Improved Convergence Guarantees for Shallow Neural Networks

Alexander Razborov*

December 6, 2022

Abstract

We continue a long line of research aimed at proving convergence of depth 2 neural networks, trained via gradient descent, to a global minimum. Like in many previous works, our model has the following features: regression with quadratic loss function, fully connected feedforward architecture, RelU activations, Gaussian data instances and network initialization, adversarial labels. It is more general in the sense that we allow both layers to be trained simultaneously and at different rates.

Our results improve on state-of-the-art [OS20] (training the first layer only) and [Ngu21, Section 3.2] (training both layers with Le Cun's initialization). We also report several simple experiments with synthetic data. They strongly suggest that, at least in our model, the convergence phenomenon extends well beyond the "NTK regime".

1. Introduction

Deep neural networks are remarkably successful at defying common wisdom inherited from the statistical learning theory and especially the theory of PAC learning. It would perhaps not entirely be an exaggeration to say that the wonder at their ability to generalize in a (highly) overparameterized regime is expressed in at least every other paper on the theory of deep learning

^{*}University of Chicago, USA, razborov@uchicago.edu, Toyota Technological Institute at Chicago and Steklov Mathematical Institute, Moscow, Russia.

so we confine ourselves to citing two expository papers [BHMM19, ZBH+21] specifically devoted to this phenomenon. On more concrete side, recent works [BFT17, NBS18, AGNZ18, LL18, ADH+19b, SY19] offer explanations of the "unbelievably good generalization" from three pairwise disjoint perspectives.

Since there does not seem to be any reasons to expect a learning system to perform on the test examples better than on the training examples, an even more basic theoretical question is why an algorithm, which in this area almost always means gradient descent or its variants, achieves zero empirical risk. It would be fair to say that this question has met with a slightly better degree of success, particularly in the situation described by various authors as the "lazy training regime" (see e.g. [COB19]) or the "neural tangent (kernel) regime" (see e.g. [MZ20]). Notwithstanding the name, this is the behavior described by the following (very much related) properties:

- only "a few" instance-neuron pairs change their activation during training;
- network parameters do not change "much" (say, in the Frobenius norm) during training.

These properties in particular imply that the gradient of the loss function and hence the NTK¹ matrix H_t also do not change much from their value at initialization. This implies that the training trajectory is almost linear and is driven, in the error space, by the value H_0 of the NTK matrix at the initialization. In the following discussion, we will adopt the terminology "NTK regime" to (loosely) describe this kind of behavior.

As a by-side remark, this regime has been widely criticized for being irrelevant to practice, and (as far as we understand) one of the main reasons is that the above property essentially implies that the network does not learn any useful features during training. While this criticism does not look to us entirely unfounded, we would like to remark that, at least for regression problems, the evidence of convergence outside the NTK regime is either fragmentary or empirical at best (we will make a modest contribution to the latter in Section 4). From the mathematical point of view, the best, and often the only, way of approaching a difficult problem is to patiently build up

¹Neural Tangent Kernel; we remind the precise definition below but essentially it is $J_f J_f^{\mathsf{T}}$, where J_f is the Jacobian of the map f recording the network's output on m data instances

necessary tools by understanding its toy cases which in our context amounts to the NTK regime.

Previous work on the subject displays quite a diverse array of various assumptions on the activation and loss functions, data, initialization, as well as (also quite diverse) set of parameters featuring in the final results. They are not easy to compare to each other. In order to get a more coherent picture, we find it convenient to start with a "theoretical benchmark". It seems the following captures a great deal of previous work (in the sense that their message is not distorted much when reduced to this scenario); we will separately mention several notable exceptions in Section 1.4.

1.1. Set-up

Feedforward fully connected neural networks of depth 2 (that is, with one hidden layer). RelU activation function. Regression tasks with single output and the squared loss function. Data points $(X^1, \ldots X^m)$ are sampled uniformly and independently from the unit sphere. Weights at initialization are Gaussian at the input level and Rademacher ($\{\pm 1\}$) at the output level. Parameters:

- n the problem dimension;
- S the number of neurons at the middle level ("size" of the network);
- m the number of training examples.

When $m \leq n$, the data is linearly separable. We will not exclude this case from statements and proofs but in Section 1 we will assume for simplicity that $m \gg n$. We will also tacitly assume $m \leq n^{O(1)}$; this assumption is indispensable in any work using NTK methods, including ours.

1.2. Previous work

The over-parameterized region is $m \ll nS$ and, clearly, no convergence or even representability is possible when $m \gg nS$. Several authors have loosely conjectured that perhaps this is actually tight and the (stochastic) gradient descent does converge in this whole (or at least close to it) region. This conjecture is not as far-fetched as it may seem since under this assumption

alone one can make, with overwhelming probability, quite strong conclusions at initialization like:

[SH18] The activation matrix $A \in \{0,1\}^{S \times m}$ at initialization possesses an intrinsic and purely combinatorial property ensuring the following. The corresponding NTK matrix is non-singular for *almost all* (in the Lebesque sense) choices of data $X \in \mathbb{R}^{n \times m}$.

[XLS, OS20, ADH⁺19a, NMM21, MZ20] The NTK matrix w.r.t. the same data X that were used for computing the activation matrix A has the minimum eigenvalue² $\Omega(S)$, i.e. well-separated from 0.

As for the dynamics of the gradient descent, the known results are much less conclusive. Before briefly surveying them, let us again stress that many of these results pertain to more general situations such as deeper networks, less restrictive assumptions on data, more general loss or activation functions etc. In our simplified treatment we have deliberately reduced them to the "common denominator" outlined in Section 1.1.

Let us first consider the case when only the first layer is trained. [DZPS19] proved that convergence takes place whenever³ $m \leq \tilde{o}\left(S^{1/6}\right)$ (plug $m \mapsto S$; $n \mapsto m$; $\lambda_0, \delta \mapsto O(1)$ in their Theorem 4.1); apparently, this can be approved to $m \leq \tilde{o}\left(n^{1/6}S^{1/6}\right)$ by examining their proof more carefully. The paper [WDW19] improved the bound on the learning rate from [DZPS19] using more sophisticated (adaptive) gradient methods. The improvement $m \leq \tilde{o}\left(n^{1/2}S^{1/4}\right)$ can be relatively easily extracted from [SY19] (plug $m \mapsto S$; $n \mapsto m$; $\lambda \mapsto O(1)$; $\alpha \mapsto \tilde{O}(m/n)$; $\theta \mapsto \tilde{O}(m^{1/2}/n^{1/2})$ in their Theorem 1.6). Finally, the paper [OS20] achieved the state-of-the-art result: convergence takes place whenever

$$m \le \widetilde{o}\left(n^{3/4}S^{1/4}\right). \tag{1}$$

Our description may look somewhat uneventful but we would like to stress that every new achievement along these lines required introducing new ingenious and technically advanced ideas and methods.

 $^{^{2}}$ We do not normalize by S.

³We use the notation $\widetilde{O}, \widetilde{\Omega}, \widetilde{\Theta}, \widetilde{o}, \widetilde{\omega}$ to hyde factors that are poly-logarithmic in n, S, and their tildeless versions to hyde constant factors. See Remark 1 below for clarifying remarks and examples.

Less work seems to have been done on the case when both layers are trained simultaneously, and now the answer may depend on the normalization of weights at the initialization even if we disregard the learning rate. The original paper [DZPS19] gives the same estimate $m \leq \tilde{o}\left(S^{1/6}\right)$ under the "standard" normalization, when the initial weights are supposed to be of the same order (say, $\tilde{O}(1)$) at both layers. The paper [Ngu21] is mostly devoted to deep networks but in Section 3.2 it also translates the main result to the shallow case under the so-called LeCun's initialization. After performing slightly more careful calculations, in our notational system this amount to the condition

$$m \le \min\left(\widetilde{o}\left(\frac{S}{n}\right), \ \widetilde{o}\left(S^{1/2}\right)\right).$$
 (2)

1.3. Our contributions

We work in the natural multi-rate setting in which the first layer is trained at a rate $\eta_w \geq 0$ and the second layer is trained at possibly another rate $\eta_z \geq 0$. This view appears to be very instructive and allows us to easily translate results from one normalization of initial weights to another. In particular, we choose to work in the standard normalization, when all weights are of order $\tilde{O}(1)$.

As a by-side remark, the recent paper [VL22] expressed surprise (that we completely share) that this simple idea seems to have been largely overlooked in the theory of deep learning and sought to change that by providing, among other things, significant experimental evidence. One difference between our paper and [VL22] is that the latter restricts the set of steps in which a "slow" set of parameters is changed while we treat all steps uniformly. But this difference appears to be cosmetic.

Skipping a few inessential details (see Theorem 2.3 for the official statement), our main contribution is as follows.

Theorem 1.1 (informal) In the set-up described above, assume that $n, S \ge \widetilde{\omega}(1)$, $m \le n^{O(1)}$ and that $m \le \widetilde{o}(S)$. Assume additionally that one of the following happens:

$$m \le \tilde{o}\left(nS^{1/4}\right), \quad \frac{\eta_z}{\eta_w} \le \tilde{o}\left(\frac{S}{m}\right),$$
 (3)

$$m \le \tilde{o}\left(nS^{1/4}\right), \quad \frac{\eta_z}{\eta_w} \ge \tilde{\omega}\left(\frac{m^2}{nS}\right)$$
 (4)

$$\frac{\eta_z}{\eta_w} \ge \widetilde{\omega} \left(1 + \frac{m}{n} \right). \tag{5}$$

Assume also that η_w, η_z are sufficiently small. Then with probability $\geq 1 - (nS)^{-\omega(1)}$ the gradient descent converges to a global minimum.

Remark 1 This paper uses the asymptotic notation $\widetilde{O}, \widetilde{\Omega}, \widetilde{o}, \widetilde{\omega}$ (and their tildeless versions) perhaps more systematically and heavily than most papers on the subject. Thus it would be appropriate to pause for a moment and reflect on what exactly it means.

We argue about infinite sequences $\left\{\left(n^{(k)},m^{(k)},S^{(k)},\eta_w^{(k)},\eta_z^{(k)}\right)\right\}_{k\geq 1}$ of parameter values. Then (say) $m\leq \tilde{o}\left(nS^{1/4}\right)$ means that for any fixed C>0,

$$\lim_{k \to \infty} \frac{m^{(k)} \log(n^{(k)} S^{(k)})^C}{n^{(k)} (S^{(k)})^{1/4}} = 0$$
(6)

while the error bound $(nS)^{-\omega(1)}$ simply states that it is asymptotically smaller than any fixed polynomial in⁴ n, S. The condition $n, S \geq \widetilde{\omega}(1)$ says that S and n are (very mildly) balanced: $S \geq (\log n)^{\omega(1)}$ and $n \geq (\log S)^{\omega(1)}$; it appeared in most of the previous work.

For a bit of further discussion see also Remark 7 below.

Remark 2 The three cases (3), (5), and (4) are in fact three separate statements in which we set our hopes to converge fast on the first, second or both layers, respectively. We have combined them into one "mega-theorem" since the proofs share many steps.

Let us now compare our result with the previous work cited above. The case $\eta_z = 0$ (that is, training the first layer only) automatically implies the second condition in (3), and the only remaining assumptions are

$$m \le \min\left(\tilde{o}(S), \ \tilde{o}\left(nS^{1/4}\right)\right)$$
 (7)

which is an improvement on (1).

The condition $\frac{\eta_z}{\eta_w} \leq \tilde{o}\left(\frac{S}{m}\right)$ in (3) is automatically satisfied when $\eta_z = \eta_w$ (as $m \leq \tilde{o}(S)$), i.e. for the "ordinary" gradient descent. Thus our result

⁴We do not include m as we are assuming $m \leq n^{O(1)}$ throughout the paper.

improves on [DZPS19]. As an easy calculation shows, re-normalizing from LeCun's initialization to the standard one gives the ratio $\frac{\eta_z}{\eta_w} = \frac{S}{m}$. Thus the condition (5) holds whenever $m \leq \min\left(\tilde{o}(S), \ \tilde{o}\left(n^{1/2}S^{1/2}\right)\right)$. This gives an improvement on (2).

Finally, note that when $m \leq \tilde{o}\left(n^{1/3}S^{2/3}\right)$, the ranges for the relative learning rate $\frac{\eta_z}{\eta_w}$ in (3) and (4) overlap. Thus, under the assumptions

$$m \leq \min \left(\widetilde{o}(S), \ \widetilde{o}\left(n^{1/3}S^{2/3}\right), \ \widetilde{o}\left(nS^{1/4}\right) \right),$$

convergence of the gradient descent is guaranteed *regardless* of the relative learning rate; we only have to make sure the rates are individually small enough. In particular, this holds regardless of the choice of normalization at the initialization.

Our improvements are admittedly not dramatic but our main motivation was to develop new techniques that might turn out useful for pushing the bound further or in other similar situations. This is what we consider our main technical contributions.

- Most previous work used, explicitly or implicitly, bounds on the spectral norm ||X||, where X is the data matrix. We will heavily exploit that when X is random, analogous bounds also hold for all its sub-matrices X^J , where J is a subset of data instances (see (20)). We will also employ a dual lower bound on $\sigma_{\min}\left((X^J)^{\top}\right)$ for all sufficiently large J (see (21) and compare with [HY20, Assumption 2.2]).
- We will also need to extend the lower bound $\lambda_{\min}(H_0) \geq \Omega(S)$ to the uniform lower bound $\lambda_{\min}(H_0|_{\Gamma}) \geq \Omega(S)$, where $\Gamma \subseteq [S]$ is a sufficiently large but otherwise arbitrary set of neurons and $H_0|_{\Gamma}$ is obtained from H_0 by disregarding neurons not in Γ . It turned out surprisingly difficult and required a version of the Lipschitz concentration inequality (see the second part of Appendix C).
- For the first two parts (3), (4), we also need to do finer analysis of the changes the activation matrix A incurs during training. Informally speaking, there are two different ways the network may try to escape the NTK regime: by accumulating a large overall number of changes or a relatively large number concentrated on a small set of inputs J. Attempts of the first kind will be prevented by the assumption $m \leq$

 $\tilde{o}\left(nS^{1/4}\right)$ while the second kind will be taken care of by $m \leq \tilde{o}(S)$. See Appendix D for details.

• Since in case of RelU activations the gradient is discontinuous, we have to take care of the unpleasant situation when A and hence the Jacobian J_F of the feature map significantly change during one step. The previous work (see e.g. [OS20, Appendix G]) did it by comparing both Jacobians to the Jacobian at initialization but it does not work any longer with our relaxed assumptions. Instead, we adopt to our purposes the beautiful invariant from [ACH18, DHL18] concerning the behavior of weights associated to an individual neuron. See Appendix F for details.

In Section 4, we will also report several simple experiments with synthetic data. We defer further discussion to that section; for now let us just remark that one strong conclusion we draw from these experiments is that there should be very significant room for further theoretical improvements, possibly (almost) all the way up to the representability barrier $m \sim nS$.

1.4. Related research

Our attempt at the uniformization has inevitably left behind some very interesting and somewhat related work. Very briefly:

Deeper networks. Several papers (notably [ALS19, ZG19]) are specifically devoted to the NTK regime in the context of deeper networks. When scaled down to shallow networks, however, the results do not seem to be as good as those obtained with methods specifically tailored to that situation.

Classification problems. Quite a bit of similar and generally stronger results have been obtained for classification problems, both binary and multiclassification. See e.g. [BGMS18, ZCZG20, JT20, CCZG21, FCB22, LZ22]. While some methods and approaches are shared with the regression problems, many techniques seem to be specific to classification.

Smooth activation. This case tends to be easier than the case of RelU, and many technical difficulties disappear. For strong results in that direction not covered above see e.g. [HY20, SRP⁺21, BAM22].

2. Preliminaries and the main result

2.1. Notation

We let $[n] \stackrel{\text{def}}{=} \{1, 2, \dots, n\}$ and let $\binom{V}{n}$ be the collection of all n-element subsets of V. For a real number t, $\lfloor t \rfloor$ is its integer part and $\{t\} \stackrel{\text{def}}{=} t - \lfloor t \rfloor$ is its fractional part. All vectors and matrices are real, and all vectors are column vectors unless otherwise noted. Matrices are generally denoted by upper case Latin letters like A, B, H, W, X and vectors are typically denoted by lower case letters. We let $\mathbb{R}^{n \times m}$ be the space of $n \times m$ matrices. For $M \in \mathbb{R}^{n \times m}$ and $j \in [m]$, we let M^j denote its jth column vector, and for $i \in [n], M_i$ is its ith row vector. More generally, for $J \subseteq [m], M^J$ is the corresponding $n \times |J|$ matrix, and similarly for M_I . For $z \in \mathbb{R}^S$, $\operatorname{diag}(z) \in \mathbb{R}^{S \times S}$ is the corresponding diagonal matrix.

