

Neural Networks Efficiently Learn Low-Dimensional Representations with SGD

Alireza Mousavi-Hosseini

University of Toronto & Vector Institute

Sejun Park

Korea University

Manuela Girotti

Saint Mary's University

Ioannis Mitliagkas

Université de Montréal & Mila Institute

Murat A. Erdogdu

University of Toronto & Vector Institute

MOUSAVI@CS.TORONTO.EDU

SEJUN.PARK000@GMAIL.COM

MANUELA.GIROTTI@SMU.CA

IOANNIS@IRO.UMONTREAL.CA

ERDOGDU@CS.TORONTO.EDU

Abstract

We study the problem of training a two-layer neural network (NN) of arbitrary width using stochastic gradient descent (SGD) where the input $\mathbf{x} \in \mathbb{R}^d$ is Gaussian and the target $y \in \mathbb{R}$ follows a multiple-index model, i.e., $y = g(\langle \mathbf{u}_1, \mathbf{x} \rangle, \dots, \langle \mathbf{u}_k, \mathbf{x} \rangle)$ with a noisy link function g . We prove that the first-layer weights of the NN converge to the k -dimensional *principal subspace* spanned by the vectors $\mathbf{u}_1, \dots, \mathbf{u}_k$ of the true model, when online SGD with weight decay is used for training. This phenomenon has several important consequences when $k \ll d$. First, by employing uniform convergence on this smaller subspace, we establish a generalization error bound of $\mathcal{O}(\sqrt{kd/T})$ after T iterations of SGD, which is independent of the width of the NN. We further demonstrate that, SGD-trained ReLU NNs can learn a single-index target of the form $y = f(\langle \mathbf{u}, \mathbf{x} \rangle) + \epsilon$ by recovering the principal direction, with a sample complexity linear in d (up to log factors), where f is a monotonic function with at most polynomial growth, and ϵ is the noise. This is in contrast to the known $d^{\Omega(p)}$ sample requirement to learn any degree p polynomial in the kernel regime, and it shows that NNs trained with SGD can outperform the neural tangent kernel at initialization.

1. Introduction

The task of learning an unknown statistical (teacher) model using data is fundamental in many areas of learning theory. There has been a considerable amount of research dedicated to this task, especially when the trained (student) model is a neural network (NN), providing precise and non-asymptotic guarantees in various settings [1, 2, 7, 8, 23, 34, 62, 66, 74, 75]. As evident from these works, explaining the remarkable learning capabilities of NNs requires arguments beyond the classical learning theory [72].

The connection among NNs and kernel methods has been particularly useful towards this expedition [22, 38]. In particular, a two-layer NN with randomly initialized and untrained weights is an example of a random features model [58], and regression on the second layer captures several interesting phenomena that NNs exhibit in practice [45, 47], e.g. *cusp* in the learning curve. However, NNs also inherit favorable characteristics from the optimization procedure [3, 31, 44, 59, 69],

which cannot be captured by associating NNs with regression on random features. Indeed, recent works have established a separation between NNs and kernel methods, relying on the emergence of representation learning as a consequence of gradient-based training [1, 8, 9, 23], which often exhibits a natural bias towards low-complexity models.

A theme that has emerged repeatedly in modern learning theory is the implicit regularization effect provided by the training dynamics [51]. Specifically, [64] has inspired an abundance of recent works focusing on the implicit bias of gradient descent favoring, in some sense, *low-complexity* models, e.g. by achieving min-norm and/or max-margin solutions despite the lack of any explicit regularization [21, 33, 36, 39, 43, 56]. However, these works mainly consider linear models or unrealistically wide NNs, and the notion of reduced complexity as well as its implications on generalization varies. A concrete example in this domain is *compressibility* and its connection to generalization [5, 65]. Indeed, when a trained NN can be compressed into a smaller NN with similar prediction behavior, the resulting models exhibit similar generalization performance, while the latter is classically linked to better generalization.

In this paper, we demonstrate the emergence of low-complexity structures during the training procedure. More specifically, we consider training a two-layer student NN with arbitrary width m where the input $\mathbf{x} \in \mathbb{R}^d$ is Gaussian and the target $y \in \mathbb{R}$ follows a multiple-index teacher model, i.e. $y = g(\langle \mathbf{u}_1, \mathbf{x} \rangle, \dots, \langle \mathbf{u}_k, \mathbf{x} \rangle; \epsilon)$ with a link function g and a noise ϵ independent of the input. In this setting, we prove that the first-layer weights trained by online stochastic gradient descent (SGD) with weight decay converge to the k -dimensional subspace $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$, which we refer to as the *principal subspace*. Our primary focus is the case where the target values depend only on a few important directions along the input, i.e. $k \ll d$, which induces a low-dimensional structure on the SGD-trained first-layer weights, whose impact on generalization is profound. We summarize our contributions as follows.

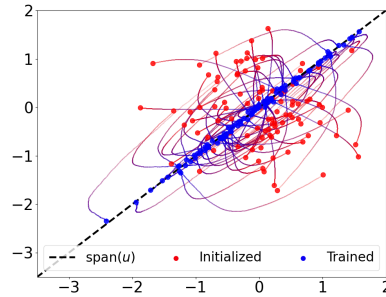


Figure 1: Two-layer ReLU network with $m = 1000$, $d = 2$ is trained to recover a tanh single-index model via SGD with weight decay. Initial neurons (red) converge to the principal subspace. 10% of student neurons are visualized.

- We show in Theorem 2 that NNs learn low-dimensional representations by proving that the iterates of online SGD on the first layer of a two-layer NN with width m converge to $\sqrt{m}\epsilon$ neighborhood of the principal subspace after $\mathcal{O}(d/\epsilon^2)$ iterations, with high probability. The error tolerance of $\sqrt{m}\epsilon$ is sufficient to guarantee that the risk of SGD iterates and that of its orthogonal projection to the principal subspace are within $\mathcal{O}(\epsilon)$ distance.
- We demonstrate the impact of learning low-dimensional representations with two applications.
 - For a single-index target $y = f(\langle \mathbf{u}, \mathbf{x} \rangle) + \epsilon$ with certain link functions f , we prove in Theorem 4 that ReLU networks of width m can learn this target with an SGD-based procedure (Algorithm 1) with number of samples T and an excess risk estimate of $\tilde{\mathcal{O}}(\sqrt{d/T} + 1/m)$, with high probability (see the illustration in Figure 1). In particular, the sample complexity is (almost) linear in d , even when f is a monotonic polynomial of any (fixed) odd degree p .
 - Based on a uniform convergence argument on the principal subspace, we prove in Theorem 5 that T iterations of SGD will produce a model with generalization error of $\mathcal{O}(\sqrt{kd/T})$, with

high probability. Remarkably, this rate is independent of the width m of the NN, even in the case $k \asymp d$ where the target is any function of the input, and not necessarily low-dimensional.

The rest of the paper is organized as follows. We discuss the notation and the related work in the remainder of this section. We describe the problem formulation and preliminaries in Section 2. Our main result on SGD is presented in Section 3. We discuss two implications of our main theorem in Section 4, where we provide results on learnability and generalization gap in Sections 4.1 and 4.2 respectively. We finally conclude with a brief discussion in Section 5.

Notation. For a loss function $\ell : \mathbb{R}^2 \rightarrow \mathbb{R}$, let $\partial_i \ell$ and $\partial_{ij}^2 \ell$ denote its partial derivatives with respect to i th and j th inputs for $i, j \in \{1, 2\}$. For quantities a and b , $a \lesssim b$ implies $a \leq Cb$ for an absolute constant C , and $a \asymp b$ implies both $a \gtrsim b$ and $a \lesssim b$. Finally, $\text{Unif}(A)$ denotes the uniform distribution over a set A and $\mathcal{N}(0, \mathbf{I}_d)$ denotes the d -dimensional isotropic Gaussian distribution.

1.1. Related work

Training dynamics of NNs. Several works have demonstrated learnability in a special case of teacher-student setting where the teacher model is *similar* to the student NN being trained [16, 42, 73–75]. This setting has also been studied through the lens of loss landscape [61] and optimization over measures [2]. We stress that our results work under misspecification and hold for generic teacher models that are not necessarily NNs with similar architecture to the student. Additional related works on different scaling regimes of wide networks is provided in Appendix A.

Feature learning with multiple-index teacher models. The task of learning a target of an unknown low-dimensional function of the input is fundamental in statistics [41]. Several recent works in learning theory literature have also focused on this problem, with an aim to demonstrate NNs can learn useful feature representations, outperforming kernel methods [12, 32]. In particular, [1] studies the necessary and sufficient conditions for learning with sample complexity linear in d with inputs on the hypercube, in the mean-field limit. Closer to our setting are the recent works [8, 9, 23] which demonstrate a clear separation between NNs and kernel methods, leveraging the effect of representation learning. However, their analysis considers a single (full) gradient step on the first-layer weights followed by training the second-layer parameters. In contrast, in our learnability result, we consider training both layers with SGD, which induces essentially different learning dynamics.

Generalization bounds for SGD. A popular algorithm-dependent approach for studying generalization is through algorithmic stability [14, 15, 28], which has been used to study the generalization behavior of gradient-based methods in various settings [11, 27, 37, 40]. Other approaches include studying the low-dimensional structure of the trajectory [55, 63] or the invariant measure of continuous-time approximations of SGD [17], and employing information-theoretic tools [50]. Among these, [10] also shows that SGD yields compressible networks; however, they assume the mean-field approximation holds and the SGD iterates converge to a heavy-tailed distribution.

2. Preliminaries: Neural Networks and the Principal Subspace

For an input $\mathbf{x} \in \mathbb{R}^d$, we consider training a two-layer neural network (NN) with m neurons

$$\hat{y}(\mathbf{x}; \mathbf{W}, \mathbf{a}, \mathbf{b}) = \sum_{i=1}^m a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle + b_i), \tag{2.1}$$

where σ is the activation function, $\{\mathbf{w}_i\}_{1 \leq i \leq m}$ are the first-layer weights collected in the rows of the matrix $\mathbf{W} \in \mathbb{R}^{m \times d}$, $\mathbf{b} \in \mathbb{R}^m$ is the bias, and $\mathbf{a} \in \mathbb{R}^m$ is the second-layer weights. We assume $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$ and the target is generated from a multiple-index (teacher) model given by

$$y = g(\langle \mathbf{u}_1, \mathbf{x} \rangle, \dots, \langle \mathbf{u}_k, \mathbf{x} \rangle; \epsilon), \quad (2.2)$$

for a weakly differentiable link function $g : \mathbb{R}^{k+1} \rightarrow \mathbb{R}$ and a noise ϵ independent of \mathbf{x} . The Gaussian input is a rather standard assumption in the literature, especially in recent works that consider the student-teacher setup; see e.g. [23, 61, 75]. The multiple-index teacher model (2.2) can encode a broad class of input-output relations through the non-linear link function, including a multi-layer fully-connected NN with arbitrary depth and width and weakly differentiable activations. While our results remain valid regardless of how k and d compare, they are most insightful when $k \ll d$; thus, we specifically consider this regime when interpreting the results. We also collect the teacher weights $\{\mathbf{u}_i\}_{1 \leq i \leq k}$ in the rows of the matrix $\mathbf{U} \in \mathbb{R}^{k \times d}$ and use $y = g(\mathbf{U}\mathbf{x}; \epsilon)$ for simplicity.

For a given loss function $\ell(\hat{y}, y)$, we consider the population and the empirical risks

$$R(\mathbf{W}, \mathbf{a}, \mathbf{b}) = \mathbb{E}[\ell(\hat{y}(\mathbf{x}; \mathbf{W}, \mathbf{a}, \mathbf{b}), y)] \quad \text{and} \quad \hat{R}(\mathbf{W}, \mathbf{a}, \mathbf{b}) = \frac{1}{T} \sum_{t=0}^{T-1} \ell(\hat{y}(\mathbf{x}^{(t)}; \mathbf{W}, \mathbf{a}, \mathbf{b}), y^{(t)}),$$

where the expectation is over the data distribution. Similarly, for some $\tau \geq 1$, the truncated loss is defined as $\ell_\tau(\hat{y}, y) := \ell(\hat{y}, y) \wedge \tau$ with the corresponding risks R_τ and \hat{R}_τ . We define the L_2 -regularized population risk with a penalty parameter $\lambda \geq 0$ as

$$\mathcal{R}_\lambda(\mathbf{W}, \mathbf{a}, \mathbf{b}) := R(\mathbf{W}, \mathbf{a}, \mathbf{b}) + \frac{\lambda}{2} \|\mathbf{W}\|_{\text{F}}^2. \quad (2.3)$$

To minimize (2.3), we use stochastic gradient descent (SGD) over the first-layer weights, where we are interested in the convergence of iterates to the *principal subspace* defined by the teacher weights

$$\mathbb{S}(\mathbf{U}) := \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)^m = \{\mathbf{C}\mathbf{U} : \mathbf{C} \in \mathbb{R}^{m \times k}\}.$$

Notice that the principal subspace satisfies $\mathbb{S}(\mathbf{U}) \subseteq \mathbb{R}^{m \times d}$, and its dimension is mk as opposed to the ambient dimension of md , with any matrix in this subspace having rank at most k . For any vector $\mathbf{v} \in \mathbb{R}^d$, we let \mathbf{v}_\parallel denote the orthogonal projection of \mathbf{v} onto $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ and $\mathbf{v}_\perp := \mathbf{v} - \mathbf{v}_\parallel$. Similarly, for a matrix $\mathbf{W} \in \mathbb{R}^{m \times d}$, we define \mathbf{W}_\parallel and \mathbf{W}_\perp by applying the projection to each row.

The smoothness properties of the activation σ play an important role in our analysis. As such, we consider two scenarios, with different requirements on the loss function.

Assumption 1.A (Smooth activation) *The activation function σ satisfies $|\sigma(z)|, |\sigma'(z)|, |\sigma''(z)| \leq 1$ for all $z \in \mathbb{R}$, the loss is $\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$ for simplicity, and y satisfies $|y| \leq K$ almost surely.*

Assumption 1.B (ReLU activation) *The activation function σ is $\sigma(z) = \max(z, 0)$ for $z \in \mathbb{R}$. The loss satisfies $0 \leq \partial_1^2 \ell(\hat{y}, y) \leq 1$, $|\partial_1 \ell(\hat{y}, y)| \leq 1$, and $|\partial_{12}^2 \ell(\hat{y}, y)| \leq 1$.*

Commonly used activations such as sigmoid and tanh satisfy Assumption 1.A. For ReLU activation in Assumption 1.B, we choose $\sigma'(z) = \mathbf{1}(z \geq 0)$ as its weak derivative. We highlight that Assumption 1.B is satisfied by common Lipschitz and convex loss functions such as the Huber loss

$$\ell_{\text{H}}(\hat{y} - y) := \begin{cases} \frac{1}{2}(\hat{y} - y)^2 & \text{if } |\hat{y} - y| \leq 1 \\ |\hat{y} - y| - \frac{1}{2} & \text{if } |\hat{y} - y| > 1, \end{cases} \quad (2.4)$$

as well as the logistic loss $\ell_{\text{L}}(\hat{y}, y) := \log(1 + e^{-\hat{y}y})$, up to appropriate scaling constants.

3. Convergence of Stochastic Gradient Descent

We now consider stochastic gradient descent (SGD) in the online setting where at each iteration t , we have access to a new data point $(\mathbf{x}^{(t)}, y^{(t)})$ drawn independently of the previous samples from the same distribution. We update the first-layer weights \mathbf{W}^t with a time varying step size η_t and weight decay, according to the update rule

$$\mathbf{W}^{t+1} = (1 - \eta_t \lambda) \mathbf{W}^t - \eta_t \nabla_{\mathbf{W}} \ell(\hat{y}(\mathbf{x}^{(t)}; \mathbf{W}^t, \mathbf{a}, \mathbf{b}), y^{(t)}). \quad (3.1)$$

The above algorithm can be used to minimize the population risk (2.3) in practice [57], even in certain non-convex landscapes [70]. We will use the following initialization for SGD

Assumption 2 (Initialization) For all $1 \leq i, j \leq m$, we initialize the NN weights and biases with $\sqrt{d} W_{ij}^0 \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$, $ma_j^0 \stackrel{\text{iid}}{\sim} \text{Unif}([-1, 1])$, and $b_j^0 \stackrel{\text{iid}}{\sim} \text{Unif}(\{-1, 1\})$.

This initialization is standard in the mean-field regime as it allows feature learning. However, we only use it to simplify the exposition. Indeed, we can initialize \mathbf{W} and \mathbf{a} with any scheme that guarantees $\|\mathbf{W}\|_F \lesssim \sqrt{m}$ and $\|\mathbf{a}\|_\infty \lesssim m^{-1}$ with high probability. Further, initialization of \mathbf{b} mostly matters in the analysis of ReLU activation.

Next, we show that the population gradient admits a certain decomposition which plays a central role in our analysis. For smooth activations, the below result is a remarkable consequence of Stein's lemma, which provides a certain alignment between the true statistical model (teacher) and the model being trained (student), which has profound impact on the learning dynamics. We generalize this result for ReLU through a sequence of smooth approximations (see Appendix C.1 for details).

Lemma 1 Under Assumption 1.A or 1.B, the gradient of the population risk can be written as

$$\nabla_{\mathbf{W}} \mathcal{R}_\lambda(\mathbf{W}) = (\mathcal{H}(\mathbf{W}) + \lambda \mathbf{I}_m) \mathbf{W} + \mathcal{D}(\mathbf{W}) \mathbf{U}, \quad (3.2)$$

for some $\mathcal{H}(\mathbf{W}) \in \mathbb{R}^{m \times d}$ and $\mathcal{D}(\mathbf{W}) \in \mathbb{R}^{k \times d}$ (with explicit forms provided in Appendix C.1).

The following result, proved in Appendix D via a recursion on the moment generating function of $\|\mathbf{W}_\perp^t\|_F$, demonstrates the algorithmic implications of Lemma 1, and shows that the iterates converges to the principal subspace.

Theorem 2 Consider running T SGD iterations (3.1) when the activation and loss either satisfy Assumption 1.A or Assumption 1.B, and the initialization satisfies Assumption 2. Let $\zeta := \mathbb{E}[|y|] + 1$ under Assumption 1.A and $\zeta := 2\sqrt{2/e\pi}$ under Assumption 1.B. Choose the decreasing step size $\eta_t = m \frac{2(t+t^*)+1}{\gamma(t+t^*+1)^2}$, $\tilde{\lambda} \geq \gamma + \zeta$, and $t^* \asymp \frac{\tilde{\lambda}}{\gamma}$ for any $\gamma > 0$. Then, for $\lambda = \frac{\tilde{\lambda}}{m}$, with probability at least $1 - \delta$,

$$\frac{\|\mathbf{W}_\perp^T\|_F}{\sqrt{m}} \lesssim \sqrt{\frac{d + \log(1/\delta)}{\gamma^2 T}}, \quad (3.3)$$

whenever $m \gtrsim \log(1/\delta)$ and $T \gtrsim \frac{\tilde{\lambda}^2}{d + \log(1/\delta)}$.

Remark 3 The above result provides a convergence rate of $\mathcal{O}(\sqrt{d/T})$ for the normalized quantity $\|\mathbf{W}_\perp^T\|_F / \sqrt{m}$. We justify this exposition by noting that since the risk is locally Lipschitz with a constant $\mathcal{O}(1/\sqrt{m})$, the condition $\|\mathbf{W}_\perp^T\|_F \lesssim \sqrt{m}\varepsilon$ is sufficient to guarantee that the risk of \mathbf{W}^T and that of its orthogonal projection to the principal subspace \mathbf{W}_\parallel^T are within $\mathcal{O}(\varepsilon)$ distance.

The above result states that, with a number of samples linear in the input dimension d , SGD is able to learn (approximately) low-dimensional weights, exhibiting an implicit bias towards low-complexity models. The interplay between two forces is in effect here. The most important one is the linear relationship between the first-layer weights and the input in both student and teacher models together with the input distribution. The alignment described in Lemma 1 yields *sparsified* weights in a basis defined by the teacher network, effectively reducing the input dimension from d to k . The second force is the explicit L_2 -regularization. We emphasize that L_2 -regularization does not play the main role in this sparsification; even though it may provide shrinkage to zero, L_2 penalty will in general produce non-sparse solutions. However, it is still required as the explicit L_2 regularization ensures that SGD avoids critical points outside of the principal subspace. We refer to Section 5 for further discussion.

