
Global Convergence of Over-parameterized Deep Equilibrium Models

Anonymous Author
Anonymous Institution

Abstract

A deep equilibrium model (DEQ) is implicitly defined through an equilibrium point of an infinite-depth weight-tied model with an input-injection. Instead of infinite computations, it solves an equilibrium point directly with root-finding and computes gradients with implicit differentiation. In this paper, the training dynamics of over-parameterized DEQs are investigated, and we propose a novel probabilistic framework to overcome the challenge arising from the weight-sharing and the infinite depth. By supposing a condition on the initial equilibrium point, we prove that the gradient descent converges to a globally optimal solution at a linear convergence rate for the quadratic loss function. We further perform a fine-grained non-asymptotic analysis about random DEQs and the corresponding weight-untied models, and show that the required initial condition is satisfied via mild over-parameterization. Moreover, we show that the unique equilibrium point always exists during the training.

1 Introduction

Deep equilibrium models (DEQs) [1] have recently emerged as a new neural network design paradigm. A DEQ is equivalent to an infinite-depth weight-tied model with input-injection. Different from conventional (explicit) neural networks, DEQs generate features by directly solving equilibrium points of implicit equations. DEQs also have the remarkable advantage that the gradients can be computed analytically by backpropagation only through the equilibrium point with implicit differentiation. Therefore, training a DEQ only requires constant memory.

DEQs have achieved impressive performance in various ap-

plications such as computer vision [2, 3], natural language processing [1], and inverse problems [4]. Although the empirical success of DEQs has been observed in many recent studies, theoretical understandings of DEQs are still limited compared to conventional models. In this paper, we aim to establish the global convergence of the gradient descent (GD) associated with an over-parameterized DEQ, as a step towards understanding general DEQs.

A large body of work [5, 6, 7] has validated the effectiveness of over-parameterization in optimizing feedforward neural networks. The main idea is to investigate the property at initialization and bound the traveling distance of GD from the initialization [8]. However, it remains unclear whether these results can be directly applied to DEQs. The implicit weight-sharing is the key challenge. Most standard concentration tools used in previous studies fail in DEQs. This is because these analyses rely on the independence of initial random weights and features, which is no longer the case in DEQs. Moreover, it remains unknown whether over-parameterization is sufficient to guarantee the well-posedness [9, 10] of the implicit mapping of a DEQ, which is crucial to the stability of the training process, *e.g.* [11] uses an extra softmax layer to resolve the well-posedness issue and achieves a global linear convergent rate. However, this result holds only for linear DEQs and it is difficult to be extended to nonlinear activations.

We start with the gradient analysis. We observe that, in the case of DEQs, the least singular value of the equilibrium points plays a key role in the gradient dynamic. Specifically, if the least singular values of the equilibrium points at all iterations can be lower bounded by a positive constant, then one can establish a corresponding version of the Polyak-Lojasiewicz inequality for DEQs, and thus the global convergence of GD can be obtained.

Our main results are based on the following observations. Firstly, we prove the global convergence of GD by supposing an initial condition on the lower bound of the least singular value of the initial equilibrium points. The perturbation of the weight matrices is small enough to ensure that the Lipschitz constant of the implicit layer transformation is smaller than 1. This means that, the unique equilibrium point always exists throughout the training. Our second observation is that, the required initial condition holds with

a high probability for general Gaussian initialization. Note that the weight of a DEQ is implicitly re-used across layers. Thus, standard concentration inequalities, which are based on the independence of weight matrices and features, cannot be employed directly in this scenario of random DEQs analysis. In order to overcome the technical difficulty, we propose a novel probabilistic framework to approximate the empirical Gram matrix of the equilibrium point with a population Gram matrix induced from a weight-untied random network with infinite depth.

1.1 Related Work

Finite-depth Over-parameterized Feedforward Networks. Recently, over-parameterization has attracted much attention due to its effectiveness in optimizing finite-depth neural networks. [12] shows that, for smooth activation and infinite wide neural networks, the trajectory of the gradient descent (GD) method could be well-captured by a kernel called the neural tangent kernel (NTK). For a finite-width feed-forward neural network with smooth activation, [5, 6, 7] prove that the neural networks’ dynamics are strongly related to a Gram matrix. For a finite-width feed-forward neural network with ReLU activation, [13, 14, 15, 16, 17] estimate the changes of the activation patterns, and show that GD can converge to a global minimum despite the non-smoothness of activation and non-convexity of the objective function. The only condition that needs to satisfy is that the width of each layer is a polynomial of the number of training samples and the number of layers. In particular, [8] provides an alternative framework that only requires tracking the evolution of the last hidden layer rather than the activation pattern. In all previous works, a non-asymptotic analysis at initialization plays a fundamental role. Their random analysis relies on the independence between the initial random weights and features. However, the weight matrix is implicitly shared in a DEQ. Thus, previous results do not directly apply to DEQs.

Finite-depth Over-parameterized Weight-tied Neural Networks. For over-parameterized weight-tied models, [18, 19, 20] investigate NTKs of the recurrent neural networks (RNNs) with infinite width by leveraging the “Gaussian Conditioning Trick” [21, 22]. The similar technique is also used in the study on infinitely wide weight-tied autoencoders [23]. These works reveal that infinitely wide Gaussian weight-tied neural networks are essentially Gaussian processes. However, their results do not apply to the regime of finite width. [24, 25] show that, for a RNN with finite width, GD can converge to a global minimum if the width of each layer is a polynomial of the number of training samples and the number of layers. However, their results do not apply to DEQs. This is because most upper bounds for the norm of the hidden units fail as the depth approaches infinity. Thus, a new analysis framework for

the convergence of the implicit models is urgently needed.

Over-parameterized DEQs. The investigation of over-parameterized DEQs is still in the initial stage. [26] considers the NTK of DEQs with infinite width. This study claims that the NTK of DEQs is equivalent to the corresponding weight-untied models in the regime of infinite width. This study also shows that DEQs have non-degenerate NTKs even in the infinite depth. This observation is similar with ours. However, the analysis on finite-width DEQs is not studied in [26]. To conduct non-asymptotic analysis on finite-width DEQs, one of the key challenges is to estimate the least singular value of the equilibrium points. To the best of our knowledge, this problem has not been addressed in general settings. Two very recent works [27] and [28] consider simple cases in which the problem can be transferred into the estimation of the least singular value of features obtained from a single layer explicit network. Their results hold for specific predictions or special initialization. In contrast, we propose a novel probabilistic framework to estimate the least singular value of the equilibrium points in general settings. Please see detailed comparisons in the discussion of Section 4.

1.2 Contributions

Our contributions are summarized as follows.

- (1) We propose a novel probabilistic framework to analyze DEQs. Our framework addresses the technical challenges arising from the weight-sharing and the infinite depth. To the best of our knowledge, this is the first time a fine-grained non-asymptotic analysis on DEQs has been performed in a general setting.
- (2) We analyze the gradient dynamics of DEQs with the quadratic loss function. Under an initial condition on the least singular value of the initial equilibrium points, we prove that the gradient descent converges to a global optimum at a linear rate. Based on our initial analysis, we show that the required initial condition is satisfied via mild over-parameterization.
- (3) We show that the unique equilibrium point of an over-parameterized DEQ always exists throughout the training process, even without using any normalization or re-parameterization method.