The symbol \circ stands for the Halliard (entryism) product of vectors/matrices of the same size. The symbol * will denote the *column-wise* Khatri-Rao product: for $A \in \mathbb{R}^{S \times m}$ and $X \in \mathbb{R}^{n \times m}$, $(A * X) \in \mathbb{R}^{Sn \times m}$ and $(A * X)_{(\nu i),j} \stackrel{\text{def}}{=} A_{\nu j} X_{ij}$. We will denote the standard Euclidean norm by ||x||. The unit sphere $S^{n-1} \subseteq \mathbb{R}^n$ is given by $S^{n-1} \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n \mid ||x|| = 1\}$; more generally, $r \cdot S^{n-1} \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n \mid ||x|| = r\}$. The spectral norm ||M|| of $M \in \mathbb{R}^{n \times m}$ is $||M|| \stackrel{\text{def}}{=} \max_{\xi \in S^{m-1}} ||M\xi||$, its minimum singular value is $\sigma_{\min}(M) \stackrel{\text{def}}{=} \min_{\xi \in S^{m-1}} ||M\xi||$ and the Frobenius norm is $||M||_{\mathbb{F}} \stackrel{\text{def}}{=} \left(\sum_{i \in [n]} \sum_{j \in [m]} M_{ij}\right)^{1/2}$. We let $||x||_{\infty}$, $||M||_{\infty}$ denote the maximum absolute value of an entry.

When H is a PSD (positive semi-definite) matrix, we will usually write $\lambda_{\min}(H)$ for $\sigma_{\min}(H)$; thus, for any real matrix M, $\sigma_{\min}(M) = \lambda_{\min}(M^\top M)^{1/2}$. For two symmetric matrices M, N of the same size, $M \succeq N$ means that M-N is PSD.

We let $\mathbf{E}[*]$ be the expectation of a real random variable, $\mathrm{Var}(*)$ be the variance of a (one-dimensional) random variable and $\mathbf{P}[E]$ be the probability of an event E. $\mathbf{1}(E)$ is the characteristic function of E. When we randomize over a part of the sample space, this will be indicated by the corresponding subscript, like $\mathbf{E}_{W_0}[*]$ or $\mathbf{P}_z[*]$. Let $||\zeta||_{\psi_2}$, $||\zeta||_{\psi_1}$ be the sub-gaussian and the sub-exponential norms of a one-dimensional ζ defined as

$$||\zeta||_{\psi_p} \stackrel{\text{def}}{=} \inf \left\{ s > 0 \, \middle| \, \mathbf{E} \left[e^{(\zeta/s)^p} \right] \le 2 \right\}.$$

2.2. Useful facts

Shur's theorem and inequalities. The Halliard product $M \circ N$ of two PSD $m \times m$ matrices is again PSD and satisfies

$$\begin{split} ||M \circ N|| & \leq & \max_{j \in [m]} |M_{jj}| \cdot ||N||; \\ \lambda_{\min}(M \circ N) & \geq & \min_{j \in [m]} |M_{jj}| \cdot \lambda_{\min}(N). \end{split}$$

If $A \in \mathbb{R}^{S \times m}$ and $X \in \mathbb{R}^{n \times m}$ then

$$(A*X)^{\top}(A*X) = (A^{\top}A) \circ (X^{\top}X).$$

In particular, applying Shur's inequalities with $M \mapsto A^{\top}A$, $N \mapsto X^{\top}X$, we get

$$||A * X|| \leq \max_{j \in [m]} ||A^{j}|| \cdot ||X||;$$

$$\sigma_{\min}(A * X) \geq \min_{j \in [m]} ||A^{j}|| \cdot \sigma_{\min}(X).$$
(8)

Norm bounds. For $z \in \mathbb{R}^S$, the norm $||(\operatorname{diag}(z)A) * X||$ can be bound in two different (dual) ways:

$$||(\operatorname{diag}(z)A) * X|| \leq ||z||_{\infty} \cdot ||A * X||;$$

$$||(\operatorname{diag}(z)A) * X|| \leq ||z|| \cdot \max_{\nu \in S} ||X \operatorname{diag}(A_{\nu})|| \leq ||z|| \cdot ||A||_{\infty} \cdot ||X||.(10)$$

Another handy inequality is

$$||MN||_{F} \le ||M||_{F} \cdot ||N||.$$
 (11)

 ϵ -nets [Ver12, Chapter 5.2.2]. An ϵ -net in S^{n-1} is any subset $\mathcal{N} \subseteq S^{n-1}$ such that $\forall v \in S^{n-1} \exists \xi \in \mathcal{N}(||v - \xi|| \le \epsilon)$.

Fact 2.1 ([Ver12, Lemma 5.2]) For any ϵ, n there exists an ϵ -net \mathcal{N} in S^{n-1} with $|\mathcal{N}| \leq \left(1 + \frac{2}{\epsilon}\right)^n$.

Fact 2.2 (similar to [Ver12, Lemma 5.3]) Let $M \in \mathbb{R}^{k \times n}$, $\sigma > 0$ and \mathcal{N} be an $\frac{\sigma}{2||M||}$ -net in S^{n-1} such that $\forall \xi \in \mathcal{N}(||M\xi|| \geq \sigma)$. Then $\sigma_{\min}(M) \geq \frac{\sigma}{2}$.

2.3. The model

Let $n, m, S \ge 1$. We are given m data instances $X^1, \ldots, X^m \in \mathbb{R}^n$ arranged as a matrix⁵ $X \in \mathbb{R}^{n \times m}$ and a label vector $y \in \mathbb{R}^m$. A (shallow) neural network is given by a pair $\theta = (W, z)$, where $W \in \mathbb{R}^{S \times n}$ and $z \in \mathbb{R}^S$. It computes as

$$F(W) \stackrel{\text{def}}{=} \sigma(WX) \in \mathbb{R}^{S \times m};$$

 $f(\theta) \stackrel{\text{def}}{=} F(W)^{\top} z \in \mathbb{R}^{m},$

where, as usual, the RelU activation function $\sigma(x) \stackrel{\text{def}}{=} \max(x,0)$ is applied to WX entryism. We define the $error\ vector$ as

$$e(\theta) \stackrel{\text{def}}{=} f(\theta) - y \in \mathbb{R}^m$$

and the (quadratic) loss function as

$$\ell(\theta) \stackrel{\text{def}}{=} \frac{1}{2} ||e(\theta)||^2.$$

Let us call $\theta = (W, z)$ regular if WX does not contain zero entries and singular otherwise. All functions just defined are polynomial in a neighborhood of any regular θ , hence the gradient $\nabla \ell(\theta)$ is well-defined for regular θ . We split it as $\nabla \ell(\theta) = (\nabla^w \ell(\theta), \nabla^z \ell(\theta))$ in the obvious way: $\nabla^z \ell(\theta) \in \mathbb{R}^S$ and $\nabla^w \ell(\theta)$ will be treated, depending on the context, either as a long vector in \mathbb{R}^{nS} or as a matrix in $\mathbb{R}^{S \times n}$.

Assume now that we are also given initial values $\theta_0 = (W_0, z_0)$ and the learning rates $\eta_w, \eta_z \geq 0$. We define the gradient descent $\theta_\tau = (W_\tau, z_\tau)$ ($\tau \geq 1$ an integer) as

$$W_{\tau} \stackrel{\text{def}}{=} W_{\tau-1} - \eta_w \nabla^w \ell(\theta_{\tau-1});$$

$$z_{\tau} \stackrel{\text{def}}{=} z_{\tau-1} - \eta_z \nabla^z \ell(\theta_{\tau-1}).$$
(12)

There is a small caveat here since this definition makes sense only when all points θ_{τ} are regular. Fortunately, it is easy to take care of with a simple measure-theoretical argument that we will present in Appendix A.

⁵In many works on the subject, data points are arranged as *rows* of the data matrix. We have found that, while also imperfect, the opposite convention leads to a cleaner and more natural notation.

2.4. Data, initialization and the main result

We pick the data X^1, \ldots, X^m uniformly (under the Haar measure) and independently from S^{n-1} . All entries in W_0 are chosen from $\mathcal{N}(0,1)$, independently of each other. We prefer to allow for a slightly better flexibility in the choice of z_0 , primarily since the literature alternates between Gaussian and Rademacher initializations.

Recall that the notation \widetilde{O} , $\widetilde{\Omega}$, \widetilde{o} , $\widetilde{\omega}$ hides factors that are poly-logarithmic in $\log(nS)$; see Remark 1 for more details. We require that z_0 is independent of W_0 , and that the entries in z_0 are S i.i.d. copies of a one-dimensional distribution ζ such that:

$$\mathbf{E}[\zeta] \stackrel{\text{def}}{=} 0; ||\zeta||_{\psi_2} \leq \tilde{O}(1);$$
 (13)

$$\operatorname{Var}(\zeta) \geq \widetilde{\Omega}(1).$$
 (14)

We are now ready to state our main result.

Theorem 2.3 Let $m \ge 1$ and $n, S \ge \widetilde{\omega}(1)$ be parameters such that $m \le n^{O(1)}$ and $m \le \widetilde{o}(S)$. Let $\eta_w, \eta_z \ge 0$, $\eta_w + \eta_z > 0$ be such that

$$\eta_w \le \min\left(\widetilde{o}\left(\frac{1}{S^2}\right), \widetilde{o}\left(\frac{n}{mS^2}\right)\right), \quad \eta_z \le \widetilde{o}\left(\frac{1}{mS^2}\right)$$
(15)

and that at least one of the following three conditions holds:

$$m \le \tilde{o}\left(nS^{1/4}\right), \quad \frac{\eta_z}{\eta_w} \le \tilde{o}\left(\frac{S}{m}\right),$$
 (16)

$$m \le \tilde{o}\left(nS^{1/4}\right), \quad \frac{\eta_z}{\eta_w} \ge \tilde{\omega}\left(\frac{m^2}{nS}\right),$$
 (17)

$$\frac{\eta_z}{\eta_w} \ge \widetilde{\omega} \left(1 + \frac{m}{n} \right). \tag{18}$$

Then with probability $\geq 1 - (nS)^{-\omega(1)}$ w.r.t. the choice of X, W_0, z_0 as above the following holds. For almost all (in the Lebesque sense) $y \in \mathbb{R}^m$ the gradient descent (12) avoids singular points and if additionally

$$||y|| \le \tilde{O}(m^{1/2}S^{1/2})$$

then it converges to a global minimum.

Before embarking on the proof of this theorem, let us make a few remarks, in addition to those that were already made in Section 1.3.

Remark 3 The choice of labels is completely adversarial (they are only required to be of "right" order), and the adversary is also allowed to see the initialization, not only the data. This aligns well with the prominent experiments in [ZBH⁺21] strongly suggesting that the choice of labels should *not* be a defining factor in convergence.

Remark 4 As in all previous work, the error probability can be decreased to exponential by examining the proofs in Appendices B, C.

Remark 5 Since $m \leq S$, the bound in (18) implies the second bound in (17). Thus, we always have either $m \leq \tilde{o}\left(nS^{1/4}\right)$ or $\frac{\eta_z}{\eta_w} \geq \tilde{\omega}\left(\frac{m^2}{nS}\right)$ and, in fact, the refinement (18) of this dichotomy will not be needed until Appendix E.

3. Proof of the main result

As noted before, we defer the simple proof that the gradient descent avoids singular points a.e. to Appendix A and for now simply assume that $y \in \mathbb{R}^m$ is chosen in such a way that this is true, and such that $||y|| \leq \tilde{O}(m^{1/2}S^{1/2})$.

We begin the proof of Theorem 2.3 with fixing some useful notation and reminding some basic facts.

First of all, it will be convenient to extend the trajectory (12) to continuous time in the piecewise linear manner. That is, for any $t \ge 0$ we set

$$\begin{array}{ccc} W_t & \stackrel{\mathrm{def}}{=} & W_{\lfloor t \rfloor} - \eta_w \{t\} \nabla^w \ell(\theta_{\lfloor t \rfloor}); \\ \\ z_t & \stackrel{\mathrm{def}}{=} & z_{\lfloor t \rfloor} - \eta_z \{t\} \nabla^z \ell(\theta_{\lfloor t \rfloor}). \end{array}$$

Since θ_{τ} is regular for all $\tau \in \mathbb{N}$, the interval $[\tau, \tau + 1]$ may contain only finitely many t for which θ_t is singular. Hence those t can and will be ignored in our estimates based on integration.

We abbreviate $F(W_t)$, $f(\theta_t)$, $e(\theta_t)$ etc. to F_t , f_t , e_t respectively. Whenever θ_t is regular, we let

$$A_t \stackrel{\text{def}}{=} \sigma'(W_t X)$$
 (the activation matrix).

If θ_t is singular, we let $A_t \stackrel{\text{def}}{=} A_{t+0}$ which makes all our constructs upper semi-continuous.

Next, let

$$B_t \stackrel{\text{def}}{=} \operatorname{diag}(z_t) A_t$$
 (the weighted activation matrix),

then $(B_t * X, F_t)$ is the transposed Jacobian $J_f(\theta_t)$ and hence

$$\nabla^{w} \ell(\theta_{t}) = (B_{t} * X) e_{t};$$

$$\nabla^{z} \ell(\theta_{t}) = F_{t} e_{t}.$$

We also let

$$H_t \stackrel{\text{def}}{=} (B_t * X)^{\top} (B_t * X) = (X^{\top} X) \circ (B_t^{\top} B_t),$$

$$G_t \stackrel{\text{def}}{=} F_t^{\top} F_t;$$

these are (for integer t) the two components of the NTK matrix.

We will frame the rest of the proof according to the paradigm known in combinatorics as quasirandomness (see [CGW89] for graphs and [CR21] for arbitrary combinatorial objects). The idea is to split a logically elaborated argument involving random objects into two totally independent parts. At the first stage we accumulate a list of so-called "quasirandom properties", that is completely deterministic facts that hold for our random objects with overwhelming probability. At the second stage, the argument proceeds completely deterministically, on the base of these properties only. This allows us to avoid unwanted serious complications, or even sheer mistakes, caused by the fact that quantifiers and randomization do not get along well in one argument (see e.g. [ABR+21] where we adopted this approach). In the context of deep learning, having a list of properties sufficient for convergence in a convenient form might have another benefit: it may help to facilitate a discussion of what exactly differentiates the data and initialization occurring in practical problems from random synthetic data.

We would like to finish this brief discourse with acknowledging that this approach is well known in the area under various names like "assumptions", "meta-theorems" (see e.g. [OS20]) etc. The main difference is perhaps that we do it a bit more systematically and on a larger scale.

3.1. Quasirandom properties

We claim that with our choice of X, W_0, z_0 (Section 2.4), the following are satisfied with probability $\geq 1 - (nS)^{-\omega(1)}$. The proofs are deferred to Appendix B (straightforward proofs) and Appendix C (not so straightforward).

3.1.1. Properties of data

Almost orthogonality.

$$\max_{j \neq j' \in [m]} \langle X^j, X^{j'} \rangle \le \widetilde{O}(n^{1/2}). \tag{19}$$

In the following property, as well as (28), the constants assumed in the \tilde{O} notation do not depend on k and j, R respectively.

Uniform bounds on norm.

$$\forall k \in [m] \left(\max_{J \in \binom{[m]}{k}} ||X^J|| \le \tilde{O} \left(1 + \frac{k^{1/2}}{n^{1/2}} \right) \right).$$
 (20)

In the following property, as well as in (24), the bounds in the existential quantifiers are completely explicit, see the proofs in Appendix B.

Dual bound on the minimum singular value.

$$\exists n^* \le \widetilde{O}(n) \left(\min_{J \in \binom{[m]}{n^*}} \sigma_{\min} \left((X^J)^\top \right) \ge \frac{n}{m} \right). \tag{21}$$

3.1.2. Properties of initialization

Rows of W_0 are large.

$$\min_{\nu \in [S]} ||(W_0)_{\nu}|| \ge \widetilde{\Omega}(n^{1/2}). \tag{22}$$

The following two properties are entirely obvious under either Rademacher or Gaussian initialization of z_0 .

Entries of θ_0 are small.

$$||W_0||_{\infty}, ||z_0||_{\infty} \le \tilde{O}(1).$$
 (23)

 z_0 contains sufficiently many large entries.

$$\exists \zeta_0 \ge \widetilde{\Omega}(1) \left(| \left\{ \nu \in [S] \, | \, |(z_0)_{\nu}| \ge \zeta_0 \right\} | \ge \widetilde{\Omega}(S) \right). \tag{24}$$

3.1.3. Smooth properties

Regularity.

$$\theta_0$$
 is regular. (25)

Data is almost orthogonal to initialization.

$$||W_0X||_{\infty} \le \tilde{O}(1). \tag{26}$$

Right order of the output at initialization.

$$||f_0||_{\infty} \le \widetilde{O}(S^{1/2}). \tag{27}$$

"Good behavior" for any data instance.

$$\forall j \in [m] \forall R \ge 0 \left(|\{ \nu \in [S] \mid |(W_0 X)_{\nu j}| \le R \} | \le \tilde{O}(SR + 1) \right). \tag{28}$$

3.1.4. NTK properties at initialization

NTK: second layer.

$$\lambda_{\min}(G_0) \ge \Omega(S). \tag{29}$$

In order to formulate our last (and the most difficult) property, let us fix $\zeta_0 \geq \widetilde{\Omega}(1)$ as in (24), and set

$$\Gamma_0 \stackrel{\text{def}}{=} \{ \nu \in [S] \mid |(z_0)_{\nu}| \geq \zeta_0 \};$$

thus, $|\Gamma_0| \geq \widetilde{\Omega}(S)$.

