+1,i,1}^t &= \Phi_{*,L+1}^v(\bar{\Theta}_i)(t), \quad i \in [m], t \in [0, T]. \end{aligned}$$

We also compare the discrete and the continuous losses denoted by  $\hat{\mathcal{L}}^k := \frac{1}{n} \sum_{n=1}^N \phi(\hat{\beta}_{L+1,1}^k(n), y^n)$  and  $\mathcal{L}^t := \frac{1}{N} \sum_{n=1}^N \phi(\beta_{L+1}(\Phi_*, t)(n), y^n)$ , respectively. We have the theorem below.

**Theorem 6** *Under Assumption 1, suppose  $\varepsilon_1 \leq \mathcal{O}(1)$  and  $m \geq \tilde{\Omega}(\varepsilon_1^{-2})$ , and treat the parameters in assumptions and  $T$  as constants. Consider the actual process from an  $\varepsilon_1$ -independent initialization in Definition 5 with step size  $\eta \leq \tilde{\mathcal{O}}(\varepsilon_1)$ . Then, the following holds with probability  $1 - \delta$ :*

(1) *The two processes are close to each other:*

$$\begin{aligned} &\sup_{k \in [0:K]} \left\{ \sup_{i \in [m]} \left\| \hat{v}_{1,i}^k - \bar{v}_{1,i}^{k\eta} \right\|_{\infty}, \quad \sup_{\ell \in [2:L], i, j \in [m]} \left| \hat{v}_{\ell,i,j}^k - \bar{v}_{\ell,i,j}^{k\eta} \right| \right\} \leq \tilde{\mathcal{O}}(\varepsilon_1), \\ &\sup_{k \in [0:K], i \in [m]} \left\{ \left| \hat{v}_{L+1,i,1}^k - \bar{v}_{L+1,i,1}^{k\eta} \right|, \quad \sup_{\ell \in [2:L]} \left\| \hat{\alpha}_{\ell,i}^k - \bar{\alpha}_{\ell,i}^{k\eta} \right\|_{\infty}, \quad \sup_{\ell \in [L]} \left\| \hat{\beta}_{\ell,i}^k - \bar{\beta}_{\ell,i}^{k\eta} \right\|_{\infty} \right\} \leq \tilde{\mathcal{O}}(\varepsilon_1). \end{aligned}$$

---

**Algorithm 2** Initialize a Discrete Res-Net.
 

---

- 1: Input the data  $\mathbf{X}$ , variance  $\sigma_1$ , and a constant  $C_3$ .
  - 2: Independently draw  $\hat{v}_{1,i,j} \sim p_0 = \mathcal{N}(0, d\sigma_1^2)$  for  $i \in [d]$  and  $j \in [m]$ .
  - 3: Set  $\hat{\beta}_{1,j} = \frac{1}{d} \sum_{i=1}^d \hat{v}_{1,i,j} \hat{\beta}_{0,i}$  where  $j \in [m]$ . ◇ Standard Initialization for layer 1
  - 4: **for**  $\ell = 2, \dots, L$  **do**
  - 5:   Independently draw  $\tilde{v}_{\ell,i,j} \sim \mathcal{N}(0, m\sigma_1^2)$  for  $i, j \in [m]$ .
  - 6:   Set  $\hat{\alpha}_{\ell,j} = \frac{1}{m} \sum_{i=1}^m \tilde{v}_{\ell,i,j} \dot{h}_1(\hat{\beta}_{\ell-1,i})$  where  $j \in [m]$ .
  - 7:   Set  $\hat{\beta}_{\ell,j} = \hat{\beta}_{\ell-1,j} + \dot{h}_2(\hat{\alpha}_{\ell,j})$  for  $j \in [m]$ . ◇ Standard Initialization for layer  $\ell$
  - 8: **end for**
  - 9: Set  $\hat{v}_{L+1,i,1} = C_3$  where  $i \in [m]$ . ◇ Simply initialize  $\{\hat{v}_{L+1,i,1}\}_{i=1}^m$  by a constant
  - 10: **for**  $\ell = 2, \dots, L$  **do**
  - 11:   **for**  $j = 1, \dots, m$  **do**
  - 12:     Solve convex optimization problem: ◇ Perform  $\ell_2$ -regression to reduce redundancy
- $$\min_{\{\hat{v}_{\ell,i,j}\}_{i=1}^m} \frac{1}{m} \sum_{i=1}^m (\hat{v}_{\ell,i,j})^2, \quad \text{s.t. } \hat{\alpha}_{\ell,j} = \frac{1}{m} \sum_{i=1}^m \hat{v}_{\ell,i,j} \dot{h}_1(\hat{\beta}_{\ell-1,i}). \quad (16)$$
- 13:   **end for**
  - 14: **end for**
  - 15: Output the discrete Res-Net parameters  $(\hat{v}, \hat{\alpha}, \hat{\beta})$ .
- 

(2) *The training losses are also close to each other:*

$$\sup_{k \in [0:K]} \left| \hat{\mathcal{L}}^k - \mathcal{L}^{k\eta} \right| \leq \tilde{\mathcal{O}}(\varepsilon_1).$$

Note in the discrete Res-Net, even though the connecting weights are independently initialized,  $\hat{\alpha}_{\ell,j}$  are not mutually independent since they all depend on a *common* set of random outputs from the previous layer. Therefore, Definition 5 restricts the skip-connected paths of the discrete Res-Net  $\{\hat{v}_{1,i}, \hat{\alpha}_{2,i}, \dots, \hat{\alpha}_{L,i}\}_{i=1}^m$  are nearly independent, which makes it possible to construct an ideal initialization with independent skip-connected paths to approximate the discrete one. Then, using a “propagation of chaos” argument (Sznitman, 1991), we obtain Theorem 6.

Now we demonstrate real examples that can achieve our assumptions. We consider the Res-Net initialized by Algorithm 2, which is composed of a standard initialization (Glorot and Bengio, 2010; He et al., 2015) and an additional regression procedure while preserving all initial features.<sup>4</sup> The standard initialization strategy scales the weights to  $\sqrt{m}$ , which diverges in the mean-field limit. We perform the simple  $\ell_2$ -regression to reduce the redundancy of the weights. The result in Theorem 7 shows that Algorithm 2 can produce an  $\varepsilon_1$ -independent initialization when  $m$  is sufficiently large.

**Theorem 7** *Under Assumptions 1 and 2, treat the parameters in assumptions as constants. With probability at least  $1 - \delta$ , Algorithm 2 produces an  $\varepsilon_1$ -independent initial discrete Res-Net with  $\varepsilon_1 \leq \tilde{\mathcal{O}}(\frac{1}{\sqrt{m}})$ .*

---

4. In Algorithm 2, the weights in the last layer  $\{\hat{v}_{L+1,i,1}\}_{i=1}^m$  can also be initialized by the standard initialization followed by an  $\ell_2$ -regression. The  $\ell_2$ -regression (16) can be replaced by a soft version

$$\min_{\{\hat{v}_{\ell,i,j}\}_{i=1}^m} \frac{\lambda_m}{m} \sum_{i=1}^m (\hat{v}_{\ell,i,j})^2 + \left\| \hat{\alpha}_{\ell,j} - \frac{1}{m} \sum_{i=1}^m \hat{v}_{\ell,i,j} \dot{h}_1(\hat{\beta}_{\ell-1,i}) \right\|^2.$$

### 5.3. Finding Global Minimum

We study the convergence of neural feature flow. We show in Theorem 8 that the neural feature flow always finds a globally optimal solution when it converges.

**Theorem 8** *Under Assumptions 1 and 2, assume that the loss function  $\phi$  is convex in the first argument. Let  $(\{v_\ell\}_{\ell=2}^{L+1}, p)$  be the initial continuous Res-Net that satisfies Assumption 3 and  $\Phi_*$  and  $\mathcal{L}^t$  be the solution and loss of the neural feature flow, respectively. If  $\Phi_{*,L}^\beta(\Theta)(t)$  converges in  $\ell_\infty(p)$  and  $\Phi_{*,L+1}^v(\Theta)(t)$  converges in  $\ell_1(p)$  as  $t \rightarrow \infty$ , where  $\Theta \sim p$ , then we have*

$$\lim_{t \rightarrow \infty} \mathcal{L}^t = \sum_{n=1}^N \left[ \min_{y'} \phi(y', y^n) \right].$$

Theorem 8 is an important application of our mean-field framework, which shows that neural feature flow can find a global minimizer after it converges<sup>5</sup>. We prove that the distribution of the weights in the first layer always have a full support in any finite time by Brouwer’s fixed-point theorem. Then, using a similar argument to Chizat and Bach (2018); Nguyen and Pham (2020), we show that all bad local minima are unstable. Our global convergence holds for Res-Nets with arbitrary (finite) depth. Before us, the global optimality was proved for three-layer ones (Nguyen and Pham, 2020) under a similar convergence assumption on the weights in the second layer.

## 6. Discussion

This paper proposed a new mean-field framework for DNNs where features in hidden layers have non-vanishing variance. We constructed a continuous dynamic called neural feature flow that captures the evolution of sufficiently over-parametrized Res-Nets trained by Gradient Descent. Furthermore, the neural feature flow reaches a globally optimal solution after it converges. We hope our new analytical tool pioneers better understandings for DNN training.

There are many interesting questions under our framework to be further investigated. First, the current work only focuses on the training part of NNs, and it still remains to study generalization. The generalization error in the mean-field regime has been studied for two-level NNs, e.g., Wei et al. (2018). Using our modeling, there are two potential directions: we may incorporate suitable regularizers on the DNNs to control the model complexity; implicit regularization is often observed in practice, which is hopefully preserved in our neural feature flow. However, a full treatment is left as a future study. Second, although our mean-field framework is applicable to standard DNNs (see Appendix E), it is still not answered how to prove that Gradient Descent can achieve the global minimum. Third, the implications of our theory in practice can be studied empirically as separate works. For example, this paper proposes a new interesting initialization strategy for DNNs and uses the scaled Gradient Descent to optimize DNNs. In Appendix F, we perform a toy simulation to validate the feasibility of this new training strategy. It is interesting to design new practical algorithms based on our learning strategy for large-scale data.

5. Remark: for  $\ell_2$  loss, i.e.  $\phi(y', y) = \|y' - y\|^2$ , as an example, Theorem 8 indicates converging to 0 loss.

## References

- Zeyuan Allen-Zhu and Yuanzhi Li. Can sgd learn recurrent neural networks with provable generalization? *arXiv:1902.01028*, 2019.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. *arXiv:1811.04918*, 2018.
- Dyego Araújo, Roberto I Oliveira, and Daniel Yukimura. A mean-field limit for certain deep neural networks. *arXiv:1906.00193*, 2019.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, 2019.
- Rajendra Bhatia. *Matrix analysis*, volume 169. Springer Science & Business Media, 2013.
- Zixiang Chen, Yuan Cao, Quanquan Gu, and Tong Zhang. Mean-field analysis of two-layer neural networks: Non-asymptotic rates and generalization bounds. *arXiv:2002.04026*, 2020.
- Lenaïc Chizat and Francis Bach. On the global convergence of gradient descent for overparameterized models using optimal transport. In *Advances in neural information processing systems*, pages 3036–3046, 2018.
- Lenaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In *Advances in Neural Information Processing Systems*, pages 2933–2943, 2019.
- Donald A Dawson. Multilevel mutation-selection systems and set-valued duals. *Journal of mathematical biology*, 76(1-2):295–378, 2018.
- Donald A Dawson, Kenneth J Hochberg, et al. Wandering random measures in the Fleming-Viot model. *The Annals of Probability*, 10(3):554–580, 1982.
- Xialiang Dou and Tengyuan Liang. Training neural networks as learning data-adaptive kernels: Provable representation and approximation benefits. *arXiv:1901.07114*, 2019.
- Simon S Du, Jason D Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, 2019a.
- Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representation*, 2019b.
- Cong Fang, Hanze Dong, and Tong Zhang. Over parameterized two-level neural networks can learn nearoptimal feature representations. *arXiv:1910.11508*, 2019.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256, 2010.

- Andrzej Granas and James Dugundji. *Fixed point theory*. Springer Science & Business Media, 2013.
- Moritz Hardt and Tengyu Ma. Identity matters in deep learning. In *International Conference on Learning Representation*, 2016.
- Philip Hartman. *Ordinary differential equations*. Wiley, 1964.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, 2018.
- Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems*, 2018.
- Yiping Lu, Chao Ma, Yulong Lu, Jianfeng Lu, and Lexing Ying. A mean-field analysis of deep resnet and beyond: Towards provable optimization via overparameterization from depth. In *International Conference on Machine Learning*, 2020.
- Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018. ISSN 0027-8424.
- Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Annual Conference on Learning Theory*, 2019.
- Phan-Minh Nguyen and Huy Tuan Pham. A rigorous framework for the mean field limit of multi-layer neural networks. *arXiv:2001.11443*, 2020.
- Huy Tuan Pham and Phan-Minh Nguyen. A note on the global convergence of multilayer neural networks in the mean field regime. *arXiv preprint arXiv:2006.09355*, 2020.
- Grant M Rotskoff and Eric Vanden-Eijnden. Neural networks as interacting particle systems: Asymptotic convexity of the loss landscape and universal scaling of the approximation error. *arXiv:1805.00915*, 2018.
- Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of neural networks: A central limit theorem. *Stochastic Processes and their Applications*, 2019a.
- Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of deep neural networks. *arXiv:1903.04440*, 2019b.



- Alain-Sol Sznitman. Topics in propagation of chaos. In *Ecole d’été de probabilités de Saint-Flour XIX—1989*, pages 165–251. Springer, 1991.
- Joel A Tropp. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.
- Mihai Turinici. Abstract comparison principles and multivariable gronwall-bellman inequalities. *Journal of Mathematical Analysis and Applications*, 117(1):100–127, 1986.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv:1011.3027*, 2010.
- Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and optimization of neural nets v.s. their induced kernel. *arXiv:1810.05369*, 2018.
- Gilad Yehudai and Ohad Shamir. On the power and limitations of random features for understanding neural networks. In *Advances in Neural Information Processing Systems*, pages 6594–6604, 2019.
- Hongyi Zhang, Yann N Dauphin, and Tengyu Ma. Fixup initialization: Residual learning without normalization. *arXiv:1901.09321*, 2019.
- Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep relu networks. In *Advances in neural information processing systems*, 2018.

## Appendix A. Overview

The appendix is sketched as follows. Appendix B provides the key proofs of this paper. Especially, we will show how the initial condition (Theorem 7) and Assumption 2 (Lemma 3) can be realized. Besides, we will prove Theorem 8. Appendix C presents the rest proofs. Especially, we will follow the technique of Picard iterations to prove Theorem 4 and a continuity argument for neural feature flow (Theorem 13). Note that the latter will be used in Theorem 8. Moreover, we will follow a “propagation of chaos” argument (Sznitman, 1991) to prove Theorem 6. Appendix D presents some basic properties for sub-gaussian distributions. Finally, Appendix E introduces the extension of our mean-field framework to fully-connected DNNs.

In our proofs, we fix a set of training data and treat the parameters in the assumptions as constants. We use  $C$  to denote a generic constant; the value of  $C$  may change from line to line.

## Appendix B. Key Proofs

### B.1. Proof of Theorem 7

In this subsection, we prove Theorem 7 which states that Algorithm 2 produces an  $\varepsilon_1$ -independent initial discrete Res-Net with  $\varepsilon_1 \leq \tilde{O}(1/\sqrt{m})$ . By Definition 5, this entails the construction of an initial distribution  $p$ , weight functions  $\{v_\ell\}_{\ell=2}^{L+1}$ , and an ideal discrete Res-Net satisfying the properties in Definition 5. We specify  $p$ ,  $\{v_\ell\}_{\ell=2}^{L+1}$ , and the ideal discrete Res-Net below, and then we verify the properties in Theorem 7.

**Initial distribution.** We first define the distribution  $p$ :

- (1) At the first layer,  $\beta_1 \sim p_1^\beta = \mathcal{N}(\mathbf{0}^N, \sigma_1^2 \mathbf{K}_0)$ , where  $\mathbf{K}_0 := \frac{1}{d} \mathbf{X} \mathbf{X}^\top$ .
- (2) At the layer  $\ell \in [L-1]$ , let  $\mathbf{K}_\ell^\beta := \int \dot{h}_1(\beta_\ell) \dot{h}_1(\beta_\ell)^\top dp_\ell^\beta(\beta_\ell)$ . We define the distribution of the residuals at layer  $\ell+1$  as

$$p_{\ell+1}^\alpha = \mathcal{N}(\mathbf{0}^N, \sigma_1^2 \mathbf{K}_\ell^\beta). \quad (17)$$

Defining the mapping  $\tilde{f}_{\ell+1}(\beta_\ell, \alpha_{\ell+1}) := \beta_\ell + \dot{h}_2(\alpha_{\ell+1})$ , the features at layer  $\ell+1$  is defined as the pushforward measure by  $\tilde{f}_{\ell+1}$ :

$$p_{\ell+1}^\beta = \tilde{f}_{\ell+1} \# (p_\ell^\beta \times p_{\ell+1}^\alpha).$$

Finally, let  $p$  be a multivariate Gaussian distribution of the form

$$p(\mathbf{v}_1, \alpha_1, \alpha_2, \dots, \alpha_L) := p_1^v(\mathbf{v}_1) \times p_2^\alpha(\alpha_2) \times p_3^\alpha(\alpha_3) \times \dots \times p_L^\alpha(\alpha_L). \quad (18)$$

**Weight functions.** Now we define the weight functions  $\{v_\ell\}_{\ell=2}^{L+1}$ . Note that those gram matrices  $\mathbf{K}_\ell^\beta$  are all positive definite under Assumption 2 (see Lemma 9) and thus are invertible. For  $\Theta = (\mathbf{v}_1, \alpha_2, \dots, \alpha_L)$ ,  $\Theta' = (\mathbf{v}'_1, \alpha'_2, \dots, \alpha'_L)$ , we define the connecting weights between consecutive layers by

$$v_\ell(\Theta, \Theta') = \dot{h}_1(\beta_{\ell-1}(\Theta))^\top \left[ \mathbf{K}_{\ell-1}^\beta \right]^{-1} \alpha'_\ell, \quad \ell \in [2:L], \quad (19)$$

where  $\beta_\ell(\Theta) = \frac{1}{d} \mathbf{X} \mathbf{v}_1 + \sum_{i=2}^\ell \dot{h}_2(\alpha_i)$  will satisfy (9). The weights at the output layer are initialized as a constant  $c$ . The forward propagation constraint (8) will also be satisfied by (19) and the definitions of  $\mathbf{K}_\ell^\beta$ . Therefore  $(\{v_\ell\}_{\ell=2}^{L+1}, p)$  constitutes a feasible continuous Res-Net.

**Ideal discrete Res-Net.** Finally we construct the initialization  $(\bar{\mathbf{v}}, \bar{\alpha}, \bar{\beta})$  of the ideal discrete Res-Net. Recall Algorithm 2 that the corresponding variables are initialized as  $(\hat{\mathbf{v}}, \hat{\alpha}, \hat{\beta})$ . Let  $\bar{v}_{1,i} := \hat{v}_{i,1}$  for  $i \in [m]$ . For  $\ell \in [L-1]$ , define the empirical Gram matrix as

$$\hat{\mathbf{K}}_\ell^\beta = \frac{1}{m} \sum_{i=1}^m \dot{h}_1(\hat{\beta}_{\ell,i}) \dot{h}_1(\hat{\beta}_{\ell,i})^\top.$$

Let  $\bar{\alpha}_{\ell+1,j} := \left( \mathbf{K}_\ell^\beta \right)^{1/2} \left( \hat{\mathbf{K}}_\ell^\beta \right)^{-1/2} \hat{\alpha}_{\ell+1,j}$  for all  $j \in [m]$  when  $\hat{\mathbf{K}}_\ell^\beta$  is invertible, and otherwise let  $\bar{\alpha}_{\ell+1,j} \stackrel{\text{i.i.d.}}{\sim} p_{\ell+1}^\alpha$ . We use Definition 5 (2) for the values of  $\bar{\beta}_\ell$  for  $\ell \in [L]$  and  $\bar{v}_\ell$  for all in  $[2:L+1]$ .

**Proof [Proof of Theorem 7]** We first show that the continuous Res-Net satisfies Assumptions 3. By definition  $p$  is a multivariate Gaussian distribution. By Lemma 9, we have  $\left\| \left( \mathbf{K}_\ell^\beta \right)^{-1} \right\|_2 \leq C$  for a constant  $C$ . Since the activation function  $h_1$  is bounded and Lipschitz continuous, the continuity conditions in Assumption 3 are all satisfied.

Now we consider the ideal discrete Res-Net. We first verify the independence. By definition,  $\bar{\alpha}_{\ell+1,j}$  are determined by the outputs of previous layer  $\hat{\beta}_{\ell,i}$  and the connecting weights  $\hat{v}_{\ell+1,i,j}$ .

Thus they are conditionally independent of  $\bar{v}_{1,i}$  and  $\bar{\alpha}_{2,i}, \dots, \bar{\alpha}_{\ell,i}$  for  $i \in [m]$  given  $\{\hat{\beta}_{\ell,i}\}_{i \in [m]}$ . Since  $\hat{v}_{\ell+1,i,j}$  are independent Gaussian,  $\hat{\alpha}_{\ell+1,j}$  and thus  $\bar{\alpha}_{\ell+1,j}$  are conditionally independent Gaussian given  $\{\hat{\beta}_{\ell,i}\}_{i \in [m]}$ . Furthermore, the conditional distribution of  $\bar{\alpha}_{\ell+1,j}$  given  $\{\hat{\beta}_{\ell,i}\}_{i \in [m]}$  is  $\mathcal{N}(\mathbf{0}^N, \sigma_1^2 \mathbf{K}_\ell^\beta) = p_{\ell+1}^\alpha$ . Therefore, marginally  $\bar{\alpha}_{\ell+1,j} \stackrel{\text{i.i.d.}}{\sim} p_{\ell+1}^\alpha$  and they are independent of  $\bar{v}_{1,i}$  and  $\bar{\alpha}_{2,i}, \dots, \bar{\alpha}_{\ell,i}$  for  $i \in [m]$ . So  $\{\bar{\Theta}_i\}_{i \in [m]} \stackrel{\text{i.i.d.}}{\sim} p$ . Since  $p$  is a product distribution, all  $\bar{v}_{1,i}, \bar{\alpha}_{2,i}, \dots, \bar{\alpha}_{L,i}$  are all mutually independent.

