

On the approximation of functions by tanh neural networks

T. De Ryck and S. Lanthaler and S. Mishra

Research Report No. 2021-14
April 2021

Seminar für Angewandte Mathematik
Eidgenössische Technische Hochschule
CH-8092 Zürich
Switzerland

ON THE APPROXIMATION OF FUNCTIONS BY TANH NEURAL NETWORKS

TIM DE RYCK, SAMUEL LANTHALER, AND SIDDHARTHA MISHRA

ABSTRACT. We derive bounds on the error, in high-order Sobolev norms, incurred in the approximation of Sobolev-regular as well as analytic functions by neural networks with the hyperbolic tangent activation function. These bounds provide explicit estimates on the approximation error with respect to the size of the neural networks. We show that tanh neural networks with only two hidden layers suffice to approximate functions at comparable or better rates than much deeper ReLU neural networks.

1. INTRODUCTION

Deep learning, relying on the use of deep artificial neural networks for regression and classification, has been very successful in different contexts in science and engineering in recent years [34]. These include image recognition, natural language understanding, machine translation, game intelligence, robotics, autonomous systems and protein folding.

Deep learning is also being increasingly used in scientific computing, particularly in the numerical solution of partial differential equations (PDEs). A very incomplete list of examples for the successful use of deep learning in this context includes the solution of high-dimensional linear and semi-linear parabolic partial differential equations [17, 21] and references therein, the solution of parametric partial differential equations that arise in many-query problems like uncertainty quantification (UQ), PDE constrained optimization and (Bayesian) inverse problems [54, 30, 31, 42, 41] and in infinite-dimensional operator learning frameworks [39, 33, 36]. Another avenue for the application of deep neural networks in scientific computing is provided by *physics-informed neural networks* (PINNs) [32, 57, 58, 44, 45], which serve as replacements for traditional numerical methods for both forward as well as inverse problems for PDEs.

The question of why deep neural networks are so successful at many diverse tasks in very different fields eludes a definitive answer. A very partial explanation may lie in the fact that artificial neural networks are *universal approximators* i.e., any continuous (even measurable) mapping can be approximated by artificial neural networks to arbitrarily high accuracy [1, 13, 24] and references therein. However, such universality results only imply the existence of a (shallow) neural network and do not provide any quantitative information (bounds) on the width of the underlying neural networks.

The task of quantitatively relating the size and architecture of neural networks to their expressivity i.e., accuracy in approximating functions of a certain hypothesis class, has received considerable attention in the literature in the last few years. A seminal work in the direction is [65], where the author derived explicit estimates

(T. De Ryck, S. Lanthaler and S. Mishra) SEMINAR FOR APPLIED MATHEMATICS, ETH ZÜRICH, RÄMISTRASSE 101, 8092 ZÜRICH, SWITZERLAND

E-mail addresses: tim.deryck@sam.math.ethz.ch, samuel.lanthaler@sam.math.ethz.ch, siddhartha.mishra@sam.math.ethz.ch.

on the size (width and depth) of a neural network with a ReLU activation function for approximating Lipschitz functions to any given accuracy in the L^∞ -norm. Expressivity results for such ReLU neural networks in Sobolev norms were presented in [19, 22, 54] and references therein, see also [35, 49, 53, 60, 66] and references therein for further approximation results for ReLU and related ReQU and RePU activation functions.

Despite the fact that several quantitative results on the expressivity of neural networks have been obtained in recent years, we highlight some of the lacunae of the current state of the art in this direction,

- Most of the available results are on the expressivity and approximation properties of ReLU neural networks. Although ReLU activations are very common in practical applications of deep learning, there is a large number of areas where other activation functions are employed. One of the most popular activations is the tanh (hyperbolic tangent) activation function and the related sigmoid or logistic function (a scaled and shifted tanh). These activation functions are the basis of heavily used recurrent neural network (RNN) architectures such as LSTM [23] and GRU [8]. Other areas where smooth activation functions such as tanh are preferred over ReLU is in physics-informed neural networks (PINNs) for solving forward and inverse problems for PDEs [57, 58, 44, 45] and references therein, and in the use of quasi-random training points [47, 38]. Although the approximation abilities of general smooth activation functions have been investigated in [10, 9, 20, 52, 55, 62] and references therein, it is fair to say that the level of detail in existing results for the expressivity of ReLU neural networks, is not yet available for tanh neural networks.
- Moreover, most of the approximation results for smooth activation functions, with the exception of the recent paper [20], measure error in L^p -norms. However, it is essential to measure errors in higher-order Sobolev norms for many applications, such as PINNs where the neural network needs to be differentiated in order to evaluate the underlying PDE residual.
- A persistent focus of approximation results for neural networks has been to highlight the role of depth of the neural network, see [56] for a review and further references. In particular, one wishes to prove that very deep neural networks are, in some sense, more expressive than shallower networks and use this to explain the superior performance of deep neural networks in many applications. The empirical superiority of deep networks over their shallower counterparts has indeed been observed in many applications in computer science. However, in the context of scientific computing, empirical experience has revealed that shallower but wider networks result in superior performance over deep and narrow neural networks, see [40] and references therein. A reason for this observation lies in the fact that deeper networks might be harder to train in the relatively data poor regime of scientific computing. Some theoretical understanding of this deterioration of performance for deeper networks, at least in the context of ReLU networks is provided in [18]. However, most of the available approximation theory results trade width for depth and there is little theoretical understanding of why relatively shallow networks can perform well in some contexts.
- Most of the available results on expressivity focus on asymptotic approximation rates i.e., the complexity of the network as the approximation error $\epsilon \rightarrow 0$. However, the fundamental question is how large a neural network

should be to provide a certain accuracy of this approximation. This requires going beyond asymptotic approximation rates and providing explicit bounds on the underlying constants. With the exception of [1, 2], such explicit bounds are mostly unavailable.

- The approximation error is only one component of the total error of neural networks, with optimization and generalization errors being the other components [61, 12]. In particular, standard approaches to estimate the generalization error such as covering number estimates [3] or Rademacher complexity [61] require explicit estimates on the weights of the underlying neural networks, in addition to bounds on their width and depth. Such estimates on weights of the best approximations of functions in the class of neural networks are rarely available in the current literature.

The main objective of this paper is to address some of the afore-mentioned deficiencies in the literature on approximation properties of neural networks. We will focus on the expressivity of neural networks with the very popular tanh activation function and will aim to prove error and complexity bounds in high-order Sobolev norms for such tanh neural networks in approximating functions belonging to Sobolev spaces as well as C^k -spaces. We go beyond the usual practice of proving only asymptotic convergence rates and will provide explicit approximation error bounds for explicit network architectures in order to answer the question of “*How large should a neural network be to approximate a specified function to some chosen accuracy $\epsilon > 0$?*”. All our results will be for tanh neural networks with at most two hidden layers.

Moreover, we also consider the case of analytic functions and will prove that a two hidden layer tanh neural network suffices to approximate an analytic function at an exponential rate, in terms of the network width, even in Sobolev norms. This result provides an improvement over available results for the approximation of analytic functions by ReLU neural networks [64, 54, 22] and also neural networks with smooth activation functions [43] and further illustrate the powers of rather shallow tanh networks at approximating smooth functions. Finally, we also derive explicit bounds on the width of the tanh neural networks as well as asymptotic bounds on their weights, thus paving the way for bounds on the generalization error for these neural networks.

