
A Convergence Analysis of Gradient Descent on Graph Neural Networks

Pranjal Awasthi
Google Research
pranjalawasthi@google.com

Abhimanyu Das
Google Research
abhidas@google.com

Sreenivas Gollapudi
Google Research
sgollapu@google.com

Abstract

Graph Neural Networks (GNNs) are a powerful class of architectures for solving learning problems on graphs. While many variants of GNNs have been proposed in the literature and have achieved strong empirical performance, their theoretical properties are less well understood. In this work we study the convergence properties of the gradient descent algorithm when used to train GNNs. In particular, we consider the realizable setting where the data is generated from a network with unknown weights and our goal is to study conditions under which gradient descent on a GNN architecture can recover near optimal solutions. While such analysis has been performed in recent years for other architectures such as fully connected feed-forward networks, the message passing nature of the updates in a GNN poses a new challenge in understanding the nature of the gradient descent updates. We take a step towards overcoming this by proving that for the case of deep linear GNNs gradient descent provably recovers solutions up to error ϵ in $O(\log(1/\epsilon))$ iterations, under natural assumptions on the data distribution. Furthermore, for the case of one-round GNNs with ReLU activations, we show that gradient descent provably recovers solutions up to error ϵ in $O(\frac{1}{\epsilon^2} \log(\frac{1}{\epsilon}))$ iterations.

1 Introduction

In the last decade, deep neural networks have been successfully used for a variety of machine learning tasks. In particular, optimization of various deep neural network architectures using gradient descent has been empirically (and in some cases, theoretically) shown to work surprisingly well, with an ability to obtain small training error on a variety of problems and datasets. In this paper, we focus on a particular class of neural network architectures known as Graph Neural Networks (GNNs).

Graph Neural Networks (GNNs) have become a popular choice for machine learning tasks with graph structured data [Hamilton et al., 2017, Kipf and Welling, 2016, Veličković et al., 2017]. Computation in GNNs is performed by each node sending and receiving messages along the edges of the graph, and aggregating messages from its neighbors to update its embedding vector. After a few rounds of message passing, the computed node embeddings from all the nodes are aggregated to compute the final output [Gilmer et al., 2017]. This leads to a simple and elegant architecture for a variety of graph-related problems in practice.

While theoretically understanding the optimization, learning, and generalization properties of deep neural networks is a challenging task in itself, our current understanding of GNNs lags significantly behind their more popular counterparts such as fully connected feed-forward networks and convolutional neural networks (CNNs). In the context of GNNs recent theoretical works have focused on questions such as the representation power of various GNN architectures [Xu et al., 2019b,a, Loukas, 2020a,b] and establishing generalization bounds [Garg et al., 2020]. However, the optimization properties of GNN architectures are less explored. These involve identifying conditions under which algorithms like gradient descent can be shown to be provably effective in optimizing GNNs.

Such analysis has been carried out in recent years for fully connected feed-forward networks and CNNs, mainly restricted to depth-2 non-linear networks or deep linear networks [Ge et al., 2017, Andoni et al., 2014, Bartlett et al., 2017, Arora et al., 2018a, Hardt and Ma, 2016]. However, the message passing nature of GNNs makes it harder to theoretically analyze the optimization problem in GNNs, compared to other neural network architectures.

We provide the first convergence analysis of gradient descent for GNNs, and provably show that gradient descent can recover near optimal solutions. We first consider the case of a one-round (or one hidden-layer) GNN with ReLU activations. This scenario already captures much of the complexity of the problem, in particular allowing for multiple local optima. Following standard assumptions made in the case of analyzing feed-forward networks, we show that when the input data distribution is the standard Gaussian distribution, and the training objective is to minimize least square regression loss, gradient descent converges to the population error of ϵ in $O(\frac{1}{\epsilon^2} \log(\frac{1}{\epsilon}))$ iterations. In particular, our analysis applies to the standard practice of initializing the network with weights from a Gaussian distribution, and we do not perform any complicated pre-processing step to initialize the network.

We then extend our result to the multi-round (or multi-layer) case, but for the case of linear activations (analogous to the setting for deep linear networks [Arora et al., 2018a, Bartlett et al., 2018]). In this case, we show a stronger convergence result; gradient descent converges to an error ϵ in $O(\log(1/\epsilon))$ iterations. All our results are for the realizable setting, i.e, we assume that the data is indeed generated by an (unknown) underlying GNN, as is typical in analysis of gradient descent for other network architectures (see the discussion of related work in Section A). The rest of the paper is organized as follows. In Section 2 we set up notation and discuss the model of GNNs that we consider. In Section 3 we present our first main convergence theorem for the case of one layer/round GNNs with ReLU activations. This is followed by our second main result on analyzing deep linear GNNs in Section 4. We conclude with discussion and open problems in Section 6. We discuss the most relevant related works in individual sections (Section 3 and Section 4), and provide a more comprehensive survey of related work in Section A.

2 Preliminaries

Graph Neural Networks operate via passing and aggregating messages among the nodes of a graph. Given an undirected unweighted graph $G = (V, E)$, let $n = |V|$ be the number of nodes in the graph and $m = |E|$ be the number of edges. Initially each node $i \in [n]$ is initialized with its own private embedding $x_i \in \mathbb{R}^r$, where r denotes the embedding size. Then the computation in a GNN proceeds in a number of rounds where at round ℓ , the new embedding x_i^ℓ for vertex i is obtained via a combination of *aggregate* and *combine* steps Gilmer et al. [2017] as shown below.

$$a_i^{(\ell)} = \text{AGGREGATE}(\{x_j^{\ell-1} : j \in N(i)\}) \quad (1)$$

$$x_i^{(\ell)} = \text{COMBINE}(x_i^{(\ell-1)}, a_i^{(\ell)}). \quad (2)$$

Here $N(i)$ is the neighborhood of vertex i . When comparing to standard architectures such as feed-forward networks, it is useful to think of a round of message passing as computation performed by one layer of a more standard architecture. Hence an ℓ -round GNN should be thought of as a network with ℓ hidden layers. Different choices of the aggregate and combine operations lead to different versions of GNNs such as graph convolutional networks (GCNs) [Kipf and Welling, 2016], GraphSAGE [Hamilton et al., 2017], and graph isomorphism networks (GINs) [Xu et al., 2019b] to name a few. The aggregate operation is typically a simple pooling operation such as the sum or average, and the combine operation is implemented via a low depth neural network. Furthermore, in the most popular implementation, i.e, GCNs, the network parameters are shared across the different rounds. This setting will also be the focus of study in our work. In particular we consider two problem settings. In the first case (Section 3) we assume that the number of rounds in the GNN is one (i.e. one hidden layer) and that the combine operation is a depth one network with ReLU activations. In the second setting (Section 4) we consider arbitrary round L GNNs but restrict the combine operation to be a linear network. In both the settings we will assume that the aggregate operation is a sum, and that the initial input embeddings for the nodes are drawn from the standard Gaussian distributions, i.e., $N(0, I_{r \times r})$. Finally, our analysis for the case of ReLU activations will rely crucially on the notion of dual activations and their properties that we recap below.

Definition 1 ([Daniely et al., 2016]). *The dual activation of σ is the function $\hat{\sigma} : [-1, 1] \mapsto \mathbb{R}$ defined as*

$$\hat{\sigma}(\rho) = \mathbb{E}[\sigma(X)\sigma(Y)], \quad (3)$$

where X and Y are jointly Gaussian random variables with mean, zero variance one, and covariance ρ .

The work of Daniely et al. [2016] showed that dual activations satisfy many nice properties such as continuity in $[-1, 1]$ and the fact that they are convex in $[0, 1]$. For a more extensive list of the properties of the dual activations please refer to Lemma 11 in Daniely et al. [2016].

3 One round GNNs with ReLU activations

In this section we present our main result on convergence of gradient descent for learning an unknown one round GNN with ReLU activations. We consider a graph $G = (V, E)$ with n nodes and maximum degree d . Furthermore, we assume that there is an unknown GNN generating outcomes as

$$y = \sum_{i=1}^n \sigma(W^* \bar{x}_i), \quad (4)$$

where $\bar{x}_i = \sum_{j \in N(i)} x_j$ and each x_j is drawn i.i.d. from $N(0, I_{r \times r})$, and $\sigma(t) = \sqrt{2} \max(t, 0)$. In other words, the aggregate operation is a sum, and the combine operation is a depth-1 neural network with ReLU activations that produces an h -dimensional embedding for each node, where h is the number of hidden units (inner dimensionality of W^*). Finally, the embeddings of all the nodes are summed up to produce an h -dimensional output for the graph G .

The $\sqrt{2}$ factor multiplication to the ReLU function is for technical convenience and does not affect our results. The ground truth matrix W^* is an $h \times r$ matrix where we denote w_1^*, \dots, w_h^* to be the rows of W^* . Without loss of generality we will assume that $\|w_j^*\| = 1$, and again our results extend easily to the case when $\|w_j\|$ is bounded for all j (see Appendix B). We train another one round GNN with unknown parameter matrix W to minimize the following loss via gradient descent.

$$L(W) = \frac{1}{2} \mathbb{E}[\|\hat{y} - y\|^2] \quad (5)$$

$$= \frac{1}{2} \mathbb{E}[\|\sum_{i=1}^n \sigma(W \bar{x}_i) - \sum_{i=1}^n \sigma(W^* \bar{x}_i)\|^2] \quad (6)$$

$$= \sum_{j=1}^h L_j(w_j) := \sum_{j=1}^h \frac{1}{2} \mathbb{E}[(\sum_{i=1}^n \sigma(w_j^\top \bar{x}_i) - \sum_{i=1}^n \sigma(w_j^{*\top} \bar{x}_i))^2]. \quad (7)$$

We will analyze the following gradient descent updates.

$$W_{t+1} = W_t - \eta \nabla L(W_t). \quad (8)$$

Here $\nabla L(W)$ denotes the gradient of the population loss and as is common in practice, we assume that the entries of the matrix W are initialized i.i.d. from a Gaussian distribution. We are now ready to state our main theorem below.

Theorem 1. *Let W^* be the unknown parameter for a one-round GNN in Eq. (4) with $\|w_j^*\| = 1$ for $j \in [h]$, and let $L(W)$ denote the population loss at W as defined in Eq. (5). If the degree d of the graph is $o(\sqrt{n})$, then for any $\epsilon \in (0, 1)$, if $W_0 \sim N(0, I)$, with probability at least $1 - h \cdot e^{-c'r}$ we have that $L(W_T) \leq \epsilon^2$ provided $T \geq c \cdot \frac{n^4 h^2 (d+1)}{\epsilon^2} \log(\frac{nh}{\epsilon})$ for absolute constants $c, c' > 0$.*

Several remarks are in order regarding the above theorem. Note that the convergence rate is $O(\frac{1}{\epsilon^2} \log(\frac{1}{\epsilon}))$ as a function of the desired error, and is polynomial in the other problem parameters such as the number of nodes, the maximum degree of the graph, and the dimensionality of the W^* matrix. In the context of feedforward neural networks, gradient descent based algorithms have been analyzed primarily in the Neural Tangent Kernel (NTK) regime where an unknown target is learned by performing gradient descent on a highly over-parameterized network [Jacot et al., 2018,

[Du et al., 2018, Daniely, 2017, Allen-Zhu et al., 2019]. The advantage of such an analysis is that it can be carried out for large depth neural networks (at least for smooth activations). However, due to the massive overparameterization the gradient descent updates move extremely slowly and the dynamics correspond to performing kernel regression in a very high dimensional space [Lee et al., 2019b, Arora et al., 2019b]. This is quite far from the behavior of gradient descent in realistic settings. We, on the other hand are interested in the setting where there is an unknown network generating the data and one would like to learn it via gradient descent with only mild over-parameterization. In particular, even the simple problem setting we consider in this section captures important aspects of the complexity of the general problem, in particular allowing for multiple local minima.

There have been recent efforts in going beyond the NTK regime for the case of feed-forward networks. Similar to our setting, these works assume that the input is generated from the Gaussian distribution. Furthermore, these works still either need some over-parameterization to argue convergence [Li et al., 2020], or need to add appropriate regularizers to the squared loss objective [Ge et al., 2017].