Lastly we show the  $\tilde{O}(1/\sqrt{m})$ -closeness specified in Definition 5 (3). By Lemma 10 we have the following events with probability  $1 - \delta$ :

$$\|\bar{\alpha}_{\ell+1,i}\|_2 \leq \tilde{O}(1), \quad (20)$$

$$\left\| \hat{\mathbf{K}}_\ell^\beta - \mathbf{K}_\ell^\beta \right\|_2 \leq \varepsilon_2, \quad (21)$$

$$\|\bar{\alpha}_{\ell+1,i} - \hat{\alpha}_{\ell+1,i}\|_2 \leq \varepsilon_2 \|\bar{\Theta}_i\|_2, \quad (22)$$

$$\left\| \bar{\beta}_{\ell+1,i} - \hat{\beta}_{\ell+1,i} \right\|_2 \leq \varepsilon_2 \|\bar{\Theta}_i\|_2, \quad (23)$$

where  $\varepsilon_2 \leq \tilde{O}(1/\sqrt{m})$ . Under (21) the matrix  $\hat{\mathbf{K}}_\ell^\beta$  is invertible, and it follows from Lemma 11 that

$$\hat{v}_{\ell+1,i,j} = \hat{\alpha}_{\ell+1,j}^\top \left[ \hat{\mathbf{K}}_\ell^\beta \right]^{-1} \dot{h}_1(\hat{\beta}_{\ell,i}), \quad \ell \in [L-1], i, j \in [m].$$

By the triangle inequality,

$$\begin{aligned} & |\bar{v}_{\ell+1,i,j} - \hat{v}_{\ell+1,i,j}| \\ &= \left\| \bar{\alpha}_{\ell+1,j}^\top \left[ \mathbf{K}_\ell^\beta \right]^{-1} \dot{h}_1(\bar{\beta}_{\ell,i}) - \hat{\alpha}_{\ell+1,j}^\top \left[ \hat{\mathbf{K}}_\ell^\beta \right]^{-1} \dot{h}_1(\hat{\beta}_{\ell,i}) \right\|_2 \\ &\leq \|\bar{\alpha}_{\ell+1,j}\|_2 \left\| \left[ \mathbf{K}_\ell^\beta \right]^{-1} \right\|_2 \left\| \dot{h}_1(\bar{\beta}_{\ell,i}) - \dot{h}_1(\hat{\beta}_{\ell,i}) \right\|_2 + \|\bar{\alpha}_{\ell+1,j}\|_2 \left\| \left[ \mathbf{K}_\ell^\beta \right]^{-1} - \left[ \hat{\mathbf{K}}_\ell^\beta \right]^{-1} \right\|_2 \left\| \dot{h}_1(\hat{\beta}_{\ell,i}) \right\|_2 \\ &\quad + \|\bar{\alpha}_{\ell+1,j} - \hat{\alpha}_{\ell+1,j}\|_2 \left\| \left[ \hat{\mathbf{K}}_\ell^\beta \right]^{-1} \right\|_2 \left\| \dot{h}_1(\hat{\beta}_{\ell,i}) \right\|_2. \end{aligned}$$

We upper bound three terms separately. By the Lipschitz continuity of  $h_1$  and (23), the first term is at most  $\tilde{O}(1/\sqrt{m}) \|\bar{\Theta}_j\|_2$ ; the second term is also at most  $\tilde{O}(1/\sqrt{m}) \|\bar{\Theta}_j\|_2$  since  $\|\bar{\alpha}_{\ell+1,j}\|_2 \leq \|\bar{\Theta}_j\|_2$ ,

$$\left\| \left[ \mathbf{K}_\ell^\beta \right]^{-1} - \left[ \hat{\mathbf{K}}_\ell^\beta \right]^{-1} \right\|_2 \leq \left\| \left[ \mathbf{K}_\ell^\beta \right]^{-1} \right\|_2 \left\| \mathbf{K}_\ell^\beta - \hat{\mathbf{K}}_\ell^\beta \right\|_2 \left\| \left[ \hat{\mathbf{K}}_\ell^\beta \right]^{-1} \right\|_2 \leq C \varepsilon_2,$$

and  $h_1$  is bounded; the third term is at most  $\tilde{O}(1/\sqrt{m}) \|\bar{\Theta}_j\|_2$  by (22).  $\blacksquare$

### B.1.1. PROOF OF ADDITIONAL LEMMAS

**Lemma 9**  $\min_{\ell \in [L-1]} \lambda_{\min}(\mathbf{K}_\ell^\beta) \geq C > 0$ .

**Proof** Fix  $\ell \in [L - 1]$ . For  $(\mathbf{v}_1, \boldsymbol{\alpha}_2, \dots, \boldsymbol{\alpha}_L) \in \text{supp}(p)$ , given  $\boldsymbol{\alpha}_2, \dots, \boldsymbol{\alpha}_L$ , we have  $\boldsymbol{\Theta} = \boldsymbol{\Theta}(\mathbf{v}_1)$  and

$$\left\| \boldsymbol{\beta}_\ell(\boldsymbol{\Theta}) - \frac{1}{d} \mathbf{X} \mathbf{v}_1 \right\|_\infty = \left\| \sum_{\ell_1=2}^{\ell} \dot{h}_2(\boldsymbol{\alpha}_{\ell_1}) \right\|_\infty \leq LL_1. \quad (24)$$

Note that  $\mathbf{v}_1$  is independent of  $\boldsymbol{\alpha}_2, \dots, \boldsymbol{\alpha}_L$ , and  $\mathbf{v}_1 \sim \mathcal{N}(0, d\sigma_1^2 \mathbf{I}^d)$  which is equivalent to the standard Gaussian distribution. By Assumption 2 with  $f_2(\mathbf{v}_1) \equiv \sum_{\ell_1=2}^{\ell} \dot{h}_2(\boldsymbol{\alpha}_{\ell_1})$  and the constant  $C_B = LL_1$ , we have

$$\mathbb{E} \left[ \dot{h}_1(\boldsymbol{\beta}_\ell(\boldsymbol{\Theta})) \dot{h}_1^\top(\boldsymbol{\beta}_\ell(\boldsymbol{\Theta})) \mid \boldsymbol{\alpha}_2, \dots, \boldsymbol{\alpha}_L \right] \succeq C \mathbf{I}^N.$$

Taking full expectation, we obtain Lemma 9. ■

In the sequel, we set  $\bar{\lambda}_1 := \min_{\ell \in [L-1]} \lambda_{\min}(\mathbf{K}_\ell^\beta)$  that is strictly bounded away from zero.

**Lemma 10** *Let  $\varepsilon_2 \leq \tilde{\mathcal{O}}(1/\sqrt{m})$ . With probability  $1 - \delta$ , for all  $\ell \in [L - 1]$  and  $i \in [m]$ ,*

$$\begin{aligned} \left\| \hat{\mathbf{K}}_\ell^\beta - \mathbf{K}_\ell^\beta \right\|_2 &\leq \varepsilon_2, & \left\| \bar{\boldsymbol{\alpha}}_{\ell+1,i} - \hat{\boldsymbol{\alpha}}_{\ell+1,i} \right\|_2 &\leq \varepsilon_2 \left\| \bar{\boldsymbol{\Theta}}_i \right\|_2, \\ \left\| \bar{\boldsymbol{\alpha}}_{\ell+1,i} \right\|_2 &\leq \tilde{\mathcal{O}}(1), & \left\| \bar{\boldsymbol{\beta}}_{\ell+1,i} - \hat{\boldsymbol{\beta}}_{\ell+1,i} \right\|_2 &\leq \varepsilon_2 \left\| \bar{\boldsymbol{\Theta}}_i \right\|_2. \end{aligned}$$

**Proof** In the proof of Lemma 7 we verified that  $\bar{\mathbf{v}}_{1,i}, \bar{\boldsymbol{\alpha}}_{2,i}, \dots, \bar{\boldsymbol{\alpha}}_{L,i}$  for all  $i \in [m]$  are independent. Therefore,  $\bar{\boldsymbol{\beta}}_{\ell,i} \stackrel{\text{i.i.d.}}{\sim} p_\ell^\beta$  by the definitions of  $\bar{\boldsymbol{\beta}}_{\ell,i}$  and  $p_\ell^\beta$ . Consider auxiliary random matrices

$$\bar{\mathbf{K}}_\ell^\beta := \frac{1}{m} \sum_{i=1}^m \dot{h}_1(\bar{\boldsymbol{\beta}}_{\ell,i}) \dot{h}_1^\top(\bar{\boldsymbol{\beta}}_{\ell,i}),$$

Since  $h_1$  is bounded, by the matrix Bernstein inequality (Tropp, 2015), with probability  $1 - \frac{\delta}{3}$ ,

$$\max_{\ell \in [L-1]} \left\| \bar{\mathbf{K}}_\ell^\beta - \mathbf{K}_\ell^\beta \right\|_2 \leq \varepsilon_3 = \tilde{\mathcal{O}}(1/\sqrt{m}). \quad (25)$$

Due to the sub-gaussianness of  $p$ , we have  $\left\| \bar{\boldsymbol{\alpha}}_{\ell+1,i} \right\|_2 \leq C \sqrt{\log(m/\delta)} = \tilde{\mathcal{O}}(1)$  with probability  $1 - \delta/3$  (see Lemma 27 (1)). We will also use the following upper bound that happens with probability  $1 - \delta/3$  by the sub-gaussianness of  $p$ :

$$\frac{1}{m} \sum_{i=1}^m \left\| \bar{\boldsymbol{\Theta}}_i \right\|_2 \leq \beta_1 = \tilde{\mathcal{O}}(1), \quad (26)$$

which can be obtained by the concentration inequality in Lemma 28).

Next we inductively prove that, for  $\ell \in [L - 1]$ ,

$$\left\| \mathbf{K}_\ell^\beta - \hat{\mathbf{K}}_\ell^\beta \right\|_2 \leq (C\beta_1)^{\ell-1} C\varepsilon_3, \quad (27)$$

$$\left\| \hat{\boldsymbol{\alpha}}_{\ell+1,i} - \bar{\boldsymbol{\alpha}}_{\ell+1,i} \right\|_2 \leq (C\beta_1)^{\ell-1} C\varepsilon_3 \left\| \bar{\boldsymbol{\Theta}}_i \right\|_2, \quad (28)$$

$$\left\| \hat{\boldsymbol{\beta}}_{\ell+1,i} - \bar{\boldsymbol{\beta}}_{\ell+1,i} \right\|_2 \leq (C\beta_1)^{\ell-1} C\varepsilon_3 \left\| \bar{\boldsymbol{\Theta}}_i \right\|_2. \quad (29)$$

For  $\ell = 1$ , by definition  $\bar{\mathbf{K}}_1^\beta = \hat{\mathbf{K}}_1^\beta$ . The upper bound of  $\left\| \left[ \hat{\mathbf{K}}_1^\beta \right]^{1/2} - \left[ \mathbf{K}_1^\beta \right]^{1/2} \right\|_2$  is achieved by matrix calculus (Bhatia, 2013, Section V.3). Since  $\left\| \hat{\mathbf{K}}_1^\beta - \mathbf{K}_1^\beta \right\|_2 \leq \frac{\bar{\lambda}_1}{2}$ , then the eigenvalues of  $\hat{\mathbf{K}}_1^\beta$  are at least  $\frac{\bar{\lambda}_1}{2}$ . Let  $f(x) := \sqrt{x}$ . Then  $|f'(x)| \geq \frac{1}{\sqrt{2\bar{\lambda}_1}}$  when  $x$  is the eigenvalue of  $\hat{\mathbf{K}}_1^\beta$ . Applying (Bhatia, 2013, (V.20)) yields that

$$\left\| \left[ \hat{\mathbf{K}}_1^\beta \right]^{1/2} - \left[ \mathbf{K}_1^\beta \right]^{1/2} \right\|_2 \leq \frac{N}{\sqrt{2\bar{\lambda}_1}} \left\| \hat{\mathbf{K}}_1^\beta - \mathbf{K}_1^\beta \right\|_2 \leq C\varepsilon_3, \quad (30)$$

and

$$\left\| \hat{\alpha}_{2,i} - \bar{\alpha}_{2,i} \right\|_2 = \left\| \left( \left[ \hat{\mathbf{K}}_1^\beta \right]^{1/2} \left[ \mathbf{K}_1^\beta \right]^{-1/2} - \mathbf{I}^N \right) \bar{\alpha}_{2,i} \right\|_2 \leq C\varepsilon_3 \left\| \bar{\alpha}_{2,i} \right\|_2 \leq C\varepsilon_3 \left\| \bar{\Theta}_i \right\|_2. \quad (31)$$

Then by the Lipschitz continuity of  $h_2$ , we have

$$\left\| \hat{\beta}_{2,i} - \bar{\beta}_{2,i} \right\|_2 = \left\| \left[ h_2(\hat{\alpha}_{2,i}) - h_2(\bar{\alpha}_{2,i}) \right] \right\|_2 \leq C\varepsilon_3 \left\| \bar{\Theta}_i \right\|_2. \quad (32)$$

For  $\ell \in [2 : L - 1]$ , suppose that

$$\left\| \hat{\beta}_{\ell,i} - \bar{\beta}_{\ell,i} \right\|_2 \leq (C\beta_1)^{\ell-2} C\varepsilon_3 \left\| \bar{\Theta}_i \right\|_2. \quad (33)$$

Then, by the boundedness of  $h_1$  and the triangle inequality, we have

$$\left\| \hat{\mathbf{K}}_\ell^\beta - \bar{\mathbf{K}}_\ell^\beta \right\|_2 \leq \frac{C}{m} \sum_{i=1}^m \left\| \dot{h}_1(\bar{\beta}_{\ell,i}) - \dot{h}_1(\hat{\beta}_{\ell,i}) \right\|_2.$$

Applying the Lipschitz continuity of  $h_1$  and (33) yields that

$$\left\| \hat{\mathbf{K}}_\ell^\beta - \bar{\mathbf{K}}_\ell^\beta \right\|_2 \leq \frac{(C\beta_1)^{\ell-2} C\varepsilon_3}{m} \sum_{i=1}^m \left\| \bar{\Theta}_i \right\|_2 \leq (C\beta_1)^{\ell-1} C\varepsilon_3. \quad (34)$$

where in the last inequality we used (26). Then we obtain (27) by triangle inequality from (25) and (34). The upper bound in (28) for  $\ell + 1$  follows from a similar argument of (30) and (31). Finally (29) for  $\ell + 1$  follows from (28) and

$$\left\| \hat{\beta}_{\ell+1,i} - \bar{\beta}_{\ell+1,i} \right\|_2 = \left\| \sum_{j=2}^{\ell+1} \left[ \dot{h}_2(\hat{\alpha}_{j,i}) - \dot{h}_2(\bar{\alpha}_{j,i}) \right] \right\|_2 \leq C(C\beta_1)^{\ell-1} \varepsilon_3 \left\| \bar{\Theta}_i \right\|_2.$$

We finish the induction. Since  $\beta = \tilde{O}(1)$  and  $\varepsilon_3 = \tilde{O}(1/\sqrt{m})$ , we complete the proof.  $\blacksquare$

**Lemma 11** *If  $\hat{\mathbf{K}}_\ell^\beta$  is invertible, then*

$$\hat{v}_{\ell+1,i,j} = \hat{\alpha}_{\ell+1,j}^\top \left[ \hat{\mathbf{K}}_\ell^\beta \right]^{-1} \dot{h}_1(\hat{\beta}_{\ell,i}), \quad \ell \in [L - 1], \quad i, j \in [m].$$

**Proof** For a given layer  $\ell$  and  $j$ , the  $\ell_2$ -regression problem in Algorithm 2 can be equivalently written as

$$\begin{aligned} \min_{\hat{\mathbf{v}}} \quad & \frac{1}{2} \|\hat{\mathbf{v}}\|^2 \\ \text{s.t.} \quad & \frac{1}{m} \hat{\mathbf{H}} \hat{\mathbf{v}} = \hat{\boldsymbol{\alpha}}_{\ell+1,j}, \end{aligned} \quad (35)$$

where  $\hat{\mathbf{v}} = (\hat{v}_{\ell+1,1,j}, \dots, \hat{v}_{\ell+1,m,j})^\top$  and  $\hat{\mathbf{H}} = [\hat{h}_1(\hat{\boldsymbol{\beta}}_{\ell,1}), \dots, \hat{h}_1(\hat{\boldsymbol{\beta}}_{\ell,m})]$ . Decompose  $\hat{\mathbf{v}}$  as

$$\hat{\mathbf{v}} = \hat{\mathbf{H}}^\top \mathbf{z} + \hat{\mathbf{v}}',$$

where  $\mathbf{z} \in \mathbb{R}^m$  and  $\hat{\mathbf{H}} \hat{\mathbf{v}}' = 0$ . Then (35) is equivalent to

$$\begin{aligned} \min_{\mathbf{z}, \hat{\mathbf{v}}'} \quad & \frac{1}{2} \|\hat{\mathbf{H}}^\top \mathbf{z}\|_2^2 + \frac{1}{2} \|\hat{\mathbf{v}}'\|_2^2 \\ \text{s.t.} \quad & \frac{1}{m} \hat{\mathbf{H}} \hat{\mathbf{H}}^\top \mathbf{z} = \hat{\boldsymbol{\alpha}}_{\ell+1,j}. \end{aligned}$$

Since  $\frac{1}{m} \hat{\mathbf{H}} \hat{\mathbf{H}}^\top = \hat{\mathbf{K}}_\ell^\beta$  is invertible, the optimal solution is  $\mathbf{z} = [\hat{\mathbf{K}}_\ell^\beta]^{-1} \hat{\boldsymbol{\alpha}}_{\ell+1,j}$  and  $\hat{\mathbf{v}}' = \mathbf{0}^N$ . ■

## B.2. Proof of Theorem 8

**Proof** [Proof of Theorem 8] In the proof we use the following abbreviated notations: for  $t \in [0, \infty)$  and  $\boldsymbol{\Theta} \in \text{supp}(p)$ , let

$$\begin{aligned} \boldsymbol{\beta}_\ell^t(\boldsymbol{\Theta}) &= \Phi_{*,\ell}^\beta(\boldsymbol{\Theta})(t), \quad \ell \in [L], \\ \boldsymbol{\alpha}_\ell^t(\boldsymbol{\Theta}) &= \Phi_{*,\ell}^\alpha(\boldsymbol{\Theta})(t), \quad \ell \in [2 : L], \\ \mathbf{v}_1^t(\boldsymbol{\Theta}) &= \Phi_{*,1}^v(\boldsymbol{\Theta})(t), \\ \mathbf{v}_{L+1}^t(\boldsymbol{\Theta}) &= \Phi_{*,L+1}^v(\boldsymbol{\Theta})(t). \end{aligned}$$

From the convergence assumptions, it is clear that  $\boldsymbol{\beta}_{L+1}^t$  converges as  $t \rightarrow \infty$ . Indeed, the convergence assumptions imply that, for any  $\varepsilon_2 > 0$ , there exists  $T$ , for any  $t \geq T$ ,

$$\|\boldsymbol{\beta}_L^t(\boldsymbol{\Theta}) - \boldsymbol{\beta}_L^\infty(\boldsymbol{\Theta})\|_\infty \leq \varepsilon_2 \quad (36)$$

holds  $p$ -almost surely and

$$\int |v_{L+1}^t(\boldsymbol{\Theta}) - v_{L+1}^\infty(\boldsymbol{\Theta})| p(\boldsymbol{\Theta}) \leq \varepsilon_2. \quad (37)$$

Then, since  $h_1$  is bounded and Lipschitz continuous, we have

$$\begin{aligned} & \|\boldsymbol{\beta}_{L+1}^t - \boldsymbol{\beta}_{L+1}^\infty\|_\infty \\ &= \left\| \int v_{L+1}^t(\boldsymbol{\Theta}) \dot{h}_1(\boldsymbol{\beta}_L^t(\boldsymbol{\Theta})) - v_{L+1}^\infty(\boldsymbol{\Theta}) \dot{h}_1(\boldsymbol{\beta}_L^\infty(\boldsymbol{\Theta})) dp(\boldsymbol{\Theta}) \right\|_\infty \\ &\leq \int |v_{L+1}^\infty(\boldsymbol{\Theta})| \left\| \dot{h}_1(\boldsymbol{\beta}_L^t(\boldsymbol{\Theta})) - \dot{h}_1(\boldsymbol{\beta}_L^\infty(\boldsymbol{\Theta})) \right\|_\infty dp(\boldsymbol{\Theta}) \\ &\quad + \int |v_{L+1}^t(\boldsymbol{\Theta}) - v_{L+1}^\infty(\boldsymbol{\Theta})| \left\| \dot{h}_1(\boldsymbol{\beta}_L^t(\boldsymbol{\Theta})) \right\|_\infty dp(\boldsymbol{\Theta}) \\ &\leq C\varepsilon_2. \end{aligned} \quad (38)$$

The goal of the proof is to show that

$$\left\| \dot{\phi}'_1(\beta_{L+1}^\infty) \right\|_2 = 0. \quad (39)$$

To this end, for any  $\varepsilon > 0$ , we will construct a function

$$f_\varepsilon(\mathbf{v}_1) := \dot{\phi}'_1(\beta_{L+1}^\infty)^\top \dot{h}_1 \left( \frac{1}{d} \mathbf{X} \mathbf{v}_1 + g_\varepsilon(\mathbf{v}_1) \right), \quad (40)$$

where the functions  $g_\varepsilon$  is uniformly bounded, such that  $|f_\varepsilon(\mathbf{v}_1)| < \varepsilon$ . Then it follows from (40) that

$$\dot{\phi}'_1(\beta_{L+1}^\infty) = \mathbf{K}^{-1} \int f_\varepsilon(\mathbf{v}_1) \dot{h}_1 \left( \frac{1}{d} \mathbf{X} \mathbf{v}_1 + g_\varepsilon(\mathbf{v}_1) \right) d\tilde{p}_1(\mathbf{v}_1),$$

where  $\tilde{p}_1 = \mathcal{N}(\mathbf{0}^d, \mathbf{I}^d)$  and  $\mathbf{K} := \int \dot{h}_1 \left( \frac{1}{d} \mathbf{X} \mathbf{v}_1 + g_{\varepsilon, \eta}(\mathbf{v}_1) \right) \dot{h}_1^\top \left( \frac{1}{d} \mathbf{X} \mathbf{v}_1 + g_{\varepsilon, \eta}(\mathbf{v}_1) \right) d\tilde{p}_1(\mathbf{v}_1)$  whose minimum eigenvalue is at least  $\bar{\lambda}_1 > 0$  by Assumption 2. The boundedness of  $\dot{h}_1$  yields that

$$\left\| \dot{\phi}'_1(\beta_{L+1}^\infty) \right\|_2 \leq C \bar{\lambda}_1^{-1} \varepsilon.$$

Since  $\bar{\lambda}_1$  is independent of  $\varepsilon$ , by letting  $\varepsilon \rightarrow 0$ , we obtain (39).

