

How Many Neurons Does it Take to Approximate the Maximum?

Itay Safran¹, Daniel Reichman², and Paul Valiant¹

¹Purdue University

²Worcester Polytechnic Institute

Abstract

We study the size of a neural network needed to approximate the maximum function over d inputs, in the most basic setting of approximating with respect to the L_2 norm, for continuous distributions, for a network that uses ReLU activations. We provide new lower and upper bounds on the width required for approximation across various depths. Our results establish new depth separations between depth 2 and 3, and depth 3 and 5 networks, as well as providing a depth $\mathcal{O}(\log(\log(d)))$ and width $\mathcal{O}(d)$ construction which approximates the maximum function. Our depth separation results are facilitated by a new lower bound for depth 2 networks approximating the maximum function over the uniform distribution, assuming an exponential upper bound on the size of the weights. Furthermore, we are able to use this depth 2 lower bound to provide tight bounds on the number of neurons needed to approximate the maximum by a depth 3 network. Our lower bounds are of potentially broad interest as they apply to the widely studied and used *max* function, in contrast to many previous results that base their bounds on specially constructed or pathological functions and distributions.

1 Introduction

How and why *depth* helps neural networks to excel in applications is one of the central challenges in the quest to understand deep learning. Both classical circuit complexity and modern deep learning theory is guided by the intuition that a modest increase in depth often leads to drastic—and often exponential—improvements in the expressive power of the circuit or network described, along with concomitant improvements in key measures of performance, including efficiency as measured by width or neuron count, and approximation accuracy. Despite this firm intuition, and much recent encouraging evidence of the *practical* expressive power of deep networks (and hence also deep circuits), theoretical insight to illuminate these phenomena remains scarce; and each additional layer of depth adds often prohibitive new theoretical challenges. This paper contributes to this area by providing several new lower and upper bounds across a range of depths for a *natural*

function, with respect to a meaningful notion of approximation, in a realistic circuit/network model with realistic parameters.

While it is known that any continuous function can be approximated arbitrarily well by a shallow (depth 2) network [2, 14], these constructions of depth 2 approximations typically require exponential size, as a function of the input dimension in the worst case. This raises the natural question: when and how can increasing the depth beyond 2 allow us to drastically improve the width, so as to be polynomial in the input dimension. Indeed, significant study has been directed towards understanding the benefit of depth in expressing functions. For example, there are several constructions of functions which can be approximated by a depth 3 and polynomial width network, but require exponentially many neurons to be approximated by a network of depth 2 [4, 3, 20, 21, 25, 8, 19]. It is also known that highly oscillatory or even sufficiently non-linear functions may require exponentially fewer neurons when approximated by networks whose depth scales with the problem’s parameters (e.g. the target accuracy) than by constant depth networks [23, 26, 15, 20]. While such results establish the expressive superiority of depth over width, the functions used to demonstrate this are at times somewhat pathological, the configuration of the network which approximates them efficiently is highly stylized, or the assumed data distribution is quite complex. Therefore we are motivated to study the effect of depth on efficiency for “natural”, well-studied functions that commonly arise in machine learning tasks.

We study the effect of depth on the quality of approximation for the maximum function

$$f_d(x_1, \dots, x_d) := \max\{x_1, \dots, x_d\},$$

where $(x_1, \dots, x_d) \in \mathbb{R}^d$. This function enjoys many favorable properties: It is a fundamental mathematical function; it has a very simple structure; it is easy to compute in linear time (assuming comparisons between real numbers can be done in time $\mathcal{O}(1)$), and it is used explicitly in popular neural network architectures (e.g. max pooling [12, 22, 9]). Computing the maximum is important in reinforcement learning tasks (choosing an action maximizing the expected reward) and has received attention in theoretical neuroscience [11]. Several works have studied how to compute the maximum efficiently using a neural network [1, 7, 17, 5]. However, most known results only deal with L_∞ approximations or even exact computation of this function, which is a far more stringent notion of approximation than the L_2 approximation which is often the metric chosen for machine learning applications, and is the metric we study in this paper (see related work subsection and Sec. 2 for further discussion).

We provide new lower and upper bounds for approximating the maximum function with respect to a broad class of distributions including the uniform distribution over the unit hypercube, and the Gaussian distribution (see Assumption 3.2). We show that for any natural $k \leq \mathcal{O}(\log(\log(d)))$,¹ ReLU networks of depth $2k+1$ and width $\mathcal{O}\left(d^{1+1/(2^k-1)}\right)$ can approximate the maximum function to arbitrary accuracy given sufficiently large weights. In particular, this implies a depth 3 and width $\mathcal{O}(d^2)$ approximation; a depth 5 and width $\mathcal{O}(d^{4/3})$ approximation; and a depth $\mathcal{O}(\log(\log(d)))$ and width $\mathcal{O}(d)$ approximation. Moreover, assuming an exponential upper bound on the size of the weights, we show corresponding lower bounds for approximating the maximum using ReLU

¹Unless stated otherwise all logarithms are base 2.

Table 1: New results in this paper for approximating the maximum function using ReLU neural networks uniformly over the domain $[0, 1]^d$. We provide a polynomial separation between depth 2 and 3, and for the same target function, we provide a polynomial separation between depth 3 and 5, requiring widths of $\Omega(d^2)$ and $\mathcal{O}(d^{4/3})$, respectively. We also provide a novel upper bound which only requires a perhaps surprising network depth of $\mathcal{O}(\log(\log(d)))$ for approximating the maximum using linear width. This in contrast to the previously best known approximation which requires width and overall size linear in d but depth $\mathcal{O}(\log(d))$ to compute the maximum exactly [1, 17]. The theta notation in the target accuracy column hides an absolute constant, and our depth 2 lower bound holds for any value of the parameter $\ell \geq 1$. An asterisk (*) denotes that, beyond the ReLU activation, the bound applies to any activation satisfying a mild assumption (Assumption 3.2); and a dagger ([†]) denotes that an exponential upper bound on the magnitude of the weights is assumed.

Depth	Target accuracy	Width lower bound	Width upper bound (any accuracy)
2	$d^{-\theta(\ell)}$, any $\ell \geq 1$	$\Omega(d^\ell)$ (Thm. 4.2) (*,†)	
3	$d^{-\theta(1)}$	$\Omega(d^2)$ (Thm. 4.3) ([†])	$\mathcal{O}(d^2)$ (Thm. 3.3)
5	$d^{-4.5}$	d (Thm. 4.4) (*)	$\mathcal{O}(d^{4/3})$ (Thm. 3.4)
$\mathcal{O}(\log(\log(d)))$	$d^{-4.5}$	d (Thm. 4.4) (*)	$\mathcal{O}(d)$ (Thm. 3.4)

networks. Specifically, we show a depth 2 lower bound requiring width d^ℓ for any $\ell \geq 1$ for obtaining a target accuracy of $d^{-c\ell}$ for constant $c > 0$, and a depth 3 lower bound requiring width $\Omega(d^2)$, establishing the tightness of our depth 3 upper bound.

We remark that by performing a change of variables in our lower bounds, one can rescale the domain of approximation as desired at the cost of rescaling the target accuracy (see Lemma D.6 for a formal statement). For example, rescaling the domain to $[0, d]$ where the maximum function has constant variance will multiply the target accuracy by a factor of d^2 in all of our lower bounds. Moreover, by further rescaling, our lower bounds show that the maximum function cannot be approximated to better than constant accuracy if the domain of approximation scales polynomially with the input dimension d . Due to the fact that our upper bounds on the required width are independent of the target accuracy (better accuracy is obtained by increasing the magnitude of the weights of the approximating network), scaling the domain of approximation does not affect the width requirement in our upper bounds. Therefore, our results also establish several new depth-based separations for approximating the maximum function to better than constant accuracy. See Table 1 for a comparison of our bounds.

It is interesting to compare lower bounds for continuous neural networks over continuous domains to lower bounds for discrete neural networks such as threshold circuits over $\{0, 1\}^d$. Devising superlinear lower bounds for depth three threshold circuits (with polynomial weights on the output gate) has been obtained relatively recently after decades of research for a family of complicated functions [10] which cannot be computed by circuits with $o(d^{3/2}/\log^3(d))$ threshold gates. Our tight quadratic lower bound for depth 3 networks with ReLU gates approximating the maxi-

mum function suggests that proving lower bounds for the continuous case is a more amenable task and that further superlinear lower bounds for bounded depth networks might be achievable over the continuous domain.

The remainder of this paper is structured as follows: After presenting our contributions in this paper in more detail below, we discuss related work in the literature. In Sec. 2 we present the notation used throughout this paper. Sec. 3 details our positive approximation results (upper bounds) and Sec. 4 details our negative inapproximability results (lower bounds). Lastly, Sec. 5 summarizes and discusses potential future work directions.

Our contributions

- We exhibit a construction to approximate the maximum function arbitrarily well—in the L_2 sense—using a depth 3 and width $d(d + 1)$ ReLU network (Thm. 3.3). Interestingly, to increase the accuracy of this construction, we increase the size of the weights, but do not need to change the network architecture. This construction arises from a piecewise-linear decomposition of the maximum function.
- We “compose” the above construction with itself, so as to enable different depth-width trade-offs. This reinterpretation of Thm. 3.3 provides new upper bounds for expressing the maximum function across odd depths, with the required width dropping rapidly as larger depths are used (Thm. 3.4).
- By contrast, we show that these constructions at depths 3 and higher are impossible at depth 2: the width of a depth 2 network approximating the maximum function *must* depend on the desired approximation accuracy. This thus shows a polynomial separation between depths 2 and 3 for approximating the maximum function. Our analysis relates to the seminal technique developed in Eldan and Shamir [4], analyzing the Fourier spectrum of the maximum function, to show that, assuming exponential upper bounds on the size of the weights, we show a lower bound for approximating the maximum function on a compact domain (Thm. 4.2).
- We show a polynomial separation between depth 3 and depth 5 neural networks. Using a combinatorial argument, we show that depth 3 ReLU networks with first hidden layer of width at most $d^2/5$ cannot capture the full structure of the maximum function on the hypercube, reducing the approximation error lower bound to the accuracy achieved by the second hidden layer. Using our previous lower bound in Thm. 4.2, this implies an approximation lower bound for depth 3 ReLU networks (Thm. 4.3). Together with Thm. 3.3, this establishes a tight bound of $\Theta(d^2)$ for approximation of the maximum using depth 3; and when combined with our Thm. 3.4 which implies a depth 5 and width $\mathcal{O}(d^{4/3})$ approximation, this provides a polynomial separation between depth 3 and 5.
- Lastly, we observe that any neural network (regardless of depth or activation function) approximating the maximum must have at least d neurons in its first hidden layer. Thus, known

upper bounds on the number of ReLUs needed to compute the maximum precisely which require size $\mathcal{O}(d)$ (e.g. Arora et al. [1], Matoba et al. [17]) are optimal up to a constant factor.

Related work

Exact computation and approximation of the maximum function Quite a few recent works have studied the problem of exact computation of the maximum function using ReLU neural networks. Arora et al. [1] establish that any d -dimensional, piecewise-linear function can be expressed exactly using a depth $\lceil \log(d+1) \rceil$ ReLU network. Observe that we are able to obtain a more depth efficient implementation of the maximum function, although in contrast to Arora et al. [1] we construct a network *approximating* the function rather than computing it exactly, and our overall size requirement is slightly larger. The construction of Arora et al. [1] implies that $\max\{x_1, x_2\}$ is computable by depth 2 ReLU networks, and $\max\{x_1, \dots, x_4\}$ is computable by depth 3 ReLU networks. In contrast, Mukherjee and Basu [18] show that the function $(x_1, x_2) \mapsto \max\{0, x_1, x_2\}$ (which can be shown to be equivalent to computing the maximum over three inputs) cannot be computed exactly by a depth 2 ReLU network regardless of its width. Hertrich et al. [7] conjecture an analogous impossibility result for computing $\max\{0, x_1, \dots, x_4\}$ exactly using a depth 3 network, and partially resolve it by assuming a certain restriction on the structure of the computing network. Haase et al. [5] further show the uncomputability of $\max\{0, x_1, \dots, x_4\}$ assuming the computing network has depth 3 and integral weights, which also implies that the exact computation of $f_d(\cdot)$ using integer weights requires depth $\Omega(d)$. Since our approximation of $f_d(\cdot)$ can be done using integer weights, this shows a depth separation between exact and approximate computation of $f_d(\cdot)$ using integer weights, as the former requires depth $\Omega(\log(d))$ whereas for the latter depth $\mathcal{O}(\log(\log(d)))$ suffices.

Despite these efforts, this conjecture is still open. Furthermore, it is currently unknown whether there exists *any* piecewise-linear function which cannot be computed exactly using depth 3 ReLU networks. We stress that the aforementioned lower bounds are concerned with *exact* computation, whereas our notion of approximation is markedly different since we consider lower bounds with respect to the L_2 norm rather than requiring zero L_∞ loss. This is a less stringent approximation requirement in the sense that an L_2 lower bound implies an L_∞ lower bound, but not vice versa.

In contrast to the exact computation requirement discussed above, Matoba et al. [17] consider approximations of the maximum function with respect to the L_∞ norm. They show a family of networks that increase in accuracy as the size of the approximating network increases, whereas we provide a construction of a network achieving arbitrarily good accuracy with fixed width (by increasing the size of the weights). Additionally, under the restriction of the approximating network to have a certain symmetry structure, they also show approximation lower bounds for the maximum function, however these do not hold in general if we relax the symmetry assumption, whereas our lower bounds hold only under the (mild) assumption that the approximating network has exponentially bounded weights.

Separations between depth 2 and 3 The seminal work of Eldan and Shamir [4] was the first to establish the existence of a (continuous) function that can be approximated efficiently using

networks of depth 3, whereas any network of depth 2 would require width exponential in the input dimension to achieve better than constant accuracy. Later, Daniely [3] showed a separation using a different technique which applies to a compactly supported distribution, but requires an exponential upper bound on the magnitude of the weights of the approximating depth 2 network. Following these works, additional separation results between depth 2 and 3 were shown (e.g. [25, 8]), including reductions to the results of Eldan and Shamir [4] and Daniely [3] that however hold for much simpler functions than the ones originally used (e.g. [20, 21, 19]). Nevertheless, due to the reduction proof technique used, these results inevitably inherit the arguably more complicated distributions over the data used in Eldan and Shamir [4] and Daniely [3]. In contrast, our separation between depth 2 and 3 holds for both the simple maximum function and for the uniform distribution over a hypercube, albeit providing a polynomial separation and requiring an exponential upper bound on the magnitude of the weights.

Limitations of deeper architectures Moving beyond depth 2, there are known constructions of functions that can be approximated by a small sized network (with no restriction on its depth), whereas an approximation to similar accuracy using constant depth networks may require exponentially many more neurons. Such lower bounds, however, are based on two main arguments and suffer from certain drawbacks making them incomparable to our results. We discuss these lower bounds and their limitations in more detail below.

Region-counting-based depth separations The seminal work of Telgarsky [23] first established depth separations between deep architectures. It is shown that a deep ReLU network can realize a one-dimensional, rapidly oscillating sawtooth function, whereas a shallower architecture cannot generate sufficiently many linear segments to be able to approximate this function efficiently. If one wishes, however, to learn this efficient representation of this function using the deeper architecture, then it is known that this cannot be done efficiently using standard techniques such as the gradient descent algorithm [16]. Different lower bounds exist that build on this region counting proof technique, but focus on smooth and non-linear target functions that may be more prone to be learned efficiently using gradient descent [26, 15, 20]. Nevertheless, there's some theoretical work which shows that when initializing deep ReLU networks, the expected number of linear regions our initialization will have is merely linear in the size of the network, and further empirical evidence suggests that this number does not tend to increase significantly, indicating that depth will impart no practical benefit for approximating these target functions [6, 24]. On the other hand, our results which focus on the maximum function and do not rely on region counting, still leave open the possibility of an optimization-based result to be shown which will demonstrate this separation in a more practical setting.

Size lower bounds and connections to circuit complexity A different, less direct approach for showing approximation lower bounds for neural networks relies on the connection between threshold circuits and neural networks. Mukherjee and Basu [18] derive sub-linear size lower bounds for neural networks by showing reductions to known threshold circuit lower bounds. Vardi et al. [24] use communication complexity to provide linear size lower bounds for approximating

a smoothed version of the binary IP mod 2 function. These results provide a different result than ours since they imply a linear width lower bound for approximating various functions using depth 3 networks, while we show a quadratic lower bound for depth 3 ReLU networks. Moreover, the lower bounds in the aforementioned papers are for the size of the network, rather than the required width for some given depth. For this reason such results cannot establish the superiority of depth over width, since these two quantities can be traded off evenly in such lower bounds, whereas our results establish a non-symmetric trade-off which indicates that depth is more efficient for approximating the maximum function.