2 Preliminaries

Notations. We use $\mathcal{N}(0, I)$ to denote the standard Gaussian distribution. We let $[n] \triangleq [1, \dots, n]$. For a vector v , $\|v\|_2$ is the Euclidean norm of v . For a matrix A , we use A_{ij} denote its (i, j) -th entry. We use $\|A\|_F$ to denote the Frobenius norm and $\|A\|_2$ to denote the operator norm. If a matrix is positive semi-definite, we use $\lambda_{\min}(A)$

and $\sigma_{\min}(\mathbf{A})$ to denote its least eigenvalue and singular value, respectively. We let $\mathcal{O}(\cdot)$, $\Theta(\cdot)$ and $\Omega(\cdot)$ denote standard Big-O, Big-Theta, and Big-Omega notations, respectively. We use $\phi(\cdot)$ to denote the ReLU function, namely $\phi(x) = \max(x, 0)$.

2.1 Problem Setup

We define a vanilla deep equilibrium model (DEQ) with the transform at the l -th layer as

$$\mathbf{Z}^{(l)} = \phi(\mathbf{W}\mathbf{Z}^{(l-1)} + \mathbf{U}\mathbf{X}), \quad (1)$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ denotes the training inputs, $\mathbf{U} \in \mathbb{R}^{m \times d}$ and $\mathbf{W} \in \mathbb{R}^{m \times m}$ are trainable weight matrices, and $\mathbf{Z}^{(l)} \in \mathbb{R}^{m \times n}$ is the output feature at the l -th hidden layer. The output of the last hidden layer is defined by $\mathbf{Z}^* \triangleq \lim_{l \rightarrow \infty} \mathbf{Z}^{(l)}$. Therefore, instead of running the infinitely deep layer-by-layer forward propagation, \mathbf{Z}^* can be calculated by directly solving the equilibrium point of the following equation

$$\mathbf{Z}^* = \phi(\mathbf{W}\mathbf{Z}^* + \mathbf{U}\mathbf{X}). \quad (2)$$

Let $\mathbf{y} = [y_1, \dots, y_n] \in \mathbb{R}^n$ denote the labels, and $\hat{\mathbf{y}}(\boldsymbol{\theta}) = \mathbf{a}^\top \mathbf{Z}^*$ be the prediction function with $\mathbf{a} \in \mathbb{R}^m$ being a trainable vector and $\boldsymbol{\theta} = \text{vec}(\mathbf{W}, \mathbf{U}, \mathbf{a})$. The object of our interest is the empirical risk minimization problem with the quadratic loss function

$$\Phi(\boldsymbol{\theta}) = \frac{1}{2} \|\hat{\mathbf{y}}(\boldsymbol{\theta}) - \mathbf{y}\|_2^2.$$

To do so, we consider the gradient descent (GD) update $\boldsymbol{\theta}_{\tau+1} = \boldsymbol{\theta}_\tau - \eta \nabla \Phi(\boldsymbol{\theta}_\tau)$, where η is the learning rate and $\boldsymbol{\theta}_\tau = \text{vec}(\mathbf{W}(\tau), \mathbf{U}(\tau), \mathbf{a}(\tau))$ is the parameter we optimize over at step τ . For notational simplicity, we omit the superscribe and denote \mathbf{Z} to be the equilibrium \mathbf{Z}^* when it is clear from the context. Moreover, the Gram matrix of the equilibrium point is defined by $\mathbf{G}(\tau) \triangleq \mathbf{Z}(\tau)^\top \mathbf{Z}(\tau)$ and we denote its least eigenvalue as $\lambda_\tau = \lambda_{\min}(\mathbf{G}(\tau))$.

In this paper, we make the following assumptions on the random initialization and the input data.

Assumption 1 (Random initialization). Take $\sigma_w^2 < 1/8$. We assume that \mathbf{W} is initialized with an $m \times m$ matrix with i.i.d. entries $\mathbf{W}_{ij} \sim \mathcal{N}(0, 2\sigma_w^2/m)$, \mathbf{U} is initialized with an $m \times d$ matrix with i.i.d. entries $\mathbf{U}_{ij} \sim \mathcal{N}(0, 2/d)$, \mathbf{a} is initialized with a random vector with i.i.d. entries $\mathbf{a} \sim \mathcal{N}(0, 1/m)$.

Assumption 2 (Input data). We assume that (i) $\|\mathbf{x}_i\|_2 = \sqrt{d}$, for all $i \in [n]$, and $\mathbf{x}_i \not\parallel \mathbf{x}_j$, for each pair $i \neq j \in [n]$, (ii) the labels satisfy $|y_i| = \mathcal{O}(1)$ for all $i \in [n]$.

2.2 Well-Posedness and Gradients

Well-Posedness. For the stability of the training of the DEQs, it is crucial to guarantee the existence and uniqueness of the equilibrium points [9, 29]. It is equivalent to

guarantee the well-posedness of the transformation defined in Eq. (1). In order to ensure the well-posedness, it suffices to take $\|\mathbf{W}(\tau)\|_2 < 1$, for all $\tau \geq 0$, with which Eq. (1) becomes a *contractive mapping*. From Lemma 1, we know that $\|\mathbf{W}(0)\|_2 < 1$ holds with a high probability under Assumption 1. Lemma 1 is a consequence of standard bounds concerning the singular values of Gaussian random matrices [30].

Lemma 1. Let \mathbf{W} be an $m \times m$ random matrix with i.i.d. entries $\mathbf{W}_{ij} \sim \mathcal{N}(0, \frac{2\sigma_w^2}{m})$. With probability at least $1 - \exp\{-\Omega(m)\}$, it holds that $\|\mathbf{W}\|_2 \leq 2\sqrt{2}\sigma_w$.

In Section 3, we show that the condition of $\|\mathbf{W}(\tau)\|_2 < 1$ always holds for $\tau \geq 0$. It is worth mentioning that the constraint on the spectral norm of $\mathbf{W}(\tau)$ can be lightened by that on the spectral radius of $\mathbf{W}(\tau)$ through special reparameterization methods [9, 10]. However, in this paper, we do not make extra assumptions on specific structures of weight matrices, and thus our constraint on the spectral norm of $\mathbf{W}(\tau)$ is in general mild.

Gradients. The gradients of conventional neural networks are usually computed via backpropagation through all the intermediate layers. On the contrary, the gradients w.r.t. parameters of a DEQ are computed analytically via backpropagation only through the equilibrium point \mathbf{Z} by applying the *implicit function theorem*. Specifically, note that the equilibrium point of Eq. (2) is the root of the function $F(\tau) \triangleq \mathbf{Z}(\tau) - \phi(\mathbf{W}(\tau)\mathbf{Z}(\tau) + \mathbf{U}(\tau)\mathbf{X}(\tau))$. Let $\mathbf{J}(\tau) \triangleq \partial \text{vec}(F(\tau)) / \partial \text{vec}(\mathbf{Z}(\tau))$ denote the Jacobian matrix. One can derive that

$$\mathbf{J}(\tau) = \mathbf{I}_{mn} - \mathbf{D}(\tau) (\mathbf{I}_n \otimes \mathbf{W}(\tau)),$$

where $\mathbf{D}(\tau) \triangleq \text{diag}[\text{vec}(\phi'(\mathbf{W}(\tau)\mathbf{Z}(\tau) + \mathbf{U}(\tau)\mathbf{X}(\tau)))]$. Using the Lipschitz property of ReLU, it is easy to check that $\mathbf{J}(\tau)$ is invertible if $\|\mathbf{W}(\tau)\|_2 < 1$. The gradient of each trainable parameter is given by the following lemma¹.