Next we construct  $g_\varepsilon$  and  $f_\varepsilon$  in (40). Let  $T$  be the time such that (36) and (37) hold with  $\varepsilon_2 \leq c\varepsilon$  for a constant  $c$  to be specified. Note that  $\mathbf{v}_1^T$  is surjective by Lemma 12. Let  $\tilde{g} : \mathbb{R}^d \rightarrow \text{supp}(p)$  be the inverse function such that  $\mathbf{v}_1^T(\tilde{g}(\mathbf{v}_1)) = \mathbf{v}_1$ . Define

$$g_\varepsilon(\mathbf{v}_1) = \sum_{\ell=2}^L \dot{h}_2(\alpha_\ell^T(\tilde{g}(\mathbf{v}_1))), \quad f_\varepsilon(\mathbf{v}_1) = \dot{\phi}'_1(\beta_{L+1}^\infty)^\top \dot{h}_1(\beta_L^T(\tilde{g}(\mathbf{v}_1))),$$

where  $g_\varepsilon$  is uniformly bounded by the boundedness of  $h_2$ . Suppose on the contrary that there exists  $\mathbf{v}'_1$  such that  $|f_\varepsilon(\mathbf{v}'_1)| > \varepsilon$ . Let  $\Theta' = \tilde{g}(\mathbf{v}'_1)$ . Since  $\Theta \mapsto \dot{\phi}'_1(\beta_{L+1}^\infty)^\top \dot{h}_1(\beta_L^T(\Theta))$  is continuous by Theorem 13 (see Appendix C.1), there exists a ball around  $\Theta'$  denoted by  $S$  such that  $p(S) > 0$  and  $|\dot{\phi}'_1(\beta_{L+1}^\infty)^\top \dot{h}_1(\beta_L^T(\Theta))| > \varepsilon/2$  with the same sign for all  $\Theta \in S$ . However, for  $t > T$ ,

$$\begin{aligned} & \int |v_{L+1}^t(\Theta) - v_{L+1}^T(\Theta)| dp(\Theta) \\ & \geq \frac{1}{N} \int \mathcal{I}_S \left| \int_T^t \dot{\phi}'_1(\beta_{L+1}^t)^\top \dot{h}_1(\beta_L^t(\Theta)) dt \right| dp(\Theta) \\ & \geq \frac{1}{N} \int \mathcal{I}_S \left( \left| \int_T^t \dot{\phi}'_1(\beta_{L+1}^\infty)^\top \dot{h}_1(\beta_L^T(\Theta)) dt \right| - \int_T^t C\varepsilon_2 dt \right) dp(\Theta), \end{aligned} \quad (41)$$

where in the last step we used  $\|\beta_L^t(\Theta) - \beta_L^T(\Theta)\|_\infty \leq 2\varepsilon_2$  from (36), (38), and the boundedness and Lipschitz continuity of  $\dot{\phi}'_1$  and  $\dot{h}_1$ . Let  $c = \frac{1}{4C}$ . The lower bound in (41) diverges with  $t$ , which contradicts (37).

Finally from (39) we show the convergence statement. Since  $\phi$  is convex on the first argument, we obtain

$$\sum_{n=1}^N \phi(\beta_{L+1}^\infty(n), y^n) = \sum_{n=1}^N \left[ \min_{y'} \phi(y', y^n) \right].$$

Since  $\beta_{L+1}^t \rightarrow \beta_{L+1}^\infty$  and  $\phi$  is continuous, we obtain that

$$\lim_{t \rightarrow \infty} L^t = \sum_{n=1}^N \phi(\beta_{L+1}^\infty(n), y^n) = \sum_{n=1}^N \left[ \min_{y'} \phi(y', y^n) \right],$$

which completes the proof.  $\blacksquare$

**Lemma 12** *The function  $t < \infty, \mathbf{v}_1^t : \text{supp}(p) \rightarrow \mathbb{R}^d$  is a surjection.*

**Proof** Recall that at the initialization let  $\Theta(\mathbf{v}) = (\mathbf{v}, \mathbf{0}^{D-d}) \in \mathbb{R}^D = \text{supp}(p)$ , where the equality follows from Assumption 3 that  $p$  has a full support. Given  $t < \infty$ , consider  $f_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$  as

$$f_t(\mathbf{v}) = \mathbf{v}_1^t(\Theta(\mathbf{v})).$$

It suffices to show that  $f_t$  is surjective. Note that  $f_t$  is continuous since  $\Theta \mapsto \mathbf{v}_1^t(\Theta)$  is continuous by Theorem 13. Furthermore, for any  $\mathbf{v} \in \mathbb{R}^d$ , by Lemma 21 which states that the gradient of the weights  $\underline{\mathcal{G}}_1^v$  is bounded (see Appendix C.1), we have

$$\|f_t(\mathbf{v}) - \mathbf{v}\|_\infty = \int_0^t \left\| \underline{\mathcal{G}}_1^v \left( [\mathbf{v}; \mathbf{0}^{D-d}], \Phi_*, s \right) \right\|_\infty ds \leq Ct.$$

For any  $\mathbf{x} \in \mathbb{R}^d$ , consider  $g(\mathbf{v}) := \mathbf{x} - (f_t(\mathbf{v}) - \mathbf{v})$  which continuously maps  $\mathcal{B}_\infty(\mathbf{x}, Ct)$  to itself. By the Brouwer's fixed-point theorem (see, e.g. Granas and Dugundji (2013)), there exists  $\mathbf{v}_* \in \mathcal{B}_\infty(\mathbf{x}, Ct)$  such that  $g(\mathbf{v}_*) = \mathbf{v}_*$ ; equivalently, we have  $f_t(\mathbf{v}_*) = \mathbf{x}$ .  $\blacksquare$

### B.3. Proof of Lemma 3

**Proof** [Proof of Lemma 3] We first note the following results in Du et al. (2019a, Lemma F.1): suppose that  $C_B = 0$  and the support of a random vector  $\mathbf{V} \in \mathbb{R}^d$  denoted by  $R$  has positive Lebesgue measure. Moreover,  $h$  is an analytic non-polynomial function on  $R$ . Then

$$\min_{\|\mathbf{a}\|_2=1} \mathbb{E} \left\| \sum_{i=1}^N a_i h(\mathbf{x}_i \circ \mathbf{V}) \right\|_2^2 = \lambda > 0,$$

where  $\mathbf{a} = (a_1, \dots, a_N)$ . Lemma 3 shows that, for  $\mathbf{V}' \sim p = \mathcal{N}(\mathbf{0}^d, \mathbf{I}^d)$ , the same result holds with a constant perturbation of the functions  $h_1$ ; namely, by letting  $g_i(\mathbf{v}) = h_1(\mathbf{x}_i \circ \mathbf{v} + C_i(\mathbf{v}))$  where  $\|C_i\|_\infty \leq C_B$ ,

$$\min_{\|\mathbf{a}\|_2=1} \mathbb{E} \left\| \sum_{i=1}^N a_i g_i(\mathbf{V}') \right\|_2^2 = \lambda' > 0, \quad (42)$$

where  $\lambda'$  is uniform over all perturbations  $\|C_i\|_\infty \leq C_B$ . It suffices to prove (42) for  $\mathbf{V}' \sim q = \text{Uniform}(R')$  where  $R'$  is determined by  $h$  and  $C_B$ , as the Radon–Nikodym derivative  $\frac{dq}{dp}$  is bounded.



We first prove (i). Consider a compact region  $R$  such that, for  $\mathbf{V} \sim \text{Uniform}(R)$  and any unit vector  $\mathbf{a}$ ,

$$\mathbb{E} \left\| \sum_{i=1}^N a_i g(\mathbf{x}_i \circ \mathbf{V}) \right\|_2^2 \geq \lambda_R > 0.$$

Then for  $\beta > 0$ , since  $g$  is bounded and Lipschitz continuous, we have

$$\mathbb{E} \left\| \sum_{i=1}^N a_i g(\mathbf{x}_i \circ \mathbf{V} + \beta C_i(\mathbf{V}/\beta)) \right\|_2^2 \geq \lambda_R - CC_B \beta \geq \frac{\lambda_R}{2},$$

when  $\beta \leq \frac{\lambda_R}{2CC_B}$ . Let  $h_1(x) = g(\beta x)$ . Then  $\mathbb{E} \left\| \sum_{i=1}^N a_i h_1(\mathbf{x}_i \circ \mathbf{V}/\beta + C_i(\mathbf{V}/\beta)) \right\|_2^2 \geq \frac{\lambda_R}{2}$ . We achieve (42) by letting  $\mathbf{V}' = \mathbf{V}/\beta$ .

For (ii), consider  $R = \{\mathbf{v} : 1/2 \leq \|\mathbf{v}\|_2 \leq 1\}$ . Then, for  $\mathbf{V} \sim \text{Uniform}(R)$  and any unit vector  $\mathbf{a}$ ,

$$\mathbb{E} \left\| \sum_{i=1}^N a_i h_1(\mathbf{x}_i \circ \mathbf{V}) \right\|_2^2 \geq \lambda_R > 0.$$

Since  $h_1(\beta x) = \beta^\alpha x$  for any  $\beta > 0$ , then we have  $\mathbb{E} \left\| \sum_{i=1}^N a_i h_1(\mathbf{x}_i \circ \beta \mathbf{V}) \right\|_2^2 \geq \beta^{2\alpha} \lambda_R$ . Note that  $|\mathbf{x}_i \circ \beta \mathbf{V}| = \Theta(\beta)$ . For  $x = \Theta(\beta)$ , we have  $|h_1(x)| \leq C\beta^\alpha$  and  $h_1$  is  $C\beta^{\alpha-1}$ -Lipschitz continuous for a constant  $C$ . Therefore,

$$\mathbb{E} \left\| \sum_{i=1}^N a_i h_1(\mathbf{x}_i \circ \beta \mathbf{V} + C_i(\beta \mathbf{V})) \right\|_2^2 \geq \beta^{2\alpha} \lambda_R - C' \beta^{2\alpha-1} C_B \geq (C' C_B)^{2\alpha} \left( \frac{2}{\lambda_R} \right)^{2\alpha-1}.$$

for a constant  $C'$  when  $\beta = \frac{2C' C_B}{\lambda_R}$ . We achieve (42) by letting  $\mathbf{V}' = \beta \mathbf{V}$ .

For (iii), we only consider  $h_1(x) = c|x|^{-\alpha}$ . The case for  $h_1(x) = c(x)_+^{-\alpha}$  can be obtained by a similar argument. We first show that there exists a compact set  $R$  such that, for all  $\mathbf{v} \in R$  and  $\mathbf{x}_i$ ,

$$|\mathbf{x}_i \circ \mathbf{v}| \geq c'. \quad (43)$$

This can be done by a simple probabilistic argument. Let  $\mathbf{v}$  be drawn from the uniform distribution on  $S^{d-1}$ , for any fixed  $\mathbf{x} \in \mathbb{R}^d$ , we have

$$P\{|\mathbf{v}^\top \mathbf{x}| < t \|\mathbf{x}\|_2\} = \frac{2\pi^{\frac{d-1}{2}}/\Gamma(\frac{d-1}{2})}{2\pi^{\frac{d}{2}}/\Gamma(\frac{d}{2})} \int_{-t}^t (1-u^2)^{\frac{d-3}{2}} du < t\sqrt{d}.$$

By a union bound, we have  $|\mathbf{x}_i \circ \mathbf{v}| \geq \frac{\|\mathbf{x}_i\|_2}{2N\sqrt{d}}$  with probability 0.5. Denote this set of  $\mathbf{v} \in S^{d-1}$  by  $S'$ . Since  $\min_i \|\mathbf{x}_i\|_2 := C_{mx} > 0$ , we obtain (43) with  $R = \{t\mathbf{v} : \mathbf{v} \in S', \frac{2c'N\sqrt{d}}{C_{mx}} \leq t \leq \frac{4c'N\sqrt{d}}{C_{mx}}\}$ . Then, for  $\mathbf{V} \sim \text{Uniform}(R)$  and any unit vector  $\mathbf{a}$ ,

$$\mathbb{E} \left\| \sum_{i=1}^N a_i h_1(\mathbf{x}_i \circ \mathbf{V}) \right\|_2^2 \geq \lambda_R > 0.$$

Then, for any  $\beta > 0$ , we have  $\mathbb{E}\|\sum_{i=1}^N a_i h_1(\mathbf{x}_i \circ \beta \mathbf{V})\|_2^2 \geq \beta^{-2\alpha} \lambda_R$ . For  $x = \Theta(\beta)$  we have  $|h_1(x)| \leq C\beta^{-\alpha}$  and  $h_1$  is  $C\beta^{-\alpha-1}$ -Lipschitz continuous for a constant  $C$ . Therefore,

$$\mathbb{E} \left\| \sum_{i=1}^N a_i h_1(\mathbf{x}_i \circ \beta \mathbf{V} + C_i(\beta \mathbf{V})) \right\|_2^2 \geq \beta^{-2\alpha} \lambda_R - C' \beta^{-2\alpha-1} C_B \geq (C' C_B)^{-2\alpha} \left( \frac{2}{\lambda_R} \right)^{-2\alpha-1}.$$

for a constant  $C'$  when  $\beta = \frac{2C' C_B}{\lambda_R}$ . We achieve (42) by letting  $\mathbf{V}' = \beta \mathbf{V}$ .  $\blacksquare$

## Appendix C. Rest Proofs

### C.1. Proofs of Theorems 4 and 13

We first demonstrate Theorem 13 which shows that the neural feature flow  $\Phi_*$  is a continuous mapping on  $\Theta$  given time  $t$ .

**Theorem 13 (Property of  $\Phi_*$ )** *Under Assumptions 1 and 3, let  $\Phi_*$  be the neural feature flow, there exist constants  $C, C' \geq 0$  such that for all  $t \in [0, T]$ ,  $\Theta_1 \in \text{supp}(p)$  and  $\tilde{\Theta}_1 \in \text{supp}(p) \cap \mathcal{B}_\infty(\Theta, 1)$ ,  $\Theta_2 \in \text{supp}(p)$ , and  $\tilde{\Theta}_2 \in \text{supp}(p) \cap \mathcal{B}_\infty(\Theta_2, 1)$ , we have*

$$\begin{aligned} \left\| \Phi_{*,\ell}^\beta(\Theta_1)(t) - \Phi_{*,\ell}^\beta(\tilde{\Theta}_1)(t) \right\|_\infty &\leq C e^{C't} (\|\Theta_1\|_\infty + 1) \left\| \Theta_1 - \tilde{\Theta}_1 \right\|_\infty, \quad \ell \in [L], \\ \left\| \Phi_{*,\ell}^\alpha(\Theta_1)(t) - \Phi_{*,\ell}^\alpha(\tilde{\Theta}_1)(t) \right\|_\infty &\leq C e^{C't} (\|\Theta_1\|_\infty + 1) \left\| \Theta_1 - \tilde{\Theta}_1 \right\|_\infty, \quad \ell \in [2 : L], \\ \left\| \Phi_{*,1}^v(\Theta_1)(t) - \Phi_{*,1}^v(\tilde{\Theta}_1)(t) \right\|_\infty &\leq C e^{C't} (\|\Theta_1\|_\infty + 1) \left\| \Theta_1 - \tilde{\Theta}_1 \right\|_\infty, \\ \left| \Phi_{*,\ell}^v(\Theta_1, \Theta_2)(t) - \Phi_{*,\ell}^v(\tilde{\Theta}_1, \Theta_2)(t) \right| &\leq C e^{C't} (\|\Theta_1\|_\infty + \|\Theta_2\|_\infty + 1) \left\| \Theta_1 - \tilde{\Theta}_1 \right\|_\infty, \quad \ell \in [2 : L], \\ \left| \Phi_{*,\ell}^v(\Theta_1, \Theta_2)(t) - \Phi_{*,\ell}^v(\Theta_1, \tilde{\Theta}_2)(t) \right| &\leq C e^{C't} (\|\Theta_1\|_\infty + \|\Theta_2\|_\infty + 1) \left\| \Theta_2 - \tilde{\Theta}_2 \right\|_\infty, \quad \ell \in [2 : L], \\ \left| \Phi_{*,L+1}^v(\Theta_1)(t) - \Phi_{*,L+1}^v(\tilde{\Theta}_1)(t) \right| &\leq C e^{C't} (\|\Theta_1\|_\infty + 1) \left\| \Theta_1 - \tilde{\Theta}_1 \right\|_\infty. \end{aligned}$$

In the proof, we fix the initial continuous Res-Net  $(\{v_\ell\}_{\ell=2}^{L+1}, p)$ , which is assumed to satisfy the Assumption 3.

#### C.1.1. PROOF OF THEOREM 4

We first show that our neural feature flow in Definition 2 necessarily satisfies several continuity properties in Lemma 15, which allows us to narrow down the search space for the solution. Then we construct a contraction mapping (also known as Picard iteration) to show the existence of uniqueness of solution in that search space. Recall that a trajectory  $\Phi$  consists of trajectories of weights  $\Phi_\ell^v$  for  $\ell \in [L+1]$ , features  $\Phi_\ell^\beta$  for  $\ell \in [L]$ , and residuals  $\Phi_\ell^\alpha$  for  $\ell \in [2 : L]$ . For  $\Theta, \tilde{\Theta} \in \text{supp}(p)$ , we also abbreviate the notations for individual trajectories as

$$\Phi_\ell^v(\mathbf{u}_\ell)(t) = v_\ell^t(\mathbf{u}_\ell), \quad \Phi_\ell^\beta(\Theta)(t) = \beta_\ell^t(\Theta), \quad \Phi_\ell^\alpha(\Theta)(t) = \alpha_\ell^t(\Theta),$$

where  $\mathbf{u}_\ell$  stands for  $\Theta, (\Theta, \tilde{\Theta}), \Theta$  for  $\ell = 1, 2 \leq \ell \leq L, \ell = L+1$ , respectively.

Throughout the proof, we fix  $T$  as a constant. We define the set of continuous restricted trajectories below.

**Definition 14 (C-Continuous Restricted Trajectory)** Given  $\mathbf{C} := (\mathbf{C}_1, \dots, \mathbf{C}_{L+1}) \in \mathbb{R}_+^{L+1}$ , we say  $\Phi$  is a  $\mathbf{C}$ -continuous restricted trajectory if  $\Phi_\ell^v(\mathbf{u}_\ell)(t)$  is  $\mathbf{C}_\ell$ -Lipschitz continuous in  $t \in [0, T]$  for  $\ell \in [L+1]$ , and  $\Phi_\ell^\alpha(\mathbf{u}_\ell)(t)$  and  $\Phi_\ell^\beta(\mathbf{u}_\ell)(t)$  are determined by the forward-propagation process, i.e.,  $\beta_1^t(\Theta) = \frac{1}{d} \mathbf{X} v_1^t(\Theta)$ ,  $\alpha_{\ell+1}^t(\Theta) = \int v_{\ell+1}^t(\Theta, \bar{\Theta}) \dot{h}_1(\beta_\ell^t(\bar{\Theta})) dp(\bar{\Theta})$ ,  $\beta_{\ell+1}^t(\Theta) = \beta_\ell^t(\Theta) + \dot{h}_2(\alpha_{\ell+1}^t(\Theta))$  for  $\ell \in [L-1]$  and  $\Theta \in \text{supp}(p)$ . The set of all  $\mathbf{C}$ -continuous restricted trajectories is denoted as  $\Phi^{\mathbf{C}}$ .

We can find that given the trajectories of weights, the trajectories of residuals and features are determined by the forward-propagation process. Lemma 15 below shows that it suffices to consider a restricted search space.

**Lemma 15** There exists a constant vector  $\mathbf{C} \in \mathbb{R}_+^{L+1}$  such that every solution  $\Phi$  of the neural feature flow is a  $\mathbf{C}$ -continuous restricted trajectory.

In the remaining of the proof we let  $\mathbf{C}$  be the constant vector in Lemma 15, and let  $\Phi := \Phi^{\mathbf{C}}$ , which will serve as the search space. The solution can be equivalently characterized as the fixed-point of a mapping from  $\Phi$  to itself that we introduce next:

**Definition 16** Define  $F : \Phi \rightarrow \Phi$  as follows: for all  $t \in [0, T]$ ,

(1) for all  $\ell \in [L+1]$  and all  $\mathbf{u}_\ell$ ,

$$F(\Phi)_\ell^v(\mathbf{u}_\ell)(t) = v_\ell(\mathbf{u}_\ell) - \int_0^t \underline{\mathcal{G}}_\ell^v(\mathbf{u}_\ell; \Phi, s) ds,$$

(2) for all  $\Theta$ ,

$$F(\Phi)_1^\beta(\Theta)(t) = \frac{1}{d} [\mathbf{X} F(\Phi)_1^v(\Theta)(t)],$$

(3) for all  $\ell \in [L-1]$  and  $\Theta$ ,

$$\begin{aligned} F(\Phi)_{\ell+1}^\alpha(\Theta)(t) &= \int F(\Phi)_{\ell+1}^v(\Theta, \bar{\Theta})(t) \dot{h}_1(\Phi_\ell^\beta(\bar{\Theta})(t)) dp(\bar{\Theta}), \\ F(\Phi)_{\ell+1}^\beta(\Theta)(t) &= \dot{h}_2(F(\Phi)_{\ell+1}^\alpha(\Theta)(t)) + F(\Phi)_\ell^\beta(\Theta)(t). \end{aligned}$$

Following the same argument as Lemma 15, we have that the image of  $\Phi$  under  $F$  is indeed contained in  $\Phi$ . We then show in Lemma 18 the contraction property of  $F$  under an appropriate metric defined below:

**Definition 17** For a pair  $\Phi_1, \Phi_2 \in \Phi$ , we define the normalized distance between each trajectories over  $[0, t]$  as

$$D^{[0,t]}(\Phi_1, \Phi_2) := \sup_{s \in [0,t], \ell \in [L+1], \mathbf{u}_\ell} \frac{\|\Phi_{1,\ell}^v(\mathbf{u}_\ell)(s) - \Phi_{2,\ell}^v(\mathbf{u}_\ell)(s)\|_\infty}{1 + \|\mathbf{u}_\ell\|_\infty}.$$

**Lemma 18** There exists a constant  $C$  such that

$$D^{[0,t]}(F(\Phi_1), F(\Phi_2)) \leq C \int_0^t D^{[0,s]}(\Phi_1, \Phi_2) ds.$$

**Proof** [Proof of Theorem 4] Firstly,  $\Phi$  contains the constant trajectory and thus is nonempty. Applying Lemma 18, the proof of existence and uniqueness follows from a similar argument of Picard–Lindelöf theorem. Specifically, iteratively applying Lemma 18 yields that

$$D^{[0,T]}(F^m(\Phi_1), F^m(\Phi_2)) \leq \frac{(CT)^m}{m!} D^{[0,T]}(\Phi_1, \Phi_2).$$

Let  $\Phi$  be the constant trajectory, for any  $\tilde{\Phi} \in \Phi$ , by the upper bounds of  $\underline{\mathcal{G}}_\ell^v$  in Lemma 21 and the Definition of  $D^{[0,T]}$  in Definition 17, there is a constant  $C$  such that

$$D^{[0,T]}(F(\tilde{\Phi}), \Phi) \leq CT < \infty.$$

We first show the uniqueness. For two fixed points of  $F$  denoted by  $\Phi_1$  and  $\Phi_2$ , we have

$$D^{[0,T]}(\Phi_1, \Phi_2) = D^{[0,T]}(F^m(\Phi_1), F^m(\Phi_2)) \leq \frac{(CT)^{m-1}}{(m-1)!} D^{[0,T]}(F(\Phi_1), F(\Phi_2)),$$

By the triangle inequality  $D^{[0,T]}(F(\Phi_1), F(\Phi_2)) \leq D^{[0,T]}(F(\Phi_1), \Phi) + D^{[0,T]}(F(\Phi_2), \Phi) < \infty$ , hence the right-hand side of the above inequality vanishes as  $m$  diverges. For the existence, we can consider the sequence  $\{F_2^i(\Phi) : i \geq 0\}$  that satisfies

$$D^{[0,T]}(F_2^{m+1}(\Phi), F_2^m(\Phi)) \leq \frac{(CT)^m}{m!} D^{[0,T]}(F(\Phi), \Phi),$$

Because  $D^{[0,T]}(F(\Phi), \Phi) < \infty$ ,  $\{F^i(\Phi) : i \geq 0\}$  is a Cauchy sequence. Since  $\Phi$  is complete under  $D^{[0,T]}$  by Lemma 22, the limit point  $\Phi_* \in \Phi$ , which is a fixed-point of  $F$ . Finally, by dominated convergence theorem, we can directly verify that  $\Phi_*$  is the solution of neural feature flow.  $\blacksquare$

### C.1.2. PROOF OF THEOREM 13

Theorem 13 is a Grönwall-type of result. However, it is not straightforward to directly derive a simple differential inequality due to the involved relations among the parameters of deep neural networks. Again we turn to the technique of Picard iterations used in the proof of Theorem 4. This approach has also been used to prove the abstract Grönwall inequality in Turinici (1986).