In contrast, an interesting aspect of our result is that we do not need any additional over-parameterization and learn the unknown GNN using another network of the same size. The recent work of Zhang et al. [2020] studied a model for GNNs similar to our setting in Eq. (4) and designed a learning algorithm that first initializes the network weights based on a tensor decomposition subroutine, followed by an accelerated gradient descent procedure. We on the other hand do not perform any special initialization based on techniques such as spectral methods and analyze gradient descent updates starting from Gaussian initialization. The model of GNN we consider in Eq. (4) also bears similarity to the model of overlapping patches considered in the work of Du et al. [2017] in the context of convolutional neural networks. The authors analyze the convergence of gradient descent under certain strong assumptions on the correlation among the different patch distributions. We on the other hand make no such assumptions. Furthermore, even in the case of feed-forward networks, current analysis in the non-NTK regime do not extend beyond one hidden layer. Hence our work brings the understanding of gradient descent for GNNs on par with that of feed-forward networks.

We next present the main ideas and key lemmas behind the proof of Theorem 1.

3.1 Expressions for loss and gradients

In this section we first compute simplified expressions for the loss and the gradients that will be useful in subsequent analysis. We first begin by writing an equivalent expression for the population loss. We have

$$L_j(w_j) = \frac{1}{2} \mathbb{E} \left[\left(\sum_{i=1}^n \sigma(w_j^\top \bar{x}_i) \right)^2 \right] + \frac{1}{2} \mathbb{E} \left[\left(\sum_{i=1}^n \sigma(w_j^{*\top} \bar{x}_i) \right)^2 \right] - \mathbb{E} \left[\sum_{i,j=1}^n \sigma(w_j^\top \bar{x}_i) \sigma(w_j^{*\top} \bar{x}_j) \right].$$

Notice that each w_j evolves independently and hence we can simply focus on the convergence of $L_j(w_j)$ to $L_j(w_j^*)$. To simplify the above expression we will make use of the dual activation function of $\sigma(\cdot)$ from the work of Daniely et al. [2016] re-stated in Definition 1. For simplicity in this section we assume that the degree of each node is exactly d . The case when degree is at most d follows with minimal changes and is presented in Appendix B. Notice that each \bar{x}_i is a sum of $d+1$ messages, one involving the node i and the other d messages involving neighbors of i . Hence, \bar{x}_i is a random variable distributed as $N(0, (d+1)I)$. Furthermore, for any $i \neq j$ define $d_{i,j}$ to be the number of common messages between \bar{x}_i and \bar{x}_j . If $i = j$, then we define $d_{i,j} = d_{i,i} = d+1$. We next compute expressions for each of the three terms above.

$$\begin{aligned} \mathbb{E} \left[\left(\sum_{i=1}^n \sigma(w_j^\top \bar{x}_i) \right)^2 \right] &= \mathbb{E} \left[\sum_{i,j=1}^n \sigma(w_j^\top \bar{x}_i) \sigma(w_j^\top \bar{x}_j) \right] \\ &= \mathbb{E} \left[(d+1) \|w_j\|^2 \sum_{i,j=1}^n \sigma \left(\frac{w_j^\top \bar{x}_i}{\sqrt{(d+1)} \|w_j\|} \right) \sigma \left(\frac{w_j^\top \bar{x}_j}{\sqrt{(d+1)} \|w_j\|} \right) \right] \\ &= (d+1) \|w_j\|^2 \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \right) \quad (\text{definition of the dual activation in Eq. (3)}). \end{aligned} \tag{9}$$

Similarly we get that,

$$\mathbb{E} \left[\left(\sum_{i=1}^n \sigma(w_j^{*\top} \bar{x}_i) \right)^2 \right] = (d+1) \|w_j^*\|^2 \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \right). \quad (10)$$

Finally we simplify the last term.

$$\begin{aligned} \mathbb{E} \left[\sum_{i,j=1}^n \sigma(w_j^\top \bar{x}_i) \sigma(w_j^{*\top} \bar{x}_j) \right] &= \mathbb{E} \left[(d+1) \|w_j\| \|w_j^*\| \sum_{i,j=1}^n \sigma \left(\frac{w_j^\top \bar{x}_i}{\sqrt{(d+1)} \|w_j\|} \right) \sigma \left(\frac{w_j^{*\top} \bar{x}_j}{\sqrt{(d+1)} \|w_j^*\|} \right) \right] \\ &= (d+1) \|w_j\| \|w_j^*\| \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \frac{w_j^\top w_j^*}{\|w_j\| \|w_j^*\|} \right). \end{aligned} \quad (11)$$

Combining the (9), (10), and (11), we have

$$\begin{aligned} L(w_j) &= \frac{1}{2} (d+1) \|w_j\|^2 \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \right) + \frac{1}{2} (d+1) \|w_j^*\|^2 \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \right) \\ &\quad - (d+1) \|w_j\| \|w_j^*\| \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \frac{w_j^\top w_j^*}{\|w_j\| \|w_j^*\|} \right). \end{aligned} \quad (12)$$

It is easy to see that if $w_{0,j}$ is the initial value of w_j then each subsequent iteration will be a linear combination of $w_{0,j}$ and w_j^* . Hence we can assume that $w_j = \alpha w_j^* + \beta w_j^\perp$, where w_j^\perp is a fixed unit vector (depending on the initialization) orthogonal to w_j^* . Then re-writing the loss in terms of α, β and recalling that $\|w_j^*\| = 1$ we get the simplified expression:

$$\begin{aligned} L(\alpha, \beta) &= \frac{1}{2} (d+1) (\alpha^2 + \beta^2) \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \right) + \frac{1}{2} (d+1) \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \right) \\ &\quad - (d+1) \sqrt{\alpha^2 + \beta^2} \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}} \right). \end{aligned} \quad (13)$$

In the rest of the section we will analyze the evolution of the updates of α and β . Furthermore we will use α_t, β_t to denote the parameters α, β associated with the iterate $w_{j,t}$ at time t . We first compute the gradient of the objective w.r.t. w or equivalently w.r.t. α, β .

$$\begin{aligned} \frac{\partial L(\alpha, \beta)}{\partial \alpha} &= \alpha(d+1) \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \right) - \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}} (d+1) \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}} \right) \\ &\quad - \frac{\beta^2}{\alpha^2 + \beta^2} (d+1) \sum_{i,j=1}^n \frac{d_{i,j}}{d+1} \hat{\sigma}' \left(\frac{d_{i,j}}{d+1} \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}} \right) \\ &= \alpha(d+1) \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \right) - \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}} (d+1) \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}} \right) \\ &\quad - \frac{\beta^2}{\alpha^2 + \beta^2} (d+1) \sum_{i,j=1}^n \frac{d_{i,j}}{d+1} \hat{\sigma}_{\text{step}} \left(\frac{d_{i,j}}{d+1} \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}} \right). \end{aligned} \quad (14)$$

Here in the last equality we have used the fact that $\hat{\sigma}' = \hat{\sigma}'$ and that $\sigma'(x) = \sqrt{2} \mathbf{1}(x \geq 0) = \sigma_{\text{step}}(x)$, where $\sigma_{\text{step}}(x)$ is the step function. See [Daniely et al., 2016] for a proof. Similarly we get that

$$\begin{aligned} \frac{\partial L(\alpha, \beta)}{\partial \beta} &= \beta(d+1) \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \right) - \frac{\beta}{\sqrt{\alpha^2 + \beta^2}} (d+1) \sum_{i,j=1}^n \hat{\sigma} \left(\frac{d_{i,j}}{d+1} \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}} \right) \\ &\quad + \frac{\alpha\beta}{\alpha^2 + \beta^2} (d+1) \sum_{i,j=1}^n \frac{d_{i,j}}{d+1} \hat{\sigma}_{\text{step}} \left(\frac{d_{i,j}}{d+1} \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}} \right). \end{aligned} \quad (15)$$

3.2 Bounding the iterates.

In this section we show that if the initial weight $w_{0,j}$ is drawn from $N(0, \sigma^2 I)$ then with high probability, the iterates remain bounded for all subsequent time steps. We first analyze how the length of w , i.e., $\|w_t\|^2$ behaves over a period of time. In the rest of the proof we will often use the shorthand A, B_t, C_t to denote the following key quantities that depend on the structure of the graph.

$$\begin{aligned} A &= \sum_{i,j=1}^n \hat{\sigma}\left(\frac{d_{i,j}}{d+1}\right) \\ B_t &= \sum_{i,j=1}^n \hat{\sigma}\left(\frac{d_{i,j}}{d+1} \frac{\alpha_t}{\ell_t}\right) \\ C_t &= \sum_{i,j=1}^n \frac{d_{i,j}}{d+1} \hat{\sigma}_{\text{step}}\left(\frac{d_{i,j}}{d+1} \frac{\alpha_t}{\ell_t}\right). \end{aligned}$$

Here ℓ_t denotes the length of the iterate w_t at time t , i.e., $\ell_t = \|w_t\| = \sqrt{\alpha_t^2 + \beta_t^2}$. The above quantities satisfy useful inequalities that will be used throughout the analysis. For example, we note that A, B_t are always non-negative when $d = o(\sqrt{n})$ and $\alpha_t \geq 0$. We also have $\frac{n^2}{4\pi} \leq |B_t| \leq A$ and $|C_t| = o(A)$. See Appendix B for the proof.

Lemma 1. *If $w_{0,j} \sim N(0, \sigma^2 I)$ and $\eta \in [\frac{1}{16\pi(d+1)A}, \frac{1}{6\pi(d+1)A}]$ then with probability at least $1 - e^{-O(r)}$, it holds that for all $t \geq 0$, $\ell_t \leq 2\sigma\sqrt{r} + 4\pi + 1$, and $\ell_t \geq \frac{1}{36\pi^2}$ for all $t \geq 1$.*

3.3 Establishing smoothness and analyzing initial updates

We first show that if the condition in Lemma 1 holds, then throughout the trajectory of the iterates, the loss function is smooth. This is formalized in the lemma below.

Lemma 2. *If the degree d of the graph is $o(\sqrt{n})$ and the conditions in Lemma 1 hold, then for all $t \geq 1$ and any $\gamma \in [0, 1]$ such that $(\alpha, \beta) = (1 - \gamma)(\alpha_t, \beta_t) + \gamma(\alpha_{t+1}, \beta_{t+1})$, we have that*

$$\lambda_{\max}(\nabla^2 L(\alpha, \beta)) \leq 4(d+1)A \left(1 + \sqrt{2\sigma\sqrt{r} + 4\pi + 1} + o(1)\right).$$

Here λ_{\max} is the maximum eigenvalue of the population Hessian denoted by $\nabla^2 L(\alpha, \beta)$. Our overall proof strategy is to show that the Polyak-Łojasiewicz (PL) condition [Polyak, 1963] holds, namely that the squared norm of the gradient lower bounds the loss value at any iterate. This will let us easily analyze convergence of gradient descent via standard arguments. However a challenge is that since we are starting from random initialization, the PL condition does not hold true at the beginning. As a result we separately analyze an initial phase of the algorithm and show that gradient descent very quickly escapes a region where the PL condition does not hold. This is formalized below.

Lemma 3. *If $w_{0,j} \sim N(0, \sigma^2 I)$ and $\eta \in [\frac{1}{16\pi(d+1)A}, \frac{1}{6\pi(d+1)A}]$ then with at least $1 - 1/h^2$, it holds that for all $t \geq c \log(\sigma \log h)$, $\alpha_t \geq -\frac{1}{100}$ and $\ell_t \geq 1 - o(1)$, where $c > 0$ is an absolute constant.*

3.4 Establishing the PL-condition

Finally, using the above lemma we show that under the conditions of Lemma 1 the iterates $w_t = (\alpha_t, \beta_t)$ satisfy the Polyak-Łojasiewicz inequality [Polyak, 1963].

Lemma 4. *If the degree d of the graph is $o(\sqrt{n})$ and the conditions in Lemma 1 hold then there is an absolute constant $c > 0$, such that for all $t \geq c \log(\sigma \log h)$ and $\epsilon \in (0, 1)$, either $|\beta_t| \leq \frac{\epsilon}{4hn}$ and $\|\ell_t - 1\| \leq \frac{\epsilon}{4hn}$ or we have that*

$$\|\nabla L(\alpha_t, \beta_t)\|^2 \geq \mu^* L(\alpha_t, \beta_t),$$

where $\mu^* \geq \frac{\epsilon^2}{380(d+1)h^2\pi n^2}$.

3.5 Putting everything together

Finally, we combine the analysis from the previous sections to prove the main theorem, i.e., Theorem 1.