2 Preliminaries and notation

Notation and terminology We let $[n]$ be shorthand for the set $\{1, \dots, n\}$. We denote vectors using bold-faced letters (e.g. \mathbf{x}) and matrices or random variables using upper-case letters (e.g. X). Multivariate random variables are denoted using bold-faced upper-case letters (e.g. \mathbf{X}). Given a vector $\mathbf{x} = (x_1, \dots, x_d)$, we let $\|\mathbf{x}\|_p$ denote its ℓ_p norm which is given by $\left(\sum_{i=1}^d |x_i|^p\right)^{1/p}$, where the case $p = \infty$ is defined as $\|\mathbf{x}\|_\infty = \max_{i \in [d]} |x_i|$. Throughout, we use the notation $f_d(\mathbf{x}) := \max\{x_1, \dots, x_d\}$ for the maximum function, $[x]_+ = \max\{0, x\}$ for the ReLU activation function, and for any natural $k \geq 1$ we use the notation $\beta(k) := \frac{1}{2^k - 1}$. A function $f : D \rightarrow \mathbb{R}$ defined in some domain $D \subseteq \mathbb{R}^d$ is (*continuous*) *piecewise-linear* if there exists a finite partition $D = \cup_i D_i$ such that f is linear on D_i for all i , where each D_i is a closed set which is referred to as a *linear region* of f . We let $\mathcal{U}(A)$ denote the uniform distribution on a set $A \subseteq \mathbb{R}^d$.

Neural networks We consider fully connected, feed-forward neural networks, computing functions from \mathbb{R}^d to \mathbb{R} . A σ neural network consists of layers of neurons. In every layer except for the output neuron, an affine function of the inputs is computed, followed by a computation of the non-linear activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$. The single output neuron simply computes an affine transformation of its inputs. Each layer with a non-linear activation is called a *hidden layer*, and the *depth* of a network is defined as the number of hidden layers plus one. The *width* of a network is defined as the number of neurons in the largest hidden layer which we generally denote by k , and the *size* of the network is the total number of neurons across all layers.

Approximation error Since we consider a regression setting in which a neural network $\mathcal{N} : \mathbb{R}^d \rightarrow \mathbb{R}$ computes a real function of its input, we define our approximation error with respect to an underlying distribution \mathcal{D} on \mathbb{R}^d and we consider the square loss. Formally, given a predictor \mathcal{N} , a target function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and an underlying distribution \mathcal{D} , our approximation error is defined as

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [(\mathcal{N}(\mathbf{x}) - f(\mathbf{x}))^2].$$

In words, the L_2 approximation defined above corresponds to the expected square error when sampling an instance from the underlying distribution \mathcal{D} , labelling it using the target function f we are trying to approximate, and making a prediction using a given neural network hypothesis

\mathcal{N} . This makes this notion of approximation a natural choice for showing approximation lower and upper bounds, as a lower bound for certain class of neural network predictors implies the existence of a particular learning problem where the architecture being considered is unable to achieve population loss better than a certain quantity, whereas a network \mathcal{N} which achieves small loss implies that this class of networks can express a good predictor.

3 Deep ReLU approximations

In this section, we focus on positive approximation results for the maximum function. We now begin with stating our assumptions on the underlying distribution generating the data, but first we would need the following definition.

Definition 3.1. *Given some $\delta > 0$, we say that a vector $\mathbf{x} = (x_1, \dots, x_d)$ is δ -separated if for all $i \neq j$, and $x_j \neq 0$ we have that*

$$\frac{x_i}{x_j} \notin [1 - \delta, 1 + \delta].$$

We denote the set of δ -separated vectors in d -dimensional space by

$$S_\delta := \{ \mathbf{x} \in \mathbb{R}^d : \mathbf{x} \text{ is } \delta\text{-separated} \}.$$

The above definition essentially guarantees that each pair of coordinates in \mathbf{x} have a ratio whose distance from 1 is at least δ for some real $\delta > 0$. Since our construction is sensitive to instances where there are coordinates that are extremely close, we would need to make sure that points that violate δ -separateness are sufficiently scarce. To this end, we make the following assumption on the distribution of the data.

Assumption 3.2. *The distribution \mathcal{D} satisfies the following:*

1.

$$\mathbb{E}_{\mathbf{X} \sim \mathcal{D}} [\|\mathbf{X}\|_\infty^2] < \infty.$$

2.

$$\lim_{\delta \rightarrow 0} \mathbb{P}_{\mathbf{X} \sim \mathcal{D}} [\mathbf{X} \notin S_\delta] = 0.$$

Item 1 merely requires that the tail of \mathbf{X} is sufficiently well-behaved in the sense of having a finite second moment for its infinity norm, and Item 2 requires that it becomes increasingly unlikely to draw an instance from \mathcal{D} which isn't δ -separated as δ decreases. These hold, for example, when the coordinates of \mathbf{X} are i.i.d. and follow any absolutely continuous distribution with a bounded density and a finite second moment, or when a certain continuous noise is added to a sufficiently concentrated random variable \mathbf{X} (see Appendix A for formal examples).

We now turn to formally state our positive approximation result for approximating the maximum function using depth 3 ReLU neural networks.

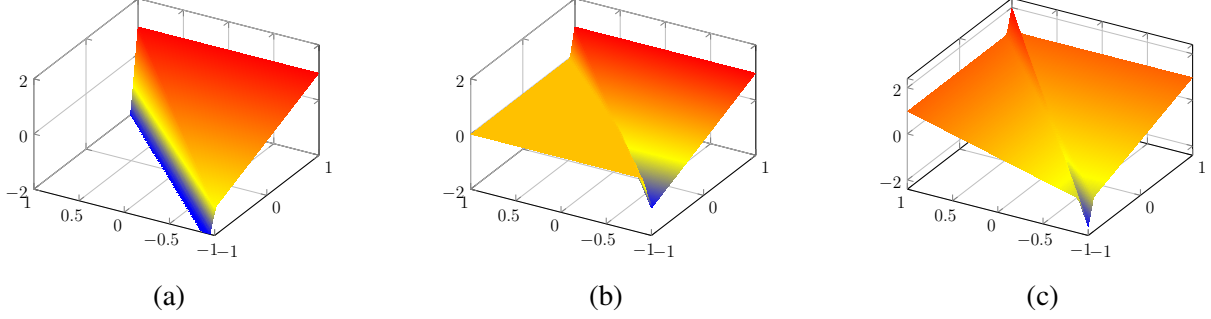


Figure 1: Three-step polytope approximation of $(x, y) \mapsto \max\{x, y\}$. Subfigure 1(a) plots the depth 2 network $[x]_+ - [10y - 10x]_+ - [-x]_+ - [10y - 10x]_+$ which computes $\max\{x, y\}$ on the polytope $\{(x, y) \in [-1, 1]^2 : x \geq y\}$. In Subfigure 1(b), a second layer of ReLUs is utilized to clip negative values that are outside of the polytope to zero, plotting the depth 3 network $\mathcal{N}(x, y) := [[x]_+ - [10y - 10x]_+]_+ - [[-x]_+ - [10y - 10x]_+]_+$. Lastly, Subfigure 1(c) plots the depth 3 network $\mathcal{N}(x, y) + \mathcal{N}(y, x)$ which is an effective approximation of $\max\{x, y\}$. We remark that while $\max\{x, y\}$ can be computed exactly using depth 2 ReLU networks, the figure is intended for illustration purposes of our construction idea used in Thm. 3.3, which generalizes to any input dimension d . Best viewed in color.

Theorem 3.3. *Let \mathcal{D} be any distribution satisfying Assumption 3.2. Then for any $\varepsilon > 0$ and natural $d \geq 2$, there exists a ReLU network \mathcal{N} of depth 3 and width $d(d + 1)$ such that*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] \leq \varepsilon.$$

The proof of the above theorem, which appears in Appendix B, relies on the observation that the structure of the maximum function is such that its surface consists of d linear regions (corresponding to the subsets of \mathbb{R}^d where each coordinate is maximal). Since each such region has exactly $d - 1$ faces (corresponding to the hyperplanes where one coordinate overtakes the other), we can use the first hidden layer to compute the linear function $\mathbf{x} \mapsto x_i$ and “peel off” the surface of the function at the relevant faces using a ReLU neuron with a very large negative slope. We then use the second hidden layer to truncate negative values. By adding such “polytope functions”, we are able to obtain a good approximation of the maximum function at points where the coordinates are sufficiently distant from each other. We refer the reader to Definition 3.6 for the formal construction and Fig. 1 for an illustration.

It is interesting to note that the width of the approximating network in our result only scales with the input dimension d , and does not scale with the desired target accuracy. Rather, by increasing the magnitude of the weights of the approximating network we can control the accuracy of the approximation. This is in contrast to many other approximation regimes where an improvement in the approximation accuracy necessitates an increase in width. E.g., when approximating the maximum function using depth 2 (see Proposition C.1 in the appendix) or when approximating non-linear functions using ReLU networks (see Safran and Shamir [20]).

Our result allows the approximation of the maximum function using a network of size $\mathcal{O}(d^2)$. It is known, however, that the maximum function can be computed exactly using a smaller network

of size $\mathcal{O}(d)$ if we allow depth $\mathcal{O}(\log(d))$. It is therefore natural to ask whether by utilizing depth, we can obtain more efficient approximations of the maximum function using our construction. Perhaps surprisingly, we are able to approximate the maximum function using linear width but by only requiring the depth to scale as $\mathcal{O}(\log(\log(d)))$. More formally, we have the following theorem.

Theorem 3.4. *Let \mathcal{D} be any distribution satisfying Assumption 3.2. Then for any $\varepsilon > 0$ and naturals $d \geq 58$ and $1 \leq k \leq \lceil \log(\log(d) + 1) \rceil$, there exists a ReLU network \mathcal{N} of depth $2k + 1$ and width at most $20d^{1+\frac{1}{2^k-1}}$ such that*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] \leq \varepsilon.$$

In particular, we have a ReLU network of width $40d$ and depth $2\lceil \log(\log(d) + 1) \rceil + 1$ which approximates $f_d(\mathbf{x})$ to accuracy ε with respect to the distribution \mathcal{D} .

The proof of the above theorem, which appears in Appendix B, exploits the key observation that the maximum taken over sub-vectors of maxima is the maximum of the vector. We can thus partition our input into smaller batches and use Thm. 3.3 to compute the maximum over each of these batches. Since we may vary the size of the batches across layers, we can gradually take larger and larger batches as our computation propagates deeper in the network, while keeping the width roughly the same across all layers. This enables a double exponential decay in the number of maxima computed at each layer, requiring depth of only $\mathcal{O}(\log(\log(d)))$ for approximating the maximum using width linear in d (see Definition 3.7 for the formal construction and Fig. 2 for an illustration). While our result is not directly comparable to previous approximations of the maximum function [1, 17] since these allow its exact computation using depth and network size linear in d , it does provide an alternative construction which allows the approximation of max using a network of size $d \log(\log(d))$, but with a much smaller depth of only $\log(\log(d))$.

Of particular interest is the following corollary, which provides an approximation guarantee in the case where the data are sampled uniformly from a hypercube. Most importantly, in such a case we can guarantee an approximation to accuracy $\varepsilon > 0$ using a network with weights that scale polynomially with d and linearly with $1/\varepsilon$. This property will turn out to be useful in the next section where we will show lower bounds for approximating the maximum function with respect to the uniform distribution.

Corollary 3.5. *For any $\varepsilon > 0$ and naturals $d \geq 58$ and $1 \leq k \leq \lceil \log(\log(d) + 1) \rceil$, there exists a ReLU network \mathcal{N} of depth $2k + 1$, width at most $20d^{1+\beta(k)}$ and weights of magnitude $\mathcal{O}\left(\frac{d^4 R^2}{\varepsilon}\right)$ such that*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}([0, R]^d)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] \leq \varepsilon.$$

Having stated our main positive approximation results, we now turn to specify the constructions used to achieve them. Beginning with defining the depth 3 network used in Thm. 3.3, we have the following architecture.

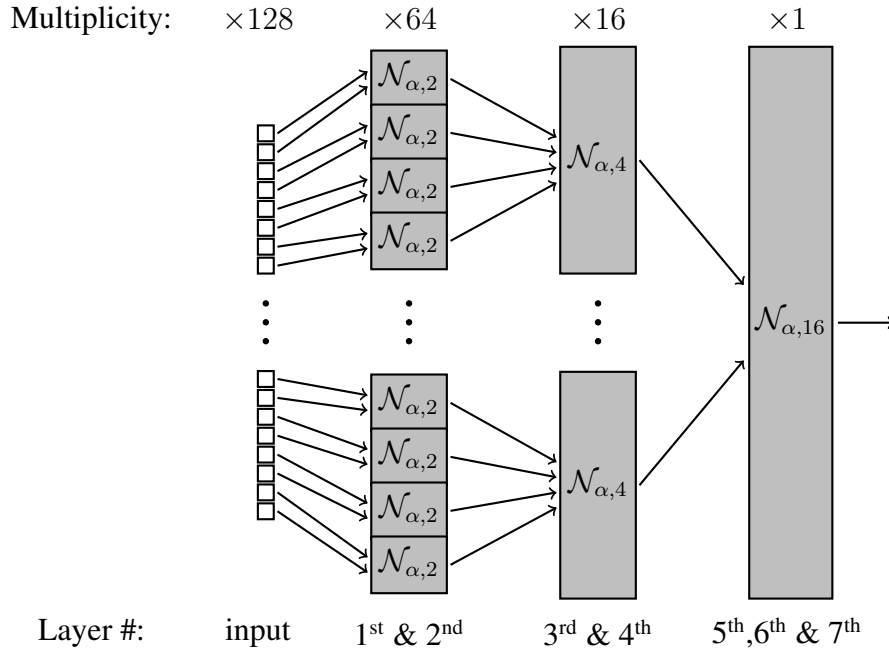


Figure 2: Sketch of the architecture $\mathcal{N}_{\alpha,128}^3$ which approximates $f_{128}(\cdot)$ using depth 7 and width $\approx 128^{8/7} = 256$. The multiplicity row at the top counts the number of components in each layer, and the layer # row at the bottom indicates which layers participate in the computation of each component. The height of each $\mathcal{N}_{\alpha,\cdot}$ component is roughly proportional to its width. Each pair of hidden layers increases the batch size on which the maxima are computed quadratically, while maintaining the width of the network roughly the same across all hidden layers. This results in a double exponential decay of the batch size, allowing an approximation of the maximum with depth $\mathcal{O}(\log(\log(d)))$.

Definition 3.6. [Depth 3 approximation] Given a weight upper bound $\alpha > 0$ and input dimension d , we define the following depth 3 width $d(d+1)$ network which approximates $f_d(\mathbf{x})$:

$$\mathcal{N}_{\alpha,d}(\mathbf{x}) := \sum_{i=1}^d \left(\left[\left[x_i \right]_+ - \sum_{\substack{j=1 \\ j \neq i}}^d [\alpha x_j - \alpha x_i]_+ \right] - \left[\left[-x_i \right]_+ - \sum_{\substack{j=1 \\ j \neq i}}^d [\alpha x_j - \alpha x_i]_+ \right] \right)_+.$$

We remark that we occasionally omit the dimension subscript whenever clear from context, and we note that the above architecture can be realized using a width $d(d+1)$ ReLU network since the inner sum terms are identical, and thus computing both requires only $d-1$ neurons.

Next, we define the architecture which approximates the function $f_d(\mathbf{x})$ using depth $2k+1$ and width at most $20d^{1+\beta(k)}$, achieving the approximating result stated in Thm. 3.4.

Definition 3.7. [Depth $2k+1$ approximation] Given a weight upper bound $\alpha > 0$ and input dimension d , we define the following depth $2k+1$ width at most $20d^{1+\beta(k)}$ network $\mathcal{N}_{\alpha,d}^k$ which approximates $f_d(\mathbf{x})$ in a recursive manner:

- For $k=1$, we have $\mathcal{N}_{\alpha,d}^1 \equiv \mathcal{N}_{\alpha,d}$.
- For integer $k > 1$, we partition the input into $\lceil d^{1-\beta(k)} \rceil$ batches, each of size at most $\lceil d^{\beta(k)} \rceil$. For each batch, the first two hidden layers compute the function $\mathcal{N}_{\alpha, \lceil d^{\beta(k)} \rceil}$. The output over all the batches is then fed into the sub-network $\mathcal{N}_{\alpha, \lceil d^{1-\beta(k)} \rceil}^{k-1}$ which consists the remaining layers of the network $\mathcal{N}_{\alpha,d}^k$.

4 Approximation lower bounds

Having shown a positive approximation result for the maximum function in the previous section, we now turn to complement our approximation upper bounds with lower bounds in this section.