Recall the set  $\Phi$  in the proof of Theorem 4, and the mapping  $F : \Phi \mapsto \Phi$  in Definition 16. It is shown that  $F$  is a contraction mapping and thus there exists a unique solution  $\Phi_* \in \Phi$ . We will construct a closed nonempty subset  $\tilde{\Phi} \subseteq \Phi$  with the desired properties in Theorem 13 such that  $F(\tilde{\Phi}) \subseteq \tilde{\Phi}$ . Then by the same argument as the proof of Theorem 4, the Picard iteration guarantees the solution in  $\tilde{\Phi}$ , thereby proving  $\Phi_* \in \tilde{\Phi}$ .

We introduce the set of  $b$ -locally Lipschitz trajectories with the desired properties in Theorem 13. We use similar notations as in the proof of Theorem 4 by letting  $\mathbf{u}_\ell$  denote  $\Theta$ ,  $(\Theta, \Theta)$ ,  $\Theta$  for  $\ell = 1, 2 \leq \ell \leq L, \ell = L + 1$ , respectively.

**Definition 19 ( $b$ -Locally Lipschitz Trajectory)** Recall the constants  $C_2$  and  $C_4$  in Assumption 3 for the locally Lipschitz continuity at  $t = 0$ . We say  $\Phi$  is  $b$ -locally Lipschitz if for all  $t \in [0, T]$ ,  $\Theta_1$ ,

$\bar{\Theta}_1 \in \mathcal{B}_\infty(\Theta_1, 1)$ ,  $\bar{\Theta}_2$ , and  $\bar{\Theta}_2 \in \mathcal{B}_\infty(\bar{\Theta}_2, 1)$ , we have

$$\|\Phi_1^v(\Theta_1)(t) - \Phi_1^v(\bar{\Theta}_1)(t)\|_\infty \leq e^{bt}(\|\Theta_1\|_\infty + 1)\|\Theta_1 - \bar{\Theta}_1\|_\infty, \quad (44a)$$

$$|\Phi_{L+1}^v(\Theta_1)(t) - \Phi_{L+1}^v(\bar{\Theta}_1)(t)| \leq (1 + C_4)e^{bt}(\|\Theta_1\|_\infty + 1)\|\Theta_1 - \bar{\Theta}_1\|_\infty, \quad (44b)$$

$$|\Phi_\ell^v(\mathbf{u}_\ell)(t) - \Phi_\ell^v(\bar{\Theta}_1, \Theta_2)(t)| \leq (1 + C_2)e^{bt}(\|\mathbf{u}_\ell\|_\infty + 1)\|\Theta_1 - \bar{\Theta}_1\|_\infty, \quad (44c)$$

$$|\Phi_\ell^v(\mathbf{u}_\ell)(t) - \Phi_\ell^v(\Theta_1, \bar{\Theta}_2)(t)| \leq (1 + C_2)e^{bt}(\|\mathbf{u}_\ell\|_\infty + 1)\|\Theta_2 - \bar{\Theta}_2\|_\infty, \quad (44d)$$

for  $\ell \in [2 : L]$ . Denote the set of all  $b$ -locally Lipschitz trajectories as  $\Phi_b$ .

**Lemma 20** *There exists a constant  $b_*$  such that  $F(\Phi \cap \Phi_{b_*}) \subseteq \Phi_{b_*}$ .*

**Proof** [Proof of Theorem 13] Let  $b_*$  be the constant in Lemma 20 and  $\Phi' := \Phi \cap \Phi_{b_*} \subseteq \Phi$ , which clearly contains the constant trajectory and thus is nonempty. It follows from Lemma 20 that  $F(\Phi') \subseteq \Phi'$ . Since  $F$  is a contraction mapping by Lemma 18 and  $\Phi'$  is a closed set by Lemma 23, by the same argument as the proof of Theorem 4, there exists a unique solution in  $\Phi'$ , which is necessarily  $\Phi_*$  by the uniqueness of the solution in Theorem 4. ■

### C.1.3. PROOFS OF LEMMAS

**Proof** [Proof of Lemma 15] We first prove the Lipschitz continuity of  $\Phi$  for weight. It suffices to show upper bounds of  $\mathcal{G}_\ell^v$  for each layer  $\ell$ . We use the backward equations to inductively upper bound  $\underline{\mathcal{D}}_\ell^\beta$  and  $\underline{\mathcal{D}}_\ell^\alpha$ , which immediately yield upper bounds  $\|\underline{\mathcal{G}}_\ell^v\|_\infty \leq C_\ell$  for constants  $C_\ell$ .

For the top layer  $\ell = L + 1$ , by Assumption 1 that  $|\phi'_1| \leq L_4$ , we have

$$\|\underline{\mathcal{D}}_{L+1}(\Phi, t)\|_\infty \leq L_4 := \tilde{C}_{L+1}.$$

At layer  $\ell = L$ , since  $|h'_1| \leq L_2$ ,

$$\left\| \underline{\mathcal{D}}_L^\beta(\Theta; \Phi, t) \right\|_\infty \leq \underbrace{\left\| \dot{h}'_1(\beta_L^t(\Theta)) \right\|_\infty}_{\leq L_2} \underbrace{\left\| \underline{\mathcal{D}}_{L+1}(\Phi, t) \right\|_\infty}_{\leq \tilde{C}_{L+1}} \underbrace{|v_{L+1}^t(\Theta)|}_{\leq C_3 + T C_{L+1}} \leq \tilde{C}_L, \quad (45)$$

where  $\tilde{C}_L := (C_3 + C_{L+1}T)L_2\tilde{C}_{L+1}$  and  $|v_{L+1}^t| \leq C_3 + C_{L+1}T$  by the upper bound of initialization in Assumption 3 and the  $C_{L+1} := L_1\tilde{C}_{L+1}$ -Lipschitz continuity of  $v_{L+1}^t$  in  $t$ . For each  $\ell = L - 1, \dots, 1$ , suppose  $\underline{\mathcal{D}}_{\ell+1}^\beta$  is uniformly bounded by  $\tilde{C}_{\ell+1}$ . Then it follows from (11c) and (12b) that

$$\left\| \underline{\mathcal{D}}_{\ell+1}^\alpha(\Theta; \Phi, t) \right\|_\infty \leq \left\| \underline{\mathcal{D}}_{\ell+1}^\beta(\Theta; \Phi, t) \right\|_\infty \left\| \dot{h}'_2(\alpha_{\ell+1}^t(\Theta)) \right\|_\infty \leq \tilde{C}_{\ell+1}L_2 := \tilde{C}'_{\ell+1},$$

and

$$\left| \underline{\mathcal{G}}_{\ell+1}^v(\Theta, \bar{\Theta}; \Phi, t) \right| \leq \left\| \underline{\mathcal{D}}_{\ell+1}^\alpha(\Theta; \Phi, t) \right\|_\infty \left\| \dot{h}'_1(\beta_\ell^t(\Theta)) \right\|_\infty \leq L_1 C'_{\ell+1} := C_{\ell+1}.$$

We then similarly apply the upper bounds of initialization in Assumption 3 and the  $C_{\ell+1}$ -Lipschitz continuity of  $v_{\ell+1}^t$  in  $t$  and obtain that

$$\begin{aligned}
 & \left\| \int v_{\ell+1}^t(\Theta, \bar{\Theta}) \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi, t) dp(\bar{\Theta}) \right\|_\infty \\
 & \leq \tilde{C}'_{\ell+1} \int |v_{\ell+1}^t(\Theta, \bar{\Theta})| dp(\bar{\Theta}) \\
 & \leq \tilde{C}'_{\ell+1} \left( C_{\ell+1}t + \int |v_{\ell+1}(\Theta_\ell, \Theta_{\ell+1})| dp_{\ell+1}(\Theta_{\ell+1}) \right) \\
 & \leq \tilde{C}'_{\ell+1}(C_{\ell+1}T + C'),
 \end{aligned} \tag{46}$$

for a constant  $C'$ , where in the last inequality we used the upper bound of  $v_{\ell+1}$  in (15), the sub-gaussian property of  $p$ , and Corollary 31. Consequently,

$$\begin{aligned}
 & \left\| \underline{\mathcal{D}}_\ell^\beta(\Theta; \Phi, t) \right\|_\infty \\
 & \leq \left\| \underline{\mathcal{D}}_{\ell+1}^\beta(\Theta; \Phi, t) \right\|_\infty + \left\| \dot{h}'_1(\beta_\ell^t(\Theta)) \circ \int v_{\ell+1}^t(\Theta, \bar{\Theta}) \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi, t) dp(\bar{\Theta}) \right\|_\infty \\
 & \leq \tilde{C}_{\ell+1} + (\tilde{C}_{\ell+1}(C_{\ell+1}T + C'))L_2 := \tilde{C}_\ell.
 \end{aligned} \tag{47}$$

Finally, denoting the infinity norm of the data  $\mathbf{X}$  by  $C_x$ , (12c) gives the upper bound of  $|\underline{\mathcal{G}}_1^v(\Theta; \Phi, t)|$  by  $C_x \tilde{C}_1 := C_1$ .

Now we turn to the forward steps. We prove that there is a constant  $C$  such that for  $\ell \in [L]$  and  $\Theta$ ,

$$\left\| \underline{\mathcal{G}}_\ell^\beta(\Theta; \Phi, t) \right\|_\infty \leq C(\|\Theta\|_\infty + 1), \tag{48}$$

and for all  $\ell \in [2 : L]$  and  $\Theta$ ,

$$\left\| \underline{\mathcal{G}}_\ell^\alpha(\Theta; \Phi, t) \right\|_\infty \leq C(\|\Theta\|_\infty + 1). \tag{49}$$

Once we obtain (48) and (49), because  $p$  has bounded finite moment (Corollary 31), the dominated convergence theorem directly implies that  $\Phi^\alpha$  and  $\Phi^\beta$  satisfy the forward equations in Definition 14, which is our desired result.

For the first layer  $\ell = 1$ , since  $\mathbf{X}$  is bounded, it follows from (13a) that

$$\left\| \underline{\mathcal{G}}_1^\beta(\Theta; \Phi, t) \right\|_\infty \leq C''_1 \leq C'_1(\|\Theta\|_\infty + 1).$$

Suppose that at layer  $\ell \in [L - 1]$ , we have  $\left\| \underline{\mathcal{G}}_\ell^\beta(\Theta; \Phi, t) \right\|_\infty \leq C''_\ell(\|\Theta\|_\infty + 1)$ . By a similar argument to (47), we have

$$\begin{aligned}
 & \left\| \int v_{\ell+1}^t(\Theta, \bar{\Theta}) \underline{\mathcal{G}}_\ell^\beta(\bar{\Theta}; \Phi, t) dp(\bar{\Theta}) \right\|_\infty \\
 & \leq C''_\ell \int (C_{\ell+1}T + |v_{\ell+1}(\Theta, \bar{\Theta})|) (\|\bar{\Theta}\|_\infty + 1) dp(\bar{\Theta}) \\
 & \stackrel{(15)}{\leq} C''_\ell(C_{\ell+1}T + C_1(\|\Theta\|_\infty + 1)) \int (\|\bar{\Theta}\|_\infty + 1) dp(\bar{\Theta}) \\
 & \leq \tilde{C}''_{\ell+1}(\|\Theta\|_\infty + 1),
 \end{aligned}$$

for some constant  $\tilde{C}_{\ell+1}''$ . Therefore, applying (13b) yields that

$$\begin{aligned} \left\| \underline{\mathcal{G}}_{\ell+1}^{\alpha}(\Theta; \Phi, t) \right\|_{\infty} &\leq L_2 \tilde{C}_{\ell+1}'' (\|\Theta\|_{\infty} + 1) + \int \underbrace{\left\| \dot{h}_1(\beta_{\ell}^t) \right\|_{\infty}}_{\leq L_1} \underbrace{\left| \underline{\mathcal{G}}_{\ell+1}^v(\Theta, \bar{\Theta}; \Phi, t) \right|}_{\leq C_{\ell}} dp(\bar{\Theta}) \\ &\leq \bar{C}_{\ell+1} (\|\Theta\|_{\infty} + 1), \end{aligned} \quad (50)$$

for some constant  $\bar{C}_{\ell+1}$ . We obtain

$$\begin{aligned} &\left\| \underline{\mathcal{G}}_{\ell+1}^{\beta}(\Theta; \Phi, t) \right\|_{\infty} \\ &\leq \left\| \dot{h}'_2(\alpha_{\ell+1}^t(\Theta)) \circ \underline{\mathcal{G}}_{\ell+1}^{\alpha}(\Theta; \Phi, t) \right\|_{\infty} + \left\| \underline{\mathcal{G}}_{\ell}^{\beta}(\Theta; \Phi, t) \right\|_{\infty} \\ &\leq (L_2 \bar{C}_{\ell+1} + C_{\ell}'') (\|\Theta\|_{\infty} + 1), \end{aligned}$$

which suggests to pick  $C_{\ell+1}'' = (L_2 \bar{C}_{\ell+1} + C_{\ell}'')$ . We achieve Lemma 15.  $\blacksquare$

Before proving Lemma 18, we first present in Lemma 21 properties of  $\Phi \in \Phi$  that will be used to prove the contraction lemma. The proof is exactly the same as Lemma 15 and is omitted.

**Lemma 21 (Property of  $\Phi$ )** *There exists a generic constant  $C$  such that, for any  $\Phi \in \Phi$ , we have*

- $\|\underline{\mathcal{D}}_{L+1}(\Phi, t)\|_{\infty} \leq C$  and  $\|\underline{\mathcal{D}}_{\ell}^{\beta}(\Theta; \Phi, t)\|_{\infty} \leq C$  for  $\ell \in [L]$ ;
- $\|\underline{\mathcal{D}}_{\ell}^{\alpha}(\Theta; \Phi, t)\|_{\infty} \leq C$  for  $\ell \in [2 : L]$ ;
- $\|\underline{\mathcal{G}}_{\ell}^v(\mathbf{u}_{\ell}; \Phi, t)\|_{\infty} \leq C$  and  $\|v_{\ell}^t(\mathbf{u}_{\ell})\|_{\infty} \leq \|v_{\ell}^0(\mathbf{u}_{\ell})\|_{\infty} + C t$  for  $\ell \in [L + 1]$ ;
- $\|\underline{\mathcal{G}}_{\ell}^{\beta}(\Theta; \Phi, t)\|_{\infty} \leq C (\|\Theta\|_{\infty} + 1)$  for  $\ell \in [L]$ ;
- $\|\underline{\mathcal{G}}_{\ell}^{\alpha}(\Theta; \Phi, t)\|_{\infty} \leq C (\|\Theta\|_{\infty} + 1)$  for  $\ell \in [2 : L]$ .

**Proof** [Proof of Lemma 18] The proof entails upper bounds of the gradient differences  $\|\underline{\mathcal{G}}_{\ell}^v(\mathbf{u}_{\ell}; \Phi_1, t) - \underline{\mathcal{G}}_{\ell}^v(\mathbf{u}_{\ell}; \Phi_2, t)\|_{\infty}$  in terms of the differences  $|v_{1,\ell}^t - v_{2,\ell}^t|$  for  $\ell \in [L + 1]$ , which can be further upper bounded in terms of  $d_t := D^{[0,t]}(\Phi_1, \Phi_2)$ , that is by Definition 17:

$$\|v_{1,\ell}^t(\mathbf{u}_{\ell}) - v_{2,\ell}^t(\mathbf{u}_{\ell})\|_{\infty} \leq (\|\mathbf{u}_{\ell}\|_{\infty} + 1) d_t, \quad \ell \in [L + 1]. \quad (51)$$

We will use the forward equations to inductively upper bound the differences between  $\beta_{\ell}$  and  $\alpha_{\ell}$ . Especially, for some constant  $C$ , we prove

$$\|\beta_{1,\ell}^t(\Theta) - \beta_{2,\ell}^t(\Theta)\|_{\infty} \leq C (\|\Theta\|_{\infty} + 1) d_t, \quad \ell \in [L], \quad (52)$$

$$\|\alpha_{1,\ell}^t(\Theta) - \alpha_{2,\ell}^t(\Theta)\|_{\infty} \leq C (\|\Theta\|_{\infty} + 1) d_t, \quad \ell \in [2 : L]. \quad (53)$$

We then use the backward equations to upper bound the difference between  $\underline{\mathcal{D}}_{\ell}^{\beta}$ ,  $\underline{\mathcal{D}}_{\ell}^{\alpha}$ ,  $\underline{\mathcal{G}}_{\ell}^v$ . Namely, we prove

$$\left\| \underline{\mathcal{D}}_{\ell}^{\beta}(\Theta; \Phi_1, t) - \underline{\mathcal{D}}_{\ell}^{\beta}(\Theta; \Phi_2, t) \right\|_{\infty} \leq C(1 + \|\Theta\|_{\infty}) d_t, \quad \ell \in [L], \quad (54)$$

$$\left\| \underline{\mathcal{D}}_{\ell}^{\alpha}(\Theta; \Phi_1, t) - \underline{\mathcal{D}}_{\ell}^{\alpha}(\Theta; \Phi_2, t) \right\|_{\infty} \leq C(1 + \|\Theta\|_{\infty}) d_t, \quad \ell \in [2 : L], \quad (55)$$

$$\left\| \underline{\mathcal{G}}_{\ell}^v(\mathbf{u}_{\ell}; \Phi_1, t) - \underline{\mathcal{G}}_{\ell}^v(\mathbf{u}_{\ell}; \Phi_2, t) \right\|_{\infty} \leq C(1 + \|\mathbf{u}_{\ell}\|_{\infty}) d_t, \quad \ell \in [L + 1], \quad (56)$$

Then the conclusion, i.e., Lemma 18 directly follows from (56), and the definitions of  $F$  and  $D^{[0,t]}$  in Definitions 16 and 17, respectively.

We consider forward steps. When  $\ell = 1$ , because  $\mathbf{X}$  is bounded, by the Definition 16 (2) and (51), we have  $\|\beta_{1,1}^t(\Theta) - \beta_{2,1}^t(\Theta)\|_\infty \leq C(\|\Theta\|_\infty + 1)d_t$ . Consider at layer  $\ell \in [L-1]$ , we obtain  $\|\beta_{1,\ell}^t(\Theta) - \beta_{2,\ell}^t(\Theta)\|_\infty \leq C(\|\Theta\|_\infty + 1)d_t$ . For layer  $\ell + 1$ , by the boundedness and Lipschitz continuity of  $h_1$ , we have

$$\begin{aligned}
 & \|\alpha_{1,\ell+1}^t(\Theta) - \alpha_{2,\ell+1}^t(\Theta)\|_\infty \\
 \leq & \left\| \int \dot{h}_1(\beta_{1,\ell}^t(\bar{\Theta})) v_{1,\ell+1}^t(\bar{\Theta}, \Theta) - \dot{h}_1(\beta_{2,\ell}^t(\bar{\Theta})) v_{2,\ell+1}^t(\bar{\Theta}, \Theta) dp(\bar{\Theta}) \right\|_\infty \\
 \leq & \int \underbrace{\left\| \dot{h}_1(\beta_{1,\ell}^t(\bar{\Theta})) - \dot{h}_1(\beta_{2,\ell}^t(\bar{\Theta})) \right\|_\infty}_{\leq L_2 C(\|\bar{\Theta}\|_\infty + 1)d_t} |v_{1,\ell+1}^t(\bar{\Theta}, \Theta)| dp(\bar{\Theta}) \\
 & \quad + \int \underbrace{\left\| \dot{h}_1(\beta_{2,\ell}^t(\bar{\Theta})) \right\|_\infty}_{\leq L_1} \underbrace{|v_{1,\ell+1}^t(\bar{\Theta}, \Theta) - v_{2,\ell+1}^t(\bar{\Theta}, \Theta)|}_{\leq (\|\bar{\Theta}\|_\infty + \|\Theta\|_\infty + 1)d_t} dp(\bar{\Theta}) \\
 \leq & C'(\|\Theta\|_\infty + 1)d_t,
 \end{aligned}$$

for a constant  $C'$ , where the last step is due to the sub-gaussianness of  $p$ , Corollary 31, and the upper bound of  $v_{\ell+1}^t$  in Lemma 21. Consequently, the Lipschitz continuity of  $L_2$  gives that

$$\begin{aligned}
 & \|\beta_{1,\ell+1}^t(\Theta) - \beta_{2,\ell+1}^t(\Theta)\|_\infty \\
 \leq & \left\| \dot{h}_2(\alpha_{1,\ell+1}^t(\Theta)) - \dot{h}_2(\alpha_{2,\ell+1}^t(\Theta)) \right\|_\infty + \|\beta_{1,\ell}^t(\Theta) - \beta_{2,\ell}^t(\Theta)\|_\infty \\
 \leq & C''(\|\Theta\|_\infty + 1)d_t
 \end{aligned}$$

for a constant  $C''$ . We achieve (52) and (53).