Notice that while Theorem 2 does not have any implications on the convergence behavior of the orthogonal projection \mathbf{W}_\parallel , in the next section, we show that the implied low-dimensional structure is sufficient to provide guarantees on the generalization error and learnability of SGD. Moreover, we showcase by a non-convex example in Appendix F that the regularization in this Theorem does not in general imply (strong) convexity in the population landscape.

4. Implications of Low-Dimensionality

4.1. Learning Single-Index Targets

An essential characteristic of NNs is their ability to learn useful representations, which allows them to adapt to the underlying misspecified statistical model. Although this fundamental property has been the guiding principle in all empirical studies, it was mathematically proven only recently for gradient-based training [1, 8, 9, 23, 29]; see also a survey of prior works in [46]. Our results in the previous section are in the same spirit, establishing the convergence of SGD to the principal subspace which is indeed a span of useful directions associated with the target function being learned. As such, we leverage the learned low-dimensional representations to demonstrate that SGD is capable of learning a target function of the form $y = f(\langle \mathbf{u}, \mathbf{x} \rangle) + \epsilon$ with a number of samples linear in d (up to logarithmic factors). For simplicity, we work with the Huber loss below; however, our analysis can accommodate any Lipschitz and convex loss at the expense of a more detailed analysis.

Algorithm 1 Training a two-layer ReLU network with SGD.

Input: $\mathbf{a}^0, \mathbf{b}^0 \in \mathbb{R}^m$, $\mathbf{W}^0 \in \mathbb{R}^{m \times d}$, $\{(\mathbf{x}^{(t)}, y^{(t)})\}_{0 \leq t \leq T-1}$, $(\eta_t)_{t \geq 0}$, $(\eta'_t)_{t \geq 0}$, $\lambda, \lambda', \Delta$.

- 1: **for** $t = 0, \dots, T - 1$ **do**
- 2: $\mathbf{W}^{t+1} = (1 - \eta_t \lambda) \mathbf{W}^t - \eta_t \nabla_{\mathbf{W}} \ell(\hat{y}(\mathbf{x}^{(t)}; \mathbf{W}^t, \mathbf{a}^0, \mathbf{b}^0), y^{(t)})$.
- 3: **end for**
- 4: Let $b_j \sim \text{Unif}(-\Delta, \Delta)$ for $1 \leq j \leq m$.
- 5: **for** $t = 0, \dots, T' - 1$ **do**
- 6: Sample $i_t \sim \text{Unif}\{0, \dots, T - 1\}$.
- 7: $\mathbf{a}^{t+1} = (1 - \eta'_t \lambda') \mathbf{a}^t - \eta'_t \nabla_{\mathbf{a}} \ell(\hat{y}(\mathbf{x}^{(i_t)}; \mathbf{W}^T, \mathbf{a}^t, \mathbf{b}), y^{(i_t)})$
- 8: **end for**
- 9: **return** $(\mathbf{W}^T, \mathbf{a}^{T'}, \mathbf{b})$.

The following Theorem states the learning guarantee of Algorithm 1, and is proved in Appendix E.2 via recovering the principal direction according to Theorem 2, along with the generalization gap of Theorem 5, and a universal approximation argument similar to [23]. We highlight that the recent works [8, 9, 23] perform only one gradient step on the first layer weights, whereas in Algorithm 1, we train the entire NN with SGD.

Theorem 4 *Suppose that the data is from a single-index model $y = f(\langle \mathbf{u}, \mathbf{x} \rangle) + \epsilon$ with a monotone differentiable f and 1-sub-Gaussian noise ϵ , and Assumption 1.B holds. Further, let $\|\mathbf{u}\|_2 = 1$, $|f(0)| < 1$, and consider the Huber loss (2.4) for simplicity. Consider running Algorithm 1 with the initialization $0 < a_j^0 = a \lesssim 1/m$, $0 < b_j^0 = b \lesssim 1$, and $\mathbf{w}_j^0 = \mathbf{w}^0 \sim \mathcal{N}(0, \frac{1}{d}\mathbf{I}_d)$ for all j with the hyper-parameters $\lambda = \frac{\tilde{\lambda}}{m} = \frac{\gamma}{m} + \frac{2a}{b} \sqrt{\frac{2}{e\pi}}$ for any $\gamma \asymp 1$, $\eta_t = m \frac{2(t^*+t)+1}{\gamma(t^*+t+1)^2}$ with $t^* \asymp \gamma^{-1}$, $\eta'_t = \frac{2t+1}{\lambda'(t+1)^2}$, and $\Delta \asymp \sqrt{\log(T/\delta)}$. Then, for $T \gtrsim (d + \log(\frac{1}{\delta})) \vee (\frac{\tilde{\lambda}}{\gamma d} \log(\frac{m}{\delta}))$, some $\lambda' > 0$ (see (E.7)), and sufficiently large T' (see (E.8)), with probability at least $1 - \delta$,*

$$R_\tau(\mathbf{W}^T, \mathbf{a}^{T'}, \mathbf{b}) - \mathbb{E}[\ell_H(\epsilon)] \lesssim \Delta_*^2 \left\{ \sqrt{\frac{\log(1/\delta)}{m}} + \sqrt{\frac{d + \log(1/\delta)}{T}} \right\}, \quad (4.1)$$

where $\Delta_* := \Delta \sup_{|z| \lesssim \Delta} |f''(z)|$ which is $\text{poly}(\log(T/\delta))$ when f'' has at most polynomial growth.

Notice that T determines the sample complexity of running Algorithm 1 as it only requires T data samples. As such, the above result implies that a ReLU NN trained with SGD can learn any monotone polynomial with a sample complexity linear in the input dimension d , up to logarithmic factors. In comparison, [23] considers training the first-layer weights with one gradient descent step with a carefully chosen weight decay, and obtains a sample complexity of d^2 to learn any unknown degree p polynomial multiple-index target. Note that learning any degree p polynomial using rotationally invariant kernels requires $d^{\Omega(p)}$ samples for a variety of input distributions including isotropic Gaussian [24]; thus, our result shows that SGD is able to efficiently learn a target function where kernel methods cannot. It is worth emphasizing that in Theorem 4, the width of the network m grows with $\text{poly} \log(T/\delta)$, which is in contrast to the neural tangent kernel (NTK) regime where m is required to grow with $\text{poly}(T)$.

On a separate note, [18] proposes a method that can train NNs to learn a polynomial of a few directions with a number of samples linear in the input dimension d ; yet, the proposed algorithm is not a simple variant of SGD and requires a non-trivial *warm-start* initialization. We refer to [26] for a review of classical literature on learning single and multiple-index models. Specifically, they provide a procedure with sample complexity at least linear in d to recover the principal direction of the single-index model with Gaussian inputs.

4.2. Generalization Gap

For a given learning algorithm, the gap between its empirical and population risks is termed as the *generalization gap* (not to be confused with excess risk), and establishing convergence estimates for this quantity is a fundamental problem in learning theory. Classical results rely on uniform convergence over the feasible domain containing the weights; thus, they apply to any learning algorithm including SGD [52]. However, these bounds often diverge with the width of the NN, yielding vacuous estimates in the overparameterized regime [72]. To alleviate this, recent works considered establishing estimates for a specific learning algorithm; see e.g. [37, 55, 64, 71].

Here, we are interested in deriving an estimate for the generalization gap over the SGD-trained first-layer weights, which holds uniformly over the second layer weights and biases. More specifically, we study, after T iterations of SGD (3.1) initialized with $(\mathbf{W}^0, \mathbf{a}^0, \mathbf{b}^0)$, the following quantity $\mathcal{E}(\mathbf{W}^T) := \sup_S R_\tau(\mathbf{W}^T, \mathbf{a}, \mathbf{b}) - \hat{R}_\tau(\mathbf{W}^T, \mathbf{a}, \mathbf{b})$ with $S := \{\mathbf{a}, \mathbf{b} \in \mathbb{R}^m : \|\mathbf{a}\|_2 \leq \frac{r_a}{\sqrt{m}}, \|\mathbf{b}\|_\infty \leq r_b\}$,

where the scaling ensures $\hat{y} = \mathcal{O}(1)$ when $\|\mathbf{w}_j\|_2 \asymp 1$, which is the setting considered in Theorem 4. We state the following bound on $\mathcal{E}(\mathbf{W}^T)$; the proof is provided in Appendix E.1, and it is based on a covering argument over the smaller dimensional principal subspace implied by Theorem 2.

Theorem 5 *Consider the setting of Theorem 2. For any $\delta > 0$, if $T \gtrsim (d + \log(1/\delta)) \vee (\frac{\tilde{\lambda}}{\gamma d} \log(m/\delta))$, then with probability at least $1 - \delta$,*

$$\mathcal{E}(\mathbf{W}^T) \lesssim \tau r_a \left\{ \sqrt{\frac{d + \log(1/\delta)}{\gamma^2 T}} + (r_b + \tilde{\lambda}^{-1}) \sqrt{\frac{dk}{T}} \right\}, \quad (4.2)$$

The above bound is independent of the width m of the NN, and only grows with the dimension of the input space d and that of the principal subspace k ; thus, producing non-vacuous estimates in the overparametrized regime where m is large. Further, the bound is stable in the number of SGD iterations T , that is, it converges to zero as $T \rightarrow \infty$. We remark that generalization bounds that rely on algorithmic stability are optimal for strongly convex objectives [37]; yet, they lead to unstable diverging bounds in non-convex settings as $T \rightarrow \infty$. As such, these techniques often require *early stopping*, which is clearly not needed in our result.

5. Conclusion

We studied the dynamics of SGD with weight decay on two-layer NNs, and proved that under a multiple-index teacher model, the first-layer weights converge to the principal subspace, i.e. the span of the weights of the teacher. This phenomenon is of particular interest when the target depends on the input along a few important directions. In this setting, we proved novel generalization bounds for SGD via uniform convergence on the low-dimensional principal subspace. Further, we proved that two-layer ReLU networks can learn a single-index target with a monotone link that has at most polynomial growth, using online SGD, with a number of samples almost linear in d . Thus, as an implication of low-dimensionality, we established a separation between kernel methods and trained NNs where the former suffers from the curse of dimensionality.

Two principal forces are responsible for the emergence of the low-dimensional structure. The main one is the linear interaction between the Gaussian input and the first-layer weights in both student and teacher models. The secondary one is the weight decay which allows SGD to avoid critical points outside of principal subspace. Figure 2 shows the convergence behavior in absence of weight decay. Understanding more precisely the range of λ that implies convergence to the principal subspace, as well as investigating the possibility of learning multiple-index models using this convergence, are left as important directions for future studies.

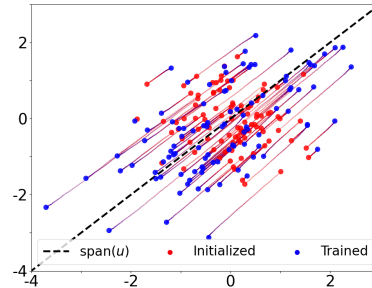


Figure 2: Neurons fail to converge to the principal subspace without weight decay, in the same experimental setup of Figure 1.

References

- [1] Emmanuel Abbe, Enric Boix Adsera, and Theodor Misiakiewicz. The merged-staircase property: a necessary and nearly sufficient condition for sgd learning of sparse functions on two-layer neural networks. In *Conference on Learning Theory*, 2022.
- [2] Shunta Akiyama and Taiji Suzuki. On Learnability via Gradient Method for Two-Layer ReLU Neural Networks in Teacher-Student Setting. In *International Conference on Machine Learning*, 2021.
- [3] Zeyuan Allen-Zhu and Yuanzhi Li. What Can ResNet Learn Efficiently, Going Beyond Kernels? In *Advances in Neural Information Processing Systems*, 2019.
- [4] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A Convergence Theory for Deep Learning via Over-Parameterization. In *International Conference on Machine Learning*, 2019.
- [5] Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger Generalization Bounds for Deep Nets via a Compression Approach. In *International Conference on Machine Learning*, 2018.
- [6] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-Grained Analysis of Optimization and Generalization for Overparameterized Two-Layer Neural Networks. In *International Conference on Machine Learning*, 2019.
- [7] Jimmy Ba, Murat Erdogdu, Taiji Suzuki, Denny Wu, and Tianzong Zhang. Generalization of two-layer neural networks: An asymptotic viewpoint. In *International Conference on Machine Learning*, 2019.
- [8] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation. *arXiv preprint arXiv:2205.01445*, 2022.
- [9] Boaz Barak, Benjamin L Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang. Hidden Progress in Deep Learning: SGD Learns Parities Near the Computational Limit. *arXiv preprint arXiv:2207.08799*, 2022.
- [10] Melih Barsbey, Milad Sefidgaran, Murat A Erdogdu, Gael Richard, and Umut Simsekli. Heavy tails in SGD and compressibility of overparametrized neural networks. In *Advances in Neural Information Processing Systems*, 2021.
- [11] Raef Bassily, Vitaly Feldman, Cristóbal Guzmán, and Kunal Talwar. Stability of stochastic gradient descent on nonsmooth convex losses. In *Advances in Neural Information Processing Systems*, 2020.
- [12] Benedikt Bauer and Michael Kohler. On deep learning as a remedy for the curse of dimensionality in nonparametric regression. *The Annals of Statistics*, 2019.
- [13] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *J. Mach. Learn. Res.*, 22:106–1, 2021.

- [14] Olivier Bousquet and André Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2002.
- [15] Olivier Bousquet, Yegor Klochkov, and Nikita Zhivotovskiy. Sharper bounds for uniformly stable algorithms. In *Conference on Learning Theory*, 2020.
- [16] Alon Brutzkus and Amir Globerson. Globally Optimal Gradient Descent for a ConvNet with Gaussian Inputs. In *International Conference on Machine Learning*, 2017.
- [17] Alexander Camuto, George Deligiannidis, Murat A Erdogdu, Mert Gurbuzbalaban, Umut Simsekli, and Lingjiong Zhu. Fractal structure and generalization properties of stochastic optimization algorithms. In *Advances in Neural Information Processing Systems*, 2021.
- [18] Sitan Chen and Raghu Meka. Learning polynomials in few relevant dimensions. In *Conference on Learning Theory*, 2020.
- [19] Lénaïc Chizat. Mean-field langevin dynamics: Exponential convergence and annealing. *arXiv preprint arXiv:2202.01009*, 2022.
- [20] Lenaic Chizat and Francis Bach. On the Global Convergence of Gradient Descent for Over-parameterized Models using Optimal Transport. In *Advances in Neural Information Processing Systems*, 2018.
- [21] Lénaïc Chizat and Francis Bach. Implicit Bias of Gradient Descent for Wide Two-layer Neural Networks Trained with the Logistic Loss. In *Conference on Learning Theory*, 2020.
- [22] Lenaic Chizat, Edouard Oyallon, and Francis Bach. On Lazy Training in Differentiable Programming. In *Advances in Neural Information Processing Systems*, 2019.
- [23] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural Networks can Learn Representations with Gradient Descent. In *Conference on Learning Theory*, 2022.
- [24] Konstantin Donhauser, Mingqi Wu, and Fanny Yang. How rotational invariance of common kernels prevents generalization in high dimensions. In *International Conference on Machine Learning*, 2021.
- [25] Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient Descent Provably Optimizes Over-parameterized Neural Networks. In *International Conference on Learning Representations*, 2019.
- [26] Rishabh Dudeja and Daniel Hsu. Learning single-index models in gaussian space. In *Conference On Learning Theory*, 2018.
- [27] Tyler Farghly and Patrick Rebeschini. Time-independent generalization bounds for SGLD in non-convex settings. In *Advances in Neural Information Processing Systems*, 2021.
- [28] Vitaly Feldman and Jan Vondrak. Generalization bounds for uniformly stable algorithms. *Advances in Neural Information Processing Systems*, 31, 2018.
- [29] Spencer Frei, Niladri S Chatterji, and Peter L Bartlett. Random feature amplification: Feature learning and generalization in neural networks. *arXiv preprint arXiv:2202.07626*, 2022.

- [30] Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy learning in deep neural networks: an empirical study. *arXiv preprint arXiv:1906.08034*, 2020.
- [31] B. Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Limitations of Lazy Training of Two-layers Neural Networks. In *Advances in Neural Information Processing Systems*, 2019.
- [32] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When Do Neural Networks Outperform Kernel Methods? In *Advances in Neural Information Processing Systems*, 2020.
- [33] Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient dynamics in linear neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- [34] Sebastian Goldt, Madhu Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborová. Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup. In *Advances in Neural Information Processing Systems*, 2019.
- [35] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik. SGD: General analysis and improved rates. In *International Conference on Machine Learning*, 2019.
- [36] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit Bias of Gradient Descent on Linear Convolutional Networks. In *Advances in Neural Information Processing Systems*, 2018.
- [37] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning*, 2016.
- [38] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural Tangent Kernel: Convergence and Generalization in Neural Networks. In *Advances in Neural Information Processing Systems*, 2018.
- [39] Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning Theory*, 2019.
- [40] Leo Kozachkov, Patrick M Wensing, and Jean-Jacques Slotine. Generalization in supervised learning through riemannian contraction. *arXiv preprint arXiv:2201.06656*, 2022.
- [41] Ker-Chau Li and Naihua Duan. Regression Analysis Under Link Violation. *The Annals of Statistics*, 1989.
- [42] Yuanzhi Li and Yang Yuan. Convergence Analysis of Two-layer Neural Networks with ReLU Activation. In *Advances in Neural Information Processing Systems*, 2017.
- [43] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic Regularization in Overparameterized Matrix Sensing and Neural Networks with Quadratic Activations. In *Conference on Learning Theory*, 2018.

- [44] Yuanzhi Li, Tengyu Ma, and Hongyang R Zhang. Learning over-parametrized two-layer neural networks beyond NTK. In *Conference on Learning Theory*, 2020.
- [45] Cosme Louart, Zhenyu Liao, and Romain Couillet. A random matrix approach to neural networks. *The Annals of Applied Probability*, 2018.
- [46] Eran Malach, Pritish Kamath, Emmanuel Abbe, and Nathan Srebro. Quantifying the benefit of using differentiable learning over tangent kernels. In *International Conference on Machine Learning*, 2021.
- [47] Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*, 2022.
- [48] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory*, 2019.
- [49] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- [50] Gergely Neu, Gintare Karolina Dziugaite, Mahdi Haghifam, and Daniel M. Roy. Information-theoretic generalization bounds for stochastic gradient descent. In *Conference on Learning Theory*, 2021.
- [51] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. *arXiv preprint arXiv:1412.6614*, 2014.
- [52] Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. The role of over-parametrization in generalization of neural networks. In *International Conference on Learning Representations*, 2019.
- [53] Atsushi Nitanda, Denny Wu, and Taiji Suzuki. Convex analysis of the mean field langevin dynamics. In *International Conference on Artificial Intelligence and Statistics*, pages 9741–9757. PMLR, 2022.
- [54] Samet Oymak and Mahdi Soltanolkotabi. Toward moderate overparameterization: Global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 2020.
- [55] Sejun Park, Umut Simsekli, and Murat A. Erdogdu. Generalization Bounds for Stochastic Gradient Descent via Localized ε -Covers. *arXiv preprint arXiv:2209.08951*, 2022.
- [56] Scott Peme, Loucas Pillaud-Vivien, and Nicolas Flammarion. Implicit Bias of SGD for Diagonal Linear Networks: a Provable Benefit of Stochasticity. In *Advances in Neural Information Processing Systems*, 2021.
- [57] Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992.

