Self-Regularity of Non-Negative Output Weights for Overparameterized Two-Layer Neural Networks

David Gamarnik*

Eren C. Kızıldağ[†]

Ilias Zadik[‡]

March 3, 2021

Abstract

We consider the problem of finding a two-layer neural network with sigmoid, rectified linear unit (ReLU), or binary step activation functions that "fits" a training data set as accurately as possible as quantified by the training error; and study the following question: does a low training error guarantee that the norm of the output layer (outer norm) itself is small? We answer affirmatively this question for the case of non-negative output weights. Using a simple covering number argument, we establish that under quite mild distributional assumptions on the input/label pairs; any such network achieving a small training error on polynomially many data necessarily has a well-controlled outer norm. Notably, our results (a) have a polynomial (in d) sample complexity, (b) are independent of the number of hidden units (which can potentially be very high), (c) are oblivious to the training algorithm; and (d) require quite mild assumptions on the data (in particular the input vector $X \in \mathbb{R}^d$ need not have independent coordinates). We then leverage our bounds to establish generalization guarantees for such networks through fat-shattering dimension, a scale-sensitive measure of the complexity class that the network architectures we investigate belong to. Notably, our generalization bounds also have good sample complexity (polynomials in d with a low degree), and are in fact near-linear for some important cases of interest.

^{*}MIT; e-mail: gamarnik@mit.edu. Research supported by the NSF grants DMS-2015517.

[†]MIT; e-mail: kizildag@mit.edu.

[‡]NYU; e-mail: zadik@nyu.edu. Research supported by a CDS Moore-Sloan Postdoctoral Fellowship.

Contents

| 1 | Introduction | 3 |
|---|---|-----------|
| 2 | Outer Norm Bounds2.1 Self-Regularity for the Sigmoid Networks2.2 Self-Regularity for the ReLU Networks2.3 Self-Regularity for the Step Networks | 8 |
| 3 | Generalization Guarantees via Outer Norm Bounds 3.1 The Learning Setting | 11 |
| 4 | Conclusion and Future Directions | 15 |
| 5 | Proofs | 17 |
| | 5.1 Proof of Theorem 2.1 | 17 |
| | 5.2 Proof of Theorem 2.3 | 18 |
| | 5.3 Proof of Theorem 2.4 | 18 |
| | 5.4 Proof of Theorem 3.1 | 20 |

1 Introduction

Neural network (NN) architectures achieved a great deal of success in practice. An ever-growing list of their applications includes image recognition [HZRS16], image classification [KSH12], speech recognition [MDH11], natural language processing [CW08], game playing [SSS+17] and more. Despite this great empirical success, however, a rigorous understanding of these networks is still an ongoing quest.

A common paradigm in classical statistics is that overparameterized models, that is, models with more parameters than necessary, pick on the idiosyncrasies of the training data itself—dubbed as overfitting; and as a consequence, tend to predict poorly on the unseen data—called poor generalization. The aforementioned success of the NN architectures, however, stands in the face of this conventional wisdom; and a growing body of recent literature, starting from [ZBH+16], has demonstrated exactly the opposite effect for a broad class of NN models: even though the number of parameters, such as the number of hidden units (neurons), of a NN significantly exceeds the sample size, and a perfect (zero) in-training error is achieved (commonly called as data interpolation); they still retain a good generalization ability. Some partial and certainly very incomplete list of references to this point are found in [DZPS18, LL18, GLSS18, GAS+19, BHMM19, ADH+19a]. Defying statistical intuition even further, it was established empirically in [BHMM19] that beyond a certain point, increasing the number of parameters increases out of sample accuracy.

Explaining this conundrum is arguably one of the most vexing current problems in the field of theoretical machine learning. Standard Vapnik-Chervonenkis (VC) theory do not help explaining the good generalization ability of overparameterized NN models, since the VC-dimension of these networks grows (at least) linearly in the number of parameters [HLM17, BHLM19]. These findings fueled significant research efforts aiming at understanding the generalization ability of such networks. One such line of research is the algorithm-independent front; and is through the lens of controlling the norm of the matrices carrying weights [NTS15, BFT17, LPRS17, GRS17, DR17], PAC-Bayes theory [NBS17, NBMS17], and compression-based bounds [AGNZ18], among others. A major drawback of these approaches, however, is that they require certain norm constraints on the weights considered; therefore making their guarantees a posteriori in nature: whether or not the weights of the NN are bounded (hence a good generalization holds) can be determined only after the training process is complete. An alternative line of research (detailed below) focuses on the end results of the algorithms, and potentially yields a priori guarantees: for instance, relatively recently, Arora et al. gave in [ADH+19b] a priori guarantees for the solution found by the qradient descent algorithm under random initialization.

A predominant explanation of the aforementioned phenomenon (that the overparameterization does not hurt the generalization ability of the NN architectures) which has emerged recently is based on the idea of self-regularization. Specifically, it is argued that even though there is an abundance of parameter choices perfectly fitting (interpolating) the data (and thus achieving zero in-training error); the algorithms used in training the models, such as the gradient descent and its many variants such as stochastic gradient descent, mirror descent, etc., tend to find solutions which are regularized according to some additional criteria, such as small norms, thus introducing algorithm dependent inductive bias. Namely, the algorithms implemented for minimizing training error "prefer" certain kinds of solutions. The use of these solutions for model building in particular is believed to result in low generalization errors. Thus a significant research effort

(as was partially mentioned above) was devoted to the analysis of the end results of the implementation of such algorithms. This line of research include the analysis of the end results of the gradient descent [BG17, FCG19], stochastic gradient descent [HRS16, BGMSS17, LL18, CG19], as well as the stochastic gradient Langevin dynamics [MWZZ18].

In this paper, we consider two-layer NN models (1)—also known as shallow architectures—consisting of an arbitrary number $\overline{m} \in \mathbb{N}$ of hidden units and sigmoid, rectified linear unit (ReLU), or binary step activations—activations that are arguably among the most popular practical choices—and investigate the following question: to what extent a low training error itself places a restriction on the weights of the learned NN? We take an algorithm-independent route; and establish the following "picture", under the assumption that the output weights $a=(a_i:1\leq i\leq \overline{m})\in\mathbb{R}^{\overline{m}}$ of the "learned" NN are non-negative. When the number N of training samples is at least an explicit (low-degree) polynomial function in d, $N=d^{O(1)}$, the norm $||a||_1$ of the output weights $a\in\mathbb{R}^{\overline{m}}_{\geq 0}$ of any NN model achieving a small training error is well-controlled: $||a||_1=O(1)$, with high probability over the training data set. In particular, for the ReLU and step networks, we obtain a near-linear sample complexity bound, $N=\Theta(d\log d)$ for such a result to hold. Note that a condition such as the non-negativity of a_i is necessary in a strict sense for such a bound on $||a||_1$. Indeed, notice that by growing the width \overline{m} arbitrarily and appropriately choosing alternating signs for the new weights a_i ; one can introduce cancellations and make $||a||_1$ to explode; while keeping the training error unchanged.

Our results are established using elementary tools, in particular through an ϵ -net argument (Definition 1.2). Notably, our results (a) are independent of the number \overline{m} of the hidden units (which can potentially be quite large), (b) are oblivious to the way the training is done (that is, independent of the choice of the training algorithm); and (c) are valid under quite mild distributional assumptions on the input/label pairs $(X,Y) \in \mathbb{R}^d \times \mathbb{R}$. In particular, the coordinates of X need not be independent.

Moreover, a bounded outer norm for such network models implies a well-controlled fat-shattering dimension (FSD) [Bar98]—a measure of the complexity of the model class achieving a low training error. In Section 3, we leverage our outer norm bounds and the FSD to establish generalization guarantees for the networks that we investigate. The current paper presents significantly strengthened versions and extensions of some results appeared in our preprint [EGKZ20].

Preliminaries

We commence this section with a list of notational convention that we follow throughout.

Notation. The set of reals, non-negative reals, and positive integers are denoted respectively by \mathbb{R} , $\mathbb{R}_{\geq 0}$, and \mathbb{N} . For any set S, |S| denotes its cardinality. For any $N \in \mathbb{N}$, $[N] \triangleq \{1, 2, ..., N\}$. For any $v \in \mathbb{R}^n$, its ℓ_p norm is denoted by $||v||_p$. For $u, v \in \mathbb{R}^n$, their Euclidean inner product is denoted by u^Tv . For any $r \in \mathbb{R}$, $\exp(r)$ denotes e^r ; and $\ln(r)$ denotes the logarithm of r base e. For any "event" E; $\mathbb{1}\{E\} = 1$ when E is true; and $\mathbb{1}\{E\} = 0$ when E is false. SGM(x) denotes the sigmoid activation function, $1/(1+\exp(-x))$; ReLU(x) denotes the ReLU activation function, $\max\{x,0\}$; and Step(x) denotes the (binary) step activation, $\mathbb{1}\{x\geq 0\}$. $X \stackrel{d}{=} \mathcal{N}(0,\Sigma)$ if X is a zero-mean multivariate normal vector with covariance Σ . A random variable U is symmetric around zero if U and U have the same distribution, that is $U \stackrel{d}{=} U$. For any random variable

U, (if finite) its moment generating function (MGF) at $s \in \mathbb{R}$, $\mathbb{E}[\exp(sU)]$, is denoted by $M_U(s)$. Finally, $\Theta(\cdot)$, $O(\cdot)$, $O(\cdot)$ are the standard asymptotic order notations.

Setup. A two-layer NN $(a, W) \in \mathbb{R}^{\overline{m}} \times \mathbb{R}^{\overline{m} \times d}$ with \overline{m} hidden units (neurons) computes, for each $X \in \mathbb{R}^d$,

$$\sum_{1 < j < \overline{m}} a_j \sigma\left(w_j^T X\right). \tag{1}$$

Here, $\sigma(\cdot)$ is the activation; $w_j \in \mathbb{R}^d$, the j^{th} row of W, carries the weights of neuron j; and $a = (a_j : 1 \le j \le \overline{m}) \in \mathbb{R}^{\overline{m}}$ carries the output weights. $||a||_1$ is referred to as the *outer norm*. We assume $a_j \ge 0$ for $j \in [\overline{m}]$. This non-negativity assumption appears often in the theoretical study of this model: see [GLM17, DKKZ20, LMZ20] for generic $a \in \mathbb{R}^{\overline{m}}_{\ge 0}$; and [DL18, SS18, ZYWG19, GKM18] for the case a_j are equal to the same positive number.

Our study of NN models under the non-negativity assumption is also partly motivated from an applied point of view, in that, non-negativity is inherent to many data sets appearing in practice, including audio data and data on muscular activity [SV17, Wik] and allow interpretability. Furthermore, non-negativity is also a commonly used assumption in the context of matrix factorization, termed as the non-negative matrix factorization problem (NMF): given a matrix $M \in \mathbb{R}^{n \times m}$ with non-negative entries and an integer $r \geq 1$, the goal of the NMF is to find matrices $A \in \mathbb{R}^{n \times r}$ and $W \in \mathbb{R}^{r \times m}$ with non-negative entries such that the product AW is as "close" to Mas possible; as quantified, e.g., by the Frobenius norm. This problem is a fundamental problem appearing in many practical applications, including information retrieval, document clustering, image segmentation, demography and chemometrics, see [AGKM16] and the references therein. Moreover, NMF is also related to the neural network models that we consider herein with a non-negative activation $\sigma(\cdot)$: observe that in the context of NN models we consider, given data $(X_i,Y_i), 1 \leq i \leq N$, the goal of the learner is to find a $(a,W) \in \mathbb{R}^{\overline{m}} \times \mathbb{R}^{\overline{m} \times d}$ such that Y_i and $a^T \sigma(WX_i)$ are as close as possible, as quantified by the ℓ_2 norm (here, σ acts coordinate-wise to the vector WX). See also [Gab19, Section 6] for a more rigorous connection between shallow NN models, matrix factorization and message passing algorithms. In addition to its key role in the NMF problem; the non-negativity was also argued as a natural assumption for representing objects in the seminal papers by Lee and Seung [LS99, LS01]; and also has roots in biology, in particular in the context of neuronal firing rates, see [Hoy02], and the references therein.

In the sequel, $d \in \mathbb{N}$ is reserved for the input dimension; and $\overline{m} \in \mathbb{N}$ is reserved for the number of neurons. We consider herein two-layer NN models with sigmoid, SGM(x); rectified linear unit, ReLU(x); and binary step, Step(x), activation functions. We refer to these as sigmoid, ReLU; and step networks, respectively. The sigmoid and the ReLU are arguably among the most popular practical choices. The step function, on the other hand, is one of the initial activations considered in the NN literature, and is inspired from a biological point of view: it resembles the firing pattern of a neuron, an initial motivation for studying NN architectures.

Given the data $(X_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$, $1 \leq i \leq N$, consider the problem of finding a two-layer NN $(a, W) \in \mathbb{R}^{\overline{m}} \times \mathbb{R}^{\overline{m} \times d}$ which "fits" the data as accurately as possible. This is achieved by solving the so-called *empirical risk minimization* problem, where the accuracy is quantified by the *training error*

$$\widehat{\mathcal{L}}(a, W) \triangleq \frac{1}{N} \sum_{1 \le i \le N} \left(Y_i - \sum_{1 \le j \le \overline{m}} a_j \sigma\left(w_j^T X_i \right) \right)^2.$$
 (2)

One then runs a training algorithm, e.g., the gradient descent algorithm or one of its variants (such as stochastic gradient descent or mirror descent), to find an (a, W) with a small $\widehat{\mathcal{L}}(a, W)$.

Distributional assumption. We study the case where the input/label pairs (X_i, Y_i) , $1 \le i \le N$, are i.i.d. samples of a distribution on $\mathbb{R}^d \times \mathbb{R}$ (which is potentially unknown to the learner). For our outer norm bounds, we assume that their distribution satisfies the following.

- We assume the input $X \in \mathbb{R}^d$ satisfies $\mathbb{P}(||X||_2^2 \le Cd) \ge 1 \exp(-\Theta(d))$ for some constant C > 0.
- We assume the label Y is such that $\mathbb{E}[|Y|] \triangleq M < \infty$.

Later in Section 3 when we study generalization guarantees, we consider a stronger assumption on labels: we assume the labels Y are bounded, that is, for some M > 0, $|Y| \le M$ almost surely.

These assumptions are quite mild. For instance, $X \in \mathbb{R}^d$ need not have i.i.d. coordinates. Moreover, most real data sets indeed have bounded labels [DZPS18]; and this bounded label assumption is employed extensively in literature, see e.g. [GWZ19, ADH+19b, DLL+19, GK19, LZA20]. Our next assumption regards the number N of samples.

Assumption 1.1. Throughout, we assume that the sample size N satisfies $N \leq \exp(cd)$ for some c > 0.

Assumption 1.1 is required for technical reasons: observe that since $\mathbb{P}(||X_i||_2^2 > Cd) \le \exp(-\Theta(d))$, it holds, by a union bound, that

$$\mathbb{P}\left(||X_i||_2^2 \le Cd, 1 \le i \le N\right) \ge 1 - N \exp(-\Theta(d)).$$

For this bound to be non-vacuous, N should at most be $\exp(cd)$ for a small enough c > 0. This assumption, again, is very benign due to obvious practical reasons. Moreover, it suffices to have $N \ge \operatorname{poly}(d)$ for our results to hold.