NTK: first layer.

$$m \leq \widetilde{o}(nS^{1/2}) \Longrightarrow \exists S^* \geq \widetilde{\Omega}\left(\frac{n^2S}{n^2 + m}\right)$$

$$\left(\min_{\Gamma \in \binom{\Gamma_0}{|\Gamma_0| - S^*}} \lambda_{\min}\left(((A_0)_{\Gamma} * X)^{\top}((A_0)_{\Gamma} * X)\right) \geq \Omega(S)\right).$$
(30)

Remark 6 We have not been able to remove the annoying term m in the denominator of the bound on S^* ; we believe it can be done with a better analysis than ours.

3.2. Deterministic part

According to the plan outlined above, from now on we are assuming that X, θ_0 are chosen arbitrarily in such a way that all properties in Section 3.1 are satisfied. The labels $y \in \mathbb{R}^m$ are also arbitrary, as long as $||y|| \leq \tilde{O}(m^{1/2}S^{1/2})$ and the gradient descent avoids singular points. The learning rates η_w, η_z are assumed to satisfy (15) and one of (16), (17), (18). All the way until Appendix E, however, we will only need that either $m \leq \tilde{o}(nS^{1/4})$ or $\frac{\eta_z}{\eta_w} \geq \tilde{\omega}\left(\frac{m^2}{nS}\right)$.

Remark 7 Before we begin, let us make one technical remark. The definition (6) of the assumption $m \leq \tilde{o}(nS^{1/4})$ can be equivalently re-written as

$$\lim_{k \to \infty} \frac{\log \left(\frac{n^{(k)} (S^{(k)})^{1/4}}{m^{(k)}} \right)}{\log \log \left(n^{(k)} S^{(k)} \right)} = \infty.$$
(31)

Likewise, $m \geq \widetilde{\Omega}(nS^{1/4})$ means that the quantity in (31) is bounded. Since from every infinite sequence we can extract an infinite subsequence for which one of the two is true, we can freely assume w.l.o.g. that either $m \leq \widetilde{o}(nS^{1/4})$ or $m \geq \widetilde{\Omega}(nS^{1/4})$ holds. The same dichotomy applies to $\frac{\eta_z}{\eta_w} \geq \widetilde{\omega}\left(\frac{m^2}{nS}\right)$ vs. $\frac{\eta_z}{\eta_w} \leq \widetilde{O}\left(\frac{m^2}{nS}\right)$.

We are going to prove the following upper bounds on the distance travelled

in the parameter space:

$$m \le \tilde{o}\left(nS^{1/4}\right) \implies \sup_{T>0} ||W_T - W_0||_{F} < \tilde{O}(m^{1/2});$$
 (32)

$$m \le \tilde{o}\left(nS^{1/4}\right) \implies \sup_{T>0} ||z_T - z_0||_F < \tilde{O}\left(\frac{\eta_z^{1/2}}{\eta_w^{1/2}}m^{1/2}\right);$$
 (33)

$$\frac{\eta_z}{\eta_w} \ge \widetilde{\omega} \left(\frac{m^2}{nS} \right) \implies \sup_{T>0} ||W_T - W_0||_{\mathcal{F}} < \widetilde{O} \left(\frac{\eta_w^{1/2}}{\eta_z^{1/2}} m^{1/2} \right); \tag{34}$$

$$\frac{\eta_z}{\eta_w} \ge \widetilde{\omega} \left(\frac{m^2}{nS} \right) \implies \sup_{T>0} ||z_T - z_0||_{\mathcal{F}} < \widetilde{O} \left(m^{1/2} \right). \tag{35}$$

With a slight abuse of notation⁶, we prove (32)-(35) by induction on T; that is, we fix T > 0 and assume that these bounds hold for all t < T. We then claim that for all t < T we also have

$$m \le \tilde{o}\left(nS^{1/4}\right) \implies \lambda_{\min}(H_t) \ge \tilde{\Omega}(S);$$
 (36)

$$\frac{\eta_z}{\eta_w} \ge \widetilde{\omega} \left(\frac{m^2}{nS} \right) \implies \lambda_{\min}(G_t) \ge \widetilde{\Omega}(S).$$
(37)

Let us prove (37) since it is easier. Given (29), it is sufficient to show that

$$||F_t - F_0|| \le \tilde{o}(S^{1/2}).$$
 (38)

For that we perform the calculation

$$||F_{t} - F_{0}|| \leq ||F_{t} - F_{0}||_{F} = ||\sigma(W_{t}X) - \sigma(W_{0}X)||_{F}$$

$$\leq ||(W_{t} - W_{0})X||_{F} \leq ||W_{t} - W_{0}||_{F} \cdot ||X||,$$
(39)

where the third inequality hold since σ is 1-Lipschitz. We now check that (39) indeed implies (38).

First,

$$||W_t - W_0||_{\mathcal{F}} \stackrel{(34)}{\leq} \widetilde{O}\left(\frac{\eta_w^{1/2}}{\eta_z^{1/2}} m^{1/2}\right) \leq \widetilde{o}\left(\frac{n^{1/2} S^{1/2}}{m^{1/2}}\right),\tag{40}$$

where for the second inequality we used the assumption in (37).

 $^{^6\}mathrm{To}$ be completely impeccable, we should first fix the constants assumed in the right-hand sides.

Next, as noted in Remark 7, we can assume w.l.o.g. that either $m \ge \tilde{\Omega}\left(nS^{1/4}\right)$ or $m \le \tilde{o}\left(nS^{1/4}\right)$.

In the first case, the bound (20) for k = m simplifies to $||X|| \leq \tilde{O}\left(\frac{m^{1/2}}{n^{1/2}}\right)$. This, along with (40), implies (38).

If, on the other hand, $m \leq \tilde{o}(nS^{1/4})$, we can also employ (32) and refine (40) to

$$||W_t - W_0||_{\mathrm{F}} \le \min\left(\widetilde{O}(m^{1/2}), \ \widetilde{o}\left(\frac{n^{1/2}S^{1/2}}{m^{1/2}}\right)\right) \le \min\left(\widetilde{o}(S^{1/2}), \ \widetilde{o}\left(\frac{n^{1/2}S^{1/2}}{m^{1/2}}\right)\right).$$

Multiplying this with the full version $||X|| \leq \tilde{O}\left(1 + \frac{m^{1/2}}{n^{1/2}}\right)$ of (20) again gives us (38). Thus, the proof of (38) and hence also of (37) is completed.

The proof of (36) is more elaborate and is deferred to Appendix D.

The dynamics of the gradient descent in the error space is ruled by the differential equation

$$\dot{\ell}_t = \langle \nabla \ell(\theta_t), \dot{\theta}_t \rangle,$$

where

$$\dot{\theta}_t = -\left(\eta_w \nabla^w \ell(\theta_{\lfloor t \rfloor}), \eta_z \nabla^z \ell(\theta_{\lfloor t \rfloor})\right).$$

Hence

$$\frac{d}{dt}||e_t|| = \frac{\frac{d}{dt}||e_t||^2}{2||e_t||} = \frac{\dot{\ell}_t}{||e_t||} = \frac{\langle \nabla \ell(\theta_t), \dot{\theta}_t \rangle}{||e_t||} = \frac{\langle \nabla \ell(\theta_t), \dot{\theta}_{\lfloor t \rfloor} \rangle}{||e_t||}.$$

Let

$$\delta_t \stackrel{\text{def}}{=} \nabla \ell(\theta_t) - \nabla \ell(\theta_{\lfloor t \rfloor})$$

(this measures how much the gradient may change during one step). We now have

$$\langle \nabla \ell(\theta_t), \dot{\theta}_{|t|} \rangle = \langle \nabla \ell(\theta_{|t|}), \dot{\theta}_{|t|} \rangle + \langle \delta_t, \dot{\theta}_{|t|} \rangle = -\eta_w ||\nabla^w \ell(\theta_{|t|})||^2 - \eta_z ||\nabla^z \ell(\theta_{|t|})||^2 + \langle \delta_t, \dot{\theta}_{|t|} \rangle.$$

In Appendices E, F we will prove that

$$\left| \langle \delta_t, \dot{\theta}_{\lfloor t \rfloor} \rangle \right| < \frac{1}{2} \langle \nabla \ell(\theta_{\lfloor t \rfloor}), \dot{\theta}_{\lfloor t \rfloor} \rangle \tag{41}$$

(the proof for the z-part is pretty straightforward, given the absolute bounds (15), but will require a new idea for the W-part due to the discontinuity of $\nabla^w \ell$).

Once we have that, we know that $||e_t||$ is decreasing or, more specifically,

$$\frac{d}{dt}||e_t|| \le \frac{\langle \nabla \ell(\theta_{\lfloor t \rfloor}), \dot{\theta}_{\lfloor t \rfloor} \rangle}{2||e_t||} \le \frac{\langle \nabla \ell(\theta_{\lfloor t \rfloor}), \dot{\theta}_{\lfloor t \rfloor} \rangle}{2||e_{|t|}||},\tag{42}$$

where the last inequality holds since $||e_t||$ is decreasing.

Note that $||f_0|| \leq \tilde{O}(m^{1/2}S^{1/2})$ by (27) and, therefore, $||e_0|| \leq \tilde{O}(m^{1/2}S^{1/2})$ due to our assumption on y. Integrating (42) from 0 to T gives us

$$\int_0^T \frac{\langle \nabla \ell(\theta_{\lfloor t \rfloor}), -\dot{\theta}_{\lfloor t \rfloor} \rangle}{||e_{\lfloor t \rfloor}||} dt \le 2 \int_0^T \left(-\frac{d}{dt} ||e_t|| \right) dt = 2(||e_0|| - ||e_T||) \le \widetilde{O}(m^{1/2} S^{1/2}),$$

where the equality holds since $||e_t||$ is decreasing.

Using Cauchy-Schwartz, we now estimate

$$||W_{T} - W_{0}||_{F} \leq \eta_{w} \int_{0}^{T} ||\nabla^{w}\ell(\theta_{\lfloor t \rfloor})|| dt$$

$$\leq \eta_{w} \left(\int_{0}^{T} \frac{||\nabla^{w}\ell(\theta_{\lfloor t \rfloor})||^{2}}{||e_{\lfloor t \rfloor}||} \right)^{1/2} \cdot \left(\int_{0}^{T} ||e_{\lfloor t \rfloor}|| dt \right)^{1/2}$$

$$\leq \eta_{w} \left(\int_{0}^{T} \frac{\langle \nabla \ell(\theta_{\lfloor t \rfloor}), -\dot{\theta}_{t} \rangle}{\eta_{w} ||e_{\lfloor t \rfloor}||} dt \right)^{1/2} \cdot \left(\int_{0}^{T} ||e_{\lfloor t \rfloor}|| dt \right)^{1/2}$$

$$\leq \tilde{O} \left(\eta_{w}^{1/2} m^{1/4} S^{1/4} \cdot \left(\int_{0}^{T} ||e_{\lfloor t \rfloor}|| dt \right)^{1/2} \right)$$

$$(43)$$

and, likewise,

$$||z_T - z_0|| \le \widetilde{O}\left(\eta_z^{1/2} m^{1/4} S^{1/4} \cdot \left(\int_0^T ||e_{\lfloor t\rfloor}||dt\right)^{1/2}\right).$$
 (44)

In order to estimate the last remaining term $\int_0^T ||e_{\lfloor t\rfloor}|| dt$, let us first assume that $m \leq \tilde{o}(nS^{1/4})$. In that case we have (36) and then we can continue (42) as follows:

$$\frac{d}{dt}||e_t|| \le -\widetilde{\Omega}\left(\frac{\langle \nabla \ell(\theta_{\lfloor t \rfloor}), -\dot{\theta}_{\lfloor t \rfloor} \rangle}{||e_{\lfloor t \rfloor}||}\right) \le -\widetilde{\Omega}\left(\eta_w \frac{||\nabla^w \ell(\theta_{\lfloor t \rfloor})||^2}{||e_{\lfloor t \rfloor}||}\right)
\le -\widetilde{\Omega}(\eta_w \lambda_{\min}(H_{\lfloor t \rfloor})||e_{\lfloor t \rfloor}||) \le -\widetilde{\Omega}(\eta_w S||e_{\lfloor t \rfloor}||).$$

This (sub)differential equation solves to

$$m \leq \tilde{o}(nS^{1/4}) \Longrightarrow ||e_t|| \leq ||e_0|| \cdot \exp(-\tilde{\Omega}(\eta_w tS))$$

$$\leq \tilde{O}\left(m^{1/2}S^{1/2}\exp(-\tilde{\Omega}(\eta_w tS))\right). \tag{45}$$

Integrating from 0 to T gives us

$$m \le \tilde{o}(nS^{1/4}) \Longrightarrow \int_0^T ||e_{\lfloor t\rfloor}|| dt \le \tilde{O}\left(\frac{m^{1/2}}{\eta_w S^{1/2}}\right)$$
 (46)

and, likewise (but applying (37) this time),

$$\frac{\eta_z}{\eta_w} \ge \tilde{\omega} \left(\frac{m^2}{nS} \right) \implies ||e_t|| \le \tilde{O} \left(m^{1/2} S^{1/2} \exp(-\tilde{\Omega}(\eta_z t S)) \right);$$
 (47)

$$\frac{\eta_z}{\eta_w} \ge \widetilde{\omega} \left(\frac{m^2}{nS} \right) \implies \int_0^T ||e_{\lfloor t \rfloor}|| dt \le \widetilde{O} \left(\frac{m^{1/2}}{\eta_z S^{1/2}} \right). \tag{48}$$

Plugging (46), (48) into (43), (44) gives us all four inequalities (32)-(35). This completes their proof by joint induction on t.

In particular, in the first case (16) we have (45), in the third case (18) we have (47) (and in the second case (17) we have both). This "completes" the proof of Theorem 2.3, with exponential speed of convergence.

4. Experiments

Since our paper does not claim any direct relevance to practical data, we have confined our experiments to synthetic data generated precisely as in Section 2.4, with $\zeta \in_R \{\pm 1\}$. In most experiments, the vector y of label data was chosen from $\mathcal{N}(0,S)^m$. We also performed a few experiments with other natural choices:

low spectrum $e_0 = f_0 - y$ is the eigenvector of the matrix H_0 corresponding to the smallest eigenvalue and scaled in such a way that $||e_0|| = m^{1/2}S^{1/2}$.

high spectrum $e_0 = n^{1/2} S^{1/2}(X^\top X^1)$. This is a "natural choice" for which $||H_0 e_0||$ is large.

⁷The quotation marks indicate the fact that all serious work is delegated to Appendix.

local
$$e_0 = \begin{pmatrix} m^{1/2}S^{1/2} \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$
. The initial error is completely concentrated on the

first data instance (and is huge).

We considered only the case $\eta_z=0$, i.e., training the first layer. All our experiments were equipped with the "safety valve": they terminated once the situation $||e_{\tau}|| > ||e_{\tau-1}||$ was encountered. A computation was declared successful and stopped once $||e_{\tau}|| < 10^{-3}$. We remark that due to our choice of normalization, the initial value was roughly of order 10^2-10^4 so our stopping condition corresponds to the loss reduction by a factor of $10^{10}-10^{14}$.

Like any other theoretical paper, ours is plugged with asymptotic notation (we would like to note in passing that attempts to replace it with explicit constants are usually not very instructive or even readable). The original purpose of our experiments was rather limited: verify that our theoretical findings are still reasonably relevant for relatively small values of parameters. What they have showed instead is that the crucial quantity $\lambda_{\min}(H_t)$ still behaves smoothly and nicely well beyond what could be optimistically called the "NTK regime", almost all the way up to the representability barrier $m \sim nS$. In simpler words, there should be a significant room for improving theoretical results on convergence with methods yet to be developed. In retrospect, this is not surprising at all: the bounds like the first inequality in (43) are hopelessly pessimistic. But it is nice to have an experimental encouragement.

Let us now be more specific. In all experiments the dimension was set to n := 100 and, unless otherwise noted, the learning rate was $\eta_w := 10^{-3}$. We did not try to optimize on the latter since this is not the main focus of the paper. But several sporadic experiments showed that for smaller values of S, m it can be significantly increased.

We kept track of the following control quantities (here T is the stopping time):

• $\kappa_H \stackrel{\text{def}}{=} \frac{\lambda_{\min}(H_T)}{\lambda_{\min}(H_0)}$. Computing this quantity is computationally costly so we did not track the evolution of $\lambda_{\min}(H_\tau)$ systematically. But a spike-like behavior $(\lambda_{\min}(H_\tau)$ goes down and then up again) was never observed in several sporadic experiments in which it was tracked.

- |D|, where $D \subseteq [S] \times [m]$ is the set of all pairs (ν, j) that change their activation at least once during training;
- $||W_T W_0||_F$. Same remark as in the first item applies (except that it is not very costly).