- [58] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, 2007.
- [59] Maria Refinetti, Sebastian Goldt, Florent Krzakala, and Lenka Zdeborová. Classifying high-dimensional Gaussian mixtures: Where kernel methods fail and neural networks succeed. In *International Conference on Machine Learning*, 2021.
- [60] Grant M Rotskoff and Eric Vanden-Eijnden. Neural networks as Interacting Particle Systems: Asymptotic convexity of the Loss Landscape and Universal Scaling of the Approximation Error. *arXiv preprint arXiv:1805.00915*, 2018.
- [61] Itay M Safran, Gilad Yehudai, and Ohad Shamir. The Effects of Mild Over-parameterization on the Optimization Landscape of Shallow ReLU Neural Networks. In *Conference on Learning Theory*, 2021.
- [62] Stefano Sarao Mannelli, Eric Vanden-Eijnden, and Lenka Zdeborová. Optimization and generalization of shallow neural networks with quadratic activation functions. In *Advances in Neural Information Processing Systems*, 2020.
- [63] Umut Simsekli, Ozan Sener, George Deligiannidis, and Murat A Erdogdu. Hausdorff dimension, heavy tails, and generalization in neural networks. In *Advances in Neural Information Processing Systems*, 2020.
- [64] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The Implicit Bias of Gradient Descent on Separable Data. *Journal of Machine Learning Research*, 2018.
- [65] Taiji Suzuki, Hiroshi Abe, and Tomoaki Nishimura. Compression based bound for non-compressed network: unified generalization error analysis of large compressible deep neural network. In *International Conference on Learning Representations*, 2020.
- [66] Rodrigo Veiga, Ludovic Stephan, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Phase diagram of Stochastic Gradient Descent in high-dimensional two-layer neural networks. *arXiv preprint arXiv:2202.00293*, 2022.
- [67] Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018.
- [68] Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, 2019.
- [69] Gilad Yehudai and Ohad Shamir. On the Power and Limitations of Random Features for Understanding Neural Networks. In *Advances in Neural Information Processing Systems*, 2019.
- [70] Lu Yu, Krishna Balasubramanian, Stanislav Volgushev, and Murat A Erdogdu. An Analysis of Constant Step Size SGD in the Non-convex Regime: Asymptotic Normality and Bias. In *Advances in Neural Information Processing Systems*, 2021.

- [71] Chulhee Yun, Shankar Krishnan, and Hossein Mobahi. A Unifying View on Implicit Bias in Training Linear Neural Networks. In *International Conference on Learning Representations*, 2021.
- [72] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 2021.
- [73] Xiao Zhang, Yaodong Yu, Lingxiao Wang, and Quanquan Gu. Learning One-hidden-layer ReLU Networks via Gradient Descent. In *International Conference on Artificial Intelligence and Statistics*, 2019.
- [74] Kai Zhong, Zhao Song, Prateek Jain, Peter L. Bartlett, and Inderjit S. Dhillon. Recovery Guarantees for One-hidden-layer Neural Networks. In *International Conference on Machine Learning*, 2017.
- [75] Mo Zhou, Rong Ge, and Chi Jin. A Local Convergence Theory for Mildly Over-Parameterized Two-Layer Neural Network. In *Conference on Learning Theory*, 2021.

Appendix A. Further Related Work

Two scaling regimes of neural networks have seen a surge of recent interest. In the regime of lazy training [22], the parameters hardly move from initialization and the NN does not learn useful features, behaving like a kernel method [4, 6, 25, 38, 54]. However, many works have shown that deep learning is more powerful than kernel models [30, 32, 69], establishing a clear separation between them; thus, several important characteristics of NNs cannot be captured in the lazy training regime [31]. In the other scaling regime, gradient descent on infinitely wide NNs reduces to Wasserstein gradient flow, which is known as the mean-field regime where feature learning is possible [19, 20, 48, 53, 60]. However, these results mostly hold for infinite or very wide NNs, and quantitative guarantees are difficult to obtain in this regime. The setting we consider in the current paper is different from both of these regimes, e.g. we may allow for NNs of arbitrary width which are not necessarily excessively overparameterized; yet, we still use the mean-field scaling when initializing the weights of the student NN.

On a separate note, we point out that our learnability result of Theorem 4 is consistent with [13]; they establish a sharp sample complexity of $\tilde{O}(d^{1 \vee (\mathcal{I}-2)})$ to learn a target with online SGD using the same activation f in the student network, where \mathcal{I} is the *information exponent* ($\mathcal{I} = 1$ in the above case due to the monotonicity of f). Despite assuming the link function f is known, we highlight that their setting covers $\mathcal{I} \geq 1$, whereas Theorem 4 is a proof concept to demonstrate the learnability implications of convergence to the principal subspace, even when f is unknown.

Appendix B. Additional Notations

For vectors \mathbf{v} and \mathbf{u} , we use $\langle \mathbf{v}, \mathbf{u} \rangle$ and $\mathbf{v} \circ \mathbf{u}$ to denote their Euclidean inner product and the element-wise product, and we use $\|\mathbf{v}\|_p$ and $\text{diag}(\mathbf{v})$ to denote the L_p -norm and the diagonal matrix whose diagonal entries are \mathbf{v} . For matrices \mathbf{V} and \mathbf{W} , we use $\langle \mathbf{V}, \mathbf{W} \rangle_{\text{F}}$, $\|\mathbf{V}\|_{\text{F}}$, and $\|\mathbf{V}\|_2$ to denote the Frobenius inner product, Frobenius norm, and the operator norm, respectively. For an activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, σ' and σ'' denote its first and second (weak) derivatives, which are applied element-wise for vector inputs. The symbol $\nabla \cdot$ is reserved for the gradient operator, and we frequently use $\nabla \ell$ to denote $\nabla_{\mathbf{W}} \ell$ when it is clear from the context.

In the appendix, we will prove the statements of the main text in a more general formulation. In particular, for smooth activations, we assume $\sup|\sigma'| \leq \beta_1$ and $\sup|\sigma''| \leq \beta_2$ for some $\beta_1, \beta_2 \in \mathbb{R}_+$, and we denote $\sup|\sigma| \leq \beta_0$, $\beta_0 \in (0, \infty]$. We will consider the following general case for the bias vector $\mathbf{b} \in \mathbb{R}^m$: $b_j \stackrel{\text{iid}}{\sim} \mathcal{D}_b$, such that $|b_j| \geq b^* > 0$, for some $b^* > 0$. This setting clearly covers the case of $b_j = \pm 1$ from the initialization of Assumption 2. Throughout the appendix, C will denote a generic positive absolute constant (e.g. 10), whose value may change from line to line.

We use the shorthand notation $\sigma_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x})$ to denote $\mathbf{a} \circ \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$, and similarly for $\sigma'_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x})$ and $\sigma''_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x})$. We use the notations $\text{vec}(\mathbf{A}) \in \mathbb{R}^{mn}$ for the vectorized representation of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, and $\mathbf{A} \otimes \mathbf{B}$ for the Kronecker product of two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{p \times q}$; we recall that the Kronecker product is an $mp \times nq$ block matrix comprised of $m \times n$ blocks of shape $p \times q$, where block (i, j) is given by $A_{ij}\mathbf{B}$.

Appendix C. Population Gradient Decomposition

C.1. Proof of Lemma 1

In what follows, ∇^\top is the Jacobian matrix and ∇ is the transpose of Jacobian for vector valued functions, which is the same as gradient for real-valued functions.

When σ is twice differentiable (Assumption 1.A), standard matrix calculations yield

$$\begin{aligned}
 \nabla_{\mathbf{W}} \mathbb{E}[R(\mathbf{W})] &\stackrel{(a)}{=} \mathbb{E}[\nabla_{\mathbf{W}} \ell(\hat{y}(\mathbf{x}; \mathbf{W}), y)] \\
 &= \mathbb{E}[\partial_1 \ell(\hat{y}(\mathbf{x}; \mathbf{W}), y) \nabla_{\mathbf{W}} \hat{y}(\mathbf{x}; \mathbf{W})] \\
 &= \mathbb{E}\left[\partial_1 \ell(\hat{y}(\mathbf{x}; \mathbf{W}), y) \sigma'_{a,b}(\mathbf{W}\mathbf{x}) \mathbf{x}^\top\right] \\
 &= \mathbb{E}\left[\mathbb{E}\left[\partial_1 \ell(\hat{y}(\mathbf{x}; \mathbf{W}), y) \sigma'_{a,b}(\mathbf{W}\mathbf{x}) \mathbf{x}^\top \mid \epsilon\right]\right] \\
 &\stackrel{(b)}{=} \mathbb{E}\left[\mathbb{E}\left[\nabla_{\mathbf{x}}^\top \{\partial_1 \ell(\hat{y}(\mathbf{x}; \mathbf{W}), g_\epsilon(\mathbf{U}\mathbf{x})) \sigma'_{a,b}(\mathbf{W}\mathbf{x})\} \mid \epsilon\right]\right] \\
 &= \mathbb{E}\left[\partial_1^2 \ell \sigma'_{a,b}(\mathbf{W}\mathbf{x}) \nabla_{\mathbf{x}}^\top \hat{y}(\mathbf{x}; \mathbf{W}) + \partial_1 \ell \nabla_{\mathbf{x}}^\top \sigma'_{a,b}(\mathbf{W}\mathbf{x}) + \partial_{12}^2 \ell \sigma'_{a,b}(\mathbf{W}\mathbf{x}) \nabla_{\mathbf{x}}^\top g_\epsilon(\mathbf{U}\mathbf{x})\right] \\
 &= \mathbb{E}\left[\left\{\partial_1^2 \ell \sigma'_{a,b}(\mathbf{W}\mathbf{x}) \sigma'_{a,b}(\mathbf{W}\mathbf{x})^\top + \partial_1 \ell \text{diag}(\sigma''_{a,b}(\mathbf{W}\mathbf{x}))\right\} \mathbf{W}\right] + \\
 &\quad + \mathbb{E}\left[\partial_{12}^2 \ell \sigma'_{a,b}(\mathbf{W}\mathbf{x}) \nabla g_\epsilon(\mathbf{U}\mathbf{x})^\top \mathbf{U}\right] \\
 &= \mathcal{H}(\mathbf{W}) \mathbf{W} + \mathcal{D}(\mathbf{W}) \mathbf{U}, \tag{C.1}
 \end{aligned}$$

where (a) follows from the dominated convergence theorem and (b) follows from the Stein's lemma, and ∇g_ϵ is the weak derivative of g_ϵ w.r.t. its inputs. Recall that Stein's lemma (Gaussian integration by parts) states that for $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$ and weakly differentiable $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we have $\mathbb{E}[\nabla f(\mathbf{x})] = \mathbb{E}[\mathbf{x} f(\mathbf{x})]$. Combining the above calculations with the gradient of the regularization term, with

$$\mathcal{D}(\mathbf{W}) = \mathbb{E}\left[\partial_{12}^2 \ell(\hat{y}, y) (\mathbf{a} \circ \sigma'(\mathbf{W}\mathbf{x} + \mathbf{b})) \nabla g_\epsilon^\top\right], \tag{C.2}$$

where ∇g_ϵ is the weak derivative of g_ϵ w.r.t. its inputs, and

$$\mathcal{H}(\mathbf{W}) = \mathbb{E}\left[(\mathbf{a} \circ \sigma'(\mathbf{W}\mathbf{x} + \mathbf{b})) (\mathbf{a} \circ \sigma'(\mathbf{W}\mathbf{x} + \mathbf{b}))^\top\right] + \mathbb{E}[(\hat{y} - y) \text{diag}((\mathbf{a} \circ \sigma''(\mathbf{W}\mathbf{x} + \mathbf{b})))], \tag{C.3}$$

the proof is complete for smooth activations.

For ReLU activations and ℓ satisfying Assumption 1.B, we introduce the following smooth approximation

$$\sigma_\iota(z) = \frac{1}{\iota} \log(1 + e^{\iota z}), \quad \iota > 0.$$

Then we have

$$\begin{aligned}
 \mathcal{H}_\iota(\mathbf{W}) &= \mathbb{E}\left[\partial_1^2 \ell (\mathbf{a} \circ \sigma'_\iota(\mathbf{W}\mathbf{x} + \mathbf{b})) (\mathbf{a} \circ \sigma'_\iota(\mathbf{W}\mathbf{x} + \mathbf{b}))^\top\right] + \mathbb{E}[\partial_1 \ell \text{diag}(\mathbf{a} \circ \sigma''_\iota(\mathbf{W}\mathbf{x} + \mathbf{b}))] \\
 &\succeq -\|\mathbf{a}\|_\infty \max_{1 \leq j \leq m} \mathbb{E}[\sigma''_\iota(\langle \mathbf{w}_j, \mathbf{x} \rangle + b_j)] \mathbf{I}_m.
 \end{aligned}$$

As $\sigma''_\tau \geq 0$, the critical step is to show $\lim_{\iota \rightarrow \infty} \mathbb{E}[\sigma''_\iota(\langle \mathbf{w}, \mathbf{x} \rangle + b)] < \infty$, uniformly for all \mathbf{w} . Let $z = \langle \mathbf{w}, \mathbf{x} \rangle + b$. Then $z \sim \mathcal{N}(b, \|\mathbf{w}\|_2^2)$, and

$$\begin{aligned}
 \int_0^\infty \sigma''_\iota(z) \frac{e^{-\frac{(z-b)^2}{2\|\mathbf{w}\|_2^2}}}{\sqrt{2\pi}\|\mathbf{w}\|_2} dz &\leq \iota \int_0^\infty \frac{e^{-\iota z - \frac{(z-b)^2}{2\|\mathbf{w}\|_2^2}}}{\sqrt{2\pi}\|\mathbf{w}\|_2} dz \\
 &= \iota e^{-\frac{b^2}{2\|\mathbf{w}\|_2^2} + \frac{(\iota\|\mathbf{w}\|_2 - \frac{b}{\|\mathbf{w}\|_2})^2}{2}} \int_0^\infty \frac{e^{-\frac{1}{2}(\frac{z}{\|\mathbf{w}\|_2} + \iota\|\mathbf{w}\|_2 - \frac{b}{\|\mathbf{w}\|_2})^2}}{\sqrt{2\pi}\|\mathbf{w}\|_2} dz. \\
 &= \iota e^{-\frac{b^2}{2\|\mathbf{w}\|_2^2} + \frac{(\iota\|\mathbf{w}\|_2 - \frac{b}{\|\mathbf{w}\|_2})^2}{2}} (1 - \Phi(\iota\|\mathbf{w}\|_2 - \frac{b}{\|\mathbf{w}\|_2})). \\
 &\stackrel{(a)}{\leq} \frac{\iota e^{-\frac{b^2}{2\|\mathbf{w}\|_2^2}}}{\sqrt{2\pi}\|\mathbf{w}\|_2(\iota - \frac{b}{\|\mathbf{w}\|_2^2})} \\
 &\stackrel{(b)}{\leq} \sqrt{\frac{2}{\pi}} \frac{e^{-\frac{b^2}{2\|\mathbf{w}\|_2^2}}}{\|\mathbf{w}\|_2} \\
 &\stackrel{(c)}{\leq} \frac{1}{|b|} \sqrt{\frac{2}{e\pi}},
 \end{aligned}$$

where (a) follows from the Gaussian tail bound $1 - \Phi(x) \leq \frac{e^{-x^2/2}}{\sqrt{2\pi}x}$, where Φ is the standard Gaussian CDF; (b) holds for large enough ι ; and (c) holds by considering supremum over $\|\mathbf{w}\|_2$. Thus $\mathbb{E}[\sigma''_\iota(\langle \mathbf{w}_j, \mathbf{x} \rangle + b_j)] \leq \frac{2}{|b_j|} \sqrt{\frac{2}{e\pi}}$ and consequently,

$$\frac{-2\|\mathbf{a}\|_\infty}{b^*} \sqrt{\frac{2}{e\pi}} \mathbf{I}_m \preceq \mathcal{H}_\iota(\mathbf{W}) \preceq \left(\|\mathbf{a}\|_2^2 + \frac{2\|\mathbf{a}\|_\infty}{b^*} \sqrt{\frac{2}{e\pi}} \right) \mathbf{I}_m$$

where $b^* = \min_{1 \leq j \leq m} |b_j|$. Moreover, as $\sigma'_\iota(\mathbf{W}\mathbf{x} + \mathbf{b})$ converges a.s. (i.e. except when $\langle \mathbf{w}_j, \mathbf{x} \rangle + b_j = 0$ for some j) to $\sigma'(\mathbf{W}\mathbf{x} + \mathbf{b})$, by the dominated convergence theorem,

$$\nabla R(\mathbf{W}) = \lim_{\iota \rightarrow \infty} \mathcal{H}_\iota(\mathbf{W})\mathbf{W} + \lim_{\iota \rightarrow \infty} \mathcal{D}_\iota(\mathbf{W})\mathbf{U}$$

We can immediately observe from the dominated convergence theorem that $\mathcal{D}_\iota(\mathbf{W}) \rightarrow \mathcal{D}(\mathbf{W})$ as $\iota \rightarrow \infty$ with $\mathcal{D}(\mathbf{W})$ given in (C.2). Moreover, we let $\mathcal{H}(\mathbf{W}) = \lim_{\iota \rightarrow \infty} \mathcal{H}_\iota(\mathbf{W})$, and observe that

$$\frac{-2\|\mathbf{a}\|_\infty}{b^*} \sqrt{\frac{2}{e\pi}} \mathbf{I}_m \preceq \mathcal{H}(\mathbf{W}) \preceq \left(\|\mathbf{a}\|_2^2 + \frac{2\|\mathbf{a}\|_\infty}{b^*} \sqrt{\frac{2}{e\pi}} \right) \mathbf{I}_m. \quad (\text{C.4})$$

This finishes the proof of Lemma 1. ■

In the case of smooth activations (Assumption 1.A), the following bounds will be useful.

Lemma 6 *Let $R(\mathbf{W}) := \mathbb{E}[\ell(\hat{y}(\mathbf{x}; \mathbf{W}, \mathbf{a}, \mathbf{b}), y)]$ be the unregularized population risk. Under Assumption 1.A we have*

$$-\beta_2\|\mathbf{a}\|_\infty \sqrt{2R(\mathbf{W})} \mathbf{I}_m \preceq \mathcal{H}(\mathbf{W}) \preceq \left\{ \beta_1^2 \|\mathbf{a}\|_2^2 + \beta_2 \|\mathbf{a}\|_\infty \sqrt{2R(\mathbf{W})} \right\} \mathbf{I}_m. \quad (\text{C.5})$$

Proof Assumption 1.A requires $\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$. Hence by definition of \mathcal{H} ,

$$\mathcal{H}(\mathbf{W}) = \mathbb{E} \left[\sigma'_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x}) \sigma'_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x})^\top \right] + \mathbb{E} \left[(\hat{y}(\mathbf{x}; \mathbf{W}) - y) \text{diag}(\sigma''_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x})) \right].$$

The first term is positive semi-definite and it can be easily bounded:

$$0 \leq \mathbf{v}^\top \mathbb{E} \left[\sigma'_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x}) \sigma'_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x})^\top \right] \mathbf{v} \leq \mathbb{E} \left[\|\sigma'_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x})\|_2 \|\mathbf{v}\|_2^2 \right] \leq \beta_1^2 \|\mathbf{a}\|_2^2 \|\mathbf{v}\|_2^2$$

for an arbitrary vector $\mathbf{v} \in \mathbb{R}^m$. For the second term, we have

$$-\beta_2 \|\mathbf{a}\|_\infty \mathbb{E} \left[|\hat{y} - y| \right] \mathbf{I}_m \preceq \mathbb{E} \left[(\hat{y}(\mathbf{x}; \mathbf{W}) - y) \text{diag}(\sigma''_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x})) \right] \preceq \beta_2 \|\mathbf{a}\|_\infty \mathbb{E} \left[|\hat{y} - y| \right] \mathbf{I}_m$$

and $\mathbb{E} \left[|\hat{y}(\mathbf{x}; \mathbf{W}) - y| \right] \leq \sqrt{2R(\mathbf{W})}$ by Jensen's inequality. \blacksquare

Appendix D. Proofs of Section 3

We begin by characterizing the tail behavior of the stochastic gradient noise in the SGD updates (3.1) through the following lemma.