The rest of the paper is organized as follows: in Section 2, we introduce the notation for the rest of the paper. Our main results, presented in Section 5, rely on the uniform approximation of polynomials by tanh neural networks, discussed in Section 3, and on the approximation of a partition of unity, presented in Section 4. In Section 6, we discuss the contents of this paper and distinguish them from other related papers.

2. PRELIMINARIES

We start by providing an overview of all the notation and the definitions that will be used frequently throughout the paper.

2.1. Multi-index notation. For $d \in \mathbb{N}$, we call a d -tuple of non-negative integers $\alpha \in \mathbb{N}_0^d$ a multi-index. We write $|\alpha| = \sum_{i=1}^d \alpha_i$, $\alpha! = \prod_{i=1}^d \alpha_i!$ and, for $x \in \mathbb{R}^d$, we denote by $x^\alpha = \prod_{i=1}^d x_i^{\alpha_i}$ the corresponding multinomial. Given two multi-indices $\alpha, \beta \in \mathbb{N}_0^d$, we say that $\alpha \leq \beta$ if, and only if, $\alpha_i \leq \beta_i$ for all $i = 1, \dots, d$. For a multi-index α , we define the following multinomial coefficient

$$(1) \quad \binom{|\alpha|}{\alpha} = \frac{|\alpha|!}{\alpha!},$$

and, given $\alpha \leq \beta$, we define a corresponding multinomial coefficient by

$$(2) \quad \binom{\beta}{\alpha} = \prod_{i=1}^d \binom{\beta_i}{\alpha_i} = \frac{\beta!}{\alpha!(\beta - \alpha)!}.$$

For $\Omega \subseteq \mathbb{R}^d$ and f a function that is at least $|\alpha|$ times continuously differentiable on Ω , we define

$$(3) \quad D^\alpha f = \frac{\partial^{|\alpha|} f}{\partial x_1^{\alpha_1} \cdots \partial x_d^{\alpha_d}}.$$

We will frequently encounter the set $P_{n,d} = \{\alpha \in \mathbb{N}_0^d : |\alpha| = n\}$ (notation as in [48]). In particular, we will need estimates on its cardinality. This is the subject of the following lemma.

Lemma 2.1. *Let $n \in \mathbb{N}$, $d \in \mathbb{N}_{\geq 2}$ and let $P_{n,d} = \{\alpha \in \mathbb{N}_0^d : |\alpha| = n\}$. Then*

$$(4) \quad |P_{n,d}| = \binom{n+d-1}{n} \leq \sqrt{\pi} \min\{e^{d-1}n^{d-1}, e^n(d-1)^n\} \quad \text{and} \quad |P_{d,d}| \leq 5^d.$$

Proof. It is well known that $|P_{n,d}| = \binom{n+d-1}{n}$. We use Stirling's approximation,

$$\begin{aligned} |P_{n,d}| &= \binom{n+d-1}{n} \leq \frac{e(n+d-1)^{n+d-1/2}}{2\pi n^{n+1/2}(d-1)^{d-1/2}} \\ &\leq \frac{e}{2\pi} \left(\frac{n+d-1}{d-1}\right)^{d-1} \left(\frac{n+d-1}{n}\right)^n \sqrt{\frac{n+d-1}{n(d-1)}} \\ &\leq \frac{e}{\sqrt{2\pi}} \left(1 + \frac{n}{d-1}\right)^{d-1} \left(1 + \frac{d-1}{n}\right)^n. \end{aligned}$$

To estimate the last term, we note that there are two possible approximations: for $a, b \geq 1$ it holds that $(1+a/b)^b \leq ea^b$ and also $(1+a/b)^b \leq e^a$. Using the fact that $e^2/\sqrt{2\pi} \leq \sqrt{\pi}$, we obtain

$$(5) \quad |P_{n,d}| \leq \sqrt{\pi} e^{d-1} n^{d-1} \quad \text{and} \quad |P_{n,d}| \leq \sqrt{\pi} e^n (d-1)^n.$$

Setting $n = d$ for $\lambda \in \mathbb{N}$, we also find that

$$(6) \quad |P_{d,d}| \leq \frac{e}{\sqrt{2\pi}} \left(1 + \frac{d}{d-1}\right)^d \left(1 + \frac{d-1}{d}\right)^d \leq \left(3 + \frac{d}{d-1}\right)^d \leq 5^d,$$

since $\frac{e}{\sqrt{2\pi}} \leq 1$ and $\frac{x}{x-1} \leq 2$ for $x \geq 2$. □

2.2. Sobolev spaces. Let $d \in \mathbb{N}$, $1 \leq p \leq \infty$ and let $\Omega \subseteq \mathbb{R}^d$ be open. We denote by $L^p(\Omega)$ the usual Lebesgue space and for $k \in \mathbb{N}_0$ we define the Sobolev space $W^{k,p}(\Omega)$ as

$$(7) \quad W^{k,p}(\Omega) = \{f \in L^p(\Omega) : D^\alpha f \in L^p(\Omega) \text{ for all } \alpha \in \mathbb{N}_0^d \text{ with } |\alpha| \leq k\}.$$

For $p < \infty$, we define the following seminorms on $W^{k,p}(\Omega)$,

$$(8) \quad |f|_{W^{m,p}(\Omega)} = \left(\sum_{|\alpha|=m} \|D^\alpha f\|_{L^p(\Omega)}^p \right)^{1/p} \quad \text{for } m = 0, \dots, k,$$

and for $p = \infty$ we define

$$(9) \quad |f|_{W^{m,\infty}(\Omega)} = \max_{|\alpha|=m} \|D^\alpha f\|_{L^\infty(\Omega)} \quad \text{for } m = 0, \dots, k.$$

Based on these seminorms, we can define the following norm for $p < \infty$,

$$(10) \quad \|f\|_{W^{k,p}(\Omega)} = \left(\sum_{m=0}^k |f|_{W^{m,p}(\Omega)}^p \right)^{1/p},$$

and for $p = \infty$ we define the norm

$$(11) \quad \|f\|_{W^{k,\infty}(\Omega)} = \max_{0 \leq m \leq k} |f|_{W^{m,\infty}(\Omega)}.$$

The space $W^{k,p}(\Omega)$ equipped with the norm $\|\cdot\|_{W^{k,p}(\Omega)}$ is a Banach space.

2.3. Neural networks. In this paper, we will consider function approximation using feedforward artificial neural networks where only connections between neighbouring layers are allowed. In the following, we formally introduce our definition of a neural network and the related terminology.

Let $L \in \mathbb{N}$ and $l_0, \dots, l_L \in \mathbb{N}$. Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an *activation function* and define the parameter space

$$(12) \quad \Theta = \bigcup_{L \in \mathbb{N}} \bigcup_{l_0, \dots, l_L \in \mathbb{N}} \bigtimes_{k=1}^L \left(\mathbb{R}^{l_k \times l_{k-1}} \times \mathbb{R}^{l_k} \right).$$

For $\theta \in \Theta$, we define $\theta_k := (W_k, b_k)$ and $\mathcal{A}_k : \mathbb{R}^{l_{k-1}} \rightarrow \mathbb{R}^{l_k} : x \mapsto W_k x + b_k$ for $1 \leq k \leq L$ and we denote by $\Psi_\theta : \mathbb{R}^{l_0} \rightarrow \mathbb{R}^{l_L}$, $x \mapsto \Psi_\theta(x)$, the function

$$(13) \quad \Psi_\theta(x) = \begin{cases} \mathcal{A}_1(x) & L = 1, \\ (\mathcal{A}_L \circ \sigma \circ \mathcal{A}_{L-1} \circ \sigma \circ \dots \circ \sigma \circ \mathcal{A}_1)(x) & L \geq 2, \end{cases}$$

where σ is applied element-wise. We refer to Ψ_θ as the realization of the *neural network* associated to the parameter θ with L layers and widths (l_0, l_1, \dots, l_L) . We refer to the first $L - 1$ layers as *hidden layers*. For $1 \leq k \leq L$, we say that layer k has width l_k and we refer to W_k and b_k as the *weights and biases* corresponding to layer k . The width of Ψ_θ is defined as $\max(l_0, \dots, l_L)$. If $L = 2$, we say that Ψ_θ is a *shallow neural network*; if $L \geq 3$, we say that Ψ_θ is a *deep neural network*. Hence, a shallow neural network has exactly one hidden layer whereas deep neural networks can have two or more hidden layers.