In the main series of our experiments, the vector $y \in \mathbb{R}^m$ was chosen from $\mathcal{N}(0,S)^m$, and we considered four values of S: $S=100,\ 200,\ 500,\ 1000$. The number of input data m ranged from 100 to 1000, with step $\frac{S}{10}$. For every triple (n,S,m) we performed 10 experiments, with fresh values of (X,y,W_0,z_0) each time. All these experiments resulted in convergence as defined above but for $S \in \{500,1000\},\ m \geq 900$ we had to decrease the learning rate to $\eta_w = 5 \cdot 10^{-4}$ when S = 500 and to $\eta_w = 2 \cdot 10^{-4}$ when S = 1000.

The ranges for our control parameters are partially⁸ reported in Table 1, and their average values are depicted on Figures 1, 2 and 3; for further comments and explanations see the respective captions.

In less systematic way, for S = 100 we also tried really large values (note that nS = 10000) m = 1000, 2000, 3000, 4000, 5000, with 10, 10, 6, 3 and 1 experiments, respectively. All of them have converged, Table 2 tabulates the results of these experiments; note that it quite smoothly extends the part of Table 1 for S = 100.

We also ran a few experiments (again, for S=100) with other choices of the labels y mentioned above. No significant discrepancies have been found with the high spectrum case and lower spectrum case.

⁸To make it comprehensible, we confine ourselves to $m = 100, 200, \ldots, 1000$ for all four values of S. Our main purpose for presenting results in the tabular form is to demonstrate that there is sufficiently sharp concentration in 10 experiments performed for every pair (S, m).

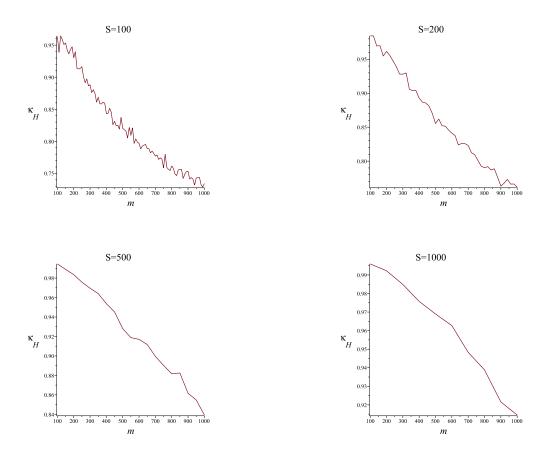


Figure 1: The dependence of κ_H on m.

S	m	T	κ_H	D	$ W_T-W_0 _F$
100	100	468-624; 528	0.89- 1.00; 0.96	657-1024; 805	15.91-22.86; 18.70
100	200	664-772; 699	0.88- 0.96; 0.93	2091-2482; 2232	27.10-31.74; 28.59
100	300	825-905; 861	0.84- 0.91; 0.89	3801-4449; 4074	34.81-39.73; 37.14
100	400	931-1054; 991	0.82- 0.88; 0.84	5763-6952; 6072	40.81-48.48; 43.33
100	500	1091-1176; 1129	0.80- 0.84; 0.82	8133-8802; 8397	47.43-51.65; 49.81
100	600	1216-1388; 1282	0.74- 0.84; 0.80	10858-12893; 11311	55.23-63.03; 57.30
100	700	1362-1492; 1418	0.76- 0.82; 0.78	13496-15722; 14282	60.32-68.55; 63.96
100	800	1489-1623; 1547	0.72- 0.78; 0.76	15974-17881; 16747	61.82-70.55; 67.37
100	900	1591-1816; 1685	0.71- 0.80; 0.75	19126-21581; 19958	68.92-78.38; 72.85
100	1000	1784-1915; 1840	0.71- 0.75; 0.73	22411-24654; 23255	76.52-81.95; 79.26
200	100	235-264; 245	0.95- 1.00; 0.98	888-1451; 1084	14.68-22.89; 17.77
200	200	306-330; 314	0.92- 0.98; 0.96	2785-3308; 3045	26.09-29.24; 27.51
200	300	362-384; 368	0.91- 0.94; 0.93	5163-5971; 5530	32.50-38.23; 34.80
200	400	406-441; 425	0.87- 0.92; 0.89	8116-9286; 8740	39.43-45.76; 42.61
200	500	446-505; 480	0.83- 0.90; 0.86	11173-13078; 12373	43.74-52.91; 49.63
200	600	508-538; 524	0.82- 0.86; 0.84	15004-17063; 16056	51.39-58.14; 54.92
200	700	547-590; 570	0.81- 0.84; 0.82	19367-21210; 20345	57.03-63.41; 60.18
200	800	610-654; 632	0.77- 0.82; 0.79	23235-25827; 24723	61.05-67.69; 65.17
200	900	674-707; 691	0.75- 0.77; 0.76	28437-31031; 29753	68.31-74.35; 70.89
200	1000	706-785; 728	0.72- 0.78; 0.76	32606-36192; 34320	70.12-80.08; 74.74
500	100	91-98; 95	0.98- 1.00; 0.99	1654-2065; 1822	16.23-20.71; 18.64
500	200	114-122; 117	0.98- 0.99; 0.98	4356-5386; 4869	24.75-30.71; 27.61
500	300	129-137; 133	0.95- 0.98; 0.97	8453-9347; 8898	32.96-37.27; 35.12
500	400	144-150; 148	0.94- 0.97; 0.95	13192-14614; 13684	39.68-43.53; 41.36
500	500	159-166; 163	0.92- 0.94; 0.93	18258-20345; 19450	44.44-49.85; 47.62
500	600	172-178; 175	0.90- 0.93; 0.92	24445-26661; 25699	48.71-55.22; 52.66
500	700	182-196; 187	0.87- 0.93; 0.90	30777-35985; 32429	54.55-65.46; 57.59
500	800	193-209; 201	0.86- 0.91; 0.88	36625-42393; 39962	57.16-66.52; 62.20
500	900	427-459; 437	0.84- 0.88; 0.86	44111-48943; 46445	62.97-70.43; 66.44
500	1000	463-478; 470	0.83- 0.85; 0.84	52602-58394; 55638	68.06-75.92; 72.29
1000	100	94-101; 96	0.99- 1.00; 1.00	2123-3066; 2546	15.78-22.89; 18.33
1000	200	112-119; 115	0.98- 1.00; 0.99	5831-7786; 6688	22.94-31.31; 26.70
1000	300	125-131; 130	0.98- 1.00; 0.98	11194-14368; 12620	31.81-40.68; 35.31
1000	400	137-146; 142	0.96- 0.99; 0.98	17651-19942; 18908	37.78-43.71; 40.45
1000	500	147-156; 152	0.95- 0.99; 0.97	25194-27879; 26485	43.62-49.31; 45.75
1000	600	158-166; 162	0.95- 0.98; 0.96	33262-37519; 35391	48.37-53.29; 51.03
1000	700	168-180; 174	0.92- 0.98; 0.95 ¿	43723-48665; 45819	54.62-59.54; 56.94
1000	800	180-189; 185	0.92- 0.95; 0.94 ²	53662-60491; 56434	57.42-66.78; 61.55
1000	900	487-503; 498	0.91- 0.94; 0.92	62774-67471; 64791	63.01-68.03; 65.22
1000	1000	511-532; 520	0.90- 0.93; 0.91	74727-78658; 76410	67.13-70.88; 69.24

Table 1: Average values are shown in bold.

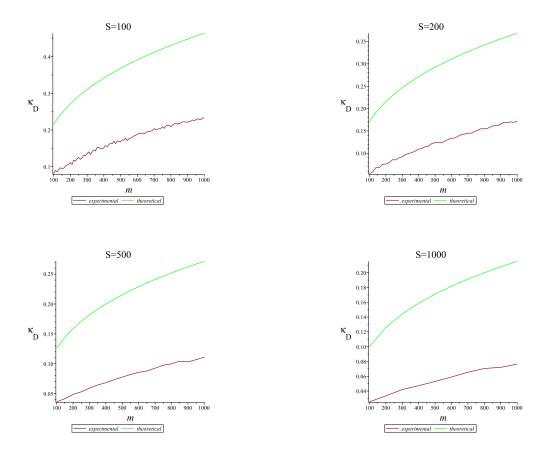


Figure 2: $\kappa_D \stackrel{\text{def}}{=} \frac{|D|}{|mS|}$ is the *density* of changes. The green line shows the theoretical prediction $\left(\frac{m}{nS}\right)^{1/3}$ obtained by dividing (62) by S.

m	T	κ_H	D	$ W_T-W_0 _F$
1000	1739 - 1948; 1859	0.7 - 0.76; 0.73	21835 - 24793; 23321	74.62 - 83.61; 79.73
2000	3363 - 3731; 3581	0.62 - 0.68; 0.64	57838 - 61459; 59157	125.67 - 136.35; 131.03
3000	6482 - 7088; 6689	0.53 - 0.57; 0.56	99676 - 102597; 100896	186 - 192.5; 189.2
4000	11558 - 12199; 11984	0.48 - 0.5; 0.49	146339 - 146489; 146410	245 - 255.3; 251.8
5000	28217	0.36	210147	360.8

Table 2: $S=100,\,m$ is large. Average values are shown in bold. The values for m=1000 were recomputed anew.

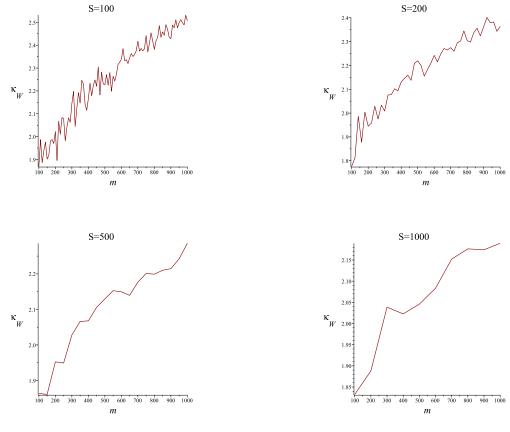


Figure 3: $\kappa_W \stackrel{\text{def}}{=} \frac{||W_T - W_0||_F}{m^{1/2}}$ (cf. (32))

The situation appears more interesting in the local case, that is when the initial error is placed on just one data point; in fact, this is the only situation where we have discovered anything that looks like a sharp transition. To report on the most striking experiment, for $n=S=100,\ m=1000,\ \eta_w=10^{-3}$ the gradient descent converged 19 times out of 20. While for m=2000, even after decreasing η_w to 10^{-4} , it aborted (according to the stopping rule $||e_\tau|| > ||e_{\tau-1}||$) in 10 (out of 10) attempts. In yet another experiment, decreasing the learning rate even further did not help convergence. But since these effects seem to be happening well beyond the range $m \ll S$ covered by our (and all previous) paper, we have not tried to investigate this systematically but rather defer this to future research.

5. Conclusion

In this paper we have improved on previously known convergence guarantees of the plain gradient descent in the case of shallow networks. The real significance of this whole line of work depends on the yet elusive properties of data and labels that allow for the generalization. In the optimistic scenario (from convergence to generalization) this research will be eventually extended into the practical region and then it might turn out to be helpful to explain why certain choices of labels allow us not only to converge but also to generalize (again, cf. [ZBH⁺21]). In the other scenario, generalization will be explained in such a way that the implied convergence will clearly follow from this ad-hoc, label-dependent explanation. That might render the whole line or research this paper belongs to obsolete.

But while the jury is out, let us briefly sketch what, in our view, are worthy directions for future work assuming the first scenario.

First and foremost, it would be nice to be able to extend our knowledge for depth 2 networks to networks of larger constant depth, without significantly weakening the results.

It would be extremely nice (and probably very difficult, too) to prove convergence outside of the NTK regime, i.e. find more intelligent ways of controlling the learning trajectory than simply by the distance travelled. A significant improvement along these lines would be to remove or relax the assumption $m \ll S$ ubiquitously present in our and all previous work. The invariant from [ACH18, DHL18] that we will explore in Appendix F seems to be a paradigmatic example of the tools that are to be developed for the

purpose.

It would be interesting to pinpoint those "quasirandom properties" from Section 3.1 (as well as the previous work) that are most problematic from the practical viewpoint. That might also give a good indication of what might be promising directions to generalize known results.

Finally, many (if not most) practical applications do not fit the set-up in Section 1.1: they either deal with (multi)-classification problems or with more sophisticated network architectures, like CNN or Resent, or both. A great deal of interesting work was done in these directions; some of it was cited in Section 1. But our impression is that less effort has been invested into working out more general united theory that would be less prone to this kind of changes in the model. That, in our view, is another interesting, and potentially more accessible, goal.

References

- [ABR⁺21] A. Atserias, I. Bonacina, S. Rezende, M. Lauria, J. Nordström, and A. Razborov. Clique is hard on average for regular resolution. *Journal of the ACM*, 68(4), 2021.
- [ACH18] Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, pages 244–253, 2018.
- [ADH⁺19a] Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Proceedings of the 33rd Conference on Neural Information Processing Systems* (NeurIPC), pages 8139–8148, 2019.
- [ADH⁺19b] Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *Proceedings of the 36th International Conference on Machine Learning* (*ICML*), pages 322–332, 2019.
- [AGNZ18] Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression

- approach. In Proceedings of the 35th International Conference on Machine Learning (ICML), pages 254–263, 2018.
- [ALS19] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *Proceedings of the 36th International Conference on Machine Learning* (*ICML*), pages 242–252, 2019.
- [BAM22] Simone Bombari, Mohammad Hossein Amani, and Marco Mondelli. Memorization and optimization in deep neural networks with minimum over-parameterization. Technical Report 2205.10217 [stat.ML], arxiv e-print, 2022.
- [BFT17] Peter L. Bartlett, Dylan J. Foster, and Matus Telgarsky. Spectrally-normalized margin bounds for neural networks. In *Proceedings of the 31st Conference on Neural Information Processing Systems (NeurIPC)*, pages 6240–6249, 2017.
- [BGMS18] Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz. SGD learns over-parameterized networks that provably generalize on linearly separable data. In *Proceedings of the 6th International Conference on Learning Representations* ((ICLR)), 2018.
- [BHMM19] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias-variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [CCZG21] Zixiang Chen, Yuan Cao, Difan Zou, and Quanquan Gu. How much over-parameterization is sufficient to learn deep ReLU networks? In *Proceedings of the 9th International Conference on Learning Representations ((ICLR))*, 2021.
- [CGW89] F. Chung, R. Graham, and R. Wilson. Quasi-random graphs. Combinatorica, 9:345–362, 1989.
- [COB19] Lénaïc Chizat, Edouard Oyallon, and Francis R. Bach. On lazy training in differentiable programming. In *Proceedings of the 33rd Conference on Neural Information Processing Systems* (NeurIPC), pages 2933–2943, 2019.

- [CR21] L. Coregliano and A. Razborov. Natural quasirandom properties. Technical Report 2012.11773 [math.CO], arxiv e-print, 2021.
- [DHL18] Simon S. Du, Wei Hu, and Jason D. Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. In *Proceedings of the 32nd Conference on Neural Information Processing Systems (NeurIPC)*, pages 382–393, 2018.
- [DZPS19] Simon S. Du, Xiyu Zhai, Barnabás Póczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *Proceedings of the 7th International Conference on Learning Representations ((ICLR))*, 2019.
- [FCB22] Spencer Frei, Niladri S. Chatterji, and Peter L. Bartlett. Benign overfitting without linearity: Neural network classifiers trained by gradient descent for noisy linear data. In *Proceedings of the* 35th Conference on Learning Theory (COLT), pages 2668–2703, 2022.
- [HY20] Jiaoyang Huang and Horng-Tzer Yau. Dynamics of deep neural networks and neural tangent hierarchy. In *Proceedings of the* 37th International Conference on Machine Learning (ICML), pages 4542–4551, 2020.
- [JT20] Ziwei Ji and Matus Telgarsky. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow ReLU networks. In *Proceedings of the 8th International Conference on Learning Representations ((ICLR))*, 2020.
- [LL18] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Proceedings of the 32nd Conference on Neural Information Processing Systems* (NeurIPC), pages 8168–8177, 2018.
- [LZ22] Bochen Lyu and Zhanxing Zhu. Implicit bias of adversarial training for deep neural networks. In *Proceedings of the 10th International Conference on Learning Representations ((ICLR))*, 2022.