Lemma 7 For any fixed $\mathbf{W} \in \mathbb{R}^{m \times k}$, let

$$\mathbf{\Gamma} := \nabla \ell(\hat{y}(\mathbf{x}; \mathbf{W}), y) - \mathbb{E}[\nabla \ell(\hat{y}(\mathbf{x}; \mathbf{W}), y)]$$

denote the zero-mean stochastic noise in the gradient of the loss function ℓ when $(\mathbf{x}, y) \sim \mathcal{P}$, and recall that

$$\nabla \ell(\hat{y}(\mathbf{x}; \mathbf{W}), y) = \partial_1 \ell(\hat{y}(\mathbf{x}; \mathbf{W}), y) \sigma'_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x}) \mathbf{x}^\top.$$

Suppose $\sup_{\hat{y}, y} |\partial_1 \ell(\hat{y}, y)| \leq \varkappa$. Then for any $\mathbf{V} \in \mathbb{R}^{m \times d}$, the zero-mean random variable $\langle \mathbf{V}, \mathbf{\Gamma} \rangle_{\mathbb{F}}$ is $C\beta_1 \varkappa \|\mathbf{a}\|_2 \|\mathbf{V}\|_{\mathbb{F}}$ -sub-Gaussian.

Proof We use the shorthand notation $\nabla \ell := \nabla_{\mathbf{W}} \ell(\mathbf{W}\mathbf{x}, y)$ and $\nabla R := \nabla_{\mathbf{W}} R(\mathbf{W})$. We compute the following

$$\begin{aligned} \mathbb{E} \left[|\langle \mathbf{V}, \nabla \ell - \nabla R \rangle_{\mathbb{F}}|^p \right]^{\frac{1}{p}} &\stackrel{(a)}{\leq} \mathbb{E} \left[|\langle \mathbf{V}, \nabla \ell \rangle_{\mathbb{F}}|^p \right]^{\frac{1}{p}} + \mathbb{E} \left[|\langle \mathbf{V}, \nabla R \rangle_{\mathbb{F}}|^p \right]^{\frac{1}{p}} \\ &\stackrel{(b)}{\leq} 2 \mathbb{E} \left[|\langle \mathbf{V}, \nabla \ell \rangle_{\mathbb{F}}|^p \right]^{\frac{1}{p}} \\ &\leq 2\varkappa \mathbb{E} \left[\left| \left\langle \mathbf{V}, \sigma'_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x}) \mathbf{x}^\top \right\rangle_{\mathbb{F}} \right|^{2p} \right]^{\frac{1}{2p}}. \end{aligned}$$

where (a) and (b) follow from the Minkowski and Jensen inequalities respectively. Furthermore, we have

$$\begin{aligned} \mathbb{E} \left[\left| \left\langle \mathbf{V}, \sigma'_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x}) \mathbf{x}^\top \right\rangle_{\mathbb{F}} \right|^{2p} \right]^{\frac{1}{2p}} &= \mathbb{E} \left[|\langle \mathbf{V}\mathbf{x}, \sigma'_{\mathbf{a}, \mathbf{b}}(\mathbf{W}\mathbf{x}) \rangle|^{2p} \right]^{\frac{1}{2p}} \\ &\leq \beta_1 \|\mathbf{a}\|_2 \mathbb{E} \left[\|\mathbf{V}\mathbf{x}\|_2^{2p} \right]^{\frac{1}{2p}} \\ &\leq \beta_1 \|\mathbf{a}\|_2 (\|\mathbf{V}\|_{\mathbb{F}} + C\|\mathbf{V}\|_2 \sqrt{p}), \end{aligned}$$

where the last inequality follows from Gaussianity of $\mathbf{V}\mathbf{x}$ and Lemma 25. Hence

$$\mathbb{E}[|\langle \mathbf{V}, \nabla \ell - \nabla R \rangle_{\mathbb{F}}|^p]^{\frac{1}{p}} \leq C\beta_1 \varkappa \|\mathbf{a}\|_2 \|\mathbf{V}\|_{\mathbb{F}} \sqrt{p}.$$

Invoking Lemma 21 implies sub-Gaussianity of $\langle \mathbf{V}, \nabla \ell - \nabla R \rangle_{\mathbb{F}}$ and completes the proof. \blacksquare

We proceed by presenting a lemma which constitutes the main part of the proof of Theorem 2 via establishing a recursive bound on the moment generating function (MGF) of $\|\mathbf{W}_{\perp}^t\|_{\mathbb{F}}^2$, which will in turn be used to prove high probability statements for $\|\mathbf{W}_{\perp}^t\|_{\mathbb{F}}^2$.

Lemma 8 *Consider running the iterates of SGD (3.1), under either Assumption 1.A or 1.B, with step size sequence $\{\eta_t\}_{t \geq 0}$ that is either constant $\eta_t = \eta$ or decreasing $\eta_t = m \frac{2(t^*+t)+1}{\gamma(t^*+t+1)^2}$ (cf. [35, Theorem 3.2]). Let $\varkappa := \sup |\partial_1 \ell(\hat{y}, y)|$, $\kappa := \beta_1 \|\mathbf{a}\|_2 \varkappa$, $\zeta := \beta_2 (\mathbb{E}[|y|] + \beta_0 \|\mathbf{a}\|_1)$ under Assumption 1.A and $\zeta := 2\sqrt{2/(e\pi)}$ under Assumption 1.B, and finally $\tilde{\varrho} := \lambda + \beta_1 \|\mathbf{a}\|_2^2 + \zeta \|\mathbf{a}\|_{\infty}$. Suppose $\eta_0 \lesssim \tilde{\varrho}^{-1}$. Let \mathcal{F}_t denote the sigma algebra generated by $\{\mathbf{W}^j\}_{j=0}^t$, and let $\{A_t\}_{t \geq 0}$ be a sequence of decreasing events (i.e. $A_{t+1} \subseteq A_t$), such that $A_t \in \mathcal{F}_t$ and on A_t we have $\mathcal{H}(\mathbf{W}^t) + \lambda \mathbf{I}_m \succeq \frac{\gamma}{m} \mathbf{I}_m$. Then, for every $t \geq 0$, with probability at least $\mathbb{P}(A_t) - \delta$,*

$$\|\mathbf{W}_{\perp}^t\|_{\mathbb{F}}^2 \lesssim \prod_{j=0}^{t-1} (1 - \frac{\eta_j \gamma}{m}) \|\mathbf{W}_{\perp}^0\|_{\mathbb{F}}^2 + \frac{m\eta_t \kappa^2 (d + \log(1/\delta))}{\gamma}. \quad (\text{D.1})$$

Proof Let \mathcal{F}_t denote the sigma algebra generated by $\{\mathbf{W}^j\}_{j=0}^t$. Recall from Lemma 7 that we define

$$\mathbf{\Gamma}^t = \nabla \ell(\hat{y}(\mathbf{x}^{(t)}; \mathbf{W}^t), y^{(t)}) - \mathbb{E}[\nabla \ell(\hat{y}(\mathbf{x}^{(t)}; \mathbf{W}^t), y^{(t)})]$$

with

$$\nabla \ell(\hat{y}(\mathbf{x}^{(t)}; \mathbf{W}^t), y^{(t)}) = \partial_1 \ell(\hat{y}(\mathbf{x}^{(t)}; \mathbf{W}^t), y^{(t)}) \sigma'_{a,b}(\mathbf{W}^t \mathbf{x}^{(t)}) (\mathbf{x}^{(t)})^{\top}.$$

Then for the SGD updates we have

$$\mathbf{W}^{t+1} = \mathbf{W}^t - \eta_t \nabla \mathcal{R}_{\lambda}(\mathbf{W}^t) - \eta_t \mathbf{\Gamma}^t.$$

By projecting the iterates onto the orthogonal complement of the principal subspace,

$$\mathbf{W}_{\perp}^{t+1} = (\mathbf{I}_m - \eta_t (\mathcal{H}(\mathbf{W}^t) + \lambda \mathbf{I}_m)) \mathbf{W}_{\perp}^t - \eta_t \mathbf{\Gamma}_{\perp}^t.$$

Let $\mathbf{M}_t := \mathbf{I}_m - \eta_t (\mathcal{H}(\mathbf{W}^t) + \lambda \mathbf{I}_m)$. Then, by observing that $\mathbf{1}_{A_{t+1}} \leq \mathbf{1}_{A_t}$, for any $0 \leq s \lesssim \frac{\gamma}{m\eta_t \kappa^2}$ we have

$$\begin{aligned} \mathbb{E} \left[\mathbf{1}_{A_{t+1}} e^{s \|\mathbf{W}_{\perp}^{t+1}\|_{\mathbb{F}}^2} \mid \mathcal{F}_0 \right] &\leq \mathbb{E} \left[\mathbf{1}_{A_t} e^{s \|\mathbf{M}_t \mathbf{W}_{\perp}^t\|_{\mathbb{F}}^2 + s \eta_t^2 \|\mathbf{\Gamma}_{\perp}^t\|_{\mathbb{F}}^2 + \langle -2s\eta_t \mathbf{M}_t \mathbf{W}_{\perp}^t, \mathbf{\Gamma}_{\perp}^t \rangle_{\mathbb{F}}} \mid \mathcal{F}_0 \right] \\ &= \mathbb{E} \left[\mathbf{1}_{A_t} e^{s \|\mathbf{M}_t \mathbf{W}_{\perp}^t\|_{\mathbb{F}}^2} \mathbb{E} \left[e^{s \eta_t^2 \|\mathbf{\Gamma}_{\perp}^t\|_{\mathbb{F}}^2} e^{\langle -2s\eta_t \mathbf{M}_t \mathbf{W}_{\perp}^t, \mathbf{\Gamma}_{\perp}^t \rangle_{\mathbb{F}}} \mid \mathcal{F}_t \right] \mid \mathcal{F}_0 \right] \\ &\leq \mathbb{E} \left[\mathbf{1}_{A_t} e^{s \|\mathbf{M}_t \mathbf{W}_{\perp}^t\|_{\mathbb{F}}^2} \mathbb{E} \left[e^{2s\eta_t^2 \|\mathbf{\Gamma}_{\perp}^t\|_{\mathbb{F}}^2} \mid \mathcal{F}_t \right]^{\frac{1}{2}} \mathbb{E} \left[e^{\langle -4s\eta_t \mathbf{M}_t \mathbf{W}_{\perp}^t, \mathbf{\Gamma}_{\perp}^t \rangle_{\mathbb{F}}} \mid \mathcal{F}_t \right]^{\frac{1}{2}} \mid \mathcal{F}_0 \right], \end{aligned} \quad (\text{D.2})$$

where the last inequality follows from the Cauchy-Schwartz inequality for conditional expectation.

Moreover, it is straightforward to observe that $\|\nabla\ell(\hat{y}(\mathbf{x}^{(t)}; \mathbf{W}^t), y^{(t)})\|_{\mathbb{F}}^2 \leq \kappa^2 \|\mathbf{x}\|_2^2$, hence

$$\mathbb{E}\left[\|\nabla\ell(\hat{y}(\mathbf{x}^{(t)}; \mathbf{W}^t), y^{(t)})\|_{\mathbb{F}}^2\right] \leq \kappa^2 d.$$

Note that by Jensen's inequality

$$\|\mathbf{\Gamma}_{\perp}^t\|_{\mathbb{F}}^2 \leq 2\|\nabla\ell(\mathbf{W}^t \mathbf{x}^{(t)}, y^{(t)})\|_{\mathbb{F}}^2 + 2\mathbb{E}\left[\|\nabla\ell(\mathbf{W}^t \mathbf{x}^{(t)}, y^{(t)})\|_{\mathbb{F}}^2\right].$$

Consequently

$$\begin{aligned} \mathbb{E}\left[\exp(2s\eta_t^2 \|\mathbf{\Gamma}_{\perp}^t\|_{\mathbb{F}}^2) \mid \mathcal{F}_t\right] &\leq \exp(4s\eta_t^2 \kappa^2 d) \mathbb{E}\left[\exp(4s\eta_t^2 \kappa^2 \|\mathbf{x}\|_2^2) \mid \mathcal{F}_t\right] \\ &\leq \exp(4s\eta_t^2 \kappa^2 d) \exp(8s\eta_t^2 \kappa^2 d), \end{aligned}$$

where the second inequality follows from Lemma 26 for $4s\eta_t^2 \kappa^2 \leq 1/4$. Since $s \lesssim \frac{\gamma}{m\eta_t \kappa^2}$, in order to satisfy the condition of Lemma 26 we need to ensure $\eta_t \gamma / m \lesssim 1$, which is guaranteed by our $\eta_t \tilde{\varrho} \lesssim 1$ assumption for a suitably small absolute constant, as $\gamma / m \leq \lambda \leq \tilde{\varrho}$.

Next, we bound the last term in (D.2). Let $\mathbf{V} := -4s\eta_t \mathbf{M}_t \mathbf{W}_{\perp}^t$. Then by Lemma 7 we have

$$\mathbb{E}\left[\exp(\langle \mathbf{V}, \mathbf{\Gamma}_{\perp}^t \rangle_{\mathbb{F}}) \mid \mathcal{F}_t\right] \leq \exp(Cs^2 \eta_t^2 \kappa^2 \|\mathbf{M}_t \mathbf{W}_{\perp}^t\|_{\mathbb{F}}^2)$$

Putting things back together in (D.2) and using the tower property of expectation, we have

$$\mathbb{E}\left[\mathbf{1}_{A_{t+1}} e^{s\|\mathbf{W}_{\perp}^{t+1}\|_{\mathbb{F}}^2} \mid \mathcal{F}_0\right] \leq \mathbb{E}\left[\mathbf{1}_{A_t} e^{s(1+Cs\eta_t^2 \kappa^2)\|\mathbf{M}_t \mathbf{W}_{\perp}^t\|_{\mathbb{F}}^2 + Cs\eta_t^2 \kappa^2 d} \mid \mathcal{F}_0\right]. \quad (\text{D.3})$$

Next, we bound $\|\mathbf{M}_t\|_2$. By definition of A_t , we can already ensure $\mathcal{H}(\mathbf{W}^t) \succeq \frac{\gamma}{m} \mathbf{I}_m$ in (D.3). Recall the definition of $\mathcal{H}(\mathbf{W}^t)$

$$\mathcal{H}(\mathbf{W}^t) = \mathbb{E}\left[\partial_1^2 \ell(\hat{y}(\mathbf{x}; \mathbf{W}^t), y) \sigma'_{a,b}(\mathbf{W}^t \mathbf{x}) \sigma'_{a,b}(\mathbf{W}^t \mathbf{x})^{\top} + \partial_1 \ell(\hat{y}(\mathbf{x}; \mathbf{W}^t), y) \text{diag}(\sigma''_{a,b}(\mathbf{W}^t \mathbf{x}))\right].$$

Notice that $0 \leq \partial_1^2 \ell(\hat{y}(\mathbf{x}; \mathbf{W}), y) \leq 1$ under either Assumption 1.A or Assumption 1.B. Moreover we have, $|\partial_1 \ell(\hat{y}, y)| \leq \varkappa$. Thus,

$$\mathcal{H}(\mathbf{W}^t) + \lambda \mathbf{I}_m \preceq (\lambda + \beta_1^2 \|\mathbf{a}\|_2^2 + \zeta \|\mathbf{a}\|_{\infty}) \mathbf{I}_m = \tilde{\varrho} \mathbf{I}_m.$$

Therefore,

$$0 \preceq \mathbf{I}_m - \eta_t (\mathcal{H}(\mathbf{W}^t) + \lambda \mathbf{I}_m) \preceq (1 - \frac{\eta_t \gamma}{m}) \mathbf{I}_m.$$

As a result $\|\mathbf{M}_t\|_2 \leq 1 - \frac{\eta_t \gamma}{m}$. Combined with (D.3) we have

$$\begin{aligned} \mathbb{E}\left[\mathbf{1}_{A_{t+1}} \exp(s\|\mathbf{W}_{\perp}^{t+1}\|_{\mathbb{F}}^2) \mid \mathcal{F}_0\right] &\leq \mathbb{E}\left[\mathbf{1}_{A_t} \exp(s(1+Cs\eta_t^2 \kappa^2)(1 - \frac{\eta_t \gamma}{m})^2 \|\mathbf{W}_{\perp}^t\|_{\mathbb{F}}^2 + Cs\eta_t^2 d \kappa^2) \mid \mathcal{F}_0\right] \\ &\leq \exp(Cs\eta_t^2 \kappa^2 d) \mathbb{E}\left[\mathbf{1}_{A_t} \exp(s(1 - \frac{\eta_t \gamma}{m}) \|\mathbf{W}_{\perp}^t\|_{\mathbb{F}}^2) \mid \mathcal{F}_0\right] \end{aligned} \quad (\text{D.4})$$

where the second inequality holds by the fact that $Cs\eta_t^2 \kappa^2 \leq \eta_t \gamma / m$, which in turn holds when a small enough absolute constant is chosen in $0 \leq s \lesssim \frac{\gamma}{m\eta_t \kappa^2}$. Also notice that for decreasing step size,

$$1 - \frac{\gamma \eta_t}{m} = \frac{(t+t^*)^2}{(t+t^*+1)^2} \leq \frac{1 - \frac{(t+t^*)^2}{(t+t^*+1)^2}}{1 - \frac{(t+t^*-1)^2}{(t+t^*)^2}} = \frac{\eta_t}{\eta_{t-1}}, \quad (\text{D.5})$$

(and the above holds trivially for constant step size), thus when $s \leq \frac{C_1\gamma}{\eta_t\kappa^2}$ for some absolute constant C_1 , we have $s(1 - \eta_t\gamma) \leq \frac{C_1\gamma}{\eta_{t-1}\kappa^2}$ with the same absolute constant. Hence we are allowed to expand the recursion (D.4), which implies

$$\mathbb{E}[\mathbf{1}_{A_t} \exp(s\|\mathbf{W}_\perp^t\|_{\mathbb{F}}^2) \mid \mathcal{F}_0] \leq \exp\left(s \prod_{j=0}^{t-1} \left(1 - \frac{\eta_j\gamma}{m}\right) \|\mathbf{W}_\perp^0\|_{\mathbb{F}}^2 + Cs\kappa^2d \sum_{i=0}^{t-1} \eta_i^2 \prod_{j=i+1}^{t-1} \left(1 - \frac{\eta_j\gamma}{m}\right)\right)$$

for all $0 \leq s \lesssim \frac{\gamma}{m\eta_{t-1}\kappa^2}$. Moreover, direct calculation implies that with both constant and decreasing step sizes of Lemma 8, we have $\sum_{i=0}^{t-1} \eta_i^2 \prod_{j=i+1}^{t-1} \left(1 - \frac{\eta_j\gamma}{m}\right) \leq \frac{Cm\eta_t}{\gamma}$ (with $C = 1$ for constant step size). Thus, for all $0 \leq s \lesssim \frac{\gamma}{m\eta_{t-1}\kappa^2}$

$$\mathbb{E}[\mathbf{1}_{A_t} \exp(s\|\mathbf{W}_\perp^t\|_{\mathbb{F}}^2) \mid \mathcal{F}_0] \leq \exp\left(s \prod_{j=0}^{t-1} \left(1 - \frac{\eta_j\gamma}{m}\right) \|\mathbf{W}_\perp^0\|_{\mathbb{F}}^2 + \frac{Csm\eta_t\kappa^2d}{\gamma}\right).$$

Finally, we can apply a Chernoff bound to obtain

$$\mathbb{P}(A_t \cap \{\|\mathbf{W}_\perp^t\|_{\mathbb{F}}^2 \geq \varepsilon\} \mid \mathcal{F}_0) \leq \exp\left(s \left\{ \prod_{j=0}^{t-1} \left(1 - \eta_j\gamma\right) \|\mathbf{W}_\perp^0\|_{\mathbb{F}}^2 + \frac{Cm\eta_t\kappa^2d}{\gamma} - \varepsilon \right\}\right)$$

By choosing

$$\varepsilon = \prod_{j=0}^{t-1} \left(1 - \frac{\eta_j\gamma}{m}\right) \|\mathbf{W}_\perp^0\|_{\mathbb{F}}^2 + \frac{Cm\eta_t\kappa^2(d + \log(1/\delta))}{\gamma}.$$

and the largest possible $s \lesssim \frac{\gamma}{m\eta_t\kappa^2}$, we obtain

$$\mathbb{P}(\|\mathbf{W}_\perp^t\|_{\mathbb{F}}^2 \geq \varepsilon \mid \mathcal{F}_0) \leq \mathbb{P}(A_t \cap \{\|\mathbf{W}_\perp^t\|_{\mathbb{F}} \geq \varepsilon\}) + \mathbb{P}(A_t^C) \leq \delta + \mathbb{P}(A_t^C).$$

Taking another expectation to remove conditioning on initialization completes the proof. \blacksquare

The proof of Theorem 2 for decreasing step size follows by a direct computation of the quantities in Lemma 8 and is presented below. We remark that proving the same result with a constant step size is essentially similar.