Proof of Theorem 1. We will analyze an arbitrary $j \in [h]$ and the evolution of the corresponding w_j vector. By setting $\sigma = 1$ we have from Lemma 2 that the smoothness parameter of the loss function is

$$L \leq 4(d+1)A\left(1 + \sqrt{4\pi+3} + o(1)\right).$$

Hence we get that for any $t \geq 0$,

$$\begin{aligned} L_j(w_{j,t+1}) &\leq L_j(w_{j,t}) + \nabla L_j(w_{j,t})(w_{j,t+1} - w_{j,t}) + \frac{L}{2}\|w_{j,t+1} - w_{j,t}\|^2 \\ &\leq L_j(w_{j,t}) - \eta\|\nabla L_j(w_{j,t})\|^2 + \frac{\eta^2 L}{2}\|\nabla L_j(w_{j,t})\|^2 \\ &= L_j(w_{j,t}) - \eta\|\nabla L_j(w_{j,t})\|^2\left(1 - \frac{\eta L}{2}\right). \end{aligned} \quad (16)$$

From the range of η in Lemma 1 we get that $\eta L \leq 1$. Furthermore, using Lemma 4 we can write

$$\begin{aligned} L_j(w_{j,t+1}) &\leq L_j(w_{j,t}) + \nabla L_j(w_{j,t})(w_{j,t+1} - w_{j,t}) + \frac{L}{2}\|w_{j,t+1} - w_{j,t}\|^2 \\ &\leq L_j(w_{j,t})(1 - \eta\mu^*) \\ &\leq L_j(w_{j,0})(1 - \eta\mu^*)^t. \end{aligned} \quad (17)$$

Hence after $T \geq c \cdot \frac{n^4 h^2 (d+1)}{\epsilon^2} \log(\frac{nh}{\epsilon})$ time steps we will either have $L_j(w_{j,t}) \leq \epsilon^2/h$, or that $|\beta_t| \leq \frac{\epsilon}{4hn}$ and $\|\ell_t - 1\| \leq \frac{\epsilon}{4hn}$. The latter implies that

$$\|w_{j,t} - w_j^*\|^2 \leq \frac{\epsilon^2}{2hn^2}.$$

Furthermore, it is easy to see that

$$L_j(w_{j,t}) \leq n(n-1)\|w_{j,t} - w_j^*\|^2.$$

Hence if the latter happens then again, $L_j(w_{j,t}) \leq \epsilon^2/h$ thereby implying that $L(W_T) \leq \epsilon^2$. \square

4 Deep linear GNNs

In the previous section we analyzed one round GNNs with ReLU activations. Analyzing the dynamics of gradient descent beyond one layer (in the non NTK regime) with non-linear activations is a challenging problem, even for standard network architectures. Hence, in this section we consider deep linear GNNs in order to analyze the dynamics of gradient descent in deeper message passing architectures. Linearity has been a commonly studied setting in recent years for analyzing deep feed-forward networks [Hardt and Ma, 2016, Bartlett et al., 2018, Arora et al., 2018a].

We assume the existence of an unknown L -round GNN characterized by positive-definite symmetric matrices $W_1^*, W_2^* \in \mathbb{R}^{r \times r}$. Given node inputs $x_1, \dots, x_n \sim N(0, I_{r \times r})$ the embedding of node i at round $\ell \geq 1$ is defined as

$$a_i^{(\ell)} = \text{AGGREGATE}(\{x_j^{(\ell-1)} : j \in N(i)\}) = W_2^* \sum_{j \in N(i)} x_j^{(\ell-1)}. \quad (18)$$

$$x_i^{(\ell)} = \text{COMBINE}(x_i^{(\ell-1)}, a_i^{(\ell)}) = W_1^* x_i^{(\ell-1)} + a_i^{(\ell)}. \quad (19)$$

Notice that this is a more general setting than the case of ReLU network studied in the previous section where the aggregate was a simple summation operation. In this section we can analyze the case where both the aggregate and combine operations have their own set of parameters W_1^* and W_2^* . Finally, the output of the network on input $\vec{x} = \{x_1, \dots, x_n\}$ is defined as

$$y = \sum_{i=1}^n x_i^{(L)}.$$

As in the case of ReLU activations from the previous section, we will analyze the non-overparameterized case where another architecture of the same size is used to learn the unknown ground truth network. Given another L -round GNN defined by parameters W_1, W_2 we define the population loss as

$$L(W_1, W_2) = \frac{1}{2} \mathbb{E} [\|\hat{y} - y\|^2] \quad (20)$$

$$= \frac{1}{2} \mathbb{E} [\|\sum_{i=1}^n \hat{x}_i^{(L)} - \sum_{i=1}^n x_i^{(L)}\|^2]. \quad (21)$$

We will analyze the gradient descent updates defined as

$$W_{1,t+1} = W_{1,t} - \eta \nabla L(W_{1,t}, W_{2,t}) \quad (22)$$

$$W_{2,t+1} = W_{2,t} - \eta \nabla L(W_{1,t}, W_{2,t}). \quad (23)$$

For ease of interpretation, we state our main theorem below for the case when each node in the graph has degree exactly d . The general theorem is stated in Appendix C and achieves similar rate of convergence. The theorem below shows that if the network parameters are initialized at identity, i.e., $W_{1,0}, W_{2,0} = I$, then gradient descent minimizes the population loss at a polylogarithmic rate.

Theorem 2. *Let the initialization of $W_{1,0}, W_{2,0}$ be identity. Then, $L(W_{1,T}, W_{2,T}) \leq \epsilon^2$ if,*

$$T \geq \max_{i \in [n]} \frac{1}{\eta_0 L^2 (d^2 + 1) \ell_i^{2L-2}} \log\left(\frac{r L u_i ((1 - \sigma_i^*)^2)}{\epsilon}\right) \quad (24)$$

$$\eta := \eta_0 \leq \min_i \frac{\ell_i^{2L-1} \min(1, \frac{1}{(\sigma_i^*)^{2L-2}})}{2n L^2 (\delta_i)^2 (d^2 + 1) u_i^{2L-3} (1 + u_i)^{L-1}}. \quad (25)$$

Here σ_i^* is the i th smallest singular value of $W_1^* + dW_2^*$, and u_i, ℓ_i and δ_i depend only on σ_i^* .

We make a few comments about the theorem above. Notice that unlike the case of ReLU networks we do not assume any bound on the maximum degree of the graph and the theorem applies generally. Furthermore, the convergence rate is logarithmic in $\frac{1}{\epsilon}$ and exponential in the depth L of the network. The dependence on the number of nodes in the graph is linear and the dependence on the maximum degree d is polynomial and appears via the bound on the maximum singular value of the matrix $W_1^* + dW_2^*$. Our analysis for the case of equal degrees proceeds by showing that it is enough to track the evolution of the singular values $\sigma_{i,t}$ of the matrix $M_t = W_{1,t} + dW_{2,t}$ at each time step. We then show that if η is sufficiently small then the singular values $\sigma_{i,t}$ will converge to the singular values σ_i^* of the true unknown matrix $W_1^* + dW_2^*$. Handling the case of unequal degrees is along similar lines but is technically more challenging as we need to separately track the evolution of the singular values of $W_{1,t}$ and $W_{2,t}$ separately. See Appendix C for details.

5 Experiments

In this section we verify our theoretical results via simulations. Our first goal is to understand whether the $\sim \frac{1}{\epsilon^2}$ rate of convergence obtained for the case of one round ReLU GNNs is tight, or whether in practice we can obtain $\text{polylog}(\frac{1}{\epsilon})$ rates even for the ReLU setting. Secondly, we understand via experiments the regime where the degree of the graph d exceeds \sqrt{n} , and hence our theoretical results for the case of ReLU GNNs do not hold. In order to do this we generate random d -regular graphs over $n = 100$ nodes. For the case of one round GNNs with ReLU activations we set the embedding size $r = 10$, and $h = 10$ (number of hidden units in the ReLU GNN). We use the same value of r for the case of deep linear GNNs, where r equals the input dimensionality and also the dimensionality of the matrices W_1^* , and W_2^* .

We generate an unknown ground truth network in the case of ReLU GNNs by choosing each column of W^* to be a random unit length vector. For the case of linear networks we generate positive definite matrices W_1^*, W_2^* by picking random Gaussian entries, and then adding a small multiplicative factor of 0.001 times the identity matrix. As dictated by our theory, we initialize network weights from a standard Gaussian for the case of ReLU networks, and to the identity matrix for the linear case.

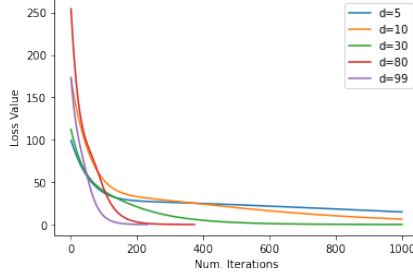


Figure 1: Loss vs. number of iterations of gradient descent for one round GNN with ReLU activations.

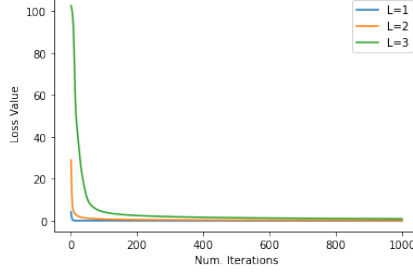


Figure 2: Loss vs. number of iterations of gradient descent for deep linear GNNs.

We simulate population gradient descent and implement our networks using the JAX programming language [Bradbury et al., 2018]. Our experiments are run using one GPU. See Section D for further details.

For the case of ReLU networks we run gradient descent on instances with varying values of the degree d . In particular we present results for $d \in \{5, 10, 30, 80, 99\}$. For the case of linear networks we vary the depth L in $\{1, 2, 3\}$. Finally, in each case we plot the loss value (ℓ_2 loss) vs. the number of iterations. The results are shown in Figure 1 and Figure 2. We make a few remarks about the experiments. Note that the convergence rate for the case of linear networks is indeed much better than the case of ReLU GNNs and hence we, in general, do not expect significantly better rates of convergence than the $\sim \frac{1}{\epsilon^2}$ bound that we prove. An interesting observation is that as d increases and is close to the regime of $\Theta(n)$, the rates of convergence of gradient descent for the case of a one round GNN with ReLU activations, become much better. This is along expected lines, since if $d = n - 1$, i.e., the graph is a complete graph, then the network in Eq. (4) boils down to h independent single ReLU units with standard Gaussian inputs. It is known for this extreme setting that gradient descent has a logarithmic convergence in $\frac{1}{\epsilon}$ [Soltanolkotabi, 2017]. Hence, extending our existing theoretical analysis beyond the \sqrt{n} degree setting is an interesting open question.

6 Conclusions

The main question left open by our work is to analyze the gradient descent updates for multi-round GNNs with ReLU activations. Given the difficulty of analyzing gradient descent for multi-layer networks beyond the NTK setting, the above seems like a challenging problem. It would also be interesting to explore whether the convergence rate for the case of one round GNNs can be improved via over-parameterization. Finally, it would also be of interest to extend our results to more general distributions going beyond the Gaussian setting. Another fascinating question is to analyze gradient descent in the *agnostic* setup. If the best one round GNN with ReLU activations (or a deep linear GNN) achieves an error of OPT on the data distribution, can gradient descent recover a network of error $f(\text{OPT})$? This style of analysis has recently been carried out in the case of a single ReLU unit [Frei et al., 2020].

Acknowledgments and Disclosure of Funding

Funding in direct support of this work: Google Research.

References

- Z. Allen-Zhu and Y. Li. Backward feature correction: How deep learning performs deep learning. *arXiv preprint arXiv:2001.04413*, 2020.
- Z. Allen-Zhu, Y. Li, and Z. Song. A convergence theory for deep learning via over-parameterization. *arXiv preprint arXiv:1811.03962*, 2018.
- Z. Allen-Zhu, Y. Li, and Z. Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252. PMLR, 2019.
- A. Andoni, R. Panigrahy, G. Valiant, and L. Zhang. Learning polynomials with neural networks. In *International conference on machine learning*, pages 1908–1916. PMLR, 2014.
- S. Arora, N. Cohen, N. Golowich, and W. Hu. A convergence analysis of gradient descent for deep linear neural networks. *arXiv preprint arXiv:1810.02281*, 2018a.
- S. Arora, N. Cohen, and E. Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In *ICML*, pages 244–253, 2018b.
- S. Arora, N. Cohen, W. Hu, and Y. Luo. Implicit regularization in deep matrix factorization. In *Advances in Neural Information Processing Systems*, pages 7411–7422, 2019a.
- S. Arora, S. Du, W. Hu, Z. Li, and R. Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pages 322–332. PMLR, 2019b.
- S. Arora, S. S. Du, W. Hu, Z. Li, R. Salakhutdinov, and R. Wang. On exact computation with an infinitely wide neural net. In *NeurIPS*, 2019c.
- A. Bakshi, R. Jayaram, and D. P. Woodruff. Learning two layer rectified neural networks in polynomial time. In *Conference on Learning Theory*, pages 195–268. PMLR, 2019.
- P. Bartlett, D. Helmbold, and P. Long. Gradient descent with identity initialization efficiently learns positive definite linear transformations by deep residual networks. In *International conference on machine learning*, pages 521–530. PMLR, 2018.
- P. L. Bartlett, D. J. Foster, and M. J. Telgarsky. Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems*, pages 6240–6249, 2017.
- J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- S. Chen, A. R. Klivans, and R. Meka. Learning deep relu networks is fixed-parameter tractable. *arXiv preprint arXiv:2009.13512*, 2020.
- L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *arXiv preprint arXiv:1805.09545*, 2018.
- A. Daniely. Sgd learns the conjugate kernel class of the network. *arXiv preprint arXiv:1702.08503*, 2017.
- A. Daniely, R. Frostig, and Y. Singer. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity. *arXiv preprint arXiv:1602.05897*, 2016.
- I. Diakonikolas, D. M. Kane, V. Kontonis, and N. Zarifis. Algorithms and sq lower bounds for pac learning one-hidden-layer relu networks. In *Conference on Learning Theory*, pages 1514–1539. PMLR, 2020.