We turn to the backward steps. We focus on the upper bound of the difference between  $\underline{\mathcal{D}}_\ell^\alpha$  and  $\underline{\mathcal{D}}_\ell^\beta$ . Since both  $h_1$  and  $\mathbf{X}$  are bounded,  $h_1$  is Lipschitz continuous, and  $\underline{\mathcal{D}}_\ell^\alpha$  is bounded by Lemma 21, the upper bound of the difference between  $\underline{\mathcal{G}}_\ell^\nu$  follows immediately. To begin with, we introduce

$$\underline{\mathcal{D}}_\ell^\gamma(\Theta; \Phi, t) := \int v_\ell^t(\Theta, \bar{\Theta}) \underline{\mathcal{D}}_\ell^\alpha(\bar{\Theta}; \Phi, t) dp(\bar{\Theta}) \circ \dot{h}_1(\beta_{\ell-1}^t(\Theta)), \quad \ell \in [2:L]. \quad (57)$$

Then we have

$$\underline{\mathcal{D}}_{\ell-1}^\beta(\Theta; \Phi, t) = \underline{\mathcal{D}}_\ell^\beta(\Theta; \Phi, t) + \underline{\mathcal{D}}_\ell^\gamma(\Theta; \Phi, t), \quad \ell \in [2:L].$$

For the top layer  $\ell = L+1$ , the Lipschitz continuity of  $\phi_1^t$  in Assumption 1 implies that,

$$\begin{aligned}
 & \|\underline{\mathcal{D}}_{L+1}(\Phi_1, t) - \underline{\mathcal{D}}_{L+1}(\Phi_2, t)\|_\infty \\
 \leq & L_5 \|\beta_{1,L+1}^t - \beta_{2,L+1}^t\|_\infty \\
 \leq & L_5 \int \|h_1(\beta_{1,L}^t) v_{1,L+1}^t(\Theta) - h_1(\beta_{2,L}^t) v_{2,L+1}^t(\Theta)\|_\infty dp(\Theta).
 \end{aligned} \quad (58)$$



Since  $h_1$  is bounded and Lipschitz continuous,  $v_{i,L+1}^t$  is bounded for  $t \leq T$  by Lemma 21, and (52), we have

$$\|\underline{\mathcal{D}}_{L+1}(\Phi_1, t) - \underline{\mathcal{D}}_{L+1}(\Phi_2, t)\|_\infty \leq C d_t,$$

for a constant  $C$ . At layer  $\ell = L$ , recall that

$$\underline{\mathcal{D}}_L(\Theta; \Phi, t) = v_{L+1}^t(\Theta) \underline{\mathcal{D}}_{L+1}(\Phi, t) \circ h_1'(\beta_L^t).$$

Since the three terms in the product are all bounded, and  $h_1'$  is  $L_3$ -Lipschitz continuous, we have

$$\|\underline{\mathcal{D}}_L^\beta(\Theta; \Phi_1, t) - \underline{\mathcal{D}}_L^\beta(\Theta; \Phi_2, t)\|_\infty \leq C'(1 + \|\Theta\|_\infty) d_t, \quad (59)$$

for a constant  $C'$ . For each  $\ell = L - 1, \dots, 1$ , suppose there is a constant  $C$ , such that

$$\left\| \underline{\mathcal{D}}_{\ell+1}^\beta(\Theta; \Phi_1, t) - \underline{\mathcal{D}}_{\ell+1}^\beta(\Theta; \Phi_2, t) \right\|_\infty \leq C (\|\Theta\|_\infty + 1) d_t.$$

As  $h_2'$  is  $L_3$ -Lipschitz continuous, the boundednesses of  $h_2$  and  $\underline{\mathcal{D}}_{\ell+1}^\beta$  in Lemma 21 and (53) give that

$$\left\| \underline{\mathcal{D}}_{\ell+1}^\alpha(\Theta; \Phi_1, t) - \underline{\mathcal{D}}_{\ell+1}^\alpha(\Theta; \Phi_2, t) \right\|_\infty \leq C' (\|\Theta\|_\infty + 1) d_t,$$

for a constant  $C'$ . Moreover,

$$\begin{aligned} & \int \|v_{1,\ell+1}^t(\Theta, \bar{\Theta}) \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi_1, t) - v_{2,\ell+1}^t(\Theta, \bar{\Theta}) \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi_2, t)\|_\infty dp(\bar{\Theta}) \\ & \leq \int \underbrace{|v_{1,\ell+1}^t(\Theta, \bar{\Theta}) - v_{2,\ell+1}^t(\Theta, \bar{\Theta})|}_{\leq (\|\Theta\|_\infty + \|\bar{\Theta}\|_\infty + 1) d_t} \left\| \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi_1, t) \right\|_\infty \\ & \quad + \underbrace{|v_{2,\ell+1}^t(\Theta, \bar{\Theta})|}_{\leq C'(\|\bar{\Theta}\|_\infty + 1) d_t} \left\| \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi_1, t) - \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi_2, t) \right\|_\infty dp(\bar{\Theta}) \\ & \leq C'' (\|\Theta\|_\infty + 1) d_t, \end{aligned} \quad (60)$$

for a constant  $C''$ , where the last step is due to the sub-gaussianness of  $p$ , Corollary 31, and the upper bounds of  $\underline{\mathcal{D}}_{\ell+1}^\alpha$  and  $v_{\ell+1}^t$  in Lemma 21. Consequently,

$$\left\| \underline{\mathcal{D}}_{\ell+1}^\gamma(\Theta; \Phi_1, t) - \underline{\mathcal{D}}_{\ell+1}^\gamma(\Theta; \Phi_2, t) \right\|_\infty \leq \bar{C} (\|\Theta\|_\infty + 1) d_t,$$

for a constant  $\bar{C}$ , which further implies

$$\begin{aligned} & \left\| \underline{\mathcal{D}}_\ell^\beta(\Theta; \Phi_1, t) - \underline{\mathcal{D}}_\ell^\beta(\Theta; \Phi_2, t) \right\|_\infty \\ & \leq \left\| \underline{\mathcal{D}}_{\ell+1}^\beta(\Theta; \Phi_1, t) - \underline{\mathcal{D}}_{\ell+1}^\beta(\Theta; \Phi_2, t) \right\|_\infty + \left\| \underline{\mathcal{D}}_{\ell+1}^\gamma(\Theta; \Phi_1, t) - \underline{\mathcal{D}}_{\ell+1}^\gamma(\Theta; \Phi_2, t) \right\|_\infty \\ & \leq \bar{C}' (\|\Theta\|_\infty + 1) d_t, \end{aligned}$$

for a constant  $\bar{C}'$ . We finish the proof. ■

**Lemma 22**  $\Phi$  is complete under  $D^{[0,T]}$ .

**Proof** Let  $\{\Phi_n : n \geq 0\}$  be a Cauchy sequence under  $D^{[0,T]}$ . Then  $\frac{\Phi_{n,\ell}^v(\mathbf{u}_\ell)(t)}{1+\|\mathbf{u}_\ell\|_\infty}$  converges uniformly under the  $\ell_\infty$ -norm. Let  $\Phi_{*,\ell}^v(\mathbf{u}_\ell)(t) = \lim_{n \rightarrow \infty} \Phi_{n,\ell}^v(\mathbf{u}_\ell)(t)$  for  $\ell \in [L+1]$ . Since the Lipschitz continuity is preserved under the pointwise convergence, we have  $\Phi_{*,\ell}^v$  is  $C$ -Lipschitz continuous in  $t$ . Let

$$\begin{aligned}\Phi_{*,1}^\beta(\Theta)(t) &= \frac{1}{d} \mathbf{X} \Phi_{*,1}^v(\Theta)(t), \\ \Phi_{*,\ell+1}^\alpha(\Theta)(t) &= \int \Phi_{*,\ell+1}^v(\Theta, \bar{\Theta})(t) \dot{h}_1 \left( \Phi_{*,\ell}^\beta(\bar{\Theta})(t) \right) dp(\bar{\Theta}), \quad \ell \in [L], \\ \Phi_{*,\ell+1}^\beta(\Theta)(t) &= \Phi_{*,\ell}^\beta(\Theta)(t) + \dot{h}_2 \left( \Phi_{*,\ell+1}^\alpha(\Theta)(t) \right), \quad \ell \in [L].\end{aligned}$$

By the dominated convergence theorem, we have  $\Phi_{*,\ell}^\beta(\Theta)(t) = \lim_{n \rightarrow \infty} \Phi_{n,\ell}^\beta(\Theta)(t)$  and  $\Phi_{*,\ell}^\alpha(\Theta)(t) = \lim_{n \rightarrow \infty} \Phi_{n,\ell}^\alpha(\Theta)(t)$ . Then  $\Phi_*$  is a limit point of  $\{\Phi_n : n \geq 0\}$  under  $D^{[0,T]}$  and  $\Phi_* \in \Phi$ .  $\blacksquare$

Next we prove lemmas for Theorem 13.

**Proof** [Proof of Lemma 20] Analogous to the notation of  $\mathbf{u}_\ell$ , for the convenience of presenting continuity of  $\Phi_\ell^v$ , we introduce notations  $\bar{\mathbf{u}}$  and  $\bar{\mathbf{u}}'_\ell$  by letting

$$\bar{\mathbf{u}}_\ell = \begin{cases} \bar{\Theta}_1, \\ (\bar{\Theta}_1, \Theta_2), \\ \bar{\Theta}_1 \end{cases}, \quad \bar{\mathbf{u}}'_\ell = \begin{cases} \bar{\Theta}_1, & \ell = 1, \\ (\Theta_1, \bar{\Theta}_2), & \ell \in [2 : L], \\ \bar{\Theta}_1 & \ell = L + 1. \end{cases}$$

We also abbreviate the notations for the individual trajectories as:

$$\Phi_\ell^v(\mathbf{u}_\ell)(t) = v_\ell^t(\mathbf{u}_\ell), \quad \Phi_\ell^v(\bar{\mathbf{u}}_\ell)(t) = v_\ell^t(\bar{\mathbf{u}}_\ell), \quad \Phi_\ell^v(\bar{\mathbf{u}}'_\ell)(t) = v_\ell^t(\bar{\mathbf{u}}'_\ell),$$

$\Phi_{\ell_1}^\beta(\Theta_1)(t) = \beta_{\ell_1}^t(\Theta_1)$ ,  $\Phi_{\ell_1}^\beta(\bar{\Theta}_1)(t) = \beta_{\ell_1}^t(\bar{\Theta}_1)$ ,  $\Phi_{\ell_2}^\alpha(\Theta_1)(t) = \alpha_{\ell_2}^t(\Theta_1)$ ,  $\Phi_{\ell_2}^\alpha(\bar{\Theta}_1)(t) = \alpha_{\ell_2}^t(\bar{\Theta}_1)$ , for  $\ell \in [L+1]$ ,  $\ell_1 \in [L]$ , and  $\ell_2 \in [2 : L]$ , respectively.

We first investigate the set  $F(\Phi \cap \Phi_b)$  for a general  $b$ . We follow similar steps as the proof of Lemma 18. We first consider forward steps and inductively show upper bound for the differences between  $\beta_\ell$  and  $\alpha_\ell$ . Namely, we prove for any  $\Phi \in \Phi \cap \Phi_b$ ,

$$\|\beta_\ell^t(\Theta_1) - \beta_\ell^t(\bar{\Theta}_1)\|_\infty \leq C e^{bt} (\|\Theta_1\|_\infty + 1) \|\Theta_1 - \bar{\Theta}_1\|_\infty, \quad \ell \in [L], \quad (61)$$

$$\|\alpha_\ell^t(\Theta_1) - \alpha_\ell^t(\bar{\Theta}_1)\|_\infty \leq C e^{bt} (\|\Theta_1\|_\infty + 1) \|\Theta_1 - \bar{\Theta}_1\|_\infty, \quad \ell \in [2 : L], \quad (62)$$

for a constant  $C$ . Then we study the backward steps, and prove that there is a constant  $\tilde{C}$  independent of  $b$  such that for any  $\Phi \in \Phi \cap \Phi_b$ ,

$$\|\underline{\mathcal{D}}_\ell^\beta(\Theta_1; \Phi, t) - \underline{\mathcal{D}}_\ell^\beta(\bar{\Theta}_1; \Phi, t)\|_\infty \leq \tilde{C} e^{bt} (\|\Theta_1\|_\infty + 1) \|\Theta_1 - \bar{\Theta}_1\|_\infty, \quad \ell \in [L], \quad (63)$$

$$\|\underline{\mathcal{D}}_\ell^\alpha(\Theta_1; \Phi, t) - \underline{\mathcal{D}}_\ell^\alpha(\bar{\Theta}_1; \Phi, t)\|_\infty \leq \tilde{C} e^{bt} (\|\Theta_1\|_\infty + 1) \|\Theta_1 - \bar{\Theta}_1\|_\infty, \quad \ell \in [2 : L], \quad (64)$$

$$\|\underline{\mathcal{G}}_\ell^v(\mathbf{u}_\ell; \Phi, t) - \underline{\mathcal{G}}_\ell^v(\bar{\mathbf{u}}_\ell; \Phi, t)\|_\infty \leq \tilde{C} e^{bt} (1 + \|\mathbf{u}_\ell\|_\infty) \|\mathbf{u}_\ell - \bar{\mathbf{u}}_\ell\|_\infty, \quad \ell \in [L+1], \quad (65)$$

$$\|\underline{\mathcal{G}}_\ell^v(\mathbf{u}_\ell; \Phi, t) - \underline{\mathcal{G}}_\ell^v(\bar{\mathbf{u}}'_\ell; \Phi, t)\|_\infty \leq \tilde{C} e^{bt} (1 + \|\mathbf{u}_\ell\|_\infty) \|\mathbf{u}_\ell - \bar{\mathbf{u}}'_\ell\|_\infty, \quad \ell \in [L+1]. \quad (66)$$

In forward steps, we prove (61) and (62). For the 1-st layer, because  $\mathbf{X}$  is bounded, applying (44a) with the formula of  $\beta_1^t$  as Definition 16 (2) yields (61) with  $\ell = 1$ . Suppose at layer  $\ell \in [L-1]$ , we have

$$\|\beta_\ell^t(\Theta_1) - \beta_\ell^t(\bar{\Theta}_1)\|_\infty \leq C' e^{bt} (\|\Theta_1\|_\infty + 1) \|\Theta_1 - \bar{\Theta}_1\|_\infty$$

for a constant  $C'$ . It follows that

$$\begin{aligned} \|\alpha_{\ell+1}^t(\Theta_2) - \alpha_{\ell+1}^t(\bar{\Theta}_2)\|_\infty &\leq \int \underbrace{\|\dot{h}_1(\beta_\ell^t(\Theta_1))\|_\infty}_{\leq L_1} \underbrace{|v_{\ell+1}^t(\Theta_1, \Theta_2) - v_{\ell+1}^t(\Theta_1, \bar{\Theta}_2)|}_{(1+C_2)e^{bt}(\|\Theta_1\|_\infty + \|\Theta_2\|_\infty + 1)\|\Theta_2 - \bar{\Theta}_2\|_\infty} dp(\Theta_1) \\ &\leq C'' e^{bt}(\|\Theta_2\|_\infty + 1)\|\Theta_2 - \bar{\Theta}_2\|_\infty, \end{aligned}$$

for a constant  $C''$ , where we use sub-gaussianness of  $p$  and Corollary 31. We conclude that

$$\begin{aligned} \|\beta_{\ell+1}^t(\Theta_1) - \beta_{\ell+1}^t(\bar{\Theta}_1)\|_\infty &\leq \left\| \dot{h}_2(\alpha_{\ell+1}^t(\Theta_1)) - \dot{h}_2(\alpha_{\ell+1}^t(\bar{\Theta}_1)) \right\|_\infty + \|\beta_\ell^t(\Theta_1) - \beta_\ell^t(\bar{\Theta}_1)\|_\infty \\ &\leq C''' e^{bt}(\|\Theta_1\|_\infty + 1)\|\Theta_1 - \bar{\Theta}_1\|_\infty. \end{aligned}$$

Therefore, (61) and (62) is achieved.

We turn to backward process. Again, we focus on the difference between  $\underline{\mathcal{D}}_\ell^\beta$  and  $\underline{\mathcal{D}}_\ell^\alpha$ , i.e., (63) and (64). Then the upper bound for the difference between  $\underline{\mathcal{G}}_\ell^v$ , i.e., (65) and (66) follows immediately. For example, for the top layer  $\ell = L + 1$ , because of the boundedness of  $h_1$ ,  $\underline{\mathcal{D}}_{L+1}$ , applying (63) with the formula of  $\underline{\mathcal{G}}_{L+1}^v$  (12a) gives that

$$\left| \underline{\mathcal{G}}_{L+1}^v(\Theta_1; \Phi, t) - \underline{\mathcal{G}}_{L+1}^v(\bar{\Theta}_1; \Phi, t) \right| \leq C e^{bt}(\|\Theta_1\|_\infty + 1)\|\Theta_1 - \bar{\Theta}_1\|_\infty.$$

Other layers can be analogously obtained. At layer  $\ell = L$ , recall that

$$\underline{\mathcal{D}}_L(\Theta; \Phi, t) = v_{L+1}^t(\Theta) \underline{\mathcal{D}}_{L+1}(\Phi, t) \circ \dot{h}'_1(\beta_L^t).$$

For,  $\Phi \in \Phi_b$ , the upper bound for  $|v_{L+1}^t(\Theta_1) - v_{L+1}^t(\bar{\Theta}_1)|$  is given in (44b). Then applying the Lipschitz continuity of  $h'_1$ , (63), and the boundednesses of  $v_{L+1}^t$ ,  $\underline{\mathcal{D}}_{L+1}$  in Lemma 21 and  $h'_1$  yields

$$\left\| \underline{\mathcal{D}}_L^\beta(\Theta_1; \Phi, t) - \underline{\mathcal{D}}_L^\beta(\bar{\Theta}_1; \Phi, t) \right\|_\infty \leq C' e^{bt}(\|\Theta_1\|_\infty + 1)\|\Theta_1 - \bar{\Theta}_1\|_\infty, \quad (67)$$

for a constant  $C'$ .

For each  $\ell = L - 1, L - 2, \dots, 1$ , suppose we have

$$\left\| \underline{\mathcal{D}}_{\ell+1}^\beta(\Theta_1; \Phi, t) - \underline{\mathcal{D}}_{\ell+1}^\beta(\bar{\Theta}_1; \Phi, t) \right\|_\infty \leq C e^{bt}(\|\Theta_1\|_\infty + 1)\|\Theta_1 - \bar{\Theta}_1\|_\infty,$$

for a constant  $C$ . Combining the Lipschitz continuity of  $h'_2$  with the boundednesses of  $h'_2$ ,  $\underline{\mathcal{D}}_{\ell+1}^\beta$  in Lemma 21 gives that

$$\begin{aligned} &\left\| \underline{\mathcal{D}}_{\ell+1}^\alpha(\Theta_1; \Phi, t) - \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}_1; \Phi, t) \right\|_\infty \\ &= \left\| \dot{h}'_2(\alpha_{\ell+1}^t(\Theta_1)) \circ \underline{\mathcal{D}}_{\ell+1}^\beta(\Theta_1; \Phi, t) - \dot{h}'_2(\alpha_{\ell+1}^t(\bar{\Theta}_1)) \circ \underline{\mathcal{D}}_{\ell+1}^\beta(\bar{\Theta}_1; \Phi, t) \right\|_\infty \\ &\leq C' e^{bt}(\|\Theta_1\|_\infty + 1)\|\Theta_1 - \bar{\Theta}_1\|_\infty, \end{aligned}$$

for a constant  $C'$ . Moreover,

$$\begin{aligned} &\int \underbrace{|v_{\ell+1}^t(\Theta_1, \Theta_2) - v_{\ell+1}^t(\bar{\Theta}_1, \Theta_2)|}_{\leq C e^{bt}(\|\Theta_1\|_\infty + \|\Theta_2\|_\infty + 1)\|\Theta_1 - \bar{\Theta}_1\|_\infty} \left\| \underline{\mathcal{D}}_{\ell+1}^\alpha(\Theta_2; \Phi, t) \right\|_\infty dp(\Theta_2) \\ &\leq C' e^{bt}(\|\Theta_1\|_\infty + 1)\|\Theta_1 - \bar{\Theta}_1\|_\infty, \end{aligned} \quad (68)$$

where in the last step we used the sub-gaussianness of  $p$ , Corollary 31, and upper bound of  $\overline{\mathcal{D}}_{\ell+1}^\alpha$  in Lemma 21. Then, by the upper bound in (46), boundedness and Lipschitz continuity of  $h'_1$ , we obtain from (57) that

$$\|\underline{\mathcal{D}}_{\ell+1}^\gamma(\Theta_1; \Phi, t) - \underline{\mathcal{D}}_{\ell+1}^\gamma(\bar{\Theta}_1; \Phi, t)\|_\infty \leq C'' e^{bt} (\|\Theta_1\|_\infty + 1) \|\Theta_1 - \bar{\Theta}_1\|_\infty,$$

for a constant  $C''$ , which further yields

$$\|\underline{\mathcal{D}}_\ell^\beta(\Theta_1; \Phi, t) - \underline{\mathcal{D}}_\ell^\beta(\bar{\Theta}_1; \Phi, t)\|_\infty \leq C''' e^{bt} (\|\Theta_1\|_\infty + 1) \|\Theta_1 - \bar{\Theta}_1\|_\infty.$$

We achieve (63) and (64).