D.1. Proof of Theorem 2

This part is directly implied by Lemma 8. The following argument holds on an event where $\|\mathbf{W}\|_{\mathbb{F}} \lesssim \sqrt{m}$, which happens with probability at least $1 - \mathcal{O}(\delta)$. In order to see this connection, we will first present an improved statement over Lemma 6 for the case of smooth activations. Recall the definition of $\mathcal{H}(\mathbf{W})$ for the squared error loss $\ell(\hat{y}, y) = \frac{(\hat{y} - y)^2}{2}$,

$$\mathcal{H}(\mathbf{W}) = \mathbb{E}\left[\sigma'_{\mathbf{a},\mathbf{b}}(\mathbf{W}\mathbf{x})\sigma'_{\mathbf{a},\mathbf{b}}(\mathbf{W}\mathbf{x})^\top\right] + \mathbb{E}\left[(\hat{y}(\mathbf{x}; \mathbf{W}, \mathbf{a}, \mathbf{b}) - y) \text{diag}(\sigma''_{\mathbf{a},\mathbf{b}}(\mathbf{W}\mathbf{x}))\right].$$

Notice that under Assumption 1.A we have $|\hat{y}| \leq \beta_0\|\mathbf{a}\|_1$. Then basic matrix algebra similar to that of Lemma 6 along with the triangle inequality shows

$$-\beta_2\|\mathbf{a}\|_\infty(\beta_0\|\mathbf{a}\|_1 + \mathbb{E}[|y|])\mathbf{I}_m \prec \mathcal{H}(\mathbf{W}) \preceq (\beta_1^2\|\mathbf{a}\|_2^2 + \beta_2\|\mathbf{a}\|_\infty(\beta_0\|\mathbf{a}\|_1 + \mathbb{E}[|y|]))\mathbf{I}_m.$$

Therefore, with $\lambda \geq \gamma/m + \beta_2 \|\mathbf{a}\|_\infty (\|\mathbf{a}\|_1 \beta_0 + \mathbb{E}[|y|])$, we have $\mathcal{H}(\mathbf{W}) + \lambda \mathbf{I}_m \succeq \gamma/m \mathbf{I}_m$ for all \mathbf{W} . In addition, $|\partial_1 \ell(\hat{y}, y)| \leq \beta_0 \|\mathbf{a}\|_1 + K$ by the triangle inequality. Thus we can invoke Lemma 8 with $\eta_t = \frac{m}{\gamma} \left(1 - \frac{(t^*+t)^2}{(t^*+t+1)^2}\right)$, $\varkappa = \beta_0 \|\mathbf{a}\|_1 + K$, $\zeta = \beta_2 \varkappa$ and $\mathbf{1}_{A_t} = 1$. Recall that in the statement of the theorem, $\beta_0 = \beta_1 = \beta_2 = 1$, $K \lesssim 1$, $\|\mathbf{a}\|_\infty \leq 1/m$, $\|\mathbf{a}\|_2 \leq 1/\sqrt{m}$, and $\|\mathbf{a}\|_1 \leq 1$, hence $\tilde{\rho} \asymp \lambda$, and with $t^* \asymp \frac{\lambda}{\gamma}$ we can guarantee $\eta_t \lambda \lesssim 1$. As the step size condition of Lemma 8 is satisfied, the desired result follows.

Similarly, for ReLU we have $|\partial_1 \ell(\hat{y}, y)| \leq 1$ by Assumption 1.B, and for $\lambda \geq \gamma/m + \frac{2\|\mathbf{a}\|_\infty}{b^*} \sqrt{\frac{2}{e\pi}}$, we have $\mathcal{H}(\mathbf{W}) + \lambda \mathbf{I}_m \succeq \gamma/m \mathbf{I}_m$. Hence this time, we can invoke Lemma 8 with the same decreasing η_t and $\mathbf{1}_{A_t} = 1$, $\varkappa = 1$, and $\zeta = \frac{2}{b^*} \sqrt{\frac{2}{e\pi}}$ (recall $b^* = 1$ in the statement of the theorem). \blacksquare

Appendix E. Proofs of Section 4

E.1. Proof of Theorem 5

As our arguments are based on the Rademacher complexity of a two-layer neural network, we require the knowledge of the norm of \mathbf{W}^t . We prove a high probability bound for this norm in the following lemma.

Lemma 9 *Under Assumption 1.A or 1.B with either decreasing or constant step size as in Theorem 2, let $\varkappa = \sup_{\hat{y}, y} |\partial_1 \ell(\hat{y}, y)| < \infty$ and $\kappa_\infty := \beta_1 \|\mathbf{a}\|_\infty \varkappa$. Then for any $t \geq 1$, with probability at least $1 - m \exp\left(\frac{-\gamma t d}{2m\lambda}\right)$ we have for all $1 \leq j \leq m$*

$$\|\mathbf{w}_j^t\|_2 \leq \prod_{i=0}^{t-1} (1 - \eta_i \lambda) \|\mathbf{w}_j^0\|_2 + \frac{3\kappa_\infty \sqrt{d}}{\lambda}. \quad (\text{E.1})$$

Proof First, we prove that for any $t \geq 0$ and $0 \leq s \leq \frac{2\sqrt{d}}{\kappa_\infty \eta_t}$, we have

$$\mathbb{E}[\exp(s \|\mathbf{w}_j^t\|_2) \mid \mathbf{W}^0] \leq \exp\left(s \prod_{i=0}^{t-1} (1 - \eta_i \lambda) \|\mathbf{w}_j^0\|_2 + \frac{2s\kappa_\infty \sqrt{d}}{\lambda}\right), \quad (\text{E.2})$$

The base case of $t = 0$ is trivial, and for the induction step we have

$$\begin{aligned} \mathbb{E}\left[\exp\left(s \|\mathbf{w}_j^{t+1}\|_2\right) \mid \mathbf{W}^0\right] &= \mathbb{E}\left[\exp\left(s \|(1 - \eta_t \lambda) \mathbf{w}_j^t - \eta_t \nabla_{\mathbf{w}_j} \ell(\hat{y}(\mathbf{x}; \mathbf{W}^t), y)\|_2\right) \mid \mathbf{W}^0\right] \\ &\leq \mathbb{E}\left[\exp\left(s(1 - \eta_t \lambda) \|\mathbf{w}_j^t\|_2 + s \eta_t \|\nabla_{\mathbf{w}_j} \ell(\hat{y}(\mathbf{x}; \mathbf{W}^t), y)\|_2\right) \mid \mathbf{W}^0\right] \\ &= \mathbb{E}\left[\exp\left(s(1 - \eta_t \lambda) \|\mathbf{w}_j^t\|_2 + s \eta_t \kappa_\infty \|\mathbf{x}\|_2\right) \mid \mathbf{W}^0\right] \\ &= \mathbb{E}\left[\exp\left(s(1 - \eta_t \lambda) \|\mathbf{w}_j^t\|_2\right) \mathbb{E}\left[\exp\left(s \eta_t \kappa_\infty \|\mathbf{x}\|_2\right) \mid \mathbf{W}^t, \mathbf{W}^0\right] \mid \mathbf{W}^0\right] \\ &\stackrel{(a)}{\leq} \mathbb{E}\left[\exp\left(s(1 - \eta_t \lambda) \|\mathbf{w}_j^t\|_2\right) \exp\left(s \eta_t \kappa_\infty \sqrt{d} + \frac{s^2 \kappa_\infty^2 \eta_t^2}{2}\right) \mid \mathbf{W}^0\right] \\ &\stackrel{(b)}{\leq} \exp\left(s \prod_{i=0}^t (1 - \eta_i \lambda) \|\mathbf{w}_j^0\|_2 + \frac{2s\kappa_\infty \sqrt{d}}{\lambda}\right) \end{aligned}$$

where (a) holds since $\|\mathbf{x}\|_2$ is a 1-Lipschitz function of a standard Gaussian random vector, thus it is sub-Gaussian with parameter 1 (Lemma 22) and additionally $\mathbb{E}[\|\mathbf{x}\|_2] \leq \sqrt{d}$, and (b) holds by the induction hypothesis (notice that for decreasing step size $s(1 - \eta_t \lambda) \leq \frac{2\sqrt{d}}{\kappa_\infty \eta_{t-1}}$ by (D.5)). Next, we apply the following Chernoff bound,

$$\mathbb{P}\left(\|\mathbf{w}_j^t\|_2 > \prod_{i=0}^{t-1} (1 - \eta_i \lambda) \|\mathbf{w}_j^0\|_2 + \frac{3\kappa_\infty \sqrt{d}}{\lambda} \mid \mathbf{W}^0\right) \leq \exp\left(-\frac{s\kappa_\infty \sqrt{d}}{\lambda}\right),$$

which holds for any $0 \leq s \leq \frac{2\sqrt{d}}{\kappa_\infty \eta_t}$. Choosing the largest s possible and noting that $\eta_t \leq \frac{4m}{\gamma t}$ yields an $\exp\left(\frac{-\gamma t d}{2m\lambda}\right)$ upper bound on the conditional probability, which followed by taking expectation removes the randomness of conditioning on \mathbf{w}_j^0 . Finally applying a union bound gives us the desired bound. \blacksquare

In addition, we would like to approximate $R_\tau(\mathbf{W}^T)$ and $\hat{R}_\tau(\mathbf{W}^T)$ with $R_\tau(\mathbf{W}_\parallel^T)$ and $\hat{R}_\tau(\mathbf{W}_\parallel^T)$ respectively. As a result, we will investigate the Lipschitzness of the population and empirical risk in the next lemma.

Lemma 10 *Under either Assumption 1.A or 1.B, the truncated risk $\mathbf{W} \mapsto R_\tau(\mathbf{W})$ is $\sqrt{2}\tau\beta_1\|\mathbf{a}\|_2$ -Lipschitz. Moreover, for $T \geq d + \log(1/\delta)$ with probability at least $1 - \delta$ over the stochasticity of $\{\mathbf{x}^{(t)}\}_{0 \leq t \leq T-1}$, the truncated empirical risk $\mathbf{W} \mapsto \hat{R}_\tau(\mathbf{W})$ is $C\tau\beta_1\|\mathbf{a}\|_2$ -Lipschitz for some absolute constant C .*

Proof We begin by the simple observation that $\hat{y} \mapsto \ell(\hat{y}, y) \wedge \tau$ is $\sqrt{2}\tau$ -Lipschitz when $\ell(\hat{y}, y) = 1/2(\hat{y} - y)^2$ and 1-Lipschitz when $|\partial_1 \ell(\hat{y}, y)| \leq 1$. As $\tau \geq 1$, we can consider both of them as $\sqrt{2}\tau$ Lipschitz. Thus by Jensen's inequality

$$\begin{aligned} |R_\tau(\mathbf{W}) - R_\tau(\mathbf{W}')| &\leq \sqrt{2}\tau \mathbb{E}[|\hat{y}(\mathbf{x}; \mathbf{W}) - \hat{y}(\mathbf{x}; \mathbf{W}')|] \\ &\leq \sqrt{2}\tau \mathbb{E}\left[\left(\sum_{j=1}^m a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle + b_j) - \sum_{j=1}^m a_j \sigma(\langle \mathbf{w}'_j, \mathbf{x} \rangle + b_j)\right)^2\right]^{\frac{1}{2}} \\ &\stackrel{(a)}{\leq} \sqrt{2}\tau \|\mathbf{a}\|_2 \sqrt{\sum_{j=1}^m \mathbb{E}\left[\left(\sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle + b_j) - \sigma(\langle \mathbf{w}'_j, \mathbf{x} \rangle + b_j)\right)^2\right]} \\ &\leq \sqrt{2}\tau\beta_1 \|\mathbf{a}\|_2 \sqrt{\sum_{j=1}^m \mathbb{E}\left[\langle \mathbf{w}_j - \mathbf{w}'_j, \mathbf{x} \rangle^2\right]} \\ &\leq \sqrt{2}\tau\beta_1 \|\mathbf{a}\|_2 \|\mathbf{W} - \mathbf{W}'\|_F \end{aligned} \tag{E.3}$$

where (a) follows from the Cauchy-Schwartz inequality. Note that Equation (E.3) also holds for $|\hat{R}_\tau(\mathbf{W}) - \hat{R}_\tau(\mathbf{W}')|$ when expectation is over the empirical distribution given by the training samples, meaning

$$|\hat{R}_\tau(\mathbf{W}) - \hat{R}_\tau(\mathbf{W}')| \leq \sqrt{2}\tau\beta_1 \|\mathbf{a}\|_2 \sqrt{\sum_{j=1}^m (\mathbf{w}_j - \mathbf{w}'_j)^\top \left(\frac{1}{T} \sum_{t=0}^{T-1} \mathbf{x}^{(t)} \mathbf{x}^{(t)\top}\right) (\mathbf{w}_j - \mathbf{w}'_j)}. \tag{E.4}$$

By Example 6.3 of [68], with probability at least $1 - \delta$ we have

$$\left\| \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{x}^{(t)} \mathbf{x}^{(t)\top} - \mathbf{I}_d \right\|_2 \leq C \left(\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{d + \log(1/\delta)}{n} \right),$$

which completes the proof. \blacksquare

Lemma 11 *Suppose either Assumptions 1.A or 1.B hold. Denote the loss with $\ell(\hat{y}, y) = \ell(\hat{y} - y)$,*

$$\tilde{S} = \left\{ \tilde{\mathbf{W}} \in \mathbb{R}^{m \times k}, \mathbf{a}, \mathbf{b} \in \mathbb{R}^m : \|\mathbf{a}\|_2 \leq \frac{r_a}{\sqrt{m}}, \quad \|\mathbf{b}\|_\infty \leq r_b, \quad \|\tilde{\mathbf{w}}_j\|_2 \leq r_{\tilde{\mathbf{w}}}, \forall 1 \leq j \leq m \right\}$$

and

$$\mathcal{G} = \left\{ (\tilde{\mathbf{x}}, y) \mapsto \ell(\hat{y}(\tilde{\mathbf{x}}; \tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}), y) \wedge \tau : (\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) \in \tilde{S} \right\}$$

for $\tilde{\mathbf{x}} \in \mathbb{R}^k$ and $y \in \mathbb{R}$. Let $\mathfrak{R}(G)$ denote the Rademacher complexity of the function class \mathcal{G} (see Lemma 11 for definition). Then with $\tilde{\mathbf{x}} \sim \mathcal{N}(0, \mathbf{U})$ we have

$$\mathfrak{R}(\mathcal{G}) \leq 2\tau\beta_1(r_{\tilde{\mathbf{w}}}\|\mathbf{U}\|_F + r_b)r_a\sqrt{\frac{2}{T}},$$

where T is the number of samples.

Proof Let $\mathcal{F} = \{(\tilde{\mathbf{x}}, y) \mapsto f_{\mathbf{a}, \tilde{\mathbf{W}}}(\tilde{\mathbf{x}}, y) : (\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) \in \tilde{S}\}$ for $f_{\mathbf{a}, \tilde{\mathbf{W}}}(\tilde{\mathbf{x}}, y) = \hat{y}(\tilde{\mathbf{x}}; \tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) - y$. Define $g(z) := \ell(z) \wedge \tau$, and notice $\mathcal{G} = \{(\tilde{\mathbf{x}}, y) \mapsto g(f_{\mathbf{a}, \tilde{\mathbf{W}}}(\tilde{\mathbf{x}}, y)) : f_{\mathbf{a}, \tilde{\mathbf{W}}} \in \mathcal{F}\}$, and that g is a $\sqrt{2\tau}$ -Lipschitz (thus $\sqrt{2\tau}$ -Lipschitz as well, for $\tau > 1$) function. Then by Talagrand's contraction principle we have $\mathfrak{R}(\mathcal{G}) \leq \sqrt{2\tau}\mathfrak{R}(\mathcal{F})$. Moreover, let $\{\xi_i\}_{1 \leq i \leq n}$ be a sequence of i.i.d. Rademacher

random variables. Then similar to the Rademacher bound of [23]

$$\begin{aligned}
 \mathfrak{R}(\mathcal{F}) &= \mathbb{E} \left[\sup_{(\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) \in \tilde{\mathcal{S}}} \frac{1}{T} \sum_{t=0}^{T-1} \xi_i \left(\mathbf{a}^\top \sigma(\tilde{\mathbf{W}} \tilde{\mathbf{x}}^{(t)} + \mathbf{b}) - y^{(i)} \right) \right] \\
 &= \mathbb{E} \left[\sup_{(\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) \in \tilde{\mathcal{S}}} \frac{1}{T} \sum_{t=0}^{T-1} \xi_i \mathbf{a}^\top \sigma(\tilde{\mathbf{W}} \tilde{\mathbf{x}}^{(t)} + \mathbf{b}) \right] \\
 &\stackrel{(a)}{\leq} \frac{r_a}{T} \mathbb{E} \left[\sup_{(\tilde{\mathbf{W}}, \mathbf{b}) \in \tilde{\mathcal{S}}} \left\| \sum_{t=0}^{T-1} \xi_i \sigma(\tilde{\mathbf{W}} \tilde{\mathbf{x}}^{(t)} + \mathbf{b}) \right\|_\infty \right] \\
 &\leq \frac{r_a}{T} \mathbb{E} \left[\sup_{\|\tilde{\mathbf{w}}\|_2 \leq r_{\tilde{\mathbf{w}}}, |\tilde{b}| \leq r_b} \left| \sum_{t=0}^{T-1} \xi_i \sigma \left(\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}^{(t)} \rangle + \tilde{b} \right) \right| \right] \\
 &\stackrel{(b)}{\leq} \frac{2\beta_1 r_a}{n} \mathbb{E} \left[\sup_{\|\tilde{\mathbf{w}}\|_2 \leq r_{\tilde{\mathbf{w}}}, |\tilde{b}| \leq r_b} \left| \sum_{t=0}^{T-1} \xi_i \left(\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}^{(t)} \rangle + \tilde{b} \right) \right| \right] \\
 &\leq \frac{2\beta_1 r_a}{T} \mathbb{E} \left[\sup_{\|\tilde{\mathbf{w}}\|_2 \leq r_{\tilde{\mathbf{w}}}} \left| \sum_{t=0}^{T-1} \xi_i \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}^{(t)} \rangle \right| + \sup_{|\tilde{b}| \leq r_b} \left| \sum_{t=0}^{T-1} \xi_i \tilde{b} \right| \right] \\
 &\leq \frac{2\beta_1 r_a}{T} \left(r_{\tilde{\mathbf{w}}} \mathbb{E} \left[\left\| \sum_{t=0}^{T-1} \xi_i \tilde{\mathbf{x}}^{(t)} \right\|_2 \right] + r_b \sqrt{n} \right) \\
 &\leq \frac{2\beta_1 (r_{\tilde{\mathbf{w}}} \|\mathbf{U}\|_F + r_b) r_a}{\sqrt{T}},
 \end{aligned}$$

where (a) holds by Hölder's inequality and the fact that $\|\mathbf{a}\|_1 \leq \sqrt{m} \|\mathbf{a}\|_2 \leq r_a$, and (b) follows from the fact that σ is β_1 Lipschitz, thus another application of Talagrand's contraction principle. ■