In this work, we will focus on neural networks which use the hyperbolic tangent as activation function, defined by

$$(14) \quad \sigma(x) := \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad \text{for } x \in \mathbb{R}.$$

We will refer to these networks as *tanh neural networks*. Even though our ideas can be carried over to other smooth activation functions, focusing on a particular activation will allow us to prove precise and explicit bounds without sacrificing the clarity of our arguments. In particular, we note that all results of this work directly apply to the sigmoid or logistic activation function, which is simply a shifted and scaled version of the hyperbolic tangent.

We close this section by recalling some basic properties of neural network calculus which we will use throughout, without explicitly referring to them.

Proposition 2.2 (Parallelization of neural networks). *Let $L \in \mathbb{N}$, $l_0, l'_0, \dots, l_L, l'_L \in \mathbb{N}$ and $\theta, \vartheta \in \Theta$ such that Ψ_θ is a neural network with widths (l_0, \dots, l_L) and Ψ_ϑ is a neural network with widths (l'_0, \dots, l'_L) . Then there exists $\eta \in \Theta$ such that Ψ_η is a neural network with widths $(l_0 + l'_0, \dots, l_L + l'_L)$ for which it holds that $\Psi_\eta(x) = (\Psi_\theta((x_1, \dots, x_{l_0})), \Psi_\vartheta((x_{l_0+1}, \dots, x_{l_0+l'_0})))$ for all $x \in \mathbb{R}^{l_0+l'_0}$.*

Proposition 2.3 (Composition of neural networks). *Let $L, L' \in \mathbb{N}$, $l_0, \dots, l_L = l'_0, \dots, l'_{L'} \in \mathbb{N}$ and $\theta, \vartheta \in \Theta$ such that Ψ_θ is a neural network with widths (l_0, \dots, l_L) and Ψ_ϑ is a neural network with widths $(l'_0, \dots, l'_{L'})$. Then there exists $\eta \in \Theta$ such that Ψ_η is a neural network with widths $(l_0, \dots, l_L = l'_0, \dots, l'_{L'})$ for which it holds that $\Psi_\eta = \Psi_\vartheta \circ \Psi_\theta$.*

3. UNIFORM APPROXIMATION OF POLYNOMIALS

The first step in our strategy for deriving bounds on approximation error for tanh neural networks is to provide uniform bounds in Sobolev norms, on the error for approximating polynomials by shallow tanh neural networks. We do so in the current section.

The observation that a shallow neural network of fixed size can approximate monomials to arbitrary accuracy in the supremum norm was already observed in [55]. A generalization of this approximation result to Sobolev norms was proven in e.g. [20]. In the current section, we present a *novel* generalization that allows us to obtain explicit error estimates for the *uniform* approximation of all polynomials of a certain maximal degree, which will be crucial for the efficient approximation of analytic functions.

3.1. Univariate polynomials. We first describe how to approximate univariate polynomials of any degree with tanh neural networks. We introduce the p -th order central finite difference operator δ_h^p for any $f \in C^{p+2}([a, b])$ for some $p \in \mathbb{N}$ by

$$(15) \quad \delta_h^p[f](x) = \sum_{i=0}^p (-1)^i \binom{p}{i} f\left(x + \left(\frac{p}{2} - i\right)h\right).$$

Next we define for any $p \in \mathbb{N}$, $q \in 2\mathbb{N} - 1$ and $M > 0$ the monomials $f_p : [-M, M] \rightarrow \mathbb{R}$ and the tanh neural networks $\hat{f}_{q,h} : [-M, M] \rightarrow \mathbb{R}$ as

$$(16) \quad f_p(y) := y^p \quad \text{and} \quad \hat{f}_{q,h}(y) := \frac{\delta_{hy}^q[\sigma](0)}{\sigma^{(q)}(0)h^q}.$$

We first prove that these neural networks are accurate approximations to monomials with odd degree.

Lemma 3.1. *Let $k \in \mathbb{N}_0$ and $s \in 2\mathbb{N} - 1$. Then it holds that for all $\epsilon > 0$ there exists $h > 0$ and a shallow tanh neural network $\Psi_{s,h} : [-M, M] \rightarrow \mathbb{R}^{\frac{s+1}{2}}$ of width $\frac{s+1}{2}$ such that*

$$(17) \quad \max_{\substack{p \leq s, \\ p \text{ odd}}} \|f_p - (\Psi_{s,h})_p\|_{W^{k,\infty}} \leq \epsilon,$$

Moreover, the weights of $\Psi_{s,h}$ scale as $O\left(\epsilon^{-s/2}(2(s+2)\sqrt{2M})^{s(s+3)}\right)$ for small ϵ and large s .

Proof. Let $p \leq s$ be odd and let $0 < h < 2/pM$. Let $0 \leq m \leq \min\{k, p+1\}$. Then Taylor's theorem guarantees the existence of $\xi_{x,i}$ such that

$$\begin{aligned} \frac{d^m}{dx^m} \delta_{hx}^p[\sigma](0) &= \sum_{i=0}^p (-1)^i \binom{p}{i} \left(\frac{p}{2} - i\right)^m h^m \cdot \sigma^{(m)}\left(\left(\frac{p}{2} - i\right)hx\right) \\ &= \sum_{i=0}^p (-1)^i \binom{p}{i} \left(\frac{p}{2} - i\right)^m h^m \left(\sum_{l=m}^{p+1} \frac{\sigma^{(l)}(0)}{(l-m)!} \left(\frac{p}{2} - i\right)^{l-m} (hx)^{l-m}\right) \\ &\quad + \sum_{i=0}^p (-1)^i \binom{p}{i} \left(\frac{p}{2} - i\right)^m h^m \frac{\sigma^{(p+2)}(\xi_{x,i})}{(p+2-m)!} \left(\frac{p}{2} - i\right)^{p+2-m} (hx)^{p+2-m}. \end{aligned}$$