Lemma 2. If $\mathbf{J}(\tau)$ is invertible, the gradient of the objective function $\Phi(\tau)$ w.r.t. each trainable parameters is given by

$$\begin{cases} \text{vec}(\nabla_{\mathbf{W}} \Phi(\tau)) = (\mathbf{Z}(\tau) \otimes \mathbf{I}_m) \mathbf{R}(\tau)^\top (\hat{\mathbf{y}}(\tau) - \mathbf{y}) \\ \text{vec}(\nabla_{\mathbf{U}} \Phi(\tau)) = (\mathbf{X}(\tau) \otimes \mathbf{I}_m) \mathbf{R}(\tau)^\top (\hat{\mathbf{y}}(\tau) - \mathbf{y}) \\ \nabla_{\mathbf{a}} \Phi(\tau) = \mathbf{Z}(\tau) (\hat{\mathbf{y}}(\tau) - \mathbf{y}) \end{cases}, \quad (3)$$

where $\mathbf{R}(\tau) = (\mathbf{a}(\tau) \otimes \mathbf{I}_n) \mathbf{J}(\tau)^{-1} \mathbf{D}(\tau)$.

2.3 Polyak-Lojasiewicz Inequalities

Polyak-Lojasiewicz (PL) inequality [31] is a commonly used recipe to prove linear convergence of GD algorithms [17, 8]. In order to obtain a corresponding version

¹To simplify the notation, we omit the parameter $\boldsymbol{\theta}$ and write just $\Phi(\tau)$ and $\hat{\mathbf{y}}(\tau)$.

of PL inequalities for DEQs, our starting observation is that

$$\|\nabla_{\theta}\Phi(\tau)\|_2^2 \geq 2\lambda_{\min}(\mathbf{H}(\tau))\Phi(\tau), \quad (4)$$

where $\mathbf{H}(\tau) = \mathbf{H}_1(\tau) + \mathbf{H}_2(\tau) + \mathbf{H}_3(\tau)$ is a sum of three positive semi-definite matrices defined as

$$\begin{cases} \mathbf{H}_1(\tau) = \mathbf{G}(\tau) \\ \mathbf{H}_2(\tau) = \mathbf{R}(\tau)(\mathbf{G}(\tau) \otimes \mathbf{I}_m)\mathbf{R}(\tau)^\top, \\ \mathbf{H}_3(\tau) = \mathbf{R}(\tau)(\mathbf{X}^\top \mathbf{X} \otimes \mathbf{I}_m)\mathbf{R}(\tau)^\top \end{cases},$$

Eq. (4) is a direct application of Lemma 2. It suggests that if $\lambda_{\min}(\mathbf{H}(\tau))$ can be lower bounded away from zero, both at initialization and throughout the training, then one can establish a PL inequality that holds for the loss function, and thus GD converges to a global minimum. However, it is technically difficult to directly estimate the lower bound of $\lambda_{\min}(\mathbf{H}(\tau))$ because (1) for $\tau = 0$, $\mathbf{H}(0)$ involves both the sum and multiplication of random matrices with complex structures, and (2) for $\tau > 0$, $\mathbf{H}(\tau)$ involves the derivatives of ReLU at activation neurons which requires the estimation related to the changes of the activation patterns [32, 33].

To make the problem tractable, we further observe that $\lambda_{\min}(\mathbf{H}(\tau)) \geq \lambda_\tau$, i.e. the least eigenvalue of the Gram matrix of the equilibrium point. Applying this observation to Eq. (4), one obtains

$$\|\nabla_{\theta}\Phi(\tau)\|_2^2 \geq 2\lambda_\tau\Phi(\tau). \quad (5)$$

This means that, in order to obtain a PL-like inequality for DEQs, it suffices to bound the changes of $\mathbf{G}(\tau)$ throughout the training if λ_0 is bounded away from zero at initialization. In Section 3, we show that it holds $\lambda_\tau \geq \frac{1}{2}\lambda_0$ for every $\tau > 0$, and we have a PL inequality for DEQs: $\|\nabla_{\theta}\Phi(\tau)\|_2^2 \geq \lambda_0\Phi(\tau)$. Based on this, we prove that GD converges to a global optimum at a linear rate. The result of Section 3 is built upon an initial condition on λ_0 , and we further demonstrate that such an initial condition can be satisfied with mild over-parameterization.

2.4 Challenges in initial analysis

The initial condition on the lower bound of λ_0 plays a fundamental role in our convergence result. It is hard to estimate λ_0 directly. A common way is to build a concentration inequality between the initial *empirical* Gram matrix \mathbf{G} and the corresponding *population* Gram matrix with easily estimated least eigenvalue. In the case of DEQs, we consider a population Gram matrix \mathbf{K} defined as follows.

Definition 1. We define the population Gram matrices

$\mathbf{K}^{(l)}$ of each layer recursively as

$$\begin{aligned} \mathbf{K}^{(0)} &= 0 \\ \Lambda_{ij}^{(l)} &= \begin{bmatrix} \sigma_w^2 \mathbf{K}_{ii}^{(l-1)} + 1 & \sigma_w^2 \mathbf{K}_{ij}^{(l-1)} + \frac{1}{d} \mathbf{x}_i^\top \mathbf{x}_j \\ \sigma_w^2 \mathbf{K}_{ji}^{(l-1)} + \frac{1}{d} \mathbf{x}_j^\top \mathbf{x}_i & \sigma_w^2 \mathbf{K}_{jj}^{(l-1)} + 1 \end{bmatrix}, \\ \mathbf{K}_{ij}^{(l)} &= 2\mathbb{E}_{(\mathbf{u}, \mathbf{v})^\top \sim \mathcal{N}(0, \Lambda_{ij}^{(l)})} [\phi(\mathbf{u})\phi(\mathbf{v})] \end{aligned} \quad (6)$$

for $l \geq 1$ and $(i, j) \in [n] \times [n]$. Letting $l \rightarrow \infty$, we define $\mathbf{K} \triangleq \mathbf{K}^{(\infty)}$ and $\lambda_* \triangleq \lambda_{\min}(\mathbf{K})$.

The population Gram matrix \mathbf{K} is induced from an *infinite-depth weight-untied* model. The convergence of $\mathbf{K}^{(l)}$ for $l \rightarrow \infty$ and the positive definiteness of \mathbf{K} are deferred to Section 4.1. The least eigenvalue of \mathbf{K} is the fundamental quantity that determines the lower bound of λ_0 .

Non-asymptotic analysis on DEQs is more difficult than explicit models. The weight-sharing is the key technical challenge, and one cannot resort to standard concentration tools. Specifically, note that initial \mathbf{G}_{ij} is implicitly defined as $\mathbf{z}_i^\top \mathbf{z}_j = \phi([\mathbf{W}, \mathbf{U}][\mathbf{z}_i^\top, \mathbf{x}_i^\top]^\top) \phi([\mathbf{W}, \mathbf{U}][\mathbf{z}_j^\top, \mathbf{x}_j^\top]^\top)^2$. On one hand, one cannot directly apply the standard inequality like previous works. This is because they rely on the independence between initial weight matrices and features, which is no longer the case in DEQs. On the other hand, one cannot directly use the standard ε -net argument [30]. Note that $[\mathbf{W}, \mathbf{U}]$ is a *short* matrix, i.e. $[\mathbf{W}, \mathbf{U}] \in \mathbb{R}^{m \times (m+d)}$, and the size of ε -net for $[\mathbf{z}_i^\top, \mathbf{x}_i^\top]$ is too large for us to derive a mild over-parameterization condition. In order to overcome the technical challenge, we propose a novel probabilistic framework in Section 4 by introducing the “new fresh randomness” [24] and perform a fine-grained non-asymptotic analysis on random DEQs.