- [MZ20] Andrea Montanari and Yiqiao Zhong. The interpolation phase transition in neural networks: Memorization and generalization under lazy training. Technical Report 2007.12826 [stat.ML], arxiv e-print, 2020.
- [NBS18] Behnam Neyshabur, Srinadh Bhojanapalli, and Nathan Srebro. A PAC-bayesian approach to spectrally-normalized margin bounds for neural networks. In *Proceedings of the 6th International Conference on Learning Representations ((ICLR))*, 2018.
- [Ngu21] Quynh Nguyen. On the proof of global convergence of gradient descent for deep ReLU networks with linear widths. In *Proceedings of the 38th International Conference on Machine Learning* (ICML), pages 8056–8062, 2021.
- [NMM21] Quynh Nguyen, Marco Mondelli, and Guido F. Montúfar. Tight bounds on the smallest eigenvalue of the neural tangent kernel for deep ReLU networks. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, pages 8119–8129, 2021.
- [OS20] 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.
- [SH18] Daniel Soudry and Elad Hoffer. Exponentially vanishing suboptimal local minima in multilayer neural networks. In *Proceed*ings of the 6th International Conference on Learning Representations ((ICLR)), 2018.
- [SRP+21] Chaehwan Song, Ali Ramezani-Kebrya, Thomas Pethick, Armin Eftekhari, and Volkan Cevher. Subquadratic overparameterization for shallow neural networks. In *Proceedings of the 35th Conference on Neural Information Processing Systems* (NeurIPC), pages 11247–11259, 2021.
- [SY19] Zhao Song and Xin Yang. Quadratic suffices for overparametrization via matrix Chernoff bound. Technical Report 1906.03593 [cs.LG], arxiv e-print, 2019.

- [Tro12] Joel Tropp. User-friendly tail bounds for sums of random matrices. Foundations of Computational Mathematics, 12:389–434, 2012.
- [Ver12] Roman Vershinin. Introduction to the non-asymptotic analysis of random matrices, chapter 5, pages 210–268. Cambridge University Press, 2012.
- [Ver18] Roman Vershinin. High-Dimensional Probability: An Introduction with Applications in Data Science. Cambridge University Press, 2018.
- [VL22] Tiffany J. Vlaar and Benedict J. Leimkuhler. Multirate training of neural networks. In *Proceedings of the 39th International Conference on Machine Learning (ICML)*, pages 22342–22360, 2022.
- [WDW19] Xiaoxia Wu, Simon S. Du, and Rachel Ward. Global convergence of adaptive gradient methods for an over-parameterized neural network. Technical Report 1902.07111 [cs.LG], arxiv e-print, 2019.
- [XLS] Bo Xie, Yingyu Liang, and Le Song. Diverse neural network learns true target functions. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, (AISTATS).
- [ZBH⁺21] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- [ZCZG20] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-parameterized deep ReLU networks. Machine Learning, 109(3):467–492, 2020.
- [ZG19] Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural networks. In *Proceedings of the 33rd Conference on Neural Information Processing Systems* (NeurIPC), pages 2053–2062, 2019.

A. Gradient descent avoids singular points

Fix X, W_0, z_0 such that θ_0 is regular, that is, W_0X does not contain zero entries. Let $\eta_w, \eta_z \geq 0$ be arbitrary. Our goal in this section is to show that for Lebesque almost all $y \in \mathbb{R}^m$, all points θ_τ ($\tau = 1, 2, 3, \ldots$) are also regular.

Due to countable additivity, it is sufficient too prove that

$$\lambda\left(\left\{y\mid\theta_0(y),\theta_1(y),\ldots\theta_{\tau-1}(y)\text{ are regular, }\theta_\tau(y)\text{ is singular}\right\}\right)=0$$

for any fixed positive integer τ . By the same token, it is sufficient to show that for any fixed collection $A_0, A_1, \ldots, A_{\tau-1}$ of 0-1 $(S \times m)$ matrices and for any fixed $\nu \in [S], j \in [m]$,

$$\lambda(\{y|\theta_0(y), \theta_1(y), \dots, \theta_{\tau-1}(y) \text{ are regular},$$

 $A_0(y) = A_0, \ A_1(y) = A_1, \dots, A_{\tau-1}(y) = A_{\tau-1}, \ (W_{\tau}(y)X)_{\nu j} = 0\}) = 0.$

But once the activation matrices $A_0, A_1, \ldots, A_{\tau-1}$ are fixed, a simple induction on τ shows that all our mappings, including $(W_{\tau}(y)X)_{\nu j}$, become polynomial in y. Moreover, if we initialize $y := f_0$ then θ_0 becomes the global minimum, $\theta_0(y) = \theta_1(y) = \ldots = \theta_{\tau}(y)$ and hence $(W_{\tau}(y)X)_{\nu j} = (W_0X)_{\nu j} \neq 0$. Hence this polynomial is not identically zero and therefore the set of its zeros has Lebesque measure 0.

B. Easy quasirandom properties

In this section we prove that all properties in Section 3.1, except for (29), (30), hold with probability $\geq 1 - (nS)^{-\omega(1)}$; all proofs in this section are routine exercises. As a preliminary remark, note that the nature of the error probability $(nS)^{-\omega(1)}$ (along with $m \leq n^{O(1)}$) allows us to freely use the union bound over sets whose cardinality is polynomial in n, m, S.

(19) This is essentially [Ver18, Remark 3.2.5] but let us do estimates slightly more carefully.

By the remark just made, we can assume that j,j' are fixed. Let $x,y \sim \mathcal{N}\left(0,\frac{1}{n}I_n\right)$ be two independent samples. Then, due to isotropy, X^j and $X^{j'}$ can be alternately represented as $X^j \sim \frac{x}{||x||}$, $X^{j'} \sim \frac{y}{||y||}$ and hence $\left|\langle X^j, X^{j'} \rangle \right| = \frac{\langle x,y \rangle}{||x|| \cdot ||y||}$.

Now, for any $i \in [n]$, the random variables $x_i^2 - 1/n$, $y_i^2 - 1/n$, $x_i y_i$ are centered and $||x_i^2 - 1/n||_{\psi_1}$, $||y_i^2 - 1/n||_{\psi_1}$, $||x_i y_i||_{\psi_1} \leq O(1)$. [Ver18,

Example 2.5.8(i); Lemma 2.7.7; Exercise 2.7.10]. Hence by Bernstein's inequality ([Ver18, Theorem 2.8.1]; plug $N\mapsto n,\ t\mapsto \frac{\log(nS)}{n^{1/2}}$), with probability $\geq 1-\exp(-\Omega(\log(nS))^2)\geq 1-(nS)^{-\omega(1)}$ we have $||x||^2,\ ||y||^2\geq 1-\frac{\log(nS)}{n^{1/2}}\geq \frac{1}{2}$ and $|\langle x,y\rangle|\leq \frac{\log(nS)}{n^{1/2}}\leq \widetilde{O}(n^{1/2})$. The bound (19) follows.

(20) We can treat each value of k individually (and then apply the union bound).

Let us first consider the case k = m. By [Ver12, Theorem 5.39] (plug $N \mapsto m$, $A \mapsto n^{1/2} \cdot X^{\top}$), there are absolute constants c, C > 0 such that for every t > 0,

$$\mathbf{P}[||n^{1/2} \cdot X|| \le C(m^{1/2} + n^{1/2}) + t] \ge 1 - 2\exp(-ct^2).$$

Set $t \stackrel{\text{def}}{=} m^{1/2} \log(nS)$. Then we get that $||X|| \leq \widetilde{O}\left(1 + \frac{m^{1/2}}{n^{1/2}}\right)$ with probability $\geq 1 - \exp(-\Omega(m\log(nS)^2)) \geq 1 - (nS)^{-\omega(1)}$.

Let now $k \in [m]$ be arbitrary and $J \in {[m] \choose k}$. Then, plugging in the above argument $m \mapsto k$, we see that $||X^J|| \leq \widetilde{O}\left(1 + \frac{k^{1/2}}{n^{1/2}}\right)$ with probability $\geq 1 - \exp(-\Omega(k\log(nS)^2))$, and $\exp(\Omega(k\log(nS)^2))$ dominates the number ${m \choose k} \leq \exp(k\log m)$ of choices of J (recall again that $m \leq n^{O(1)}$). Hence we can apply the union bound to complete the proof of (20).

(21) This time we have not been able to find a convenient off-the-shelf inequality so we will do a simple ad hoc argument using ϵ -nets.

We set $n^* \stackrel{\text{def}}{=} n(\log n)^2$ and $\epsilon \stackrel{\text{def}}{=} \frac{1}{m}$. By Fact 2.1, we can find an ϵ -net \mathcal{N} in S^{n-1} with $|\mathcal{N}| \leq \exp(O(n\log n))$. Then it is sufficient to prove that for any fixed $\xi \in \mathcal{N}$ we have

$$\mathbf{P}\left[\forall J \in {\binom{[m]}{n^*}} | |(X^J)^\top \xi|| \ge \frac{2n}{m}\right] \ge 1 - \exp(-\omega(n\log n)) \tag{49}$$

since then we can apply the union bound over all $\xi \in \mathcal{N}$, followed by an application of Fact 2.2 (with $\sigma := n/m$; note that $||X^J|| \leq \tilde{O}(1)$ by the already proven (20)).

In order to prove (49), by isotropy we can assume that $\xi = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$,

that is $X^{\top}\xi = X_1^{\top}$. Let $R \stackrel{\text{def}}{=} \frac{n^{1/2}}{m}$, then (by an argument similar to the

one used in the proof of (19)) $\mathbf{P}[|X_{1j}| \leq R] \leq 3Rn^{1/2} \leq \frac{3n}{m}$ and hence $\mathbf{E}[|\{j \mid |X_{1j}| \leq R\}|] \leq 3n \leq \frac{n^*}{3}$. Hence, by the plain Chernoff bound,

$$\mathbf{P}\Big[|\{j \mid |X_{1j}| \le R\}| \le \frac{n^*}{2}\Big] \ge 1 - \exp(-\Omega(n^*)) \ge 1 - \exp(-\omega(n\log n)).$$

On the other hand, this event logically implies that for every $J \in {[m] \choose n^*}$, $(X^J)^{\top}\xi$ contains at least $\frac{n^*}{2}$ entries X_{1j} with $|X_{1j}| \geq R$. Hence $||(X^J)^{\top}\xi|| \geq R\left(\frac{n^*}{2}\right)^{1/2} \geq \frac{2n}{m}$. This completes the proof of (49) and hence of (21) as well.

(22) Fix $\nu \in [S]$. Like in the proof of (19), for every $i \in [n]$, $(W_0)_{\nu i}^2 - 1$ is a centered distribution with $||(W_0)_{\nu i}^2 - 1||_{\psi_1} \leq O(1)$. Applying to it Bernstein's inequality ([Ver18, Theorem 2.8.1], plug $N \mapsto m$, $X_i \mapsto (W_0)_{\nu i}^2 - 1$, $t \mapsto n/2$), we get

$$\mathbf{P}[||(W_0)_{\nu}||^2 \le n/2] \ge 1 - \exp(-\Omega(n)) \ge 1 - (nS)^{-\omega(1)}.$$

(23) The inequality $||W_0||_{\infty} \leq \tilde{O}(1)$ is obvious. For $||z_0||_{\infty} \leq \tilde{O}(1)$, (13) implies that for sufficiently large C,

$$\mathbf{P}[|\zeta| \le (\log nS)^C] \ge 1 - (nS)^{-\omega(1)}. \tag{50}$$

- (24) The bound (14) on the variance allows us to conclude that for some $\zeta_0 \geq \widetilde{\Omega}(1)$, $\mathbf{E}[|\zeta|^2 \cdot \mathbf{1}(|\zeta| \geq \zeta_0)] \geq \widetilde{\Omega}(1)$ and then the bound (50) implies $\mathbf{P}[|\zeta| \geq \zeta_0] \geq \widetilde{\Omega}(1)$. The desired bound (24) now follows from the plain Chernoff inequality.
- (25) is obvious (in fact, it holds with probability 1).
- (26) Due to isotropy, every individual entry $(W_0X)_{\nu j}$ is a standard Gaussian.
- (27) By the just proven (26), we also have $||F_0||_{\infty} \leq \tilde{O}(1)$ and then, for any fixed F_0 , the jth entry of $f_0 = F_0^{\top} z_0$ is of the form $a_1 \zeta_1 + \ldots + a_S \zeta_S$, where $a_{\nu} \stackrel{\text{def}}{=} (F_0)_{\nu j}$ is fixed, $|a_{\nu}| \leq \tilde{O}(1)$ and ζ_1, \ldots, ζ_S are S independent copies of ζ . Now we only have to apply the sub-Gaussian Hoeffding inequality [Ver18, Theorem 2.6.2].
- (28) Fix $j \in [m]$. Due to the term +1 in RS + 1, we can assume w.l.o.g. that $R \ge \frac{1}{S}$; we can also assume that $R \le 1$ as otherwise the bound is trivial. By replacing R with the nearest rational of the form 2^{-h} (h = 0, 1, 2, ...), we

can assume that R takes on only $\widetilde{O}(1)$ different values. Hence we can apply the union bound on R and assume that $R \in \left[\frac{1}{S}, 1\right]$ is fixed.

Now, for an fixed $j \in [m]$, the values $(W_0X)_{\nu j}$ are i.i.d. from $\mathcal{N}(0,1)$. Hence $\mathbf{P}[|(W_0X)_{\nu j}| \leq R] \leq \widetilde{O}(R)$ and now (28) follows from the plain Chernoff bound.

C. NTK matrix at initialization

In this section we prove (29) and (30).

Let $w \sim \mathcal{N}(0, I_n)$. Recall that the *limit NTK matrices* (that we will denote by $H^w(X)$ and $H^z(X)$, respectively) are defined as the following expectations:

$$H^{w}(X) \stackrel{\text{def}}{=} \mathbf{E}_{w} \left[(X^{\top}X) \circ \left(\sigma'(X^{\top}w) \sigma'(w^{\top}X) \right) \right]$$

$$H^{z}(X) \stackrel{\text{def}}{=} \mathbf{E}_{w} \left[\sigma(X^{\top}w) \sigma(w^{\top}X) \right].$$

The first step is to show that (19) implies $\lambda_{\min}(H^w(X))$, $\lambda_{\min}(H^z(X)) \ge \Omega(S)$. The argument is essentially the same as in [NMM21] but since we are working in way less general set-up, we prefer to present a self-contained (and simpler) proof.

The entries $H^w(X)_{jj'}$, $H^z(X)_{jj'}$ of those matrices depend only on $\langle X^j, X^{j'} \rangle$ and this dependence is provided by analytical (near 0) functions f^w, f^z , respectively such that all coefficients in their Taylor expansions are nonnegative.

Indeed, by isotropy we can assume that

$$X^{j} = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad X^{j'} = \begin{pmatrix} \cos \phi \\ \sin \phi \\ \vdots \\ 0 \end{pmatrix},$$

where $\phi \stackrel{\text{def}}{=} \arccos\left(\langle X^j, X^{j'}\rangle\right)$. Next, (w_1, w_2) is distributed as $(r\cos\psi, r\sin\psi)$, where r is the squared χ -distribution with 2 degrees of freedom, $\psi \sim_R [0, 2\pi]$ and r, ψ are independent. We now compute

$$(H^{w}(X))_{jj'} = (X^{\top}X)_{jj'} \cdot \frac{1}{2\pi} \int_{\phi-\pi/2}^{\pi/2} 1d\psi = \langle X^{j}, X^{j'} \rangle \left(\frac{1}{2} - \frac{\phi}{2\pi}\right)$$

that is,

$$f^w(\gamma) = \gamma \left(\frac{1}{2} - \frac{\arccos \gamma}{2\pi}\right)$$

(this computation already appeared several times in the literature). Also, assuming $\phi \in [0, \pi]$,

$$(H^{z}(X))_{jj'} = \mathbf{E}[r^{2}] \cdot \frac{1}{2\pi} \int_{\phi-\pi/2}^{\pi/2} \cos(\phi - \psi) d\psi = \frac{\cos(\phi) \cdot (\pi - \phi) + \sin(\phi)}{2\pi},$$

that is

$$f^{z}(\gamma) = \frac{\gamma(\pi - \arccos(\gamma)) + \sqrt{1 - \gamma^{2}}}{2\pi}.$$

These functions have the following explicit Taylor expansions:

$$f^{w}(\gamma) = \frac{1}{4}\gamma + \frac{1}{2\pi} \cdot \sum_{r=0}^{\infty} \frac{(2r)!}{4^{r}(r!)^{2}(2r+1)} \gamma^{2r+2},$$

$$f^{z}(\gamma) = \frac{1}{2\pi} + \frac{1}{4}\gamma + \frac{1}{48\pi}\gamma^{2} + \frac{1}{2\pi} \sum_{r=2}^{\infty} \frac{(2r-3)!!}{(2r)!!} \gamma^{2r}$$

which verifies our claim.

Now, since $H^w(X)$ is obtained from the ordinary kernel matrix $X^\top X$ by entryism applications of f^w , we have

$$H^{w}(X) = \frac{1}{4}X^{\top}X + \frac{1}{2\pi} \cdot \sum_{r=0}^{\infty} \frac{(2r)!}{4^{r} (r!)^{2} (2r+1)} \underbrace{(X^{\top}X \circ \cdots \circ X^{\top}X)}_{2r+2 \text{ times}}$$

and, by Shur's theorem, all summands in the right-hand side are positive semi-definite. Hence for every fixed $r \geq 1$, $H^w(X) \succeq \Omega\left(\underbrace{(X^\top X \circ \cdots \circ X^\top X)}_{2r \text{ times}}\right)$. Also, for a fixed r > 0 all off-diagonal entries in $\underbrace{(X^\top X \circ \cdots \circ X^\top X)}_{2r \text{ times}}$ are bounded as $\widetilde{O}(n^{r/2})$ (due to (19)) which is o(m) if r > 0 is large enough. On the other hand, all diagonal entries are still $\Omega(1)$. Hence $\lambda_{\min}\left(\underbrace{(X^\top X \circ \cdots \circ X^\top X)}_{2r \text{ times}}\right) \geq \Omega(1)$ which completes the proof of $\lambda_{\min}(H^w(X)) \geq \Omega(1)$. The argument for $\lambda_{\min}(H^z(X)) \geq \Omega(1)$ is identical.