From [27, Theorem 1] it follows that

$$(18) \quad \sum_{i=0}^p (-1)^i \binom{p}{i} \left(\frac{p}{2} - i\right)^l = p! \delta(l-p) = \begin{cases} p!, & (l=p), \\ 0, & (l \neq p). \end{cases}$$

for $l = 0, \dots, p$. We observe that (18) remains true also for $l = p+1$, since all summands change sign when i is replaced by $p-i$. Using this fact, we can then rewrite the first term as

$$(19) \quad \begin{aligned} & \sum_{i=0}^p (-1)^i \binom{p}{i} \left(\frac{p}{2} - i\right)^m h^m \left(\sum_{l=m}^{p+1} \frac{\sigma^{(l)}(0)}{(l-m)!} \left(\frac{p}{2} - i\right)^{l-m} (hx)^{l-m} \right) \\ &= h^m \sum_{l=m}^{p+1} \frac{\sigma^{(l)}(0)}{(l-m)!} (hx)^{l-m} \sum_{i=0}^p (-1)^i \binom{p}{i} \left(\frac{p}{2} - i\right)^l \\ &= \begin{cases} h^m \frac{\sigma^{(p)}(0)}{(p-m)!} (hx)^{p-m} p!, & 0 \leq m \leq p \\ 0, & m = p+1 \end{cases} = h^p \sigma^{(p)}(0) f_p^{(m)}(x). \end{aligned}$$

Combining the previous results, it thus follows that we have

$$\hat{f}_{p,h}^{(m)}(x) - f_p^{(m)}(x) = \sum_{i=0}^p (-1)^i \binom{p}{i} \frac{1}{(p+2-m)!} \frac{\sigma^{(p+2)}(\xi_{x,i})}{\sigma^{(p)}(0)} \left(\frac{p}{2} - i\right)^{p+2} h^2 x^{p+2-m}.$$

Together with the lower and upper bounds on the derivatives of σ from Lemma A.1 and Lemma A.4, this yields for $m \leq \min(k, p+1)$:

$$(20) \quad \begin{aligned} \left| f_p - \hat{f}_{p,h} \right|_{W^{m,\infty}} &\leq \sum_{i=0}^p \binom{p}{i} \frac{\left| \sigma^{(p+2)}(\xi_{x,i}) \right|}{\left| \sigma^{(p)}(0) \right|} \left| \frac{p}{2} - i \right|^{p+2} h^2 M^{p+2} \\ &\leq 2^p \frac{(2(p+2))^{p+3}}{1} \left(\frac{p}{2}\right)^{p+2} h^2 M^{p+2} \\ &\leq (2(p+2)pM)^{p+3} h^2. \end{aligned}$$

If $k \leq p+1$, then this shows that

$$(21) \quad \left\| f_p - \hat{f}_{p,h} \right\|_{W^{k,\infty}} \leq (2(p+2)pM)^{p+3} h^2.$$

If $k > p+1$, let $p+2 \leq m \leq k$. In this case, $f_p^{(m)} = 0$, therefore it suffices to bound $\hat{f}_{p,h}^{(m)}$. We see that for $0 < h < 1$,

$$(22) \quad \begin{aligned} \left| \hat{f}_{p,h}^{(m)}(x) \right| &= \left| \frac{1}{h^p \sigma^{(p)}(0)} \sum_{i=0}^p (-1)^i \binom{p}{i} \left(\frac{p}{2} - i\right)^m h^m \cdot \sigma^{(m)} \left(\left(\frac{p}{2} - i\right) hx \right) \right| \\ &\leq 2 \sum_{i=0}^p \binom{p}{i} \left| \frac{p}{2} - i \right|^m h^2 (2m)^{m+1} \\ &\leq 2^{p+1} \left(\frac{p}{2}\right)^k (2k)^{k+1} h^2 \leq (2pk)^{k+1} h^2. \end{aligned}$$

We thus obtain, for arbitrary $k \in \mathbb{N}$:

$$(23) \quad \left\| f_p - \hat{f}_{p,h} \right\|_{W^{k,\infty}} \leq \left((2(p+2)pM)^{p+3} + (2pk)^{k+1} \right) h^2 =: \epsilon.$$

Furthermore observe that the weights scale as $O\left(\max_i \binom{p}{i} h^{-p}\right)$. For $\epsilon \rightarrow 0$ and large p , it holds that $O(h^{-p}) = O\left(\epsilon^{-p/2} ((p+2)\sqrt{2M})^{p(p+3)}\right)$, where the implied

constant depends on k . Next, we find using Stirling's approximation that for $0 \leq i \leq p$ it holds that

$$(24) \quad \binom{p}{i} \leq \binom{p}{\frac{p-1}{2}} \leq \frac{ep^{p+1/2}}{2\pi \left(\frac{p-1}{2}\right)^{\frac{p}{2}} \left(\frac{p+1}{2}\right)^{\frac{p}{2}+1}} = O\left(\frac{2^p}{\sqrt{p}}\right).$$

The weights therefore scale as $O\left(\epsilon^{-p/2}(2(p+2)\sqrt{2M})^{p(p+3)}\right)$.

Regarding the network architecture, note that the neurons needed for all $\hat{f}_{p,h}$ are already available in the network $\hat{f}_{s,h}$. This allows us to define the shallow tanh neural network $\Psi_{s,h}$ by $(\Psi_{s,h})_p = \hat{f}_{p,h}$ such that it only has $\frac{s+1}{2}$ neurons in its hidden layer. The width follows directly from its definition and the fact that σ is an odd function. \square

We would like to state that the above proof is largely inspired by [20, Proposition 4.7] but differs at some crucial points. In particular, we take into account the fact that $\xi_{x,i}$'s are functions of x and derivatives with respect to x have to take this into account.

We now extend the previous result to monomials with even degree. To this end, we rely on the observation that for $n \in \mathbb{N}$ and $\alpha > 0$, it holds that

$$(25) \quad y^{2n} = \frac{1}{2\alpha(2n+1)} \left((y+\alpha)^{2n+1} - (y-\alpha)^{2n+1} - 2 \sum_{k=0}^{n-1} \binom{2n+1}{2k} \alpha^{2(n-k)+1} y^{2k} \right).$$

This formula allows us to construct recursively defined tanh neural network approximations of even powers of y . The following lemma quantifies the uniform approximation accuracy of these networks in the Sobolev norm.

Lemma 3.2. *Let $k \in \mathbb{N}_0$, $s \in 2\mathbb{N} - 1$ and $M > 1$. For every $\epsilon > 0$, there exists $h > 0$ and a shallow tanh neural network $\psi_{s,h} : [-M+1, M-1] \rightarrow \mathbb{R}^s$ of width $\frac{3(s+1)}{2}$ such that*

$$(26) \quad \max_{p \leq s} \|f_p - (\psi_{s,h})_p\|_{W^{k,\infty}} \leq \epsilon.$$

Furthermore, the weights scale as $O\left(\epsilon^{-s/2}(\sqrt{M}(s+2))^{3s(s+3)/2}\right)$ for small ϵ and large s .