3 Main Result

3.1 Global Convergence under an Initial Condition

Let δ be any positive constant such that $\rho_w(0) + \delta < 1$. We define the following quantities:

$$\bar{\rho}_w = \|\mathbf{W}(0)\|_2 + \delta, \bar{\rho}_u = \|\mathbf{U}(0)\|_2 + \delta, \bar{\rho}_a = \|\mathbf{a}(0)\|_2 + \delta.$$

We first present the global convergence of GD by supposing the following condition on the least eigenvalue λ_0 of the initial Gram matrix $\mathbf{G}(0)$.

Condition 1. At initialization, it holds that

$$\lambda_0 \geq \frac{4}{\delta} \max(c_w, c_u, c_a) \|\mathbf{X}\|_F \|\hat{\mathbf{y}}(0) - \mathbf{y}\|_2, \quad (7)$$

$$\lambda_0^{\frac{3}{2}} \geq 4(2 + \sqrt{2})\bar{\rho}_a^{-1} (c_w^2 + c_u^2) \|\mathbf{X}\|_F^2 \|\hat{\mathbf{y}}(0) - \mathbf{y}\|_2, \quad (8)$$

$$\lambda_0 \geq 4(c_w^2 + c_u^2) \|\mathbf{X}\|_F^2, \quad (9)$$

where $c_w = \frac{\bar{\rho}_u \bar{\rho}_a}{(1 - \bar{\rho}_w)^2}$, $c_u = \frac{\bar{\rho}_a}{1 - \bar{\rho}_w}$, $c_a = \frac{\bar{\rho}_u}{1 - \bar{\rho}_w}$.

²For notational simplicity, we denote $\mathbf{W}(0)$, $\mathbf{U}(0)$ and $\mathbf{G}(0)$ by \mathbf{W} , \mathbf{U} and \mathbf{G} , respectively.

The convergence result under Condition 1 is presented as follows.

Theorem 1. *Consider a DEQ defined in Eq. (2). Suppose that Condition 1 holds at initialization. If the learning rate satisfies*

$$\eta < \min \left(\frac{2}{\lambda_0}, \frac{2(c_w^2 + c_u^2)}{(c_w^2 + c_u^2 + c_a^2)^2 \|\mathbf{X}\|_F^2} \right), \quad (10)$$

for every $\tau \geq 0$, the following holds

- (i) $\|\mathbf{W}(\tau)\|_2 < 1$, i.e. the equilibrium points always exist,
- (ii) $\lambda_\tau > \frac{1}{2}\lambda_0$, and thus the PL condition holds as $\|\nabla_{\boldsymbol{\theta}}\Phi(\tau)\|_2^2 \geq \lambda_0\Phi(\tau)$,
- (iii) the loss converges to a global minimum as

$$\Phi(\tau) \leq \left(1 - \eta \frac{\lambda_0}{2}\right)^\tau \Phi(0).$$

Theorem 1 shows that the unique equilibrium point always exists during the training, and GD converges to a global optimum under Condition 1. The proof of this part is inspired by the framework proposed in [8, Theorem 2.2]. The complete proof is deferred to the supplementary material. Next, we discuss how these initial conditions can be fulfilled via over-parameterization.

3.2 Initial Condition

In this section, we aim to show that, Condition 1 holds under Assumptions 1 and 2 via over-parameterization. To do so, it suffices to derive a lower bound on λ_0 , and upper bounds on $\max(c_w, c_u, c_a)$, $\bar{\rho}_a^{-1}$ and the initial loss $\Phi(0)$, and put these bounds into Eq. (7)-Eq. (9) to obtain a specific condition on width m .

Firstly, we present the lower bound of λ_0 as follows.

Theorem 2. *If $m = \Omega\left(\frac{n^2}{\lambda_*^2} \left(\log \frac{n}{\lambda_* t}\right)\right)$, with probability at least $1 - t$, it holds that*

$$\lambda_0 \geq \frac{m}{2} \lambda_*.$$

The proof of Theorem 2 is based on a concentration inequality between the empirical Gram matrix \mathbf{G} and the population Gram matrix \mathbf{K} . Since that the weight matrices are implicitly reused in DEQs for infinite times (see Eq. (1)), standard concentration tools are not applicable. In order to overcome the challenge arising from the weight-sharing and the infinite depth. We present a novel probabilistic framework. It is one of our core contributions. Detailed analysis is presented in Section 4.

Secondly, by standard bounds on the operator norm of Gaussian matrices [30], we have $w.p. \geq 1 - \exp\{-\Omega(m)\}$

that, $\|\mathbf{W}(0)\|_2 = \mathcal{O}(1)$, $\|\mathbf{U}(0)\|_2 = \mathcal{O}(\sqrt{m/d})$ under Assumption 1. Thus, $w.p. \geq 1 - \exp\{-\Omega(m)\}$, it holds that

$$\bar{\rho}_w = \mathcal{O}(1), \quad \bar{\rho}_u = \mathcal{O}\left(\sqrt{\frac{m}{d}}\right)$$

which implies that

$$\max(c_w, c_u, c_a) = \mathcal{O}\left(\sqrt{\frac{m}{d}}\right).$$

By standard bounds on the norm of Gaussian vector [30], we have $w.p. \geq 1 - \exp\{-\Omega(m)\}$ that, $\bar{\rho}_a^{-1} = \mathcal{O}(1)$.

Thirdly, by using the property of the contractive mapping Eq. (1) and using the standard concentration argument, it is easy to show that $w.p. \geq 1 - t$, it holds that

$$\Phi(0) = \mathcal{O}(n).$$

Putting all these bounds into Eq. (7)-Eq. (9), one can show that $w.p. \geq 1 - t$, Condition 1 is satisfied for $m = \Omega\left(\frac{n^2}{\lambda_*^2} \left(\log \frac{n}{\lambda_* t}\right)\right)$.

4 Analysis at Initialization

For notational simplicity, we denote $\mathbf{W}(0)$ and $\mathbf{U}(0)$ by \mathbf{W} and \mathbf{U} in this section. The Gram matrices of the equilibrium point $\mathbf{Z}^{(l)}$ of the l -th layer are defined as $\mathbf{G}^{(l)} \triangleq (\mathbf{Z}^{(l)})^\top \mathbf{Z}^{(l)}$, for $l \geq 1$. Without loss of generality, we assume that $\mathbf{Z}^{(0)} = \mathbf{0}$.

Our main idea is to establish the concentration inequality between the empirical Gram matrix \mathbf{G} and the population Gram matrix \mathbf{K} . A simple case is to initialize DEQs as single-layer explicit models. Specifically, take $\sigma_w^2 = 0$ and $\mathbf{Z} = \phi(\mathbf{U}\mathbf{X})$. Using the standard Bernstein inequality, one can show that $w.p. \geq 1 - t$, $\lambda_0 \geq m\tilde{\lambda}/2$, as long as $m = \Omega\left(\tilde{\lambda}^{-2} n^2 \log(n/t)\right)$ where $\tilde{\lambda} = \lambda_{\min}(\mathbf{K}^{(1)})$. This case has been studied in many previous works [17, 16].