Finally, let  $b_* = \tilde{C}$  in (65) and (66). It remains to verify that  $F(\Phi) \in \Phi_{b_*}$  for any  $\Phi \in \Phi \cap \Phi_{b_*}$ , that is, to verify the conditions in Definition 19. For  $F(\Phi)_1^v$ , we have

$$\begin{aligned} & \|F(\Phi)_1^v(\Theta_1)(t) - F(\Phi)_1^v(\bar{\Theta}_1)(t)\|_\infty \\ & \leq \|\Theta_1 - \bar{\Theta}_1\|_\infty + \int_0^t \|\underline{\mathcal{G}}_1^v(\Theta_1; \Phi, s) - \underline{\mathcal{G}}_1^v(\bar{\Theta}_1; \Phi, s)\|_\infty ds \\ & \leq \|\Theta_1 - \bar{\Theta}_1\|_\infty + \int_0^t \tilde{C} (\|\Theta_1\|_\infty + 1) e^{b_* s} \|\Theta_1 - \bar{\Theta}_1\|_\infty ds \\ & \leq e^{b_* t} (\|\Theta_1\|_\infty + 1) \|\Theta_1 - \bar{\Theta}_1\|_\infty. \end{aligned} \quad (69)$$

The verification of other cases are entirely analogous and is omitted. Thus we obtain Lemma 20. ■

**Lemma 23**  $\Phi \cap \Phi_b$  is a closed set.

**Proof** Given a convergent sequence  $\{\Phi_n : n \geq 0\} \subseteq \Phi \cap \Phi_b$ , it follows from Lemma 22 that the limit point  $\Phi_* \in \Phi$ . Since Lipschitz property is preserved under pointwise convergence, we also have  $\Phi_* \in \Phi_b$ . ■

## C.2. Proof of Theorem 6

In the proof, we fix  $\Phi_*$  and the initialization  $\{\bar{\Theta}_i\}_{i=1}^m$ . We first write down the exact formula for the gradient of the actual discrete DNN, i.e.  $\hat{\mathcal{G}}_{\ell,i,j}^k$ . Especially, define intermediate variables in the back-propagation as

$$\begin{aligned} \hat{\mathcal{D}}_{L+1,1}^k & := N \frac{\partial \hat{\mathcal{L}}^k}{\partial \hat{\beta}_{L+1}} = \left[ \phi'_1 \left( \hat{\beta}_{L+1}^k(1), y^1 \right), \phi'_1 \left( \hat{\beta}_{L+1}^k(2), y^2 \right), \dots, \phi'_1 \left( \hat{\beta}_{L+1}^k(N), y^N \right) \right]^\top, \\ \hat{\mathcal{D}}_{L,i}^{\beta,k} & := N \frac{\partial \hat{\mathcal{L}}^k}{\partial \hat{\beta}_{L,i}} = \frac{1}{m} \left[ \hat{v}_{L+1,i,1}^k \hat{\mathcal{D}}_{L+1,1}^k \right] \circ \dot{h}'_1 \left( \hat{\beta}_{L,i}^k \right), \quad i \in [m], \\ \hat{\mathcal{D}}_{L,i}^{\alpha,k} & := N \frac{\partial \hat{\mathcal{L}}^k}{\partial \hat{\alpha}_{L,i}} = \hat{\mathcal{D}}_{L,i}^{\beta,k} \circ \dot{h}'_2 \left( \hat{\alpha}_{L,i}^k \right), \quad i \in [m], \\ \hat{\mathcal{D}}_{\ell,i}^{\beta,k} & := N \frac{\partial \hat{\mathcal{L}}^k}{\partial \hat{\beta}_{\ell,i}} = \frac{1}{m} \left[ \sum_{j=1}^m \hat{v}_{\ell+1,i,j}^k \hat{\mathcal{D}}_{\ell+1,j}^{\alpha,k} \right] \circ \dot{h}'_1 \left( \hat{\beta}_{\ell,i}^k \right) + \hat{\mathcal{D}}_{\ell+1,i}^{\beta,k}, \quad \ell \in [L-1], i \in [m], \\ \hat{\mathcal{D}}_{\ell,i}^{\alpha,k} & := N \frac{\partial \hat{\mathcal{L}}^k}{\partial \hat{\alpha}_{\ell,i}} = \hat{\mathcal{D}}_{\ell,i}^{\beta,k} \circ \dot{h}'_2 \left( \hat{\alpha}_{\ell,i}^k \right), \quad \ell \in [2 : L-1], i \in [m]. \end{aligned}$$

We have

$$\begin{aligned}\hat{\mathcal{G}}_{L+1,i,1}^k &= \frac{1}{Nm} \left[ \hat{\mathcal{D}}_{L+1}^k \right]^\top \dot{h}_1 \left( \hat{\beta}_{\ell,i}^k \right), \quad i \in [m], \\ \hat{\mathcal{G}}_{\ell+1,i,j}^k &= \frac{1}{Nm} \left[ \hat{\mathcal{D}}_{\ell+1,j}^{\alpha,k} \right]^\top \dot{h}_1 \left( \hat{\beta}_{\ell,i}^k \right), \quad \ell \in [L-1], i, j \in [m], \\ \hat{\mathcal{G}}_{1,i,j}^k &= \frac{1}{Nd} \left[ \hat{\mathcal{D}}_{1,j}^{\beta,k} \right]^\top \hat{\beta}_{0,i}, \quad i \in [d], j \in [m].\end{aligned}$$

To compare the discrete and continuous trajectories on the same time scale, we normalize discrete gradients by

$$\widehat{\mathcal{N}}\mathcal{D}_{\ell,i}^{\alpha,k} = [m] \widehat{\mathcal{D}}_{\ell,i}^{\alpha,k}, \quad \ell \in [2:L], \quad \widehat{\mathcal{N}}\mathcal{D}_{\ell,i}^{\beta,k} = [m] \widehat{\mathcal{D}}_{\ell,i}^{\beta,k}, \quad \ell \in [L]$$

and

$$\widehat{\mathcal{N}}\mathcal{G}_{\ell,i,j}^k = [m_{\ell-1}m_\ell] \hat{\mathcal{G}}_{\ell,i,j}^k, \quad \ell \in [L+1].$$

Moreover, recalling the definition of  $\underline{\mathcal{D}}_\ell^\gamma$  in (57), we also introduce

$$\widehat{\mathcal{N}}\mathcal{D}_{\ell,i}^{\gamma,k} = \frac{1}{m} \sum_{j=1}^m \left[ \hat{v}_{\ell,i,j}^k \widehat{\mathcal{N}}\mathcal{D}_{\ell,i}^{k,\alpha} \right] \circ \dot{h}'_1 \left( \hat{\beta}_{\ell-1,i}^k \right), \quad \ell \in [2:L],$$

For the ideal process, similar to the notation  $\mathbf{u}_\ell$  in the proof of Theorem 4, we introduce the notations  $\bar{\mathbf{u}}_{\ell,i,j}$  that stands for  $\bar{\Theta}_j$ ,  $(\bar{\Theta}_i, \bar{\Theta}_i)$ ,  $\bar{\Theta}_i$  for  $\ell = 1, 2 \leq \ell \leq L, \ell = L+1$ , respectively.

We also abbreviate the gradients of the ideal process as

$$\begin{aligned}\underline{\mathcal{D}}_{\ell,i}^{\beta,t} &:= \underline{\mathcal{D}}_\ell^\beta(\bar{\Theta}_i, \Phi_*, t), \quad \underline{\mathcal{D}}_{\ell,i}^{\alpha,t} := \underline{\mathcal{D}}_\ell^\alpha(\bar{\Theta}_i, \Phi_*, t), \\ \underline{\mathcal{D}}_{\ell,i}^{\gamma,t} &:= \underline{\mathcal{D}}_\ell^\gamma(\bar{\Theta}_i, \Phi_*, t), \quad \underline{\mathcal{G}}_{\ell,i,j}^v = \underline{\mathcal{G}}_\ell^v(\bar{\mathbf{u}}_{\ell,i,j}; \Phi_*, t).\end{aligned}$$

We use a common notation  $\bar{v}_{\ell,i,j}^t$  to the weights at layer  $\ell$ ; for  $\ell = 1$  let  $\bar{v}_{1,i,j}^t = \bar{v}_{1,j}^t$ .

When  $m$  is finite, the forward and backward propagation for the ideal process is no long exact. Nevertheless, for sufficiently large  $m$ , those propagations relations approximately holds by the following events that happen with high probability:

$$\left\| \frac{1}{m} \sum_{i=1}^m \left[ \bar{v}_{\ell+1,i,j}^{k\eta} \dot{h}_1 \left( \bar{\beta}_{\ell,i}^{k\eta} \right) \right] - \bar{\alpha}_{\ell+1,j}^{k\eta} \right\|_\infty \leq (\|\bar{\Theta}_j\|_\infty + 1) \varepsilon_1, \quad \ell \in [L-1], k \in [0:K], j \in [m], \quad (70)$$

$$\left\| \frac{1}{m} \sum_{j=1}^m \left[ \bar{v}_{\ell,i,j}^{k\eta} \underline{\mathcal{D}}_{\ell,j}^{\alpha,k\eta} \right] \circ \dot{h}'_1 \left( \bar{\beta}_{\ell-1,i}^{k\eta} \right) - \underline{\mathcal{D}}_{\ell,i}^{\gamma,k\eta} \right\|_\infty \leq \varepsilon_1, \quad \ell \in [2:L], k \in [0:K], i \in [m], \quad (71)$$

$$\|\bar{\Theta}_i\|_\infty \leq C \sqrt{\log\left(\frac{m}{\delta}\right)}, \quad i \in [m], \quad (72)$$

$$\frac{1}{m} \sum_{i=1}^m \|\bar{\Theta}_i\|_\infty^j \leq C, \quad j \in [2], \quad (73)$$

for a constant  $C$ . In the proofs of this section, we condition on those events.

**Lemma 24** *The events (70) – (73) happen with probability  $1 - \delta$ .*

The proof consists of the deviation of the actual discrete trajectory from the ideal trajectory over the iteration  $k \in [0 : K]$ . We will upper bound the deviation by induction on  $k$ . For  $k = 0$ , we have the deviation of weights  $\|\bar{v}_{\ell,i,j}^0 - \hat{v}_{\ell,i,j}^0\|_\infty$  from the initial conditions in Definition 5. The induction proceeds as follows. In Lemma 25, we first upper bound the deviation of features using the forward propagation, and then upper bound the deviation of gradients using the backward propagation. Note that

$$\left\| \bar{v}_{\ell,i,j}^{(k+1)\eta} - \hat{v}_{\ell,i,j}^{k+1} \right\|_\infty \leq \left\| \bar{v}_{\ell,i,j}^{k\eta} - \hat{v}_{\ell,i,j}^k \right\|_\infty + \int_{k\eta}^{(k+1)\eta} \left\| \underline{\mathcal{G}}_{\ell,i,j}^s - \widehat{\mathcal{N}}\underline{\mathcal{G}}_{\ell,i,j}^k \right\|_\infty ds. \quad (74)$$

Combining with the Lipschitz continuity of  $\underline{\mathcal{G}}_{\ell,i,j}^t$  in Lemma 26, we complete the inductive step.

**Lemma 25** *Given  $k \in [0 : K]$  and  $\varepsilon < 1$ . Suppose*

$$\left\| \bar{v}_{\ell,i,j}^{k\eta} - \hat{v}_{\ell,i,j}^k \right\|_\infty \leq (\|\bar{\mathbf{u}}_{\ell,i,j}\|_\infty + 1)\varepsilon, \quad \forall \ell \in [L+1], i \in [m_{\ell-1}], j \in [m_\ell]. \quad (75)$$

*Then there exists a constant  $C$  such that*

$$\left\| \bar{\beta}_{L+1,1}^{k\eta} - \hat{\beta}_{L+1,1}^k \right\|_\infty \leq C(\varepsilon + \varepsilon_1), \quad (76)$$

$$\left\| \bar{\beta}_{\ell,i}^{k\eta} - \hat{\beta}_{\ell,i}^k \right\|_\infty \leq C(\|\bar{\Theta}_i\|_\infty + 1)(\varepsilon + \varepsilon_1), \quad \forall \ell \in [L], i \in [m], \quad (77)$$

$$\left\| \bar{\alpha}_{\ell,i}^{k\eta} - \hat{\alpha}_{\ell,i}^k \right\|_\infty \leq C(\|\bar{\Theta}_i\|_\infty + 1)(\varepsilon + \varepsilon_1), \quad \forall \ell \in [2 : L], i \in [m], \quad (78)$$

$$\left\| \underline{\mathcal{G}}_{\ell,i,j}^{k\eta} - \widehat{\mathcal{N}}\underline{\mathcal{G}}_{\ell,i,j}^k \right\|_\infty \leq C(\|\bar{\mathbf{u}}_{\ell,i,j}\|_\infty + 1)(\varepsilon + \varepsilon_1), \quad \forall \ell \in [L+1], i \in [m_{\ell-1}], j \in [m_\ell]. \quad (79)$$

**Lemma 26** *There exists a constant  $C$  such that, for all  $\ell \in [L+1]$ ,  $t_1, t_2 \in [0, T]$ , and  $\mathbf{u}_\ell$ ,*

$$\left\| \underline{\mathcal{G}}_{\ell,i,j}^{t_1} - \underline{\mathcal{G}}_{\ell,i,j}^{t_2} \right\|_\infty \leq C(\|\bar{\mathbf{u}}_{\ell,i,j}\|_\infty + 1)|t_1 - t_2|.$$

**Proof** [Proof of Theorem 6] By Lemma 24, the events in (70) – (73) happen with probability  $1 - \delta$ . Conditioned on those events, we prove by induction on  $k \in [0 : K]$  that

$$\left\| \bar{v}_{\ell,i,j}^{k\eta} - \hat{v}_{\ell,i,j}^k \right\|_\infty \leq (\|\bar{\mathbf{u}}_{\ell,i,j}\|_\infty + 1)e^{Ck\eta}\varepsilon_1, \quad \forall \ell \in [L+1], i \in [m_{\ell-1}], j \in [m_\ell], \quad (80)$$

for some constant  $C$  to be specified. The base case  $k = 0$  follows from Definition 5. Suppose that (80) holds for  $k \in [0 : K - 1]$ . By Lemmas 25 and 26, for  $s \in [k\eta, (k+1)\eta]$ ,

$$\left\| \underline{\mathcal{G}}_{\ell,i,j}^s - \widehat{\mathcal{N}}\underline{\mathcal{G}}_{\ell,i,j}^k \right\|_\infty \leq C'(\|\bar{\mathbf{u}}_{\ell,i,j}\|_\infty + 1)\left(e^{Ck\eta}\varepsilon_1 + \varepsilon_1 + s - k\eta\right).$$

Applying (74) yields that

$$\begin{aligned} \left\| \bar{v}_{\ell,i,j}^{(k+1)\eta} - \hat{v}_{\ell,i,j}^{k+1} \right\|_\infty &\leq \left\| \bar{v}_{\ell,i,j}^{k\eta} - \hat{v}_{\ell,i,j}^k \right\|_\infty + \int_{k\eta}^{(k+1)\eta} \left\| \underline{\mathcal{G}}_{\ell,i,j}^s - \widehat{\mathcal{N}}\underline{\mathcal{G}}_{\ell,i,j}^k \right\|_\infty ds \\ &\leq (\|\bar{\mathbf{u}}_{\ell,i,j}\|_\infty + 1)\left(e^{Ck\eta}\varepsilon_1 + 2C'e^{Ck\eta}\varepsilon_1\eta + C'\frac{\eta^2}{2}\right) \\ &\leq (\|\bar{\mathbf{u}}_{\ell,i,j}\|_\infty + 1)e^{Ck\eta}\varepsilon_1(1 + C''\eta), \end{aligned}$$

for a constant  $C''$ . By letting  $C = C''$ , we arrive at (80) for  $k + 1$  using  $1 + C\eta \leq e^{C\eta}$ . Note that  $k\eta \leq T$  for  $k \in [0 : K]$ ,  $\varepsilon_1 \leq \tilde{\mathcal{O}}(1/\sqrt{m})$ , and  $\|\bar{\mathbf{u}}_{\ell,i,j}\|_\infty \leq \mathcal{O}(\log(m))$  from (72). The conclusion follows from Lemma 25 and the Lipschitz continuity of  $\phi$ .  $\blacksquare$

## C.2.1. PROOFS OF LEMMAS

**Proof** [Proof of Lemma 24] We prove each of the four events happens with probability  $1 - \frac{\delta}{4}$  by standard concentration inequalities. Both (72) and (73) happen with probability  $1 - \frac{\delta}{4}$  by the concentration of sub-gaussian random variables; in particular, (72) follows from Lemma 27 and (73) follows from Lemmas 28 and 29.

For (70) with a given  $k, \ell, j, n$ , consider random vectors

$$\zeta_i := \frac{\bar{v}_{\ell,i,j}^{k\eta} h_1 \left( \bar{\beta}_{\ell-1,i}^{k\eta}(n) \right)}{\|\bar{\Theta}_j\|_\infty + 1},$$

which are bounded by a constant  $C'$  due to the upper bound of  $\bar{v}_\ell$  in Lemma 21. Conditioned on  $\bar{\Theta}_j$ , when  $i \neq j$ ,  $\zeta_i$  are independent and  $\mathbb{E}[\zeta_i | \bar{\Theta}_j] = \frac{\bar{\alpha}_{\ell,j}^{k\eta}(n)}{\|\bar{\Theta}_j\|_\infty + 1}$ . By Hoeffding's inequality, we have

$$\left| \frac{1}{m-1} \sum_{i=1, i \neq j}^m \zeta_i - \frac{\bar{\alpha}_{\ell,j}^{k\eta}(n)}{\|\bar{\Theta}_j\|_\infty + 1} \right| < \varepsilon_1/2,$$

with probability  $1 - \frac{\delta}{4mL(K+1)N}$ . On the other hand, when  $i = j$ , we also have

$$\frac{1}{m} \left| \zeta_j - \frac{\bar{\alpha}_{\ell,j}^{k\eta}(n)}{\|\bar{\Theta}_j\|_\infty + 1} \right| \leq \tilde{C}' \varepsilon_1^2 \leq \varepsilon_1/2,$$

where we use the upper bound of  $\bar{\alpha}_\ell$  in Lemma 21. Therefore, applying the union bound over  $k \in [0 : K]$ ,  $\ell \in [L]$ ,  $j \in [m]$  and  $n \in [N]$ , we have (70) with probability  $1 - \frac{\delta}{4}$ .

For (71) with a given  $k, \ell, i, n$ , consider the random vectors

$$\zeta'_j := [\bar{v}_{\ell+1,i,j}^{k\eta} \underline{\mathcal{D}}_{\ell+1,j}^{\alpha,k\eta}(n)] h'_1 \left( \beta_{\ell,i}^{k\eta}(n) \right).$$

Conditioned on  $\Theta_i$ , when  $i \neq j$ ,  $\zeta'_j$  are independent and  $\mathbb{E}[\zeta'_j | \bar{\Theta}_i] = \underline{\mathcal{D}}_{\ell+1,i}^{\gamma,k\eta}(n)$ . By the boundedness of  $h'_1$  and the upper bound of  $\underline{\mathcal{D}}_{\ell+1}$  in Lemma 21, we have  $\zeta'_j$

$$|\zeta'_j| \leq C' |\bar{v}_{\ell+1,i,j}^{k\eta}| \leq C(1 + \|\Theta_j\|_\infty),$$

and thus  $\zeta'_j$  is sub-gaussian. Applying Lemma 28, we obtain that

$$\left| \frac{1}{m-1} \sum_{j=1, j \neq i}^m \zeta'_j - \underline{\mathcal{D}}_{\ell+1,i}^{\gamma,k\eta}(n) \right| < \varepsilon_1/2,$$

with probability  $1 - \frac{\delta}{4mL(K+1)N}$ . On the other hand, under event (72), we have

$$\frac{1}{m} \left| \zeta'_i - \underline{\mathcal{D}}_{\ell+1,i}^{\gamma,k\eta}(n) \right| \leq \tilde{\mathcal{O}}(\varepsilon_1^2) \leq \varepsilon_1/2$$

Therefore, applying the union bound again over  $k \in [0 : K]$ ,  $\ell \in [L]$ ,  $j \in [m]$ , and  $n \in [N]$ , we have (71) with probability  $1 - \frac{\delta}{4}$ .

■

**Proof** [Proof of Lemma 25] We first consider the forward propagation and prove (76), (77) and (78). For  $\ell = 1$ , since  $\mathbf{X}$  is bounded,

$$\left\| \bar{\beta}_{1,i}^{k\eta} - \hat{\beta}_{1,i}^k \right\|_{\infty} \leq C \|\bar{v}_{1,i}^{k\eta} - \hat{v}_{1,i}^k\|_{\infty} \leq C(\|\bar{\Theta}_i\|_{\infty} + 1)\varepsilon.$$

For  $\ell \in [2 : L]$ , by the triangle inequality,

$$\begin{aligned} & \left\| \bar{\alpha}_{\ell+1,j}^{k\eta} - \hat{\alpha}_{\ell+1,j}^k \right\|_{\infty} \\ & \leq \left\| \bar{\alpha}_{\ell+1,j}^{k\eta} - \frac{1}{m} \sum_{i=1}^m \left[ \bar{v}_{\ell+1,i,j}^{k\eta} \dot{h}_1 \left( \bar{\beta}_{\ell,i}^{k\eta} \right) \right] \right\|_{\infty} + \left\| \frac{1}{m} \sum_{i=1}^m \left[ \bar{v}_{\ell+1,i,j}^{k\eta} \dot{h}_1 \left( \bar{\beta}_{\ell,i}^{k\eta} \right) - \hat{v}_{\ell+1,i,j}^k \dot{h}_1 \left( \hat{\beta}_{\ell,i}^k \right) \right] \right\|_{\infty}. \end{aligned} \quad (81)$$

The first term is approximately the forward propagation that is at most  $(\|\bar{\Theta}\|_{\infty} + 1)\varepsilon_1$  by (70). For the second term, since  $h_1$  is bounded and Lipschitz continuous and the weights  $\bar{v}_{\ell,i,j}$  are upper bounded by Lemma 21 and Assumption 3, we have a further upper bound

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m \underbrace{\left| \bar{v}_{\ell+1,i,j}^{k\eta} \right|}_{\leq C(\|\bar{\Theta}_j\|_{\infty} + 1)} \left\| \dot{h}_1 \left( \bar{\beta}_{\ell,i}^{k\eta} \right) - \dot{h}_1 \left( \hat{\beta}_{\ell,i}^k \right) \right\|_{\infty} + \frac{1}{m} \sum_{i=1}^m \underbrace{\left| \bar{v}_{\ell+1,i,j}^{k\eta} - \hat{v}_{\ell+1,i,j}^k \right|}_{(75)} \left\| \dot{h}_1 \left( \hat{\beta}_{\ell,i}^k \right) \right\|_{\infty} \\ & \leq C(\|\bar{\Theta}_j\|_{\infty} + 1)(\varepsilon + \varepsilon_1), \end{aligned}$$

where in the last step we used (73). We have

$$\left\| \bar{\alpha}_{\ell+1,i}^{k\eta} - \hat{\alpha}_{\ell+1,i}^k \right\|_{\infty} \leq C'(\|\bar{\Theta}_i\|_{\infty} + 1)\varepsilon,$$

for a constant  $C'$ , which gives that

$$\begin{aligned} & \left\| \bar{\beta}_{\ell+1,i}^{k\eta} - \hat{\beta}_{\ell+1,i}^k \right\|_{\infty} \\ & \leq \left\| \bar{\beta}_{\ell,i}^{k\eta} - \hat{\beta}_{\ell,i}^k \right\|_{\infty} + \left\| \dot{h}_2 \left( \bar{\alpha}_{\ell+1,i}^{k\eta} \right) - \dot{h}_2 \left( \hat{\alpha}_{\ell+1,i}^k \right) \right\|_{\infty} \\ & \leq C''(\|\bar{\Theta}_i\|_{\infty} + 1)\varepsilon. \end{aligned}$$

The output layer  $\ell = L + 1$  is similar by applying the upper bound of  $v_{L+1}$  in Assumption 3.