Nets and Covering Numbers. The crux of our proofs is the so-called ϵ -net argument [Ver10, Ver18]. This (rather elementary) argument is also known as the covering number argument; and has been employed extensively in the literature; including compressed sensing, machine learning and probability theory.

Definition 1.2. Let $\epsilon > 0$. Given a metric space (X, ρ) , a subset $\mathcal{N}_{\epsilon} \subset X$ is called an ϵ -net of X if, for every $x \in X$, there is a $y \in \mathcal{N}_{\epsilon}$ such that $\rho(x, y) \leq \epsilon$. The smallest cardinality of such an \mathcal{N}_{ϵ} , if finite, is called the covering number of X, denoted by $\mathcal{N}(X, \epsilon)$.

The next result, verbatim from [Ver18, Corollary 4.2.13], is an upper bound on the covering number of the Euclidean ball.

Theorem 1.3. Let $B_2(0,R) \triangleq \{x \in \mathbb{R}^d : ||x||_2 \leq R\}$. Then for $R \geq 1$ and any $\epsilon > 0$

$$\mathcal{N}(B_2(0,R),\epsilon) \leq (3R/\epsilon)^d$$
.

Paper organization. The rest of the paper is organized as follows. Our main results on the self-regularity of output weights are presented in Section 2. In particular, see Sections 2.1, 2.2, and 2.3 for the cases of sigmoid, ReLU, and step networks, respectively. By leveraging our outer norm bounds and employing earlier results on the fat shattering dimension, we establish in Section 3 generalization guarantees. We outline several future directions in Section 4. Finally, we present our proofs in Section 5.

2 Outer Norm Bounds

In this section, we establish the self-regularity of the output weights for the aforementioned networks. That is, we establish that the outer norms of sigmoid, ReLU, and step networks with non-negative output weights achieving a small training error (2) on polynomially many data is O(1).

2.1 Self-Regularity for the Sigmoid Networks

Our first focus in on the sigmoid networks. This object, for each $X \in \mathbb{R}^d$, computes the function (1) with $\sigma = \text{SGM}(\cdot) = (1 + \exp(-x))^{-1}$. Our first main result establishes an outer norm bound for this architecture.

Theorem 2.1. Let $\delta, M, R > 0$; and $(X_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$, $i \in [N]$ be i.i.d. data with $\mathbb{E}[|Y_i|] = M < \infty$; where N satisfies Assumption 1.1. For any $\overline{m} \in \mathbb{N}$, define

$$\mathcal{S}\left(\overline{m}, \delta, R\right) = \left\{(a, W) \in \mathbb{R}^{\overline{m}}_{\geq 0} \times \mathbb{R}^{\overline{m} \times d} : \max_{1 \leq j \leq \overline{m}} ||w_j||_2 \leq R, \ \widehat{\mathcal{L}}\left(a, W\right) \leq \delta^2 \right\},$$

where $\widehat{\mathcal{L}}(\cdot)$ is defined in (2) with $\sigma(\cdot) = SGM(\cdot)$. Suppose, in addition, that the random variable $w^TX \in \mathbb{R}$ is symmetric around zero for every $w \in \mathbb{R}^d$. Then,

$$\mathbb{P}\left(\sup_{(a,W)\in\mathcal{S}(\delta,R)}||a||_1 \le 3(1+e)(\delta+2M)\right) \ge 1 - \left(3R\sqrt{Cd}\right)^d \exp\left(-\Theta(N)\right) - N\exp\left(-\Theta(d)\right) - o_N(1),\tag{3}$$

where $S(\delta, R) \triangleq \bigcup_{\overline{m} \in \mathbb{N}} S(\overline{m}, \delta, R)$.

Corollary 2.2. Let $R = \exp(d^{O(1)})$. Then, under the assumptions of Theorem 2.1; it holds w.h.p. that $\sup_{(a,W)\in\mathcal{S}(\delta,R)} ||a||_1 \leq 3(1+e)(\delta+2M)$, provided $N \geq d^{O(1)}$.

The proof of Theorem 2.1 is provided in Section 5.1.

Above, $o_N(1)$ is a function which depends only on the distribution of Y and N; and tends to zero as $N \to \infty$. Several remarks are now in order. Theorem 2.1 states that any two-layer sigmoid NN which (a) consists of internal weights w_j bounded in norm by an exponentially large (in d) quantity and non-negative output weights; and (b) achieves a small training error on a sufficiently large data set, has a well-controlled outer norm. It is worth noting that Theorem 2.1 is oblivious to how the training is done: this result not only applies to the weights obtained, say, via the gradient descent algorithm; but applies to any weights (subject to the aforementioned assumptions) achieving a small training loss.

Moreover, the upper bound established in Theorem 2.1 is also oblivious to the number \overline{m} of the neurons of the NN used for fitting. In particular, adopting a teacher/student setting as in [GAS+19] where the input/label pairs (X_i, Y_i) are generated by a teacher NN; the output norm of any student NN—which may potentially be significantly overparameterized with respect to the teacher NN—is still well-controlled, provided the assumptions of Theorem 2.1 are satisfied. The extra requirement that $w^T X$ is symmetric is quite mild: it holds for many data distributions, e.g., for $X \stackrel{d}{=} \mathcal{N}(0, \Sigma)$ where Σ is an arbitrary positive semidefinite matrix.

The $o_N(1)$ term is due to a certain high probability event \mathcal{E}_0 , see (10) in the proof. The probability of this event is controlled through the weak law of large numbers; and the $o_N(1)$ term can be improved explicitly (a) to O(1/N) if $\mathbb{E}[Y^2] < \infty$; and (b) to $\exp(-\Theta(N))$ if Y_i satisfy the large deviations bounds (which holds, for instance, when the moment generating function of Y_i exists in a neighbourhood around zero). Moreover, if Y is (almost surely) bounded (which holds for real data sets, as noted earlier), then it can be dropped altogether.

Furthermore, Corollary 2.2—which follows immediately from Theorem 2.1—asserts that even under the mild assumption $R = \exp(d^{O(1)})$ (i.e., the weights w_j are unbounded from a practical perspective), $\sum_j a_j$ is still O(1), provided that that the number N of data is polynomial in d.

Moreover, an inspection of the proof of Theorem 2.1 reveals the following. The constant 3(1+e) can be improved to any constant greater than four with slightly more work. Moreover, the thesis of Theorem 2.1 still remains valid (with appropriately modified constants) for any non-negative activation which is continuous at the origin and whose value at the origin is positive. This includes the softplus activation $\ln(1+e^x)$ [GBB11], the Gaussian activation, $\exp(-x^2)$; among others.

2.2 Self-Regularity for the ReLU Networks

Our next focus is on the ReLU networks. This object, for each input $X \in \mathbb{R}^d$, computes the function (1) with $\sigma(x) = \text{ReLU}(x) = \max\{x, 0\} = \frac{1}{2}(x + |x|)$.

We first observe that the ReLU function is positive homogeneous: for any $c \geq 0$ and $x \in \mathbb{R}$, $\text{ReLU}(cx) = c \cdot \text{ReLU}(x)$. For this reason, we may assume, without loss of generality, that $||w_j||_2 = 1$ for $1 \leq j \leq \overline{m}$. Indeed, if $w_j \neq 0$, one can simply "push" its norm outside; whereas if $w_j = 0$, then one can replace it with any unit norm vector and set $a_j = 0$ instead.

It is worth noting that since the ReLU case requires no explicit assumptions on $||w_j||_2$, an outer bound for this case is a somewhat stronger conclusion than an outer bound for the case of sigmoid activation.

Equipped with this, we now present our next result.

Theorem 2.3. Let $\delta, M > 0$; and $(X_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$, $i \in [N]$ be i.i.d. data with $\mathbb{E}[|Y_i|] = M < \infty$; where N satisfies Assumption 1.1. For any $\overline{m} \in \mathbb{N}$, define

$$\mathcal{G}\left(\overline{m},\delta\right) = \left\{(a,W) \in \mathbb{R}^{\overline{m}}_{\geq 0} \times \mathbb{R}^{\overline{m} \times d} : ||w_j||_2 = 1, 1 \leq j \leq \overline{m}; \ \widehat{\mathcal{L}}\left(a,W\right) \leq \delta^2\right\},\,$$

where $\widehat{\mathcal{L}}(\cdot)$ is defined in (2) with $\sigma(\cdot) = \text{ReLU}(\cdot)$. Suppose, in addition, that for $Y_w \triangleq w^T X$, (a) there exists a $\mu^* > 0$ such that $\mathbb{E}[\text{ReLU}(Y_w)] \geq \mu^*$ for any $w \in B_2(0,1)$; and (b) for some s > 0, $M_1(s)$ and $M_2(s)$ are independent of d and are finite; where $M_1(s) \triangleq \sup_{w:||w||_2=1} M_{Y_w}(s)$ and

 $M_2(s) \triangleq \sup_{w:||w||_2=1} M_{Y_w}(-s)$. Then,

$$\mathbb{P}\left(\sup_{(a,W)\in\mathcal{G}(\delta)}||a||_{1} \leq 4(\delta+2M)(\boldsymbol{\mu}^{*})^{-1}\right) \geq 1 - \left(\frac{12\sqrt{Cd}}{\boldsymbol{\mu}^{*}}\right)^{d} \exp\left(-\Theta(N)\right) - N\exp\left(-\Theta(d)\right) - o_{N}(1),$$

$$where \ \mathcal{G}(\delta) \triangleq \bigcup_{\overline{m}\in\mathbb{N}} \mathcal{G}\left(\overline{m},\delta\right).$$
(4)

The proof of Theorem 2.3 is provided in Section 5.2.

In particular, it suffices to have a near-linear number of samples, $N = \Theta(d \log d)$, to obtain a good, uniform, control over $||a||_1$. As mentioned above, we managed to bypass the dependence on the term R that appears in Theorem 2.1 by leveraging the fact that ReLU is a positive homogeneous function.

Analogous to Theorem 2.1, the bound established in Theorem 2.3 is also oblivious to (a) how the training is done, and (b) the number \overline{m} of neurons. In particular, even potentially overparameterized networks have a well-controlled outer norm; provided that they achieve a small training error on a sufficient number N of data. The additional distributional requirements are still mild. For instance, when $X \stackrel{d}{=} \mathcal{N}(0, I_d)$, $w^T X \stackrel{d}{=} \mathcal{N}(0, 1)$ for any w with $||w||_2 = 1$; and μ^* can be taken to be $1/\sqrt{2\pi}$. The requirement (b) ensures the existence of the moment generating function in a neighborhood around zero, hence the large deviations bounds are applicable. The same remarks on $o_N(1)$ term following Theorem 2.1 also apply here: it can be improved to O(1/N) or $\exp(-\Theta(N))$ under slightly stronger assumptions on Y_i .

2.3 Self-Regularity for the Step Networks

Our final focus is on the step networks. This object, for each $X \in \mathbb{R}^d$, computes (1) with $\sigma(x) = \text{Step}(x) = \mathbb{1}\{x \geq 0\}$.

Like the ReLU case, Step(x) is also homogeneous: for every $c \ge 0$, Step(cx) = Step(x). For this reason, we assume, without loss of generality, $||w_j||_2 = 1$, $1 \le j \le \overline{m}$.

Theorem 2.4. Let $\delta, M > 0$; and $(X_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$, $i \in [N]$ be i.i.d. data with $\mathbb{E}[|Y_i|] = M < \infty$; where N satisfies Assumption 1.1. For any $\overline{m} \in \mathbb{N}$, define

$$\mathcal{H}\left(\overline{m},\delta\right) = \left\{ (a,W) \in \mathbb{R}^{\overline{m}}_{\geq 0} \times \mathbb{R}^{\overline{m} \times d} : ||w_j||_2 = 1, 1 \leq j \leq \overline{m}; \ \widehat{\mathcal{L}}\left(a,W\right) \leq \delta^2 \right\},\,$$

with $\widehat{\mathcal{L}}(\cdot)$ as in (2) with $\sigma(\cdot) = Step(\cdot)$. Moreover, assume that for some $\eta > 0$, $\inf_{w:||w||_2=1} \mathbb{P}\left(w^TX \geq \eta\right) \geq \eta$. Then,

$$\mathbb{P}\left(\sup_{(a,W)\in\mathcal{H}(\delta)}||a||_1 \le 2(\delta+2M)\eta^{-1}\right) \ge 1 - \left(\frac{6\sqrt{Cd}}{\eta}\right)^d \exp\left(-\Theta(N)\right) - N\exp\left(-\Theta(d)\right) - o_N(1)$$

where $\mathcal{H}(\delta) \triangleq \bigcup_{\overline{m} \in \mathbb{N}} \mathcal{H}(\overline{m}, \delta)$.

The proof of Theorem 2.4 is provided in Section 5.3.

Main remarks following Theorems 2.1 and 2.3—in particular, independence from \overline{m} as well as the training algorithm—apply here, as well.

The extra condition on the distribution ensures that the collection $\{\mathbb{P}(w^TX \geq \eta) : ||w||_2 = 1\}$ is uniformly bounded away from zero. This is again quite mild, as demonstrated by the following example. Suppose $Y_w \triangleq w^TX$ is centered and equidistributed for w with $||w||_2 = 1$. (Observe that this is indeed the case, e.g. when $X \stackrel{d}{=} \mathcal{N}(0, I_d)$.) Then as long as $\operatorname{Var}(Y_w) > 0$ the extra requirement per Theorem 2.4 is satisfied. Indeed, for this case $\mathbb{P}(Y_w > 0) > 0$. Hence, using the continuity of probabilities

$$\mathbb{P}(Y_w > 0) = \mathbb{P}(w^T X > 0) = \lim_{t \to \infty} \mathbb{P}(w^T X > t^{-1}) > 0,$$

one ensures the existence of such an η . In the case where $X \stackrel{d}{=} \mathcal{N}(0, I_d)$, one can concretely take $\eta = 0.3$.

3 Generalization Guarantees via Outer Norm Bounds

3.1 The Learning Setting

In this section, we leverage the outer norm bounds we established in Theorems 2.1-2.4 to provide generalization guarantees for the neural network architectures having non-negative output weights that we investigated.

Our approach is through a quantity called the *fat-shattering dimension* (FSD) of such networks introduced by Kearns and Schapire [KS94]. This quantity is essentially a scale-sensitive measure of the complexity of the "class" (appropriately defined) that the network architecture being considered belongs to. We introduce the FSD formally in Definition 5.1 found in Section 5.4. For more information on the FSD, we refer the interested reader to the original paper by Kearns and Schapire [KS94]; as well as earlier papers by Bartlett, Long, and Williamson [BLW96], and Bartlett [Bar98].