- S. S. Du, J. D. Lee, and Y. Tian. When is a convolutional filter easy to learn? *arXiv preprint arXiv:1709.06129*, 2017.
- S. S. Du, X. Zhai, B. Póczos, and A. Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- S. Frei, Y. Cao, and Q. Gu. Agnostic learning of a single neuron with gradient descent. *arXiv preprint arXiv:2005.14426*, 2020.
- V. K. Garg, S. Jegelka, and T. Jaakkola. Generalization and representational limits of graph neural networks. *ICML*, 2020.
- R. Ge, J. D. Lee, and T. Ma. Learning one-hidden-layer neural networks with landscape design. *arXiv preprint arXiv:1711.00501*, 2017.
- R. Ge, R. Kuditipudi, Z. Li, and X. Wang. Learning two-layer neural networks with symmetric inputs. *arXiv preprint arXiv:1810.06793*, 2018.
- J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. *arXiv preprint arXiv:1704.01212*, 2017.
- S. Goel and A. R. Klivans. Learning neural networks with two nonlinear layers in polynomial time. In *Conference on Learning Theory*, pages 1470–1499. PMLR, 2019.
- W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs. In *Advances in neural information processing systems*, pages 1024–1034, 2017.
- M. Hardt and T. Ma. Identity matters in deep learning. *arXiv preprint arXiv:1611.04231*, 2016.
- A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *arXiv preprint arXiv:1806.07572*, 2018.
- M. Janzamin, H. Sedghi, and A. Anandkumar. Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods. *arXiv preprint arXiv:1506.08473*, 2015.
- T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- J. Lee, L. Xiao, S. Schoenholz, Y. Bahri, R. Novak, J. Sohl-Dickstein, and J. Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. In *NeurIPS*, pages 8570–8581. 2019a.
- J. Lee, L. Xiao, S. S. Schoenholz, Y. Bahri, R. Novak, J. Sohl-Dickstein, and J. Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *arXiv preprint arXiv:1902.06720*, 2019b.
- Y. Li and Y. Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *NeurIPS*, pages 8168–8177, 2018.
- Y. Li and Y. Yuan. Convergence analysis of two-layer neural networks with relu activation. In *Advances in neural information processing systems*, pages 597–607, 2017.
- Y. Li, T. Ma, and H. R. Zhang. Learning over-parametrized two-layer relu neural networks beyond ntk. *arXiv preprint arXiv:2007.04596*, 2020.
- A. Loukas. What graph neural networks cannot learn: depth vs width. *ICLR*, 2020a.
- A. Loukas. How hard is to distinguish graphs with graph neural networks? Technical report, 2020b.
- E. Malach and S. Shalev-Shwartz. Is deeper better only when shallow is good? In *Advances in Neural Information Processing Systems*, pages 6426–6435, 2019.
- S. Mei, A. Montanari, and P.-M. 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.

- S. Oymak and M. Soltanolkotabi. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *CoRR*, abs/1902.04674, 2019.
- B. T. Polyak. Gradient methods for minimizing functionals. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki*, 3(4):643–653, 1963.
- G. Rotskoff and E. Vanden-Eijnden. Parameters as interacting particles: long time convergence and asymptotic error scaling of neural networks. In *Advances in neural information processing systems*, pages 7146–7155, 2018.
- J. Sirignano and K. Spiliopoulos. Mean field analysis of neural networks. *arXiv preprint arXiv:1805.01053*, 2018.
- M. Soltanolkotabi. Learning relus via gradient descent. *arXiv preprint arXiv:1705.04591*, 2017.
- M. Soltanolkotabi, A. Javanmard, and J. D. Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory*, 65(2): 742–769, 2018.
- L. Su and P. Yang. On learning over-parameterized neural networks: A functional approximation perspective. In *NeurIPS*, pages 2637–2646, 2019.
- P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- S. Vempala and J. Wilmes. Gradient descent for one-hidden-layer neural networks: Polynomial convergence and sq lower bounds. *arXiv preprint arXiv:1805.02677*, 2018.
- C. Wei, J. D. Lee, Q. Liu, and T. Ma. Regularization matters: Generalization and optimization of neural nets vs their induced kernel. In *Advances in Neural Information Processing Systems*, pages 9709–9721, 2019.
- K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? *ICLR*, 2019a.
- K. Xu, J. Li, M. Zhang, S. S. Du, K.-i. Kawarabayashi, and S. Jegelka. What can neural networks reason about? *arXiv preprint arXiv:1905.13211*, 2019b.
- S. Zhang, M. Wang, S. Liu, P.-Y. Chen, and J. Xiong. Fast learning of graph neural networks with guaranteed generalizability: One-hidden-layer case. In *International Conference on Machine Learning*, pages 11268–11277. PMLR, 2020.
- D. Zou and Q. Gu. An improved analysis of training over-parameterized deep neural networks. In *NeurIPS*, 2019.

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? **[Yes]** See Section ??.
- Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#) See Section 3 and Section 4.
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section 5 and Section 6.
 - (c) Did you discuss any potential negative societal impacts of your work? [\[No\]](#) Our work is theoretical in nature and we do not anticipate any negative societal impact.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
- (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) See Section 3 and Section 4.
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) See the Appendix.
3. If you ran experiments...
- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) See supplementary material.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) See Section D.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[No\]](#) We ran gradient descent on population loss and the variance due to random initialization was negligible.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See Section 5.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#) See Section 5.
 - (b) Did you mention the license of the assets? [\[Yes\]](#) See Section D.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[N/A\]](#)
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[N/A\]](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)

A Related Work

There is a vast amount of literature on theory of deep learning. Some recent works [Ge et al., 2017, Bakshi et al., 2019, Ge et al., 2018, Janzamin et al., 2015, Vempala and Wilmes, 2018] provide polynomial time algorithms for learning depth-2 feedforward ReLU networks using techniques based on tensor decompositions and spectral methods. Other works focus on designing learning algorithms without assumptions on the linear independence of columns of W . These results incur an exponential dependence on either the input dimensionality or the number of parameters in the unknown network [Diakonikolas et al., 2020, Chen et al., 2020]. Polynomial time algorithms beyond depth 2 can be designed (under some assumptions) for other activations such as the sigmoid function [Goel and Klivans, 2019].

Properties of Gradient Descent The works of Bartlett et al. [2018] and Arora et al. [2018a] study the convergence of gradient descent on deep linear feed-forward networks. Recent extensions include studying the behavior of gradient descent on overparameterized deep linear networks for solving problems such as matrix factorization [Arora et al., 2019a, 2018b] and also exploring distributions that are hard for gradient descent [Malach and Shalev-Shwartz, 2019].

There is also a large body of work in analyzing the convergence of gradient descent and stochastic gradient descent (SGD) on over parameterized neural networks. The work of Andoni et al. [2014] concerns learning low degree polynomials via gradient descent on depth-2 networks with quadratic activations. There have also been works analyzing the convergence properties of gradient descent on depth-2 feedforward networks with identity mapping [Li and Yuan, 2017, Soltanolkotabi et al., 2018, Li and Liang, 2018] and certain assumptions on the data distribution.

Going beyond assumptions on the data distribution, there has been a recent surge in analyzing the convergence of gradient descent in the so called “NTK regime”, i.e., gradient descent on massively overparameterized networks [Daniely, 2017, Daniely et al., 2016, Allen-Zhu et al., 2018, Du et al., 2018, Arora et al., 2019b, Lee et al., 2019b, Chizat and Bach, 2018, Jacot et al., 2018, Oymak and Soltanolkotabi, 2019, Zou and Gu, 2019, Su and Yang, 2019].

The NTK or the neural tangent kernel as proposed in the work of Jacot et al. [2018] is a kernel defined by the gradient of the network parameters at random initialization. Recent works have shown that the NTK effectively captures the dynamic of gradient descent in the above mentioned overparameterized settings [Lee et al., 2019a, Arora et al., 2019c]. An alternate line of work focused on studying the mean field dynamics of SGD on infinite width neural networks [Mei et al., 2018, Chizat and Bach, 2018, Rotskoff and Vanden-Eijnden, 2018, Sirignano and Spiliopoulos, 2018]. Very recent works have also explored analyzing gradient descent beyond the NTK regime [Wei et al., 2019, Allen-Zhu and Li, 2020].

B Proof from Section 3

Properties of A, B_t, C_t . Recall the following key quantities that we will use throughout the proof.

$$\begin{aligned} A &= \sum_{i,j=1}^n \hat{\sigma}\left(\frac{d_{i,j}}{d+1}\right) \\ B_t &= \sum_{i,j=1}^n \hat{\sigma}\left(\frac{d_{i,j}}{d+1} \frac{\alpha_t}{\ell_t}\right) \\ C_t &= \sum_{i,j=1}^n \frac{d_{i,j}}{d+1} \hat{\sigma}_{\text{step}}\left(\frac{d_{i,j}}{d+1} \frac{\alpha_t}{\ell_t}\right). \end{aligned}$$

We first establish some useful relations among the above quantities.

Lemma 5. *If the degree of the graph $d = o(\sqrt{n})$, then we have that*

1. $A \geq \frac{n^2}{\pi}, \frac{n^2}{4\pi} \leq B_t \leq A$, and $|C_t| = o(A)$.
2. When $\alpha_t \geq -\frac{1}{100}$ and $\ell_t \geq 1 - o(1)$, then $C_t \geq \frac{A}{2n}$.

We prove the lemma above at the end of the section and first prove the main claims needed for the proof of Theorem 1.

Lemma 6 (Restatement of Lemma 1.). *If $w_{0,j} \sim N(0, \sigma^2 I)$ and $\eta \in [\frac{1}{16\pi(d+1)A}, \frac{1}{6\pi(d+1)A}]$ then with probability at least $1 - e^{-O(r)}$, it holds that for all $t \geq 0$, $\ell_t \leq 2\sigma\sqrt{r} + 4\pi + 1$, and $\ell_t \geq \frac{1}{36\pi^2}$ for all $t \geq 1$.*

Proof. From (14) and (15) we have

$$\begin{aligned} \ell_{t+1}^2 &:= \|w_{t+1}\|^2 \\ &= \alpha_{t+1}^2 + \beta_{t+1}^2 \\ &= \left(\alpha_t - \eta \frac{\partial L(\alpha, \beta)}{\alpha} \Big|_{\alpha_t, \beta_t}\right)^2 + \left(\beta_t - \eta \frac{\partial L(\alpha, \beta)}{\beta} \Big|_{\alpha_t, \beta_t}\right)^2 \\ &= \left(\alpha_t - \eta \alpha_t (d+1)A + \eta \frac{\alpha_t}{\ell_t} (d+1)B_t + \eta \frac{\beta_t^2}{\ell_t^2} (d+1)C_t\right)^2 \\ &\quad + \left(\beta_t - \eta \beta_t (d+1)A + \eta \frac{\beta_t}{\ell_t} (d+1)B_t - \eta \frac{\alpha_t \beta_t}{\ell_t^2} (d+1)C_t\right)^2 \\ &= \ell_t^2 + \eta^2 (d+1)^2 A^2 \ell_t^2 + \eta^2 (d+1)^2 B_t^2 + \eta^2 \frac{\beta_t^2}{\ell_t^2} (d+1)^2 C_t^2 \\ &\quad - 2\eta (d+1) A \ell_t^2 + 2\eta (d+1) \ell_t B_t - 2\eta^2 (d+1)^2 A B_t \ell_t. \end{aligned}$$

Hence setting $\eta \in [\frac{1}{16\pi(d+1)A}, \frac{1}{6\pi(d+1)A}]$ we get that

$$\ell_{t+1}^2 \leq \left(\ell_t \left(1 - \frac{1}{16\pi}\right)\right)^2 + \frac{1}{16}.$$