However, when $\sigma_w^2 > 0$, previous random analyses on explicit networks cannot be directly applied to DEQs. This is due to the fact that these analyses rely on the independence of initial random weights and features, which is no longer the case in DEQs. In order to overcome these technical difficulties, we propose a novel probabilistic framework based on the following observations: $\mathbf{K}_{ij}^{(l)}$, $\frac{1}{m}\mathbf{G}_{ij}^{(l)}$, and $\mathbf{G}_{ij}^{(l)}$ converge to \mathbf{K}_{ij} , $\mathbf{K}_{ij}^{(l)}$, and \mathbf{G}_{ij} at an exponential rate, respectively. Moreover, \mathbf{K} is strictly positive definite. These observations imply that it suffices to take sufficiently large l and m to ensure that $\|\frac{1}{m}\mathbf{G} - \mathbf{K}\|_F$ is smaller than λ_* . Thus, one can lower bound λ_0 by invoking Weyl's inequality. Detailed analysis is given as follows.

4.1 Bound between infinite-depth and finite-depth weight-untied models

In the case of $\phi = \text{ReLU}$, given any positive definite matrix $\mathbf{A} = \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}$ with $|x| \leq 1$, and two random variables $(\mathbf{u}, \mathbf{v})^\top \sim \mathcal{N}(0, \mathbf{A})$, as shown in [34], it holds that

$$\mathbb{E}[\phi(\mathbf{u})\phi(\mathbf{v})] = \frac{1}{2}Q(x), \quad (11)$$

where $Q(x) \triangleq \frac{\sqrt{1-x^2} + (\pi - \arccos x)x}{\pi}$.

Combining Eq. (11) with the homogeneity of ReLU, we can have more precise expressions of $\mathbf{K}_{ij}^{(l)}$ as follows.

Lemma 3. Let $\cos \theta_{ij}^{(l)} = \frac{\sigma_w^2 \mathbf{K}_{ij}^{(l-1)} + d^{-1} \mathbf{x}_i^\top \mathbf{x}_j}{\sigma_w^2 \mathbf{K}_{ij}^{(l-1)} + 1}$, and $\rho^{(l)} = \mathbf{K}_{ii}^{(l)}$. For $l \geq 1$ and $(i, j) \in [n] \times [n]$, $\mathbf{K}_{ij}^{(l)}$ defined in Definition 1 can be written as

$$\mathbf{K}_{ij}^{(l)} = \rho^{(l)} Q(\cos \theta_{ij}^{(l)}), \quad (12)$$

with $\rho^{(l)} = \frac{1 - \sigma_w^{2l}}{1 - \sigma_w^2}$ and

$$\cos \theta_{ij}^{(l)} = \left(1 - \frac{1}{\rho^{(l-1)}}\right) Q(\cos \theta_{ij}^{(l-1)}) + \frac{1}{\rho^{(l-1)}} \frac{\mathbf{x}_i^\top \mathbf{x}_j}{d}. \quad (13)$$

Theorem 3. Under Assumptions 1 and 2, it holds that

(i) $\|\mathbf{K} - \mathbf{K}^{(l)}\|_F = \mathcal{O}\left(nl^{\frac{1}{2}}\sigma_w^l\right)$, which implies that, for $l \rightarrow \infty$, $\mathbf{K}^{(l)}$ converges to \mathbf{K} with each entry

$$\mathbf{K}_{ij} = \frac{1}{1 - \sigma_w^2} Q(\cos \theta_{ij}), \quad (14)$$

where $\cos \theta_{ij} = \sigma_w^2 Q(\cos \theta_{ij}) + (1 - \sigma_w^2) \frac{\mathbf{x}_i^\top \mathbf{x}_j}{d}$.

(ii) \mathbf{K} is strictly positive definite, i.e. $\lambda_* > 0$.

Sketch of Proof. (i) By Eq. (6) and Eq. (12), one can show that $|\mathbf{K}_{ij}^{(l+1)} - \mathbf{K}_{ij}^{(l)}| \leq \sigma_w^2 |\mathbf{K}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l-1)}| + 2\sigma_w^{2l}$, which implies that $|\mathbf{K}_{ij} - \mathbf{K}_{ij}^{(l)}| = \mathcal{O}(l\sigma_w^l)$.

Note that $\sigma_w < 1$. Thus $\mathbf{K}^{(l)}$ converges to a deterministic Gram matrix \mathbf{K} , and one can easily obtain Eq. (14) from Eq. (13) by letting $l \rightarrow \infty$.

(ii) The proof of this part is similar with [16, 17]. By performing the Hermite analysis on Eq. (14), one can show that \mathbf{K} is strictly positive definite if $|\cos \theta_{ij}| < 1$ for all $i \neq j$. Using the fact that $|Q(x)| \leq 1$, Eq. (14) implies that $|\cos \theta_{ij}| < 1$ under Assumption 2. Please see the complete proof in the supplementary material. \square

4.2 Bound between infinite-depth and finite-depth weight-tied models

By leveraging the contractility of the transformation defined in Eq. (1) and by invoking the standard Bernstein inequality, we establish a concentration inequality between the Gram matrix of the l -th layer's output and that of the initial equilibrium point as follows.

Theorem 4. Under Assumptions 1 and 2, with probability at least $1 - n^2 \exp\{-\Omega(m)\}$, it holds

$$\frac{1}{m} \|\mathbf{G} - \mathbf{G}^{(l)}\|_F = \mathcal{O}\left(n \left(2\sqrt{2}\sigma_w\right)^l\right).$$

Sketch of Poof. Using the contractility of the implicit layer, it is easy to have $w.p. \geq 1 - \exp\{-\Omega(m)\}$, $\|\mathbf{z}_i^{(l)}\|_2 = \mathcal{O}(\|\mathbf{z}_i^{(1)}\|_2)$, and $\|\mathbf{z}_i - \mathbf{z}_i^{(l)}\|_2 = \mathcal{O}\left((2\sqrt{2}\sigma_w)^l \|\mathbf{z}_i^{(1)}\|_2\right)$. Moreover, using Bernstein inequality, we have $w.p. \geq 1 - \exp\{-\Omega(mt^2)\}$, $|\frac{1}{m}(\mathbf{z}_i^{(1)})^\top \mathbf{z}_i^{(1)} - 1| \leq t$. Let t be an absolute positive constant. Theorem 4 can be proved by applying the simple union bound. Please see the complete proof in the supplementary material. \square

4.3 Bound between weight-tied and weight-untied models with finite-depth

Due to the implicit weight-tied structure of DEQs, standard concentration inequalities built upon the independence cannot be directly applied. In order to overcome technical difficulties, we build necessary probabilistic tools for DEQs. We first introduce Lemma 4 which provides the ‘‘fresh new randomness’’ for our analysis. The construction method in Lemma 4 is inspired by the previous work [24] on RNNs.

Lemma 4. For $l \geq 1$, $\mathbf{G}_{ij}^{(l+1)}$ can be reconstructed as $\mathbf{G}_{ij}^{(l+1)} = \phi(\mathbf{M}\mathbf{h})^\top \phi(\mathbf{M}\mathbf{h}')$ such that

(i) $\mathbf{h}^\top \mathbf{h}' = \frac{\sigma_w^2}{m} \mathbf{G}_{ij}^{(l)} + \frac{1}{d} \mathbf{x}_i^\top \mathbf{x}_j$,

(ii) $\mathbf{M} \in \mathbb{R}^{m \times (2l+d+2)}$ is a rectangle matrix, and the entries of \mathbf{M} are i.i.d. from $\mathcal{N}(0, 2)$ conditioning on previous layers.