In what follows, we prove our promised generalization guarantee (Theorem 3.1 below) by combining the prior results on the FSD of such networks with our outer norm bounds. Bartlett provides in [Bar98] upper bounds on the FSD of certain function classes H. He then leverages these bounds to give good generalization guarantees. One of the classes he studies is precisely the class of two-layer NN with a **bounded outer norm** (as we do). In particular, he establishes in [Bar98, Corollary 24] (which is restated as Theorem 5.2 below) that the class of two-layer networks with **bounded outer norm** has a well-controlled FSD: informally, it has "low complexity". He then leverages the FSD bounds to devise good generalization guarantees for the architectures that he investigates. It is worth noting, however, that he establishes this link in the context of classification setting, $Y \in \{\pm 1\}$. Since we assume $a_j \geq 0$, and the activations we study are non-negative, this does not apply to our case: the outputs of the networks we study are always non-negative. Nevertheless, we by-pass this by combining our outer norm bounds (Theorems 2.1-2.4), Theorem 5.2, as well as building upon several other prior results tailored for the regression setting.

We next recall the learning setting for convenience. Let \mathcal{D} be a distribution on $\mathbb{R}^d \times \mathbb{R}$ for the input/label pairs (X,Y); and let $(X_i,Y_i) \sim \mathcal{D}$, $1 \leq i \leq N$, be the i.i.d. training data. The goal of the learner is to find a NN $(a,W) \in \mathbb{R}^{\overline{m}} \times \mathbb{R}^{\overline{m} \times d}$ with \overline{m} hidden units (neurons) and activation $\sigma(\cdot)$ which "explains" the data (X_i,Y_i) , $1 \leq i \leq N$, as accurately as possible, often by solving the *empirical risk minimization* problem, $\min_{a,W} \widehat{\mathcal{L}}(a,W)$ (2). The "learned"

network is then used for predicting the unseen data. The generalization ability of the "learned" network $(a, W) \in \mathbb{R}^{\overline{m}} \times \mathbb{R}^{\overline{m} \times d}$ is quantified by the so-called *generalization error* (also known as the *population risk*)

$$\mathcal{L}(a, W) \triangleq \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\left(Y - \sum_{1 < j < \overline{m}} a_j \sigma \left(w_j^T X \right) \right)^2 \right]. \tag{5}$$

Here, the expectation is taken w.r.t. to a fresh sample $(X,Y) \sim \mathcal{D}$, which is independent of the training data. The "gap"

$$\left|\widehat{\mathcal{L}}\left(a,W\right) - \mathcal{L}(a,W)\right| = \left|\frac{1}{N} \sum_{1 \leq i \leq N} \left(Y_i - \sum_{1 \leq j \leq \overline{m}} a_j \sigma\left(w_j^T X_i\right)\right)^2 - \mathbb{E}_{(X,Y) \sim \mathcal{D}}\left[\left(Y - \sum_{1 \leq j \leq \overline{m}} a_j \sigma\left(w_j^T X\right)\right)^2\right]\right|$$

between the training error and the generalization error is called the generalization gap.

In what follows, we focus our attention on the generalization ability (5) of the learned networks (a, W) that achieved a small training error, $\widehat{\mathcal{L}}(a, W) \leq \delta^2$ (2), on a polynomial (in d) number of data. The details of the training process (such as the algorithm used for training) are immaterial to us; and our results apply to **any** NN (a, W) provided it achieved a small training error, $\widehat{\mathcal{L}}(a, W) \leq \delta^2$.

In this section, we also assume that the labels Y are bounded: \mathcal{D} is such that for some M > 0, $|Y| \leq M$ almost surely. This is necessary, as the prior results we employ from Haussler [Hau92] and Bartlett, Long, and Williamson [BLW96] (in particular, see Theorem 5.4) apply only to the case where the labels are bounded. For this reason, the $o_N(1)$ terms present in Theorems 2.1-2.3 disappear, see the remarks following each theorem.

3.2 The Generalization Guarantees

Equipped with our outer norm bounds (Theorems 2.1-2.4) and Theorem 5.2, we now establish the promised generalization guarantees for the aforementioned networks whose output weights a_i are non-negative. To that end, let $\alpha, M, M, A > 0$ be certain parameters (elaborated below); and set

$$\xi\left(\alpha, M, \mathcal{M}, A\right) \triangleq \frac{2}{\ln 2} \cdot \frac{c \cdot 128^2 \cdot \mathcal{M}^6 A^6 \cdot \max\{\mathcal{M}A, 2M\}^2}{\alpha^2} \cdot \ln\left(\frac{128\mathcal{M}^3 A^3 \max\{\mathcal{M}A, 2M\}}{\alpha}\right), (6)$$

where c > 0 is the absolute constant appearing in Theorem 5.2. Our result is as follows.

Theorem 3.1. Let $\alpha, \delta, M, R > 0$, and $(X_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$, $1 \leq i \leq N$, be i.i.d. samples drawn from an arbitrary distribution \mathcal{D} on $\mathbb{R}^d \times \mathbb{R}$ with $|Y| \leq M$ almost surely; where N satisfies Assumption 1.1. For the ξ term defined in (6), set

$$\zeta(\alpha, M, A, N) \triangleq \exp\left(\xi(\alpha, M, 2, A) \cdot d \cdot \ln^2\left(\frac{2304 \cdot N \cdot A^2 \cdot \max\{2A, M\}}{\alpha}\right) - \frac{\alpha^2 \cdot N}{64 \cdot \max\{2A, M\}^2}\right).$$

(a) (Sigmoid Networks) Under the assumptions of Theorem 2.1, with probability at least

$$1 - \zeta(\alpha, M, 3(1+e)(\delta+2M), N) - \left(3R\sqrt{Cd}\right)^d \exp\left(-\Theta(N)\right) - N\exp\left(-\Theta(d)\right)$$

over $(X_i, Y_i) \sim \mathcal{D}$, $1 \leq i \leq N$, it holds that

$$\sup_{(a,W)\in\mathcal{S}(\delta,R)}\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\left(Y-\sum_{1\leq j\leq\overline{m}}a_{j}\mathit{SGM}\left(w_{j}^{T}X\right)\right)^{2}\right]\leq\alpha+\delta^{2},$$

provided

$$N \geq c \cdot 2^{21} \cdot \frac{A^6 \cdot \max\{A, M\}^2}{\alpha^2} \cdot d \quad and \quad \alpha \leq 2^{11} \cdot A^3 \cdot \max\{A, M\}$$

with $A = 3(1 + e)(\delta + 2M)$. Here, $S(\delta, R)$ is the set introduced in Theorem 2.1.

(b) (ReLU Networks) Under the assumptions of Theorem 2.3 and assuming additionally $\mu^* = \exp(o(d))$, with probability at least

$$1 - \zeta\left(\alpha, M, \frac{4\sqrt{Cd}(\delta + 2M)}{\boldsymbol{\mu^*}}, N\right) - \left(\frac{12\sqrt{Cd}}{\boldsymbol{\mu^*}}\right)^d \exp\left(-\Theta(N)\right) - N\exp\left(-\Theta(d)\right),$$

over $(X_i, Y_i) \sim \mathcal{D}$, $1 \leq i \leq N$, it holds that

$$\sup_{(a,W)\in\mathcal{G}(\delta)}\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\left(Y-\sum_{1\leq j\leq \overline{m}}a_{j}\textit{ReLU}\left(w_{j}^{T}X\right)\right)^{2}\right]\leq \alpha+\delta^{2}+e^{-\Theta(d)},$$

provided

$$N \ge c \cdot 2^{21} \cdot \frac{A^6 \cdot \max\{A, M\}^2}{\alpha^2} \cdot d \quad and \quad \alpha \le 2^{11} \cdot A^3 \cdot \max\{A, M\}$$

with $A = \frac{4\sqrt{Cd}(\delta + 2M)}{\mu^*}$. Here, $\mathcal{G}(\delta)$ is the set introduced in Theorem 2.3.

(c) (Step Networks) Under the assumptions of Theorem 2.4, with probability at least

$$1 - \zeta\left(\alpha, M, \frac{2(\delta + 2M)}{\eta}, N\right) - \left(\frac{6\sqrt{Cd}}{\eta}\right)^d \exp\left(-\Theta(N)\right) - N\exp\left(-\Theta(d)\right)$$

over $(X_i, Y_i) \sim \mathcal{D}$, $1 \leq i \leq N$, it holds that

$$\sup_{(a,W)\in\mathcal{H}(\delta)}\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\left(Y-\sum_{1\leq j\leq \overline{m}}a_{j}\textit{Step}\left(w_{j}^{T}X\right)\right)^{2}\right]\leq \alpha+\delta^{2},$$

provided

$$N \ge c \cdot 2^{21} \cdot \frac{A^6 \cdot \max\{A, M\}^2}{\alpha^2} \cdot d \quad and \quad \alpha \le 2^{11} \cdot A^3 \cdot \max\{A, M\}$$

with $A = \frac{2(\delta + 2M)}{\eta}$. Here, $\mathcal{H}(\delta)$ is the set introduced in Theorem 2.4.

Theorem 3.1 is established by combining various individual results established in separate works [Hau92, BLW96, ABDCBH97, Bar98] together with our outer norm bounds. See Section 5.4 for its proof.

We next comment on the performance parameters appearing in Theorem 3.1. The parameter α controls the so-called generalization gap: the gap between the training error and the generalization error. The parameter δ controls the training error: we study those (a, W) with $\widehat{\mathcal{L}}(a, W) \leq \delta^2$. The parameter M is an (almost sure) upper bound on the labels; whereas R is an (quite mild) upper bound on internal weights required for the technical reasons, only for the case of sigmoid networks, see Theorem 2.1, Corollary 2.2; and the remarks following them.

The term $\zeta(\alpha, M, A, N)$ is a probability term appearing in the uniform convergence result (Proposition 5.3) that we employ. This proposition provides a control for the generalization gap uniformly over all two-layer neural networks with *bounded outer norm* like we investigate herein.

It is worth noting that in the regime $d \to \infty$ (which is a legitimate assumption for many existing guarantees in the field of machine learning) and for $N \geq d^{O(1)}$; $\xi(\alpha, M, A) = O(1)$ (with respect to d), provided that A = O(1) like we establish earlier. Namely, this object is simply a constant in d. Furthermore, while we made no attempts in simplifying it, it can potentially be improved. In the sigmoid and step cases, the value of A that we consider is indeed O(1). For the ReLU case, however, the situation is more involved; and a certain scaling which makes A = poly(d) is necessary, as we elaborate soon. Soon in Section 3.3, we investigate the probability term $\zeta(\alpha, M, A, N)$ appearing in (7). We show that provided N is sufficiently large (while remaining polynomial in d), the ζ term behaves like $\exp\left(-d^{O(1)}\right)$, thus it is indeed $o_d(1)$. Moreover, our analysis will also reveal that the dependence of N on d is quite mild; and is in fact near-linear in some important cases of interest.

In particular, the probability term $\zeta(\alpha, M, A, N)$ is $o_d(1)$ provided

Theorem 5.2 as well as the uniform generalization gap guarantee, Proposition 5.3, apply to activations with a bounded output; whereas the output of ReLU is potentially unbounded. In our proof, we bypass this by considering an auxiliary activation S-ReLU(·), which is a "saturated" version of the ReLU. Specifically, we let S-ReLU(x) = 0 for $x \le 0$, S-ReLU(x) = x for $0 < x \le 1$; and S-ReLU(x) = 1 for $x \ge 1$. We then rescale w_j to have $||w_j||_2 = 1/\sqrt{Cd}$ and multiply A by \sqrt{Cd} (we therefore consider $A = 4\sqrt{Cd}(\delta + 2M)/\mu^*$, \sqrt{Cd} times the bound appearing in Theorem 2.3). Note that this step is indeed valid due to the homogeneity of the ReLU activation, see also Section 2.2. Since $||X||_2 \le \sqrt{Cd}$ with probability at least $1 - \exp(-\Theta(d))$ and since $|w_j^T X| \le 1$ for $||w_j||_2 = 1/\sqrt{Cd}$ and $||X||_2 \le \sqrt{Cd}$ by Cauchy-Schwarz inequality; the output of this activation will, w.h.p., coincide with that of the ReLU activation. We then control the difference between the generalization errors for a pair of two-layer neural networks having the same architecture, the same number $\overline{m} \in \mathbb{N}$ of hidden units, the same weights (a, W); but different activations (one with ReLU(·) and the other with S-ReLU(·)). This done by a conditioning argument. See the proof for further details.

Similar to what we have noted previously for our outer norm bounds, Theorem 3.1 is also oblivious to (a) how the training is done and (b) the number \overline{m} of hidden units as long as $a_i \geq 0$, and $\widehat{\mathcal{L}}(a, W) \leq \delta^2$ for the learned network. Moreover, similar to prior cases, the extra conditional expectation requirement (30) is quite mild.

Our next focus is on the sample complexity required by Theorem 3.1. We show that they are indeed polynomial in d. Furthermore for some very important cases, they are even near-linear.

3.3 Sample Complexity Analysis

While the required sample complexity N can simply be inferred from Theorem 3.1, we spell out the implied scaling analysis below for convenience. In what follows, all asymptotic notations are w.r.t. the natural parameter d (namely the dimension) of the problem in the regime $d \to \infty$; and our goal is to ensure that the corresponding probability term is $1 - o_d(1)$ for an appropriate function $o_d(1)$. (It is worth noting though that our bounds will be in fact much stronger, e.g. $1 - \exp(-d^{O(1)})$.)

To that end, recall the term (6) with $\mathcal{M}=2$ appearing in Theorem 3.1:

$$\xi(\alpha, M, 2, A) = \frac{2^{23} \cdot c \cdot A^6 \cdot \max\{A, M\}^2}{\ln 2 \cdot \alpha^2} \cdot \ln\left(\frac{2^{11} \cdot A^3 \cdot \max\{A, M\}}{\alpha}\right). \tag{8}$$

Sigmoid and Step Networks

First, the outer norm bounds we establish indicate A = O(1). Hence, the "A parameter" considered in parts (a) and (c) of Theorem 3.1 are O(1). Moreover, M = O(1) (since it is not sound for the real-valued label Y to grow with dimension d). Treating α as a constant in d, we then obtain $\xi(\alpha, M, A) = O(1)$ for the term appearing in (8). Hence, in order to ensure that the probability term ζ appearing in (7) is $o_d(1)$, a necessary and sufficient condition is $N = \Omega(d \ln^2 N)$. We claim that it suffices to have

$$N = \Omega \left(d \ln^2 d \right). \tag{9}$$

Indeed, if N satisfies (9), then provided N remains polynomial in d, N = poly(d), it holds that

$$\ln^2 N = O(\ln^2 d) \implies d \ln^2 N = O(d \ln^2 d) = O(N).$$

We now investigate the sample complexity required by the corresponding outer norm bounds for the case of sigmoid and step networks.

Sigmoid networks. Note, in this case, that the dominant contribution to the probability term appearing in Theorem 2.1/Theorem 3.1(a) (other than ξ term) is $(3R\sqrt{Cd})^d \exp(-\Theta(N))$. Suppose first that $R = d^K$ where K = O(1) (namely R remains polynomial in d). Then

$$(3R\sqrt{Cd})^d \exp(-\Theta(N)) = \exp\left(-\Theta(N) + d\left(K + \frac{1}{2}\right)\ln d + d\ln(3\sqrt{C})\right)$$
$$= \exp\left(-\Theta(N) + \Theta(d\ln d) + o(d\ln d)\right).$$

provided $N = \Omega(d \ln d)$, this bound is indeed $o_d(1)$. Taking the maximum between this and (9), we obtain that it suffices to have $N = \Omega(d \ln^2 d)$, which is near-linear.