4.1 Improving accuracy requires increasing width for depth 2 networks

Before presenting our main theorem for this subsection, we first state the following very mild assumption that we use which is adopted from Eldan and Shamir [4]:

Assumption 4.1 (Polynomially-bounded activation). *The activation function σ is Lebesgue measurable and satisfies*

$$|\sigma(x)| \leq C_\sigma (1 + |x|^{\alpha_\sigma})$$

for all $x \in \mathbb{R}$ and for some constants $C_\sigma, \alpha_\sigma > 0$.

Our depth 2 lower bound is the following:

Theorem 4.2. *Let $\ell \geq 1$ be arbitrary and suppose that σ satisfies Assumption 4.1. Then there exist constants $c_1, c_2 > 0$ which depend solely on σ such that for all dimensions $d \geq c_1$, a σ depth 2 neural network \mathcal{N} of width at most d^ℓ and with weights bounded by $\mathcal{O}(\exp(\mathcal{O}(d)))$ must satisfy*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}([0,1]^d)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] > \Omega(d^{-c_2 \cdot \ell}).$$

The proof of the above theorem, which appears in Appendix C.1, builds on the proof technique introduced in Eldan and Shamir [4]. Roughly speaking, they build on the important observation that a neural network \mathcal{N} approximates a function f if and only if the Fourier transform of \mathcal{N} approximates the Fourier transform of f . Our main technical contribution here is the computation of the Fourier transform of the maximum function and showing that it has sufficient L_2 mass far away from the origin which is sufficiently spread. This, in turn, shows that the support of the Fourier transform of a neural network (which is a linear combination of the Fourier transform of its activation function) must be contained inside a d -dimensional union of ‘Gaussian tubes’. Namely, under the assumption of exponentially bounded weights, the approximation contribution of each neuron is negligible outside a union of tubes with bounded radius. This entails that to approximate the Fourier transform of the maximum function, one must use sufficiently many neurons in order to be able to capture its non-trivial structure which is sufficiently spread across the domain of approximation.

While the above theorem establishes a polynomial rather than exponential separation between depths 2 and 3, such a gap is nevertheless significant since it provides a compelling practical example where depth is more beneficial than width: Modern machine learning problems are often high-dimensional, hence even such polynomial gaps quickly translate into a significant practical advantage in the size of the required network. We further remark that our assumption that the approximating network has exponentially bounded weights is mild and very reasonable. This follows from the fact that approximations with weights that have exponential magnitude are known to be difficult to learn using stable gradient descent [19], so from a more practical perspective having exponentially bounded weights and having unbounded weights is equivalent. Lastly, as we also pointed out earlier, our result also implies a lower bound with constant accuracy by rescaling our domain of approximation polynomially with d . Performing such a manipulation is justified since our upper bounds from the previous section are not sensitive to scaling of the domain and merely require that the weights of the approximating network would also scale appropriately (Corollary 3.5). This enables us to show a separation in which a depth 2 network cannot approximate the maximum to better than constant accuracy, whereas a depth 3 network with a fixed width of $d(d+1)$ that does not depend on the desired accuracy can achieve arbitrarily good accuracy. In contrast, known results in the literature do require that the width of the depth 3 network would scale with the accuracy parameter (e.g. Eldan and Shamir [4], Daniely [3] and results that build on their technique – see related work subsection).

4.2 Depth 3 requires $\Omega(d^2)$ width

In this subsection, we show that approximating the maximum function using depth 3 ReLU networks with weights bounded by $\mathcal{O}(\exp(\mathcal{O}(d)))$ requires width at least $\Omega(d^2)$. Our main result in this subsection is the following.

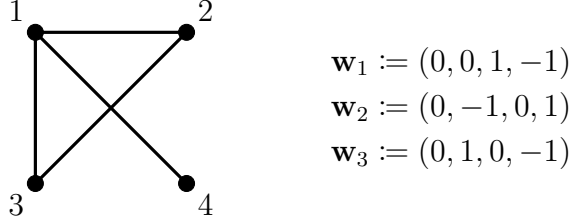


Figure 3: An informal simplification of the main proof idea behind Thm. 4.3. To identify a region in the domain of approximation where $f_d(\cdot)$ is poorly approximated, we construct a graph G as follows: We begin with the complete graph on d vertices (a vertex for each input), and we remove an edge (i, j) if there exists a neuron in the first hidden layer whose weight vector is all zeros except for the coordinates i, j . In the example portrayed in the figure, for weight vectors $\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3$ we remove the edge $(3, 4)$ due to \mathbf{w}_2 , and the edge $(2, 4)$ due to either of $\mathbf{w}_2, \mathbf{w}_3$. Loosely speaking, whenever the first hidden layer contains a neuron with an all-zero weight vector except for two coordinates that equal -1 and 1 , then this neuron is able to capture the non-linearity of $f_d(\cdot)$ associated with the two non-zero coordinates when one overtakes the other in value. Since the width of the network is at most $d^2/5$, this guarantees that G contains at least $\binom{d}{2} - d^2/5 > d^2/4$ edges. By Mantel’s theorem (see Thm. D.1), G must contain a triangle, which implies the existence of a 3-dimensional sub-cube where the non-linearities in the first hidden layer are redundant. This effectively reduces the approximation of $f_d(\cdot)$ using a depth 3 network to the approximation of $f_3(\cdot)$ using a depth 2 network.

Theorem 4.3. *Suppose that \mathcal{N} is a depth 3 ReLU network of width at most $\frac{d^2}{5}$ and with weights bounded by $\mathcal{O}(\exp(\mathcal{O}(d)))$. Then there exist absolute constants $c_1, c_2 > 0$ such that for all $d \geq c_1$*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}([0,1]^d)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] > \Omega(d^{-c_2}).$$

The proof of the above theorem, which appears in Appendix C.2, exploits the structure of the maximum function which computes a lower dimensional version of itself on every subset of its inputs. Using a combinatorial argument, we show that with fewer than $d^2/5$ neurons in the first hidden layer, our approximating network must be able to approximate the maximum over three inputs well on a non-negligible subset of its domain with its remaining layers (see Fig. 3 for a more detailed explanation of the proof technique). Since in the previous subsection we have shown a lower bound on the approximation capabilities of depth 2 ReLU networks for the maximum function, this implies that if the second hidden layer is also at most $d^2/5$, then we cannot obtain a good approximation.

We remark that together with Thm. 3.3, we establish tight bounds on the capability of depth 3 ReLU networks to approximate the maximum function (up to constant factors). It is also interesting to note that other existing lower bound techniques in the literature such as region counting arguments (c.f. Telgarsky [23]) when applied to the maximum function yield a far weaker lower bound for depth 3 networks, since the maximum function consists of d different linear regions, a number which is attainable by merely using $\log(d)$ neurons. In contrast, our lower bound of $\Omega(d^2)$ highlights an inherent limitation of depth 3 ReLU networks for capturing the particular structure

of the maximum function. When combined with Thm. 3.4 for $k = 2$, our theorem also implies a polynomial depth separation between depths 3 and 5 where the former requires width $\Omega(d^2)$, yet the latter can approximate with width $\mathcal{O}(d^{4/3})$.

4.3 Width of at least d is necessary

Having shown lower bounds for depth 2 and 3 ReLU networks, we now turn to show a general width d lower bound requirement for approximating the maximum function.

Theorem 4.4. *Let \mathcal{N} be neural network employing any activation function and having first hidden layer width of at most $d - 1$. Then*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}([0,1]^d)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] \geq \Omega(d^{-4.5}).$$

The proof of the above theorem, which appears in Appendix C.3, relies on the observation that having fewer than d neurons in the first hidden layer implies that the linear transformation defined by them has a non-trivial kernel, and therefore establishes the existence of some direction in the domain of approximation where the function computed by the network is constant. Since the maximum function typically does not remain constant in most directions, this results in a non-trivial approximation error.

We remark that the exponent of -4.5 in the target accuracy can possibly be improved somewhat, but in any case it must be strictly positive, since a network with a single neuron which computes the constant value function $1 - 1/d$ will achieve better than constant accuracy over the domain $[0, 1]^d$. Moreover, we note that this result also immediately implies a size d lower bound. Together with Thm. 3.4, this establishes size bounds for approximating the maximum function using ReLU networks that are tight up to a factor of $\mathcal{O}(\log(\log(d)))$.

5 Summary

We have shown that the maximum function can be gradually approximated more efficiently by increasing the depth of the approximating ReLU network. This holds under an appropriate (but mostly mild) assumption on the distribution of the data (Assumption 3.2). Interestingly, the width in our positive approximation results does not scale with the desired target accuracy, but rather by increasing the magnitude of the weights of the approximating network we can obtain an arbitrarily good approximation. Assuming exponentially bounded weights, we show a polynomial lower bound on the required width when approximating the maximum function using depth 2, and a quadratic lower bound on the width required for approximating using depth 3. Additionally, we also provide a general width d lower bound for approximating the maximum function using neural networks of any depth or with any activation function. Our results establish a partial depth hierarchy for approximating a simple target function and with respect to the uniform distribution on a hypercube, which provides a more grounded example for the benefits of depth compared to previous results which make more stylized assumptions on the problem.

Our analysis leaves several important open questions. First, our lower bound for depth 2 is (inverse) polynomial in the desired accuracy, which becomes polynomial in the input dimension d if we scale the domain with d . However, it is not clear what is the optimal rate at which a depth 2 network can approximate the maximum function, and whether this quantity is polynomial or rather exponential in the input dimension. Second, despite our tight $\Theta(d^2)$ bound on the width for approximating the maximum using depth 3, our lower bound for deeper architectures is only linear, which leaves open the question of showing superlinear width lower bounds for depths larger than 3. Moreover, our upper bounds essentially suggest that depth $2k + 1$ and depth $2k + 2$ ReLU network approximations require the same width, up to constant factors. It is therefore natural to ask whether one can improve our depth $2k + 1$ upper bounds to apply to depth $k + 1$ instead. Finally, our analysis opens an avenue for novel optimization-based separations for the maximum function. Proving that indeed deep architectures are capable of learning the representations constructed by our upper bounds (efficiently, from finite data) using standard techniques such as gradient descent is a tantalizing future research direction. Such a result holds the potential to establish an optimization-based depth hierarchy for learning the maximum function, exemplifying the superiority of depth in a simple and natural problem setting.

Acknowledgements

Part of this work was done while the second author was visiting the Simons Institute for the Theory of Computing. Their hospitality is greatly acknowledged.

References

- [1] R. Arora, A. Basu, P. Mianjy, and A. Mukherjee. Understanding deep neural networks with rectified linear units. *arXiv preprint arXiv:1611.01491*, 2016.
- [2] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- [3] A. Daniely. Depth separation for neural networks. In *Conference on Learning Theory*, pages 690–696. PMLR, 2017.
- [4] R. Eldan and O. Shamir. The power of depth for feedforward neural networks. In *Conference on learning theory*, pages 907–940, 2016.
- [5] C. Haase, C. Hertrich, and G. Loho. Lower bounds on the depth of integral relu neural networks via lattice polytopes. *arXiv preprint arXiv:2302.12553*, 2023.
- [6] B. Hanin and D. Rolnick. Complexity of linear regions in deep networks. In *International Conference on Machine Learning*, pages 2596–2604. PMLR, 2019.
- [7] C. Hertrich, A. Basu, M. Di Summa, and M. Skutella. Towards lower bounds on the depth of relu neural networks. *Advances in Neural Information Processing Systems*, 34:3336–3348, 2021.

- [8] D. Hsu, C. Sanford, R. A. Servedio, and E.-V. Vlatakis-Gkaragkounis. On the approximation power of two-layer networks of random relus. *arXiv preprint arXiv:2102.02336*, 2021.
- [9] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015.
- [10] D. M. Kane and R. Williams. Super-linear gate and super-quadratic wire lower bounds for depth-two and depth-three threshold circuits. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 633–643, 2016.
- [11] B. Kriener, R. Chaudhuri, and I. R. Fiete. Robust parallel decision-making in neural circuits with nonlinear inhibition. *Proceedings of the National Academy of Sciences*, 117(41):25505–25516, 2020.
- [12] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- [13] R. Lang. A note on the measurability of convex sets. *Archiv der Mathematik*, 47:90–92, 1986.
- [14] M. Leshno, V. Y. Lin, A. Pinkus, and S. Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural networks*, 6(6): 861–867, 1993.
- [15] S. Liang and R. Srikant. Why deep neural networks for function approximation? In *5th International Conference on Learning Representations, ICLR 2017*, 2017.
- [16] E. Malach, G. Yehudai, S. Shalev-Shwartz, and O. Shamir. The connection between approximation, depth separation and learnability in neural networks. *arXiv preprint arXiv:2102.00434*, 2021.
- [17] K. Matoba, N. Dimitriadis, and F. Fleuret. The theoretical expressiveness of maxpooling. *arXiv preprint arXiv:2203.01016*, 2022.
- [18] A. Mukherjee and A. Basu. Lower bounds over boolean inputs for deep neural networks with relu gates. *arXiv preprint arXiv:1711.03073*, 2017.
- [19] I. Safran and J. Lee. Optimization-based separations for neural networks. In *Conference on Learning Theory*, pages 3–64. PMLR, 2022.
- [20] I. Safran and O. Shamir. Depth-width tradeoffs in approximating natural functions with neural networks. In *International Conference on Machine Learning*, pages 2979–2987. PMLR, 2017.
- [21] I. Safran, R. Eldan, and O. Shamir. Depth separations in neural networks: what is actually being separated? In *Conference on Learning Theory*, pages 2664–2666. PMLR, 2019.

- [22] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [23] M. Telgarsky. Benefits of depth in neural networks. In *Conference on learning theory*, pages 1517–1539. PMLR, 2016.
- [24] G. Vardi, D. Reichman, T. Pitassi, and O. Shamir. Size and depth separation in approximating benign functions with neural networks. In *Conference on Learning Theory*, pages 4195–4223. PMLR, 2021.
- [25] L. Venturi, S. Jelassi, T. Ozuch, and J. Bruna. Depth separation beyond radial functions. *arXiv preprint arXiv:2102.01621*, 2021.
- [26] D. Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*, 94:103–114, 2017.

A Distributions that satisfy Assumption 3.2

In this appendix, we exemplify two instances of distribution that satisfy Assumption 3.2. While the examples presented here are quite broad, they are in no way exhaustive, and a far richer family of distributions can be shown to satisfy the assumption.

Theorem A.1. *Suppose that X is absolutely continuous with bounded density and has a finite second moment. Then $\mathbf{X} := (X_1, \dots, X_d)$ satisfies Assumption 3.2, where each X_i is an i.i.d X random variable.*

Proof. Beginning with Item 1, we have

$$\mathbb{E}_{\mathbf{X}} [\|\mathbf{X}\|_{\infty}^2] \leq \mathbb{E}_{\mathbf{X}} [\|\mathbf{X}\|_1^2] \leq d \sum_{i=1}^d \mathbb{E}_{\mathbf{X}} [X_i^2] = d^2 \mathbb{E}_X [X^2] < \infty.$$

In the above, the second inequality is an application of Cauchy-Schwartz to the vectors $(1, \dots, 1)$ and (X_1, \dots, X_d) and due to the linearity of expectation, and the last inequality is due to our assumption that X has a finite second moment.