Proof. (i) Let $\mathbf{V}_l \in \mathbb{R}^{m \times 2l}$ denote a column orthonormal matrix using Gram-Schmidt as

$$\mathbf{V}_l = \text{GS}\left(\mathbf{z}_i^{(1)}, \dots, \mathbf{z}_i^{(l)}, \mathbf{z}_j^{(1)}, \dots, \mathbf{z}_j^{(l)}\right).$$

For each i, j , we define $\mathbf{p} \triangleq (\mathbf{I} - \mathbf{V}_{l-1} \mathbf{V}_{l-1}^\top) \mathbf{z}_i^{(l)}$ and $\mathbf{q} \triangleq (\mathbf{I} - \mathbf{V}_{l-1} \mathbf{V}_{l-1}^\top) \mathbf{z}_j^{(l)}$. We split \mathbf{q} into two parts $\mathbf{q} = \mathbf{q}^\parallel + \mathbf{q}^\perp$, where \mathbf{q}^\parallel is parallel to \mathbf{p} and \mathbf{q}^\perp is orthogonal to \mathbf{p} as

$$\mathbf{q}^\parallel = \frac{\mathbf{p}^\top \mathbf{q}}{\|\mathbf{p}\|_2^2} \mathbf{p}, \quad \mathbf{q}^\perp = \left(\mathbf{I} - \frac{\mathbf{p} \mathbf{p}^\top}{\|\mathbf{p}\|_2^2}\right) \mathbf{q}.$$

First, we construct \mathbf{M} as $\mathbf{M} = [\mathbf{M}_1, \mathbf{M}_2, \mathbf{M}_3, \mathbf{M}_4]$ with

$$\begin{aligned}\mathbf{M}_1 &= \sigma_w^{-1} \sqrt{m} \mathbf{W} \mathbf{V}_{l-1}, & \mathbf{M}_2 &= \sigma_w^{-1} \sqrt{m} \frac{\mathbf{W} \mathbf{p}}{\|\mathbf{p}\|_2}, \\ \mathbf{M}_3 &= \sigma_w^{-1} \sqrt{m} \frac{\mathbf{W} \mathbf{q}}{\|\mathbf{q}\|_2}, & \mathbf{M}_4 &= \sqrt{d} \mathbf{U}.\end{aligned}$$

Then, we construct $\mathbf{h} = [\mathbf{h}_1^\top, \mathbf{h}_2^\top, \mathbf{h}_3^\top, \mathbf{h}_4^\top]^\top$ with

$$\begin{aligned}\mathbf{h}_1 &= \frac{\sigma_w}{\sqrt{m}} \mathbf{V}_{l-1}^\top \mathbf{z}_i^{(l)}, & \mathbf{h}_2 &= \frac{\sigma_w}{\sqrt{m}} \|\mathbf{p}\|_2, \\ \mathbf{h}_3 &= 0, & \mathbf{h}_4 &= \frac{1}{\sqrt{d}} \mathbf{x}_i,\end{aligned}$$

and $\mathbf{h}' = [\mathbf{h}'_1{}^\top, \mathbf{h}'_2{}^\top, \mathbf{h}'_3{}^\top, \mathbf{h}'_4{}^\top]^\top$ with

$$\begin{aligned}\mathbf{h}'_1 &= \frac{\sigma_w}{\sqrt{m}} \mathbf{V}_{l-1}^\top \mathbf{z}_j^{(l)}, & \mathbf{h}'_2 &= \frac{\sigma_w \mathbf{p}^\top \mathbf{q}}{\sqrt{m} \|\mathbf{p}\|_2}, \\ \mathbf{h}'_3 &= \frac{\sigma_w}{\sqrt{m}} \|\mathbf{q}\|_2, & \mathbf{h}'_4 &= \frac{1}{\sqrt{d}} \mathbf{x}_j.\end{aligned}$$

It is easy to check that $\mathbf{G}_{ij}^{(l+1)} = \phi(\mathbf{M} \mathbf{h})^\top \phi(\mathbf{M} \mathbf{h}')$ with $\mathbf{h}^\top \mathbf{h}' = \frac{\sigma_w^2}{m} \mathbf{G}_{ij}^{(l)} + \frac{1}{d} \mathbf{x}_i^\top \mathbf{x}_j$.

(ii) Let $\mathbf{V}_{l-1} \triangleq [\mathbf{v}_1, \dots, \mathbf{v}_{2(l-1)}]$. Note that \mathbf{v}_j only depends on the randomness of \mathbf{U} and $\mathbf{W}[\mathbf{v}_1, \dots, \mathbf{v}_{j-1}]$. This means that, conditioning on \mathbf{U} and $\mathbf{W}[\mathbf{v}_1, \dots, \mathbf{v}_{j-1}]$, $\mathbf{W} \mathbf{v}_j$ is still an independent Gaussian vector. Similarly, we have $\mathbf{W} \mathbf{p}$ and $\mathbf{W} \mathbf{q}$ are also independent Gaussian vectors. Consequently, we prove that the entries of \mathbf{M} are *i.i.d.* from $\mathcal{N}(0, 2)$. \square

We stress that although \mathbf{M} constructed in Lemma 4 has *i.i.d.* entries, it still depends on \mathbf{h} and \mathbf{h}' . Thus, the standard Bernstein inequality cannot be directly applied. To address this issue, we perform the standard ε -argument [30] and leverage the “fresh new randomness” [24] provided in Lemma 4. The concentration inequality between $\frac{1}{m} \mathbf{G}^{(l)}$ and $\mathbf{K}^{(l)}$ is established as follows.

Theorem 5. *Under Assumptions 1 and 2, with probability at least $1 - n^2 \exp\{-\Omega(m 8^l \sigma_w^{2l}) + \mathcal{O}(l^2)\}$, it holds that*

$$\left\| \frac{1}{m} \mathbf{G}^{(l)} - \mathbf{K}^{(l)} \right\|_F \leq n \left(2\sqrt{2} \sigma_w \right)^l.$$

Sketch of Proof. We give the main ideas of the proof.

It holds that $\left| \frac{1}{m} \mathbf{G}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l)} \right| \leq \left| \frac{1}{m} \mathbf{G}_{ij}^{(l)} - \mathbb{E} \left[\frac{1}{m} \mathbf{G}_{ij}^{(l)} \right] \right| + \left| \mathbb{E} \left[\frac{1}{m} \mathbf{G}_{ij}^{(l)} \right] - \mathbf{K}_{ij}^{(l)} \right|$ by the triangle inequality. Following Lemma 4, we reconstruct $\mathbf{G}_{ij}^{(l)} = \phi(\mathbf{M} \mathbf{h})^\top \phi(\mathbf{M} \mathbf{h}')$.