Proof. For $p \leq s$, $\alpha \leq 1$ and $y \in [-M+1, M-1]$, we define $(\psi_{s,h}(y))_p = \hat{f}_{p,h}(y)$ for p odd and, for $p = 2n$ even, we define $(\psi_{s,h}(y))_p = (\psi_{s,h}(y))_{2n}$ recursively by $(\psi_{s,h})_0(y) := 1$, and

$$(27) \quad (\psi_{s,h}(y))_{2n} = \frac{1}{2\alpha(2n+1)} \left(\hat{f}_{2n+1,h}(y+\alpha) - \hat{f}_{2n+1,h}(y-\alpha) - 2 \sum_{k=0}^{n-1} \binom{2n+1}{2k} \alpha^{2(n-k)+1} (\psi_{s,h}(y))_{2k} \right).$$

Moreover, we introduce the notation $E_p = \|f_p - (\psi_{s,h})_p\|_{W^{k,\infty}}$. We will prove the statement that for all $\epsilon > 0$, there exists $h > 0$, such that for all $p \leq s$, we have

$$(28) \quad E_p \leq E_p^* := \frac{2^{p/2}(1+\alpha)^{(p^2+p)/2}}{\alpha^{p/2}} \cdot \epsilon.$$

We first note that choosing h as in Lemma 3.1 implies that

$$(29) \quad \max_{\substack{p \leq s, \\ p \text{ odd}}} E_p \leq \epsilon,$$

which proves the statement for p odd, since $(1 + \alpha)/\alpha \geq 1$. We will now prove (28) for even p using induction. First note that

$$(30) \quad E_2 \leq \frac{1}{6\alpha} \cdot 2\epsilon \leq E_2^*,$$

which proves the base step. To prove the induction step, let $n \in \mathbb{N}$ be such that $2n + 1 \leq s$ and $n > 1$, and we assume by the induction hypothesis that $E_{2k} \leq E_{2k}^*$ for all $k < n$. It then follows from (25) and (27), that

$$(31) \quad E_{2n} \leq \frac{1}{2\alpha(2n+1)} \left(E_{2n+1} + E_{2n+1} + 2 \sum_{k=1}^{n-1} \binom{2n+1}{2k} \alpha^{2(n-k)+1} E_{2k} \right).$$

Note that by the induction hypothesis and the fact that E_{2k}^* is monotonically increasing in k , we have $E_{2k} \leq E_{2k}^* \leq E_{2(n-1)}^*$. Using also (29), and the fact that $\epsilon \leq E_{2(n-1)}^*$, this allows us to estimate (31), by

$$(32) \quad \begin{aligned} E_{2n} &\leq \frac{1}{\alpha(2n+1)} \left(\max_{\substack{p \leq s, \\ p \text{ odd}}} E_p + \sum_{k=1}^{n-1} \binom{2n+1}{2k} \alpha^{2(n-k)+1} E_{2(n-1)}^* \right) \\ &\leq \frac{1}{\alpha} \left(E_{2(n-1)}^* + (1 + \alpha)^{2n+1} E_{2(n-1)}^* \right) \\ &\leq \frac{2}{\alpha} (1 + \alpha)^{2n+1} E_{2(n-1)}^*. \end{aligned}$$

Recalling the definition of $E_{2(n-1)}^*$, we obtain

$$(33) \quad E_{2n} \leq \frac{2}{\alpha} (1 + \alpha)^{2n+1} E_{2(n-1)}^* \leq \left(\frac{2}{\alpha} (1 + \alpha)^{2n+1} \right)^n \cdot \epsilon = E_{2n}^*.$$

This proves the claimed estimate (28) also for the case where $p = 2n$ is even, and therefore concludes the proof of (28).

Next, we optimize (28) by choosing the optimal value of α . Lemma A.2 proves that the optimal choice is $\alpha = 1/s$. We conclude that for any $\epsilon > 0$, there exists a shallow tanh neural network $\psi_{s,h}$, with width independent of ϵ , such that

$$(34) \quad \max_{p \leq s} \|f_p - (\psi_{s,h})_p\|_{W^{k,\infty}} \leq \sqrt{\epsilon} (2es)^{s/2} \epsilon.$$

Replacing $\epsilon \rightarrow \epsilon/\sqrt{\epsilon} (2es)^{s/2}$ recovers the claimed error bound in the statement of this lemma. To quantify the size of the weights, we observe that equation (27) reveals that the weight bound of Lemma 3.1 needs to be multiplied with a factor

$$(35) \quad \max_k s \binom{s}{2k} \left(\frac{1}{s} \right)^{s-2k} \leq s \sum_{j=0}^s \binom{s}{j} \left(\frac{1}{s} \right)^{s-j} \leq s \left(1 + \frac{1}{s} \right)^s = O(s),$$

where we used the binomial theorem. The weight bound can seen to be equal to

$$(36) \quad O \left(\epsilon^{-s/2} s (2\sqrt[4]{2es} \sqrt{2M}(s+2))^{s(s+3)} \right) = O \left(\epsilon^{-s/2} (\sqrt{M}(s+2))^{3s(s+3)/2} \right)$$

for small ϵ and large s . This proves the weight bound stated in the lemma.

Finally, we note that the constructed approximations indeed correspond to a shallow tanh neural network of the stated size. Indeed, one can see from the fact that σ is odd, equation (15), Lemma 3.1 and equation (27) that a shallow tanh

neural network suffices, where the values of the $3(s+1)/2$ neurons in the hidden layer are given by

$$(37) \quad \sigma \left(\left(\frac{s}{2} - i \right) h(y + \beta) \right) \quad \text{where } i = 0, 1, \dots, \frac{s-1}{2} \text{ and } \beta \in \{-\alpha, 0, \alpha\}.$$

□

Remark 3.3. *Combined with the Weierstrass approximation theorem [63], the preceding results show that any continuous function can be uniformly approximated in supremum norm on a compact interval by shallow tanh neural networks to arbitrary accuracy, as was already observed by e.g. [55]. Using the constructive proof of the Weierstrass approximation theorem based on Bernstein polynomials, one can even obtain a rate of convergence in terms of the width of the neural network and the modulus of continuity of the continuous function [14].*

Remark 3.4. *Note that one can also construct monomials with even powers directly, as was done in e.g. [20]. Indeed, there exists an $x \in \mathbb{R}$ such that $\tanh^{(p)}(x) \neq 0$ for all $p \in \mathbb{N}$, allowing us to use a neural network as in (16) for even p as well. However, a key difficulty lies in explicitly finding a function $\gamma : \mathbb{N} \rightarrow (0, \infty)$ such that $|\tanh^{(p)}(x)| \geq \gamma(p)$ for all $p \in \mathbb{N}$. It is unclear if such a function, which is quite essential when proving uniform bounds as in Lemma 3.2, can be constructed directly. Instead, our construction circumvents this issue and can be readily extended to other activation functions.*

3.2. Approximating multivariate polynomials. Next, we consider the approximation of multivariate polynomials using tanh neural networks. As an application, we will also present two different approximations of the multiplication operator.

First, recall the set $P_{n,q} = \{\alpha \in \mathbb{N}_0^q : |\alpha| = n\}$ from Section 2.1. The multinomial theorem implies that for $\alpha \in P_{n,q}$ and $\omega \in \mathbb{R}^q$, it holds that

$$(38) \quad \sum_{\beta \in P_{n,q}} \binom{n}{\beta} \frac{\alpha^\beta}{n^n} \omega^\beta = \left(\sum_{i=1}^q \frac{\alpha_i}{n} \omega_i \right)^n.$$

Now let $x \in \mathbb{R}^d$, set $q = d$ and $\omega = x$. Then the set $\{\omega^\beta : \beta \in P_{n,q}\}$ corresponds to the set of all d -variate monomials of total degree *equal to* n . Similarly, one can set $q = d+1$ and $\omega = (1, x)$, such that $\{\omega^\beta : \beta \in P_{n,q}\}$ corresponds to the set of all d -variate monomials of total degree *at most* n . It is the goal of this section to approximate these monomials ω^β using tanh neural networks. Notice however that the results from the previous section already allow us to approximate the right hand side of (38), as it is merely a composition of a linear map and a univariate monomial. Writing $b_\alpha = (\sum_i \alpha_i \omega_i / n)^n$ one can interpret (38) for every $\alpha \in P_{n,q}$ as the linear equation $\sum_\beta D_{\alpha,\beta} \omega^\beta = b_\alpha$, where

$$(39) \quad D_{\alpha,\beta} = \binom{n}{\beta} \frac{\alpha^\beta}{n^n},$$

which leads us to a linear system $\{\sum_\beta D_{\alpha,\beta} \omega^\beta = b_\alpha : \alpha \in P_{n,q}\}$ with as unknowns the monomials ω^β . Since the *Dyson matrix* $D = (D_{\alpha,\beta})_{\alpha,\beta \in P_{n,q}}$, where the order of rows and columns reflects the lexicographic order on $P_{n,q}$, is invertible [48], it is possible to write every monomial as a linear combination of the b_α 's. We will exploit this fact to construct approximations of multivariate polynomials and the multiplication $\prod_{i=1}^d x_i$ in particular.