Proof [Proof of Theorem 5] Let \mathcal{E}_1 denote the event of Lemma 10. We begin with the following decomposition for generalization error which holds on \mathcal{E}_1 ,

$$\begin{aligned}
 R_\tau(\mathbf{W}^T) - \hat{R}_\tau(\mathbf{W}^T) &= R_\tau(\mathbf{W}^T) - R_\tau(\mathbf{W}_\parallel^T) + R_\tau(\mathbf{W}_\perp^T) - \hat{R}_\tau(\mathbf{W}_\parallel^T) + \hat{R}_\tau(\mathbf{W}_\parallel^T) - \hat{R}_\tau(\mathbf{W}^T) \\
 &\leq C\tau\beta_1 \|\mathbf{a}\|_2 \|\mathbf{W}_\perp^T\|_F + R_\tau(\mathbf{W}_\parallel^T) - \hat{R}_\tau(\mathbf{W}_\parallel^T).
 \end{aligned}$$

where the upper bound follows from Lemma 10. Consequently,

$$\sup_{\mathbf{a}, \mathbf{b}} R_\tau(\mathbf{W}^T, \mathbf{a}, \mathbf{b}) - \hat{R}_\tau(\mathbf{W}^T, \mathbf{a}, \mathbf{b}) \leq \frac{C\tau\beta_1 r_a}{\sqrt{m}} \|\mathbf{W}_\perp^T\|_F + \sup_{\mathbf{a}, \mathbf{b}} R_\tau(\mathbf{W}_\parallel^T, \mathbf{a}, \mathbf{b}) - \hat{R}_\tau(\mathbf{W}_\parallel^T, \mathbf{a}, \mathbf{b}). \quad (\text{E.5})$$

We begin by upper bounding the first term. From Theorem 2, on an event \mathcal{E}_2 we have with probability at least $1 - \mathcal{O}(\delta)$

$$\|\mathbf{W}_\perp^T\|_F \lesssim \kappa \sqrt{\frac{d + \log(1/\delta)}{\gamma^2 T}}.$$

Next, we bound the second term in (E.5). For each \mathbf{W} , define $\tilde{\mathbf{W}} := \mathbf{U}^\dagger \mathbf{W}_\parallel$, where \mathbf{U}^\dagger is the Moore–Penrose pseudo-inverse of \mathbf{U} . Then, since we have the representation $\mathbf{W}_\parallel = \mathbf{M}\mathbf{U}$ for some $\mathbf{M} \in \mathbb{R}^{m \times k}$,

$$\tilde{\mathbf{W}}\mathbf{U} = \mathbf{W}_\parallel \mathbf{U}^\dagger \mathbf{U} = \mathbf{M}\mathbf{U}\mathbf{U}^\dagger \mathbf{U} = \mathbf{M}\mathbf{U} = \mathbf{W}_\parallel.$$

Thus, $\mathbf{W}\mathbf{x} = \tilde{\mathbf{W}}\tilde{\mathbf{x}}$ and $\ell(\hat{y}(\mathbf{x}; \mathbf{W}, \mathbf{a}, \mathbf{b}), y) = \ell(\hat{y}(\tilde{\mathbf{x}}; \tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}), y)$ for $\tilde{\mathbf{x}} = \mathbf{U}\mathbf{x}$, when \mathbf{W} is in the principal subspace, i.e. $\mathbf{W} = \mathbf{W}_\parallel$. Let \mathcal{E}_3 denote the event of Lemma 9, on which

$$\|\mathbf{w}_j^T\|_2 \leq \prod_{i=0}^{T-1} \left(1 - \frac{\eta_i \gamma}{m}\right) \|\mathbf{w}_j^0\|_2 + \frac{3\kappa_\infty \sqrt{d}}{\lambda}$$

and consequently

$$\|\tilde{\mathbf{w}}_j^T\|_2 \leq \|\mathbf{U}^\dagger\|_2 \left(\prod_{i=0}^{T-1} \left(1 - \frac{\eta_i \gamma}{m}\right) \|\mathbf{w}_j^0\|_2 + \frac{3\kappa_\infty \sqrt{d}}{\lambda} \right)$$

for any $1 \leq j \leq m$. Define $r_{\tilde{\mathbf{w}}^T}$ as the RHS bound above. Then on \mathcal{E}_3

$$\sup_{\mathbf{a}, \mathbf{b}} R_\tau(\mathbf{W}_\parallel^T) - \hat{R}_\tau(\mathbf{W}_\parallel^T) \leq \sup_{(\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) \in \tilde{S}} R_\tau(\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) - \hat{R}_\tau(\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}),$$

where we recall

$$\tilde{S} := \left\{ \tilde{\mathbf{W}} \in \mathbb{R}^{m \times k}, \mathbf{a}, \mathbf{b} \in \mathbb{R}^m : \|\mathbf{a}\|_2 \leq \frac{r_a}{\sqrt{m}}, \|\mathbf{b}\|_\infty \leq r_b, \|\tilde{\mathbf{w}}_j\|_2 \leq r_{\tilde{\mathbf{w}}^T}, \forall 1 \leq j \leq m \right\}.$$

Additionally define

$$\mathcal{G} = \{(\tilde{\mathbf{x}}, y) \mapsto \ell(\hat{y}(\tilde{\mathbf{x}}; \tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) \wedge \tau : (\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) \in \tilde{S}\}.$$

Then Lemma 24 and Lemma 11 yield

$$\mathbb{E} \left[\sup_{(\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) \in \tilde{S}} R_\tau(\tilde{\mathbf{W}}) - \hat{R}_\tau(\tilde{\mathbf{W}}) \right] \leq 2\mathfrak{R}(\mathcal{G}) \lesssim \tau \beta_1 (r_{\tilde{\mathbf{w}}^T} + r_b) r_a \|\mathbf{U}\|_{\text{F}} \sqrt{\frac{1}{T}}.$$

Besides, as the loss is bounded by τ , by McDiarmid's inequality, on an event \mathcal{E}_4 which happens with probability at least $1 - \mathcal{O}(\delta)$ we have

$$\sup_{(\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) \in \tilde{S}} R_\tau(\tilde{\mathbf{W}}) - \hat{R}_\tau(\tilde{\mathbf{W}}) \leq \mathbb{E} \left[\sup_{(\tilde{\mathbf{W}}, \mathbf{a}, \mathbf{b}) \in \tilde{S}} R(\tilde{\mathbf{W}}) \right] + C\tau \sqrt{\frac{\log(1/\delta)}{T}}.$$

and consequently on $\cap_{i=1}^4 \mathcal{E}_i$

$$\sup_{\mathbf{a}, \mathbf{b}} R_\tau(\mathbf{W}_\parallel^T, \mathbf{a}, \mathbf{b}) - \hat{R}_\tau(\mathbf{W}_\parallel^T, \mathbf{a}, \mathbf{b}) \lesssim \tau \beta_1 (r_{\tilde{\mathbf{w}}^T} + r_b) r_a \|\mathbf{U}\|_{\text{F}} \sqrt{\frac{1}{T}} + \tau \sqrt{\frac{\log(1/\delta)}{T}}.$$

Finally, observe that $\|\mathbf{a}\|_1 \leq \sqrt{m} \|\mathbf{a}\|_2 \leq r_a$ and $\|\mathbf{U}^\dagger\|_2 = \|\mathbf{U}\|_2^{-1} \leq \sqrt{k} \|\mathbf{U}\|_{\text{F}}^{-1}$, thus with probability at least $1 - o(\delta)$,

$$\begin{aligned} \sup_{\mathbf{a}, \mathbf{b}} R_\tau(\mathbf{W}^T, \mathbf{a}, \mathbf{b}) - \hat{R}_\tau(\mathbf{W}^T, \mathbf{a}, \mathbf{b}) &\lesssim \tau \beta_1 r_a \kappa \sqrt{\frac{d + \log(1/\delta)}{\gamma^2 T}} \\ &\quad + \tau \beta_1 r_a \left\{ \left(\frac{t^*}{t^* + T} \right)^2 r_w + \frac{\kappa_\infty}{\lambda} + r_b \right\} \sqrt{\frac{dk}{T}} \\ &\quad + \tau \sqrt{\frac{\log(1/\delta)}{T}}. \end{aligned} \tag{E.6}$$

We remark that in the setting of Theorem 2 which is adapted in Theorem 5, $\|\mathbf{a}\|_\infty \lesssim m^{-1}$, thus $\kappa_\infty \lesssim m^{-1}$. Finally, we observe that $r_w \leq \sqrt{2m}$ with probability at least $1 - \mathcal{O}(\delta)$ over initialization, which completes the proof. \blacksquare

E.2. Proof of Theorem 4

Note that due to the special symmetry in the initialization of Algorithm 1, while training the first layer, all neurons have an identical value, i.e. $\mathbf{w}_j^t = \mathbf{w}^t$ for all j , and that the stochastic gradient with respect to any neuron can be denote by $\nabla \ell = a\partial_1 \ell(\hat{y}, y)\sigma'(\langle \mathbf{w}, \mathbf{x} \rangle + b)\mathbf{x}$. Furthermore, $\nabla_{\mathbf{w}_j} \mathcal{R}_\lambda(\mathbf{W})$ will also be identical for all j , which due to the population gradient formula (3.2), we denote by

$$\nabla \mathcal{R}_\lambda(\mathbf{w}) = (h(\mathbf{w}) + \lambda)\mathbf{w} + \mathfrak{d}(\mathbf{w})\mathbf{u},$$

where $h(\mathbf{w}) = \sum_{j=1}^m \mathcal{H}_{ij}(W)$ and $\mathfrak{d}(\mathbf{w}) = a\mathbb{E}[\partial_{12}^2 \ell(\hat{y}, y)\sigma'(\langle \mathbf{w}, \mathbf{x} \rangle + b)f'(\langle \mathbf{u}, \mathbf{x} \rangle)]$. Additionally, via the arguments in the proof of Lemma 1, it is not difficult to observe $\gamma/m \leq h(\mathbf{w}) + \lambda \lesssim m^{-1}$. Furthermore, similar to the arguments of Lemma 7, $\langle \nabla \ell, \mathbf{v} \rangle$ is $Ca\|\mathbf{v}\|_2$ -sub-Gaussian for any $\mathbf{v} \in \mathbb{R}^d$. Next, we will derive a lower bound for $\langle \mathbf{w}^t, \mathbf{u} \rangle$ to argue that useful features are learned, which first requires obtaining a sharper upper bound on $\|\mathbf{w}^t\|_2$ than that of Lemma 9. This improvement is possible due to considering the special case of $\mathbf{w}_j^t = \mathbf{w}^t$ here.

Lemma 12 *Suppose $t \geq d$. Then,*

$$\|\mathbf{w}^t\|_2 \leq \left(\frac{t^*}{t^* + t} \right) \|\mathbf{w}^0\|_2 + \frac{Cma}{\gamma}$$

with probability at least $1 - \exp(-C(t^* + t))$. In particular, using the union bound, we have

$$\sup_{t \geq t_0} \|\mathbf{w}^t\|_2 \leq \|\mathbf{w}^0\|_2 + \frac{Cma}{\gamma} \lesssim 1$$

with probability at least $1 - \exp(-C(t^* + t_0)) - \exp(-Cd)$.

Proof Let $\mathbf{e}^t := \nabla_{\mathbf{w}} \ell - \nabla_{\mathbf{w}} \mathcal{R}_\lambda$. Then we have

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \nabla_{\mathbf{w}} \mathcal{R}_\lambda - \eta_t \mathbf{e}^t.$$

Recall that $\langle \mathbf{e}^t, \mathbf{v} \rangle$ is $Ca\|\mathbf{v}\|_2$ -sub-Gaussian, and \mathcal{F}_t is the sigma algebra generated by $\{\mathbf{w}^j\}_{0 \leq j \leq t}$. Let $\boldsymbol{\omega}^t := \mathbf{w}^t - \eta_t \nabla_{\mathbf{w}} \mathcal{R}_\lambda$. Then, for any $0 \leq s \lesssim \frac{\gamma}{\eta_t a^2}$,

$$\begin{aligned} \mathbb{E}[\exp(s\|\mathbf{w}^{t+1}\|_2^2) | \mathcal{F}_0] &= \mathbb{E}[\exp(s\|\boldsymbol{\omega}^t\|_2^2 - 2s\eta_t \langle \boldsymbol{\omega}^t, \mathbf{e}^t \rangle + s\eta_t^2 \|\mathbf{e}^t\|_2^2) | \mathcal{F}_0] \\ &\leq \mathbb{E}[\exp(s\|\boldsymbol{\omega}^t\|_2^2) \mathbb{E}[\exp(-4s\eta_t \langle \boldsymbol{\omega}^t, \mathbf{e}^t \rangle) | \mathcal{F}_t]^{\frac{1}{2}} \mathbb{E}[\exp(2s\eta_t^2 \|\mathbf{e}^t\|_2^2) | \mathcal{F}_t]^{\frac{1}{2}} | \mathcal{F}_0]. \end{aligned}$$

By sub-Gaussianity of $\langle \boldsymbol{\omega}^t, \mathbf{e}^t \rangle$ we have $\mathbb{E}[\exp(-4s\eta_t \langle \boldsymbol{\omega}^t, \mathbf{e}^t \rangle) | \mathcal{F}_t] \leq \exp(Cs^2\eta_t^2 a^2 \|\boldsymbol{\omega}^t\|_2^2)$. Moreover, as $\|\nabla \ell\|_2 \leq |a|\|\mathbf{x}\|_2$, by Jensen's inequality

$$\|\mathbf{e}^t\|_2^2 \leq 2\|\nabla \ell\|_2^2 + 2\mathbb{E}[\|\nabla \ell\|_2^2] \leq 2a^2(\|\mathbf{x}\|_2^2 + d).$$

Thus $\mathbb{E}[\exp(2s\eta_t^2\|\mathbf{e}^t\|_2^2) \mid \mathcal{F}_t] \leq \exp(Cs\eta_t^2a^2d)$ for $s \lesssim \frac{1}{\eta_t^2a^2}$ (which holds by $s \lesssim \frac{\gamma}{\eta_t a^2}$, see the proof of Lemma 8 for more details), i.e. we have

$$\mathbb{E}[\exp(s\|\mathbf{w}^{t+1}\|_2^2) \mid \mathcal{F}_0] \leq \mathbb{E}[\exp(s(1 + Cs\eta_t^2a^2)\|\mathbf{w}^t\|_2^2 + Cs\eta_t^2a^2d) \mid \mathcal{F}_0].$$

Recall that by our choice of η_t , $0 \leq (1 - \eta_t(h(\mathbf{w}^t) + \lambda)) \leq 1 - \frac{\eta_t\gamma}{m}$ (cf. proof of Lemma 8), and $\mathbf{w}^t = (1 - \eta_t(h(\mathbf{w}^t) + \lambda))\mathbf{w}^t - \eta_t\mathfrak{d}(\mathbf{w}^t)\mathbf{u}$. As $\|\mathbf{u}\|_2 = 1$ and $|\mathfrak{d}(\mathbf{w}^t)| \lesssim |a|$, we have

$$\begin{aligned} \|\mathbf{w}^t\|_2^2 &\leq (1 - \frac{\eta_t\gamma}{m})^2\|\mathbf{w}^t\|_2^2 + Ca^2\eta_t^2 + 2\eta_tCa(1 - \frac{\eta_t\gamma}{m})\|\mathbf{w}^t\|_2 \\ &\stackrel{(a)}{\leq} (1 - \frac{\eta_t\gamma}{m})^2\|\mathbf{w}^t\|_2^2 + \eta_t\left(\frac{2Cma^2}{\gamma} + \frac{\gamma}{2m}(1 - \frac{\eta_t\gamma}{m})^2\|\mathbf{w}^t\|_2^2\right) + Ca^2\eta_t^2 \\ &\stackrel{(b)}{\leq} (1 - \frac{3\eta_t\gamma}{2m})\|\mathbf{w}^t\|_2^2 + \frac{Cm\eta_t a^2}{\gamma} + Ca^2\eta_t^2. \end{aligned}$$

where (a) holds by Young's inequality and (b) holds for $\eta_t\gamma \lesssim 1$ with a sufficiently small absolute constant. Therefore, for $s \lesssim \frac{\gamma}{m\eta_t a^2}$,

$$\mathbb{E}[\exp(s\|\mathbf{w}^{t+1}\|_2^2) \mid \mathcal{F}_0] \leq \mathbb{E}\left[\exp\left((1 - \frac{\eta_t\gamma}{m})\|\mathbf{w}^t\|_2^2 + \frac{Csm\eta_t a^2}{\gamma} + Cs\eta_t^2a^2d\right) \mid \mathcal{F}_0\right].$$

Expanding the recursion yields,

$$\mathbb{E}[\exp(s\|\mathbf{w}^t\|_2^2) \mid \mathcal{F}_0] \leq \exp\left(s\left(\frac{t^*}{t^* + t}\right)^2\|\mathbf{w}^0\|_2^2 + \frac{Csm^2a^2(t+d)}{\gamma^2(t^* + t)}\right).$$

Finally, we apply a Chernoff bound with the maximum choice of $s \lesssim \frac{\gamma}{\eta_t a^2}$, and combine it with the fact that $\|\mathbf{w}^0\|_2 \lesssim 1$ with probability at least $1 - \exp(-Cd)$. \blacksquare

Lemma 13 *Suppose $mab < 1 - |f(0)|$. Then, we have $|\langle \mathbf{w}^t, \mathbf{u} \rangle| \gtrsim 1$ with probability at least $1 - 2\exp(-Ct) - \exp(-Cd)$.*

Proof We will only prove for the case where f is increasing as the case for decreasing f is similar. We begin by proving an upper bound for $\mathfrak{d}(\mathbf{w})$ when $\|\mathbf{w}\|_2 \lesssim 1$. By the triangle inequality,

$$|\hat{y} - y| \leq |\hat{y}| + |f(0)| + |f(\langle \mathbf{u}, \mathbf{x} \rangle) - f(0)| + |\epsilon|.$$

Furthermore, $|\hat{y}| \leq ma(|\langle \mathbf{w}, \mathbf{x} \rangle| + b)$. Thus, for

$$|\langle \mathbf{w}, \mathbf{x} \rangle| \leq \left(\frac{1 - |f(0)| - mab}{2ma}\right) \wedge b,$$

$|\langle \mathbf{u}, \mathbf{x} \rangle| \lesssim 1$ and $|\epsilon| \lesssim 1$ for sufficiently small absolute constants, we have $|\hat{y} - y| \leq 1$ hence $\partial_{12}^2\ell(\hat{y}, y) = -1$. Then we have,

$$\begin{aligned} \mathfrak{d}(\mathbf{w}) &= a\mathbb{E}[\partial_{12}^2\ell(\hat{y}, y)\sigma'(\langle \mathbf{w}, \mathbf{x} \rangle + \beta)f'(\langle \mathbf{u}, \mathbf{x} \rangle)] \\ &\lesssim a\mathbb{E}[\mathbf{1}(|\epsilon| \lesssim 1)\mathbf{1}(|\langle \mathbf{w}, \mathbf{x} \rangle| \lesssim 1)\mathbf{1}(|\langle \mathbf{u}, \mathbf{x} \rangle| \lesssim 1)f'(\langle \mathbf{u}, \mathbf{x} \rangle)] \\ &= a\mathbb{E}[\mathbf{1}(|\epsilon| \lesssim 1)]\mathbb{E}[\mathbf{1}(|\langle \mathbf{w}, \mathbf{x} \rangle| \lesssim 1)\mathbf{1}(|\langle \mathbf{u}, \mathbf{x} \rangle| \lesssim 1)f'(\langle \mathbf{u}, \mathbf{x} \rangle)] \\ &\lesssim -a. \end{aligned}$$

where the last line is obtained by considering supremum over $\|\mathbf{w}\|_2 \lesssim 1$.