The proof of (29) is now easy to complete by the matrix Chernoff inequality. Fix an X satisfying (19). Then G_0 is the sum of S independent copies of the random (rank 1) PSD matrix $\sigma(w^{\top}X)^{\top}\sigma(w^{\top}X)$ while $H^z(X)$ is its expectation. Moreover,

$$||\sigma(w^{\top}X)^{\top}\sigma(w^{\top}X)|| = ||\sigma(w^{\top}X)||^2 \le ||w^{\top}X||^2.$$

Let

$$H^{z}(X, w) \stackrel{\text{def}}{=} \begin{cases} \sigma(w^{\top} X)^{\top} \sigma(w^{\top} X) & \text{if } ||w^{\top} X||^{2} \leq m(\log(nS))^{2} \\ 0 & \text{otherwise.} \end{cases}$$

Then $||H^z(X, w)|| \leq \tilde{O}(m)$ and

$$||\mathbf{E}_w[H^z(X,w)] - H^z(X)|| \le \mathbf{E}_w[||w^\top X||^2 \cdot \mathbf{1} \left(||w^\top X||^2 \ge m(\log(nS))^2\right)]$$

$$\le \mathbf{E}_w[||w^\top X||^2 \cdot \mathbf{1} \left(||w^\top X||_{\infty} \ge \log(nS)\right)] \le o(1)$$

(for the last inequality, observe that $w^{\top}X$ is a tuple of standard Gaussian, albeit not necessary independent). This in particular implies that

$$\lambda_{\min}\left(\mathbf{E}_w[H^z(X,w)]\right) \ge \Omega(1),$$

and now we can apply Matrix Chernoff Inequality ([Tro12, Theorem 1.1]; plug $d \mapsto m$, $X_k \mapsto H^z(X, w)$, $R \mapsto \widetilde{O}(m)$, $\mu_{\min} \mapsto \Omega(S)$, $\delta \mapsto 1/2$) to conclude that (29) holds with probability $\geq 1 - \exp(-\Omega(S/m)) \geq 1 - \exp(-\widetilde{\omega}(1)) \geq 1 - (nS)^{-\omega(1)}$.

The proof of (30) is tricker, and, as we indicated above, there probably should be an easier way of doing this.

We first note that Γ_0 depends only on z_0 while $(A_0 *X)^{\top}(A_0 *X)$ depends only on X, W_0 . Hence we can fix an arbitrary $\Gamma_0 \subseteq [S]$ with $|\Gamma_0| \ge \widetilde{\Omega}(S)$ and prove (30) conditioned by $\Gamma_0(z_0) = \Gamma_0$. After that we can simplify our notation by assuming w.l.o.g. that $\Gamma_0 = [S]$ (for the Rademacher initialization this holds automatically anyway).

For any $\Gamma \subseteq [S]$,

$$||(A_0)_{\Gamma} * X|| \le ||A_0 * X|| = ||(A_0^{\top} A_0) \circ (X^{\top} X)||^{1/2} \stackrel{(8)}{\le} S^{1/2} \cdot ||X|| \stackrel{(20)}{\le} \widetilde{O} \left(S^{1/2} \left(1 + \frac{m^{1/2}}{n^{1/2}} \right) \right).$$

Let now

$$\delta \stackrel{\text{def}}{=} \lambda_{\min}(H^w(X)) \ge \Omega(1);$$

$$\sigma \stackrel{\text{def}}{=} \min\left(\frac{\delta^{1/2}}{3}, \frac{1}{2}\right) S^{1/2};$$

$$\epsilon \stackrel{\text{def}}{=} \frac{\sigma}{2 \cdot ||(A_0)_{\Gamma} * X||} \ge \Omega(m^{-1/2}).$$

Fix an arbitrary ϵ -net $\mathcal{N} \subseteq S^{m-1}$ of cardinality $\exp(\tilde{O}(m))$ (by Fact 2.1). Set also

$$S^* \stackrel{\text{def}}{=} \left\lfloor \frac{n^2 S}{(n^2 + m) \log(nS)^C} \right\rfloor,$$

where $C \geq 1$ is a sufficiently large constant, also to be specified later. Note that

$$S^* \ge \widetilde{\omega}(m),\tag{51}$$

due to $S \geq \tilde{\omega}(m)$ and the assumption $m \leq \tilde{o}(nS^{1/2})$ incorporated in (30). Then it is sufficient to show that for any fixed $\xi \in \mathcal{N}$ we have

$$\mathbf{P}\left[\forall \Gamma \in \binom{[S]}{S - S^*} \mid ||((A_0)_{\Gamma} * X)\xi|| \ge \sigma\right] \ge 1 - \exp(-\widetilde{\omega}(m)) \tag{52}$$

since after that we can apply the union bound over \mathcal{N} , and then Fact 2.2.

Before proceeding with the formal proof, let us briefly explain the predicament we are facing. The bound $\lambda_{\min}(H^w(X)) \geq \Omega(1)$ we have already proven only implies that the expectation of the random variable $\xi^{\top}(X^{\top}X \circ (\sigma'(X^{\top}w)\sigma'(w^{\top}X)))\xi$ is bounded away from zero, and $\xi^{\top}(A_0*X)^{\top}(A_0*X)\xi$ is a sum of S independent copies of that variable. Since $S \geq \widetilde{\omega}(m)$, we would have been done by a simple application of Chernoff's inequality if we knew that this random variable does not behave abnormally; say, if we knew that the probability that it is separated from 0 is also $\Omega(1)$ (cf. [SY19, Assumption 1.2.2]). We can not directly apply Markov's inequality since we only have an upper bound of $\widetilde{O}(m/n)$ on the value of this random variable. Our way to rule out this pathological situation (i.e. when the large expectation is made by large values occurring with small probability) is to apply some ideas from the proof of the Lipschitz concentration inequality [Ver18, Theorem 5.1.4].

Returning to the formal argument, let us first express our random variable in more compact way:

$$\xi^{\top}(X^{\top}X \circ (\sigma'(X^{\top}w)\sigma'(w^{\top}X)))\xi = ||X(\sigma'(X^{\top}w) \circ \xi)||^2.$$

Thus,

$$\mathbf{E}_w [||X(\sigma'(X^\top w) \circ \xi)||^2] \ge \lambda_{\min}(H^w(X)) = \delta.$$

For the rest of this section, it will be more convenient to assume that $w \in_R \sqrt{n} \cdot S^{n-1}$; since σ' is invariant under positive scalings, this will not change anything. We will denote by μ the standard (Haar) measure on $\sqrt{n} \cdot S^{n-1}$.

Let

$$K \stackrel{\text{def}}{=} \frac{\delta S}{2S^*}$$

and

$$\mathcal{W} \stackrel{\text{def}}{=} \left\{ w \in \sqrt{n} \cdot S^{n-1} \, \middle| \, ||X(\sigma'(X^{\top}w) \circ \xi)||^2 \ge K \right\}.$$

We split the analysis according to whether \mathcal{W} is small or large.

Case 1. $\mu(W) \leq 1/m$.

This case is easy. Since

$$||X(\sigma'(X^{\top}w)\circ\xi)||^{2} \leq ||X||^{2} \cdot ||\sigma'(X^{\top}w)\circ\xi||^{2} \leq ||X||^{2} \stackrel{(20)}{\leq} \tilde{O}\left(1 + \frac{m}{n}\right) \leq \tilde{o}(m),$$

we have that $\mathbf{E}[||X(\sigma'(X^{\top}w) \circ \xi)||^2 \cdot \mathbf{1}(w \in \mathcal{W})] \leq \tilde{o}(1)$. Hence if we consider the truncated (and scaled by K) function

$$f(w) \stackrel{\text{def}}{=} \min \left(\frac{1}{K} ||\sigma'(X^{\top}w) \circ \xi||^2, 1 \right),$$

we will still have $\mathbf{E}[f(w)] \geq \frac{2}{3} \frac{\delta}{K} = \frac{4}{3} \frac{S^*}{S}$. Noting that $f(w) \in [0, 1]$, we have by the plain Chernoff bound:

$$\mathbf{P} \left[f(w_1) + \ldots + f(w_S) \ge \frac{5}{4} S^* \right] \ge 1 - \exp(-\Omega(S^*)) \stackrel{(51)}{\ge} 1 - \exp(-\widetilde{\omega}(m)).$$

Now, if we remove from this sum S^* terms, it will get decreased by at most S^* and hence

$$||((A_0)_{\Gamma} * X)\xi||^2 \ge K \cdot \sum_{\nu \in \Gamma} f(w_{\nu}) \ge \frac{1}{4}KS^* = \frac{1}{8}\delta S,$$

for any $\Gamma \in \binom{[S]}{S-S^*}$. This completes the analysis of Case 1. Case 2. $\mu(\mathcal{W}) \geq 1/m$.

In this case we will be able to achieve our dream goal and show that $||X(\sigma'(X^{\top}w)\circ\xi)||\geq 1$ with probability at least 1/2. For $\rho\geq 0$, let

$$\mathcal{W}_{\rho} \stackrel{\text{def}}{=} \left\{ w \in \sqrt{n} \cdot S^{n-1} \mid \exists w^* \in \mathcal{W}(||w - w^*|| \le \rho) \right\}.$$

Then, by the standard isoperimetric inequality we can fix $\rho \leq \tilde{O}(1)$ in such a way that $\mu(W_{\rho}) \geq 2/3$. Our goal is to show that for all but a negligible fraction of "bad" $w \in W_{\rho}$, we have $||X(\sigma'(X^{\top}w) \circ \xi)|| \geq 1$.

In order to identify the first set of "bad" points in W_{ρ} , set first

$$\phi \stackrel{\text{def}}{=} \frac{\log(nS)^{C_1}}{n^{1/2}},$$

where $C_1 \geq 1$ is large enough. Let

$$E(w) \stackrel{\text{def}}{=} \left\{ j \in [m] \, \middle| \, ||w^{\top} X^{j}|| \le \phi \right\},$$

and let $\chi_w \in \{0,1\}^m$ be the characteristic function of this set. Then, since $\sum_{j \in [m]} \xi_j^2 = 1$, we have the estimate

$$\mathbf{E}_{w}[||\chi_{w} \circ \xi||^{2}] = \sum_{j \in [m]} \xi_{j}^{2} \cdot \mathbf{P}_{w}[||w^{\top} X^{j}|| \leq \phi] \leq \widetilde{O}(n^{-1/2}).$$

Hence we can choose $\psi \leq \tilde{O}(n^{-1/4})$ such that

$$\mathbf{P}_w[||\chi_w \circ \xi|| \ge \psi] \le \frac{1}{12},$$

and this is our first "bad" event.

For the second "bad" event we need to identify one more quasirandom property of the data X that in a sense is a uniform version of a dual to (28). Namely, for R > 0, $w \in \sqrt{n} \cdot S^{n-1}$ and $X \in (S^{n-1})^m$, let

$$\operatorname{Bad}_{R}(w, X) \stackrel{\text{def}}{=} \left| \left\{ j \in [m] \, \middle| \, ||w^{\top} X^{j}|| \le R \right\} \right| \ge \log(nS)^{2} \cdot (mR + 1). \tag{53}$$

Similarly to the proof of (28), for any fixed w we have $\mathbf{P}_X[\operatorname{Bad}_R(w,X)] \le \exp(-\Omega(\log(nS))^2)$ and, averaging over w,

$$\mathbf{P}[\mathrm{Bad}_R(w,X)] \le \exp(-\Omega((\log(nS))^2)) \le 1 - (nS)^{-\omega(1)}.$$

On the other hand, if we set

$$\mathcal{B}_{R,X} = \left\{ w \in \sqrt{n} \cdot S^{n-1} \mid \text{Bad}_R(w,X) \right\},\,$$

then

$$\mathbf{P}[\mathrm{Bad}_R(w,X)] = \mathbf{E}_X[\mu(\mathcal{B}_{R,X})].$$

Therefore,

$$\mathbf{P}_X[\mu(\mathcal{B}_{R,X}) \le \widetilde{o}(1)] \ge 1 - (nS)^{-\omega(1)}.$$

Let now

$$\mathcal{B}_X \stackrel{\mathrm{def}}{=} \bigcup_{R>0} \mathcal{B}_{R,X}.$$

Then we still have

$$\mathbf{P}_X[\mu(\mathcal{B}_X) \le \widetilde{o}(1)] \ge 1 - (nS)^{-\omega(1)}$$

since (cf. the proof of (28)) it is sufficient to consider only $\widetilde{O}(1)$ different values of R in this union. Thus we also require that

$$\mu(\mathcal{B}_X) \le \frac{1}{12}.$$

We stress that this is the quasirandom property of the $data\ X$ only and it does not depend in any way on the actual initialization W_0 .

The set \mathcal{B}_X is our second "bad event", and now we let

$$\widetilde{\mathcal{W}}_{\rho} \stackrel{\text{def}}{=} \left\{ w \in \mathcal{W}_{\rho} \, | \, || \chi_{w} \circ \xi || \leq \psi \wedge w \notin \mathcal{B}_{X} \right\};$$

note that $\mu(\widetilde{\mathcal{W}}_{\rho}) \geq \frac{2}{3} - 2 \cdot \frac{1}{12} = \frac{1}{2}$. We claim that

$$\forall w \in \widetilde{\mathcal{W}}_{\rho} ||X(\sigma'(X^{\top}w) \circ \xi)|| \ge 1.$$
 (54)

Indeed, fix $w \in \widetilde{\mathcal{W}}_{\rho}$ and let $w^* \in \mathcal{W}$ be such that $||w - w^*|| \le \rho \le \widetilde{O}(1)$. Let

$$d \stackrel{\text{def}}{=} \sigma'(X^{\top}w) - \sigma'(X^{\top}w^*); \ d \in \{-1, 0, 1\}^m.$$

We start with the obvious estimate

$$||X(\sigma'(X^{\top}w) \circ \xi)|| > ||X(\sigma'(X^{\top}w^*) \circ \xi)|| - ||X(d \circ \xi)|| > K^{1/2} - ||X(d \circ \xi)||$$

(the second inequality holds since $w^* \in \mathcal{W}$). Let now $D \stackrel{\text{def}}{=} \sup(d)$, and split d as d = d' + d'', where $d' \stackrel{\text{def}}{=} d \circ \chi_w$ is the part corresponding to E(w) and d'' corresponds to co - E(w). Then we have the bound

$$||X(d \circ \xi)|| \le ||X(d' \circ \xi)|| + ||X(d'' \circ \xi)||$$

$$\le ||X^{D}|| \cdot ||d' \circ \xi|| + ||X^{D \setminus E(w)}|| \cdot ||d'' \circ \xi||$$

$$\le ||X^{D}|| \cdot ||\chi_{w} \circ \xi|| + ||X^{D \setminus E(w)}||$$

$$\le \psi \cdot ||X^{D}|| + ||X^{D \setminus E(w)}||.$$
(55)

Let us estimate |D| and $|D \setminus E(w)|$ (and then we will apply (20)). For any $j \in D$, we have

$$||(w^* - w)^\top X^j|| \ge ||w^\top X^j||$$
 (56)

since $w^{\top}X^{j}$ and $(w^{*})^{\top}X^{j}$ have the opposite sign. Thus, $||(w^{*}-w)^{\top}X^{D}|| \geq ||w^{\top}X^{D}||$; recalling that $||w^{*}-w|| \leq \tilde{O}(1)$, we obtain

$$||w^{\top}X^{D}|| \le \tilde{O}\left(||X^{D}||\right) \stackrel{(20)}{\le} \tilde{O}\left(1 + \frac{|D|^{1/2}}{n^{1/2}}\right).$$

On the other hand, let

$$R \stackrel{\text{def}}{=} \frac{1}{m} \left(\frac{|D|}{2\log(nS)^2} - 1 \right).$$

so that the right-hand side in (53) becomes $\frac{|D|}{2}$. Since $w \notin \mathcal{B}_{R,X}$, there are at least $\frac{|D|}{2}$ indices $j \in D$ such that $||w^{\top}X^{j}|| \geq R$ which implies

$$||w^\top X^D|| \geq \frac{1}{\sqrt{2}} R|D|^{1/2} \geq \widetilde{\Omega}\left(\frac{|D|^{3/2}}{m}\right) - \widetilde{O}\left(\frac{|D|^{1/2}}{m}\right).$$

Comparing now the upper and lower bounds on $||w^{\top}X^{D}||$, we get $\frac{|D|^{3/2}}{m} \leq \tilde{O}\left(1 + \frac{|D|^{1/2}}{n^{1/2}} + \frac{|D|^{1/2}}{m}\right)$ which solves to $|D| \leq \tilde{O}\left(m^{2/3} + \frac{m}{n^{1/2}}\right)$. By (20), this implies

$$||X^D|| \le \widetilde{O}\left(1 + \frac{m^{1/3}}{n^{1/2}} + \frac{m^{1/2}}{n^{3/4}}\right) \le \widetilde{O}\left(1 + \frac{m^{1/2}}{n^{3/4}}\right),$$
 (57)

where the second inequality holds simply because $\frac{m^{1/3}}{n^{1/2}} = \left(\frac{m^{1/2}}{n^{3/4}}\right)^{2/3}$.