Lemma 3.5. *Let $q, n \in \mathbb{N}$, $k \in \mathbb{N}_0$ and $M > 0$. Then for every $\epsilon > 0$, there exist a shallow tanh neural network $\Psi_{n,q} : [-M, M]^q \rightarrow \mathbb{R}^{|P_{n,q}|}$ of width $3 \lceil \frac{n+1}{2} \rceil |P_{n,q}|$ such that*

$$(40) \quad \max_{\beta \in P_{n,q}} \left\| \omega^\beta - (\Psi_{n,q}(\omega))_\beta \right\|_{W^{k,\infty}} \leq \epsilon.$$

Furthermore, the weights of the network scale as $O\left(\epsilon^{-n/2}(n(n+2))^{3(n+2)^2}\right)$ for small ϵ and large n .

Proof. From the previous section, we can see that approximating $b_\alpha = (\sum_i \alpha_i \omega_i / n)^n$ requires a shallow tanh subnetwork \widehat{b}_α of width $3 \lceil \frac{n+1}{2} \rceil$. As we require $|P_{n,q}|$ such subnetworks, the total network width can be summarized as $3 \lceil \frac{n+1}{2} \rceil |P_{n,q}|$. Now denote by $\widehat{\omega}^\beta$ the neural network approximation one obtains by solving the linear system $\{\sum_\beta D_{\alpha,\beta} \widehat{\omega}^\beta = \widehat{b}_\alpha : \alpha \in P_{n,q}\}$. We then set $(\Psi_{n,q}(\omega))_\beta := (\widehat{\omega}^\beta)_\beta$. Then it holds that

$$(41) \quad \left\| (\widehat{\omega}^\beta)_\beta - (\omega^\beta)_\beta \right\|_{W^{k,\infty}} \leq \left\| D^{-1} \right\|_\infty \left\| (\widehat{b}_\alpha)_\alpha - (b_\alpha)_\alpha \right\|_{W^{k,\infty}}.$$

Now define $h_\alpha(\omega) = \sum_i \alpha_i \omega_i / n$, then it holds that $\widehat{b}_\alpha - b_\alpha = (\widehat{f}_{n,h} - f_n) \circ h_\alpha$. It is easy to check that $\|h_\alpha\|_{W^{k,\infty}} \leq \max\{1, M\}$. Invoking Lemma A.7 then gives us

$$(42) \quad \left\| \widehat{b}_\alpha - b_\alpha \right\|_{W^{k,\infty}} \leq 16(e^2 k^4 q^2)^k \left\| \widehat{f}_{n,h} - f_n \right\|_{W^{k,\infty}} \max\{1, M\}^k.$$

In addition, Lemma A.3 provides us with the bound

$$(43) \quad \left\| D^{-1} \right\|_\infty \leq (n!)^3 |P_{n,q}|^2 2^n \leq \pi e^3 q^{2n} n^{3(n+1/2)},$$

where we used Stirling's approximation and Lemma 2.1. Now let $\epsilon > 0$. Combining the two obtained inequalities with Lemma 3.2 then proves the existence of $h > 0$ such that

$$(44) \quad \left\| (\widehat{\omega}^\beta)_\beta - (\omega^\beta)_\beta \right\|_{W^{k,\infty}} \leq \left\| D^{-1} \right\|_\infty \cdot 16(e^2 k^4 q^2)^k \max\{1, M\}^k \cdot \epsilon$$

where the weights of $\widehat{f}_{n,h}$ scale as $O\left(\epsilon^{-n/2}(\sqrt{M}(n+2))^{3n(n+3)/2}\right)$ for small ϵ and large n . We can now rescale ϵ such that

$$(45) \quad \left\| (\widehat{\omega}^\beta)_\beta - (\omega^\beta)_\beta \right\|_{W^{k,\infty}} \leq \epsilon.$$

As a consequence, the weights of $\Psi_{n,q}$ will scale as

$$(46) \quad O\left(\epsilon^{-n/2} \left\| D^{-1} \right\|_\infty^{n/2+1} \left(4(e k^2 q)^k \sqrt{1+M}\right)^n (\sqrt{M}(n+2))^{3n(n+2)/2}\right).$$

Note that $O\left(\left\| D^{-1} \right\|_\infty^{n/2+1}\right) = O\left((\pi e^3 q n)^{3(n+2)^2/2}\right)$ and that therefore a (conservative) upper bound of the weights of $\Psi_{n,q}$ is given by

$$(47) \quad O\left(\epsilon^{-n/2}(n(n+2))^{3(n+2)^2}\right)$$

for small ϵ and large n . \square

Corollary 3.6 (Approximation of multivariate monomials). *Let $d, s \in \mathbb{N}$, $k \in \mathbb{N}_0$ and $M > 0$. Then for every $\epsilon > 0$, there exist a shallow tanh neural network $\Phi_{s,d} : [-M, M]^d \rightarrow \mathbb{R}^{|P_{s,d+1}|}$ of width $3 \lceil \frac{s+1}{2} \rceil |P_{s,d+1}|$ such that*

$$(48) \quad \max_{\beta \in P_{s,d+1}} \sup_{x \in [-M, M]^d} \left\| x^\beta - (\Phi_{s,d}(x))_\beta \right\|_{W^{k,\infty}} \leq \epsilon.$$

Furthermore, the weights of the network scale as $O\left(\epsilon^{-s/2}(s(s+2))^{3(s+2)^2}\right)$ for small ϵ and large s .

Proof. The statement follows directly from Lemma 3.5 with $n \leftarrow s$, $q \leftarrow d + 1$ and $\omega \leftarrow (1, x)$, where $x \in [-M, M]^d$. \square

Next, we discuss how the multiplication operator can be approximated. To begin with, Lemma 3.5 shows that the multiplication of d numbers can easily be approximated using a shallow tanh neural network.

Corollary 3.7 (Shallow approximation of multiplication of d numbers). *Let $d \in \mathbb{N}$, $k \in \mathbb{N}_0$ and $M > 0$. Then for every $\epsilon > 0$, there exist a shallow tanh neural network $\widehat{\times}_d^\epsilon : [-M, M]^d \rightarrow \mathbb{R}$ of width $3 \left\lceil \frac{d+1}{2} \right\rceil |P_{d,d}|$ such that*

$$(49) \quad \left\| \widehat{\times}_d^\epsilon(x) - \prod_{i=1}^d x_i \right\|_{W^{k,\infty}} \leq \epsilon.$$

Furthermore, the weights of the network scale as $O(\epsilon^{-d/2})$ for small ϵ .