A simple calculation then shows that for all $t \geq 1$, $\ell_t \leq \ell_0 + 4\pi + 1$. Furthermore, we also have that if $B_t \geq 0$ then

$$\begin{aligned} \ell_{t+1}^2 &\geq \ell_t^2 (1 - \eta(d+1)A)^2 + 2\eta(d+1)B_t \ell_t (1 - \eta(d+1)A) + \eta^2 (d+1)^2 B_t^2 \\ &\geq \eta^2 (d+1)^2 B_t^2. \end{aligned}$$

Using the fact that $B_t \geq n^2/(4\pi)$ and the range of η we get that $\ell_t \geq \frac{1}{36\pi^2}$ for $t \geq 1$. Finally, notice that with probability at least $1 - e^{-O(r)}$, we will have $\ell_0 = O(\sigma\sqrt{r})$. \square

Lemma 7 (Restatement of Lemma 2). *If the degree d of the graph is $o(\sqrt{n})$ and the conditions in Lemma 1 hold then for all $t \geq 1$ and any $\gamma \in [0, 1]$ such that $(\alpha, \beta) = (1-\gamma)(\alpha_t, \beta_t) + \gamma(\alpha_{t+1}, \beta_{t+1})$, we have that*

$$\lambda_{\max}(\nabla^2 L(\alpha, \beta)) \leq 4(d+1)A \left(1 + \sqrt{2\sigma\sqrt{r} + 4\pi + 1} + o(1)\right).$$

Proof. We first establish upper and lower bounds on the length of the iterate, i.e., $\ell^2 = \alpha^2 + \beta^2$. By convexity of length and Lemma 1 we immediately get that $\alpha^2 + \beta^2 \leq 2\sigma\sqrt{r} + 4\pi + 1$. To obtain a lower bound on the length we have

$$\begin{aligned}\ell^2 &= ((1-\lambda)\alpha_t + \lambda\alpha_{t+1})^2 + ((1-\lambda)\beta_t + \lambda\beta_{t+1})^2 \\ &= (\alpha_t - \eta\lambda\frac{\partial L}{\partial \alpha})^2 + (\beta_t - \eta\lambda\frac{\partial L}{\partial \beta})^2\end{aligned}$$

Following the same calculation as in the proof of Lemma 1 we get that

$$\begin{aligned}\ell^2 &\geq \ell_t^2(1 - \lambda\eta(d+1)A)^2 + 2\lambda\eta(d+1)B_t\ell_t(1 - \lambda\eta(d+1)A) + \eta^2(d+1)^2B_t^2 \\ &\geq \ell_t^2(1 - \lambda\eta(d+1)A)^2 \\ &\geq \ell_t^2(1 - \frac{1}{6\pi})^2.\end{aligned}$$

Hence we get that $\ell \geq \frac{1}{36\pi^2}(1 - \frac{1}{6\pi})$. Next we upper bound $\lambda_{\max}(\nabla^2(\alpha, \beta))$. We will use the following simple upper bound on the maximum eigenvalue.

$$\lambda_{\max}(\nabla^2(\alpha, \beta)) \leq \left| \frac{\partial^2 L(\alpha, \beta)}{\partial \alpha^2} \right| + \left| \frac{\partial^2 L(\alpha, \beta)}{\partial \beta^2} \right| + \left| \frac{\partial^2 L(\alpha, \beta)}{\partial \alpha \partial \beta} \right| + \left| \frac{\partial^2 L(\alpha, \beta)}{\partial \beta \partial \alpha} \right|.$$

Taking the second derivatives we get

$$\begin{aligned}\frac{\partial^2 L(\alpha, \beta)}{\partial \alpha^2} &= (d+1)A - \frac{\alpha}{\ell}(d+1)\frac{\partial B}{\partial \alpha} - \frac{\beta^2}{\ell^{3/2}}(d+1)B - \frac{\beta^2}{\ell^2}(d+1)\frac{\partial C}{\partial \alpha} + \frac{2\alpha\beta^2}{\ell^4}(d+1)C \\ \frac{\partial^2 L(\alpha, \beta)}{\partial \beta^2} &= (d+1)A - \frac{\beta}{\ell}(d+1)\frac{\partial B}{\partial \beta} + \frac{\alpha\beta}{\ell^{3/2}}(d+1)B + \frac{\alpha\beta}{\ell^2}(d+1)\frac{\partial C}{\partial \alpha} + \frac{\alpha(\alpha^2 - \beta^2)}{\ell^4}(d+1)C \\ \frac{\partial^2 L(\alpha, \beta)}{\partial \alpha \partial \beta} &= -\frac{\beta}{\ell}(d+1)\frac{\partial B}{\partial \alpha} + \frac{\alpha\beta}{\ell^{3/2}}(d+1)B + \frac{\alpha\beta}{\ell^2}(d+1)\frac{\partial C}{\partial \alpha} + \frac{\beta(\beta^2 - \alpha^2)}{\ell^4}(d+1)C \\ \frac{\partial^2 L(\alpha, \beta)}{\partial \beta \partial \alpha} &= -\frac{\alpha}{\ell}(d+1)\frac{\partial B}{\partial \beta} + \frac{\alpha\beta}{\ell^{3/2}}(d+1)B - \frac{\beta^2}{\ell^2}(d+1)\frac{\partial C}{\partial \beta} - 2\frac{\alpha\beta^2}{\ell^4}(d+1)C.\end{aligned}$$

Next we have

$$\begin{aligned}\left| \frac{\partial B}{\partial \alpha} \right| &= \left| \sum_{i,j=1}^n \frac{d_{i,j}}{d+1} \hat{\sigma}_{\text{step}}\left(\frac{d_{i,j}}{d+1} \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}}\right) \cdot \frac{\beta^2}{\ell^{3/2}} \right| \leq \sqrt{\ell}|C| = o(A) \\ \left| \frac{\partial B}{\partial \beta} \right| &= \left| \sum_{i,j=1}^n \frac{d_{i,j}}{d+1} \hat{\sigma}_{\text{step}}\left(\frac{d_{i,j}}{d+1} \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}}\right) \cdot \frac{\alpha\beta}{\ell^{3/2}} \right| \leq \sqrt{\ell}|C| = o(A),\end{aligned}$$

where we have used the fact that $|C| \leq nd^2 = o(A)$.

To differentiate C we use the fact that $\hat{\sigma}_{\text{step}}(\rho) = 1 - \frac{\arccos(\rho)}{\pi}$ [Daniely et al., 2016] and $\arccos'(\rho) = -\frac{1}{\sqrt{1-\rho^2}}$ to get

$$\begin{aligned}\left| \frac{\partial C}{\partial \alpha} \right| &= \left| \sum_{i,j=1}^n \frac{d_{i,j}^2}{(d+1)^2} \frac{\ell}{\beta} \frac{\beta^2}{\ell^{3/2}} \right| \leq \sqrt{\ell} \sum_{i,j} \frac{d_{i,j}^2}{(d+1)^2} = o(A) \\ \left| \frac{\partial C}{\partial \beta} \right| &= \left| \sum_{i,j=1}^n \frac{d_{i,j}^2}{(d+1)^2} \frac{\ell}{\beta} \frac{\alpha\beta}{\ell^{3/2}} \right| \leq \sqrt{\ell} \sum_{i,j} \frac{d_{i,j}^2}{(d+1)^2} = o(A).\end{aligned}$$

Hence we get that

$$\begin{aligned}\left| \frac{\partial^2 L(\alpha, \beta)}{\partial \alpha^2} \right| &\leq (d+1)A + \sqrt{\ell}(d+1)A + o(A) \leq (d+1)A(1 + \sqrt{2\sigma\sqrt{r} + 4\pi + 1} + o(1)) \\ \left| \frac{\partial^2 L(\alpha, \beta)}{\partial \beta^2} \right| &\leq (d+1)A + \sqrt{\ell}(d+1)A + o(A) \leq (d+1)A(1 + \sqrt{2\sigma\sqrt{r} + 4\pi + 1} + o(1)) \\ \left| \frac{\partial^2 L(\alpha, \beta)}{\partial \alpha \partial \beta} \right| &\leq \sqrt{\ell}(d+1)A + o(A) \leq (d+1)A(\sqrt{2\sigma\sqrt{r} + 4\pi + 1} + o(1)) \\ \left| \frac{\partial^2 L(\alpha, \beta)}{\partial \beta \partial \alpha} \right| &\leq \sqrt{\ell}(d+1)A + o(A) \leq (d+1)A(\sqrt{2\sigma\sqrt{r} + 4\pi + 1} + o(1)).\end{aligned}$$

□

Lemma 8 (Restatement of Lemma 3). *If $w_{0,j} \sim N(0, \sigma^2 I)$ and $\eta \in [\frac{1}{16\pi(d+1)A}, \frac{1}{6\pi(d+1)A}]$ then with at least $1 - 1/h^2$, it holds that for all $t \geq c \log(\sigma \log h)$, $\alpha_t \geq -\frac{1}{100}$ and $\ell_t \geq 1 - o(1)$.*

Proof. Notice that due to random initialization, with probability at least $1 - \frac{1}{h^2}$ we will have that $\ell_0 = \theta(\sigma\sqrt{r})$ and $\alpha_0 \geq -c'\sigma\sqrt{\log h}$. Furthermore we have the following updates for ℓ_t and α_t :

$$\begin{aligned}\ell_{t+1}^2 &= (\ell_t(1 - \eta(d+1)A) + \eta(d+1)B_t)^2 + \eta^2 \frac{\beta_t^2}{\ell_t^2} (d+1)^2 C_t^2 \\ \alpha_{t+1} &= \alpha_t(1 - \eta(d+1)A) + \eta \frac{\alpha_t}{\ell_t} (d+1)B_t + \eta \frac{\beta_t^2}{\ell_t^2} (d+1)C_t.\end{aligned}$$

Since $B_t \geq 0$ and is bounded by A and $C_t = o(A)$, if $\alpha_t < 0$ and $\ell_t \geq 2$ we will have that

$$\begin{aligned}\ell_{t+1} &\geq \ell_t(1 - \eta(d+1)A). \\ \alpha_{t+1} &\geq \alpha_t(1 - \frac{\eta}{2}(d+1)A - \eta(d+1)o(A)).\end{aligned}$$

Hence, after $t \geq c \log(\sigma \log h)$ steps we will have that $\alpha_t \geq -\frac{1}{100}$ and $\ell_t \geq 2$. We next argue that from this point on α_t and ℓ_t continue to satisfy the conditions stated in the Lemma. From the update equation for ℓ_{t+1} above we have that

$$\ell_{t+1} \geq \ell_t - \eta(d+1)A(\ell_t - 1) + \eta(d+1)\Delta_t,$$

where

$$\Delta_t := B_t - A = \sum_{i,j:d_{i,j} \neq 0} \hat{\sigma}(\frac{d_{i,j}}{d+1} \frac{\alpha_t}{\ell_t}) - \hat{\sigma}(\frac{d_{i,j}}{d+1}).$$

Once $\alpha_t \geq -\frac{1}{100}$ we will have that

$$\left| \hat{\sigma}(\frac{d_{i,j}}{d+1} \frac{\alpha_t}{\ell_t}) - \hat{\sigma}(\frac{d_{i,j}}{d+1}) \right| \leq 2 \frac{d_{i,j}}{d+1}.$$

Hence defining

$$\Delta = \sum_{i,j:d_{i,j} \neq 0} \frac{d_{i,j}}{d+1}$$

we get that if $\alpha_t \geq -\frac{1}{100}$ then $\ell_t \geq 1 - O(\frac{\Delta}{A}) = 1 - o(1)$. Next we argue that α_t continues to be larger than $-\frac{1}{100}$. We first notice that if $\alpha_t \geq 0$ the is continues to remain so. Furthermore if $\alpha_t \in [-\frac{1}{100}, 0)$, then C_t is non negative and is at least $\Delta/4$. Hence we get that

$$\alpha_{t+1} \geq \alpha_t - \eta(d+1)A\alpha_t(1 - \frac{1}{\ell_t}) + \eta \frac{\alpha_t}{\ell_t} (d+1)\Delta_t + \eta \frac{\beta_t^2}{\ell_t^2} (d+1) \frac{\Delta}{4}.$$

Using the fact that $|\ell_t - 1| = O(\frac{\Delta}{A})$ and the fact that if $\alpha_t \in [-\frac{1}{100}, 0)$ and $\ell_t \geq 1 - o(1)$, then $|\frac{\beta_t}{\ell_t}| \geq \frac{1}{2}$ we get that

$$\alpha_{t+1} \geq \alpha_t - c\eta(d+1)|\alpha_t|\Delta + \eta(d+1)\frac{\Delta}{16} \geq \alpha_t.$$

□

Lemma 9 (Restatement of Lemma 4). *If the degree d of the graph is $o(\sqrt{n})$ and the conditions in Lemma 1 hold then for all $t \geq c \log(\sigma \log h)$, either $|\beta_t| \leq \frac{\epsilon}{4hn}$ and $\|\ell_t - 1\| \leq \frac{\epsilon}{4hn}$ or we have that*

$$\|\nabla L(\alpha_t, \beta_t)\|^2 \geq \mu^* L(\alpha_t, \beta_t),$$

where $\mu^* \geq \frac{\epsilon^2}{380(d+1)h^2\pi n^2}$.