Suppose next that $R = \exp(d^K)$, like in Corollary 2.2. Then provided K > 0,

$$(3R\sqrt{Cd})^{d} \exp(-\Theta(N)) = \exp\left(-\Theta(N) + d^{K+1} + \frac{1}{2}d\ln d + d\ln(3\sqrt{C})\right)$$
$$= \exp\left(-\Theta(N) + d^{K+1} + o\left(d^{K+1}\right)\right).$$

Hence, provided $N = \Omega(d^{K+1})$, this bound is indeed $o_d(1)$. Taking the maximum between this and (9), we obtain that it suffices to have $N = \Omega(d^{K+1})$, which is polynomial in d.

Step networks. Treating the distributional parameter η appearing in Theorem 2.4/Theorem 2.3 as a constant in d, we have

$$\exp(-\Theta(N)) \left(\frac{6\sqrt{Cd}}{\eta}\right)^d = \exp\left(-\Theta(N) + \frac{1}{2}d\ln d + d\ln\left(\frac{6\sqrt{C}}{\eta}\right)\right)$$
$$= \exp\left(-\Theta(N) + \Theta(d\ln d) + o(d\ln d)\right).$$

Thus, provided $N = \Omega(d \ln d)$, this bound is indeed $o_d(1)$. Taking the maximum between this and (9), we obtain that it suffices to have $N = \Omega(d \ln^2 d)$, which, again, is near-linear.

ReLU Networks

The situation is more involved for the case of ReLU networks. We first study the ξ term (8). Treating $M, \alpha, C, \delta, \mu^* = O(1)$ (in d),

$$\xi\left(\alpha, M, 2, \frac{4\sqrt{C}(\delta + 2M)}{\mu^*}\sqrt{d}\right) = \Theta\left(d^4 \ln d\right).$$

Hence,

$$\zeta\left(\alpha, M, \frac{4\sqrt{Cd}(\delta + 2M)}{\mu^*}, N\right) = \exp\left(\Theta\left(d^4 \cdot \ln d \cdot d \cdot \ln^2(Nd)\right) - \Theta\left(\frac{N}{d}\right)\right) \\
= \exp\left(\Theta\left(d^5 \cdot \ln d \cdot \ln^2(Nd)\right) - \Theta\left(\frac{N}{d}\right)\right) \\
= \exp\left(\Theta\left(d^5 \cdot \ln^3 d\right) - \Theta\left(\frac{N}{d}\right)\right),$$

where we used the fact $\ln(Nd) = \Theta(\ln d)$ if N = poly(d). Thus, provided $N = \Omega\left(d^6 \ln^3 d\right)$, this bound is indeed $o_d(1)$. Inspecting next the term $(12\sqrt{Cd}/\boldsymbol{\mu}^*)^d \exp(-\Theta(N))$ appearing in the probability bound, we observe as long as $N = \Omega(d \ln d)$, this term is also $o_d(1)$. Taking the maximum of these two, it suffices to have $N = \Omega\left(d^6 \ln^3 d\right)$. This, again, is a polynomial in d; albeit having a slightly worse degree (of six).

4 Conclusion and Future Directions

We have studied two-layer NN models with sigmoid, ReLU, and step activations; and established that the *outer norm* of any such NN achieving a small training loss on a polynomially (in d) many data and having non-negative output weights is well-controlled. Our results are independent of the width \overline{m} of the network and the training algorithm; and are valid under very mild distributional assumptions on input/label pairs. We then leveraged the outer norm bounds we established to obtain good generalization guarantees for the networks we investigated. Our generalization results are obtained by employing earlier results on the fat-shattering dimension of such networks, and have good sample complexity bounds as we have discussed. In particular, for certain important cases of interest, we obtain near-linear sample guarantees.

We now provide future directions. As was already mentioned, our approach operates under mild distributional requirements; and can potentially handle different distributions as well as other activations, provided (rather natural) certain properties of these objects we leveraged remain in place.

A very important question is to which extent our approach applies to deeper networks. In what follows, we give a very brief argument demonstrating that for such an extension, one needs much more stringent regularity assumptions on the internal weights. Consider, as an example, a ReLU network with three hidden layers. Observe that the outputs of the neurons at the first hidden layer are non-negative as $ReLU(x) \geq 0$ for all $x \in \mathbb{R}$. Let us now focus on its second hidden layer, which takes weighted sums of the outputs of the first hidden layer. If all the weights in the second layer are negative, then upon passing to ReLU, one obtains all zeroes, forcing the final output to be zero. Now, let us assume, instead, that the weights of the second layer are such that the input to the ReLU functions are positive, though arbitrarily close to zero (this can potentially be achieved, e.g., by taking many small negative weights and few large positive weights in a way that ensures proper cancellation). If this holds, then even if the outer norm, $||a||_1$, is very large, one still obtains a bounded output at the end of the network. As demonstrated by this conceptual example, one indeed needs more stringent assumptions on the internal weights so as to address larger depth. At the present time, we are unable to have a complete resolution of necessary and sufficient assumptions for addressing deeper architectures (while maintaining the position that these assumptions must also be sound from a practical point of view).

Yet another important direction pertains to the non-negativity of the weights, and a crucial question is whether this assumption can be relaxed. We now provide a brief argument demonstrating that in full generality, this is not necessarily the case. Namely, strictly speaking, the non-negativity assumption is necessary. We focus on the so-called "teacher/student" setting, a setting that has been quite popular recently, see, e.g. [GAS+19]. In this setting, given i.i.d. input data $X_i \in \mathbb{R}^d$, $1 \leq i \leq N$, a teacher network $(a^*, W^*) \in \mathbb{R}^{m^*} \times \mathbb{R}^{m^* \times d}$ with $m^* \in \mathbb{N}$ neurons and activation $\sigma(\cdot)$ generates the labels Y_i . That is, $Y_i = \sum_{1 < j < m^*} a_j^* \sigma((w_j^*)^T X_i)$. A student network with an $\overline{m} \in \mathbb{N}$ number of hidden units (where \overline{m} is not necessarily equal to m^*) is then "trained" by minimizing the objective function (2) on the data (X_i, Y_i) , $1 \le i \le N$; and the resulting network is then used for predicting the unseen data. We now construct a wider student network interpolating the data whose vector of output weights has arbitrarily large norm, by introducing many cancellations. Fix $z \in \mathbb{N}$, a non-zero $v \in \mathbb{R}^d$; and v > 0. Construct a new network $(\overline{a}, \overline{W})$ on $m^* + 2z$ neurons as follows. Set $\overline{a_j} = a_j^*$ and $\overline{W_j} = W_j^*$ for $1 \leq j \leq m^*$. For any $m^* + 1 \le j \le m^* + 2z$, set $\overline{a_j} = \nu$ if j is even, and $-\nu$, if j is odd. At the same time, set $\overline{W_j} = v$ for $m^* + 1 \le j \le m^* + 2z$. This network interpolates the data while $\|\overline{a}\|_1 = \|a^*\|_1 + 2z\nu$. Hence, $||\overline{a}||_1$ can be made arbitrarily large by amplifying z and/or $\nu > 0$. In particular, in full generality, such a non-negativity assumption is indeed necessary. It is worth noting, however, that the example above is a somewhat tailored one involving many dependencies/cancellations. It might still be possible to establish similar bounds for the case of potentially negative weights under more stringent constraints on them which prevent such cancellations.

5 Proofs

5.1 Proof of Theorem 2.1

Proof of Theorem 2.1. Observe that

$$\mathbb{P}(\mathcal{E}_0) \ge 1 - o_N(1) \quad \text{for} \quad \mathcal{E}_0 \triangleq \{ \sum_{i=1}^N |Y_i| \le 2MN \}, \tag{10}$$

using the weak law of large numbers [Dur19, Thm 2.2.14]. Next, let $(a, W) \in \mathcal{S}(\delta, R)$. Then there exists an $\overline{m} \in \mathbb{N}$ such that $(a, W) \in \mathcal{S}(\overline{m}, \delta, R)$. Applying Cauchy-Schwarz inequality, $\sum_{1 \leq i \leq N} \left| Y_i - \sum_{1 \leq j \leq \overline{m}} a_j \text{SGM}\left(w_j^T X_i\right) \right| \leq N\delta$. Next, by the triangle inequality and the fact $\sum_i |Y_i| \leq 2MN$ on \mathcal{E}_0 ,

$$\sum_{1 \le i \le N} \sum_{1 \le j \le \overline{m}} a_j \text{SGM}\left(w_j^T X_i\right) \le N(\delta + 2M),\tag{11}$$

on the event \mathcal{E}_0 . Now, let \mathcal{N}_{ϵ} be an ϵ -net for $B_2(0,R)$, $\epsilon > 0$ to be tuned appropriately. Using Theorem 1.3, one can ensure $|\mathcal{N}_{\epsilon}| \leq (3R/\epsilon)^d$. Next, fix any $\widehat{w} \in \mathcal{N}_{\epsilon}$, and set $\overline{Z_i} \triangleq \widehat{w}^T X_i$, $1 \leq i \leq N$. Since $\overline{Z_i}$ is symmetric, $\mathbb{P}(\overline{Z_i} \geq 0) = \mathbb{P}(-\overline{Z_i} \leq 0) = \mathbb{P}(\overline{Z_i} \leq 0)$, implying $\mathbb{P}(\overline{Z_i} \geq 0) \geq \frac{1}{2}$. Define now $Z_i \triangleq \mathbb{I}\left\{\overline{Z_i} \geq 0\right\}$. Since Z_i "stochastically dominates" Bernoulli(1/2), we have $\mathbb{P}\left(\sum_{1\leq i\leq N} Z_i \geq N/3\right) \geq \mathbb{P}\left(\text{Binomial}\left(N, 1/2\right) \geq N/3\right) \geq 1 - \exp\left(-\Theta(N)\right)$. The last inequality is due to standard large deviations bounds. Taking a union bound over the net \mathcal{N}_{ϵ} , we obtain

$$\mathbb{P}(\mathcal{E}_1) \ge 1 - (3R/\epsilon)^d \exp\left(-\Theta(N)\right), \quad \text{where} \quad \mathcal{E}_1 \triangleq \bigcap_{\widehat{w} \in \mathcal{N}_{\epsilon}} \left\{ \sum_{1 \le i \le N} \mathbb{1} \left\{ \widehat{w}^T X_i \ge 0 \right\} \ge N/3 \right\}.$$
 (12)

Furthermore, another union bound over the data yields

$$\mathbb{P}(\mathcal{E}_2) \ge 1 - N \exp\left(-\Theta(d)\right), \quad \text{where} \quad \mathcal{E}_2 \triangleq \{||X_i||_2^2 \le Cd, 1 \le i \le N\}.$$
 (13)

We now choose $\epsilon = 1/\sqrt{Cd}$. We claim that on the event $\mathcal{E}_1 \cap \mathcal{E}_2$, it is the case that for every $w \in B_2(0,R)$; $\sum_{1 \leq i \leq N} \mathbbm{1} \left\{ w^T X_i \geq -1 \right\} \geq \frac{N}{3}$. Let $w \in B_2(0,R)$, and $\widehat{w} \in \mathcal{N}_{\epsilon}$ be such that $||w - \widehat{w}||_2 \leq \epsilon = (Cd)^{-1/2}$. Using Cauchy-Schwarz inequality, $|\widehat{w}^T X_i - w^T X_i| \leq ||X_i||_2 (Cd)^{-1/2} \leq 1$, where $||X_i||_2 \leq \sqrt{Cd}$ due to the event \mathcal{E}_2 (13). In particular, if $\widehat{w}^T X_i \geq 0$, then $w^T X_i \geq -1$. Hence $\sum_{1 \leq i \leq N} \mathbbm{1} \left\{ w^T X_i \geq -1 \right\} \geq \sum_{1 \leq i \leq N} \mathbbm{1} \left\{ \widehat{w}^T X_i \geq 0 \right\} \geq \frac{N}{3}$. Using now the fact $a_j \geq 0$, and SGM(·) ≥ 0 for the sigmoid activation, we arrive at

$$\sum_{1 < j < \overline{m}} a_j \sum_{1 < i < N} \text{SGM}\left(w_j^T X_i\right) \ge \frac{N}{3} \cdot \text{SGM}(-1) \cdot \sum_{1 < j < n} a_j. \tag{14}$$

We now combine the facts $SGM(-1) = (1+e)^{-1}$, (11) and (14), to obtain that on the event $\mathcal{E}_0 \cap \mathcal{E}_1 \cap \mathcal{E}_2$,

$$\sum_{1 \le j \le \overline{m}} a_j \le 3(1+e)(\delta+2M).$$

Since the event $\mathcal{E}_0 \cap \mathcal{E}_1 \cap \mathcal{E}_2$ holds with probability at least $1 - \left(3R\sqrt{Cd}\right)^d \exp\left(-\Theta(N)\right) - N\exp\left(-\Theta(d)\right) - o_N(1)$ by a union bound, the proof is complete.

5.2 Proof of Theorem 2.3

Proof of Theorem 2.3. Recall from (10) the event $\mathcal{E}_0 = \{\sum_{1 \leq i \leq N} |Y_i| \leq 2MN \}$ where $\mathbb{P}(\mathcal{E}_0) \geq 1 - o_N(1)$.

Let $(a, W) \in \mathcal{G}(\delta)$. Then, for some $\overline{m} \in \mathbb{N}$, $(a, W) \in \mathcal{G}(\overline{m}, \delta)$. Using Cauchy-Schwarz inequality and the triangle inequality like in the beginning of the proof of Theorem 2.1; we first establish that on the event \mathcal{E}_0 , the following holds:

$$\sum_{1 \le i \le N} \sum_{1 \le j \le \overline{m}} a_j \text{ReLU}\left(w_j^T X_i\right) \le N(\delta + 2M). \tag{15}$$

Next, let \mathcal{N}_{ϵ} be a ϵ -net for $B_2(0,1)$, $\epsilon > 0$ to be tuned. Using Theorem 1.3, one can ensure $|\mathcal{N}_{\epsilon}| \leq (3/\epsilon)^d$. Fix any $\widehat{w} \in \mathcal{N}_{\epsilon}$. Consider the i.i.d. random variables $Y_{\widehat{w},i} \triangleq \text{ReLU}\left(\widehat{w}^T X_i\right), i \in [N]$. The condition (b) on the distribution of $Y_{\widehat{w},i}$ ensures that a large deviations bound is applicable. This, together with the condition (a) yield $\mathbb{P}\left(\sum_{1\leq i\leq N} Y_{\widehat{w},i} \geq \frac{1}{2} \mu^* N\right) \geq 1 - \exp\left(-\Theta(N)\right)$. Due to the distributional assumption, the lower bound is uniform in $\widehat{w} \in \mathcal{N}_{\epsilon}$.