Proof. The statement follows directly from Lemma 3.5 with $n \leftarrow d$, $q \leftarrow d$ and $\omega \leftarrow x$, where $x \in [-M, M]^d$. \square

One issue with this shallow approximation is that the width of the network grows quickly with the dimension. The next lemma shows that the same accuracy can also be obtained using a deep tanh neural network for which both width and depth scale at most linearly with the input dimension.

Lemma 3.8 (Deep approximation of multiplication of d numbers). *Let $d \in \mathbb{N}$, $k \in \mathbb{N}_0$ and $M > 0$. Then for every $\epsilon > 0$, there exist a tanh neural network $\widehat{\times}_d^\epsilon : [-M, M]^d \rightarrow \mathbb{R}$ with $\lceil \log_2(d) \rceil$ hidden layers and of width at most $3d$ such that*

$$(50) \quad \left\| \widehat{\times}_d^\epsilon(x) - \prod_{i=1}^d x_i \right\|_{W^{k,\infty}} \leq \epsilon.$$

Furthermore, the weights of the network scale as $O(\epsilon^{-1/2})$ for small ϵ .

Proof. Using the finite difference approach (15), we can approximate the quadratic function using $\delta_h^2[f](x_0)$ for some $h > 0$ and $x_0 \in [-1, 1]$ such that $\sigma^{(2)}(x_0) \neq 0$. Observing that

$$(51) \quad xy = \frac{1}{4} \left((x+y)^2 - (x-y)^2 \right)$$

then provides a recipe to approximate (in Sobolev norm) the multiplication of two numbers using a shallow tanh neural network with 6 neurons in its hidden layer. The proof is similar to that of Lemma 3.1. Moreover, Lemma 3.1 shows as well that the identity can be approximated using a shallow tanh neural network with only one neuron in its hidden layer.

The multiplication of d numbers then follows easily from the multiplication of 2 numbers. In e.g. [54, Proposition 2.36], it is proven that the multiplication of d numbers requires a neural network in the form of a binary tree of depth $\lceil \log_2(d) \rceil$ where each node computes the (approximate) multiplication of two numbers. The proof of our error bound follows from Lemma A.6 and A.7. \square

Remark 3.9. *For simplicity and motivated by its widespread use, we only focused on the hyperbolic tangent activation function here. Our approach can be generalized to any activation function ϕ for which there exist $\mathcal{P} \subseteq \mathbb{N}$ with $\sup \mathcal{P} = \infty$ and an explicitly known function $\gamma : \mathcal{P} \rightarrow (0, \infty)$ with $|\phi^{(p)}| \geq \gamma(p)$ for all $p \in \mathcal{P}$. Monomials with degree $p \in \mathcal{P}$ can be constructed as in (16), the construction of monomials with degree $p \in \mathbb{N} \setminus \mathcal{P}$ is similar to the one described for multivariate polynomials.*

4. APPROXIMATION OF PARTITION OF UNITY

Once we have approximated polynomials with shallow tanh neural networks, the next step in our construction is to approximate a suitable partition of unity. In this section, we show how one can mimic a partition of unity using tanh neural networks. We recall that a partition of unity is a set of functions $f_i : [0, 1]^d \rightarrow [0, 1]$ such that every f_i is non-zero on only a small part of $[0, 1]^d$ and such that $\sum_i f_i = 1$. For ReLU and RePU neural networks, such partitions of unity can be constructed exactly [65]. For tanh neural networks, we will prove that an approximate partition of unity can be constructed. A unifying framework for approximating partitions of unity by general neural networks has been proposed in [20].

Let $d, N \in \mathbb{N}$ and $k \in \mathbb{N}_0$. For every $j \in \mathbb{N}^d$ with $\|j\|_\infty \leq N$ we define x_j^N such that $(x_j^N)_i = j_i/N_i$. We also define

$$(52) \quad I_j^N = \bigtimes_{i=1}^d ((j_i - 1)/N, j_i/N).$$

Let $R > 0$ be such that $\sigma^{(m)}$ is decreasing on $[R, \infty)$ for every $1 \leq m \leq k$. Given $\epsilon > 0$, we first find an $\alpha = \alpha(N, \epsilon)$ large enough such that

$$(53) \quad \alpha/N \geq R, \quad 1 - \sigma(\alpha/N) \leq \epsilon, \quad \alpha^m \left| \sigma^{(m)}(\alpha/N) \right| \leq \epsilon \text{ for all } 1 \leq m \leq k.$$

This is possible because $\lim_{x \rightarrow \infty} \sigma(x) = 1$ and because of Lemma A.4. In particular, Lemma A.5 shows that a suitable choice of α is given by

$$(54) \quad \alpha = N \max \left\{ R, \ln \left(\frac{(2k)^{k+1} (Nk)^k}{e^k \epsilon} \right) \right\}.$$

For $y \in \mathbb{R}$, we then define

$$(55) \quad \rho_1^N(y) = \frac{1}{2} - \frac{1}{2} \sigma \left(\alpha \left(y - \frac{1}{N} \right) \right),$$

$$(56) \quad \rho_j^N(y) = \frac{1}{2} \sigma \left(\alpha \left(y - \frac{j-1}{N} \right) \right) - \frac{1}{2} \sigma \left(\alpha \left(y - \frac{j}{N} \right) \right) \quad \text{for } 2 \leq j \leq N-1,$$

$$(57) \quad \rho_N^N(y) = \frac{1}{2} \sigma \left(\alpha \left(y - \frac{N-1}{N} \right) \right) + \frac{1}{2}.$$

In the remainder of the paper, we will assume for simplicity that ρ_j^N is always of the second form. The calculations involving ρ_1^N and ρ_N^N can be done entirely similarly and do not change the stated results. Finally, we define for $D \leq d$ the functions

$$(58) \quad \Phi_j^{N,D}(x) = \prod_{i=1}^D \rho_{j_i}^{N_i}(x_i)$$

and the sets $\mathcal{V}_D = \{v \in \mathbb{Z}^d : \max_{1 \leq i \leq D} |v_i| \leq 1 \text{ and } v_{D+1} = \dots = v_d = 0\}$. We will prove that the functions $\Phi_j^{N,d}$ approximate a partition of unity in the sense that for every j it holds on I_j^N that,

$$(59) \quad \sum_{v \in \mathcal{V}_d} \Phi_{j+v}^{N,d} \approx 1 \quad \text{and} \quad \sum_{\substack{v \notin \mathcal{V}_d, \\ j+v \in \{1, \dots, N\}^d}} \Phi_{j+v}^{N,d} \approx 0.$$

An example for $d = 1$ and $N = 7$ is shown in Figure 1. The next two lemmas formalize this approximation. Finally, a tanh neural network approximation of

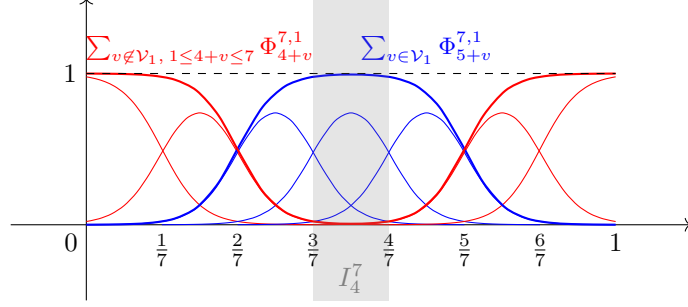


FIGURE 1. Example of an approximate partition of unity on $[0, 1]$ with $N = 7$. The thin lines represent the $\Phi_j^{7,1} = \rho_j^7$, $1 \leq j \leq 7$.