Moving on to Item 2, assuming X has density f satisfying $\sup_{x \in \mathbb{R}} f(x) \leq C$ for some $C > 0$, we have that the density of the ratio distribution between two different coordinates of \mathbf{X} satisfies

$$f_R(r) = \int_{\mathbb{R}} |x| f(r \cdot x) f(x) dx \leq C \cdot \int_{\mathbb{R}} |x| f(x) dx = C \mathbb{E}_X [|X|] \leq C \sqrt{\mathbb{E}_X [X^2]} < \infty,$$

where the penultimate inequality follows from Jensen’s inequality applied to the function $x \mapsto x^2$. The above implies that for any two coordinates X_i, X_j , we have

$$\lim_{\delta \rightarrow 0} \frac{X_i}{X_j} \in [1 - \delta, 1 + \delta] = 0,$$

therefore by taking a union bound over all $\leq d^2$ pairs of coordinates we have that

$$\lim_{\delta \rightarrow 0} \mathbb{P}_{\mathbf{X} \sim \mathcal{D}} [\mathbf{X} \notin S_\delta] = 0.$$

□

Theorem A.2. *Suppose that $\mathbf{X} := (X_1, \dots, X_d)$ satisfies Item 1 in Assumption 3.2. Then the vector $\mathbf{X} + \mathbf{Y}$ satisfies Assumption 3.2, where $\mathbf{Y} := (Y_1, \dots, Y_d)$ is an i.i.d noise vector such that Y_i is absolutely continuous, with bounded density, and $0 < \mathbb{E}[Y_i^2] < \infty$.*

Proof. To show that Item 1 holds, we have by Thm. A.1 that \mathbf{Y} satisfies Item 1. This entails that

$$\|\mathbf{X} + \mathbf{Y}\|_\infty \leq \|\mathbf{X}\|_\infty + \|\mathbf{Y}\|_\infty < \infty.$$

To show that Item 2 is satisfied, consider any coordinates $i \neq j$. Let $z_i := x_i + y_i$ denote the realizations of $Z_i := X_i + Y_i$. Then we have that for any realization x_j of X_j , it must hold that y_j falls within an interval of length at most 2δ to have that

$$z_j \in [z_i - \delta, z_i + \delta] \iff |z_j - z_i| \leq \delta.$$

Since

$$\sup_{a \in \mathbb{R}} \mathbb{P}[Y_j \in [a, a + 2\delta]] = \sup_{a \in \mathbb{R}} \int_a^{a+2\delta} f(x) dx \leq 2\delta C,$$

where f is the density of Y which satisfies $\sup_x f(x) \leq C$ by assumption, we obtain

$$\mathbb{P}[|z_j - z_i| \leq \delta] \leq 2\delta C.$$

Since Z_i has bounded first and second moments for all i by assumption, we have from Chebyshev's inequality that there exists some $M_\delta > 0$ such that $\mathbb{P}[\|\mathbf{X}\|_\infty \leq M_\delta] \geq 1 - \delta$. Since Lemma D.2 guarantees that if $|z_j - z_i| > \delta$ then $\mathbf{z} \in S_{\delta/M_\delta}$, we have from a union bound taken over all $\leq d^2$ pairs of coordinates and the event where $\|\mathbf{X}\|_\infty$ is bounded that

$$\lim_{\delta \rightarrow 0} \mathbb{P}_{\mathbf{X} \sim \mathcal{D}} [\mathbf{X} \notin S_\delta] = 0.$$

□

B Proofs from Sec. 3

To prove Thm. 3.3 and Thm. 3.4, we would first need the following lemmas and proposition:

Lemma B.1. *Given a real $\alpha > 0$ and integer $i \geq 1$, let*

$$n_{\alpha,i}^+(\mathbf{x}) := \left[[x_i]_+ - \sum_{\substack{j=1 \\ j \neq i}}^d [\alpha x_j - \alpha x_i]_+ \right]_+, \quad n_{\alpha,i}^-(\mathbf{x}) := - \left[[-x_i]_+ - \sum_{\substack{j=1 \\ j \neq i}}^d [\alpha x_j - \alpha x_i]_+ \right]_+.$$

Then

$$n_{\alpha,i}^+(\mathbf{x}) = \begin{cases} x_i, & \text{If } f_d(\mathbf{x}) = x_i \text{ and } x_i \geq 0, \\ 0, & \text{If } f_d(\mathbf{x}) > x_i \text{ and } \mathbf{x} \in S_{1/\alpha}, \end{cases}$$

and

$$n_{\alpha,i}^-(\mathbf{x}) = \begin{cases} x_i, & \text{If } f_d(\mathbf{x}) = x_i \text{ and } x_i \leq 0, \\ 0, & \text{If } f_d(\mathbf{x}) > x_i \text{ and } \mathbf{x} \in S_{1/\alpha}. \end{cases}$$

Proof. We will only focus on the proof for $n_{\alpha,i}^+$ since the analysis is symmetric for $n_{\alpha,i}^-$.

Suppose that $f_d(\mathbf{x}) = x_i$. Then $x_j < x_i$ for all $j \neq i$, implying that $\alpha x_j - \alpha x_i < 0$ and thus

$$n_{\alpha,i}^+(\mathbf{x}) = \left[[x_i]_+ - \sum_{\substack{j=1 \\ j \neq i}}^d [\alpha x_j - \alpha x_i]_+ \right]_+ = [x_i]_+.$$

Suppose that $f_d(\mathbf{x}) > x_i$. Then if $x_i \leq 0$ we immediately have that $n_{\alpha,i}^+(\mathbf{x}) = 0$. Assuming $x_i > 0$ and letting $j := \operatorname{argmax}_{i \in [d]} x_i$, we have that $x_j > x_i > 0$. Next, from the assumption $\mathbf{x} \in S_{1/\alpha}$ we obtain

$$\frac{x_j}{x_i} > 1 + \frac{1}{\alpha},$$

which entails

$$\alpha x_j - \alpha x_i > x_i,$$

implying that

$$[x_i]_+ - \sum_{\substack{i=1 \\ i \neq i}}^d [\alpha x_i - \alpha x_i]_+ \leq [x_i]_+ - [\alpha x_j - \alpha x_i]_+ < x_i - x_i = 0,$$

and thus

$$n_{\alpha,i}^+(\mathbf{x}) = 0.$$

□

Lemma B.2. *Given a real $\alpha > 0$, we have*

$$\mathcal{N}_\alpha(\mathbf{x}) = f_d(\mathbf{x}), \quad \forall \mathbf{x} \in S_{1/\alpha},$$

and

$$|\mathcal{N}_\alpha(\mathbf{x})| \leq \|\mathbf{x}\|_1, \quad \forall \mathbf{x} \in \mathbb{R}^d.$$

Proof. By the definition of \mathcal{N}_α and Lemma B.1, we have for all $\mathbf{x} \in S_{1/\alpha}$ that

$$\mathcal{N}_\alpha(\mathbf{x}) = \sum_{i=1}^d (n_{\alpha,i}^+(\mathbf{x}) + n_{\alpha,i}^-(\mathbf{x})) = \sum_{i=1}^d x_i \cdot \mathbb{1}\{f_d(\mathbf{x}) = x_i\} = f_d(\mathbf{x}).$$

Assuming any arbitrary $\mathbf{x} \in \mathbb{R}^d$, we have by the definitions of $n_{\alpha,i}^{\pm}$ that

$$|n_{\alpha,i}^{\pm}(\mathbf{x})| = \left| [\pm x_i]_+ - \sum_{\substack{j=1 \\ i \neq j}}^d [\alpha x_j - \alpha x_i]_+ \right| \leq [\pm x_i]_+ \leq |x_i|,$$

therefore by the definition of \mathcal{N}_α and the fact that at most one of $n_{\alpha,i}^{\pm}$ is non-zero, we have

$$|\mathcal{N}_\alpha(\mathbf{x})| \leq \sum_{i=1}^d |n_{\alpha,i}^+(\mathbf{x}) + n_{\alpha,i}^-(\mathbf{x})| \leq \|\mathbf{x}\|_1.$$

□

Proposition B.3. *Given any real $\alpha > 0$, we have*

$$\mathcal{N}_\alpha^k(\mathbf{x}) = f_d(\mathbf{x}), \quad \forall \mathbf{x} \in S_{1/\alpha},$$

and

$$|\mathcal{N}_\alpha^k(\mathbf{x})| \leq d \cdot f_d(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^d.$$

Proof. The proof follows by induction on k . The base case $k = 1$ follows immediately from Lemma B.2. In what follows, given a natural k , recall that we use the shorthand $\beta(k) := \frac{1}{2^{k-1}}$.

For the inductive step, assume the induction hypothesis for k , and consider the network \mathcal{N}_α^{k+1} . Since any sub-vector of a δ -separated vector is also δ -separated for all $\delta > 0$, we have from Lemma B.2 that the output of the second hidden layer of \mathcal{N}_α^{k+1} is the maximum over each of the $\lceil d^{1-\beta(k)} \rceil$ batches. Applying the inductive hypothesis on the sub-network consisting of layers 3 to $2k + 3$, the network outputs $f_d(\mathbf{x})$.

For the second part of the proposition, partition the input \mathbf{x} into $\lceil d^{1-\beta(k)} \rceil$ batches such that the vector of inputs in each batch has dimension at most $\lceil d^{\beta(k)} \rceil$. Then by Lemma B.2, each coordinate in the output of the second hidden layer of \mathcal{N}_α^{k+1} is upper bounded by

$$\left(\|\mathbf{x}_1\|_1, \dots, \|\mathbf{x}_{\lceil d^{1-\beta(k)} \rceil}\|_1 \right).$$

Applying the induction hypothesis on the sub-network consisting of layers 3 to $2k + 3$, we obtain

$$|\mathcal{N}_\alpha^k(\mathbf{x})| \leq \sum_{i=1}^{\lceil d^{1-\beta(k)} \rceil} \|\mathbf{x}_i\|_1 = \|\mathbf{x}\|_1 \leq d \cdot f_d(\mathbf{x}).$$

□

With the above lemmas and proposition, we are now ready to prove the theorems. Since the proof of Thm. 3.4 follows mainly by induction and since Thm. 3.3 consists the base case for the induction, it would be convenient to prove both theorems simultaneously.

Proofs of Thm. 3.3 and Thm. 3.4. Recall we use the shorthand $\beta(k) := \frac{1}{2^{k-1}}$ for any natural $k \geq 1$. We begin with asserting the size of the approximating network. We have that \mathcal{N}_α^k has depth $2k + 1$ and weights of magnitude at most α by definition (note that the output neuron of \mathcal{N}_α has weights of magnitude exactly 1, and therefore composing its weights with the subsequent layer's weights does not increase the magnitude).

Next, we bound the width of \mathcal{N}_α^k using induction. To this end, we will show that for all natural $k \geq 1$ we have an upper bound on the width of

$$\prod_{i=1}^k \left(1 + \frac{2}{i^3}\right)^2 d^{1+\beta(k)}.$$

By Lemma D.3, we have that $\prod_{i=1}^\infty \left(1 + \frac{2}{i^3}\right)^2 \leq 20$, thus the above implies the desired upper bound on the width.

The base case is immediate since $\mathcal{N}_\alpha^1 \equiv \mathcal{N}_\alpha$ which has width exactly $d(d+1) \leq 2d^2$. Assume the inductive hypothesis for k , and consider the network \mathcal{N}_α^{k+1} . Its first two hidden layers consist of $\lceil d^{1-\beta(k+1)} \rceil \leq d^{1-\beta(k+1)} + 1$ batches of $\mathcal{N}_{\alpha, \lceil d^{\beta(k+1)} \rceil}$, each of which having width at most

$$\lceil d^{\beta(k+1)} \rceil^2 + \lceil d^{\beta(k+1)} \rceil \leq (d^{\beta(k+1)} + 1)^2 + d^{\beta(k+1)} + 1 = d^{2\beta(k+1)} + 3d^{\beta(k+1)} + 2,$$

for a total width upper bound of

$$(d^{1-\beta(k+1)} + 1) (d^{2\beta(k+1)} + 3d^{\beta(k+1)} + 2) \leq 12d^{1+\beta(k+1)},$$

thus implying an upper bound on the width of

$$12d^{1+\beta(k+1)} \leq \prod_{i=1}^{k+1} \left(1 + \frac{2}{i^3}\right)^2 d^{1+\beta(k+1)},$$

since $k \geq 1$ and the product over the first two elements is at least 14. Moving on to bound the width of the remaining layers, we have by definition that layers 3 to $2k + 3$ is the network $\mathcal{N}_{\alpha, \lceil d^{1-\beta(k+1)} \rceil}^k$, which by the induction hypothesis has width at most

$$\begin{aligned} \prod_{i=1}^k \left(1 + \frac{2}{i^3}\right)^2 \lceil d^{1-\beta(k+1)} \rceil^{1+\beta(k)} &\leq \prod_{i=1}^k \left(1 + \frac{2}{i^3}\right)^2 (d^{1-\beta(k+1)} + 1)^{1+\beta(k)} \\ &\leq \prod_{i=1}^k \left(1 + \frac{2}{i^3}\right)^2 \left(\left(1 + \frac{2}{(k+1)^3}\right) d^{1-\beta(k+1)} \right)^{1+\beta(k)} \\ &\leq \prod_{i=1}^k \left(1 + \frac{2}{i^3}\right)^2 \left(1 + \frac{2}{(k+1)^3}\right)^2 (d^{1-\beta(k+1)})^{1+\beta(k)} \\ &= \prod_{i=1}^{k+1} \left(1 + \frac{2}{i^3}\right)^2 d^{1+\beta(k+1)}. \end{aligned}$$

In the above, the second inequality follows from Lemma D.4 by our assumption that $d \geq 58$ and $1 \leq k \leq \lceil \log(\log(d) + 1) \rceil$, and the third inequality follows from the fact that $\beta(k) \leq 1$ for all natural $k \geq 1$. We thus conclude that \mathcal{N}_α^k has width at most $20d^{1+\beta(k)}$.

Turning to bound the approximation error of \mathcal{N}_α^k , we first have from Assumption 3.2 that there exists some $\delta_0 > 0$ such that

$$\mathbb{P}_{\mathbf{x} \sim \mathcal{D}} [\mathbf{x} \notin S_{\delta_0}] \leq \frac{\varepsilon}{(d+1)^2 \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|\mathbf{x}\|_\infty^2]}. \quad (1)$$

We now compute using the law of total expectation

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[(\mathcal{N}_{1/\delta_0}^k(\mathbf{x}) - f_d(\mathbf{x}))^2 \right] &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[(\mathcal{N}_{1/\delta_0}^k(\mathbf{x}) - f_d(\mathbf{x}))^2 \mid \mathbf{x} \in S_{\delta_0} \right] \cdot \mathbb{P}_{\mathbf{x} \sim \mathcal{D}} [\mathbf{x} \in S_{\delta_0}] \\ &\quad + \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[(\mathcal{N}_{1/\delta_0}^k(\mathbf{x}) - f_d(\mathbf{x}))^2 \mid \mathbf{x} \notin S_{\delta_0} \right] \cdot \mathbb{P}_{\mathbf{x} \sim \mathcal{D}} [\mathbf{x} \notin S_{\delta_0}] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[(\mathcal{N}_{1/\delta_0}^k(\mathbf{x}) - f_d(\mathbf{x}))^2 \mid \mathbf{x} \notin S_{\delta_0} \right] \cdot \mathbb{P}_{\mathbf{x} \sim \mathcal{D}} [\mathbf{x} \notin S_{\delta_0}] \\ &\leq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [(d+1)^2 \|\mathbf{x}\|_\infty^2] \cdot \frac{\varepsilon}{(d+1)^2 \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|\mathbf{x}\|_\infty^2]} = \varepsilon, \end{aligned}$$

where the second equality follows from Proposition B.3, and the inequality also follows from Proposition B.3 and Eq. (1).

Lastly, we verify that $\lceil \log(\log(d) + 1) \rceil = \mathcal{O}(\log(\log(d)))$. We have

$$20d^{1+\frac{1}{2^k-1}} \leq 20d^{1+\frac{1}{\log(d)}} = 20d \cdot d^{\frac{1}{\log(d)}} = 20d \cdot 2^{\log(d) \frac{1}{\log(d)}} = 40d,$$

thus for this choice of k we have a network of depth $\mathcal{O}(\log(\log(d)))$ and width $\mathcal{O}(d)$ which approximates $f_d(\cdot)$, concluding the proof of the theorem. \square

Proof of Corollary 3.5. To prove the corollary, we need only show that $\mathcal{D} \sim \mathcal{U}([0, R]^d)$ satisfies Assumption 3.2 and compute the $\delta_0 > 0$ for which Eq. (1) holds. Starting with Item 1, it is trivial that

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}([0, R]^d)} [\|\mathbf{x}\|_\infty^2] \leq R^2. \quad (2)$$

Moving on to Item 2, let $\delta > 0$ be some arbitrary real number. Drawing any $x_i \sim \mathcal{U}([0, R])$ for some $i \in [d]$, we have with probability at most $\frac{2\delta}{R}$ that it is within distance of at most δ from any other x_j , $j < i$. By a union bound taken over the distances from all coordinates, any freshly sampled coordinate is within distance at least δ from all the other coordinates with probability at least $1 - \frac{2d\delta}{R}$. Taking another union bound over drawing each coordinate sufficiently far and using Lemma D.2 with the fact that $\mathbb{P}[|x_i| \leq R] = 1$ for all i , we have that

$$\mathbb{P}_{\mathbf{x} \sim \mathcal{U}([0, R]^d)} [\mathbf{x} \notin S_{\delta/R}] \leq \frac{2d^2\delta}{R},$$

which by a change of variables $\delta_0 = \delta/R$ implies

$$\mathbb{P}_{\mathbf{x} \sim \mathcal{U}([0, R]^d)} [\mathbf{x} \notin S_{\delta_0}] \leq 2d^2\delta_0.$$

It is only left to compute the δ_0 for which Eq. (1) holds. To this end, we wish to find $\delta > 0$ such that

$$2d^2\delta \leq \frac{\varepsilon}{(d+1)^2R^2} \leq \frac{\varepsilon}{(d+1)^2\mathbb{E}_{\mathbf{x}\sim\mathcal{U}([0,R]^d)}[\|\mathbf{x}\|_\infty^2]},$$

where the second inequality uses Eq. (2). Solving the above for δ , we have $\delta_0 = \Omega\left(\frac{\varepsilon}{d^4R^2}\right)$, implying a weight upper bound of $\mathcal{O}\left(\frac{d^4R^2}{\varepsilon}\right)$ and concluding the proof of the corollary. \square

C Proofs from Sec. 4

C.1 Proof of Thm. 4.2

The following proposition is key in the proof of the theorem.