Taking now a union bound over $\widehat{w} \in \mathcal{N}_{\epsilon}$, we obtain $\mathbb{P}(\mathcal{E}_1) \geq 1 - (3/\epsilon)^d \exp{(-\Theta(N))}$ where $\mathcal{E}_1 \triangleq \bigcap_{\widehat{w} \in \mathcal{N}_{\epsilon}} \left\{ \sum_{1 \leq i \leq N} \text{ReLU}\left(\widehat{w}^T X_i\right) \geq \frac{1}{2} \mu^* N \right\}$. Another union bound over data X_i , $1 \leq i \leq N$, yields that $\mathbb{P}(\mathcal{E}_2) \geq 1 - N \exp{(-\Theta(d))}$ where $\mathcal{E}_2 \triangleq \{||X_i||_2^2 \leq Cd, 1 \leq i \leq N\}$. Choose $\epsilon \triangleq \frac{\mu^*}{4\sqrt{Cd}}$, and assume in the remainder that we are on the event $\mathcal{E}_1 \cap \mathcal{E}_2$. Next, observe

Choose $\epsilon \triangleq \frac{\mu^*}{4\sqrt{Cd}}$, and assume in the remainder that we are on the event $\mathcal{E}_1 \cap \mathcal{E}_2$. Next, observe that ReLU is 1-Lipschitz: $|\text{ReLU}(x) - \text{ReLU}(y)| = \left|\frac{x+|x|}{2} - \frac{y+|y|}{2}\right| \leq |x-y|$, using triangle inequality twice. Now, fix $any \ w \in B_2(0,1)$. Let $\widehat{w} \in \mathcal{N}_{\epsilon}$ be the member of the net closest to w. Using the Lipschitz property, and the Cauchy-Schwarz, we obtain $|\text{ReLU}(w^TX_i) - \text{ReLU}(\widehat{w}^TX_i)| \leq |w^TX_i - \widehat{w}^TX_i| \leq ||w - \widehat{w}||_2 \cdot ||X_i||_2 \leq \frac{\mu^*}{4}$. Consequently, $|\text{ReLU}(w^TX_i)| \geq ||w^TX_i||_2 \leq |w^TX_i| \leq ||w - \widehat{w}||_2 \cdot ||x_i||_2 \leq \frac{\mu^*}{4}$. Summing this over $|1 \leq i \leq N|$, we have $|\sum_{1 \leq i \leq N} ||\text{ReLU}(w^TX_i)||_2 \leq \sum_{1 \leq i \leq N} ||\text{ReLU}(\widehat{w}^TX_i)||_2 = \frac{\mu^*}{4}N$. Using $|a_j| \geq 0$, we obtain by taking $|w_j|$ in place of |w|:

$$\sum_{1 \le j \le \overline{m}} a_j \sum_{1 \le i \le N} \text{ReLU}\left(w_j^T X_i\right) \ge \frac{\mu^*}{4} N \sum_{1 \le j \le \overline{m}} a_j. \tag{16}$$

Combining (15) and (16), we obtain that on the event $\mathcal{E}_0 \cap \mathcal{E}_1 \cap \mathcal{E}_2$,

$$\sum_{1 < j < \overline{m}} a_j \le 4(\delta + 2M) \left(\boldsymbol{\mu}^*\right)^{-1}.$$

Since the event $\mathcal{E}_0 \cap \mathcal{E}_1 \cap \mathcal{E}_2$ holds with probability at least $1 - \left(12\sqrt{Cd}(\boldsymbol{\mu}^*)^{-1}\right)^d \exp\left(-\Theta(N)\right) - N \exp\left(-\Theta(d)\right) - o_N(1)$ via a union bound, we complete the proof.

5.3 Proof of Theorem 2.4

Proof of Theorem 2.4. The proof is quite similar to that of the proof of Theorems 2.1/2.3, and is provided for completeness.

Again, recall from (10) the event $\mathcal{E}_0 = \{\sum_{1 \leq i \leq N} |Y_i| \leq 2MN\}$ where $\mathbb{P}(\mathcal{E}_0) \geq 1 - o_N(1)$. Then, take an $(a, W) \in \mathcal{H}(\delta)$. There exists an $\overline{m} \in \mathbb{N}$ such that $(a, W) \in \mathcal{H}(\overline{m}, \delta)$. Using

again Cauchy-Schwarz inequality and the triangle inequality like in the beginning of the proof of Theorems 2.1/2.3; we have that on the event \mathcal{E}_0 , the following holds:

$$\sum_{1 \leq i \leq N} \left| Y_i - \sum_{1 \leq j \leq \overline{m}} a_j \mathtt{Step}\left(w_j^T X_i
ight)
ight| \leq N \delta.$$

This, together with (a) the fact that the labels are bounded, $|Y_i| \leq M$; and (b) the triangle inequality; then yields

$$\sum_{1 \le i \le N} \sum_{1 \le j \le \overline{m}} a_j \operatorname{Step}\left(w_j^T X_i\right) \le N\left(\delta + M\right). \tag{17}$$

Let \mathcal{N}_{ϵ} be an ϵ -net for $B_2(0,1)$, where $\epsilon > 0$ to be tuned appropriately. Using Theorem 1.3, one can ensure $|\mathcal{N}_{\epsilon}| \leq (3/\epsilon)^d$.

Next, fix any $\widehat{w} \in \mathcal{N}_{\epsilon}$; and set $Z_i \triangleq \mathbb{I}\left\{\widehat{w}^T X_i \geq \eta\right\}$, $1 \leq i \leq N$ (where we drop the dependence of Z_i on \widehat{w} for convenience). Evidently, Z_i is an i.i.d. collection of Bernoulli random variables, with $\mathbb{E}[Z_i] \geq \eta$ (due to the assumption on the distribution of X). Hence, using standard concentration results, $\mathbb{P}\left(\sum_{1\leq i\leq N} Z_i \geq N\eta/2\right) \geq 1 - \exp\left(-\Theta(N)\right)$. Moreover, the lower bound is, again, uniform in \widehat{w} via an exact same stochastic domination argument, like in the proof of Theorem 2.1.

Taking now a union bound over the net \mathcal{N}_{ϵ} ,

$$\mathbb{P}(\mathcal{E}_1) \ge 1 - (3/\epsilon)^d \exp\left(-\Theta(N)\right), \quad \text{where} \quad \mathcal{E}_1 \triangleq \bigcap_{\widehat{w} \in \mathcal{N}_{\epsilon}} \left\{ \sum_{1 \le i \le N} \mathbb{1} \left\{ \widehat{w}^T X_i \ge \eta \right\} \ge N\eta/2 \right\}. \quad (18)$$

Furthermore, another union bound over data, $1 \le i \le N$, yields

$$\mathbb{P}(\mathcal{E}_2) \ge 1 - N \exp\left(-\Theta(d)\right), \quad \text{where} \quad \mathcal{E}_2 \triangleq \left\{ ||X_i||_2^2 \le Cd, 1 \le i \le N \right\}. \tag{19}$$

We now choose $\epsilon = \frac{\eta}{2\sqrt{Cd}}$; and assume in the remainder that we are on the event $\mathcal{E}_1 \cap \mathcal{E}_2$.

Fix any $w \in B_2(0,1)$; and let $\widehat{w} \in \mathcal{N}_{\epsilon}$ be such that $||w - \widehat{w}||_2 \leq \frac{\eta}{2\sqrt{Cd}}$. Using Cauchy-Schwarz inequality, $|\widehat{w}^T X_i - w^T X_i| \leq ||w - \widehat{w}||_2 ||X_i||_2 \leq \eta/2$, for every $i \in [N]$, since the event we are on is a subset of \mathcal{E}_2 in (19). Observe now that $\{\widehat{w}^T X \geq \eta\} \subseteq \{w^T X \geq \eta/2\}$. Thus, on the event $\mathcal{E}_1 \cap \mathcal{E}_2$, it holds that

$$\sum_{1 < i < N} \mathbb{1}\{w^T X_i \ge \eta/2\} \ge \sum_{1 < i < N} \mathbb{1}\{\widehat{w}^T X_i \ge \eta/2\} \ge N\eta/2.$$

Since $w \in B_2(0,1)$ is arbitrary, and $Step(w^T X_i) = 1$ if $w^T X_i \ge \eta/2 > 0$, we arrive at

$$\sum_{1 \le j \le \overline{m}} a_j \sum_{1 \le i \le N} \operatorname{Step}\left(w_j^T X_i\right) \ge \frac{N\eta}{2} \sum_{1 \le j \le \overline{m}} a_j. \tag{20}$$

We now combine (17) and (20) to arrive at the conclusion that on the event $\mathcal{E}_1 \cap \mathcal{E}_2$, it holds

$$\sum_{1 < j < \overline{m}} a_j \le 2(\delta + M)\eta^{-1}.$$

Finally, we combine (18) (with $\epsilon = \frac{\eta}{2\sqrt{Cd}}$) and (19) via a union bound; and arrive at the conclusion that $\mathbb{P}(\mathcal{E}_1 \cap \mathcal{E}_2) \geq 1 - \left(6\sqrt{Cd} \cdot (\eta)^{-1}\right)^d \exp\left(-\Theta(N)\right) - N \exp\left(-\Theta(d)\right) - o_N(1)$. This concludes the proof.

5.4 Proof of Theorem 3.1

In this section, we establish Theorem 3.1. We build upon earlier results by Bartlett [Bar98] and Bartlett, Long, and Williamson [BLW96]. For the results we cite from the latter, the numbers recorded below are from the version accessed at http://phillong.info/publications/fatshat.pdf¹.

The FSD of the Networks with a Bounded Outer Norm. We now recall the definition of the fat-shattering dimension (FSD), verbatim from [Bar98], for convenience.

Definition 5.1. Let X be an input space, H be a class of real-valued functions defined on X (that is, H consists of functions $f: X \to \mathbb{R}$). Fix a $\gamma > 0$, which is a certain scale parameter. We say that a sequence (x_1, x_2, \ldots, x_m) of m points from X is γ -shattered by H if there is an $r = (r_1, \ldots, r_m) \in \mathbb{R}^m$ such that, for all $b = (b_1, \ldots, b_m) \in \{-1, 1\}^m$ there is an $h \in H$ satisfying $(h(x_i) - r_i)b_i \geq \gamma$. Define the fat-shattering dimension of H as the function

$$FSD_{H}(\gamma) \triangleq \max \Big\{ m : H \ \gamma\text{-shatters some } x \in X^{m} \Big\}.$$
 (21)

We next record the following result.

Theorem 5.2. [Bar98, Corollary 24] Let $\mathcal{M} > 0$, and $\sigma : \mathbb{R} \to [-\mathcal{M}/2, \mathcal{M}/2]$ be a non-decreasing function. Define a class F of functions on \mathbb{R}^d by

$$F \triangleq \{X \mapsto \sigma\left(w^T X + w_0\right) : w \in \mathbb{R}^d, w_0 \in \mathbb{R}\}$$

and let

$$H(A) \triangleq \left\{ \sum_{1 \leq j \leq \overline{m}} a_j f_j : \overline{m} \in \mathbb{N}, f_j \in F, ||a||_1 \leq A \right\}$$

where $A \geq 1$. Then for every $\gamma \leq \mathcal{M}A$,

$$FSD_{H(A)}(\gamma) \le \frac{c\mathcal{M}^2 A^2 d}{\gamma^2} \ln\left(\frac{\mathcal{M}A}{\gamma}\right)$$

for some universal constant c > 0.

Here $a = (a_j : 1 \le j \le \overline{m}) \in \mathbb{R}^{\overline{m}}$ is the vector of output weights, $||a||_1$ is the outer norm; and $\gamma > 0$ is a certain *scale* parameter. Observe that H(A) is precisely the class of two-layer NN with activation function $\sigma(\cdot)$ whose outer norm is at most A. Per Theorem 5.2, the FSD of the class of two-layer networks with **bounded outer norm** is upper bounded by an explicit quantity.

Some Extra Notation on Covering Numbers. We next introduce several quantities verbatim from [BLW96]. Let W be an arbitrary set, and $f: W \to \mathbb{R}$ be any function. For any $w = (w_1, \ldots, w_N) \in W^N$, denote by $f|_w$ the N-tuple $(f(w_1), f(w_2), \ldots, f(w_N)) \in \mathbb{R}^N$. For a class \mathcal{C} of functions $f: W \to \mathbb{R}$, let $\mathcal{C}|_w \subseteq \mathbb{R}^N$ denotes the set

$$C|_{w} \triangleq \left\{ f|_{w} : f \in C \right\} = \left\{ \left(f(w_{1}), \dots, f(w_{N}) \right) : f \in C \right\} \subseteq \mathbb{R}^{N}.$$
 (22)

¹See the archived version at http://web.archive.org/web/20200921180645/http://phillong.info/publications/fatslif the link above is expired.

Next, recall the covering numbers from Definition 1.2. Throughout this section, and in particular the proof of Theorem 3.1, we take the metric ρ appearing in Definition 1.2 to be the normalized ℓ_1 distance: for any $w, \bar{w} \in \mathbb{R}^N$, set

$$\rho(w, \bar{w}) = \frac{1}{N} \sum_{1 \le i \le N} |w_i - \bar{w}_i|.$$

For any $U \subseteq \mathbb{R}^N$, denote by $\mathcal{N}(\epsilon, U)$ the covering number of U (at scale ϵ) with respect to the metric ρ above. That is, $\mathcal{N}(\epsilon, U)$ is the cardinality of the smallest $\mathcal{N}_{\epsilon} \subset U$ (if finite) such that for every $w \in U$, there exists a $\overline{w} \in N_{\epsilon}$ with $\rho(w, \overline{w}) = \frac{1}{N} \sum_{1 \leq i \leq N} |w_i - \overline{w}_i| \leq \epsilon$. (It is worth noting that here we flipped the order of arguments in \mathcal{N} appearing in Definition 1.2. The rationale for this is to be consistent with the notation of Bartlett et al. [BLW96].)

Throughout this section, we often consider the following special case of $\mathcal{N}(\cdot, \cdot)$: we employ $\mathcal{N}(\cdot, \mathcal{C}|_w)$ for appropriate classes \mathcal{C} of functions where w is an element of the Euclidean space \mathbb{R}^N for some N.

We now establish the following proposition which provides a control for the generalization gap uniformly over all two-layer NN models with bounded outer norm.

Proposition 5.3. Let $M, \mathcal{M}, A > 0$; $\sigma : \mathbb{R} \to [-\mathcal{M}/2, \mathcal{M}/2]$ be a non-decreasing activation function; and \mathcal{D} be an arbitrary distribution on $\mathbb{R}^d \times \mathbb{R}$ for the input/label pairs (X, Y) where $|Y| \leq M$ almost surely. Recall the class H(A) of two-layer neural networks with activation σ and outer norm at most A from Theorem 5.2; and let (X_i, Y_i) , $1 \leq i \leq N$, be i.i.d. samples drawn from \mathcal{D} . Then for any $\alpha > 0$, with probability at least

$$1 - 4\exp\left(\xi(\alpha, M, \mathcal{M}, A) \cdot d \cdot \ln^2\left(\frac{576N\mathcal{M}^2A^2\max\{\mathcal{M}A, 2M\}}{\alpha}\right) - \frac{\alpha^2N}{64\max\{\mathcal{M}A, 2M\}^2}\right)$$

over the draw of the training data (X_i, Y_i) , $1 \le i \le N$, it holds that

$$\sup_{\varphi \in H(A)} \left| \frac{1}{N} \sum_{1 < i < N} \left(\varphi(X_i) - Y_i \right)^2 - \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\left(\varphi(X) - Y \right)^2 \right] \right| \le \alpha,$$

provided

$$N \geq 64 \cdot 128c \cdot \frac{\mathcal{M}^6 A^6 \max{\{\mathcal{M}A, 2M\}^2}}{\alpha^2} d$$

Here, c, c' > 0 are absolute constants, the term ξ is introduced in (6) and the expectation is taken with respect to a fresh sample $(X, Y) \sim \mathcal{D}$ independent of (X_i, Y_i) , $1 \le i \le N$.