To bound $|D \setminus E(w)|$, we again use (56) which, along with the definition of E(w), gives us

$$||(w - w^*)^\top X^{D \setminus E(w)}|| \ge \phi \cdot |D \setminus E(w)|^{1/2}.$$

On the other hand,

$$||(w-w^*)^\top X^{D\setminus E(w)}|| \le \rho \cdot ||X^{D\setminus E(w)}|| \stackrel{(20)}{\le} \widetilde{O}\left(1 + \frac{|D\setminus E(w)|^{1/2}}{n^{1/2}}\right).$$

If the constant C_1 in the definition of ϕ is large enough (namely, exceeds the constant assumed in (20)), the second term $\tilde{O}\left(\frac{|D\setminus E(w)|^{1/2}}{n^{1/2}}\right)$ is dominated by $\phi|D\setminus E(w)|^{1/2}$ and therefore $\phi|D\setminus E(w)|^{1/2}\leq \tilde{O}(1)$. Thus $|D\setminus E(w)|\leq \tilde{O}(n)$ and (as always, by (20)),

$$||X^{D\setminus E(w)}|| \le \widetilde{O}(1). \tag{58}$$

Plugging (57), (58) and $\psi \leq \tilde{O}(n^{-1/4})$ into (55), we finally conclude that

$$||X(d \circ \xi)|| \le \tilde{O}\left(1 + \frac{m^{1/2}}{n}\right) \le K^{1/2} - 1,$$

provided the constant C in the definition of S^* is large enough. This completes the proof of (54).

The rest is easy. By plain Chernoff bound, $\left|\left\{\nu\in[S]\,\middle|\,w_{\nu}^{\top}\in\widetilde{\mathcal{W}}_{\rho}\right\}\right|\geq S/3$ with probability $\geq 1-\exp(-\Omega(S))\geq 1-\exp(-\widetilde{\omega}(m))$. Since $S^{*}\leq o(S)$, removing any S^{*} entries will still leave us with $\geq \frac{S}{4}$ neurons $\nu\in[S]$ such that $||X(\sigma'(X^{\top}w_{\nu}^{\top})\circ\xi)||\geq 1$. This completes the proof of (52) (since $\sigma\leq\frac{S^{1/2}}{2}$) and (30).

D. NTK matrix at the first layer

In this section we prove (36), assuming that (32), (33) and (35) hold (for the same t).

Let

$$\Gamma \stackrel{\text{def}}{=} \left\{ \nu \in \Gamma_0 \,\middle|\, |(z_t)_{\nu}| \ge \frac{1}{2} \zeta_0 \right\}.$$

Our first task is to check that this Γ satisfies the condition $|\Gamma_0 \setminus \Gamma| \leq \tilde{o}\left(\frac{n^2S}{n^2+m}\right)$ (and hence can be chosen in (30)). For that we note that for any $\nu \in \Gamma_0 \setminus \Gamma$,

 $|(z_t)_{\nu} - (z_0)_{\nu}| \geq \frac{1}{2}\zeta_0 \geq \widetilde{\Omega}(1)$ and hence $||z_t - z_0|| \geq \widetilde{\Omega}(|\Gamma_0 \setminus \Gamma|^{1/2})$. So we only have to check that

$$||z_t - z_0|| \le \min\left(\tilde{o}(S^{1/2}), \ \tilde{o}\left(\frac{nS^{1/2}}{m^{1/2}}\right)\right).$$
 (59)

We can assume (see Remark 7) that either $\frac{\eta_z}{\eta_w} \geq \widetilde{\omega}\left(\frac{m^2}{nS}\right)$ or $\frac{\eta_z}{\eta_w} \leq \widetilde{O}\left(\frac{m^2}{nS}\right)$. If $\frac{\eta_z}{\eta_w} \geq \widetilde{\omega}\left(\frac{m^2}{nS}\right)$, we can apply (35) and use the condition $m \leq \widetilde{o}(nS^{1/4})$ in (36).

If on the other hand $\frac{\eta_z}{\eta_w} \leq \widetilde{O}\left(\frac{m^2}{nS}\right)$ then the cases (17), (18) are ruled out and hence (16) must hold. Now we apply (33) to conclude that $||z_t - z_0|| \leq \widetilde{O}\left(\frac{m^{3/2}}{n^{1/2}S^{1/2}}\right) \leq \widetilde{O}\left(\frac{nS^{1/2}}{m^{1/2}}\right)$, where the last inequality follows from the calculation $m \leq \widetilde{O}\left((nS^{1/4})^{3/4} \cdot (S)^{1/4}\right) \leq \widetilde{O}(n^{3/4}S^{7/16}) \leq \widetilde{O}\left(n^{3/4}S^{1/2}\right)$. Alternatively, from (16) we have $\frac{\eta_z}{\eta_w} \leq \widetilde{O}\left(\frac{S}{m}\right)$ which, also by (33), gives us $||z_t - z_0|| \leq \widetilde{O}(S^{1/2})$. Thus in either case we have (59).

We can now apply (30) to our particular Γ , and we get $\sigma_{\min}((A_0)_{\Gamma} * X) \ge \Omega(S^{1/2})$. Further, since $(z_t)_{\nu} \ge \widetilde{\Omega}(1)$ whenever $\nu \in \Gamma$, we have

$$\sigma_{\min}((B_t)_{\Gamma} * X) \geq \widetilde{\Omega}\left(\sigma_{\min}((A_t)_{\Gamma} * X)\right).$$

Thus, all that remains to prove is that $\sigma_{\min}((A_t)_{\Gamma} * X) \geq \frac{1}{2}\sigma_{\min}((A_0)_{\Gamma} * X)$, and for that it is sufficient to establish

$$||(A_t - A_0) * X|| \le \tilde{o}(S^{1/2}).$$
 (60)

Let

$$D \stackrel{\text{def}}{=} \sup(A_t - A_0) = \left\{ (\nu, j) \in [S] \times [m] \mid \sigma'((W_t)_{\nu} X^j) \neq \sigma'((W_0)_{\nu} X^j) \right\}$$

and

$$D^{j} \stackrel{\text{def}}{=} \left\{ \nu \in [S] \mid (\nu, j) \in D \right\}.$$

Let $J \subseteq [m]$ consist of $\min(m, n)$ data instances $j \in [m]$ with the maximum value of $|D^j|$ and let $co - J \stackrel{\text{def}}{=} [m] \setminus J$. We will bound $||(A_t - A_0)^J * X^J||$ and $||(A_t - A_0)^{co - J} * X^{co - J}||$ as $\tilde{o}(S^{1/2})$ separately.

For the first term, (20) implies that $||X^J|| \leq O(1)$. Hence, by Shur's inequality (8), it is sufficient to prove that $\max_{j \in J} |D^j| \leq \tilde{o}(S)$. Fix $j \in J$,

and fix $R \geq \widetilde{\Omega}\left(\frac{|D^j|}{S}\right)$ such that the right-hand side in (28) is $\leq \frac{|D^j|}{2}$. Then for at least $\frac{|D^j|}{2}$ values $\nu \in D^j$ we have $|(W_0X)_{\nu j}| \geq R$ and hence

$$||(W_0)_{D^j}X^j|| \ge \Omega\left(R \cdot |D^j|^{1/2}\right) \ge \widetilde{\Omega}\left(\frac{|D^j|^{3/2}}{S}\right).$$

Also (cf. (56)) $\forall \nu \in D^j(|(W_0)_{\nu}X^j| \leq |(W_0 - W_t)_{\nu}X^j|)$. This gives us

$$||(W_0 - W_t)X^j|| \ge ||(W_0 - W_t)_{D^j}X^j|| \ge ||(w_0)_{D^j}X^j|| \ge \widetilde{\Omega}\left(\frac{|D^j|^{3/2}}{S}\right).$$
 (61)

On the other hand,

$$||(W_0 - W_t)X^j|| \le ||W_0 - W_t|| \le ||W_0 - W_t||_F \stackrel{(32)}{\le} \widetilde{O}(m^{1/2}) \le \widetilde{o}(S^{1/2}).$$

Comparing these two bounds proves $\max_{j\in J} |D^j| \leq \tilde{o}(S)$ and thus $||(A_t - A_0)^J * X^J|| \leq \tilde{o}(S^{1/2})$.

For the second term $||(A_t - A_0)^{co-J} * X^{co-J}||$, we can assume that $m \ge n$ (and hence |J| = n) as otherwise the statement is void. Let s be such that

$$\min_{j \in J} |D^j| \ge s \ge \max_{j \in co-J} |D^j|$$

(it exists due to our choice of J). We can now continue (61) as $||(W_0 - W_t)X^j|| \ge \widetilde{\Omega}\left(\frac{s^{3/2}}{S}\right)$ $(j \in J)$ and then conclude $||(W_0 - W_t)X^J||_F \ge \widetilde{\Omega}\left(\frac{s^{3/2}n^{1/2}}{S}\right)$. On the other hand,

$$||(W_0 - W_t)X^J||_{\mathbf{F}} \stackrel{(11)}{\leq} ||W_0 - W_t||_{\mathbf{F}} \cdot ||X^J|| \stackrel{(20), (32)}{\leq} \widetilde{O}(m^{1/2}).$$

Comparing these two bounds gives us

$$s \le \widetilde{O}\left(\frac{S^{2/3}m^{1/3}}{n^{1/3}}\right). \tag{62}$$

Applying now (8), we see that

$$||(A_t - A_0)^{co - J} * X^{co - J}|| \le \left(\max_{j \in co - J} |D^j|\right)^{1/2} \cdot ||X|| \stackrel{(20)}{\le} \widetilde{O}\left(\frac{s^{1/2} m^{1/2}}{n^{1/2}}\right)$$
$$\le \widetilde{O}\left(\frac{S^{1/3} m^{2/3}}{n^{2/3}}\right) \stackrel{m \le \widetilde{o}(nS^{1/4})}{\le} \widetilde{o}(S^{1/2}).$$

This completes the proof of (60) and hence also of (36).

E. Gradient does not change radically between steps

In this section we prove (41) assuming that the bounds (32)-(35) hold for all s < t. By applying another auxiliary induction⁹ we can assume that (41) also holds for all s < t. In particular, for s < t we may freely use all the conclusions made in Section 3.2. Finally, if t is an integer then the statement is trivially true. Hence we can assume w.l.o.g. that all inequalities in Section 3.2 hold for $\tau \stackrel{\text{def}}{=} |t|$.

Recall that

$$\nabla \ell(\theta_t) = ((B_t * X)e_t, F_t e_t);$$

$$\dot{\theta}_\tau = -(\eta_w, (B_\tau * X)e_\tau, \eta_z F_\tau e_\tau);$$

$$\left| \langle \nabla \ell(\theta_\tau), \dot{\theta}_\tau \rangle \right| = \eta_w ||(B_\tau * X)e_\tau||^2 + \eta_z ||F_\tau e_\tau||^2.$$

Let us introduce the uniform notation

$$\bar{\eta}_w = \begin{cases} \eta_w & \text{if either (16) or (17) holds} \\ 0 & \text{otherwise} \end{cases}$$

$$\bar{\eta}_z = \begin{cases} \eta_z & \text{if either (17) or (18) holds} \\ 0 & \text{otherwise} \end{cases}.$$

Then by (36) and (37), the right-hand side of (41) gets bounded from below as

$$\left| \langle \nabla \ell(\theta_{\tau}), \dot{\theta}_{\tau} \rangle \right| \geq \widetilde{\Omega} \left(S^{1/2} ||e_{\tau}|| (\bar{\eta}_{w} || (B_{\tau} * X) e_{\tau} || + \bar{\eta}_{z} || F_{\tau} e_{\tau} ||) \right)$$

$$\geq \widetilde{\Omega} \left(S ||e_{\tau}||^{2} (\bar{\eta}_{w} + \bar{\eta}_{z}) \right).$$
(63)

In order to upper bound the left-hand side in (41), let us first collect some useful estimates; in what follows, s < t is arbitrary.

First we have

$$||W_s - W_0||_{\mathcal{F}} \le \tilde{O}(m^{1/2}).$$
 (64)

In the cases (16), (17) it follows from (32), and in the case (18) – from (34) (since $\eta_z \ge \eta_w$ in that case).

 $^{^{9}}$ Recall that due to our convention the gradient is upper semi-continuous and hence the set of those t for which (41) fails is violated is either empty or contains the minimum element.

The following is the first part 10 of (59):

$$||z_s - z_0|| \le \tilde{o}(S^{1/2});$$
 (65)

along with (23) this gives us

$$||z_s|| \le \widetilde{O}(S^{1/2}).$$

We can now bound the gradient. Namely, by (10) we have

$$||B_s * X|| \le ||z_s|| \cdot ||X|| \le \widetilde{O}\left(S^{1/2}\left(1 + \frac{m^{1/2}}{n^{1/2}}\right)\right)$$
 (66)

and

$$||F_s - F_0||_{\text{F}} \le ||W_s X - W_0 X||_{\text{F}} \le ||W_s - W_0||_{\text{F}} \cdot ||X||$$

$$\stackrel{(64), (20)}{\le} \widetilde{O}\left(m^{1/2} + \frac{m}{n^{1/2}}\right) \le \widetilde{O}\left(m^{1/2} S^{1/2}\right)$$

which, along with (26), implies

$$||F_s||_{\mathsf{F}} \le \widetilde{O}\left(m^{1/2}S^{1/2}\right).$$

Now we can also control the evolution from τ to t, both in the parameter space and in the feature space. Namely, recalling that we have already proved that $||e_s||$ is decreasing for s < t and that $||e_s|| \le \tilde{O}\left(m^{1/2}S^{1/2}\right)$, we have

$$||W_{t} - W_{\tau}||_{F} = \eta_{w}(t - \tau)||\nabla^{w}\ell(\theta_{\tau})|| \leq \eta_{w} \cdot ||B_{\tau} * X|| \cdot ||e_{\tau}||$$

$$\leq \tilde{O}\left(\eta_{w} m^{1/2} S\left(1 + \frac{m^{1/2}}{n^{1/2}}\right)\right)$$
(67)

and thus

$$||F_t - F_\tau||_F \le ||W_t - W_\tau||_F \cdot ||X|| \le \tilde{O}\left(\eta_w m^{1/2} S\left(1 + \frac{m}{n}\right)\right).$$
 (68)

In the W-department, we claim that

$$m \le \tilde{o}(nS^{1/4}) \Longrightarrow ||(B_t - B_\tau) * X|| \le \tilde{o}(S^{1/2}). \tag{69}$$

¹⁰That part did not use the assumption $m \leq \tilde{o}(nS^{1/4})$.

Since our proof of this fact requires quite new ideas, it is postponed to Appendix F.

Next, for $s \in [\tau, t)$ we have

$$\dot{e}_s = -(\eta_w (B_s * X)^\top (B_\tau * X) + \eta_z F_s^\top F_\tau) e_\tau$$

and hence by the above bounds on $||B_s * X||, ||F_s||,$

$$||\dot{e}_{s}|| \leq \widetilde{O}\left(||e_{\tau}|| \cdot (\eta_{w}||B_{s} * X|| \cdot ||B_{\tau} * X|| + \eta_{z}||F_{s}|| \cdot ||F_{\tau}||\right))$$

$$\leq \widetilde{O}\left(S \cdot ||e_{\tau}|| \cdot \left(\eta_{w}\left(1 + \frac{m}{n}\right) + \eta_{z}m\right)\right) \stackrel{(15)}{\leq} \widetilde{o}\left(\frac{||e_{\tau}||}{S}\right).$$

Therefore,

$$||e_t - e_\tau|| \le \int_\tau^t ||\dot{e}_s|| ds \le \tilde{o}\left(\frac{||e_\tau||}{S}\right). \tag{70}$$

Equipped with all this knowledge, we can now proceed to completing the proof of (41).