Proposition C.1. *For each fixed dimension $d \geq 2$, the squared L_2 error of approximating the function $\max\{0, x_1, x_2, \dots, x_d\}$ on the unit Gaussian, using n neurons in a depth 2 σ network satisfying Assumption 4.1, with coefficients of size at most s , is at least*

$$\frac{\text{polylog}(n)}{\text{polylog}(s)} \cdot \frac{1}{n^{1+\frac{3}{d-1}}}.$$

Before we prove the proposition, however, we will first need the following definition and lemmas.

Definition C.2. *Let $M(\mathbf{x}) = \max\{0, x_1, x_2, \dots, x_d\}$ denote the max function. Let $M_1(\mathbf{x})$ denote only the portion of the max function where coordinate x_1 is biggest, namely, the function that takes value x_1 if x_1 is the largest of $0, x_2, x_3, \dots, x_d$, and 0 otherwise. Equivalently, let $q_1(\mathbf{x}) = x_1 \cdot \mathbb{1}_{[x_1 \geq 0]} \cdot \mathbb{1}_{[x_2 \leq 0]} \cdot \mathbb{1}_{[x_3 \leq 0]} \cdot \dots \cdot \mathbb{1}_{[x_d \leq 0]}$, namely the function taking value x_1 but only when x_1 is nonnegative and all the other coordinates are nonpositive; let $S_1(\mathbf{x})$ be the ‘‘skew’’ matrix such that $S_1 \cdot (x_1, x_2, x_3, \dots, x_d)^T = (x_1, x_2 - x_1, x_3 - x_1, \dots, x_d - x_1)^T$, namely subtracting x_1 from all the other coordinates; thus $M_1(\mathbf{x}) = q_1(S_1 \cdot \mathbf{x})$ since x_1 is at least some other coordinate x_j if and only if $x_j - x_1 \leq 0$.*

Lemma C.3. *The Fourier transform of the skew of a function is the inverse transpose skew of the Fourier transform of the function: for a function g , we have $\widehat{g(S_1 \cdot \mathbf{x})} = \widehat{g}(S_1^{-\top} \boldsymbol{\xi})$, using the standard notation $S_1^{-\top}$ to represent the matrix inverse transpose.*

Proof. Standard (linear) change of variables relation for the Fourier transform integral. \square

Lemma C.4. *Letting $q_1(\mathbf{x}) = x_1 \cdot \mathbb{1}_{[x_1 \geq 0]} \cdot \mathbb{1}_{[x_2 \leq 0]} \cdot \mathbb{1}_{[x_3 \leq 0]} \cdot \dots \cdot \mathbb{1}_{[x_d \leq 0]}$ and defining Dawson’s integral to be $\text{daw}(x) := \exp(-x^2) \int_0^x \exp(t^2) dt$, then the Fourier transform of $q_1(\mathbf{x}) \exp(-\|\mathbf{x}\|_2^2/4)$ equals*

$$(2 - 4\xi_1 \text{daw}(\xi_1) - 2i\sqrt{\pi}\xi_1 \exp(-\xi_1^2)) \prod_{j=2}^d (\exp(-\xi_j^2)\sqrt{\pi} + 2i \cdot \text{daw}(\xi_j)).$$

And further, this Fourier transform has magnitude at most $2\pi^{\frac{d-1}{2}}$ everywhere, and for inputs ξ each of whose coordinates is positive and has value at least $\Omega(\log(d))$, the Fourier transform has a component in the (complex) direction $-i^{d-1}$ that is at least $\frac{1}{\xi_1^2} \prod_{j=2}^d \frac{1}{\xi_j}$.

Proof. The Fourier transform is a straightforward calculation on each dimension separately, since the function $q_1(\mathbf{x}) \exp(-\|\mathbf{x}\|_2^2/4)$ is separable.

The global magnitude bound simply comes from evaluating the Fourier transform at the origin, since the Fourier transform of a nonnegative real function attains its largest magnitude at the origin.

For the final bound, we take the approximation of Dawson's integral $\text{daw}(x) = \frac{1}{2x} + \frac{1}{4x^3} + \Theta(\frac{1}{x^5})$ for inputs x away from 0. These inverse polynomial terms in ξ_j dominate the inverse exponential $\exp(-\xi_j^2)$ terms, even when d such terms are multiplied together, for $\xi = \Omega(\log d)$. Substituting in this approximation for $\text{daw}(x)$ into our Fourier transform expression and dropping lower-order terms yields $-\frac{1}{\xi_1^2} \prod_{j=2}^d \frac{i}{\xi_j}$, with the next-largest terms from the expansion of Dawson's integral contributing inverse-polynomially farther in the same direction. Thus we conclude the lemma. \square

Lemma C.5. For vector ξ with all of its coordinates positive and at least $\Omega(d)$, but less than some parameter b , the Fourier transform of $\exp(-\|\mathbf{x}\|_2^2/4) \max(0, x_1, x_2, \dots, x_d)$ evaluated at ξ has magnitude at least $b^{-(d+1)} 2^{-\mathcal{O}(d)}$.

Proof. The Fourier transform of $\exp(-\|\mathbf{x}\|_2^2/4) \max\{0, x_1, x_2, \dots, x_d\}$ can be decomposed into the sum of the contributions of the d separate components of the \max function, which are all symmetric up to relabeling the coordinates. We thus compute the contribution from the first component.

We compute the Fourier transform of $\exp(-\|\mathbf{x}\|_2^2/4) M_1(\mathbf{x})$ by expressing $M_1 = q_1 \circ s_1$ from Definition C.2, as the composition of a separable function q_1 with a volume-preserving affine transformation s_1 . We make further use of the transformation s_1 by breaking the scaling term $\exp(-\|\mathbf{x}\|_2^2/4)$ into 2 parts, one of which is a spherical Gaussian even after begin transformed by s_1 . Explicitly, we have

$$\exp(-\|\mathbf{x}\|_2^2/4) M_1(\mathbf{x}) = \exp(-\mathbf{x}^\top Q \mathbf{x}) \exp(-\|s_1(\mathbf{x})\|_2^2/(4(d+1))) q_1(s_1(\mathbf{x})) \quad (3)$$

for some symmetric positive semidefinite matrix Q with eigenvalues at most $\frac{1}{4}$.

Since the Fourier transform of a product equals the convolution of the Fourier transform of the terms, we thus have that the Fourier transform of Equation 3 equals the convolution of the Fourier transform of the Gaussian $\exp(-\mathbf{x}^\top Q \mathbf{x})$ with the Fourier transform of the expression $\exp(-\|s_1(\mathbf{x})\|_2^2/(4(d+1))) q_1(s_1(\mathbf{x}))$. Since this expression is an affine transformation of $\exp(-\|\mathbf{y}\|_2^2/(4(d+1))) q_1(\mathbf{y})$, its Fourier transform is the corresponding (inverse transpose) affine transformation of the Fourier transform of $\exp(-\|\mathbf{y}\|_2^2/(4(d+1))) q_1(\mathbf{y})$.

We bound this Fourier transform via Lemma C.4. Explicitly, let $g(\xi)$ be the Fourier transform of $\exp(-\|\mathbf{y}\|_2^2/4) q_1(\mathbf{y})$, which Lemma C.4 bounds. Then the Fourier transform of $\exp(-\|\mathbf{y}\|_2^2/(4(d+1))) q_1(\mathbf{y})$ is exactly $g(\xi \sqrt{d+1})(d+1)^{\frac{d+1}{2}}$, since replacing \mathbf{y} by $\mathbf{y}\sqrt{d+1}$ scales the function q_1 by $\sqrt{d+1}$ and thus scales its integral and hence Fourier transform by the $d+1$ power of this, as claimed. Next, transforming the inputs of a function by the affine function s_1 transforms its Fourier transform by the transpose of the inverse of the affine function. Thus the

Fourier transform of $\exp(-\|s_1(\mathbf{x})\|_2^2/(4(d+1))) q_1(s_1(\mathbf{x}))$ equals $g(\sqrt{d+1}((\sum_j \xi_j), \xi_2, \xi_3, \dots, \xi_d)) \cdot (d+1)^{\frac{d+1}{2}}$.

We now use the bounds of Lemma C.4 to characterize g . For ξ with all coordinates positive and at least $\Omega(1)$, the transformed coordinates $\sqrt{d+1}((\sum_j \xi_j), \xi_2, \xi_3, \dots, \xi_d)$ will all be at least $\Omega(\log(d))$ and thus the lemma applies. Thus we conclude that the Fourier transform has component in the direction $-i^{d-1}$ at least $\frac{1}{(\sum_j \xi_j)^2} \prod_{j=2}^d \frac{1}{\xi_j}$, where the factors of $d+1$ all cancel; by the second part of Lemma C.4, this Fourier transform has magnitude at most $\mathcal{O}(d)^{\mathcal{O}(d)}$ everywhere.

Finally, to obtain the overall Fourier transform of $\exp(-\|\mathbf{x}\|_2^2/4) M_1(\mathbf{x})$, it remains to convolve this last expression with the Fourier transform of the remaining term $\exp(-\mathbf{x}^\top Q \mathbf{x})$; namely, we convolve this last expression with the Gaussian of covariance Q , which thus has radius $\leq \frac{1}{4}$ by construction. Since all but $\exp(-\Omega(d^2))$ fraction of the mass of this Gaussian must be within radius $\mathcal{O}(d)$, we thus have that—even after this final convolution—the component of the Fourier transform of $\exp(-\|\mathbf{x}\|_2^2/4) M_1(\mathbf{x})$ in the direction $-i^{d-1}$ must be at least $2^{-\mathcal{O}(d)} \frac{1}{(\sum_j \xi_j)^2} \prod_{j=2}^d \frac{1}{\xi_j}$ provided all coordinates of ξ are at least $\Omega(d)$.

Summing this bound over all d components of the maximum function, and then pointing out that the magnitude of a complex number must be at least its component in the direction $-i^{d-1}$ yields our final bound. \square

Using the above lemmas, we now turn to the proof of the proposition.

Proof of Proposition C.1. Let $b = \omega(d)$. Consider the region in Fourier space where each coordinate ξ_j lies in $[\Omega(d), b]$. This region has volume $\Omega(b^d)$. By Lemma C.5, the Fourier transform of the maximum function, weighted by the square root of the pdf of the unit Gaussian, has magnitude at least $b^{-(d+1)} 2^{-\mathcal{O}(d)}$; denote this bound by ℓ .

However, a ReLU with bounded coefficients has a Fourier transform which is large only on a relatively small volume, which is what gives us the desired contradiction.

We will show that the linear combination of ReLU units cannot closely approximate the max function on our Gaussian via the following “accounting” scheme: consider the contribution of each neuron separately, and letting $f_k(\xi)$ be the Fourier transform of the contribution of the k^{th} neuron, then we give this neuron a “score” of $\int_{[\Omega(d), b]^d} \min\{1, \frac{1}{\ell} |f_k(\xi)|\} d\xi$. We will show that the total score over all neurons is at most half the volume of the region, which implies a squared- L_2 error of at least $\Omega(b^d \ell^2) = b^{-d-2} 2^{-\mathcal{O}(d)}$.

Consider a neuron with an activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ satisfying Assumption 4.1, and weights of magnitude at most some bound s . In the context of the neural net, σ will be applied as $w\sigma(\mathbf{x} \cdot \mathbf{v})$ where w is a weight of magnitude at most s and \mathbf{v} is a vector each of whose coordinates has magnitude at most s . We decompose the Fourier transform of $w\sigma(\mathbf{x} \cdot \mathbf{v}) \exp(-\|\mathbf{x}\|_2^2/4)$ into the Fourier transform along the direction of \mathbf{v} , and then the Fourier transform in the transverse $d-1$ dimensional space.

Since by Assumption 4.1, σ is polynomially bounded, the 1-dimensional Fourier transform along direction \mathbf{v} of this scaled version $w\sigma(\mathbf{x} \cdot \mathbf{v}) \exp(-\|\mathbf{x}\|_2^2/4)$ will be bounded by $\mathcal{O}(sd)^{d+1}$, where the parameters of the big-O expression depend on σ . We now consider the Fourier transform of along the $d-1$ dimensional space orthogonal to d : along any hyperplane orthogonal to \mathbf{v} , we

have a $d - 1$ dimensional Gaussian $\exp(-\|\mathbf{x}\|_2^2/4)$ times some number bounded by $\mathcal{O}(sd)^{d+1}$. We thus consider how much ‘‘score’’ this can contribute.

Namely, for a scaling factor $t = \mathcal{O}(sd)^{d+1}$, we bound $\int_{\mathbb{R}^{d-1}} \min \left\{ 1, \frac{t}{\ell} (2\sqrt{\pi})^{d-1} \exp(-\|\boldsymbol{\xi}\|_2^2) \right\} d\boldsymbol{\xi}$, where the expression $(2\sqrt{\pi})^{d-1} \exp(-\|\boldsymbol{\xi}\|_2^2)$ is the $d - 1$ dimensional Fourier transform of the Gaussian $\exp(-\mathbf{x}^2/4)$. This integral is straightforward to bound once we convert it to an integral over the radius. We solve for the radius r where the min function transitions from the first term to the second: $\frac{t}{\ell} (2\sqrt{\pi})^{d-1} \exp(-r^2) \leq 1$ means that $r \geq \sqrt{\log \frac{t}{\ell} (2\sqrt{\pi})^{d-1}}$. Since the surface area of a radius r ball in $d - 1$ dimensions is $\mathcal{O}(1)r^{d-1}$, we bound our integral as

$$\int_{\mathbb{R}^{d-1}} \min \left\{ 1, \frac{t}{\ell} (2\sqrt{\pi})^{d-1} \exp(-\|\boldsymbol{\xi}\|_2^2) \right\} d\boldsymbol{\xi} \leq \int_0^\infty \mathcal{O}(1)r^{d-1} \min \left\{ 1, \frac{t}{\ell} (2\sqrt{\pi})^{d-1} \exp(-r^2) \right\} dr$$

And for $\log \frac{t}{\ell} (2\sqrt{\pi})^{d-1} \geq d$ this integral is bounded by $\mathcal{O}(1) \cdot (\log \frac{t}{\ell} (2\sqrt{\pi})^{d-1})^{\frac{d}{2}}$. Substituting in our bound for t yields that, for $\ell \leq \mathcal{O}(sd)^{d+1}$, the contribution to the score per unit in the transverse direction \mathbf{v} is at most $((d + 1) \log \mathcal{O}(sd) + \log \frac{1}{\ell})^{\frac{d}{2}}$.

Since we integrate the score over the hypercube of side length b , the projection to direction \mathbf{v} has length at most $b\sqrt{d}$, and thus the total score of our neural network with n neurons is at most $nb\sqrt{d}((d + 1) \log \mathcal{O}(sd) + \log \frac{1}{\ell})^{\frac{d}{2}}$. Substituting in the definition $\ell = b^{-(d+1)}2^{-\mathcal{O}(d)}$ yields $nb\sqrt{d}((d + 1) \log \mathcal{O}(sdb))^{\frac{d}{2}}$.

As described above, we show this neural network does not closely approximate the max function by showing that this score is less than half the volume of the region of frequency space under consideration, $\frac{1}{2}(b - \Omega(d))^d$, which is true provided $b \geq n^{\frac{1}{d-1}} \frac{\text{polylog}(s)}{\text{polylog}(n)}$, where we emphasize that the base and exponents of the polylog terms may depend on d . Thus our overall L_2 -squared error bound of $b^{-d-2}2^{-\mathcal{O}(d)}$ becomes $\frac{\text{polylog}(n)}{\text{polylog}(s)} \frac{1}{n^{1+\frac{3}{d-1}}}$, as claimed. \square

With Proposition C.1 at our disposal, we can now prove Thm. 4.2.