It is worth noting that while we made no attempts for simplifying the constants appearing throughout Proposition 5.3, we believe that they can be improved.

Proof of Proposition 5.3. We first provide a result established originally in [Hau92, Theorem 3,p. 107].

Theorem 5.4. Let X, Y be sets; G be a PH-permissible class of [0,T]-valued functions defined on $Z \triangleq X \times Y$ where $T \in \mathbb{R}^+$, and P be any distribution on Z. Suppose Z_i , $1 \leq i \leq N$, are i.i.d. samples from P. Then for any $\alpha > 0$, with probability at least

$$1 - 4 \left(\sup_{z \in Z^{2N}} \mathcal{N} \left(\frac{\alpha}{16}, G \Big|_z \right) \right) \cdot \exp \left(-\alpha^2 N / 64T^2 \right)$$

over data Z_i , $1 \le i \le N$, it holds that

$$\sup_{g \in G} \left| \frac{1}{N} \sum_{1 < i < N} g(Z_i) - \mathbb{E}_{Z \sim P}[g(Z)] \right| \le \alpha,$$

where $\mathbb{E}[g(Z)]$ is taken with respect to a fresh sample (namely a sample drawn from P, and independent of Z_i).

The version we record above is verbatim from [BLW96, Theorem 13]. (The parameters M and m in [BLW96] are replaced, respectively, with the parameters T and N above.)

Here, PH-permissible refers to a rather mild measurability constraint², see [Hau92, Section 9.2]. The precise details of this technicality are immaterial to us; and it is satisfied for our purposes. Moreover, $\mathcal{N}(\cdot, \cdot)$ is the covering numbers quantity defined above.

In what follows, we take $X = \mathbb{R}^d$, Y = [0, M] (recall that the labels are bounded almost surely by M) thus $Z = \mathbb{R}^d \times [0, M]$ and we set P to simply be \mathcal{D} , the distribution from which the data are drawn. We then set

$$G \triangleq \left\{ (\varphi(X) - Y)^2 : X \in \mathbb{R}^d, Y \in [0, M], \varphi(\cdot) \in H(A) \right\}, \tag{23}$$

and take T to be $\max\{\mathcal{M}A, 2M\}^2$ (see below, in particular (26)). This is nothing but the ℓ_2 error obtained for predicting the label Y with $\varphi(X)$, with X being the input and $\varphi(\cdot)$ being the "predictor".

Upon inserting these parameters in Theorem 5.4, we obtain immediately

$$\sup_{\varphi \in H(A)} \left| \frac{1}{N} \sum_{1 \le i \le N} \left(\varphi(X_i) - Y_i \right)^2 - \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[(\varphi(X) - Y)^2 \right] \right| \le \alpha \tag{24}$$

with probability at least

$$1 - 4\left(\sup_{z \in Z^{2N}} \mathcal{N}\left(\frac{\alpha}{16}, G\Big|_z\right)\right) \cdot \exp\left(-\frac{\alpha^2 N}{64 \max\{\mathcal{M}A, 2M\}^2}\right) \tag{25}$$

over data $Z_i = (X_i, Y_i) \sim \mathcal{D}$, $1 \leq i \leq N$. Above, we used the facts (a) $|Y| \leq M$ almost surely; and (b) for any $\varphi \in H(A)$, it is the case $\varphi(X) = \sum_{1 \leq j \leq \overline{m}} a_j \sigma\left(w_j^T X\right)$ (for an $\overline{m} \in \mathbb{N}$ and $w_j \in \mathbb{R}^d$, $1 \leq j \leq \overline{m}$), where $||a||_1 \leq A$ and $\sup_{x \in \mathbb{R}} |\sigma(x)| \leq \mathcal{M}/2$. These together with the triangle inequality yield

$$-\mathcal{M}A/2 \le \varphi(X) \le \mathcal{M}A/2$$
 and $-M \le Y \le M$.

Hence,

$$-\max\{\mathcal{M}A/2,M\} \leq \varphi(X), Y \leq \max\{\mathcal{M}A/2,M\} \implies \left(\varphi(X)-Y\right)^2 \leq \max\{\mathcal{M}A,2M\}^2,$$

thus T can be taken as

$$T \triangleq \max\{\mathcal{M}A, 2M\}^2. \tag{26}$$

We next study covering number quantity $\sup_{z\in Z^{2N}} \mathcal{N}\left(\alpha/16, G|_{z}\right)$ appearing in (25). For this, we rely on the following result taken verbatim from [BLW96, Lemma 17].

The letters H and P stand, respectively, for Haussler and Pollard—who gave a preliminary version of Theorem 5.4.

Lemma 5.5. Let X be a set, and F be a set of functions from X to [0,1]. Then for any $\epsilon > 0$ and any $N \in \mathbb{N}$, if $a \leq 0$ and $b \geq 1$, we have

$$\sup_{z \in (X \times [a,b])^N} \mathcal{N}\left(\epsilon, (\ell_F) \Big|_z\right) \leq \sup_{x \in X^N} \mathcal{N}\left(\frac{\epsilon}{3|b-a|}, F\Big|_x\right).$$

Here, $\ell_f(x,y) = (f(x) - y)^2$, $\ell_F = \{\ell_f : f \in F\}$, and for $z = (z_1, \dots, z_N)$ (where $z_i = (x_i, y_i)$),

$$(\ell_F)|_z = \left\{ \left(\ell_f(x_i, y_i)^2 : 1 \le i \le N \right) : f \in F \right\},$$

which is the notation introduced in (22) with $\mathcal{C} := \ell_F$ and w := z.

We take F = H(A) and $\ell_F = G$ to arrive at

$$\sup_{z \in Z^{2N}} \mathcal{N}\left(\frac{\alpha}{16}, G\Big|_{z}\right) \le \sup_{x \in \left(\mathbb{R}^d\right)^{2N}} \mathcal{N}\left(\frac{\alpha}{32\mathcal{M}A \max\{\mathcal{M}A, 2M\}}, H(A)|_{x}\right). \tag{27}$$

Here, in addition to inserting $\alpha/16$, we also rescaled ϵ so as to reflect the fact that the functions in H(A) take values in $[0, \mathcal{M}A]$. (While all the bounds established by Bartlett et al. in [BLW96] assume the output space to be [0, 1], they extend in a straightforward manner to any output spaces of form [L, U] by rescaling corresponding parameters. This is already noted in the beginning of [BLW96, Section 6].)

We next record yet another result by Bartlett et al. [BLW96, Corollary 16].

Lemma 5.6. Let F be a class of [0,1]-valued functions defined on X, $0 < \epsilon < 1/2$ and $2N \ge FSD_F(\epsilon/4)$. Then,

$$\sup_{x \in X^N} \mathcal{N}\left(\epsilon, F|_x\right) \le \exp\left(\frac{2}{\ln 2} \mathrm{FSD}_F(\epsilon/4) \ln^2 \frac{9N}{\epsilon}\right),$$

where the quantity $FSD_F(\cdot)$ stands for the fat-shattering dimension introduced in (21).

(While we again skip the proof of this lemma, it is worth noting that it is obtained by combining two earlier results by Alon et al. [ABDCBH97, Lemmas 14,15].)

Taking now $X = \mathbb{R}^d$ and F = H(A); rescaling ϵ to $\frac{\epsilon}{\mathcal{M}A}$; and then plugging

$$\epsilon = \frac{\alpha}{32\mathcal{M}A\max\{\mathcal{M}A, 2M\}}$$

as in (27), we obtain

$$\sup_{x \in (\mathbb{R}^d)^{2N}} \mathcal{N}\left(\frac{\alpha}{32\mathcal{M}A \max\{\mathcal{M}A, 2M\}}, H(A)|_{x}\right)$$

$$\leq \exp\left(\frac{2}{\ln 2} FSD_{H(A)}\left(\frac{\alpha}{128\mathcal{M}^2 A^2 \max\{\mathcal{M}A, 2M\}}\right) \ln^2\left(\frac{576N\mathcal{M}^2 A^2 \max\{\mathcal{M}A, 2M\}}{\alpha}\right)\right). \tag{28}$$

We finally apply Theorem 5.2 above to upper bound the FSD term appearing in (28). Provided

$$\frac{\alpha}{128\mathcal{M}^2A^2\max\{\mathcal{M}A,2M\}} \le \mathcal{M}A \Leftrightarrow \alpha \le 128\mathcal{M}^3A^3\max\{\mathcal{M}A,2M\}$$

it holds that

$$FSD_{H(A)}\left(\frac{\alpha}{128\mathcal{M}^2 A^2 \max\{\mathcal{M}A, 2M\}}\right) \leq \frac{128^2 c \mathcal{M}^6 A^6 \max\{\mathcal{M}A, 2M\}^2}{\alpha^2} d \cdot \ln\left(\frac{128\mathcal{M}^3 A^3 \max\{\mathcal{M}A, 2M\}}{\alpha}\right)$$
(29)

where c > 0 is the absolute constant appearing in Theorem 5.2.

Finally, combining the chain of equations (25), (27), (28), and (29), we complete the proof. \Box

We finally provide a technical lemma to be used in the proof for the ReLU case.

Lemma 5.7. Suppose that the distribution of $X \in \mathbb{R}^d$ satisfies the assumptions of Theorem 2.3. Then,

$$\lambda(d) \triangleq \sup_{w \in \mathbb{R}^d: ||w||_2 = 1/\sqrt{Cd}} \mathbb{E}\left[\left|w^T X\right|^2 \mathbb{1}\left\{||X||_2^2 > Cd\right\}\right] \le \exp\left(-\Theta(d)\right). \tag{30}$$

The scaling $||w||_2 = 1/\sqrt{Cd}$ is required for technical reasons for the proof of the part (b) of Theorem 3.1.

Proof of Lemma 5.7. Define

$$\bar{\lambda}(d) \triangleq \sup_{w \in \mathbb{R}^{d}: ||w||_2 = 1} \mathbb{E}\Big[\left| w^T X \right|^2 \mathbb{1} \left\{ ||X||_2^2 > Cd \right\} \Big].$$

Clearly $\bar{\lambda}(d) = Cd\lambda(d)$. Since C = O(1), it suffices to prove $\bar{\lambda}(d) \leq \exp(-\Theta(d))$.

Next, fix a $w \in \mathbb{R}^d$ with $||w||_2 = 1$. Observe that using the inequality $e^x \ge 1 + x$, we obtain

$$e^{rw^TX} + e^{-rw^TX} \ge r |w^TX|$$
, for any $r \ge 0$.

Using the chain of inequalities

$$8(a^4+b^4) \ge 4(a^2+b^2)^2 \ge (a+b)^4$$

both due to Cauchy-Schwarz, we thus obtain

$$\frac{8}{r^4} \left(e^{4rw^T X} + e^{-4rw^T X} \right) \ge \left| w^T X \right|^4.$$

Now, take r = s/4 and then take the expectation of both sides to obtain

$$\frac{2048}{s^4} \Big(M_1(s) + M_2(s) \Big) \ge \mathbb{E} \Big[\Big| w^T X \Big|^4 \Big], \tag{31}$$

where $M_1(s)$ and $M_2(s)$ are defined in Theorem 2.3. Thus,

$$\mathbb{E}\left[\left|w^{T}X\right|^{2}\mathbb{1}\left\{||X||_{2}^{2} > Cd\right\}\right]^{2} \leq \mathbb{E}\left[\left|w^{T}X\right|^{4}\right]\mathbb{E}\left[\mathbb{1}\left\{||X||_{2}^{2} > Cd\right\}^{2}\right]$$
(32)

$$= \mathbb{E}\left[\left|w^T X\right|^4\right] \mathbb{P}\left(||X||_2^2 > Cd\right) \tag{33}$$

$$\leq \frac{2048}{s^4} \cdot \left(M_1(s) + M_2(s) \right) \cdot \exp\left(-\Theta(d) \right) \tag{34}$$

$$\leq \exp(-\Theta(d)),$$
 (35)

where (32) uses Cauchy-Schwarz inequality; (33) uses the fact $\mathbb{E}[\mathbb{1}\{E\}^2] = \mathbb{P}(E)$ valid for any event E; (34) uses (31) and the fact $\mathbb{P}(||X||_2^2 > Cd) \leq \exp(-\Theta(d))$; and finally (35) uses the condition (b) on the distribution of X stated in Theorem 2.3. Taking square roots and taking the supremum over all $||w||_2 = 1$, we obtain $\overline{\lambda}(d) \leq \exp(-\Theta(d))$; establishing Lemma 5.7.

Having established Proposition 5.3 and Lemma 5.7, we now complete the proof of Theorem 3.1.

Proof of Theorem 3.1. Throughout the proof, we assume that N is a sufficiently large polynomial in d and satisfies Assumption 1.1. Moreover, since the labels are bounded, $|Y| \leq M$ almost surely, the $o_N(1)$ terms in Theorems 2.1-2.4 disappear, as noted previously.

For the case of sigmoid and step activations, \mathcal{M} can be taken as 2. Thus, for the ξ term appearing in Proposition 5.3, we simply employ $\xi(\alpha, M, 2, A)$.

Part (a). Define the class

$$\overline{\mathcal{S}}(\delta,R) = \left\{ X \mapsto \sum_{1 \leq j \leq \overline{m}} a_j \mathrm{SGM}\left(w_j^T X\right) : (a,W) \in \mathcal{S}(\delta,R) \right\},$$

where $S(\delta, R)$ is introduced in Theorem 2.1. Note, by the definition of $S(\delta, R)$, that

$$\sup_{(a,W)\in\mathcal{S}(\delta,R)}\widehat{\mathcal{L}}\left(a,W\right)=\sup_{(a,W)\in\mathcal{S}(\delta,R)}\frac{1}{N}\sum_{1\leq i\leq N}\left(Y_{i}-\sum_{1\leq j\leq\overline{m}}a_{j}\mathrm{SGM}\left(w_{j}^{T}X_{i}\right)^{2}\right)\leq\delta^{2}.$$

Applying Theorem 2.1, we find that provided $N \ge \text{poly}(d)$, $\overline{\mathcal{S}}(\delta, R) \subset H(A)$ with probability bounded by (3), where H(A) is the class defined in Theorem 5.2 with $\sigma(\cdot) = \text{SGM}(\cdot)$ and $A = 3(1+e)(\delta+2M)$.

Finally, we (a) set $\mathcal{M}=2$ in Proposition 5.3; (b) then consider $\xi(\alpha,M,2,A)$; and (c) set $\zeta(\alpha,M,A,N)$ as in (7). Combining now Theorem 2.1 and Proposition 5.3 via a union bound, we establish the desired conclusion.