Proof. We have

$$\begin{aligned}\|\nabla L(\alpha_t, \beta_t)\|^2 &= \left(\alpha_t(d+1)A - \frac{\alpha_t}{\ell_t}(d+1)B_t - \frac{\beta_t^2}{\ell_t^2}C_t\right)^2 + \left(\beta_t(d+1)A - \frac{\beta_t}{\ell_t}(d+1)B_t + \frac{\alpha\beta}{\ell_t^2}C_t\right)^2 \\ &= (d+1)^2 \left((\ell_t A - B_t)^2 + \frac{\beta_t^2}{\ell_t^2(d+1)^2}C_t^2\right).\end{aligned}$$

On the other hand the loss $L(\alpha_t, \beta_t)$ can be written as

$$L(\alpha_t, \beta_t) = \frac{1}{2}(d+1)(\ell_t^2 + 1)A - (d+1)\ell_t B \leq \frac{1}{2}(d+1)(\ell_t^2 + 1)A. \quad (26)$$

Hence we get that

$$\frac{\|\nabla L(\alpha_t, \beta_t)\|^2}{L(\alpha_t, \beta_t)} \geq (d+1) \frac{(\ell_t A - B_t)^2}{A} + 2 \frac{(d+1)\beta_t^2 C_t^2}{(\ell_t^2 + 1)A}.$$

If $\ell_t - 1 > \frac{\epsilon}{4hn}$ then the first term above combined with the fact that $B_t \leq A$ leads to

$$\begin{aligned}\frac{\|\nabla L(\alpha_t, \beta_t)\|^2}{L(\alpha_t, \beta_t)} &\geq (d+1)A \frac{\epsilon^2}{16h^2n^2} \\ &\geq (d+1) \frac{\epsilon^2}{16h^2\pi}.\end{aligned} \quad (27)$$

If $|\ell_t - 1| \leq \frac{\epsilon}{4hn} \leq 2$ and $|\beta_t| > \frac{\epsilon}{4hn}$ then the second term leads to

$$\begin{aligned}\frac{\|\nabla L(\alpha_t, \beta_t)\|^2}{L(\alpha_t, \beta_t)} &\geq 2 \frac{(d+1)\beta_t^2 C_t^2}{(\ell_t^2 + 1)h^2(d+1)^2 A} \\ &\geq \frac{C_t^2}{5A(d+1)} \frac{\epsilon^2}{16h^2n^2}.\end{aligned} \quad (28)$$

Using the fact that $C_t \geq \frac{A}{2n}$ we get that

$$\frac{\|\nabla L(\alpha_t, \beta_t)\|^2}{L(\alpha_t, \beta_t)} \geq \frac{\epsilon^2}{380\pi h^2 n^2 (d+1)}. \quad (29)$$

Finally, consider the case when $\ell_t < 1 - \frac{\epsilon}{4hn}$. Furthermore, we can assume that $|B_t - \ell_t A| \leq \frac{\epsilon}{8hn} A$ since otherwise we get the same bound as in (27), up to a loss of a constant factor. In this case we will show that $|\beta_t|$ must be at least $\frac{\epsilon}{4hn}$ and the hence the bound of (28) will be applicable. To see this we use the fact that $\hat{\sigma}(\cdot)$ is convex in $[0, 1]$ to get

$$\hat{\sigma}\left(\frac{d_{i,j}}{d+1} \frac{\alpha_t}{\ell_t}\right) - \hat{\sigma}\left(\frac{d_{i,j}}{d+1}\right) \geq \frac{d_{i,j}}{d+1} \frac{\alpha_t - \ell_t}{\ell_t} \hat{\sigma}_{\text{step}}\left(\frac{d_{i,j}}{d+1}\right). \quad (30)$$

Summing over i, j , we can conclude that

$$B_t - A \geq \frac{\alpha_t - \ell_t}{\ell_t} \sum_{i,j} \frac{d_{i,j}}{d+1} \hat{\sigma}_{\text{step}}\left(\frac{d_{i,j}}{d+1}\right). \quad (31)$$

Substituting $B_t = \ell_t A \pm \frac{\epsilon A}{8hn}$ we get that

$$(\ell_t - 1)A \pm \frac{\epsilon A}{8n} \geq \frac{\alpha_t - \ell_t}{\ell_t} \sum_{i,j} \frac{d_{i,j}}{d+1} \hat{\sigma}_{\text{step}}\left(\frac{d_{i,j}}{d+1}\right). \quad (32)$$

Using the bound on ℓ_t the above implies that

$$\frac{\epsilon A}{8hn} \leq \frac{\ell_t - \alpha_t}{\ell_t} \sum_{i,j} \frac{d_{i,j}}{d+1} \hat{\sigma}_{\text{step}}\left(\frac{d_{i,j}}{d+1}\right) \quad (33)$$

$$\leq \frac{\ell_t - \alpha_t}{\ell_t} n d^2. \quad (34)$$

Noticing that $A \geq \frac{n^2}{\pi}$ we get that

$$\frac{\ell_t - \alpha_t}{\ell_t} \geq \frac{\epsilon}{8\pi h d^2}. \quad (35)$$

Since $\ell_t \geq \frac{1}{36\pi^2}$ we can further conclude that

$$\beta_t \geq \ell_t - \alpha_t \geq \frac{\epsilon}{8 \cdot 36 \cdot h \pi^3 d^2} \geq \frac{\epsilon}{4hn}, \quad (36)$$

where the last inequality uses the fact that $d = o(\sqrt{n})$. \square

Proof of Lemma 5. For property 1, we first recall that $\hat{\sigma}(x) \geq \frac{1}{\pi}$ whenever $x \geq 0$ [Daniely et al., 2016]. Hence, the bound on A follows. To bound B_t we use the fact that $|\hat{\sigma}(x)| \leq \hat{\sigma}(|x|)$, and that $\hat{\sigma}(\cdot)$ is a non-decreasing function in $[0, 1]$. Hence, the upper bound on B_t follows. For the lower bound on B_t notice that whenever $d_{i,j} = 0$, the corresponding term in B_t contributes $\hat{\sigma}(0) = \frac{1}{\pi}$. The remaining at most nd^2 terms can contribute at most -1 each. Therefore we have that $B_t \geq (n^2 - nd^2)\frac{1}{\pi} - nd^2 \geq \frac{n^2}{4\pi}$.

For the upper bound on C_t notice that each term in the summation is non-zero only when $d_{i,j} \neq 0$. Hence, there are at most nd^2 non-zero terms in the summation, and each is upper bounded by 1 since $d_{i,j} \leq d + 1$ and $|\hat{\sigma}(x)| \leq 1$. Using the fact that $d = o(\sqrt{n})$, the upper bound on C_t follows.

For the lower bound on C_t in property 2 we recall that $\hat{\sigma}_{\text{step}}(x) \geq \frac{1}{2}$ whenever $|x| \leq \frac{1}{50}$ [Daniely et al., 2016] which is ensured by the fact that $\alpha_t \geq -\frac{1}{100}$ and $\ell_t \geq 1 - o(1)$. Hence, in this case each term in the summation in the expression of C_t will be non-negative. Finally, notice that for $i = j$, each of the n terms in the summation will contribute at least $\frac{1}{200}$. Hence in this case $C_t \geq \frac{n}{2} \geq \frac{A}{2n}$ (since $A \leq n^2$). \square

Handling unequal degrees. Our main theorem (Theorem 1) continues to hold when the vertices have unequal degrees and the maximum degree is bounded by d . All our proofs easily and (almost) identically translate to the setting of unequal degrees and hence we avoid repeating them. We instead highlight the main differences in the different terms that appear in the proofs. In the case of unequal degrees we redefine the quantities A, B_t, C_t as

$$\begin{aligned} A' &= \sum_{i,j=1}^n \sqrt{(d_i + 1)(d_j + 1)} \hat{\sigma}\left(\frac{d_{i,j}}{\sqrt{(d_i + 1)(d_j + 1)}}\right) \\ B'_t &= \sum_{i,j=1}^n \sqrt{(d_i + 1)(d_j + 1)} \hat{\sigma}\left(\frac{d_{i,j}}{\sqrt{(d_i + 1)(d_j + 1)}} \frac{\alpha_t}{\ell_t}\right) \\ C'_t &= \sum_{i,j=1}^n d_{i,j} \hat{\sigma}_{\text{step}}\left(\frac{d_{i,j}}{\sqrt{(d_i + 1)(d_j + 1)}} \frac{\alpha_t}{\ell_t}\right). \end{aligned}$$

Here, d_i is the degree of node i . In other words we systematically replace $(d+1)$ in all the summations with the term $\sqrt{(d_i + 1)(d_j + 1)}$. Following this change all our Lemmas follow as is.

C Proof from Section 4

C.1 Expressions for Loss and Gradients

In this section we compute expressions for the population loss and the gradients that will be useful in subsequent analysis. We first analyze the population loss.

Lemma 10. *If G is a degree- d graph and node inputs x_i are drawn from $N(0, I)$ then for any W_1, W_2 we have that*

$$L(W_1, W_2) = \frac{n}{2} \|M^L - (M^*)^L\|_F^2, \quad (37)$$

where $M = W_1 + d \cdot W_2$ and $M^* = W_1^* + d \cdot W_2^*$.

Proof. Define $\mu_{\vec{x}} = \sum_{i=1}^n x_i$. We will prove via induction that $\sum_{i=1}^n h_i^{(L)} = M^L \mu_{\vec{x}}$. The case of $L = 1$ is easy to see since

$$\begin{aligned} \sum_{i=1}^n h_i^{(1)} &= \sum_{i=1}^n W_1 x_i + \sum_{i=1}^n W_2 \sum_{j \in N(i)} x_j \\ &= (W_1 + d \cdot W_2) \mu_{\vec{x}} \\ &= M \mu_{\vec{x}}. \end{aligned}$$

Next for any $k > 1$ we have

$$\begin{aligned} \sum_{i=1}^n h_i^{(k)} &= \sum_{i=1}^n W_1 \hat{h}_i^{(k-1)} + \sum_{i=1}^n W_2 \sum_{j \in N(i)} \hat{h}_j^{(k-1)} \\ &= (W_1 + d \cdot W_2) \sum_{i=1}^n \hat{h}_i^{(k-1)} \\ &= (W_1 + d \cdot W_2) M^{k-1} \mu_{\vec{x}} \quad (\text{by induction hypothesis.}) \\ &= M^k \mu_{\vec{x}}. \end{aligned}$$

Hence we get that

$$\begin{aligned} L(W_1, W_2) &= \frac{1}{2} \mathbb{E} [\|M^L \mu_{\vec{x}} - (M^*)^L \mu_{\vec{x}}\|^2] \\ &= \frac{1}{2} \mathbb{E} [\langle (M^L - (M^*)^L) \mu_{\vec{x}}, (M^L - (M^*)^L) \mu_{\vec{x}} \rangle] \\ &= \frac{1}{2} \mathbb{E} [\langle (M^L - (M^*)^L), (M^L - (M^*)^L) \mu_{\vec{x}} \mu_{\vec{x}}^\top \rangle] \\ &= \frac{n}{2} \|M^L - (M^*)^L\|_F^2. \end{aligned}$$

□

Next we compute the expression for the gradient of the loss w.r.t. the weights.

Lemma 11. *For the multi-round GNN as defined in Eq. (18) it holds that*

$$\frac{\partial L(W_1, W_2)}{\partial W_1} = \frac{n}{2} \sum_{k=0}^{L-1} M^k E M^{L-k-1} \quad (38)$$

$$\frac{\partial L(W_1, W_2)}{\partial W_2} = \frac{nd}{2} \sum_{k=0}^{L-1} M^k E M^{L-k-1}, \quad (39)$$

where $E = M^L - (M^*)^L$.