(i) For *fixed* \mathbf{h} and \mathbf{h}' , using the Bernstein inequality, and one can show that *w.p.* $\geq 1 - \exp\{-\Omega(m \varepsilon^2)\}$, it holds that

$$\left| \frac{1}{m} \mathbf{G}_{ij}^{(l)} - \mathbb{E} \left[\frac{1}{m} \mathbf{G}_{ij}^{(l)} \right] \right| \leq \varepsilon.$$

(ii) For *all* \mathbf{h} and \mathbf{h}' , we apply the standard ε -net argument. Note that the size of ε -net for \mathbf{h}, \mathbf{h}' is at most $\exp\{\mathcal{O}(l \log \frac{1}{\varepsilon})\}$. Thus, one can derive that *w.p.* $\geq 1 - \exp\{-\Omega(m \varepsilon^2) + \mathcal{O}(l \log \frac{1}{\varepsilon})\}$, it holds that

$$\left| \frac{1}{m} \mathbf{G}_{ij}^{(l)} - \mathbb{E} \left[\frac{1}{m} \mathbf{G}_{ij}^{(l)} \right] \right| \leq \varepsilon.$$

(iii) Substitute the choice of \mathbf{h} and \mathbf{h}' such that $\mathbf{h}^\top \mathbf{h}' = \frac{\sigma_w^2}{m} \mathbf{G}_{ij}^{(l-1)} + \frac{1}{d} \mathbf{x}_i^\top \mathbf{x}_j$. Using the fact that $\mathbb{E} \left[\frac{1}{m} \mathbf{G}_{ij}^{(l)} \right]$ is determined by $\mathbf{h}^\top \mathbf{h}'$, one can further derive that *w.p.* $\geq 1 - \exp\{-\Omega(m \varepsilon^2) + \mathcal{O}(l \log \frac{1}{\varepsilon})\}$,

$$\left| \mathbb{E} \left[\frac{1}{m} \mathbf{G}_{ij}^{(l)} \right] - \mathbf{K}_{ij}^{(l)} \right| \leq \sigma_w^2 \left| \frac{1}{m} \mathbf{G}_{ij}^{(l-1)} - \mathbf{K}_{ij}^{(l-1)} \right| + \varepsilon,$$

which implies that

$$\left| \frac{1}{m} \mathbf{G}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l)} \right| \leq \sigma_w^2 \left| \frac{1}{m} \mathbf{G}_{ij}^{(l-1)} - \mathbf{K}_{ij}^{(l-1)} \right| + 2\varepsilon.$$

Note that $\sigma_w^2 < 1/8$, and thus a simple induction argument works here.

Lastly, let $\varepsilon = (2\sqrt{2}\sigma_w)^l$, and Theorem 5 can be proved by using the simple union bound. Please see the complete proof in the supplementary material. \square

4.4 Proof of Theorem 2

Firstly, combining Theorems 3, 4 and 5, one can show that *w.p.* $\geq 1 - n^2 \exp\{-\Omega(m 8^l \sigma_w^{2l}) + \mathcal{O}(l^2)\}$, it holds

$$\begin{aligned}& \left\| \frac{1}{m} \mathbf{G} - \mathbf{K} \right\|_F \\ & \leq \frac{1}{m} \left\| \mathbf{G} - \mathbf{G}^{(l)} \right\|_F + \left\| \frac{1}{m} \mathbf{G}^{(l)} - \mathbf{K}^{(l)} \right\|_F + \left\| \mathbf{K} - \mathbf{K}^{(l)} \right\|_F \\ & = \mathcal{O} \left(n \left(2\sqrt{2} \sigma_w \right)^l \right) + \mathcal{O} \left(n \left(2\sqrt{2} \sigma_w \right)^l \right) + \mathcal{O} \left(n l^{\frac{1}{2}} \sigma_w^l \right) \\ & = \mathcal{O} \left(n \left(2\sqrt{2} \sigma_w \right)^l \right),\end{aligned}$$

where the last equality comes from the fact that $l^{\frac{1}{2}} \leq (2\sqrt{2})^l$, for $l \geq 1$.

Next, we fix l to omit the explicit dependence in l . Specifically, we take $l = \Theta(\log(\lambda_*^{-1}n)/\log(\sqrt{2}/4\sigma_w))$. The lower bound of l is large enough to ensure that $\left\| \frac{1}{m} \mathbf{G} - \mathbf{K} \right\|_F \leq \frac{\lambda_*}{2}$. Therefore, by Weyl's inequality, it holds that $\lambda_0 > \frac{m}{2} \lambda_*$. Meanwhile, the upper bound of l guarantees that the probability does not decrease exponentially. Thus, one can use a mild over-parameterization condition on m to ensure the high probability.

Consequently, one can show that, it holds *w.p.* $\geq 1 - t$, $\lambda_0 \geq \frac{m}{2} \lambda_*$, as long as $m = \Omega \left(\frac{n^2}{\lambda_*^2} \left(\log \frac{n}{\lambda_* t} \right) \right)$.

Discussion on comparisons with [27] and [28]. In [27], the prediction of an implicit network is formulated as a weighted summation of equilibrium points and explicit features *i.e.* $\hat{\mathbf{y}} = \mathbf{a}^\top \mathbf{Z} + \mathbf{b}^\top \phi(\mathbf{U} \mathbf{X})$. Therefore, it is hard to measure the contribution of the equilibrium point to the capacity of implicit models in their case. In contrast, we only use equilibrium points for predictions, and our result illustrates that arbitrary training label can be fitted only using

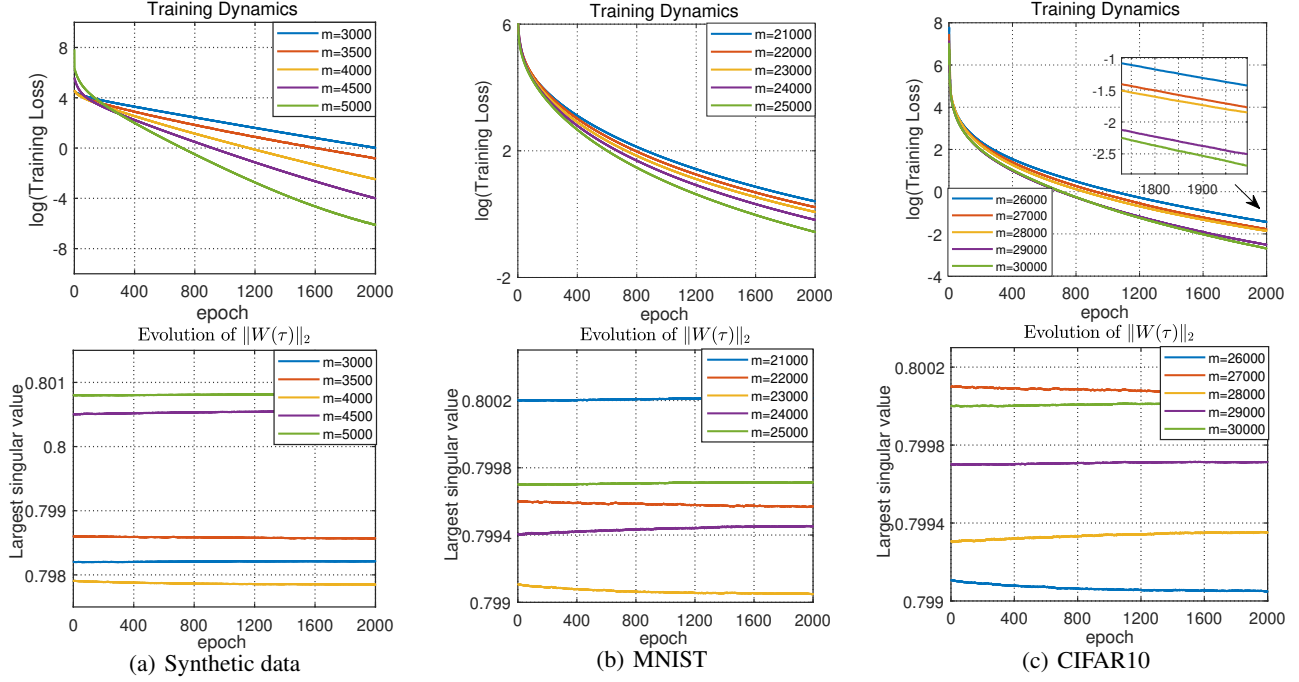


Figure 1: Results of different widths on (a) Synthetic data; (b) MNIST; (c) CIFAR10.

equilibrium points. In a concurrent work [28] on ReLU implicit networks, they consider “a subset of initialization” which requires entries of \mathbf{W} to be non-negative. In their case, one can show that $\mathbf{Z} = (\mathbf{I} - \mathbf{W})^{-1}\phi(\mathbf{W}\mathbf{X})$ at initialization. This implies that the initial equilibrium point is essentially a linear transformation of the explicit feature $\phi(\mathbf{U}\mathbf{X})$. In contrast, we consider general Gaussian initialization which is commonly used in practice. The results in [27, 28] can be built upon the lower bound of the least singular value of $\phi(\mathbf{U}\mathbf{X})$, which is given by previous works on explicit models. However, previous analyzes do not applied in our case. In order to address the technical problem, we propose a novel probabilistic framework. Therefore, our studies have distinct differences.