Part (b). As the output of the ReLU is not bounded, the situation is more involved. First, recall from Theorem 2.3 the sets

$$\mathcal{G}(\overline{m},\delta) \triangleq \left\{ (a,W) \in \mathbb{R}^{\overline{m}}_{\geq 0} \times \mathbb{R}^{\overline{m} \times d} : ||w_j||_2 = 1, 1 \leq j \leq \overline{m}, \widehat{\mathcal{L}}\left(a,W\right) \leq \delta^2 \right\} \quad \text{and} \quad \mathcal{G}(\delta) \triangleq \bigcup_{\overline{m} \in \mathbb{N}} \mathcal{G}(\overline{m},\delta).$$

By Theorem 2.3, it holds that with probability bounded by (4), for any $(a, W) \in \mathcal{G}(\delta)$, $||a||_1 \le 4(\delta + 2M)(\boldsymbol{\mu}^*)^{-1}$. Using the homogeneity of the ReLU activation, we instead rescale w_j by $1/\sqrt{Cd}$; and consider throughout the sets

$$\widetilde{\mathcal{G}}(\overline{m}, \delta) \triangleq \left\{ (a, W) \in \mathbb{R}^{\overline{m}}_{\geq 0} \times \mathbb{R}^{\overline{m} \times d} : ||w_j||_2 = \frac{1}{\sqrt{Cd}}, j \in [\overline{m}], \widehat{\mathcal{L}}(a, W) \leq \delta^2 \right\} \quad \text{and} \quad \widetilde{\mathcal{G}}(\delta) \triangleq \bigcup_{\overline{m} \in \mathbb{N}} \mathcal{G}(\overline{m}, \delta). \tag{36}$$

Then, with probability at least

$$1 - \left(\frac{12\sqrt{Cd}}{\mu^*}\right)^d \exp\left(-\Theta(N)\right) - N\exp\left(-\Theta(d)\right),\tag{37}$$

it holds that

$$\sup_{(a,W)\in\widetilde{G}(\delta)} ||a||_1 \le \frac{4\sqrt{Cd}(\delta + 2M)}{\mu^*}.$$
(38)

We now define an activation function, which is a "saturated" version of the ReLU:

$$S-ReLU(x) \triangleq \begin{cases} 0 & x < 0 \\ x & 0 \le x < 1 \\ 1 & x \ge 1 \end{cases}$$
 (39)

Next, using a union bound over data (X_i, Y_i) , $1 \le i \le N$,

$$\mathbb{P}\left(||X_i||_2^2 \le Cd, 1 \le i \le N\right) \ge 1 - N \exp\left(-\Theta(d)\right).$$

Hence by Cauchy-Schwarz inequality,

$$\mathbb{P}\left(\sup_{||w||_2 = \frac{1}{\sqrt{Cd}}} \left| w^T X_i \right| \le 1, 1 \le i \le N\right) \ge 1 - N \exp\left(-\Theta(d)\right).$$

Consequently, w.p. at least $1 - N \exp(-\Theta(d))$ over (X_i, Y_i) ; it holds that for all $(a, W) \in \widetilde{G}(\delta)$

$$\frac{1}{N} \sum_{1 \leq i \leq N} \left(Y_i - \sum_{1 \leq j \leq \overline{m}} a_j \text{ReLU}\left(w_j^T X_i\right) \right)^2 = \frac{1}{N} \sum_{1 \leq i \leq N} \left(Y_i - \sum_{1 \leq j \leq \overline{m}} a_j \text{S-ReLU}\left(w_j^T X_i\right) \right)^2 \leq \delta^2. \tag{40}$$

Define next the class

$$\overline{\mathcal{G}}(\delta) \triangleq \left\{ X \mapsto \sum_{1 < j < \overline{m}} a_j \text{S-ReLU}\left(w_j^T X\right) : (a, W) \in \widetilde{\mathcal{G}}(\delta) \right\}. \tag{41}$$

Note that, this set consists of all two-layer neural networks with (a) activation $S-ReLU(\cdot)$, the saturated version of $ReLU(\cdot)$; and (b) weights trained on the $ReLU(\cdot)$ network.

By Theorem 2.3 and (38), we find that provided $N \ge \text{poly}(d)$, $\overline{\mathcal{G}}(\delta) \subset H(A)$ with probability given by (37), where H(A) is the class defined in Theorem 5.2 with $\sigma(\cdot) = \text{S-ReLU}(\cdot)$ and $A = 4\sqrt{Cd}(\delta + 2M)/\mu^*$.

Observe that S-ReLU is a non-decreasing activation with bounded range. Hence, Proposition 5.3 applies: one can simply take $\mathcal{M}=2$. We now apply Proposition 5.3 with $\mathcal{M}=2$, and $A=4\sqrt{Cd}(\delta+2M)(\mu^*)^{-1}$ as in (38). Combining the probability bound (37) and the one in Proposition 5.3 by a union bound, we find that for every $\alpha>0$, with probability at least

$$1 - \zeta \left(\alpha, M, \frac{4\sqrt{Cd}(\delta + 2M)}{\mu^*}, N\right) - \left(\frac{12\sqrt{Cd}}{\mu^*}\right)^d \exp\left(-\Theta(N)\right) - N \exp\left(-\Theta(d)\right)$$
(42)

(where ξ is introduced in (7)) over training data (X_i, Y_i) , $1 \le i \le N$, it holds that

$$\sup_{\varphi \in \overline{\mathcal{G}}(\delta)} \left| \frac{1}{N} \sum_{1 \leq i \leq N} (Y_i - \varphi(X_i))^2 - \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\left(Y - \varphi(X) \right)^2 \right] \right| \leq \alpha,$$

for the class $\overline{\mathcal{G}}(\delta)$ introduced in (41). Recalling also (40) which holds with probability $1 - N \exp(-\Theta(d))$, we conclude that

$$\sup_{(a,W)\in\widetilde{\mathcal{G}}(\delta)} \mathbb{E}_{(X,Y)\sim\mathcal{D}} \left[\left(Y - \sum_{1\leq j\leq \overline{m}} a_j \operatorname{S-ReLU}\left(w_j^T X\right) \right)^2 \right] \leq \alpha + \delta^2, \tag{43}$$

with probability at least

$$1 - \zeta \left(\alpha, M, \frac{4\sqrt{Cd}(\delta + 2M)}{\mu^*}, N \right) - \left(\frac{12\sqrt{Cd}}{\mu^*} \right)^d \exp\left(-\Theta(N) \right) - 2N \exp\left(-\Theta(d) \right). \tag{44}$$

We next fix an $(a, W) \in \widetilde{\mathcal{G}}(\delta)$, and study the quantity

$$\Delta(a, W) \triangleq \left| \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\left(Y - \sum_{1 \leq j \leq \overline{m}} a_j \text{ReLU} \left(w_j^T X \right) \right)^2 \right] - \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\left(Y - \sum_{1 \leq j \leq \overline{m}} a_j \text{S-ReLU} \left(w_j^T X \right) \right)^2 \right] \right|. \tag{45}$$

This quantity is nothing but the difference of generalization errors between two networks of same architecture, same number \overline{m} of hidden units and same weights (a, W); but different activations, ReLU(\cdot) and S-ReLU(\cdot).

For convenience, denote

$$\varphi_{SR}(X) \triangleq \sum_{1 < j < \overline{m}} a_j \text{S-ReLU}\left(w_j^T X\right) \quad \text{and} \quad \varphi_R(X) \triangleq \sum_{1 < j < \overline{m}} a_j \text{ReLU}\left(w_j^T X\right).$$

In what follows, we employ the simple observation that since $a_j \geq 0$ and $0 \leq \text{S-ReLU}(x) \leq 1$, $0 \leq \varphi_{SR}(X) \leq ||a||_1$.

Suppressing the subscript $(X,Y) \sim \mathcal{D}$ from the expectations, we have

$$\Delta(a, W) = \left| \mathbb{E}\left[(Y - \varphi_{SR}(X))^2 \right] - \mathbb{E}\left[(Y - \varphi_R(X))^2 \right] \right|$$
(46)

$$= \left| \mathbb{E}\left[2Y\varphi_R(X) - 2Y\varphi_{SR}(X)\right] + \mathbb{E}\left[\varphi_{SR}(X)^2 - \varphi_R(X)^2\right] \right| \tag{47}$$

$$\leq \left| \mathbb{E}\left[2Y\varphi_R(X) - 2Y\varphi_{SR}(X) \right] \right| + \left| \mathbb{E}\left[\varphi_{SR}(X)^2 - \varphi_R(X)^2 \right] \right| \tag{48}$$

$$\leq \mathbb{E}\Big[\Big|2Y\varphi_R(X) - 2Y\varphi_{SR}(X)\Big|\Big] + \mathbb{E}\Big[\Big|\varphi_{SR}(X)^2 - \varphi_R(X)^2\Big|\Big]. \tag{49}$$

Above, (46) follows by the definition of $\Delta(a, W)$ per (45); (47) follows after simple algebra; (48) follows by the triangle inequality; and (49) follows by the Jensen's inequality.

We next study two individual terms appearing in (49) separately, while keeping in mind that $(a,W) \in \widetilde{\mathcal{G}}(\delta)$ implies $a_j \geq 0$ for $1 \leq j \leq \overline{m}$ and $||w_j||_2 = 1/\sqrt{Cd}$ for $1 \leq j \leq \overline{m}$. We have

$$\mathbb{E}\Big[\Big|2Y\varphi_{R}(X) - 2Y\varphi_{SR}(X)\Big|\Big] \leq 2M\mathbb{E}\Big[\Big|\varphi_{R}(X) - \varphi_{SR}(X)\Big|\Big] \qquad (50)$$

$$= 2M\Big(\mathbb{E}\Big[\Big|\varphi_{R}(X) - \varphi_{SR}(X)\Big|\Big|\|X\|_{2}^{2} \leq Cd\Big]\mathbb{P}\left(\|X\|_{2}^{2} \leq Cd\right)\Big) \qquad (51)$$

$$+ 2M\Big(\mathbb{E}\Big[\Big|\varphi_{R}(X) - \varphi_{SR}(X)\Big|\mathbb{1}\Big\{\|X\|_{2}^{2} > Cd\Big\}\Big] \qquad (52)$$

$$\leq 2M\mathbb{E}\Big[\Big|\varphi_{R}(X) - \varphi_{SR}(X)\Big|\mathbb{1}\Big\{\|X\|_{2}^{2} > Cd\Big\}\Big] \qquad (53)$$

$$\leq 2M\Big(\mathbb{E}\Big[\varphi_{R}(X)\mathbb{1}\Big\{\|X\|_{2}^{2} > Cd\Big\}\Big] + \mathbb{E}\Big[\varphi_{SR}(X)\mathbb{1}\Big\{\|X\|_{2}^{2} > Cd\Big\}\Big]\Big)$$

$$\leq 2Me^{-\Theta(d)}\|a\|_{1}\Big(\sqrt{\lambda(d)} + 1\Big). \qquad (55)$$

Here, (50) uses the fact $|Y| \leq M$ almost surely; (52) is by the law of total expectation; (53) uses the fact that on the event $||X||_2^2 \leq Cd$, $\varphi_R(X) = \varphi_{SR}(X)$ since $||w_j||_2 = 1/\sqrt{Cd}$; (54) uses the triangle inequality; and finally (55) uses the facts $0 \le S-ReLU(x) \le 1$ for every $x, a_i \ge 0$ for every $1 \le j \le \overline{m}$; ReLU $(x) \le |x|$; and

$$\mathbb{E}\Big[\left|w_j^TX\middle|\mathbb{1}\{||X||_2^2>Cd\}\right]\leq \sqrt{\mathbb{E}\Big[\left|w_j^TX\middle|^2\mathbb{1}\Big\{||X||_2^2>Cd\Big\}\Big]}\cdot\mathbb{E}\Big[\mathbb{1}\Big\{||X||_2^2>Cd\Big\}\Big]\leq e^{-\Theta(d)}\sqrt{\lambda(d)}$$

using Lemma 5.7 and Cauchy-Schwarz inequality. Here, $\lambda(d)$ is the function defined in (30). We now study the second term in (49). Observe that

$$\mathbb{E}\Big[\Big|\varphi_{SR}(X)^{2} - \varphi_{R}(X)^{2}\Big|\Big] = \mathbb{E}\Big[\Big|\varphi_{SR}(X)^{2} - \varphi_{R}(X)^{2}\Big|\Big|\|X\|_{2}^{2} \le Cd\Big]\mathbb{P}\left(\|X\|_{2}^{2} \le Cd\right)$$

$$+ \mathbb{E}\Big[\Big|\varphi_{SR}(X)^{2} - \varphi_{R}(X)^{2}\Big|\mathbb{1}\Big\{\|X\|_{2}^{2} > Cd\Big\}\Big]$$

$$= \mathbb{E}\Big[\Big|\varphi_{SR}(X)^{2} - \varphi_{R}(X)^{2}\Big|\mathbb{1}\Big\{\|X\|_{2}^{2} > Cd\Big\}\Big]$$

$$\leq \Big(\mathbb{E}\Big[\varphi_{SR}(X)^{2}\mathbb{1}\Big\{\|X\|_{2}^{2} > Cd\Big\}\Big] + \mathbb{E}\Big[\varphi_{R}(X)^{2}\mathbb{1}\Big\{\|X\|_{2}^{2} > Cd\Big\}\Big] \Big)$$

$$\leq \Big(e^{-\Theta(d)} \cdot \|a\|_{1}^{2} + \sum_{1 \le j \le \overline{m}} a_{j}^{2}\mathbb{E}\Big[\text{ReLU}\left(w_{j}^{T}X\right)^{2}\mathbb{1}\Big\{\|X\|_{2}^{2} > Cd\Big\}\Big] \Big)$$

$$+ 2 \sum_{1 \le j_{1} < j_{2} \le \overline{m}} a_{j_{1}} a_{j_{2}}\mathbb{E}\Big[\text{ReLU}\left(w_{j_{1}}^{T}X\right) \text{ReLU}\left(w_{j_{2}}^{T}X\right) \mathbb{1}\Big\{\|X\|_{2}^{2} > Cd\Big\}\Big] \Big)$$

$$\leq e^{-\Theta(d)}\Big(\|a\|_{1}^{2} + \lambda(d) \sum_{1 \le j \le \overline{m}} a_{j}^{2} + 2\lambda(d) \sum_{1 \le j_{1} < j_{2} \le \overline{m}} a_{j_{1}} a_{j_{2}}\Big)$$

$$= e^{-\Theta(d)}\|a\|_{1}^{2}\Big(\lambda(d) + 1\Big).$$

$$(63)$$

Indeed, (57) is again by the law of total expectation; (58) uses the fact that on $||X||_2^2 \leq Cd$, $\varphi_{SR}(X) = \varphi_R(X)$ since $||w_j||_2 = 1/\sqrt{Cd}$; (59) uses triangle inequality; (61) is obtained by