Proof of Thm. 4.2. We wish to find some $c_2 > 0$ such that

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}([0,1]^d)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] > \Omega\left(\frac{1}{d^{c_2 \ell}}\right).$$

To this end, we first define the sets $A := [0, 1 - 1/d]^{d-3}$ and $B = [1 - 1/d, 1]^3$. For any natural n , denote the uniform distribution over $[0, 1]^n$ by \mathcal{D}_n and compute by repeatedly using the law of total expectation

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_d} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] &= \mathbb{E}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}_1, \mathbf{x}_2) - f_d(\mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_1]] \\ &= \mathbb{E}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}_1, \mathbf{x}_2) - f_d(\mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_1 \in A, \mathbf{x}_2 \in B]] \cdot \mathbb{P}[\mathbf{x}_1 \in A, \mathbf{x}_2 \in B] \\ &\quad + \mathbb{E}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}_1, \mathbf{x}_2) - f_d(\mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_1 \notin A \text{ or } \mathbf{x}_2 \notin B]] \cdot \mathbb{P}[\mathbf{x}_1 \notin A \text{ or } \mathbf{x}_2 \notin B] \\ &\geq \mathbb{E}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}_1, \mathbf{x}_2) - f_d(\mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_1 \in A, \mathbf{x}_2 \in B]] \cdot \mathbb{P}[\mathbf{x}_1 \in A, \mathbf{x}_2 \in B] \\ &\geq \mathbb{E}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}_1, \mathbf{x}_2) - f_d(\mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_1 \in A, \mathbf{x}_2 \in B]] \cdot d^{-3} \exp(-1) \\ &= \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}'_1, \mathbf{x}_2) - f_d(\mathbf{x}'_1, \mathbf{x}_2))^2 | \mathbf{x}_2 \in B] \cdot d^{-3} \exp(-1). \end{aligned}$$

In the above, the last equality holds for some intermediate point $\mathbf{x}'_1 \in A$ whose existence is guaranteed by Lemma D.5 due to the fact that $g(\mathbf{x}) = \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}, \mathbf{x}_2) - f_d(\mathbf{x}, \mathbf{x}_2))^2 | \mathbf{x}_2 \in B]$ is continuous on $[0, 1]^{d-3}$. Since $\mathbf{x}_2 \mapsto \mathcal{N}(\mathbf{x}'_1, \mathbf{x}_2)$ defines a depth 2 σ network which we denote by $\tilde{\mathcal{N}}(\cdot)$, and by using the fact that $f_d(\mathbf{x}'_1, \mathbf{x}_2) = f_3(\mathbf{x}_2)$ for all $\mathbf{x}'_1 \in A$ and $\mathbf{x}_2 \in B$, we can let $\mathcal{N}_{(2)}$ denote the class of depth 2 networks employing a σ activation function, and lower bound the above by

$$\inf_{\tilde{\mathcal{N}} \in \mathcal{N}_{(2)}} \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{U}(B)} \left[\left(\tilde{\mathcal{N}}(\mathbf{x}_2) - f_3(\mathbf{x}_2) \right)^2 \right] \cdot d^{-3} \exp(-1).$$

It now suffices to lower bound the expectation term above by $\Omega(d^{-c' \cdot \ell})$ for some constant $c' > 0$. Focusing on the expectation term above and letting $\tilde{\mathcal{N}}$ denote arbitrary (not necessarily fixed) elements in $\mathcal{N}_{(2)}$, we once more use the law of total expectation repeatedly to obtain

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{U}(B)} \left[\left(\tilde{\mathcal{N}}(\mathbf{x}_2) - f_3(\mathbf{x}_2) \right)^2 \right] \\ &= \mathbb{E}_{\tilde{x}_1 \sim \mathcal{U}([1-\frac{1}{d}, 1])} \left[\mathbb{E}_{\tilde{\mathbf{x}}_2 \sim \mathcal{U}([1-\frac{1}{d}, 1]^2)} \left[\left(\tilde{\mathcal{N}}(\tilde{x}_1, \tilde{\mathbf{x}}_2) - f_3(\tilde{x}_1, \tilde{\mathbf{x}}_2) \right)^2 | \tilde{x}_1 \in \left[1 - \frac{1}{d}, 1 \right] \right] \right] \\ &\geq 0.5 \mathbb{E}_{\tilde{x}_1 \sim \mathcal{U}([1-\frac{1}{d}, 1])} \left[\mathbb{E}_{\tilde{\mathbf{x}}_2 \sim \mathcal{U}([1-\frac{1}{d}, 1]^2)} \left[\left(\tilde{\mathcal{N}}(\tilde{x}_1, \tilde{\mathbf{x}}_2) - f_3(\tilde{x}_1, \tilde{\mathbf{x}}_2) \right)^2 | \tilde{x}_1 \in \left[1 - \frac{3}{4d}, 1 - \frac{1}{4d} \right] \right] \right] \\ &= 0.5 \mathbb{E}_{\tilde{\mathbf{x}}_2 \sim \mathcal{U}([1-\frac{1}{d}, 1]^2)} \left[\left(\tilde{\mathcal{N}}(x_0, \tilde{\mathbf{x}}_2) - f_3(x_0, \tilde{\mathbf{x}}_2) \right)^2 \right], \end{aligned}$$

where the last equality uses Lemma D.5 to establish the existence of some intermediate point $x_0 \in [1 - \frac{3}{4d}, 1 - \frac{1}{4d}]$ satisfying the above. Writing the above expectation term in integral form, we have that it equals

$$\int_{1-\frac{1}{d}}^1 \int_{1-\frac{1}{d}}^1 \left(\tilde{\mathcal{N}}(x_0, \tilde{\mathbf{x}}_2) - f_3(x_0, \tilde{\mathbf{x}}_2) \right)^2 d^2 d\tilde{\mathbf{x}}_2.$$

By the change of variables $\tilde{\mathbf{x}}_2 = \mathbf{y} + (x_0, x_0)$, $d\tilde{\mathbf{x}}_2 = d\mathbf{y}$, the above equals

$$\begin{aligned} & \int_{1-\frac{1}{d}-x_0}^{1-x_0} \int_{1-\frac{1}{d}-x_0}^{1-x_0} \left(\tilde{\mathcal{N}}(x_0, \mathbf{y} + (x_0, x_0)) - f_3(x_0, \mathbf{y} + (x_0, x_0)) \right)^2 d^2 d\mathbf{y} \\ &= \int_{1-\frac{1}{d}-x_0}^{1-x_0} \int_{1-\frac{1}{d}-x_0}^{1-x_0} \left(\tilde{\mathcal{N}}(x_0, \mathbf{y} + (x_0, x_0)) - x_0 - f_3(0, \mathbf{y}) \right)^2 d^2 d\mathbf{y} \\ &= \int_{1-\frac{1}{d}-x_0}^{1-x_0} \int_{1-\frac{1}{d}-x_0}^{1-x_0} \left(\tilde{\mathcal{N}}(\mathbf{y}) - f_3(0, \mathbf{y}) \right)^2 d^2 d\mathbf{y} \\ &\geq \int_{-\frac{1}{4d}}^{\frac{1}{4d}} \int_{-\frac{1}{4d}}^{\frac{1}{4d}} \left(\tilde{\mathcal{N}}(\mathbf{y}) - f_3(0, \mathbf{y}) \right)^2 d^2 d\mathbf{y}, \end{aligned}$$

where the first equality uses the fact that $f_3(\mathbf{x} + (c, c, c)) = c + f_3(\mathbf{x})$ for any vector \mathbf{x} and real c , the second equality follows from the fact that $\mathcal{N}_{(2)}$ is closed under linear transformations of the input

and the output, and since we can simulate the fixed input x_0 in the first coordinate by an appropriate linear rescaling of the first hidden layer, and the inequality follows from $x_0 \in [1 - \frac{3}{4d}, 1 - \frac{1}{4d}]$ which implies that $[-0.25d^{-1}, 0.25d^{-1}] \subseteq [1 - 1/d - x_0, 1 - x_0]$. Letting $\gamma > 0$ to be determined later, we perform a second linear change of variables $\mathbf{y} = 0.5\gamma^{-1}\mathbf{z}$, $d\mathbf{y} = 0.5\gamma^{-2}d\mathbf{z}$, which entails that the above displayed equation equals

$$\begin{aligned}
& 0.5 \int_{-\frac{\gamma}{2d}}^{\frac{\gamma}{2d}} \int_{-\frac{\gamma}{2d}}^{\frac{\gamma}{2d}} \left(\tilde{\mathcal{N}}(0.5\gamma^{-1}\mathbf{z}) - f_3(0, 0.5\gamma^{-1}\mathbf{z}) \right)^2 d^2\gamma^{-2}d\mathbf{z} \\
&= 0.5 \int_{-\frac{\gamma}{2d}}^{\frac{\gamma}{2d}} \int_{-\frac{\gamma}{2d}}^{\frac{\gamma}{2d}} \left(0.5\gamma^{-1}2\gamma\tilde{\mathcal{N}}(0.5\gamma^{-1}\mathbf{z}) - 0.5\gamma^{-1}f_3(0, \mathbf{z}) \right)^2 d^2\gamma^{-2}d\mathbf{z} \\
&= 0.25\pi d^2\gamma^{-2} \int_{-\frac{\gamma}{2d}}^{\frac{\gamma}{2d}} \int_{-\frac{\gamma}{2d}}^{\frac{\gamma}{2d}} \left(\tilde{\mathcal{N}}(\mathbf{z}) - f_3(0, \mathbf{z}) \right)^2 \frac{1}{2\pi} d\mathbf{z} \\
&\geq 0.25\pi d^2\gamma^{-2} \int_{\{\mathbf{z}:\|\mathbf{z}\|_2 \leq 0.5d^{-1}\gamma\}} \left(\tilde{\mathcal{N}}(\mathbf{z}) - f_3(0, \mathbf{z}) \right)^2 \frac{1}{2\pi} \exp(-0.5\|\mathbf{z}\|_2^2) d\mathbf{z} \\
&= 0.25\pi d^2\gamma^{-2} \left(\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, I_2)} \left[\left(\tilde{\mathcal{N}}(\mathbf{z}) - f_3(0, \mathbf{z}) \right)^2 \right] \right. \\
&\quad \left. - \int_{\{\mathbf{z}:\|\mathbf{z}\|_2 \geq 0.5d^{-1}\gamma\}} \left(\tilde{\mathcal{N}}(\mathbf{z}) - f_3(0, \mathbf{z}) \right)^2 \frac{1}{2\pi} \exp(-0.5\|\mathbf{z}\|_2^2) d\mathbf{z} \right). \tag{4}
\end{aligned}$$

In the above, the first equality follows from the fact that $f_3(\alpha \cdot \mathbf{x}) = \alpha f_3(\mathbf{x})$ for all $\alpha > 0$ and $\mathbf{x} \in \mathbb{R}$, the second equality follows from the fact that $\mathcal{N}_{(2)}$ is closed under linear scaling of its input and output, and the inequality follows from $\{\mathbf{z} : \|\mathbf{z}\|_2 \leq 0.5d^{-1}\gamma\} \subseteq [-0.5d^{-1}\gamma, 0.5d^{-1}\gamma]^2$ and the fact that the maximum of a bivariate standard Gaussian is $\frac{1}{2\pi}$.

Next, we upper bound the square of the above approximation. We begin with the output of a single neuron:

$$\begin{aligned}
|\sigma(\langle \mathbf{w}_i, \mathbf{z} \rangle + b_i)| &\leq C_\sigma (1 + |\langle \mathbf{w}_i, \mathbf{z} \rangle + b_i|^{\alpha_\sigma}) \\
&\leq C_\sigma (1 + \|\mathbf{w}_i\| \cdot \|\mathbf{z}\| + |b_i|^{\alpha_\sigma}) \leq \mathcal{O}(\exp(\mathcal{O}(d)) \|\mathbf{z}\|^{\alpha_\sigma}),
\end{aligned}$$

where the second inequality follows from Cauchy-Schwartz and the last inequality follows from our assumption on the magnitude of the weights. Using the above, we can upper bound the output of the network by

$$\left| \tilde{\mathcal{N}}(\mathbf{z}) \right| = \left| \sum_{i=1}^k \sigma(\langle \mathbf{w}_i, \mathbf{z} \rangle + b_i) + b_0 \right| \leq d^\ell \cdot \mathcal{O}(\exp(\mathcal{O}(d)) \|\mathbf{z}\|^{\alpha_\sigma}),$$

implying

$$\left(\tilde{\mathcal{N}}(\mathbf{z}) - f_3(0, \mathbf{z}) \right)^2 \leq d^{2\ell} \cdot \mathcal{O}(\exp(\mathcal{O}(d)) \|\mathbf{z}\|^{2\alpha_\sigma}).$$

Substituting $\gamma = \mathcal{O}(d^2)$ in Eq. (4) and using the above, we get a lower bound of

$$0.25\pi d^{-2} \left(\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, I_2)} \left[\left(\tilde{\mathcal{N}}(\mathbf{z}) - f_3(0, \mathbf{z}) \right)^2 \right] - \mathcal{O}(\exp(-0.5d^2)) \right).$$

Lastly, using Proposition C.1 to lower bound the expectation term above with the assumed bounds on the parameters, the theorem follows. \square

C.2 Proof of Thm. 4.3

Proof. We begin with reducing the approximation error of f_d over a depth 3 ReLU network to the approximation error of f_3 over a depth 2 ReLU network. To this end, we first identify three coordinates in the domain of \mathcal{N} where a certain sub-cube B of dimension 3, and a set $A \subseteq [0, 1]^{d-3}$ exist, which satisfy the following properties:

1. $\mathbb{P}_{\mathbf{x} \sim \mathcal{D}_{d-3}} [\mathbf{x} \in A] \geq 0.1$.
2. $\mathbb{P}_{\mathbf{x} \sim \mathcal{D}_3} [\mathbf{x} \in B] = d^{-18}$.
3. $f_d(\mathbf{x}_1, \mathbf{x}_2) = f_3(\mathbf{x}_2)$ for all $\mathbf{x}_1 \in A, \mathbf{x}_2 \in B$.
4. A can be decomposed into a disjoint partition of at most 2^k convex sets A_1, A_2, \dots and a set of measure zero Δ , such that $A = \bigcup_j A_j \cup \Delta$, where for all j and $i \in [k]$, $\mathbb{P}_{\mathbf{x} \sim \mathcal{D}_{d-3}} [\mathbf{x} \in A_j] > 0$ and $\text{sign}(n_i(\mathbf{x}_1, \mathbf{x}_2))$ is fixed for all $\mathbf{x}_1 \in A_j$ and all $\mathbf{x}_2 \in B$.

Before we prove the existence of A and B , we shall first show how they imply a reduction to an approximation using depth 2. For any natural n , denote the uniform distribution over $[0, 1]^n$ by \mathcal{D}_n and compute by repeatedly using the law of total expectation

$$\begin{aligned}
\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_d} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] &= \mathbb{E}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}_1, \mathbf{x}_2) - f_d(\mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_1]] \\
&= \mathbb{E}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}_1, \mathbf{x}_2) - f_d(\mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_1 \in A, \mathbf{x}_2 \in B]] \cdot \mathbb{P}[\mathbf{x}_1 \in A, \mathbf{x}_2 \in B] \\
&\quad + \mathbb{E}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}_1, \mathbf{x}_2) - f_d(\mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_1 \notin A \text{ or } \mathbf{x}_2 \notin B]] \cdot \mathbb{P}[\mathbf{x}_1 \notin A \text{ or } \mathbf{x}_2 \notin B] \\
&\geq \mathbb{E}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}_1, \mathbf{x}_2) - f_d(\mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_1 \in A, \mathbf{x}_2 \in B]] \cdot \mathbb{P}[\mathbf{x}_1 \in A, \mathbf{x}_2 \in B] \\
&= \sum_j \mathbb{E}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}_1, \mathbf{x}_2) - f_d(\mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_2 \in B, \mathbf{x}_1] | \mathbf{x}_1 \in A_j] \cdot \mathbb{P}[\mathbf{x}_1 \in A_j, \mathbf{x}_2 \in B] \\
&= \sum_j \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}'_j, \mathbf{x}_2) - f_d(\mathbf{x}'_j, \mathbf{x}_2))^2 | \mathbf{x}_2 \in B] \cdot \mathbb{P}[\mathbf{x}_1 \in A_j, \mathbf{x}_2 \in B].
\end{aligned}$$

In the above, the penultimate equality holds despite the omission of Δ from the decomposition of A since it is a set of measure zero, and the last equality holds for a set of intermediate points $(\mathbf{x}'_1, \mathbf{x}'_2, \dots) \in \mathbb{R}^{d-3}$ whose existence is guaranteed by Lemma D.5 due to Item 4 and the fact that $g(\mathbf{x}) = \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}(\mathbf{x}, \mathbf{x}_2) - f_d(\mathbf{x}, \mathbf{x}_2))^2 | \mathbf{x}_2 \in B]$ is continuous on $[0, 1]^{d-3}$. Since for any given j , $\text{sign}(n_i(\mathbf{x}'_j, \mathbf{x}_2))$ is fixed for all $i \in [k]$ and all $\mathbf{x}_2 \in B$, we can collapse the first hidden layer of \mathcal{N} ,² obtaining a depth 2 ReLU network \mathcal{N}_j for each \mathbf{x}'_j such that $\mathcal{N}(\mathbf{x}'_j, \mathbf{x}_2) = \mathcal{N}_j(\mathbf{x}_2)$ for all

²More formally, we can obtain \mathcal{N}_j from \mathcal{N} and \mathbf{x}'_j by choosing some arbitrary $\mathbf{x}_2 \in B$ and considering the sign of $n_i(\mathbf{x}'_j, \mathbf{x}_2)$ for each $i \in [k]$. If it is negative we can discard the neuron and set its incoming weight in the second layer to 0, and if it is positive then we discard the ReLU activation and compose the obtained linear transformation with the linear transformation computed by the corresponding neuron in the second layer. In both cases, the neuron in the first layer is either canceled or is absorbed into the second layer, thus removing the first hidden layer altogether without increasing the width of the network.