5 Numerical Experiments

In this section, we implement several numerical experiments to verify our main theoretical conclusions. We first evaluate our method on various datasets including synthetic data, MNIST, and CIFAR10. For constructing synthetic data, we uniformly generate $n = 1000$ data points from a $d = 1000$ dimensional sphere with radius \sqrt{d} , and labels are generated from a one-dimensional standard Gaussian distribution. For each dataset of MNIST and CIFAR10, we randomly sample 500 images from each of class 0 and class 1 to generate the training dataset with $n = 1000$ samples. We use Gaussian initialization as suggested in Assumption 1 and σ_w^2 is set as 0.08. We normalize each data point as suggested in Assumption 2.

In the first experiment, we test how the width affects the convergence rate. As shown in Figure 1, the convergence speed becomes faster as m increases, and the final training loss becomes smaller. We believe that the reason is as m increases, Gram matrices become more stable. The second experiment verifies that $\|\mathbf{W}(\tau)\|_2$ is smaller than 1 throughout the training, which implies that the unique equilibrium point always exists.

6 Conclusion

In this paper, we analyze the gradient dynamics of DEQs with the quadratic loss function. Under a specific initial condition, we prove that the GD converges to a global optimum at a linear rate. By performing a fine-grained analysis on Gaussian initialized DEQs, we further show that the initial conditions are satisfied via mild over-parameterization. Specifically, we present a new probabilistic framework to address the challenge arising from the weight-sharing and the infinite depth. To the best of our knowledge, it is the first time to analyze the equilibrium point of a random DEQ. Moreover, we show that the unique equilibrium points always exist during the training process. Our analysis is specific to ReLU DEQs. For future research, it would be interesting to generalize the result to other nonlinear activations. Moreover, it would be interesting to explore the generalization performance of over-parameterized DEQs based on our analysis at initialization.

References

- [1] Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. Deep equilibrium models. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [2] Shaojie Bai, Vladlen Koltun, and J Zico Kolter. Multiscale deep equilibrium models. *Advances in Neural Information Processing Systems*, 2020.
- [3] Xingyu Xie, Qiu hao Wang, Zenan Ling, Xia Li, Guangcan Liu, and Zhou chen Lin. Optimization induced equilibrium networks: An explicit optimization perspective for understanding equilibrium models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [4] Davis Gilton, Gregory Ongie, and Rebecca Willett. Deep equilibrium architectures for inverse problems in imaging. *arXiv preprint arXiv:2102.07944*, 2021.
- [5] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing Systems*, pages 8141–8150, 2019.
- [6] 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*, pages 3003–3048, 2019.
- [7] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in neural information processing systems*, 2018.
- [8] Quynh Nguyen. On the proof of global convergence of gradient descent for deep relu networks with linear widths. In *International Conference on Machine Learning*, pages 8056–8062. PMLR, 2021.
- [9] Ezra Winston and J Zico Kolter. Monotone operator equilibrium networks. In *Advances in Neural Information Processing Systems*, 2020.
- [10] Max Revay, Ruigang Wang, and Ian R Manchester. Lipschitz bounded equilibrium networks. *arXiv preprint arXiv:2010.01732*, 2020.
- [11] Kenji Kawaguchi. On the theory of implicit deep learning: Global convergence with implicit layers. In *International Conference on Learning Representations*, 2020.
- [12] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, pages 8580–8589, 2018.
- [13] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via overparameterization. In *International Conference on Machine Learning*, pages 242–252. PMLR, 2019.
- [14] Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural networks. In *Advances in neural information processing systems*, pages 2055–2064, 2019.
- [15] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-parameterized deep relu networks. *Machine Learning*, 109(3):467–492, 2020.
- [16] Samet Oymak and Mahdi Soltanolkotabi. Toward moderate overparameterization: Global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 1(1):84–105, 2020.
- [17] Quynh N Nguyen and Marco Mondelli. Global convergence of deep networks with one wide layer followed by pyramidal topology. *Advances in Neural Information Processing Systems*, 33:11961–11972, 2020.
- [18] Greg Yang. Wide feedforward or recurrent neural networks of any architecture are gaussian processes. *Advances in Neural Information Processing Systems*, 32, 2019.
- [19] Greg Yang. Tensor programs ii: Neural tangent kernel for any architecture. *arXiv preprint arXiv:2006.14548*, 2020.
- [20] Greg Yang. Tensor programs iii: Neural matrix laws. *arXiv preprint arXiv:2009.10685*, 2020.
- [21] Erwin Bolthausen. An iterative construction of solutions of the tap equations for the sherrington–kirkpatrick model. *Communications in Mathematical Physics*, 325(1):333–366, 2014.
- [22] Mohsen Bayati and Andrea Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Transactions on Information Theory*, 57(2):764–785, 2011.
- [23] Ping Li and Phan-Minh Nguyen. On random deep weight-tied autoencoders: Exact asymptotic analysis, phase transitions, and implications to training. In *International Conference on Learning Representations*, 2018.
- [24] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. On the convergence rate of training recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 6676–6688, 2019.

- [25] Lifu Wang, Bo Shen, Bo Hu, and Xing Cao. On the provable generalization of recurrent neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- [26] Zhili Feng and J Zico Kolter. On the neural tangent kernel of equilibrium models. *arxiv*, 2020.
- [27] T. Gao, H. Liu, J. Liu, H. Rajan, and H. Gao. A global convergence theory for deep relu implicit networks via over-parameterization. *ICLR*, 2022.
- [28] T. Gao and H. Gao. Gradient descent optimizes infinite-depth relu implicit networks with linear widths. *arxiv*, 2022.
- [29] Laurent El Ghaoui, Fangda Gu, Bertrand Travacca, Armin Askari, and Alicia Tsai. Implicit deep learning. *SIAM Journal on Mathematics of Data Science*, 3(3):930–958, 2021.
- [30] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [31] B. T. Polyak. Gradient methods for minimizing functionals. *Zh.vychisl.mat.mat.fiz.*, 3:643–653, 1963.
- [32] Simon S Du, X. Zhai, B. Póczos, and A. Singh. Gradient descent provably optimizes over-parameterized neural networks. *ICLR*, 2018.
- [33] D. Zou and Q. Gu. An improved analysis of training over-parameterized deep neural networks. *NeurIPS*, 2019.
- [34] Amit Daniely, Roy Frostig, and Yoram Singer. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.