(63)

opening the parantheses while using $a_i \geq 0$, $0 \leq \text{S-ReLU}(x) \leq 1$; (62) uses the fact $a_j \geq 0$, Lemma 5.7 as well as the Cauchy-Schwarz inequality

$$\begin{split} \mathbb{E}\Big[\text{ReLU}\left(w_{j_1}^TX\right)\text{ReLU}\left(w_{j_2}^TX\right)\mathbbm{1}\Big\{||X||_2^2 > Cd\Big\}\Big] &\leq \sqrt{\mathbb{E}\Big[\text{ReLU}^2\left(w_{j_1}^TX\right)\mathbbm{1}\Big\{||X||_2^2 > Cd\Big\}\Big]} \times \\ &\sqrt{\mathbb{E}\Big[\text{ReLU}^2\left(w_{j_2}^TX\right)\mathbbm{1}\Big\{||X||_2^2 > Cd\Big\}\Big]} &\leq \lambda(d), \end{split}$$

since ReLU(x) $\leq |x|$. Finally, (63) is obtained by just noticing that for $a_j \geq 0$,

$$||a||_1^2 = \left(\sum_{1 \le j \le \overline{m}} a_j\right)^2 = \sum_{1 \le j \le \overline{m}} a_j^2 + 2 \sum_{1 \le j_1 < j_2 \le \overline{m}} a_{j_1} a_{j_2}.$$

We now combine (55) and (63) to upper bound the right hand side of (49) and arrive at

$$\Delta(a, W) \le 2Me^{-\Theta(d)} ||a||_1 \left(\sqrt{\lambda(d)} + 1 \right) + e^{-\Theta(d)} ||a||_1^2 \left(\lambda(d) + 1 \right).$$

Since $||a||_1 \leq 4\sqrt{Cd}(\delta + 2M)/\mu^*$ on $\widetilde{\mathcal{G}}(\delta)$ as recorded in (38), we obtain

$$\sup_{(a,W)\in\widetilde{\mathcal{G}}(\delta)} \Delta(a,W) \leq e^{-\Theta(d)} \left(\frac{8M\sqrt{C}(\delta+2M)}{\boldsymbol{\mu}^*} \sqrt{d} \left(\sqrt{\lambda(d)} + 1 \right) + \frac{16C(\delta+2M)^2}{\boldsymbol{\mu}^{*2}} d\left(\lambda(d) + 1 \right) \right). \tag{64}$$

Recall that $\lambda(d) \leq \exp(-\Theta(d))$ by (30). Note that as long as $M, C, \delta, \mu^* = \exp(o(d))$ as well, the term on the right hand side of (64) is $e^{-\Theta(d)}$.

We finally combine (43), (45); and (64) to obtain

$$\sup_{(a,W) \in \mathcal{G}(\delta)} \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\left(Y - \sum_{1 \leq j \leq \overline{m}} a_j \text{ReLU} \left(w_j^T X \right) \right)^2 \right] \leq \alpha + \delta^2 + e^{-\Theta(d)}$$

with probability at least

$$1 - \zeta\left(\alpha, M, \frac{4\sqrt{Cd}(\delta + 2M)}{\boldsymbol{\mu^*}}, N\right) - \left(\frac{12\sqrt{Cd}}{\boldsymbol{\mu^*}}\right)^d \exp\left(-\Theta(N)\right) - 2N\exp\left(-\Theta(d)\right),$$

as shown in (44). This concludes the proof of Part (b).

Part (c). This is quite similar to Part (a). Define the class

$$\overline{\mathcal{H}}(\delta) \triangleq \left\{ X \mapsto \sum_{1 < j < \overline{m}} a_j \mathtt{Step}\left(w_j^T X\right) : (a, W) \in \mathcal{H}(\delta) \right\}$$

where $\mathcal{H}(\delta)$ is introduced in Theorem 2.4. Note, by definition, that

$$\sup_{(a,W)\in\mathcal{H}(\delta)}\widehat{\mathcal{L}}\left(a,W\right) = \sup_{(a,W)\in\mathcal{H}(\delta)}\frac{1}{N}\sum_{1\leq i\leq N}\left(Y_i - \sum_{1\leq j<\overline{m}}a_j\mathrm{Step}\left(w_j^TX_i\right)^2\right) \leq \delta^2.$$

Applying Theorem 2.4, we find that provided $N \ge \text{poly}(d)$, $\overline{\mathcal{H}}(\delta) \subset H(A)$ w.h.p., where H(A) is the class defined in Theorem 5.2 with $\sigma(\cdot) = \text{Step}(\cdot)$ and $A = 2(\delta + 2M)/\eta$.

Like in the previous case, we then (a) set $\mathcal{M}=2$ in Proposition 5.3; (b) then let $\xi(\alpha,M,2)$ to be $\xi(\alpha,M,2,A)$; and (c) set $\zeta(\alpha,M,A,N)$ as in (7). Combining now Theorem 2.4 and Proposition 5.3 via a union bound, we establish the desired conclusion.

References

- [ABDCBH97] Noga Alon, Shai Ben-David, Nicolo Cesa-Bianchi, and David Haussler, Scale-sensitive dimensions, uniform convergence, and learnability, Journal of the ACM (JACM) 44 (1997), no. 4, 615–631.
- [ADH⁺19a] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang, On exact computation with an infinitely wide neural net, Advances in Neural Information Processing Systems, 2019, pp. 8139–8148.
- [ADH⁺19b] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang, Finegrained analysis of optimization and generalization for overparameterized twolayer neural networks, arXiv preprint arXiv:1901.08584 (2019).
- [AGKM16] Sanjeev Arora, Rong Ge, Ravi Kannan, and Ankur Moitra, Computing a non-negative matrix factorization—provably, SIAM Journal on Computing 45 (2016), no. 4, 1582–1611.
- [AGNZ18] Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang, Stronger generalization bounds for deep nets via a compression approach, arXiv preprint arXiv:1802.05296 (2018).
- [Bar98] Peter L Bartlett, The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network, IEEE transactions on Information Theory 44 (1998), no. 2, 525–536.
- [BFT17] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky, Spectrally-normalized margin bounds for neural networks, Advances in Neural Information Processing Systems, 2017, pp. 6240–6249.
- [BG17] Alon Brutzkus and Amir Globerson, Globally optimal gradient descent for a convnet with gaussian inputs, Proceedings of the 34th International Conference on Machine Learning-Volume 70, JMLR. org, 2017, pp. 605–614.
- [BGMSS17] Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz, Sgd learns over-parameterized networks that provably generalize on linearly separable data, arXiv preprint arXiv:1710.10174 (2017).

- [BHLM19] Peter L Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian, Nearlytight vc-dimension and pseudodimension bounds for piecewise linear neural networks., Journal of Machine Learning Research 20 (2019), no. 63, 1–17.
- [BHMM19] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal, Reconciling modern machine-learning practice and the classical bias-variance trade-off, Proceedings of the National Academy of Sciences 116 (2019), no. 32, 15849–15854.
- [BLW96] Peter L Bartlett, Philip M Long, and Robert C Williamson, Fat-shattering and the learnability of real-valued functions, journal of computer and system sciences 52 (1996), no. 3, 434–452.
- [CG19] Yuan Cao and Quanquan Gu, Generalization bounds of stochastic gradient descent for wide and deep neural networks, Advances in Neural Information Processing Systems, 2019, pp. 10836–10846.
- [CW08] Ronan Collobert and Jason Weston, A unified architecture for natural language processing: Deep neural networks with multitask learning, Proceedings of the 25th international conference on Machine learning, ACM, 2008, pp. 160–167.
- [DKKZ20] Ilias Diakonikolas, Daniel M Kane, Vasilis Kontonis, and Nikos Zarifis, Algorithms and sq lower bounds for pac learning one-hidden-layer relu networks, Conference on Learning Theory, PMLR, 2020, pp. 1514–1539.
- [DL18] Simon Du and Jason Lee, On the power of over-parametrization in neural networks with quadratic activation, International Conference on Machine Learning, PMLR, 2018, pp. 1329–1338.
- [DLL⁺19] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai, *Gradient descent finds global minima of deep neural networks*, International Conference on Machine Learning, PMLR, 2019, pp. 1675–1685.
- [DR17] Gintare Karolina Dziugaite and Daniel M Roy, Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data, arXiv preprint arXiv:1703.11008 (2017).
- [Dur19] Rick Durrett, *Probability: theory and examples*, vol. 49, Cambridge university press, 2019.
- [DZPS18] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh, *Gradient descent provably optimizes over-parameterized neural networks*, arXiv preprint arXiv:1810.02054 (2018).
- [EGKZ20] Matt Emschwiller, David Gamarnik, Eren C Kızıldağ, and Ilias Zadik, Neural networks and polynomial regression. demystifying the overparametrization phenomena, arXiv preprint arXiv:2003.10523 (2020).
- [FCG19] Spencer Frei, Yuan Cao, and Quanquan Gu, Algorithm-dependent generalization bounds for overparameterized deep residual networks, Advances in Neural Information Processing Systems, 2019, pp. 14797–14807.

- [Gab19] Marylou Gabrié, Towards an understanding of neural networks: mean-field incursions, Ph.D. thesis, Paris Sciences et Lettres, 2019.
- [GAS⁺19] 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*, Advances in Neural Information Processing Systems, 2019, pp. 6979–6989.
- [GBB11] Xavier Glorot, Antoine Bordes, and Yoshua Bengio, *Deep sparse rectifier neu*ral networks, Proceedings of the fourteenth international conference on artificial intelligence and statistics, 2011, pp. 315–323.
- [GK19] Surbhi Goel and Adam R Klivans, Learning neural networks with two nonlinear layers in polynomial time, Conference on Learning Theory, PMLR, 2019, pp. 1470–1499.
- [GKM18] Surbhi Goel, Adam Klivans, and Raghu Meka, Learning one convolutional layer with overlapping patches, International Conference on Machine Learning, PMLR, 2018, pp. 1783–1791.
- [GLM17] Rong Ge, Jason D Lee, and Tengyu Ma, Learning one-hidden-layer neural networks with landscape design, arXiv preprint arXiv:1711.00501 (2017).
- [GLSS18] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro, *Implicit bias of gradient descent on linear convolutional networks*, Advances in Neural Information Processing Systems, 2018, pp. 9461–9471.
- [GRS17] Noah Golowich, Alexander Rakhlin, and Ohad Shamir, Size-independent sample complexity of neural networks, arXiv preprint arXiv:1712.06541 (2017).
- [GWZ19] Rong Ge, Runzhe Wang, and Haoyu Zhao, Mildly overparametrized neural nets can memorize training data efficiently, arXiv preprint arXiv:1909.11837 (2019).
- [Hau92] David Haussler, Decision theoretic generalizations of the pac model for neural net and other learning applications, Information and computation **100** (1992), no. 1, 78–150.
- [HLM17] Nick Harvey, Christopher Liaw, and Abbas Mehrabian, Nearly-tight vc-dimension bounds for piecewise linear neural networks, Conference on Learning Theory, 2017, pp. 1064–1068.
- [Hoy02] Patrik O Hoyer, *Non-negative sparse coding*, Proceedings of the 12th IEEE Workshop on Neural Networks for Signal Processing, IEEE, 2002, pp. 557–565.
- [HRS16] Moritz Hardt, Ben Recht, and Yoram Singer, Train faster, generalize better: Stability of stochastic gradient descent, International Conference on Machine Learning, PMLR, 2016, pp. 1225–1234.

- [HZRS16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, *Deep residual learning for image recognition*, Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [KS94] Michael J Kearns and Robert E Schapire, Efficient distribution-free learning of probabilistic concepts, Journal of Computer and System Sciences 48 (1994), no. 3, 464–497.
- [KSH12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, *Imagenet classification* with deep convolutional neural networks, Advances in neural information processing systems, 2012, pp. 1097–1105.
- [LL18] Yuanzhi Li and Yingyu Liang, Learning overparameterized neural networks via stochastic gradient descent on structured data, Advances in Neural Information Processing Systems, 2018, pp. 8157–8166.
- [LMZ20] Yuanzhi Li, Tengyu Ma, and Hongyang R Zhang, Learning over-parametrized two-layer neural networks beyond ntk, Conference on Learning Theory, PMLR, 2020, pp. 2613–2682.
- [LPRS17] Tengyuan Liang, Tomaso Poggio, Alexander Rakhlin, and James Stokes, Fisher-rao metric, geometry, and complexity of neural networks, arXiv preprint arXiv:1711.01530 (2017).
- [LS99] Daniel D Lee and H Sebastian Seung, Learning the parts of objects by non-negative matrix factorization, Nature **401** (1999), no. 6755, 788–791.
- [LS01] Daniel Lee and H. Sebastian Seung, Algorithms for non-negative matrix factorization, Advances in Neural Information Processing Systems (T. Leen, T. Dietterich, and V. Tresp, eds.), vol. 13, MIT Press, 2001.
- [LZA20] Zhiyuan Li, Yi Zhang, and Sanjeev Arora, Why are convolutional nets more sample-efficient than fully-connected nets?, arXiv preprint arXiv:2010.08515 (2020).
- [MDH11] Abdel-rahman Mohamed, George E Dahl, and Geoffrey Hinton, *Acoustic modeling using deep belief networks*, IEEE transactions on audio, speech, and language processing **20** (2011), no. 1, 14–22.
- [MWZZ18] Wenlong Mou, Liwei Wang, Xiyu Zhai, and Kai Zheng, Generalization bounds of sgld for non-convex learning: Two theoretical viewpoints, Conference on Learning Theory, PMLR, 2018, pp. 605–638.
- [NBMS17] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro, Exploring generalization in deep learning, Advances in Neural Information Processing Systems, 2017, pp. 5947–5956.
- [NBS17] Behnam Neyshabur, Srinadh Bhojanapalli, and Nathan Srebro, A pac-bayesian approach to spectrally-normalized margin bounds for neural networks, arXiv preprint arXiv:1707.09564 (2017).

- [NTS15] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro, *Norm-based capacity control in neural networks*, Conference on Learning Theory, 2015, pp. 1376–1401.
- [SS18] Itay Safran and Ohad Shamir, Spurious local minima are common in two-layer relu neural networks, International Conference on Machine Learning, PMLR, 2018, pp. 4433–4441.
- [SSS⁺17] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al., *Mastering the game of go without human knowledge*, Nature **550** (2017), no. 7676, 354.
- [SV17] Paris Smaragdis and Shrikant Venkataramani, A neural network alternative to non-negative audio models, 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2017, pp. 86–90.
- [Ver10] Roman Vershynin, Introduction to the non-asymptotic analysis of random matrices, arXiv preprint arXiv:1011.3027 (2010).
- [Ver18] _____, High-dimensional probability: An introduction with applications in data science, vol. 47, Cambridge university press, 2018.
- [Wik] Non-negative matrix factorization, https://en.wikipedia.org/wiki/Non-negative_matrix_Accessed: 2021-03-01.
- [ZBH⁺16] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, Understanding deep learning requires rethinking generalization, arXiv preprint arXiv:1611.03530 (2016).
- [ZYWG19] Xiao Zhang, Yaodong Yu, Lingxiao Wang, and Quanquan Gu, Learning one-hidden-layer relu networks via gradient descent, The 22nd International Conference on Artificial Intelligence and Statistics, PMLR, 2019, pp. 1524–1534.