$\mathbf{x}_2 \in B$. Combining the previous argument with Item 3, the above displayed equation is equal to

$$\begin{aligned} & \sum_j \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{D}_3} [(\mathcal{N}_j(\mathbf{x}_2) - f_3(\mathbf{x}_2))^2 | \mathbf{x}_2 \in B] \cdot \mathbb{P}[\mathbf{x}_1 \in A_j, \mathbf{x}_2 \in B] \\ &= \sum_j \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{U}(B)} [(\mathcal{N}_j(\mathbf{x}_2) - f_3(\mathbf{x}_2))^2] \cdot \mathbb{P}[\mathbf{x}_1 \in A_j] \cdot \mathbb{P}[\mathbf{x}_2 \in B]. \end{aligned}$$

Letting $\mathcal{N}_{(2)}$ denote the class of depth 2 ReLU networks of width at most $d^2/5$, and letting $b > a \geq 0$ such that $B = [a, b]^3$ where $b := a + d^{-6}$, we can lower bound the above by

$$\begin{aligned} & \sum_j \inf_{\tilde{\mathcal{N}} \in \mathcal{N}_{(2)}} \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{U}(B)} \left[\left(\tilde{\mathcal{N}}(\mathbf{x}_2) - f_3(\mathbf{x}_2) \right)^2 \right] \cdot \mathbb{P}[\mathbf{x}_1 \in A_j] \cdot \mathbb{P}[\mathbf{x}_2 \in B] \\ &= \inf_{\tilde{\mathcal{N}} \in \mathcal{N}_{(2)}} \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{U}(B)} \left[\left(\tilde{\mathcal{N}}(\mathbf{x}_2) - f_3(\mathbf{x}_2) \right)^2 \right] \cdot \mathbb{P}[\mathbf{x}_1 \in A] \cdot \mathbb{P}[\mathbf{x}_2 \in B] \\ &\geq 0.1d^{-18} \cdot \inf_{\tilde{\mathcal{N}} \in \mathcal{N}_{(2)}} \mathbb{E}_{\mathbf{x}_2 \sim \mathcal{U}(B)} \left[\left(\tilde{\mathcal{N}}(\mathbf{x}_2) - f_3(\mathbf{x}_2) \right)^2 \right], \end{aligned} \quad (5)$$

where the equality is due to Item 4 and the inequality is due to Items 1 and 2. Applying Lemma D.6 and Thm. 4.2 with input dimension 3 and $\ell = 2$, the lower bound follows.

It now remains to show the existence of A and B . Starting with B , we first assume w.l.o.g. that no neuron in the first hidden layer of \mathcal{N} has an all-zero weight vector. This is justified since if such a neuron exists, it merely outputs a constant as input to the second layer which can be simulated by modifying the bias terms in the second layer, which doesn't increase the width of \mathcal{N} . Denote for all $i \in [k]$, $w_i^{\max} = w_{i, j_i}$ where $j_i = \operatorname{argmax}_{j \in [d]} |w_{i, j}|$, we define the set

$$P := \left\{ x \in \left[1 - \frac{1}{d}, 1 \right] : \left| x + \frac{b_i}{w_i^{\max}} \right| \leq d^{-3}, \quad \forall i \in [k] \right\}.$$

Note that by our assumption $k \leq \frac{d^2}{5}$, we have that P consists of at most $\frac{d^2}{5}$ connected components where each is of length at most $2d^{-3}$. Therefore, the overall length of P is no more than $\frac{2}{5d}$, and we can thus find an interval $[a, a + d^{-6}] \subseteq [1 - 1/d] \setminus P$ for some $a \in [1 - 1/d, 1 - d^{-6}]$. We can now define

$$B := [a, a + d^{-6}]^3.$$

Note that this immediately entails that $\mathbb{P}_{\mathbf{x} \sim \mathcal{D}_3}[\mathbf{x} \in B] = d^{-18}$, proving Item 2. Continuing to showing the existence of A , we first define it formally as the set given by

$$A := \left[0, 1 - \frac{1}{d} \right]^{d-3} \setminus \left\{ \mathbf{x}_1 \in [0, 1]^{d-3} : \exists \mathbf{x}_2 \in B, i \in [k] \text{ s.t. } n_i(\mathbf{x}_1, \mathbf{x}_2) = 0 \right\}.$$

Note that this immediately implies Item 3. To show Item 4, we observe that A can be defined as the set difference between a cube and the union of k closed sets, one for each neuron in the first hidden

layer of \mathcal{N} . More specifically, for $i \in [k]$, suppose that the i -th neuron has weights $\mathbf{w} = (\mathbf{w}_1, \mathbf{w}_2)$ and bias b where $\mathbf{w}_1 \in \mathbb{R}^{d-3}$ and $\mathbf{w}_2 \in \mathbb{R}^3$, and consider the set

$$\tilde{A}_i := \{\mathbf{x}_1 \in [0, 1]^{d-3} : \exists \mathbf{x}_2 \in B \text{ s.t. } n_i(\mathbf{x}_1, \mathbf{x}_2) = 0\}.$$

Observing that

$$\tilde{A}_i = \bigcup_{\mathbf{x}_2 \in B} \{\mathbf{x}_1 : \langle \mathbf{w}_1, \mathbf{x}_1 \rangle = b - \langle \mathbf{w}_2, \mathbf{x}_2 \rangle\},$$

we have that $\tilde{A}_i = \mathbb{R}^{d-3}$ or $\tilde{A}_i = \emptyset$ if $\mathbf{w}_1 = \mathbf{0}$, depending on whether $b - \langle \mathbf{w}_2, \mathbf{x}_2 \rangle$ equals zero for some $\mathbf{w}_2 \in B$ or not. Otherwise, if $\mathbf{w}_1 \neq \mathbf{0}$, we have that \tilde{A}_i can be represented as a union of parallel affine subspaces, the boundary of which is given by $\langle \mathbf{w}_1, \mathbf{x}_1 \rangle = \min_{\mathbf{x}_2 \in B} b - \langle \mathbf{w}_2, \mathbf{x}_2 \rangle$ and $\langle \mathbf{w}_1, \mathbf{x}_1 \rangle = \max_{\mathbf{x}_2 \in B} b - \langle \mathbf{w}_2, \mathbf{x}_2 \rangle$, where both the minimum and maximum are defined since B is compact and $\mathbf{x}_2 \mapsto b - \langle \mathbf{w}_2, \mathbf{x}_2 \rangle$ is continuous. Moreover, by continuity we also obtain that \tilde{A}_i is connected. We thus have in either case that we can represent $\overline{\tilde{A}_i} = \tilde{A}_{i,1} \cup \tilde{A}_{i,2}$ for some disjoint and convex sets $\tilde{A}_{i,1}, \tilde{A}_{i,2} \subseteq \mathbb{R}^{d-3}$. We now compute

$$\begin{aligned} A &= \left[0, 1 - \frac{1}{d}\right]^{d-3} \setminus \left(\bigcup_{i \in [k]} \tilde{A}_i\right) = \left[0, 1 - \frac{1}{d}\right]^{d-3} \cap \overline{\bigcup_{i \in [k]} \tilde{A}_i} = \left[0, 1 - \frac{1}{d}\right]^{d-3} \cap \left(\bigcap_{i \in [k]} \overline{\tilde{A}_i}\right) \\ &= \left[0, 1 - \frac{1}{d}\right]^{d-3} \cap \left(\bigcap_{i \in [k]} (\tilde{A}_{i,1} \cup \tilde{A}_{i,2})\right) = \bigcup_{j_1, \dots, j_k \in \{1, 2\}} \left(\left(\bigcap_{i \in [k]} \tilde{A}_{i, j_i}\right) \cap \left[0, 1 - \frac{1}{d}\right]^{d-3}\right). \end{aligned}$$

Namely, A is a union of at most 2^k disjoint and convex sets (since the intersection of convex sets is also convex). For $j_1, \dots, j_k \in \{1, 2\}$, denote

$$A_{j_1, \dots, j_k} := \left(\bigcap_{i \in [k]} \tilde{A}_{i, j_i}\right) \cap \left[0, 1 - \frac{1}{d}\right]^{d-3}.$$

By defining

$$\Delta := \left\{ \bigcup_{j_1, \dots, j_k \in \{1, 2\}} A_{j_1, \dots, j_k} : \mathbb{P}_{\mathbf{x} \sim \mathcal{D}_{d-3}} [\mathbf{x} \in A_{j_1, \dots, j_k}] = 0 \right\},$$

we can partition A into Δ and $A \setminus \Delta$ where each connected component in $A \setminus \Delta$ is not a measure zero set. Finally, we have that $\text{sign}(n_i(\mathbf{x}_1, \mathbf{x}_2))$ is fixed on each such convex component. This holds true since if otherwise, by contradiction, it holds that some neuron satisfies $n_i(\mathbf{x}_1, \mathbf{x}_2) \geq 0$ and $n_i(\mathbf{x}'_1, \mathbf{x}'_2) < 0$ for some $(\mathbf{x}_1, \mathbf{x}_2), (\mathbf{x}'_1, \mathbf{x}'_2) \in A_j \times B$ and some $j \in [2^k]$. Consider the path $p : [0, 1] \rightarrow \mathbb{R}^d$ given by

$$p(\lambda) := \lambda(\mathbf{x}_1, \mathbf{x}_2) + (1 - \lambda)(\mathbf{x}'_1, \mathbf{x}'_2).$$

Since A_j and B are convex, so is $A_j \times B$, and we have $p(\lambda) \in A_j \times B$ for all $\lambda \in [0, 1]$. Since $n_i(p(\lambda))$ is continuous in λ , by the intermediate value theorem, we can find some $\lambda_0 \in [0, 1]$ such that $n_i(p(\lambda_0)) = 0$, which contradicts the definition of A .

To show Item 1, we first show that with high probability over drawing \mathbf{x}_1 , it holds that $n_i(\mathbf{x}_1, \mathbf{x}_2) \neq 0$ for all $\mathbf{x}_2 \in B$ and $i \in [k]$. Suppose that $\mathbf{x}_1 \sim \mathcal{U}([0, 1 - 1/d]^{d-3})$. We now construct the following graph G : G has d vertices, one for each coordinate of the input dimension, and the set of edges is determined according to the values of the weights in the first hidden layer of \mathcal{N} . Specifically, there's an edge between two vertices $j_1, j_2 \in [d]$ if and only if there exists no neuron $m \in [k]$ such that $|w_{m,j_1}| > |w_{m,\ell}|$ for all $\ell \in [k] \setminus \{j_2\}$ and $|w_{m,j_2}| > |w_{m,\ell}|$ for all $\ell \in [k] \setminus \{j_1\}$. In words, there's no edge between vertices j_1 and j_2 if and only if there exists a neuron in which the coordinates j_1, j_2 have the strictly largest weights in the neuron in absolute value. Since $k \leq d^2/5$, we have that G must contain at least $\binom{d}{2} - \frac{d^2}{5}$ edges, which is strictly greater than $\frac{d^2}{4}$ for sufficiently large d , and thus by Mantel's theorem (Thm. D.1) we have that G must contain a triangle. Consider this triangle in G , and assume w.l.o.g. that it is formed on the last three coordinates. This means by the definition of G that at least one of the two largest coordinates in each neuron in the hidden layer of \mathcal{N} have an index $j \leq d-3$. Fix some neuron, and let $\mathbf{w} = (w_1, \dots, w_d)$ and b denote the weights and bias, respectively, of the neuron. Further assume w.l.o.g. that $|w_1| \geq |w_2| \geq \dots \geq |w_{d-3}|$ and that $|w_{d-2}| \leq |w_{d-1}| \leq |w_d|$. We now perform a case analysis depending on the ratio between the two largest coordinates in the weights of the neuron.

- Suppose that $|w_1| \leq |w_d|d^{-4}$. Note that this also entails $|w_d| \geq |w_1|d^4 > |w_1|$. Namely, w_d has the largest magnitude in absolute value among the weights of the neuron. By our construction of B , we have

$$\left| x + \frac{b}{w_d} \right| > d^{-3}, \quad \forall x \in [a, a + d^{-6}].$$

Simple algebra and the above imply that

$$w_d x + b \notin [-|w_d|d^{-3}, |w_d|d^{-3}], \quad \forall x \in [a, a + d^{-6}]. \quad (6)$$

We now compute

$$\left| \sum_{j=1}^{d-1} w_j x_j \right| \leq \sum_{j=1}^{d-1} |w_j| \leq \sum_{j=1}^{d-1} |w_1| < |w_d|d^{-3}.$$

In the above, the first inequality is Hölder's inequality, the second inequality is due to the fact that $|w_1|$ is the second largest in absolute value among the weights of the neuron, and the last inequality is by our assumption $|w_1| \leq |w_d|d^{-4}$. Adding the above displayed inequality with Eq. (6) we obtain

$$\sum_{j=1}^d w_j x_j + b \neq 0, \quad \forall \mathbf{x}_1 \in [0, 1 - 1/d]^{d-3}, \quad \forall \mathbf{x}_2 \in B. \quad (7)$$

- Suppose that $|w_1| > |w_d|d^{-4}$. Then in such a case, we cannot guarantee that Eq. (7) holds with probability 1. However, we can show that the randomness over drawing x_1 induces sufficient variance and therefore it holds with high probability. Define the random variables

$X := \sum_{j=1}^{d-3} w_j x_j + b$, $\tilde{X} := \sum_{j=2}^{d-3} w_j x_j + b$, where the randomness is taken over drawing $\mathbf{x}_1 \sim \mathcal{U}([0, 1 - 1/d]^{d-3})$, and let $I \subseteq \mathbb{R}$ be any interval of length $3|w_d|d^{-6}$. We compute

$$\begin{aligned} \mathbb{P} \left[X \in I | \tilde{X} = x \right] &= \mathbb{P} \left[w_1 x_1 + x \in I | \tilde{X} = x \right] \leq \int_{\mathbb{R}} \frac{1 - 1/d}{|w_1|} \mathbb{1} \{t + x \in I\} dt \\ &\leq \frac{3|w_d|d^{-6}}{|w_1|} < 3d^{-2}, \end{aligned} \quad (8)$$

where we used the fact that the density of the random variable $w_1 x_1$ is $\frac{1-1/d}{|w_1|}$ in its support, and our assumption which implies that $|w_d| < |w_1|d^4$. Next, compute using the law of total probability to obtain

$$\mathbb{P} [X \in I] = \int_{\mathbb{R}} \mathbb{P} [X \in I | \tilde{X} = x] p_{\tilde{X}}(x) dx < \int_{\mathbb{R}} 3d^{-2} p_{\tilde{X}}(x) dx = 3d^{-2},$$

where the inequality follows from Eq. (8). Since

$$\max_{\mathbf{x}_2 \in B} \sum_{i=d-2}^d w_i x_i - \min_{\mathbf{x}_2 \in B} \sum_{i=d-2}^d w_i x_i \leq 3|w_d|d^{-6},$$

we have that

$$\mathbb{P}_{\mathbf{x}_1 \sim \mathcal{U}([0, 1-1/d]^{d-3})} \left[\sum_{j=1}^d w_j x_j + b \neq 0 \right] \geq 1 - 3d^{-2}, \quad \forall \mathbf{x}_2 \in B. \quad (9)$$

Having shown that in both cases Eq. (9) holds, we proceed by taking a union bound over all the $k \leq d^2/5$ neurons in the first hidden layer of \mathcal{N} , obtaining

$$\mathbb{P}_{\mathbf{x}_1 \sim \mathcal{U}([0, 1-1/d]^{d-3})} \left[\sum_{j=1}^d w_{i,j} x_j + b \neq 0 \right] \geq \frac{2}{5}, \quad \forall i \in [k], \quad \forall \mathbf{x}_2 \in B.$$