Proof. For i and j in $[r]$,

$$\frac{\partial L(W_1, W_2)}{\partial W_1^{i,j}} = \frac{n}{2} \langle \text{vec}(E), \text{vec}(\frac{\partial M^L}{\partial W_1^{i,j}}) \rangle.$$

Furthermore, by chain rule we get that

$$\frac{\partial M^L}{\partial W_1^{i,j}} = \sum_{k=0}^{L-1} M^k R_{i,j} M^{L-k-1},$$

where $R_{i,j}$ is an $r \times r$ matrix with 1 in the (i, j) th entry and 0 everywhere else. Denoting by A_i to be the i th column of a matrix we can simplify the above as

$$\frac{\partial M^L}{\partial W_1^{i,j}} = \sum_{k=0}^{L-1} (M^k)_i \otimes (M^{L-k-1})_j.$$

Hence we get that

$$\begin{aligned}
\frac{\partial L(W_1, W_2)}{\partial W_1^{i,j}} &= \frac{n}{2} \langle \text{vec}(E), \text{vec}(\sum_{k=0}^{L-1} (M^k)_i \otimes (M^{L-k-1})_j) \rangle \\
&= \frac{n}{2} \sum_{k=0}^{L-1} \langle \text{vec}(E), \text{vec}((M^k)_i \otimes (M^{L-k-1})_j) \rangle \\
&= \frac{n}{2} \sum_{k=0}^{L-1} (M^k)_i^\top E (M^j)^{K-k-1}.
\end{aligned}$$

From the above it follows that

$$\frac{\partial L(W_1, W_2)}{\partial W_1} = \frac{n}{2} \sum_{k=0}^{L-1} M^k E M^{L-k-1}.$$

An identical calculation establishes the bound for the partial derivative w.r.t. W_2 . \square

C.2 Analyzing Gradient Updates

In this section we analyze the gradient descent updates as defined in (22) and prove our main convergence result. To begin with we will track the evolution of the matrix M_t defined as $M_t := W_{1,t} + d \cdot W_{2,t}$. From Lemma 11 we get that

$$M_{t+1} = M_t - \eta \frac{n}{2} (d^2 + 1) \sum_{k=0}^{L-1} M_t^k E_t M_t^{L-k-1}. \quad (40)$$

Let the singular value decomposition of M^* be denoted by $M^* = \sum_{i=1}^r \sigma_i v_i v_i^\top$. We will ensure that throughout the trajectory of the updates the matrix M_t remains positive definite in which case we can write $M_t = \sum_{i=1}^r \sigma_{i,t} v_i v_i^\top$. This is true for $t = 0$ since we initialize the matrices W_1, W_2 to identity. Then from (40) we get

$$\begin{aligned}
\sum_{i=1}^r \sigma_{i,t+1} v_i v_i^\top &= \sum_{i=1}^r \sigma_{i,t} v_i v_i^\top - \eta \frac{n}{2} (d^2 + 1) \sum_{k=0}^{L-1} \left(\sum_{i=1}^r \sigma_{i,t}^k v_i v_i^\top \right) \left(\sum_{i=1}^r (\sigma_{i,t}^L - (\sigma_i^*)^L) v_i v_i^\top \right) \left(\sum_{i=1}^r \sigma_{i,t}^{L-k-1} v_i v_i^\top \right) \\
&= \sum_{i=1}^r \sigma_{i,t} v_i v_i^\top - \eta \frac{n}{2} (d^2 + 1) \sum_{k=0}^{L-1} \sum_{i=1}^r \sigma_{i,t}^{L-1} (\sigma_{i,t}^L - (\sigma_i^*)^L) \\
&= \sum_{i=1}^r (\sigma_{i,t} - \eta \frac{n}{2} L (d^2 + 1) \sigma_{i,t}^{L-1} (\sigma_{i,t}^L - (\sigma_i^*)^L)) v_i v_i^\top.
\end{aligned}$$

Hence each singular value evolves as

$$\sigma_{i,t+1} := \sigma_{i,t} - \eta \frac{n}{2} g_t \quad (41)$$

$$= \sigma_{i,t} - \eta \frac{n}{2} L (d^2 + 1) \sigma_{i,t}^{L-1} (\sigma_{i,t}^L - (\sigma_i^*)^L). \quad (42)$$

It is enough to show that over a period of time the singular values σ_i are getting closer to σ_i^* . In order to do this we first prove that the singular values remain bounded throughout the trajectory of the iterates.

Lemma 12. *Let the initialization be a symmetric matrices $W_{1,0}, W_{2,0}$ such that $M_0 = W_{1,0} + d \cdot W_{2,0} = \sum_{i=1}^n \sigma_{i,0} v_i v_i^\top$, where $\sigma_{i,0} > 0$ for all i . For each $i \in [r]$ define*

$$\begin{aligned}
u_i &= \sigma_i^* \left(1 + \max\left(1, \frac{\sigma_{i,0}}{\sigma_i^*} - 1\right) \right) \\
\ell_i &= \sigma_i^* \left(1 - \max\left(\frac{1}{2}, 1 - \frac{\sigma_{i,0}}{\sigma_i^*}\right) \right).
\end{aligned}$$

Then for all $t \geq 0$ it holds that $\sigma_{i,t} \in [\ell_i, u_i]$ provided that

$$\eta \leq \eta_i := \frac{\max(1, \frac{\sigma_{i,0}}{\sigma_i^*} - 1)}{2nL^2u_i^L(d^2 + 1)(\sigma_i^*)^{2L-2}(1 + u_i)^{L-1}}.$$

Proof. We will prove the claim by induction. At $t = 0$ the bound holds by definition. Consider time t and assume that $\sigma_i^t = \sigma_i^*(1 + \delta)$ where $\delta > 0$. The case for $\delta \leq 0$ will be similar. In this case we have that

$$\begin{aligned} \sigma_{i,t+1} &= \sigma_{i,t} - \eta \frac{n}{2} L(d^2 + 1) \sigma_{i,t}^{L-1} (\sigma_{i,t}^L - (\sigma_i^*)^L) \\ &\leq \sigma_{i,t} \leq u_i, \quad (\text{by induction hypothesis.}) \end{aligned}$$

Furthermore we have

$$\begin{aligned} \sigma_{i,t+1} &= \sigma_{i,t} - \eta \frac{n}{2} L(d^2 + 1) \sigma_{i,t}^{L-1} (\sigma_{i,t}^L - (\sigma_i^*)^L) \\ &= \sigma_i^* (1 + \delta - \eta \frac{n}{2} L(d^2 + 1) (\sigma_i^*)^{2L-2} (1 + \delta)^{L-1} ((1 + \delta)^L - 1)) \\ &\geq \sigma_i^* (1 - \eta \frac{n}{2} L(d^2 + 1) (\sigma_i^*)^{2L-2} (1 + u_i)^{L-1} \cdot Lu_i^L) \quad (\text{by induction hypothesis.}) \\ &\geq \frac{1}{2} \sigma_i^*, \end{aligned}$$

where the last inequality holds provided that

$$\eta \leq \frac{1}{2nL^2u_i^L(d^2 + 1)(\sigma_i^*)^{2L-2}(1 + u_i)^{L-1}}.$$

Next consider the case when $\sigma_{i,t} = \sigma_i^*(1 - \delta)$ for $\delta > 0$. Then we have

$$\begin{aligned} \sigma_{i,t+1} &= \sigma_{i,t} - \eta \frac{n}{2} L(d^2 + 1) \sigma_{i,t}^{L-1} (\sigma_{i,t}^L - (\sigma_i^*)^L) \\ &\geq \sigma_{i,t} \geq \ell_i, \quad (\text{by induction hypothesis.}) \end{aligned}$$

Finally, we have

$$\begin{aligned} \sigma_{i,t+1} &= \sigma_{i,t} - \eta \frac{n}{2} L(d^2 + 1) \sigma_{i,t}^{L-1} (\sigma_{i,t}^L - (\sigma_i^*)^L) \\ &= \sigma_i^* (1 - \delta + \eta \frac{n}{2} L(d^2 + 1) (\sigma_i^*)^{2L-2} (1 + \delta)^{L-1} (1 - (1 - \delta)^L)) \\ &\leq \sigma_i^* (1 + \eta \frac{n}{2} L(d^2 + 1) (\sigma_i^*)^{2L-2} (1 + u_i)^{L-1} \cdot Lu_i^L) \quad (\text{by induction hypothesis.}) \\ &\leq u_i, \end{aligned}$$

where the last inequality holds provided that

$$\eta \leq \frac{\max(1, \frac{\sigma_{i,0}}{\sigma_i^*} - 1)}{2nL^2u_i^L(d^2 + 1)(\sigma_i^*)^{2L-2}(1 + u_i)^{L-1}}.$$

□

Next we use the above bounds to establish two useful properties of the gradient value g_t in (41).

Lemma 13. *Let u_i, ℓ_i be as defined in Lemma 12. Then for all $t \geq 0$ it holds that*

$$\begin{aligned} g_t^2 &\leq L^2(d^2 + 1)^2 u_i^{2L-3} (L\delta_i)^2 (\sigma_{i,t} - \sigma_i^*)^2 \\ g_t(\sigma_{i,t} - \sigma_i^*) &\geq L^2(d^2 + 1) \ell_i^{2L-2} (\sigma_{i,t} - \sigma_i^*)^2. \end{aligned}$$

Here $\delta_i := \max(\frac{u_i}{\sigma_i^*}, \frac{\sigma_i^*}{\ell_i})$.

Proof. For the upper bound on g_t notice that

$$\begin{aligned} g_t^2 &= L^2(d^2 + 1)^2 \sigma_{i,t}^{2L-2} (\sigma_{i,t}^L - (\sigma_i^*)^L)^2 \\ &\leq L^2(d^2 + 1)^2 u_i^{2L-2} (\sigma_{i,t}^L - (\sigma_i^*)^L)^2 \\ &\leq L^2(d^2 + 1)^2 u_i^{2L-3} (L\delta_i)^2 (\sigma_{i,t} - \sigma_i^*)^2. \end{aligned}$$

Next notice that

$$g_t(\sigma_{i,t} - \sigma_i^*) = L(d^2 + 1) \sigma_{i,t}^{L-1} (\sigma_{i,t}^L - (\sigma_i^*)^L) (\sigma_{i,t} - \sigma_i^*) \quad (43)$$

$$\geq L(d^2 + 1) \ell_i^{L-1} (\sigma_{i,t}^L - (\sigma_i^*)^L) (\sigma_{i,t} - \sigma_i^*) \quad (44)$$

$$\geq L^2(d^2 + 1) \ell_i^{2L-2} (\sigma_{i,t} - \sigma_i^*)^2, \quad (45)$$

where the last inequality follows from the fact that $(a^L - b^L)(a - b) \geq L \min(a, b)^{L-1} (a - b)^2$. \square

We are now ready to prove our main theorem regarding convergence of gradient descent on multi-round linear GNNs.