First, we claim that in (63) we can now replace $\bar{\eta}_w$ with η_w , i.e. that

$$\left| \langle \nabla \ell(\theta_{\tau}), \dot{\theta}_{\tau} \rangle \right| \geq \widetilde{\Omega} \left(S^{1/2} ||e_{\tau}|| (\eta_{w}|||(B_{\tau} * X)e_{\tau}|| + \bar{\eta}_{z}||F_{\tau}e_{\tau}||) \right)$$

$$\geq \widetilde{\Omega} \left(S||e_{\tau}||^{2} (\eta_{w} + \bar{\eta}_{z}) \right).$$

$$(71)$$

Indeed, we only need to consider the last case (18). But then $\bar{\eta}_z = \eta_z$,

$$|\eta_z||F_{\tau}e_{\tau}|| \stackrel{(37)}{\geq} \widetilde{\Omega}(\eta_z S^{1/2}||e_{\tau}||)$$

and

$$|\eta_w||(B_\tau * X)e_\tau|| \le \widetilde{O}\left(\eta_w S^{1/2}\left(1 + \frac{m^{1/2}}{n^{1/2}}\right)||e_\tau||\right).$$

Now, the condition in (18) implies that the former expression dominates the latter.

Second, let $\delta_t = (\delta_t^w, \delta_t^z)$, where

$$\delta_t^w = (B_t * X)e_t - (B_\tau * X)e_\tau
\delta_t^z = F_t e_t - F_\tau e_\tau;$$

we bound their contributions separately.

More specifically,

$$||\delta_t^w|| \le ||(B_\tau * X)(e_t - e_\tau)|| + ||((B_t - B_\tau) * X)e_t|| \le ||B_\tau * X|| \cdot ||e_t - e_\tau|| + ||(B_t - B_\tau) * X|| \cdot ||e_\tau||.$$

We now do some case analysis.

If $m \leq \tilde{o}(nS^{1/4})$, we can apply (66), (70) and (69) to conclude that

$$||\delta_t^w|| \le \tilde{o}\left(\frac{||e_\tau||}{S^{1/2}}\left(1 + \frac{m^{1/2}}{n^{1/2}}\right) + S^{1/2}||e_\tau||\right) \le \tilde{o}(S^{1/2}||e_\tau||).$$

Hence

$$\left| \left\langle \delta_t^w, \dot{\theta}_\tau \right\rangle \right| \leq \eta_w ||\delta_t^w|| \cdot ||(B_\tau * X) e_\tau|| \leq \tilde{o} \left(\eta_w S^{1/2} ||e_\tau|| \cdot ||(B_\tau * X) e_\tau|| \right) \leq \tilde{o} \left(\left| \left\langle \nabla \ell(\theta_\tau), \dot{\theta}_\tau \right\rangle \right| \right)$$
 by (71).

If, on the other hand, (18) takes place then $\bar{\eta}_z = \eta_z$ and we bound $||\delta_t^w||$ trivially as

$$||\delta_t^w|| \le O\left((||B_\tau * X|| + ||B_t * X||) \cdot ||e_\tau||\right) \le \widetilde{O}\left(S^{1/2}\left(1 + \frac{m^{1/2}}{n^{1/2}}\right)||e_\tau||\right)$$

and then

$$\left| \langle \delta_t^w, \dot{\theta}_\tau \rangle \right| \leq \widetilde{O} \left(\eta_w S^{1/2} \left(1 + \frac{m^{1/2}}{n^{1/2}} \right) ||e_\tau|| \cdot ||(B_\tau * X) e_\tau|| \right)$$

$$\stackrel{(66)}{\leq} \widetilde{O} \left(\eta_w ||e_\tau||^2 S \left(1 + \frac{m}{n} \right) \right) \stackrel{(18)}{\leq} \widetilde{o} (\eta_z ||e_\tau||^2 S) \stackrel{(71)}{\leq} \widetilde{o} \left(\left| \langle \nabla \ell(\theta_\tau), \dot{\theta}_\tau \rangle \right| \right).$$

This completes the analysis of the contribution of δ_t^w in (41).

The analysis of δ_t^z is analogous (and easier):

$$\begin{split} ||\delta_{t}^{z}|| &\leq ||F_{\tau}(e_{t} - e_{\tau})|| + ||(F_{t} - F_{\tau})e_{t}||_{\mathbf{F}} \leq ||F_{\tau}|| \cdot ||e_{t} - e_{\tau}|| + ||F_{t} - F_{\tau}||_{\mathbf{F}} \cdot ||e_{\tau}|| \\ &\leq \widetilde{o}\left(\frac{m^{1/2}}{S^{1/2}}||e_{\tau}|| + \eta_{w}m^{1/2}S\left(1 + \frac{m}{n}\right)||e_{\tau}||\right) \\ &\leq \widetilde{o}\left(\frac{m^{1/2}}{S^{1/2}}||e_{\tau}||\right) \end{split}$$

and then

$$\left| \left\langle \delta_t^z, \dot{\theta}_\tau \right\rangle \right| \le \widetilde{o} \left(\eta_z \frac{m^{1/2}}{S^{1/2}} ||e_\tau|| \cdot ||F_\tau e_\tau|| \right).$$

If (17) or (18) takes place then $\bar{\eta}_z = \eta_z$ and we are done by (71). If, on the other hand, (16) takes place then we can continue this estimate

$$\left| \left\langle \delta_t^z, \dot{\theta}_\tau \right\rangle \right| \le \tilde{o} \left(\eta_w \frac{S^{1/2}}{m^{1/2}} ||e_\tau|| \cdot (m^{1/2} S^{1/2}) \cdot ||e_\tau|| \right) = \tilde{o}(\eta_w S ||e_\tau||^2),$$

and we are done again since $\bar{\eta}_w = \eta_w$.

This completes the proof of (41).

F. Not too many activation changes

In this section we prove (69); as noted in the introduction, our proof uses (a discretized version of) the beautiful invariant discovered in [ACH18, DHL18].

Let us first remind the set-up: we are given a non-integer t>0 such that for all s< t we have all the facts and inequalities proven in Section 3.2 as well as Appendix D. By continuity, we also have the seed inequalities (32)-(35) for our chosen t as well. The bound (41) is not guaranteed for this t; in fact, this is exactly what we are proving. But it is used only in the integral form which means that we still have all conclusions from Section 3.2 and Appendix D for our chosen t as well. We are specifically interested in (60): $||(A_t - A_0) * X|| \le \tilde{o}(S^{1/2})$.

Let us now start the argument. First, we have

$$||(B_t - B_\tau) * X|| \le ||(\operatorname{diag}(z_t - z_\tau) A_\tau) * X|| + ||(\operatorname{diag}(z_t)(A_t - A_\tau)) * X||.$$

The estimate of the first term is immediate:

$$||(\operatorname{diag}(z_t - z_\tau) A_\tau) * X|| \stackrel{(10)}{\leq} ||z_t - z_\tau|| \cdot ||X|| \stackrel{(20)}{\leq} \widetilde{O}\left(||z_t - z_\tau|| \left(1 + \frac{m^{1/2}}{n^{1/2}}\right)\right).$$

The upper bound on $||z_t - z_\tau||$ is obtained via a computation completely analogous to (67):

$$||z_t - z_\tau|| = \eta_z(t - \tau)||\nabla^z \ell(\theta_\tau)|| \le \eta_z ||F_\tau|| \cdot ||e_\tau|| \le \tilde{O}(\eta_z m S);$$
 (72)

note for the record (we will need it later) that (67) and (72) also hold for $t \mapsto s$, $\tau \mapsto \lfloor s \rfloor$, for an arbitrary s < t. Hence

$$\left|\left|\left(\operatorname{diag}(z_t - z_\tau)A_\tau\right) * X\right|\right| \le \widetilde{O}\left(\eta_z m S\left(1 + \frac{m^{1/2}}{n^{1/2}}\right)\right) \stackrel{(15)}{\le} \widetilde{O}(1).$$

It remains to handle the second term, that is $||(\operatorname{diag}(z_t)(A_t - A_\tau)) * X||$. We can identify yet another contribution that we already know how to handle:

$$||(\operatorname{diag}(z_{0})(A_{t} - A_{\tau})) * X|| \overset{(9)}{\leq} ||z_{0}||_{\infty} \cdot ||(A_{t} - A_{\tau}) * X||$$

$$\overset{(23)}{\leq} \widetilde{O}(||(A_{t} - A_{0}) * X|| + ||(A_{\tau} - A_{0}) * X||)$$

$$\overset{(60)}{\leq} \widetilde{o}(S^{1/2}).$$

Thus, all that remains to show is

$$||(\operatorname{diag}(z_t - z_0)(A_t - A_\tau)) * X|| \le \tilde{o}(S^{1/2}),$$
 (73)

and the difficulty is that we do not have good enough bound on $||z_t - z_0||_{\infty}$. In order to circumvent this difficulty, let

$$D \stackrel{\text{def}}{=} \sup(A_t - A_\tau);$$

$$D_\nu \stackrel{\text{def}}{=} \{ j \in [m] \mid (\nu, j) \in D \}$$

and

$$\Gamma \stackrel{\text{def}}{=} \left\{ \nu \in [S] \, | \, |D_{\nu}| \le n^* \right\},\,$$

where $n^* \leq \tilde{O}(n)$ is as in (21). We split $(\operatorname{diag}(z_t - z_0)(A_t - A_\tau)) * X$ in two parts, $(\operatorname{diag}(z_t - z_0)(A_t - A_\tau))_{\Gamma} * X$ and $(\operatorname{diag}(z_t - z_0)(A_t - A_\tau))_{co-\Gamma} * X$ and bound their norms separately.

The first one follows from what we already know:

$$||(\operatorname{diag}(z_{t} - z_{0})(A_{t} - A_{\tau}))_{\Gamma} * X|| \overset{(10)}{\leq} ||z_{t} - z_{0}|| \cdot \max_{\nu \in \Gamma} ||X \operatorname{diag}((A_{t} - A_{\tau})_{\nu})||$$

$$\overset{(65)}{\leq} \tilde{o}\left(S^{1/2} \cdot \max_{\nu \in \Gamma} ||X^{D_{\nu}}||\right) \overset{(20)}{\leq} \tilde{o}(S^{1/2})$$

(recall that $|D_{\nu}| \leq \tilde{O}(n)$ for $\nu \in \Gamma$ by our choice of Γ).

We bound the last remaining term as

$$||(\operatorname{diag}(z_{t}-z_{0})(A_{t}-A_{\tau}))_{co-\Gamma} * X|| \overset{(9)}{\leq} ||(z_{t}-z_{0})_{co-\Gamma}||_{\infty} \cdot ||(A_{t}-A_{\tau}) * X||$$

$$\overset{(60)}{\leq} \widetilde{o}\left(S^{1/2}||(z_{t}-z_{0})_{co-\Gamma}||_{\infty}\right).$$

So it only remains to prove that $||(z_t - z_0)_{co-\Gamma}||_{\infty} \leq \tilde{O}(1)$ which, given (23), amounts to proving

$$||(z_t)_{co-\Gamma}||_{\infty} \le \tilde{O}(1). \tag{74}$$

In words: for every particular neuron $\nu \in [S]$ we need to prove the dichotomy: either its weight on the second layer is small or it changes activation only on a small number of input data. This is where [ACH18, DHL18] steps in.

Let us fix a neuron $\nu \in co - \Gamma$. Pick an arbitrary $\widetilde{D}_{\nu} \subseteq D_{\nu}$ of cardinality exactly n^* . Then we have the following chain of inequalities:

$$||(W_{\tau})_{\nu}|| \leq \frac{m}{n}||(W_{\tau})_{\nu}X^{\widetilde{D}_{\nu}}|| \leq \frac{m}{n}||(W_{\tau} - W_{t})_{\nu}X^{\widetilde{D}_{\nu}}|| \leq \widetilde{O}\left(\frac{m}{n}||(W_{\tau} - W_{t})_{\nu}||\right)$$

$$\leq \widetilde{O}\left(\frac{m}{n}||W_{\tau} - W_{t}||_{F}\right) \leq \widetilde{O}\left(\frac{m}{n}\eta_{w}m^{1/2}S\left(1 + \frac{m^{1/2}}{n^{1/2}}\right)\right) \leq \widetilde{O}(1).$$

The first inequality holds since $\sigma_{\min}\left(\left(X^{\widetilde{D}_{\nu}}\right)^{\top}\right) \geq \frac{n}{m}$ by (21). The second inequality holds since for any $j \in \widetilde{D}_{\nu}$, $(W_{\tau})_{\nu j}$ and $(W_{t})_{\nu j}$ have opposite signs. The third inequality is true since $||X^{\widetilde{D}_{\nu}}|| \leq \widetilde{O}(1)$ by (20), the fourth is obvious and the fifth is (67). Finally, the sixth inequality follows from $\eta_{w} \leq \widetilde{o}\left(\frac{n}{mS^{2}}\right)$ in (15). Hence by (22) we have

$$||(W_0)_{\nu}||^2 - ||(W_{\tau})_{\nu}||^2 \ge \widetilde{\Omega}(n). \tag{75}$$

For any s < t such that θ_s is regular, we have

$$\frac{d}{ds}(z_s)_{\nu}^2 = 2(z_s)_{\nu} \cdot (\dot{z}_s)_{\nu} = -2\eta_z(z_s)_{\nu} (F_{\lfloor s \rfloor})_{\nu} e_{\lfloor s \rfloor}$$

and

$$\frac{d}{ds}||(W_s)_{\nu}||^2 = 2\langle (W_s)_{\nu}, (\dot{W}_s)_{\nu} \rangle = -2\eta_w(z_{\lfloor s \rfloor})_{\nu}(W_s)_{\nu}X \left(e_{\lfloor s \rfloor} \circ (A_{\lfloor s \rfloor})_{\nu}^{\top}\right).$$

We note that

$$(W_{\lfloor s\rfloor})_{\nu}X\left(e_{\lfloor s\rfloor}\circ (A_{\lfloor s\rfloor})_{\nu}^{\top}\right)=(F_{\lfloor s\rfloor})_{\nu}e_{\lfloor s\rfloor}.$$

Hence if we let

$$R_s \stackrel{\text{def}}{=} \eta_w(z_s)_{\nu}^2 - \eta_z ||(W_s)_{\nu}||^2$$

then

$$\begin{split} ||\dot{R}_{s}|| &\leq \tilde{O}\left(\eta_{w}\eta_{z}\left(|(z_{\lfloor s\rfloor})_{\nu}|\cdot\left|(W_{\lfloor s\rfloor}-W_{s})_{\nu}X\left(e_{\lfloor s\rfloor}\circ(A_{\lfloor s\rfloor})_{\nu}^{\top}\right)\right|\right. \\ &\left. + |(z_{\lfloor s\rfloor}-z_{s})_{\nu}|\cdot|(F_{\lfloor s\rfloor})_{\nu}e_{\lfloor s\rfloor}|\right)\right) \\ &\leq \tilde{O}\left(\eta_{w}\eta_{z}||e_{\lfloor s\rfloor}||\cdot\left(||z_{\lfloor s\rfloor}||\cdot||W_{\lfloor s\rfloor}-W_{s}||\cdot||X||+||z_{\lfloor s\rfloor}-z_{s}||\cdot||F_{\lfloor s\rfloor}||\right)\right) \\ &\leq \tilde{O}\left(\eta_{w}\eta_{z}||e_{\lfloor s\rfloor}||\cdot m^{1/2}S^{3/2}\left(\eta_{w}\left(1+\frac{m}{n}\right)+\eta_{z}m\right)\right), \end{split}$$

where for the last inequality we used the bounds from Appendix E, as well as (72). Let us now integrate this.

In the case (16) we have $\eta_w \left(1 + \frac{m}{n}\right) + \eta_z m \leq \tilde{o}(\eta_w S)$, and applying (46), $|R_\tau - R_0| \leq \int_0^\tau ||\dot{R}_s|| ds \leq \tilde{O}(\eta_w \eta_z m S^2)$.

In the two remaining cases (17), (18) we have (see Remark 5) $\frac{\eta_z}{\eta_w} \geq \tilde{\omega} \left(\frac{m^2}{nS}\right)$ and then $\eta_w \left(1 + \frac{m}{n}\right) + \eta_z m \leq \tilde{O}\left(\eta_z \left(\frac{nS}{m^2} + \frac{S}{m} + m\right)\right)$. Hence (48) gives us $|R_\tau - R_0| \leq \tilde{O}\left(\eta_w \eta_z mS\left(\frac{nS}{m^2} + \frac{S}{m} + m\right)\right)$. Thus in either case we have the bound

$$|R_{\tau} - R_0| \le \widetilde{O}\left(\eta_w \eta_z S^2\left(m + \frac{n}{m}\right)\right) \stackrel{(15)}{\le} \widetilde{o}(n\eta_z). \tag{76}$$

On the other hand,

$$R_{\tau} - R_{0} = \eta_{z}(||(W_{0})_{\nu}||^{2} - ||(W_{\tau})_{\nu}||^{2}) - \eta_{w}((z_{0})_{\nu}^{2} - (z_{\tau})_{\nu}^{2})$$

$$\stackrel{(75)}{\geq} \widetilde{\Omega}(n\eta_{z}) - \eta_{w}((z_{0})_{\nu}^{2} - (z_{\tau})_{\nu}^{2}).$$

$$(77)$$

Comparing (76) and (77), we conclude¹¹ that $||(z_{\tau})_{\nu}|| \leq ||(z_{0})_{\nu}|| \leq \widetilde{O}(1)$. This concludes the proof of (74), (73) and (69).

¹¹If $\eta_w = 0$ then $D = \emptyset$ and (69) is trivial.