We now observe that

$$\begin{aligned} &\mathbb{P}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbf{x}_1 \in A] \\ &= \mathbb{P}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbf{x}_1 \in A | \mathbf{x}_1 \in [0, 1 - 1/d]^{d-3}] \cdot \mathbb{P}_{\mathbf{x}_1 \sim \mathcal{D}_{d-3}} [\mathbf{x}_1 \in [0, 1 - 1/d]^{d-3}] \\ &\geq \frac{2}{5} \left(1 - \frac{1}{d}\right)^{d-3} \geq \frac{2}{5} \exp(-1) \geq 0.1, \end{aligned}$$

which thus proves Item 1, and completes the proof of the theorem. \square

C.3 Proof of Thm. 4.4

Proof. Consider the matrix of first hidden layer weights $W \in \mathbb{R}^{k \times d}$. Since $k \leq d - 1$ by our assumption, we have that $\dim(\ker(W)) \geq 1$. Fix some vector $\mathbf{v} = (v_1, \dots, v_d) \in \ker(W)$ such that $\|\mathbf{v}\|_2 = 1$ and assume w.l.o.g. $\|\mathbf{v}\|_\infty = v_1$. Denote $X := [0, 1]^d$, we now consider the triangular matrix and vector

$$P := \begin{pmatrix} \frac{1}{d}v_1 & 0 & 0 & \cdots & 0 \\ \frac{1}{d}v_2 & 1 - \frac{2}{d} & 0 & \cdots & 0 \\ \frac{1}{d}v_3 & 0 & 1 - \frac{2}{d} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{d}v_d & 0 & 0 & \cdots & 1 - \frac{2}{d} \end{pmatrix}, \quad \mathbf{b} := \begin{pmatrix} 1 - \frac{1}{d} \\ \frac{1}{d} \\ \frac{1}{d} \\ \vdots \\ \frac{1}{d} \end{pmatrix},$$

and the set defined by

$$\mathcal{P} := \{P\mathbf{x} + \mathbf{b} : \mathbf{x} \in X\}.$$

By its definition, \mathcal{P} is a parallelotope satisfying $\mathcal{P} \subseteq X$. Moreover, we have

$$f_d(\mathbf{u}) = u_1, \quad \forall \mathbf{u} = (u_1, \dots, u_d) \in \mathcal{P}. \quad (10)$$

The above holds true since for all $\mathbf{u} \in \mathcal{P}$ there exists some $\mathbf{x} = (x_1, \dots, x_d) \in X$ such that $u_i = 1 - \frac{1}{d} + \frac{1}{d}v_i x_1$ for all $i \geq 2$ and $u_1 = 1 - \frac{1}{d} + \frac{1}{d}v_1$, and thus by our assumption that $\|\mathbf{v}\|_\infty = v_1$ we have

$$u_1 = 1 - \frac{1}{d} + \frac{1}{d}v_1 x_1 \geq 1 - \frac{1}{d} + \frac{1}{d}v_i x_1 = u_i.$$

Using the change of variables $\mathbf{u} = P\mathbf{x} + \mathbf{b}$, $d\mathbf{u} = |\det(P)| d\mathbf{x}$; the fact that $\mathcal{P} \subseteq X$; and the fact that the squared loss is non-negative, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{u} \sim \mathcal{U}(X)} [(\mathcal{N}(\mathbf{u}) - f_d(\mathbf{u}))^2] &\geq \int_{\mathcal{P}} (\mathcal{N}(\mathbf{u}) - f_d(\mathbf{u}))^2 d\mathbf{u} \\ &= \int_X (\mathcal{N}(P\mathbf{x} + \mathbf{b}) - f_d(P\mathbf{x} + \mathbf{b}))^2 |\det(P)| d\mathbf{x}. \end{aligned} \quad (11)$$

Letting \mathbf{e}_i denote the standard unit vector with coordinate $e_i = 1$, we get from $P\mathbf{x} = \frac{1}{d}\mathbf{v}x_1 + \sum_{i=2}^d (1 - \frac{2}{d})x_i \mathbf{e}_i$ and $\mathbf{v} \in \ker(W)$ that we can write $\mathcal{N}(P\mathbf{x} + \mathbf{b}) = c(x_2, \dots, x_d)$ for some function $c : \mathbb{R}^{d-1} \rightarrow \mathbb{R}$. Since P is triangular, we have $|\det(P)| = \frac{1}{d} (1 - \frac{2}{d})^{d-1} v_1 \geq \frac{1}{10d} v_1$. Moreover, since $\|\mathbf{v}\|_\infty = v_1$ and $\|\mathbf{v}\|_2 = 1$, we have that $v_1 \geq d^{-0.5}$ and we can further lower bound the above to obtain $|\det(P)| \geq \frac{1}{10d^{1.5}}$. Plugging the above and Eq. (10) back in Eq. (11), we obtain

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{U}(X)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] &\geq \frac{1}{10d^{1.5}} \int_X \left(c(x_2, \dots, x_d) - \left(1 - \frac{1}{d} + \frac{1}{d}v_1 x_1 \right) \right)^2 d\mathbf{x} \\ &= \frac{1}{10d^{1.5}} \int_{x_d} \cdots \int_{x_2} \int_{x_1} \left(c(x_2, \dots, x_d) - \left(1 - \frac{1}{d} + \frac{1}{d}v_1 x_1 \right) \right)^2 dx_1 dx_2 \cdots dx_d. \end{aligned} \quad (12)$$

It is easy to verify that the optimal constant approximation for the linear function $1 - \frac{1}{d} + \frac{1}{d}v_1x_1$ is $1 - \frac{1}{d} + \frac{1}{2d}v_1$, in which case the optimal L_2 approximation error is

$$\int_0^1 \left(1 - \frac{1}{d} + \frac{1}{2d}v_1 - \left(1 - \frac{1}{d} + \frac{1}{d}v_1x\right)\right)^2 dx = \frac{v_1^2}{d^2} \int_0^1 \left(\frac{1}{2} - x\right)^2 dx = \frac{v_1^2}{12d^2}.$$

Plugging the above back in Eq. (12) and using the fact that $v_1 \geq d^{-0.5}$ again, we obtain

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}(X)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] \geq \frac{1}{10d^{1.5}} \int_{x_d} \dots \int_{x_2} \frac{v_1^2}{12d^2} dx_2 \dots dx_d \geq \frac{1}{120d^{4.5}},$$

concluding the proof of the lemma. \square

D Technical lemmas

The following theorem is a well-known result in graph theory, which we state here for the sake of completeness.

Theorem D.1 (Mantel's theorem). *Let G be a graph with d vertices and more than $d^2/4$ edges. Then G contains a triangle.*

Lemma D.2. *Let $\mathbf{x} = (x_1, \dots, x_d)$ and $M, \delta > 0$ such that $|x_i - x_j| > \delta$ and $|x_j| \leq M$ for all i and j . Then $\mathbf{x} \in \mathcal{S}_{\delta/M}$.*

Proof. Assuming $x_j \neq 0$, we have that

$$\frac{x_i}{x_j} = 1 + \frac{x_i - x_j}{x_j} \notin \left[1 - \frac{\delta}{M}, 1 + \frac{\delta}{M}\right] \iff \left|\frac{x_i - x_j}{x_j}\right| > \frac{\delta}{M}.$$

Thus the lemma follows from

$$\left|\frac{x_i - x_j}{x_j}\right| > \frac{\delta}{|x_j|} \geq \frac{\delta}{M},$$

where the first inequality is by the assumption $|x_i - x_j| > \delta$ and the second inequality is by the assumption $|x_j| \leq M$ which implies $1/|x_j| \geq 1/M$. \square

Lemma D.3.

$$\prod_{i=1}^{\infty} \left(1 + \frac{2}{i^3}\right)^2 \leq 20.$$

Proof. Compute

$$\begin{aligned} \prod_{i=1}^{\infty} \left(1 + \frac{2}{i^3}\right)^2 &= 9 \left(\frac{5}{4}\right)^2 \prod_{i=3}^{\infty} \left(1 + \frac{2}{i^3}\right)^2 = \frac{225}{16} \exp\left(2 \sum_{i=3}^{\infty} \ln\left(1 + \frac{2}{i^3}\right)\right) \\ &\leq \frac{225}{16} \exp\left(2 \sum_{i=3}^{\infty} \frac{2}{i^3}\right) = \frac{225}{16} \exp\left(4 \left(\zeta(3) - \frac{9}{8}\right)\right) \\ &\leq \frac{225}{16} \exp\left(4 \left(1.21 - \frac{9}{8}\right)\right), \end{aligned}$$

where the first inequality follows from the inequality $\ln(1+x) < x$ for all $x > 0$, and the second inequality is a known bound $\zeta(3) \leq 1.21$ where $\zeta(\cdot)$ is Riemann's zeta function. Evaluating the above expression, the lemma follows. \square

Lemma D.4. *For all natural $d \geq 58$ and $1 \leq k \leq \lceil \log(\log(d) + 1) \rceil$, we have*

$$1 \leq \frac{2d^{1-\beta(k+1)}}{(k+1)^3}.$$

Proof. We first verify the lemma for $k = 1$. We have

$$\frac{2d^{1-\beta(2)}}{8} = \frac{1}{4}d^{\frac{2}{3}} \geq \frac{1}{4}\sqrt{d} \geq \frac{1}{4}\sqrt{100} \geq 1.$$

Next, assume $k \geq 2$ and compute

$$\frac{2d^{1-\beta(k+1)}}{(k+1)^3} \geq \frac{2d^{\frac{6}{7}}}{(k+1)^3} \geq \frac{2d^{\frac{6}{7}}}{(\lceil \log(\log(d) + 1) \rceil + 1)^3}.$$

It thus suffices to prove that

$$\lceil \log(\log(d) + 1) \rceil + 1 \leq 2^{\frac{1}{3}}d^{\frac{6}{21}}.$$

It is easy to see that this inequality holds for $d = 58$ (using any symbolic computation package), and since the left hand side is constant for all $d \in [58, 128]$ whereas the right hand side is increasing, the inequality also holds for all $d \leq 128$. By the same reasoning, we observe that the left hand side is constant on any interval of the form $[2^{2^n-1}, 2^{2^{n+1}-1}]$ for integer $n \geq 3$, and takes the value of $n + 2$. In contrast, the right hand side is lower bounded by $2^{\frac{1}{3}}(2^{2^n-1})^{\frac{6}{21}}$ on each such interval. It is thus sufficient to prove that

$$n + 2 \leq 2^{\frac{1}{3}}(2^{2^n-1})^{\frac{6}{21}}$$

for all integer $n \geq 3$. We shall show this using induction. The base case can be easily verified for $n = 3$. Assuming the induction hypothesis for n , we compute

$$\begin{aligned} 2^{\frac{1}{3}}(2^{2^{n+1}-1})^{\frac{6}{21}} &= 2^{\frac{1}{3}}(2^{2^n-1})^{\frac{6}{21}} \cdot \frac{(2^{2^{n+1}-1})^{\frac{6}{21}}}{(2^{2^n-1})^{\frac{6}{21}}} \geq (n+2)(2^{2^n})^{\frac{6}{21}} \\ &\geq (n+2)(2^8)^{\frac{6}{21}} \geq 2n+4 \geq n+3, \end{aligned}$$

where the first inequality follows from the induction hypothesis, and the second inequality follows from $n \geq 3$. \square

Lemma D.5. *Let μ denote the d -dimensional Lebesgue measure, Let $D \subseteq \mathbb{R}^d$ be compact, and suppose that $g : D \rightarrow \mathbb{R}$ is continuous and that $\Omega \subseteq [0, 1]^d$ is a convex set satisfying $\mu(\Omega) > 0$. Then there exists some $\mathbf{x}_0 \in \Omega$ such that*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}(D)} [g(\mathbf{x}) | \mathbf{x} \in \Omega] = g(\mathbf{x}_0).$$

Proof. Since $\Omega = (\partial\Omega \cap \Omega) \cup \text{int}(\Omega)$ is a disjoint union and since the boundary of a convex set in \mathbb{R}^d has measure zero [13], we have

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}(D)} [g(\mathbf{x}) | \mathbf{x} \in \Omega] = \frac{1}{\mu(\Omega)} \int_{\Omega} g(\mathbf{x}) d\mathbf{x} = \frac{1}{\mu(\Omega)} \int_{\text{int}(\Omega)} g(\mathbf{x}) d\mathbf{x}.$$

Due to the above, we may assume w.l.o.g. that $\Omega = \text{int}(\Omega)$ is thus open and Lebesgue measurable. If g is constant on Ω then the lemma holds true for all $\mathbf{x}_0 \in \Omega$. Suppose that g is not constant on Ω , then due to being continuous on the compact domain $D \supseteq \Omega$, it is bounded on Ω , and there exist $\mathbf{x}_1, \mathbf{x}_2 \in \Omega$ such that $g(\mathbf{x}_1) < g(\mathbf{x}_2)$. Denote

$$-\infty < m := \inf_{\mathbf{x} \in \Omega} g(\mathbf{x}) \leq g(\mathbf{x}_1) < g(\mathbf{x}_2) \leq \sup_{\mathbf{x} \in \Omega} g(\mathbf{x}) =: M < \infty.$$

In particular, by the continuity of g and since Ω is open, there exists some open neighborhood $U \subseteq \Omega$ containing \mathbf{x}_2 and satisfying $\mu(U) > 0$ such that $g(\mathbf{x}') \geq \frac{g(\mathbf{x}_1) + g(\mathbf{x}_2)}{2} > g(\mathbf{x}_1) \geq m$ for all $\mathbf{x}' \in U$. We now have

$$\begin{aligned} \frac{1}{\mu(\Omega)} \int_{\Omega} g(\mathbf{x}) d\mathbf{x} &= \frac{1}{\mu(\Omega)} \left(\int_{\Omega \setminus U} g(\mathbf{x}) d\mathbf{x} + \int_U g(\mathbf{x}) d\mathbf{x} \right) \\ &\geq \frac{1}{\mu(\Omega)} \left(m \cdot \mu(\Omega \setminus U) + \int_U g(\mathbf{x}) d\mathbf{x} \right) \\ &> \frac{1}{\mu(\Omega)} (m \cdot \mu(\Omega \setminus U) + m \cdot \mu(U)) = m. \end{aligned}$$

An analogous argument shows that $\frac{1}{\mu(\Omega)} \int_{\Omega} g(\mathbf{x}) d\mathbf{x} < M$, and we thus deduce that

$$m < \frac{1}{\mu(\Omega)} \int_{\Omega} g(\mathbf{x}) d\mathbf{x} < M. \quad (13)$$

Let $\{\mathbf{a}_n\}_{n=1}^{\infty}, \{\mathbf{b}_n\}_{n=1}^{\infty} \subseteq \Omega$ such that $\lim_{n \rightarrow \infty} g(\mathbf{a}_n) = m$ and $\lim_{n \rightarrow \infty} g(\mathbf{b}_n) = M$. Then for any $\varepsilon > 0$, there exists n_0 such that $g(\mathbf{a}_{n_0}) \leq m + \varepsilon$ and $g(\mathbf{b}_{n_0}) \geq M - \varepsilon$. Consider the path $p : [0, 1] \rightarrow \mathbb{R}^d$ given by

$$p(\lambda) := \lambda \mathbf{a}_{n_0} + (1 - \lambda) \mathbf{b}_{n_0}.$$

From the convexity of Ω we have that $p(\lambda) \in \Omega$ for all $\lambda \in [0, 1]$, and since $g(p(\lambda))$ is continuous in λ , for all $y \in [m + \varepsilon, M - \varepsilon]$ there exists some $\lambda \in [0, 1]$ such that

$$g(p(\lambda)) = y.$$

In particular, using Eq. (13), we can choose $\varepsilon > 0$ sufficiently small such that

$$\frac{1}{\mu(\Omega)} \int_{\Omega} g(\mathbf{x}) d\mathbf{x} \in [m + \varepsilon, M - \varepsilon],$$

and find some λ_0 satisfying

$$\frac{1}{\mu(\Omega)} \int_{\Omega} g(\mathbf{x}) d\mathbf{x} = g(p(\lambda_0)).$$

Letting $\mathbf{x}_0 := p(\lambda_0) \in \Omega$ gives the desired result. \square

Lemma D.6. Let $\mathcal{N}_{k,\ell}$ denote the class of width k and depth ℓ neural networks with an arbitrary activation function. Then there exists some $\varepsilon > 0$ such that

$$\inf_{\mathcal{N} \in \mathcal{N}_{k,\ell}} \mathbb{E}_{\mathbf{x} \sim \mathcal{U}([0,1]^d)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] \geq \varepsilon,$$

if and only if for all $R > 0$ and all $a \in \mathbb{R}$, we have

$$\inf_{\mathcal{N} \in \mathcal{N}_{k,\ell}} \mathbb{E}_{\mathbf{x} \sim \mathcal{U}([a,a+R]^d)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] \geq R^2 \varepsilon.$$

Proof. The fact that the latter implies the former is immediate by substituting $R = 1$ and $a = 0$. For the reverse implication, observing that $f_d(\mathbf{x}) = R \cdot f_d(\frac{1}{R}\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^d$, we have by [21, Theorem 9] that

$$\inf_{\mathcal{N} \in \mathcal{N}_{k,\ell}} \mathbb{E}_{\mathbf{x} \sim \mathcal{U}([0,1]^d)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] \geq \varepsilon$$

implies

$$\inf_{\mathcal{N} \in \mathcal{N}_{k,\ell}} \mathbb{E}_{\mathbf{x} \sim \mathcal{U}([0,R]^d)} [(\mathcal{N}(\mathbf{x}) - f_d(\mathbf{x}))^2] \geq R^2 \varepsilon,$$

for all $R > 0$. Writing the above expectation in integral form and performing the change of variables $\mathbf{x} = \mathbf{y} + (a, \dots, a)$, $d\mathbf{x} = d\mathbf{y}$, the lemma follows. \square