Let $A_t = \{\sup_{t_0 \leq t' \leq t} \|\mathbf{w}^{t'}\|_2 \lesssim 1\}$. Then,

$$\begin{aligned} \mathbb{E}[\exp(-s\langle \mathbf{w}^{t+1}, \mathbf{u} \rangle) \mathbf{1}_{A_{t+1}}] &\leq \mathbb{E}[\exp(-s\langle \mathbf{w}^{t+1}, \mathbf{u} \rangle) \mathbf{1}_{A_t}] \\ &\leq \mathbb{E}[\exp(-s\langle \mathbf{w}^t, \mathbf{u} \rangle + s\eta_t \langle \nabla \ell + \lambda \mathbf{w}^t, \mathbf{u} \rangle) \mathbf{1}_{A_t}] \\ &\stackrel{(a)}{\leq} \mathbb{E}[\exp(-s\langle \mathbf{w}^t, \mathbf{u} \rangle + s\eta_t \langle \nabla \mathcal{R}_\lambda, \mathbf{u} \rangle + Cs^2\eta_t^2 a^2) \mathbf{1}_{A_t}] \\ &\stackrel{(b)}{=} \mathbb{E}[\exp(-s(1 - \eta_t(h(\mathbf{w}^t) + \lambda)) \langle \mathbf{w}^t, \mathbf{u} \rangle + s\eta_t(\mathfrak{d}(\mathbf{w}^t) + Cs\eta_t)) \mathbf{1}_{A_t}] \\ &\stackrel{(c)}{\leq} \exp(-Cs\eta_t a) \mathbb{E}[\exp(-s(1 - \eta_t(h(\mathbf{w}^t) + \lambda)) \langle \mathbf{w}^t, \mathbf{u} \rangle) \mathbf{1}_{A_t}], \end{aligned}$$

where (a) follows from the sub-Gaussianity of the stochastic noise in the gradient, (b) follows since $\langle \nabla \mathcal{R}_\lambda(\mathbf{w}^t), \mathbf{u} \rangle = \mathfrak{d}(\mathbf{w}^t)$ by definition, and (c) holds for $s \lesssim \eta_t^{-1}$ with a sufficiently small absolute constant. Notice that by the condition on t^* inherited from Theorem 2, $1 - \eta_t(h(\mathbf{w}^t) + \lambda) > 0$, and since $s(1 - \eta_t(h(\mathbf{w}^t) + \lambda)) \leq s(1 - \frac{\eta_t \gamma}{m})$, we can expand the recursion,

$$\begin{aligned} \mathbb{E}[\exp(-s\langle \mathbf{w}^t, \mathbf{u} \rangle) \mathbf{1}_{A_t}] &\leq \mathbb{E} \left[\exp \left(-Csa \sum_{i=t_0}^{t-1} \eta_j \prod_{j=t-i}^{t-1} (1 - \frac{\eta_t \gamma}{m}) + s \prod_{i=t_0}^{t-1} (1 - \frac{\eta_t \gamma}{m}) |\langle \mathbf{w}^{t_0}, \mathbf{u} \rangle| \right) \mathbf{1}_{A_{t_0}} \right] \\ &\leq \mathbb{E} \left[\exp \left(-Cs \left(1 - \left(\frac{t^* + t_0}{t^* + t} \right)^2 \right) + Cs \left(\frac{t^* + t_0}{t^* + t} \right)^2 \right) \right]. \end{aligned}$$

where in the second inequality we used $a \asymp m^{-1}$ and $\gamma \asymp 1$. Applying the Chernoff bound implies that $\langle \mathbf{w}^t, \mathbf{u} \rangle \gtrsim 1$ with probability at least $1 - \mathbb{P}(A_t^C) - \exp(-Ct) \leq 1 - \exp(-C(t^* + t_0)) - \exp(-Cd) - \exp(-Ct)$. Finally the result follows by letting $t_0 = \frac{t}{2}$. \blacksquare

We have proven that $|\langle \mathbf{w}^t, \mathbf{u} \rangle| \gtrsim 1$ while $\|\mathbf{w}_\perp^t\|_2 \rightarrow 0$. This fact shows that the features learned in the first layer are useful. What remains to be shown is an approximation result, such that for a carefully constructed second layer, the network can approximate polynomials of the desired type. This type of approximation using random biases has been adopted from [23]. We first present an approximation result using infinite neurons.

Lemma 14 *Let $0 < |\alpha| \leq r$ and $b \sim \text{Unif}(-2r\Delta, 2r\Delta)$. For any smooth $f : \mathbb{R} \rightarrow \mathbb{R}$, let $\tilde{f}_\alpha : \mathbb{R} \rightarrow \mathbb{R}$ be a smooth function such that $\tilde{f}_\alpha(z) = f(z)$ for $|z| \leq \frac{r\Delta}{|\alpha|}$ and $\tilde{f}_\alpha(-\frac{2r\Delta}{\alpha}) = \tilde{f}'_\alpha(-\frac{2r\Delta}{\alpha}) = 0$. Then, for $|z| \leq \Delta$ we have*

$$\mathbb{E}_b \left[4r\Delta \tilde{f}_\alpha'' \left(-\frac{b}{\alpha} \right) \sigma(\alpha z + b) \right] = f(z).$$

Proof Using integration by parts, we have

$$\begin{aligned} \mathbb{E}_b \left[4r\Delta \tilde{f}_\alpha'' \left(-\frac{b}{\alpha} \right) \sigma(\alpha z + b) \right] &= \int_{-2r\Delta}^{2r\Delta} \tilde{f}_\alpha'' \left(-\frac{b}{\alpha} \right) (\alpha z + b) db \\ &= -\tilde{f}'_\alpha \left(-\frac{2r\Delta}{\alpha} \right) \left(z + \frac{2r\Delta}{\alpha} \right) + \int_{-\frac{2r\Delta}{\alpha}}^z \tilde{f}'_\alpha(b) db \\ &= \tilde{f}_\alpha(z) = f(z). \end{aligned}$$

■

Now, by a concentration argument, we state an approximation result with finitely many neurons.

Lemma 15 *Let $r^* \leq |\alpha_j| \leq r$ and $b_j \sim \text{Unif}(-2r\Delta, 2r\Delta)$. Then there exists $a(\alpha_j, b_j)$ such that for any fixed $z \in [-\Delta, \Delta]$, with probability at least $1 - \delta$ over the choice of (b_j) , we have*

$$\left| \sum_{j=1}^m a(\alpha_j, b_j) \sigma(\alpha_j z + b_j) - f(z) \right| \lesssim (r\Delta)^2 \sup_{|z| \leq \frac{r\Delta}{r^*}} |f''(z)| \sqrt{\frac{\log(1/\delta)}{m}}.$$

Moreover, $\|\mathbf{a}\|_2 \lesssim \frac{r\Delta}{\sqrt{m}} \sup_{|z| \leq \frac{r\Delta}{r^*}} |f''(z)|$.

Proof Let $\tilde{f}_\alpha(z)$ be a candidate in Lemma 14, which can be obtained by e.g. extending f with suitable polynomials (notice that \tilde{f}_α only needs to be twice differentiable on its domain). Now choose $a_j = 4\frac{r\Delta}{m} \tilde{f}''_{\alpha_j}(-\frac{b_j}{\alpha_j})$. Then Lemma 14 ensures that

$$\mathbb{E}_b[a(\alpha_j, b_j) \sigma(\alpha_j z + b_j)] = f(z).$$

As $\tilde{f}_\alpha(z) = f(z)$ for all $|z| \leq \frac{r\Delta}{|\alpha_j|}$, we have $\tilde{f}''_{\alpha_j}(-\frac{b_j}{\alpha_j}) = f''_{\alpha_j}(-\frac{b_j}{\alpha_j})$. It immediately follows that $\|\mathbf{a}\|_2 \leq \frac{Cr\Delta}{\sqrt{m}} \sup_{|z| \leq \frac{r\Delta}{r^*}} |f''(z)|$ and $|a_j \sigma(\alpha_j z + b_j)| \leq \frac{C(r\Delta)^2}{m} \sup_{|z| \leq \frac{r\Delta}{r^*}} |f''(z)|$. Applying the Hoeffding's inequality finishes the proof. ■

In the following lemma, we will briefly record useful properties of \mathbf{W}^T which will be of help for invoking the above approximation results and providing guarantees when the second layer is optimized by SGD. Through the rest of the proof, we will add the mild assumption that $d \gtrsim \log(1/\delta)$. Otherwise, we need to add e^{-Cd} to the probability of failure in Theorem 4.

Lemma 16 *Suppose $T \gtrsim d + \log(1/\delta)$. Then with probability at least $1 - \delta$ over the choice of $(b_j)_{1 \leq j \leq m}$ and $\{(\mathbf{x}^{(t)}, y^{(t)})\}_{t=0}^{T-1}$, the following statements hold:*

1. $\|\mathbf{w}_j^T\|_2 \asymp |\langle \mathbf{w}_j^T, \mathbf{u} \rangle| \asymp 1$ for all $1 \leq j \leq m$.
2. $\|\frac{1}{T} \sum_{t=0}^{T-1} \mathbf{x}^{(t)} (\mathbf{x}^{(t)})^\top\|_2 \lesssim 1$.
3. $\|\mathbf{W}_\perp^T\|_F \lesssim \sqrt{\frac{m(d + \log(1/\delta))}{T}}$.
4. $|\langle \mathbf{u}, \mathbf{x}^{(t)} \rangle| \lesssim \Delta$ for all $1 \leq j \leq m$ and $0 \leq t \leq T - 1$.
5. $\|\mathbf{W}^T \mathbf{x}^{(t)}\|_2 \lesssim \sqrt{m}(\sqrt{d} + \Delta)$ for all $0 \leq t \leq T - 1$.

Proof We will show that each of the events hold with probability (w.p.) at least $1 - \mathcal{O}(\delta)$. Recall from Lemma 12 that $\|\mathbf{w}_j^T\|_2 \lesssim 1$ for all j w.p. $\geq 1 - \mathcal{O}(\delta)$, which implies the same for $\langle \mathbf{w}_j^T, \mathbf{u} \rangle$. On the other hand, from Lemma 13, $|\langle \mathbf{w}_j^T, \mathbf{u} \rangle| \gtrsim 1$ for all j w.p. $\geq 1 - \mathcal{O}(\delta)$. Combining these events implies that $|\langle \mathbf{w}_j^T, \mathbf{u} \rangle| \asymp 1$. The fact that $\|\frac{1}{T} \sum_{t=0}^{T-1} \mathbf{x}^{(t)} (\mathbf{x}^{(t)})^\top\|_2 \lesssim 1$ w.p. $\geq 1 - \mathcal{O}(\delta)$ for $T \gtrsim$

$d + \log(1/\delta)$ follows from the statement of Lemma 23. Furthermore, $\|\mathbf{W}_\perp^T\|_F \lesssim \sqrt{\frac{m(d + \log(1/\delta))}{T}}$ w.p. $1 - \mathcal{O}(\delta)$ follows from Theorem 2. Note that as $\langle \mathbf{u}, \mathbf{x}^{(t)} \rangle \sim \mathcal{N}(0, 1)$, by the choice of Δ , we have $\langle \mathbf{u}, \mathbf{x}^{(t)} \rangle \gtrsim 1$ w.p. $\geq 1 - \mathcal{O}(\delta/T)$, thus $|\langle \mathbf{u}, \mathbf{x}^{(t)} \rangle| \lesssim \Delta$ w.p. $\geq 1 - \mathcal{O}(\delta)$ by a union bound. Finally, we have

$$\begin{aligned} \|\mathbf{W}^T \mathbf{x}^{(t)}\|_2 &\leq \|\mathbf{W}_\parallel^T \mathbf{x}^{(t)}\|_2 + \|\mathbf{W}_\perp^T \mathbf{x}^{(t)}\|_2 \lesssim \sqrt{m} |\langle \mathbf{u}, \mathbf{x}^{(t)} \rangle| + \sqrt{\frac{m(d + \log(1/\delta))}{T}} \|\mathbf{x}^{(t)}\|_2 \\ &\lesssim \sqrt{m} |\langle \mathbf{u}, \mathbf{x}^{(t)} \rangle| + \sqrt{m} \|\mathbf{x}^{(t)}\|_2 \end{aligned}$$

The first term is already bounded by $\sqrt{m}\Delta$ with probability at least $1 - \mathcal{O}(\delta)$. Moreover, recall that $\|\mathbf{x}^{(t)}\|_2 - \mathbb{E}[\|\mathbf{x}^{(t)}\|_2]$ is 1-sub-Gaussian, thus by the union bound $\|\mathbf{x}^{(t)}\|_2 - \sqrt{d} \lesssim \sqrt{\log(T/\delta)} \lesssim \Delta$ for all $0 \leq t \leq T-1$. Thus w.p. $\geq 1 - \mathcal{O}(\delta)$ we have $\|\mathbf{W}^T \mathbf{x}^{(t)}\|_2 \lesssim \sqrt{m}(\sqrt{d} + \Delta)$ which completes the proof. \blacksquare

From this point onwards, we will denote the Huber loss with $\ell_H(\hat{y}, y) = \ell(\hat{y} - y)$. Notice that ℓ_H is 1-Lipschitz.

Lemma 17 *Recall*

$$\hat{R}(\mathbf{W}^T, \mathbf{a}, \mathbf{b}) = \frac{1}{T} \sum_{t=0}^{T-1} \ell_H \left(\sum_{j=1}^m a_j \sigma(\langle \mathbf{w}_j^T, \mathbf{x}^{(t)} \rangle + b_j) - f(\langle \mathbf{u}, \mathbf{x}^{(t)} \rangle) - \epsilon^{(t)} \right),$$

the empirical risk of \mathbf{W}^T given by Algorithm 1. Let $\Delta \asymp \sqrt{\log(\frac{T}{\delta})}$, $\Delta_* := \Delta \sup_{|z| \lesssim \Delta} |f''(z)|$, and $b_j \stackrel{i.i.d.}{\sim} \text{Unif}(-\Delta, \Delta)$. Then, with probability at least $1 - \delta$ (over the randomness of $(b_j)_{1 \leq j \leq m}$ and $\{\mathbf{x}^{(t)}, y^{(t)}\}_{t=0}^{T-1}$ hence \mathbf{W}^T), for $T \gtrsim d + \log(1/\delta)$, there exists \mathbf{a}^* with $\|\mathbf{a}^*\|_2 \lesssim \frac{\Delta_*}{\sqrt{m}}$ such that

$$\hat{R}(\mathbf{W}^T, \mathbf{a}^*, \mathbf{b}) - \mathbb{E}[\ell_H(\epsilon)] \lesssim \Delta_* \left(\Delta \sqrt{\frac{\log(1/\delta)}{m}} + \Delta_* \sqrt{\frac{d + \log(1/\delta)}{T}} \right).$$

Proof We will condition the following discussion on the event of Lemma 16. Let $\alpha_j = \langle \mathbf{w}_j^T, \mathbf{u} \rangle$, and let \mathbf{a}^* be constructed according to Lemma 15. By the Lipschitzness of the Huber loss, for an individual sample (\mathbf{x}, y) we have

$$\begin{aligned} \ell(\hat{y}(\mathbf{x}; \mathbf{W}^T, \mathbf{a}^*, \mathbf{b}) - f(\langle \mathbf{u}, \mathbf{x} \rangle) - \epsilon) &\leq \ell_H(\epsilon) + |\hat{y}(\mathbf{x}; \mathbf{W}^T, \mathbf{a}^*, \mathbf{b}) - f(\langle \mathbf{u}, \mathbf{x} \rangle)| \\ &\leq \ell_H(\epsilon) + |\hat{y}(\mathbf{x}; \mathbf{W}^T, \mathbf{a}^*, \mathbf{b}) - \hat{y}(\mathbf{x}; \mathbf{W}_\parallel^T, \mathbf{a}^*, \mathbf{b})| \\ &\quad + |\hat{y}(\mathbf{x}; \mathbf{W}_\parallel^T, \mathbf{a}^*, \mathbf{b}) - f(\langle \mathbf{u}, \mathbf{x} \rangle)|. \end{aligned}$$

Moreover, by the Cauchy-Schwartz inequality

$$\begin{aligned} |\hat{y}(\mathbf{x}; \mathbf{W}^T, \mathbf{a}^*, \mathbf{b}) - \hat{y}(\mathbf{x}; \mathbf{W}_\parallel^T, \mathbf{a}^*, \mathbf{b})| &\leq \|\mathbf{a}^*\|_2 \sqrt{\sum_{j=1}^m \left(\sigma(\langle \mathbf{w}_j^T, \mathbf{x} \rangle + b_j) - \sigma(\langle (\mathbf{w}_j^T)_\parallel, \mathbf{x} \rangle + b_j) \right)^2} \\ &\leq \|\mathbf{a}^*\|_2 \sqrt{\sum_{j=1}^m \langle (\mathbf{w}_j^T)_\perp, \mathbf{x} \rangle^2}. \end{aligned}$$

Additionally, since $\|\frac{1}{T} \sum_{t=0}^{T-1} \mathbf{x}^{(t)}(\mathbf{x}^{(t)})^\top\|_2 \lesssim 1$, by Jensen's inequality,

$$\begin{aligned} \sum_{t=0}^{T-1} \frac{1}{T} |\hat{y}(\mathbf{x}; \mathbf{W}^T, \mathbf{a}^*, \mathbf{b}) - \hat{y}(\mathbf{x}; \mathbf{W}_\perp^T, \mathbf{a}^*, \mathbf{b})| &\leq \|\mathbf{a}^*\|_2 \sqrt{\frac{1}{T} \sum_{t=0}^{T-1} \|\mathbf{W}_\perp^T \mathbf{x}^{(t)}\|_F^2} \\ &\lesssim \|\mathbf{a}^*\|_2 \|\mathbf{W}_\perp^T\|_F \\ &\lesssim \Delta_* \sqrt{\frac{d + \log(1/\delta)}{T}} \end{aligned}$$

On the other hand, let $z^{(t)} := \langle \mathbf{u}, \mathbf{x}^{(t)} \rangle \leq \Delta$. Then, we can apply Lemma 15, which states that with probability $1 - \mathcal{O}(\delta)$ over the choice of $(b_j)_{1 \leq j \leq m}$,

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} |\hat{y}(\mathbf{x}^{(t)}; \mathbf{W}_\perp^T, \mathbf{a}^*, \mathbf{b}) - f(\langle \mathbf{u}, \mathbf{x}^{(t)} \rangle)| &\leq \frac{1}{T} \sum_{t=0}^{T-1} |\sum_{j=1}^m a_j^* \sigma(\alpha_j z^{(t)} + b_j) - f(z^{(t)})| \\ &\lesssim \Delta \Delta_* \sqrt{\frac{\log(1/\delta)}{m}}. \end{aligned}$$

Combining the events above, we have with probability at least $1 - \delta$,

$$\hat{R}(\mathbf{W}^T, \mathbf{a}^*, \mathbf{b}) - \frac{1}{T} \sum_{t=0}^{T-1} \ell_H(\epsilon^{(t)}) \lesssim \Delta_* \left(\Delta \sqrt{\frac{\log(1/\delta)}{m}} + \sqrt{\frac{d + \log(1/\delta)}{T}} \right).$$

The final step to apply a concentration for $\sum_{t=0}^{T-1} \ell_H(\epsilon^{(t)})$. Note that as $\ell_H(\epsilon) \leq |\epsilon|$, if $|\epsilon|$ is C -sub-Gaussian, then $\ell_H(\epsilon) - \mathbb{E}[\ell_H(\epsilon)]$ is also sub-Gaussian with an absolute constant as parameter (can be verified e.g. by Lemma 21). Then a sub-Gaussian concentration bound implies that $\frac{1}{T} \sum_{t=0}^{T-1} \ell_H(\epsilon^{(t)}) - \mathbb{E}[\ell_H(\epsilon)] \gtrsim \sqrt{\frac{\log(1/\delta)}{T}}$, which finishes the proof. \blacksquare

Let $\mathbb{E}_S[\cdot]$ denote expectation w.r.t. the random sampling of SGD used to train a , hence conditioned on $\{\mathbf{x}^{(t)}, y^{(t)}\}_{t=0}^{T-1}$. Also, define the stochastic noise in the gradient w.r.t. a as

$$\boldsymbol{\epsilon}_a^t = \nabla_a \ell(\hat{y}(\mathbf{x}^{(i_t)}; \mathbf{W}^T, \mathbf{a}^t, \mathbf{b}) - y^{(i_t)}) - \nabla_a \hat{R}(\mathbf{W}^T, \mathbf{a}^t, \mathbf{b}).$$

Notice that $\mathbb{E}_S[\boldsymbol{\epsilon}_a^t] = 0$.