Next we consider the backward propagation and prove (79). Since  $\mathbf{X}$  is bounded,  $h_1$  is bounded and Lipschitz continuous, and  $\underline{\mathcal{D}}_{\ell}^{\alpha}$  is bounded by Lemma 21, it suffices to prove that

$$\begin{aligned} & \left\| \underline{\mathcal{D}}_{L+1}(\Phi_*, k\eta) - \hat{\mathcal{D}}_{L+1,1}^k \right\|_{\infty} \leq C\varepsilon, \\ & \left\| \underline{\mathcal{D}}_{\ell,i}^{\beta,k\eta} - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{\ell,i}^{\beta,k} \right\|_{\infty} \leq C(1 + \|\bar{\Theta}_i\|_{\infty})\varepsilon, \quad \ell \in [L], i \in [m], \\ & \left\| \underline{\mathcal{D}}_{\ell,i}^{\alpha,k\eta} - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{\ell,i}^{\alpha,k} \right\|_{\infty} \leq C(1 + \|\bar{\Theta}_i\|_{\infty})\varepsilon, \quad \ell \in [2 : L], i \in [m], \end{aligned} \quad (82)$$

for a constant  $C$ .



At the output layer  $\ell = L + 1$ , since  $\phi'_1$  is Lipschitz continuous on the first argument,

$$\left\| \underline{\mathcal{D}}_{L+1}(\Phi_*, k\eta) - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{L+1,1}^k \right\|_\infty \leq L_5 \left\| \bar{\beta}_{L+1,1}^{k\eta} - \hat{\beta}_{L+1,1}^k \right\|_\infty \leq C' (\varepsilon + \varepsilon_1),$$

for a constant  $C'$ . At layer  $\ell = L$ , since  $h'_1$  is bounded and Lipschitz continuous and  $\underline{\mathcal{D}}_{L+1}$  is bounded by Lemma 21, applying (77) yields that

$$\left\| \underline{\mathcal{D}}_{L+1,1}^{k\eta} h'_1 \left( \bar{\beta}_{L,i}^{k\eta} \right) - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{L+1,1}^k h'_1 \left( \hat{\beta}_{L,i}^{k\eta} \right) \right\|_\infty \leq C (\|\bar{\Theta}_i\|_\infty + 1) (\varepsilon + \varepsilon_1).$$

Moreover, applying (75) and the upper bound of  $\bar{v}_{L,i}$  in Lemma 21, we obtain that

$$\left\| \underline{\mathcal{D}}_{L,i}^{\beta,k\eta} - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{L,i}^{\beta,k} \right\|_\infty \leq C (1 + \|\bar{\Theta}_i\|_\infty) \varepsilon,$$

for a constant  $C$ .

For each layer  $\ell$  from  $L - 1$  to 1, suppose we have

$$\left\| \underline{\mathcal{D}}_{\ell+1,i}^{\beta,k\eta} - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{\ell+1,i}^{\beta,k} \right\|_\infty \leq C (1 + \|\bar{\Theta}_i\|_\infty) \varepsilon.$$

Since  $h'_2$  is bounded and Lipschitz continuous and  $\underline{\mathcal{D}}_{\ell+1}^\beta$  is bounded by Lemma 21, by (78), we further have

$$\begin{aligned} & \left\| \underline{\mathcal{D}}_{\ell+1,i}^{\alpha,k\eta} - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{\ell+1,i}^{\alpha,k} \right\|_\infty \\ & \leq \left\| \underline{\mathcal{D}}_{\ell+1,i}^{\beta,k\eta} \circ h'_2 \left( \bar{\alpha}_{\ell+1,i}^{k\eta} \right) - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{\ell+1,i}^{\beta,k} \circ h'_2 \left( \hat{\alpha}_{\ell+1,i}^k \right) \right\|_\infty \\ & \leq C' (1 + \|\bar{\Theta}_i\|_\infty) \varepsilon, \end{aligned} \tag{83}$$

for a constant  $C'$ . By the triangle inequality,

$$\begin{aligned} & \left\| \underline{\mathcal{D}}_{\ell+1,i}^{\gamma,k\eta} - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{\ell+1,i}^{\gamma,k} \right\|_\infty \\ & \leq \left\| \underline{\mathcal{D}}_{\ell+1,i}^{\gamma,k\eta} - \frac{1}{m} \sum_{j=1}^m \bar{v}_{\ell+1,i,j}^{k\eta} \underline{\mathcal{D}}_{\ell+1,j}^{\alpha,k\eta} \circ h'_1 \left( \bar{\beta}_{\ell+1,i}^{k\eta} \right) \right\|_\infty \\ & + \left\| \frac{1}{m} \sum_{j=1}^m \bar{v}_{\ell+1,i,j}^{k\eta} \left[ \underline{\mathcal{D}}_{\ell+1,j}^{\alpha,k\eta} \circ h'_1 \left( \bar{\beta}_{\ell+1,i}^{k\eta} \right) \right] - \frac{1}{m} \sum_{j=1}^m \hat{v}_{\ell+1,i,j}^k \left[ \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{\ell+1,j}^{\alpha,k} \circ h'_1 \left( \hat{\beta}_{\ell+1,i}^{k\eta} \right) \right] \right\|_\infty. \end{aligned} \tag{84}$$

The first term is approximately backward propagation and is at most  $\varepsilon_1$  by (71). For the second term, note that  $h'_1$  is bounded and Lipschitz continuous,  $\bar{v}_{\ell+1,i,j}$  and  $\underline{\mathcal{D}}_{\ell+1,j}^\alpha$  are upper bounded by Lemma 21. Applying (75), (77), and (83) at layer  $\ell + 1$  yields that

$$\begin{aligned} & \left\| \underline{\mathcal{D}}_{\ell+1,j}^{\alpha,k\eta} \circ \left[ \bar{v}_{\ell+1,i,j}^{k\eta} h'_1 \left( \bar{\beta}_{\ell+1,i}^{k\eta} \right) \right] - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{\ell+1,j}^{\alpha,k} \circ \left[ \hat{v}_{\ell+1,i,j}^k h'_1 \left( \hat{\beta}_{\ell+1,i}^{k\eta} \right) \right] \right\|_\infty \\ & \leq C (\|\bar{\Theta}_j\|_\infty + 1) (\|\bar{\Theta}_j\|_\infty + \|\bar{\Theta}_i\|_\infty + 1) (\varepsilon + \varepsilon_1). \end{aligned}$$

Therefore, by (73), we have

$$\left\| \underline{\mathcal{D}}_{\ell+1,i}^{\gamma,k\eta} - \widehat{\mathcal{N}}\widehat{\mathcal{D}}_{\ell+1,i}^{\gamma,k} \right\|_\infty \leq C'' (1 + \|\bar{\Theta}_i\|_\infty) \varepsilon,$$

for a constant  $C''$ . We conclude

$$\begin{aligned}
 & \left\| \underline{\mathcal{D}}_{\ell,i}^{\beta,k\eta} - \widehat{\mathcal{N}} \widehat{\mathcal{D}}_{\ell,i}^{\beta,k} \right\|_{\infty} \\
 & \leq \left\| \underline{\mathcal{D}}_{\ell+1,i}^{\gamma,k\eta} - \widehat{\mathcal{N}} \widehat{\mathcal{D}}_{\ell+1,i}^{\gamma,k} \right\|_{\infty} + \left\| \underline{\mathcal{D}}_{\ell+1,i}^{\beta,k\eta} - \widehat{\mathcal{N}} \widehat{\mathcal{D}}_{\ell+1,i}^{\beta,k} \right\|_{\infty} \\
 & \leq C''' (1 + \|\bar{\Theta}_i\|_{\infty}) \varepsilon,
 \end{aligned}$$

for a constant  $C'''$  and obtain (82). ■

**Proof** [Proof of Lemma 26] It suffices to prove  $\underline{\mathcal{G}}_{\ell}^v(\mathbf{u}_{\ell}; \Phi_*, t)$  is continuous on  $t$  for all  $\mathbf{u}_{\ell}$ . The proof is similar to the backward steps in Lemma 25. Recalling boundedness of  $\underline{\mathcal{G}}_{\ell}^{\beta}$  and  $\underline{\mathcal{G}}_{\ell}^{\alpha}$  in Lemma 21, for any fix  $\Theta$ , we have

$$\begin{aligned}
 \|\beta_{\ell}^{t_1} - \beta_{\ell}^{t_2}\| & \leq C (\|\Theta\|_{\infty} + 1) |t_1 - t_2|, \quad \ell \in [L], \\
 \|\alpha_{\ell}^{t_1} - \alpha_{\ell}^{t_2}\| & \leq C (\|\Theta\|_{\infty} + 1) |t_1 - t_2|, \quad \ell \in [2 : L],
 \end{aligned} \tag{85}$$

where  $\beta_{\ell}^t = \Phi_{*,\ell}^{\beta}(\Theta)(t)$  and  $\alpha_{\ell}^t = \Phi_{*,\ell}^{\alpha}(\Theta)(t)$ .

By the boundedness of  $h_1$  and  $\mathbf{X}$ , it suffices to prove the Lipschitz continuity for  $\underline{\mathcal{D}}_{\ell}^{\beta}$  for  $\ell = 1$  and  $\underline{\mathcal{D}}_{\ell}^{\alpha}$  for  $\ell \in [2 : L]$ . We prove

$$\|\underline{\mathcal{D}}_{L+1}(\Phi_*, t_1) - \underline{\mathcal{D}}_{L+1}(\Phi_*, t_2)\|_{\infty} \leq C' |t_1 - t_2|, \tag{86}$$

$$\left\| \underline{\mathcal{D}}_{\ell}^{\beta}(\Theta; \Phi_*, t_1) - \underline{\mathcal{D}}_{\ell}^{\beta}(\Theta; \Phi_*, t_2) \right\|_{\infty} \leq C' (1 + \|\Theta\|_{\infty}) |t_1 - t_2|, \quad \ell \in [L], \tag{87}$$

$$\left\| \underline{\mathcal{D}}_{\ell}^{\alpha}(\Theta; \Phi_*, t_1) - \underline{\mathcal{D}}_{\ell}^{\alpha}(\Theta; \Phi_*, t_2) \right\|_{\infty} \leq C' (1 + \|\Theta\|_{\infty}) |t_1 - t_2|, \quad \ell \in [2 : L]. \tag{88}$$

At the output layer  $\ell = L + 1$ , by the Lipschitz continuity of  $\phi'$ , we have

$$\begin{aligned}
 & \left\| \underline{\mathcal{D}}_{L+1}(\Phi_*, t_1) - \underline{\mathcal{D}}_{L+1}(\Phi_*, t_2) \right\|_{\infty} \\
 & \leq L_5 \left\| \beta_{L+1}^{t_1} - \beta_{L+1}^{t_2} \right\|_{\infty} \\
 & \leq L_5 \left\| \int v_{L+1}^{t_1} \dot{h}_1(\beta_L^{t_1}) - v_{L+1}^{t_2} \dot{h}_1(\beta_L^{t_2}) dp(\Theta) \right\|_{\infty}.
 \end{aligned} \tag{89}$$

By the upper bound and Lipschitz continuity of  $v_{L+1}$  in Lemma 21, we obtain (86). At layer  $\ell = L$ , using (11b), we obtain (87) from the upper bounds and the Lipschitz continuity of  $\underline{\mathcal{D}}_{L+1}(\Theta, \Phi_*, t)$ ,  $v_{L+1}^t(\Theta)$ , and  $\dot{h}'_1(\beta_L^t)$ .

For each layer  $\ell$  from  $L - 1$  to 1, suppose we have (87) at layer  $\ell + 1$ . From the upper bounds and the Lipschitz continuity of  $\underline{\mathcal{D}}_{\ell+1}^{\beta}$  and  $\dot{h}'_2(\bar{\alpha}_{\ell+1}^t)$ , we have

$$\left\| \underline{\mathcal{D}}_{\ell+1}^{\alpha}(\Theta; \Phi_*, t_1) - \underline{\mathcal{D}}_{\ell+1}^{\alpha}(\Theta; \Phi_*, t_2) \right\|_{\infty} \leq C'' (1 + \|\Theta\|_{\infty}) |t_1 - t_2|, \tag{90}$$

for a constant  $C''$ . Moreover,

$$\begin{aligned}
 & \int \left\| v_{\ell+1}^{t_1}(\Theta, \bar{\Theta}) \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi_*, t_1) - v_{\ell+1}^{t_2}(\Theta, \bar{\Theta}) \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi_*, t_2) \right\|_\infty dp(\bar{\Theta}) \\
 & \leq \int \underbrace{\left| v_{\ell+1}^{t_1}(\Theta, \bar{\Theta}) - v_{\ell+1}^{t_2}(\Theta, \bar{\Theta}) \right|}_{\leq C|t_1-t_2|} \left\| \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi_*, t) \right\|_\infty dp(\bar{\Theta}) \\
 & \quad + \int \left| v_{\ell+1}^{t_2}(\Theta, \bar{\Theta}) \right| \underbrace{\left\| \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi_*, t_1) - \underline{\mathcal{D}}_{\ell+1}^\alpha(\bar{\Theta}; \Phi_*, t_2) \right\|_\infty}_{\leq C(\|\bar{\Theta}\|_\infty + 1)|t_1-t_2|} dp(\bar{\Theta}) \\
 & \leq C'|t_1 - t_2|, \tag{91}
 \end{aligned}$$

where in the last step we used the upper bounds of  $v_{\ell+1}$  and  $\underline{\mathcal{D}}_{\ell+1}^\alpha$  in Lemma 21, sub-gaussianness of  $p$ , and Corollary 31. Combing the above results with

$$\left\| \dot{h}_1(\beta_\ell^{t_1}) - \dot{h}_1(\beta_\ell^{t_2}) \right\| \leq C(\|\Theta\|_\infty + 1)|t_1 - t_2|, \quad \ell \in [L]$$

from (85), we obtain

$$\left\| \underline{\mathcal{D}}_{\ell+1}^\gamma(\Theta; \Phi_*, t_1) - \underline{\mathcal{D}}_{\ell+1}^\gamma(\Theta; \Phi_*, t_2) \right\|_\infty \leq C'''(1 + \|\Theta\|_\infty)|t_1 - t_2|, \tag{92}$$

for a constant  $C'''$ . Combining (90) and (92), we can achieve (87) at  $\ell$ .  $\blacksquare$

### C.3. Proof of Proposition 1

**Proof** [Proof of Proposition 1] Explicitly shown in (Nguyen and Pham, 2020, Corollary 25), in the mean-field limit that  $m \rightarrow \infty$ , the weights remain mutually independent and follow a common distribution that only depends on time  $t$  in the intermediate layers. Therefore, by the law of large numbers, the features are the same. We have Proposition 1.  $\blacksquare$

## Appendix D. Preliminary of Proofs

In this paper, we adopt the definition of sub-gaussian distributions in Vershynin (2010). Below we present properties of sub-gaussian distributions. The equivalence among those properties are given in (Vershynin, 2010, Lemma 5.5).

**Lemma 27** *Let  $\xi$  be an univariate random variable that follows a  $\sigma$ -sub-gaussian distribution. Then there exists an absolute constant  $C$  such that*

(1) *Tails*  $\mathbb{P}(|\xi| > t) \leq \exp(1 - t^2/(C\sigma)^2)$  for all  $t \geq 0$ ;

(2) *Moments*:  $(\mathbb{E}|\xi|^q)^{1/q} \leq C\sigma\sqrt{q}$  for all  $q \geq 1$ ;

(3) *If*  $\mathbb{E}[\xi] = 0$ , *then*  $\mathbb{E}[\exp(t\xi)] \leq \exp(t^2(C\sigma)^2)$  for all  $t \in \mathbb{R}$ .

**Lemma 28 (Concentration Inequality for Sub-gaussian Distributions (Vershynin, 2010, Proposition 5.10))**

Let  $\{\xi_i\}_{i=1}^m$  be independent centered  $\sigma$ -sub-gaussian random variables. Then, for an absolute constant  $C$ ,

$$\mathbb{P}\left(\left|\frac{1}{m}\sum_{i=1}^m \xi_i\right| \geq \varepsilon\right) \leq 2 \exp\left(-\frac{m\varepsilon^2}{4(C\sigma)^2}\right).$$

We say a random variable  $\xi$  is sub-exponential if

$$\sup_{q \geq 1} q^{-1}(\mathbb{E}|\xi|^q)^{1/q} < \infty.$$

A sub-exponential random variable is equivalent to the a squared sub-gaussian random variable (Vershynin, 2010, Lemma 5.14). It satisfies the following concentration inequality:

**Lemma 29 (Bernstein's Inequality for Sub-Exponential Distributions (Vershynin, 2010, Corollary 5.17))**

Let  $\{\xi_i\}_{i=1}^m$  be independent centered sub-exponential random variables such that  $(\mathbb{E}|\xi|^q)^{1/q} \leq Kq$  for all  $q \geq 1$ . Then, for an absolute constant  $c$ ,

$$\mathbb{P}\left(\left|\frac{1}{m}\sum_{i=1}^m z_i\right| \geq \varepsilon\right) \leq 2 \exp\left[-cm \min\left(\frac{\varepsilon^2}{K^2}, \frac{\varepsilon}{K}\right)\right].$$

**Lemma 30** For a  $d$ -dimensional random vector  $\boldsymbol{\xi} \in \mathbb{R}^d$ , we have:

- (1)  $\|\boldsymbol{\xi}\|_\infty$  is  $\sigma$ -sub-gaussian  $\implies \boldsymbol{\xi}$  is  $(\sigma\sqrt{d})$ -sub-gaussian;
- (2)  $\boldsymbol{\xi}$  is  $\sigma$ -sub-gaussian  $\implies \|\boldsymbol{\xi}\|_\infty$  is  $\mathcal{O}(\sigma\sqrt{\log d})$ -sub-gaussian.

**Proof** (1) For any  $\mathbf{u} \in \mathbb{S}^{d-1}$ , we have  $|\mathbf{u}^\top \boldsymbol{\xi}| \leq \sqrt{d}\|\boldsymbol{\xi}\|_\infty$ .

(2) Note that  $\|\boldsymbol{\xi}\|_\infty = \max_{i \in [d]} |e_i^\top \boldsymbol{\xi}|$ , where  $e_i$  denotes the unit vector with  $i$ -th coordinate being one. Applying Lemma 27 (1) and the union bound yields that

$$\mathbb{P}[\|\boldsymbol{\xi}\|_\infty > t] \leq \min\left\{de^{1-\frac{t^2}{(C\sigma)^2}}, 1\right\} = e^{-\left(\frac{t^2}{(C\sigma)^2} - \log(ed)\right)_+} \leq e^{1-\frac{t^2}{(C\sigma\sqrt{\log(ed)})^2}},$$

where we used the fact that  $(\frac{t^2}{a} - b)_+ \geq \frac{t^2}{ab} - 1$  when  $b \geq 1$ . Therefore,  $\|\boldsymbol{\xi}\|_\infty$  is  $\mathcal{O}(\sigma\sqrt{\log d})$ -sub-gaussian by the equivalent definition of sub-gaussian distributions in (Vershynin, 2010, Lemma 5.5). ■

**Corollary 31** For a  $\sigma$ -sub-gaussian random vector  $\boldsymbol{\xi} \in \mathbb{R}^d$ , we have

$$(\mathbb{E}\|\boldsymbol{\xi}\|_\infty^q)^{1/q} \leq \mathcal{O}(\sigma\sqrt{q\log d}), \quad q \geq 1.$$

**Proof** From Lemma 30,  $\|\boldsymbol{\xi}\|_\infty$  is  $\mathcal{O}(\sigma\sqrt{\log d})$ -sub-gaussian. Applying Lemma 27 (2), we achieve the desired result. ■

## Appendix E. Extensions on Fully-connected DNNs

### E.1. Discrete Fully-connected DNN

We first introduce the standard DNN. Let  $m_\ell$  denote the number of units at layer  $\ell$  for  $\ell \in [0 : L+1]$ . Let  $m_0 = d$  and node  $i$  outputs the value of  $i$ -th coordinate of the training data for  $i \in [d]$ . Let  $m_{L+1} = 1$  that is the unit of the final network output. For  $\ell \in [L+1]$ , the output, i.e. features, of node  $i$  in layer  $\ell$  is denoted by  $\hat{\theta}_{\ell,i} \in \mathbb{R}^N$ ; the weight that connects the node  $i$  for the  $N$  training samples at layer  $\ell - 1$  to node  $j$  at layer  $\ell$  is denoted by  $\hat{w}_{\ell,i,j} \in \mathbb{R}$ .

(1) At the input layer, for  $i \in [d]$ , let

$$\hat{\theta}_{0,i} := [\mathbf{x}^1(i), \mathbf{x}^2(i), \dots, \mathbf{x}^N(i)]^\top. \quad (93)$$

(2) We recursively define the upper layers ( $\ell \in [L]$ ) as below.

$$\hat{\theta}_{\ell,j} := \begin{cases} \frac{1}{m_0} \sum_{i=1}^{m_0} \hat{w}_{1,i,j} \hat{\theta}_{0,i}, & j \in [m_1], \ell = 1, \\ \frac{1}{m_{\ell-1}} \sum_{i=1}^{m_{\ell-1}} \hat{w}_{\ell,i,j} \hat{h}(\hat{\theta}_{\ell-1,i}), & j \in [m_\ell], \ell \in [2 : L], \end{cases} \quad (94)$$

where  $h$  is the activation function.

(3) At the output layer,

$$\hat{\theta}_{L+1,1} := \frac{1}{m_L} \sum_{i=1}^{m_L} \hat{w}_{L+1,i} \hat{h}(\hat{\theta}_{L,i}). \quad (95)$$

We collect the weights at the  $\ell$ -th layer ( $\ell \in [L+1]$ ) into a single vector denoted by  $\hat{\mathbf{w}}_\ell$  and all the weights into a single vector denoted by  $\hat{\mathbf{w}}$ . Similarly, we aggregate features at  $\ell$ -th layer ( $\ell \in [L]$ ) into a single vector denoted by  $\hat{\boldsymbol{\theta}}_\ell$  and all the features into a single vector denoted by  $\hat{\boldsymbol{\theta}}$ .

### E.2. Continuous Fully-connected DNN

We introduce our continuous DNN formulation using similar forward propagation of the the discrete DNN in Section E.1.

(1) At the input layer, let  $\mathbf{X} = [\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N]^\top \in \mathbb{R}^{N \times d}$ .