Theorem 3 (Restatement of Theorem 2.). *Let u_i, ℓ_i, δ_i be as defined in Lemma 12 and Lemma 13.*

Then we have that $L(W_{1,T}, W_{2,T}) \leq \epsilon^2$ provided $T \geq \max_i \frac{1}{\eta_0 L^2 (d^2 + 1) \ell_i^{2L-2}} \log\left(\frac{r L u_i ((1 - \sigma_i^)^2)}{\epsilon}\right)$ and*

$$\eta := \eta_0 \leq \min_i \frac{\ell_i^{2L-1} \min(1, \frac{1}{(\sigma_i^*)^{2L-2}})}{2nL^2(\delta_i)^2(d^2 + 1)u_i^{2L-3}(1 + u_i)^{L-1}}.$$

Proof. Fix an $i \in [r]$. We have

$$\begin{aligned} (\sigma_{i,t+1} - \sigma_i^*)^2 &= (\sigma_{i,t} - \eta \frac{n}{2} g_t - \sigma_i^*)^2 \\ &= (\sigma_{i,t} - \sigma_i^*)^2 + \eta_0^2 \frac{n^2}{4} g_t^2 - \eta_0 n g_t (\sigma_{i,t} - \sigma_i^*) \\ &\leq (\sigma_{i,t} - \sigma_i^*)^2 + \eta_0^2 \frac{n^2}{4} (\delta_i^*)^2 L^4 (d^2 + 1)^2 u_i^{2L-3} (\sigma_{i,t} - \sigma_i^*)^2 - \eta_0 n L^2 (d^2 + 1) \ell_i^{2L-2} (\sigma_{i,t} - \sigma_i^*)^2 \quad (\text{Lemma 13}) \\ &\leq (\sigma_{i,t} - \sigma_i^*)^2 (1 - \eta_0 n L^2 (d^2 + 1) \ell_i^{2L-2}) \\ &\leq (1 - \sigma_i^*)^2 (1 - \eta_0 n L^2 (d^2 + 1) \ell_i^{2L-2})^t. \end{aligned}$$

Hence after $T \geq \max_i \frac{1}{\eta_0 L^2 (d^2 + 1) \ell_i^{2L-2}} \log\left(\frac{nr L u_i ((1 - \sigma_i^*)^2)}{\epsilon}\right)$ time steps we will have that $|\sigma_{i,T}^L - (\sigma_i^*)^L| \leq \epsilon/(\sqrt{nr})$ for all i . This implies an ϵ^2 upper bound on the loss since we have that

$$\begin{aligned} L(W_{1,T}, W_{2,T}) &= \frac{n}{2} \|M_T^L - (M^*)^L\|_F^2 \\ &\leq nr \max_i (\sigma_{i,T}^L - (\sigma_i^*)^L)^2 \\ &\leq \epsilon^2. \end{aligned}$$

\square

Handling unequal degrees. Similar to the case of one round ReLU GNNs, handling the case of unequal (but bounded by d) vertex degrees follows from our main analysis above at the expense of more complicated expressions. We next discuss the main changes. In the case of unequal vertex degrees we have that the loss function equals:

$$L(W_1, W_2) = \frac{1}{2} \sum_{i=1}^n \|M_i^L - (M^*)^L\|_F^2. \quad (46)$$

Here $M_i = W_1 + d_i W_2$ and $M_i^* = W_1^* + d_i W_2^*$. Similarly, the gradient expressions become as stated below.

$$\frac{\partial L(W_1, W_2)}{\partial W_1} = \frac{1}{2} \sum_{i=1}^n \sum_{k=0}^{L-1} M_i^k E_i M_i^{L-k-1} \quad (47)$$

$$\frac{\partial L(W_1, W_2)}{\partial W_2} = \frac{1}{2} \sum_{i=1}^n d_i \sum_{k=0}^{L-1} M_i^k E_i M_i^{L-k-1}, \quad (48)$$

where $E_i = M_i^L - (M^*)_i^L$. In contrast to the case of equal degrees where it was enough to track the evolution of singular values of $W_1 + dW_2$, in the general case we will track the evolution of the singular values of W_1 and W_2 separately. Denote by $\alpha_{p,t}, \alpha_p^*$ the p th singular values of W_1, t and W_1^* respectively. Similarly, denote by $\beta_{p,t}, \beta_p^*$ the p th singular values of W_2, t and W_2^* respectively. Then we get the following recurrence.

$$\alpha_{p,t+1} = \alpha_{p,t} - \eta \sum_{i=1}^n g_{i,p,t} \quad (49)$$

$$\beta_{p,t+1} = \beta_{p,t} - \eta \sum_{i=1}^n d_i g_{i,p,t}. \quad (50)$$

Here $g_{i,p,t}$ is given by the following expression.

$$g_{i,p,t} = (\alpha_{p,t} + d_i \beta_{p,t})^{L-1} ((\alpha_{p,t} + d_i \beta_{p,t})^L - (\alpha_p^* + d_i \beta_p^*)^L). \quad (51)$$

Next we further simplify the above to get the following updates that are easier to analyze.

$$\alpha_{p,t+1} - \alpha_p^* = (\alpha_{p,t} - \alpha_p^*) \left(1 - \eta \sum_{i=1}^n r_i\right) - \eta (\beta_{p,t} - \beta_p^*) \sum_{i=1}^n d_i r_i \quad (52)$$

$$\beta_{p,t+1} - \beta_p^* = (\beta_{p,t} - \beta_p^*) \left(1 - \eta \sum_{i=1}^n d_i^2 r_i\right) - \eta (\alpha_{p,t} - \alpha_p^*) \sum_{i=1}^n d_i r_i. \quad (53)$$

Here r_i is defined as

$$r_i = (\alpha_{p,t} + d_i \beta_{p,t})^{L-1} \cdot (\alpha_p^* + d_i \beta_p^*)^{L-1} \gamma_i, \quad (54)$$

and γ_i equals

$$\gamma_i = \sum_{j=0}^{L-1} \binom{L}{t} s_i^j. \quad (55)$$

Finally, δ equals

$$s_i = \frac{\alpha_{p,t} + d_i \beta_{p,t}}{\alpha_p^* + d_i \beta_p^*} - 1. \quad (56)$$

Next we prove the following analog of Lemma 12.

Lemma 14. *There exist ℓ_α, u_α and ℓ_β, u_β such that for all $t \geq 0$ and p , $\alpha_{p,t} \in [\ell_\alpha, u_\alpha]$ and $\beta_{p,t} \in [\ell_\beta, u_\beta]$, provided that*

$$\eta := \eta_0 \leq \min \left(\frac{\ell_\alpha^{2L-1} \min(1, \frac{1}{(\sigma_i^*)^{2L-2}})}{2nL^2 u_\alpha^{2L-3} (1 + u_\alpha)^{L-1} (\max_i \delta_{i,\alpha})^2 (d^2 + 1)}, \frac{\ell_\beta^{2L-1} \min(1, \frac{1}{(\sigma_i^*)^{2L-2}})}{2nL^2 u_\beta^{2L-3} (1 + u_\beta)^{L-1} (\max_i \delta_{i,\beta})^2 (d^2 + 1)} \right),$$

where $\delta_{i,\alpha} := \max(\frac{u_\alpha}{\alpha_i^*}, \frac{\alpha_i^*}{\ell_\alpha})$, and $\delta_{i,\beta} := \max(\frac{u_\beta}{\beta_i^*}, \frac{\beta_i^*}{\ell_\beta})$.

Proof. We will prove the lemma by a case analysis and show that provided η is small enough there exists a choice of values $\ell_\alpha, u_\alpha, \ell_\beta, u_\beta$. We will focus on a particular p since the analysis is similar for all p . Hence in the rest of the lemma we will omit the subscript on p for notational convenience.

1. $\beta_t \geq \beta^*$, and $\alpha_t \geq \alpha^*$, $\sum_{i=1}^n g_{i,t} \geq 0$ and $\sum_{i=1}^n d_i g_{i,t} \geq 0$. In this case we have

$$\begin{aligned}\alpha_{t+1} &\geq \alpha^* - \eta n G_u \\ \alpha_{t+1} &\geq \beta^* - \eta n d G_u.\end{aligned}$$

Here G_u is a uniform upper bound on $\sum_{i=1}^n g_{i,t}$ and $\sum_{i=1}^n d_i g_{i,t}$ (See Eq. (49)). Then for the updates to remain bounded we need

$$\begin{aligned}\eta &\leq \frac{\alpha^* - \ell_\alpha}{n G_u} \\ \eta &\leq \frac{\beta^* - \ell_\beta}{n d G_u}.\end{aligned}$$

2. $\beta_t \leq \beta^*$, $\alpha_t \leq \alpha^*$, $\sum_{i=1}^n g_{i,t} \leq 0$ and $\sum_{i=1}^n d_i g_{i,t} \leq 0$. In this case we have

$$\begin{aligned}\alpha_{t+1} &\leq \alpha^* + \eta n G_u \\ \beta_{t+1} &\leq \beta^* + \eta n d G_u.\end{aligned}$$

For the updates to remain bounded we need

$$\begin{aligned}\eta &\leq \frac{u_\alpha - \alpha^*}{n G_u} \\ \eta &\leq \frac{u_\beta - \beta^*}{n d G_u}.\end{aligned}$$

3. $\beta_t \leq \beta^*$, $\alpha_t \geq \alpha^*$, $\sum_{i=1}^n g_{i,t} \geq 0$ and $\sum_{i=1}^n d_i g_{i,t} \geq 0$. This in particular implies that $\alpha_t - \alpha^* \geq \beta_t - \beta^*$. In this case we have

$$\begin{aligned}\alpha_{t+1} &\geq \alpha^* - \eta n G_u \\ \beta_{t+1} &\geq \beta_t - \eta n d G_u \\ &\geq \beta^* + \alpha^* - u_\alpha - \eta n d G_u\end{aligned}$$

For the updates to remain bounded we need

$$\eta \leq \frac{\beta^* + \alpha^* - u_\alpha - \ell_\beta}{n d G_u}.$$

4. $\beta_t \geq \beta^*$, $\alpha_t \leq \alpha^*$, $\sum_{i=1}^n g_{i,t} \geq 0$ and $\sum_{i=1}^n d_i g_{i,t} \geq 0$. This in particular implies that $\alpha_t - \alpha^* \leq \beta_t - \beta^*$. In this case we have

$$\begin{aligned}\beta_{t+1} &\geq \beta^* - \eta n G_u \\ \alpha_{t+1} &\geq \alpha_t - \eta n G_u \\ &\geq \beta^* + \alpha^* - u_\beta - \eta n G_u\end{aligned}$$

For the updates to remain bounded we need

$$\eta \leq \frac{\beta^* + \alpha^* - u_\beta - \ell_\alpha}{n G_u}.$$

5. $\beta_t \geq \beta^*$, $\alpha_t \leq \alpha^*$, $\sum_{i=1}^n g_{i,t} \leq 0$ and $\sum_{i=1}^n d_i g_{i,t} \leq 0$. This in particular implies that $\alpha_t - \alpha^* \geq \beta_t - \beta^*$. In this case we have

$$\begin{aligned}\alpha_{t+1} &\leq \alpha^* + \eta n G_u \\ \beta_{t+1} &\leq \beta^* + \alpha^* - \ell_\alpha + \eta n d G_u\end{aligned}$$

For the updates to remain bounded we need

$$\eta \leq \frac{\beta^* + \alpha^* - u_\beta}{n d G_u}.$$

The other cases not covered above lead to invalid configurations. It is then easy to see that given the chosen value of η_0 , there is a setting of $\ell_\alpha, u_\alpha, \ell_\beta, u_\beta$ that satisfies all the constraints above. \square

Given the above lemma we obtain, analogously to Theorem 2, the following general theorem.

Theorem 4. *Let $\ell_\alpha, u_\alpha, \ell_\beta, u_\beta$ and η_0 be as defined in Lemma 14. Then we have that $L(W_{1,T}, W_{2,T}) \leq \epsilon^2$ provided that*

$$T \geq \max \left(\max_i \frac{1}{\eta_0 L^2 (d^2 + 1) \ell_\alpha^{2L-2}} \log \left(\frac{r L u_\alpha ((1 - \alpha_i^*)^2)}{\epsilon} \right), \max_i \frac{1}{\eta_0 L^2 (d^2 + 1) \ell_\beta^{2L-2}} \log \left(\frac{r L u_\beta ((1 - \beta_i^*)^2)}{\epsilon} \right) \right). \quad (57)$$

Proof. The proof is analogous to the proof of Theorem 2. For notational convenience we define $P = 1 - \eta \sum_{i=1}^n r_i$, $Q = \sum_{i=1}^n d_i r_i$, and $R = 1 - \eta \sum_{i=1}^n d_i^2 r_i$. Then we have

$$\begin{aligned} (\alpha_{p,t+1} - \alpha_p^*)^2 + (\beta_{p,t+1} - \beta_p^*)^2 &= (\alpha_{p,t} - \alpha_p^*)^2 (P^2 + Q^2) + (\beta_{p,t} - \beta_p^*)^2 (R^2 + Q^2) \\ &\quad - 2(\alpha_{p,t} - \alpha_p^*)(\beta_{p,t} - \beta_p^*)(PQ + RQ) \\ &\leq (\alpha_{p,t} - \alpha_p^*)^2 (P^2 + Q^2 + Q(P + R)) \\ &\quad + (\beta_{p,t} - \beta_p^*)^2 (P^2 + R^2 + Q(P + R)) \\ &\leq (P^2 + R^2 + Q^2 + Q(P + R)) ((\alpha_{p,t} - \alpha_p^*)^2 + (\beta_{p,t} - \beta_p^*)^2). \end{aligned} \quad (58)$$

(59)

Next it is easy to verify that from our choice of η , the above decays at a geometric rate and the statement of the theorem follows. \square

D Further details on experiments

We simulate population gradient descent and implement our networks using the JAX programming language [Bradbury et al., 2018]. Our experiments are run using one GPU. Below we present the hyperparameters used in our experiments.

One round ReLU GNNs. For the case of degree $d = 5$ (see Figure 1) we use a learning rate of 2×10^{-5} . For the remaining degrees, we use a learning rate of 1×10^{-6} .

Multi round linear GNNs. For the case of depth $L = 1$ (see Figure 2) we use a learning rate of 0.1. For depths 2 and 3, we use a learning rate of 0.0075.