Lemma 18 *On the event of Lemma 16 and with $(b_j) \stackrel{i.i.d.}{\sim} \text{Unif}(-\Delta, \Delta)$, consider the mapping $\mathbf{a} \mapsto \hat{\mathcal{R}}_{\lambda'}(\mathbf{a})$. Then, $\nabla_a^2 \hat{\mathcal{R}}_{\lambda'}(\mathbf{a}) \lesssim m\Delta^2 + \lambda'$, and $\|\boldsymbol{\epsilon}_a^t\|_2 \lesssim \sqrt{mr}(\sqrt{d} + \Delta)$.*

Proof For $\nabla_a^2 \hat{R}(\mathbf{a})$, and any $\mathbf{v} \in \mathbb{R}^m$ with $\|\mathbf{v}\|_2 = 1$, we have the following computation:

$$\begin{aligned} \langle \mathbf{v}, \nabla_a^2 \hat{\mathcal{R}}_{\lambda'}(\mathbf{a}) \mathbf{v} \rangle &= \frac{2}{T} \sum_{t=0}^{T-1} \partial_{12}^2 \ell(\hat{y}, y) \mathbf{v}^\top \sigma(\mathbf{W}^T \mathbf{x}^{(t)} + \mathbf{b}) \sigma(\mathbf{W}^T \mathbf{x}^{(t)} + \mathbf{b})^\top \mathbf{v} + \lambda' \\ &\leq \frac{1}{T} \sum_{t=0}^{T-1} \langle \mathbf{v}, \mathbf{W}^T \mathbf{x}^{(t)} \rangle^2 + 2\|\mathbf{b}\|_2^2 + \lambda' \\ &\stackrel{(a)}{\lesssim} \|\mathbf{W}^T\|_F + \|\mathbf{b}\|_2^2 + \lambda' \\ &\lesssim m\Delta^2 + \lambda' \end{aligned}$$

where (a) holds since $\|\frac{1}{T} \sum_{t=0}^{T-1} \mathbf{x}^{(t)} \mathbf{x}^{(t)\top}\|_2 \lesssim 1$. Thus $\nabla_{\mathbf{a}}^2 \hat{\mathcal{R}}_{\lambda'}(\mathbf{a}) \lesssim m\Delta^2 + \lambda'$. On the other hand, as $\|\mathbf{W}^T \mathbf{x}^{(t)}\|_2 \lesssim \sqrt{m}(\sqrt{d} + \Delta)$ for all $0 \leq t \leq T-1$, we have

$$\|\boldsymbol{\epsilon}_{\mathbf{a}}^t\|_2 \leq 2\|\nabla_{\mathbf{a}} \ell\|_2 \leq 2\|\mathbf{W}^T \mathbf{x}^{(t)} + \mathbf{b}\|_2 \lesssim \sqrt{m}(\sqrt{d} + \Delta).$$

■

Now we can analyze the SGD run on the second layer \mathbf{a} to give a high probability statement for $\hat{\mathcal{R}}_{\lambda'}(\mathbf{a}^T)$. As $\hat{\mathcal{R}}_{\lambda'}(\mathbf{a})$ is a smooth and strongly convex function of \mathbf{a} , we will state the following well-known elementary convergence result of SGD for smooth and strongly convex functions with bounded noise, which we present in a high-probability framework suitable for our analysis.

Lemma 19 *Let $\mathcal{R} : \mathbb{R}^m \rightarrow \mathbb{R}$ be a μ -strongly convex function satisfying $\mu \mathbf{I}_m \preceq \nabla_{\mathbf{a}}^2 \mathcal{R}(\mathbf{a}) \preceq L \mathbf{I}_m$. Suppose we run the SGD iterates $\mathbf{a}^{t+1} = \mathbf{a}^t - \eta_t \mathbf{g}^t$ with $\mathbb{E}[\mathbf{g}^t | \mathbf{a}^t] = \nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t)$ and $\|\mathbf{g}^t\|_2 \leq G$. Choose $\eta_t = \frac{2t+1}{\mu(t+1)^2}$. Then with probability at least $1 - \delta$*

$$\mathcal{R}(\mathbf{a}^T) - \mathcal{R}^* \leq \frac{\mathcal{R}^0}{T^2} + \frac{CLG^2}{\mu^2 T} + \frac{CG^2 \log(1/\delta)}{\mu T},$$

where $\mathcal{R}^* = \arg \min_{\mathbf{a}} \mathcal{R}(\mathbf{a})$.

Proof Let $\boldsymbol{\epsilon}^t = \mathbf{g}^t - \nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t)$ denote the stochastic noise. By the smoothness property of \mathcal{R} , we have

$$\begin{aligned} \mathcal{R}(\mathbf{a}^{t+1}) - \mathcal{R}^* &\leq \mathcal{R}(\mathbf{a}^t) - \mathcal{R}^* - \eta_t \langle \nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t), \nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t) + \boldsymbol{\epsilon}^t \rangle + \frac{L\eta_t^2}{2} \|\mathbf{g}^t\|_2^2 \\ &\leq \mathcal{R}(\mathbf{a}^t) - \mathcal{R}^* - \eta_t \|\nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t)\|_2^2 - \eta_t \langle \nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t), \boldsymbol{\epsilon}^t \rangle + \frac{L\eta_t^2 G^2}{2}. \end{aligned}$$

Notice that by Jensen's inequality, $\|\nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t)\|_2 \leq G$, thus $\|\boldsymbol{\epsilon}^t\|_2 \leq 2G$ and the zero-mean random variable $\|\langle \nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t), \boldsymbol{\epsilon}^t \rangle\|_2$ is $2G\|\nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t)\|_2$ -sub-Gaussian conditioned on \mathbf{a}^t . Now, we can establish the following recursive bound on the MGF of $\mathcal{R}^t := \mathcal{R}(\mathbf{a}^t) - \mathcal{R}^*$. For $0 \leq s \leq \frac{1}{4\eta_t G^2}$ we have

$$\begin{aligned} \mathbb{E}\left[e^{s\mathcal{R}^{t+1}}\right] &\leq \mathbb{E}\left[\exp\left(s\mathcal{R}^t - s\eta_t \|\nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t)\|_2^2 - s\eta_t \langle \nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t), \boldsymbol{\epsilon}^t \rangle + \frac{\eta_t^2 LG^2}{2}\right)\right] \\ &\leq \mathbb{E}\left[\exp\left(s\mathcal{R}^t - s\eta_t(1 - 2s\eta_t G^2) \|\nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a}^t)\|_2^2 + \frac{LG^2 \eta_t^2}{2}\right)\right] \\ &\stackrel{(a)}{\leq} \mathbb{E}\left[\exp\left(s(1 - \eta_t \mu) \mathcal{R}^t + \frac{LG^2 \eta_t^2}{2}\right)\right] \end{aligned}$$

where (a) follows since $\mathcal{R}(\mathbf{a})$ is strongly convex thus satisfies the Polyak-Łojasiewicz inequality $2\lambda(\mathcal{R}(\mathbf{a}) - \mathcal{R}^*) \leq \|\nabla_{\mathbf{a}} \mathcal{R}(\mathbf{a})\|_2^2$. As $s(1 - \eta_t \mu) \leq \frac{1}{4\eta_{t-1} G^2}$ (cf. (D.5)), we can expand the recursion and have

$$\mathbb{E}\left[\exp(s\mathcal{R}^t)\right] \leq \exp\left(s\left(\frac{t^*}{t^* + t}\right)^2 \mathcal{R}^0 + \frac{16LG^2}{\mu^2(t^* + t)}\right).$$

Finally, applying a Chernoff bound using $s = (4\eta_{t-1} G^2)^{-1}$ concludes the proof. ■

We are finally in a position to complete the proof of Theorem 4.

Proof [Proof of Theorem 4] We will consider the event of Lemma 16 on which from Lemma 17 we know with probability at least $1 - \delta$ over the dataset and $(b_j)_{1 \leq j \leq m}$ we have

$$\min_{\mathbf{a}: \|\mathbf{a}\|_2 \lesssim \frac{\Delta_*}{\sqrt{m}}} \hat{R}(\mathbf{W}^T, \mathbf{a}, \mathbf{b}) - \mathbb{E}[\ell_{\mathbf{H}}(\epsilon)] \lesssim \Delta_*^2 \left(\sqrt{\frac{\log(1/\delta)}{m}} + \sqrt{\frac{d + \log(1/\delta)}{T}} \right).$$

Notice that $\mathbf{a} \mapsto \hat{R}(\mathbf{W}, \mathbf{a}, \mathbf{b})$ is a convex function. Thus by strong duality, there exists $\lambda' > 0$ such that the value of the above constrained minimization problem is equal to the value of the following regularized minimization problem,

$$\min_{\mathbf{a}} \hat{\mathcal{R}}_{\lambda'}(\mathbf{W}^T, \mathbf{a}, \mathbf{b}) - \mathbb{E}[\ell_{\mathbf{H}}(\epsilon)] \lesssim \Delta_*^2 \left(\sqrt{\frac{\log(1/\delta)}{m}} + \sqrt{\frac{d + \log(1/\delta)}{T}} \right).$$

Explicitly, this λ' can be chosen such that the unique solution to

$$\nabla_{\mathbf{a}} \hat{R}(\mathbf{W}^T, \mathbf{a}^*, \mathbf{b}) + \lambda' \mathbf{a}^* = 0 \tag{E.7}$$

has $\|\mathbf{a}^*\|_2 \lesssim \frac{\Delta_*}{\sqrt{m}}$. Notice that this \mathbf{a}^* is the unique solution to $\arg \min_{\mathbf{a}} \hat{\mathcal{R}}_{\lambda'}(\mathbf{W}^T, \mathbf{a}, \mathbf{b})$.

Moreover, from Lemma 19 we have

$$\hat{\mathcal{R}}_{\lambda'}(\mathbf{W}^T, \mathbf{a}^{T'}, \mathbf{b}) - \hat{\mathcal{R}}_{\lambda'}(\mathbf{W}^T, \mathbf{a}^*, \mathbf{b}) \lesssim \frac{\hat{R}(\mathbf{W}^T, \mathbf{a}^0, \mathbf{b})}{T'^2} + \frac{(d + \Delta^2)(\Delta^2 + \lambda'/m + \log(1/\delta))}{(\lambda'/m)^2 T'},$$

and by strong convexity

$$\|\mathbf{a}^{T'} - \mathbf{a}^*\|_2^2 \leq \frac{2}{m} \left(\hat{\mathcal{R}}_{\lambda'}(\mathbf{W}^T, \mathbf{a}^{T'}, \mathbf{b}) - \hat{\mathcal{R}}_{\lambda'}(\mathbf{W}^T, \mathbf{a}^*, \mathbf{b}) \right).$$

Thus, with sufficiently large T' such that

$$\frac{\hat{R}(\mathbf{W}^T, \mathbf{a}^0, \mathbf{b})}{T'^2} + \frac{(d + \Delta^2)(\Delta^2 + \lambda'/m + \log(1/\delta))}{(\lambda'/m)^2 T'} \lesssim \Delta_*^2 \sqrt{\frac{d + \log(1/\delta)}{T}} \wedge \frac{\lambda' \Delta_*}{\sqrt{m}}, \tag{E.8}$$

we have $\|\mathbf{a}^{T'}\|_2 \lesssim \frac{\Delta_*}{\sqrt{m}}$ and

$$\hat{\mathcal{R}}_{\lambda'}(\mathbf{a}^{T'}) - \mathbb{E}[\ell_{\mathbf{H}}(\epsilon)] \lesssim \Delta_*^2 \left(\sqrt{\frac{\log(1/\delta)}{m}} + \sqrt{\frac{d + \log(1/\delta)}{T}} \right).$$

Finally, we invoke Theorem 5, to close the generalization gap and get

$$R_{\tau}(\mathbf{W}^T, \mathbf{a}^{T'}, \mathbf{b}) - \mathbb{E}[\ell_{\mathbf{H}}(\epsilon)] \lesssim \Delta_*^2 \sqrt{\frac{\log(1/\delta)}{m}} + \Delta_*^2 \sqrt{\frac{d + \log(1/\delta)}{T}}.$$

■

Appendix F. Example of Non-Convex $\mathcal{R}_\lambda(\mathbf{W})$

Here, we outline an examples for which $\mathcal{R}_\lambda(\mathbf{W})$ is non-convex on a neighborhood around $\mathbf{W} = \mathbf{0}$ while $\lambda = \frac{\tilde{\lambda}}{m}$ satisfies the condition in Theorem 2. For simplicity of exposition, we fix $\mathbf{a} = \frac{\mathbf{1}_m}{m}$ where $\mathbf{1}_m$ is the vector of all ones. It is easy to observe that the results hold with high probability when \mathbf{a} follows the initialization of Assumption 2 as well.

This example is constructed for the smooth activation case of Theorem 2. We require $\sigma(\pm 1) = -\beta_0 = -1$ (which automatically implies $\sigma'(\pm 1) = 0$ as σ attains its minimum) and $\sigma''(\pm 1) = \beta_2 = 1$. Then simplifying $\nabla^2 \mathcal{R}_\lambda(\mathbf{0})$ yields

$$\nabla^2 \mathcal{R}_\lambda(\mathbf{0}) = \frac{\mathbf{I}_m \otimes (-1 - \mathbb{E}[y\mathbf{x}\mathbf{x}^\top])}{m} + \frac{\tilde{\lambda}}{m} \mathbf{I}_{md} = \mathbf{I}_m \otimes \left(\frac{\tilde{\lambda} - 1}{m} \mathbf{I}_d - \frac{\mathbb{E}[y\mathbf{x}\mathbf{x}^\top]}{m} \right).$$

Thus, we need to show that $\frac{\tilde{\lambda} - 1}{m} \mathbf{I}_d - \frac{\mathbb{E}[y\mathbf{x}\mathbf{x}^\top]}{m}$ is not PSD. Let $y = \frac{1}{2}(1 + \tanh(\langle \mathbf{w}, \mathbf{x} \rangle^2 - \|\mathbf{w}\|_2^2))$ and $\tilde{\lambda} = 1 + \mathbb{E}[y] + \gamma$ (notice that $y \geq 0$ thus $\tilde{\lambda}$ indeed satisfies the assumption in Theorem 2) with

$$\gamma = \frac{1}{4} \mathbb{E} \left[(\langle \mathbf{w}, \mathbf{x} \rangle^2 - \|\mathbf{w}\|_2^2) \tanh(\langle \mathbf{w}, \mathbf{x} \rangle^2 - \|\mathbf{w}\|_2^2) \right] > 0.$$

Then we have

$$\begin{aligned} \mathbf{w}^\top \left(\frac{\tilde{\lambda} - 1}{m} \mathbf{I}_d - \frac{\mathbb{E}[y\mathbf{x}\mathbf{x}^\top]}{m} \right) \mathbf{w} &= \frac{\gamma \|\mathbf{w}\|_2^2 + \mathbb{E} \left[y (\|\mathbf{w}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2) \right]}{m} \\ &= \frac{\gamma \|\mathbf{w}\|_2^2 - \frac{1}{2} \mathbb{E} \left[(\langle \mathbf{w}, \mathbf{x} \rangle^2 - \|\mathbf{w}\|_2^2) \tanh(\langle \mathbf{w}, \mathbf{x} \rangle^2 - \|\mathbf{w}\|_2^2) \right]}{m} < 0. \end{aligned}$$

Therefore, we have shown that $\mathcal{R}_\lambda(\mathbf{W})$ is not convex on a neighborhood around zero.

Appendix G. Auxiliary Lemmas

In order to be explicit, we state the following definitions and lemmas that will be used in the proof of Theorem 5. We only state the next definitions and lemmas and refer the reader to [68] and [67] for proof and more details.

Definition 20 [68, Definitions 2.2 and 2.7] *A real-valued random variable z is said to be ν -sub-Gaussian if for all $s \in \mathbb{R}$ we have $\mathbb{E}[\exp(sz)] \leq \exp(s \mathbb{E}[z] + \frac{s^2 \nu^2}{2})$, and is said to be u -sub-exponential if for all $|s| \leq \frac{1}{\nu}$ we have $\mathbb{E}[\exp(sz)] \leq \mathbb{E} \left[\exp(s \mathbb{E}[z] + \frac{s^2 \nu^2}{2}) \right]$.*

Lemma 21 [67, Propositions 2.5.2 and 2.7.1] *Suppose z is a zero-mean random variable and $\mathbb{E}[|z|^p]^{\frac{1}{p}} \leq L\sqrt{p}$ for all $p \geq 1$. Then z is cL -sub-Gaussian for an absolute constant $c > 0$, i.e. $\mathbb{E}[\exp(sz)] \leq \exp(\frac{s^2 c^2 L^2}{2})$ for all $s \in \mathbb{R}$. Similarly, suppose $\mathbb{E}[|z|^p]^{\frac{1}{p}} \leq Lp$. Then z is cL -sub-exponential for an absolute constant $c > 0$.*

Lemma 22 [68, Theorem 2.26] *Let $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$. If $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -Lipschitz, then $f(\mathbf{x})$ is sub-Gaussian with parameter L .*

Lemma 23 [68, Example 6.3] *Let $\{\mathbf{x}^{(i)}\}_{1 \leq i \leq n}$ be a sequence of i.i.d. standard Gaussian random vectors $x_i \sim \mathcal{N}(0, \mathbf{I}_d)$. It holds with probability at least $1 - \delta$ that*

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} - \mathbf{I}_d \right\|_2 \leq C \left(\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{d + \log(1/\delta)}{n} \right),$$

where C is an absolute constant.

The next lemma is the well-known symmetrization argument that upper bounds the expected value of an empirical process with Rademacher complexity.

Lemma 24 [49, Theorem 3.3] *Let \mathcal{F} be a class functions $f : \mathbb{R}^p \rightarrow \mathbb{R}$ for some $p > 0$. For a number of samples T and a probability distribution \mathcal{P} on \mathbb{R}^p , define the Rademacher complexity of \mathcal{F} as*

$$\mathfrak{R}(\mathcal{F}) = \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{T} \sum_{t=0}^{T-1} \xi_t f(\mathbf{x}^{(t)}) \right], \quad (\text{G.1})$$

where $\{\mathbf{x}^{(t)}\}_{t=0}^{T-1} \stackrel{i.i.d.}{\sim} \mathcal{P}$ and $\{\xi_t\}_{t=0}^{T-1}$ are independent Rademacher random variables (i.e. ± 1 equiprobably). Then the following holds,

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{T} \sum_{t=0}^{T-1} f(\mathbf{x}^{(t)}) - \mathbb{E}[f(\mathbf{x})] \right| \right] \leq 2\mathfrak{R}(\mathcal{F}).$$

Furthermore, we have the following fact for standard normal random vectors.

Lemma 25 *Let $x \sim \mathcal{N}(0, \mathbf{I}_d)$. There exists an absolute constant $C > 0$ such that for any $V \in \mathbb{R}^{m \times k}$ and $p \geq 1$ we have*

$$\mathbb{E}[\|\mathbf{V}\mathbf{x}\|_2^p]^{\frac{1}{p}} \leq \|\mathbf{V}\|_F + C\|\mathbf{V}\|_2\sqrt{p}.$$

Proof First of all, $\|\mathbf{V}\mathbf{x}\|_2$ is a $\|\mathbf{V}\|_2$ -Lipschitz function of x , thus Lemma 22 applies and $\|\mathbf{V}\mathbf{x}\|_2$ is sub-Gaussian. Furthermore, by applying Lemma 21 to $\|\mathbf{V}\mathbf{x}\|_2 - \mathbb{E}[\|\mathbf{V}\mathbf{x}\|_2]$ and Minkowski's inequality, we have

$$\begin{aligned} C\|\mathbf{V}\|_2\sqrt{p} &\geq \mathbb{E}[\|\mathbf{V}\mathbf{x}\|_2^p]^{\frac{1}{p}} - \mathbb{E}[\|\mathbf{V}\mathbf{x}\|_2] \\ &\geq \mathbb{E}[f\|\mathbf{V}\mathbf{x}\|_2^p]^{\frac{1}{p}} - \|\mathbf{V}\|_F, \end{aligned}$$

where the last inequality follows from Jensen's inequality. ■

Lemma 26 *Let $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$. Then $\mathbb{E}[\exp(c\|\mathbf{x}\|_2^2)] \leq \exp(2cd)$ for $c \leq 1/4$.*

Proof Gaussian integration yields $\mathbb{E}[\exp(cx_i^2)] = \frac{1}{\sqrt{1-2c}}$. Furthermore, for $c \leq \frac{1}{4}$ we have $\frac{1}{\sqrt{1-2c}} \leq \exp(2c)$. ■