(2) At the first layer, each hidden node (before the activation function) is computed by a linear mapping of the input data, so each node can be indexed by the weights connecting it to the input. We introduce a probability measure  $p_1(\mathbf{w}_1) \in \mathcal{P}(\mathbb{R}^d)$  for the weights to describe the states of first layer and let<sup>6</sup>

$$\boldsymbol{\theta}_1(\mathbf{w}_1) := \frac{1}{d} (\mathbf{X} \mathbf{w}_1). \quad (96)$$

6. The state of the first layer can be equivalently characterized by either the output or the weight that are related by a linear mapping.

---

**Algorithm 3** Scaled Gradient Descent for Training a DNN

---

- 1: Input the data  $\{\mathbf{x}^i, y^i\}_{i=1}^N$ , step size  $\eta$ , and initial weights  $\hat{\mathbf{w}}^0$ .
- 2: **for**  $k = 0, 1, \dots, K - 1$  **do**
- 3:   Perform forward-propagation (94) and (95) to compute  $\hat{\boldsymbol{\theta}}_{L+1,1}^k$ .
- 4:   Perform backward-propagation to compute the gradient  $\hat{\mathcal{G}}_{\ell,i,j}^k = \frac{\partial \hat{\mathcal{L}}_N}{\partial \hat{w}_{\ell,i,j}^k}$ .
- 5:   Perform scaled Gradient Descent:

$$\hat{w}_{\ell,i,j}^{k+1} = \hat{w}_{\ell,i,j}^k - [\eta m_{\ell-1} m_\ell] \hat{\mathcal{G}}_{\ell,i,j}^k, \quad \ell \in [L + 1], i \in [m_{\ell-1}], j \in [m_\ell].$$

- 6: **end for**
  - 7: Output the weights  $\hat{\mathbf{w}}^K$ .
- 

- (3) At the second layer, recall that the output of each node, i.e., the feature, for the training samples is a  $N$ -dimensional vector. We use the features  $\boldsymbol{\theta}_2 \in \mathbb{R}^N$  to index those nodes. We introduce a probability measure  $p_2(\boldsymbol{\theta}_2) \in \mathcal{P}(\mathbb{R}^N)$  to describe the overall states of the second layer and function  $w_2 : \text{supp}(p_1) \times \text{supp}(p_2) \rightarrow \mathbb{R}$  to denote the weights on the connections from layer 1 to 2. We have for any  $\boldsymbol{\theta}_2 \in \text{supp}(p_2)$ , we have the constraint for  $w_2$  and  $p_1$  that

$$\int w_2(\mathbf{w}_1, \boldsymbol{\theta}_2) \dot{h}(\boldsymbol{\theta}_1(\mathbf{w}_1)) dp_1(\mathbf{w}_1) = \boldsymbol{\theta}_2. \quad (97)$$

- (4) Similarly, for  $\ell \in [3 : L]$ , let  $\boldsymbol{\theta}_\ell \in \mathbb{R}^N$  be the index of nodes according to the features. We introduce a probability measure  $p_\ell(\boldsymbol{\theta}_\ell) \in \mathcal{P}(\mathbb{R}^N)$  to describe the states the  $\ell$ -th layer and function  $w_\ell : \text{supp}(p_{\ell-1}) \times \text{supp}(p_\ell) \rightarrow \mathbb{R}$  to denote the weights on the connections from layer  $\ell - 1$  to  $\ell$ . We have any all  $\boldsymbol{\theta}_\ell \in \text{supp}(p_\ell)$ , we have the constraint for  $w_\ell$  and  $p_{\ell-1}$  that:

$$\int w_\ell(\boldsymbol{\theta}_{\ell-1}, \boldsymbol{\theta}_\ell) \dot{h}(\boldsymbol{\theta}_{\ell-1}) dp_{\ell-1}(\boldsymbol{\theta}_{\ell-1}) = \boldsymbol{\theta}_\ell. \quad (98)$$

- (5) Finally, let  $w_{L+1} : \text{supp}(p) \rightarrow \mathbb{R}$  be the weights in the layer  $L + 1$  and  $\boldsymbol{\theta}_{L+1}$  be the final output, and we have the constraint for  $w_{L+1}$  and  $p_L$  that

$$\int w_{L+1}(\boldsymbol{\theta}_L) \dot{h}(\boldsymbol{\theta}_L) dp_L(\boldsymbol{\theta}_L) = \boldsymbol{\theta}_{L+1}. \quad (99)$$

Finally, a loss is imposed on the top layer and our target is to minimize the objective  $\frac{1}{N} \sum_{n=1}^N \phi(\boldsymbol{\theta}_{L+1}(n), y^n)$ . We can find that a static continuous DNN in our formulation is characterized by  $\{w_\ell\}_{\ell=2}^{L+1}$  and  $\{p_\ell\}_{\ell=1}^L$ .

### E.3. Scaled Gradient Descent for Training Fully-connected DNN

We still consider the scaled Gradient Descent algorithm to optimize the DNN. Given an initial weights  $\hat{\mathbf{w}}^0$ , the meta algorithm of the scaled Gradient Descent is shown in Algorithm 3, where the gradients  $\hat{\mathcal{G}}_{\ell,i,j}^k$  can be obtained by the standard backward-propagation algorithm.

#### E.4. Neural Feature Flow for Training Continuous Fully-connected DNN

We derive the evolution of the Gradient Descent algorithm on a continuous DNN  $(\{w_\ell\}_{\ell=2}^{L+1}, \{p_\ell\}_{\ell=1}^L)$ . We first introduce the notations for the trajectories of  $\mathbf{w}_1$ ,  $\{w_\ell\}_{\ell=2}^L$ , and  $\{\theta_\ell\}_{\ell=2}^L$ :

- $\Psi_\ell^\theta : \text{supp}(p_\ell) \rightarrow C([0, T], \mathbb{R}^N)$  is the trajectory of  $\theta_\ell$  for  $\ell \in [2 : L]$ ;
- $\Psi_1^w : \text{supp}(p_1) \rightarrow C([0, T], \mathbb{R}^d)$  and  $\Psi_{L+1}^w : \text{supp}(p_L) \rightarrow C([0, T], \mathbb{R})$  are the trajectories of  $\mathbf{w}_1$  and  $w_{L+1}$ , respectively;
- $\Psi_\ell^w : \text{supp}(p_{\ell-1}) \times \text{supp}(p_\ell) \rightarrow C([0, T], \mathbb{R})$  is the trajectory of  $w_\ell$  for  $\ell \in [2 : L]$ ;
- Let  $\Psi$  be the collection of these trajectories.

The continuous gradient for the weight can be obtained from the backward-propagation algorithm. Especially, we define

$$\theta_{L+1}(\Psi, t) := \int \Psi_{L+1}^w(\theta_L)(t) \dot{h}(\Psi_L^\theta(\theta_L)(t)) dp_L(\theta_L), \quad (100a)$$

$$\begin{aligned} \bar{\mathcal{D}}_{L+1}(\Psi, t) &:= [\phi'_1(\theta_{L+1}^t(1), y^1), \phi'_1(\theta_{L+1}^t(2), y^2), \dots, \phi'_1(\theta_{L+1}^t(N), y^N)]^\top, \\ \bar{\mathcal{D}}_L(\theta_L; \Psi, t) &:= [\Psi_{L+1}^w(\theta_L)(t) \bar{\mathcal{D}}_{L+1}(\Psi, t)] \circ \dot{h}'(\Psi_L^\theta(\theta_L)(t)), \end{aligned} \quad (100b)$$

$$\bar{\mathcal{D}}_\ell(\theta_\ell; \Psi, t) := \left[ \int \Psi_{\ell+1}^w(\theta_\ell, \theta_{\ell+1})(t) \bar{\mathcal{D}}_{\ell+1}(\theta_{\ell+1}; \Psi, t) dp_{\ell+1}(\theta_{\ell+1}) \right] \circ \dot{h}'(\Psi_\ell^\theta(\theta_\ell)(t)), \quad (100c)$$

$$\bar{\mathcal{D}}_1(\mathbf{w}_1; \Psi, t) := \left[ \int \Psi_2^w(\mathbf{w}_1, \theta_2)(t) \bar{\mathcal{D}}_2(\theta_2; \Psi, t) dp_2(\theta_2) \right] \circ \dot{h}'(\theta_1(\Psi_1^w(\mathbf{w}_1)(t))), \quad (100d)$$

where in (100b),  $\theta_L \in \text{supp}(p_L)$ , in (100c),  $\ell \in [2 : L - 1]$  and  $\theta_\ell \in \text{supp}(p_\ell)$ , and in (100d),  $\mathbf{w}_1 \in \text{supp}(p_1)$  and  $\theta_1(\cdot)$  is defined by (96). Then the gradient of the weights can be written as below.

$$\bar{\mathcal{G}}_{L+1}^w(\theta_L; \Psi, t) := \frac{1}{N} [\bar{\mathcal{D}}_{L+1}(\Psi, t)]^\top \dot{h}(\Psi_L^\theta(\theta_L)(t)), \quad (101a)$$

$$\bar{\mathcal{G}}_\ell^w(\theta_{\ell-1}, \theta_\ell; \Psi, t) := \frac{1}{N} [\bar{\mathcal{D}}_\ell(\theta_\ell; \Psi, t)]^\top \dot{h}(\Psi_{\ell-1}^\theta(\theta_{\ell-1})(t)), \quad (101b)$$

$$\bar{\mathcal{G}}_2^w(\mathbf{w}_1, \theta_2; \Psi, t) := \frac{1}{N} [\bar{\mathcal{D}}_2(\theta_2; \Psi, t)]^\top \dot{h}(\theta_1(\Psi_1^w(\mathbf{w}_1)(t))), \quad (101c)$$

$$\bar{\mathcal{G}}_1^w(\mathbf{w}_1; \Psi, t) := \frac{1}{N} \mathbf{X}^\top [\bar{\mathcal{D}}_1(\mathbf{w}_1; \Psi, t)], \quad (101d)$$

where in (101a),  $\theta_L \in \text{supp}(p_L)$ , in (101b),  $\ell \in [3 : L]$ ,  $\theta_{\ell-1} \in \text{supp}(p_{\ell-1})$ , and  $\theta_\ell \in \text{supp}(p_\ell)$ , in (101c),  $\mathbf{w}_1 \in \text{supp}(p_1)$  and  $\theta_2 \in \text{supp}(p_2)$ , and in (101d),  $\mathbf{w}_1 \in \text{supp}(p_1)$ .

Moreover, we expect that the features satisfy the constraints:

$$\begin{aligned} \int \Psi_2^w(\mathbf{w}_1, \theta_2)(t) \dot{h}(\theta_1(\Psi_1^w(\mathbf{w}_1)(t))) dp_1(\mathbf{w}_1) &= \Psi_2^\theta(\theta_2)(t), \quad \theta_2 \in \text{supp}(p_2), \\ \int \Psi_\ell^w(\theta_{\ell-1}, \theta_\ell)(t) \dot{h}(\Psi_{\ell-1}^\theta(\theta_{\ell-1})(t)) dp_{\ell-1}(\theta_{\ell-1}) &= \Psi_\ell^\theta(\theta_\ell)(t), \quad \ell \in [3 : L], \theta_\ell \in \text{supp}(p_\ell). \end{aligned}$$

So the drift term for the features can be obtained by the chain rule:

$$\bar{\mathcal{G}}_1^\theta(\mathbf{w}_1; \Psi, t) := \frac{1}{d} \left[ \mathbf{X} \bar{\mathcal{G}}_1^w(\mathbf{w}_1; \Psi, t) \right], \quad (102a)$$

$$\begin{aligned} \bar{\mathcal{G}}_2^\theta(\boldsymbol{\theta}_2; \Psi, t) &:= \int \Psi_2^w(\mathbf{w}_1, \boldsymbol{\theta}_2)(t) \left[ \dot{h}'\left(\boldsymbol{\theta}_1(\Psi_1^w(\mathbf{w}_1)(t))\right) \circ \bar{\mathcal{G}}_1^\theta(\mathbf{w}_1; \Psi, t) \right] dp_1(\mathbf{w}_1) \\ &\quad + \int \dot{h}\left(\boldsymbol{\theta}_1(\Psi_1^w(\mathbf{w}_1)(t))\right) \circ \bar{\mathcal{G}}_2^w(\mathbf{w}_1, \boldsymbol{\theta}_2; \Psi, t) dp_1(\mathbf{w}_1), \end{aligned} \quad (102b)$$

$$\begin{aligned} \bar{\mathcal{G}}_\ell^\theta(\boldsymbol{\theta}_\ell; \Psi, t) &:= \int \Psi_\ell^w(\boldsymbol{\theta}_{\ell-1}, \boldsymbol{\theta}_\ell)(t) \left[ \dot{h}'\left(\Psi_{\ell-1}^\theta(\boldsymbol{\theta}_{\ell-1})(t)\right) \circ \bar{\mathcal{G}}_{\ell-1}^\theta(\boldsymbol{\theta}_{\ell-1}; \Psi, t) \right] dp_{\ell-1}(\boldsymbol{\theta}_{\ell-1}) \\ &\quad + \int \dot{h}\left(\Psi_{\ell-1}^\theta(\boldsymbol{\theta}_{\ell-1})(t)\right) \circ \bar{\mathcal{G}}_\ell^w(\boldsymbol{\theta}_{\ell-1}, \boldsymbol{\theta}_\ell; \Psi, t) dp_{\ell-1}(\boldsymbol{\theta}_{\ell-1}), \end{aligned} \quad (102c)$$

where in (102a),  $\mathbf{w}_1 \in \text{supp}(p_1)$ , in (102b),  $\boldsymbol{\theta}_2 \in \text{supp}(p_2)$ , and in (102c),  $\ell \in [3 : L]$  and  $\boldsymbol{\theta}_\ell \in \text{supp}(p_\ell)$ . Now we define the process of a continuous DNN trained by Gradient Descent, i.e., neural feature flow. It characterizes the evolution of both weights and features.

**Definition 32 (Neural Feature Flow for DNN)** *Given an initial continuous DNN  $(\{w_\ell\}_{\ell=2}^{L+1}, \{p_\ell\}_{\ell=1}^L)$  that satisfies (96) – (99) and  $T < \infty$ , we say a trajectory  $\Psi_*$  is a neural feature flow if for all  $t \in [0, T]$ ,*

(1) for all  $\ell \in [2 : L]$  and  $\boldsymbol{\theta}_\ell \in \text{supp}(p_\ell)$ ,

$$\Psi_{*,\ell}^\theta(\boldsymbol{\theta}_\ell)(t) = \boldsymbol{\theta}_\ell - \int_{s=0}^t \bar{\mathcal{G}}_\ell^\theta(\boldsymbol{\theta}_\ell; \Psi_*, s) ds,$$

(2) for all  $\mathbf{w}_1 \in \text{supp}(p_1)$ ,

$$\Psi_{*,1}^w(\mathbf{w}_1)(t) = \mathbf{w}_1 - \int_{s=0}^t \bar{\mathcal{G}}_1^w(\mathbf{w}_1; \Psi_*, s) ds,$$

(3) for all  $\mathbf{w}_1 \in \text{supp}(p_1)$  and  $\boldsymbol{\theta}_2 \in \text{supp}(p_2)$ ,

$$\Psi_{*,2}^w(\mathbf{w}_1, \boldsymbol{\theta}_2)(t) = w_2(\mathbf{w}_1, \boldsymbol{\theta}_2) - \int_{s=0}^t \bar{\mathcal{G}}_2^w(\mathbf{w}_1, \boldsymbol{\theta}_2; \Psi_*, s) ds,$$

(4) for all  $\ell \in [2 : L - 1]$ ,  $\boldsymbol{\theta}_\ell \in \text{supp}(p_\ell)$ , and  $\boldsymbol{\theta}_{\ell+1} \in \text{supp}(p_{\ell+1})$ ,

$$\Psi_{*,\ell+1}^w(\boldsymbol{\theta}_\ell, \boldsymbol{\theta}_{\ell+1})(t) = w_{\ell+1}(\boldsymbol{\theta}_\ell, \boldsymbol{\theta}_{\ell+1}) - \int_{s=0}^t \bar{\mathcal{G}}_{\ell+1}^w(\boldsymbol{\theta}_\ell, \boldsymbol{\theta}_{\ell+1}; \Psi_*, s) ds,$$

(5) for all  $\boldsymbol{\theta}_L \in \text{supp}(p_L)$ ,

$$\Psi_{*,L+1}^w(\boldsymbol{\theta}_L)(t) = w_{L+1}(\boldsymbol{\theta}_L) - \int_{s=0}^t \bar{\mathcal{G}}_{L+1}^w(\boldsymbol{\theta}_L; \Psi_*, s) ds.$$



---

**Algorithm 4** Initializing a Discrete DNN.
 

---

- 1: Input the data  $\{\hat{\theta}_{0,i}\}_{i=1}^d$  in (93), variance  $\sigma_1 > 0$ , and a constant  $C_3$ .
- 2: Independently draw  $\hat{w}_{1,i,j} \sim p_0 = \mathcal{N}(0, d\sigma_1^2)$  for  $i \in [d]$  and  $j \in [m]$ .
- 3: Set  $\hat{\theta}_{1,j} = \frac{1}{d} \sum_{i=1}^d \hat{w}_{1,i,j} \hat{\theta}_{0,i}$  where  $j \in [m]$ . ◇ Standard Initialization for layer 1
- 4: **for**  $\ell = 2, \dots, L$  **do**
- 5:   Independently draw  $\tilde{w}_{\ell,i,j} \sim \mathcal{N}(0, m\sigma_1^2)$  for  $i, j \in [m]$ .
- 6:   Set  $\hat{\theta}_{\ell,j} = \frac{1}{m} \sum_{i=1}^m \tilde{w}_{\ell,i,j} \dot{h}(\hat{\theta}_{\ell-1,i})$  where  $j \in [m]$ . ◇ Standard Initialization for layer  $\ell$
- 7: **end for**
- 8: Set  $\hat{w}_{L+1,i,1} = C$  where  $i \in [m]$ . ◇ Simply initialize  $\{\hat{w}_{L+1,i,1}\}_{i=1}^m$  by a constant
- 9: **for**  $\ell = 2, \dots, L$  **do**
- 10:   **for**  $j = 1, \dots, m$  **do**
- 11:     Solve convex optimization problem: ◇ Perform  $\ell_2$ -regression to reduce redundancy

$$\min_{\{\hat{w}_{\ell,i,j}\}_{i=1}^m} \frac{1}{m} \sum_{i=1}^m (\hat{w}_{\ell,i,j})^2, \quad \text{s.t. } \hat{\theta}_{\ell,j} = \frac{1}{m} \sum_{i=1}^m \hat{w}_{\ell,i,j} \dot{h}(\hat{\theta}_{\ell-1,i}).$$

- 12:   **end for**
  - 13: **end for**
  - 14: Output the discrete DNN parameters  $(\hat{w}, \hat{\theta})$ .
- 

### E.5. Informal Result

We still consider the DNN initialized by a standard initialization with an additional regression procedure, shown in Algorithm 4. We show in Theorem 6 that there is a neural feature flow that can capture the evolution of a DNN that is initialized by Algorithm 4 and trained by Gradient Descent, i.e., Algorithm 3.

**Theorem 33 (Informal)** *Under suitable conditions, there is an initialization  $(\{w_\ell\}_{\ell=2}^{L+1}, \{p_\ell\}_{\ell=1}^L)$  such that the continuous DNN has the following properties.*

- (1) For any  $T < \infty$ , there exists a unique neural feature flow  $\Psi_*$  satisfying Definition 32.
- (2) Suppose  $\varepsilon \leq \tilde{\mathcal{O}}(1)$ ,  $\delta \leq 1$ ,  $m \geq \tilde{\Omega}(\varepsilon^{-2})$ , the step size  $\eta \leq \tilde{\mathcal{O}}(\varepsilon)$ . Let  $T$  be a constant and  $K := \lceil T/\eta \rceil$ . Let  $\hat{\mathcal{L}}_N^k := \frac{1}{n} \sum_{n=1}^N \phi(\hat{\theta}_{L+1,1}^k(n), y^n)$  be the loss of running scaled Gradient Descent Algorithm 3 on a DNN initialized by Algorithm 4 at  $k$ -th step, and  $\mathcal{L}_N^t := \frac{1}{N} \sum_{n=1}^n \phi(\theta_{L+1}(\Psi_*, t)(n), y^n)$  be the loss of neural feature flow at time  $t$ . Then, with probability  $1 - \delta$ ,

$$\sup_{k \in [0:K]} \left| \hat{\mathcal{L}}_N^k - \mathcal{L}_N^{k\eta} \right| \leq \tilde{\mathcal{O}}(\varepsilon).$$

### Appendix F. Simulations

In this section, we perform a toy simulation to validate our theory. We consider a synthetic 1-D regression task:  $f(x) = \sin(x)$ . We randomly generate  $N = 100$  training samples uniformly from  $[-\pi, \pi]$  and experiment on a four-hidden-layer ( $L = 4$ ) NN. We choose the activation function  $h_1$  and  $h_2$  as  $\tanh(x)$  and apply  $\ell_2$  loss, i.e.,  $\phi(y', y) = |y' - y|^2$ . The weights are initialized by a

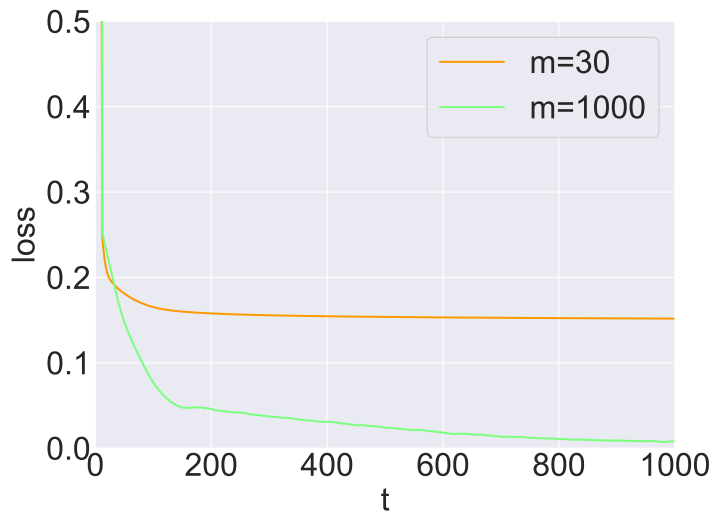


Figure 1: The experimental results on the synthetic data.  $m$  denotes the number of hidden units.

standard strategy (Glorot and Bengio, 2010) with an additional  $\ell_2$  regression in Algorithm 2. We see that the scaled GD in Algorithm 3 from the initialization achieves the global optimal solution for overparameterized DNNs.

The experimental result is shown Fig. 1. From the experiments, we can conclude that (1) the initialization and scaled GD are workable; (2) with the growth of the number of hidden units  $m$ , scaled GD achieves the global minimum.