$\Phi_j^{N,d}$ can be constructed by replacing the multiplication operator by the network from e.g. Corollary 3.7 or Lemma 3.8.

Lemma 4.1. *If $0 < \epsilon < 1/4$, then*

$$(60) \quad \left\| \sum_{v \in \mathcal{V}_d} \Phi_{j+v}^{N,d} - 1 \right\|_{W^{k,\infty}(I_j^N)} \leq 2^{dk} d \epsilon.$$

Proof. We will prove the statement holds by induction on d . We first note that, for $d = 1$, we have

$$(61) \quad \begin{aligned} \sum_{v \in \mathcal{V}_1} \Phi_{j+v}^{N,1}(x) &= \sum_{l=-1}^1 \rho_{j_1+l}^N(x_1) \\ &= \frac{1}{2} \sigma \left(\alpha \left(x_1 - \frac{j_1 - 2}{N} \right) \right) - \frac{1}{2} \sigma \left(\alpha \left(x_1 - \frac{j_1 + 1}{N} \right) \right), \end{aligned}$$

from which easily follows that

$$(62) \quad \sum_{v \in \mathcal{V}_1} \Phi_{j+v}^{N,1}(x) \leq 1.$$

Next, note that for $x \in I_j^N$

$$(63) \quad \begin{aligned} \sum_{v \in \mathcal{V}_1} \Phi_{j+v}^{N,1}(x) &= \frac{1}{2} \sigma \left(\alpha \left(x_1 - \frac{j_1 - 2}{N} \right) \right) - \frac{1}{2} \sigma \left(\alpha \left(x_1 - \frac{j_1 + 1}{N} \right) \right) \\ &\geq \sigma \left(\frac{\alpha}{N} \right) \geq 1 - \epsilon, \end{aligned}$$

where we used the definition of α on the last line. Furthermore, for $1 \leq m \leq k$, we get that

$$(64) \quad \left| \frac{d^m}{dx^m} \sum_{v \in \mathcal{V}_1} \Phi_{j+v}^{N,1}(x) \right| \leq \alpha^m \sigma^{(m)} \left(\frac{\alpha}{N} \right) \leq \epsilon,$$

where we used (61) and the monotonic decay of $\sigma^{(m)}(x)$ for $x \in [\alpha/N, \infty)$ and our choice of α (cf. equation (53)). This allows us to conclude that

$$(65) \quad \left\| \sum_{v \in \mathcal{V}_1} \Phi_{j+v}^{N,1}(x) - 1 \right\|_{W^{k,\infty}(I_j^N)} \leq \epsilon.$$

For the induction step, we assume that for some $2 \leq D \leq d$ it holds that

$$(66) \quad \left\| \sum_{v \in \mathcal{V}_{D-1}} \Phi_{j+v}^{N,D-1} - 1 \right\|_{W^{k,\infty}(I_j^N)} \leq 2^{(D-1)k} (D-1)\epsilon.$$

Using Lemma A.6, we find that for $x \in I_j^N$,

$$(67) \quad \begin{aligned} \left\| \sum_{v \in \mathcal{V}_D} \Phi_{j+v}^{N,D}(x) - 1 \right\|_{W^{k,\infty}(I_j^N)} &= \left\| \sum_{w \in \mathcal{V}_1} \rho_{j_D+w}^N(x_D) \sum_{v \in \mathcal{V}_{D-1}} \Phi_{j+v}^{N,D-1}(x) - 1 \right\|_{W^{k,\infty}(I_j^N)} \\ &\leq \left\| \sum_{w \in \mathcal{V}_1} \rho_{j_D+w}^N(x_D) - 1 \right\|_{W^{k,\infty}(I_j^N)} \\ &\quad + 2^k \left\| \sum_{w \in \mathcal{V}_1} \rho_{j_D+w}^N(x_D) \right\|_{W^{k,\infty}(I_j^N)} \left\| \sum_{v \in \mathcal{V}_{D-1}} \Phi_{j+v}^{N,D-1}(x) - 1 \right\|_{W^{k,\infty}(I_j^N)} \\ &\leq \epsilon + 2^k 2^{(D-1)k} (D-1)\epsilon \leq 2^{Dk} D\epsilon. \end{aligned}$$

This concludes the proof. \square

Lemma 4.2. *Let $k \in \mathbb{N}_0$ and $v \in \mathbb{Z}^d$ with $\|v\|_\infty \geq 2$. Then it holds that*

$$(68) \quad \left\| \Phi_{j+v}^{N,d} \right\|_{W^{k,\infty}(I_j^N)} \leq \max\{1, (2k)^{2k} \alpha^k\} \epsilon.$$

Proof. Let $x \in I_j^N$ and let $1 \leq i \leq d$ be an index such that $|v_i| \geq 2$. Using some basic equalities for the hyperbolic tangent function and the definition of α , we obtain that

$$(69) \quad \begin{aligned} \left| \rho_{j_i+v_i}^N(x_i) \right| &\leq \frac{1}{2} \sigma\left(\frac{2\alpha}{N}\right) - \frac{1}{2} \sigma\left(\frac{\alpha}{N}\right) \\ &= \frac{1}{2} \sigma\left(\frac{\alpha}{N}\right) \left(1 - \sigma\left(\frac{2\alpha}{N}\right) \sigma\left(\frac{\alpha}{N}\right)\right) \\ &\leq \frac{1}{2} \left(1 - \sigma^2\left(\frac{\alpha}{N}\right)\right) \leq \epsilon. \end{aligned}$$

In addition, for every $1 \leq \ell \leq d$, it holds that $\left| \rho_{j_\ell+v_\ell}^N(x_\ell) \right| \leq 1$. This implies that

$$(70) \quad \left\| \Phi_{j+v}^{N,d} \right\|_{L^\infty(I_j^N)} \leq \epsilon.$$

Let $1 \leq m \leq k$, then it holds that (by our choice of the index i),

$$(71) \quad \left| \frac{d^m}{dx_i^m} \rho_{j_i+v_i}^N(x_i) \right| \leq \alpha^m \left| \sigma^{(m)}\left(\frac{\alpha}{N}\right) \right| \leq \epsilon.$$

Now let $\beta \in \mathbb{N}^d$ such that $1 \leq |\beta| \leq k$. Then

$$(72) \quad \left| D^\beta \Phi_{j+v}^{N,d}(x) \right| = \left| \prod_{\ell=1}^d \frac{d^{\beta_\ell}}{dx_\ell^{\beta_\ell}} \rho_{j_\ell+v_\ell}^{N_\ell}(x_\ell) \right| \leq \epsilon \prod_{\ell=1, \ell \neq i, \beta_\ell \neq 0}^d (2\beta_\ell)^{\beta_\ell+1} \alpha^{\beta_\ell} \leq \epsilon (2k)^{2k} \alpha^k,$$

where we used the fact that $\left| \sigma^{(m)}(x) \right| \leq (2m)^{m+1}$ in the first inequality (cf. Lemma A.4). Combining (70) and (72) proves the statement. \square

5. MAIN RESULTS

5.1. Approximation of functions in Sobolev spaces. We now present the first main result of the paper. It follows from the lemma of Bramble-Hilbert (Lemma A.8) that localized Taylor polynomials can approximate a function $f \in W^{s,\infty}([0,1]^d)$. For functions $f \in C^s([0,1]^d)$, this approximation follows from Taylor's theorem (Lemma A.9). We then use the results from the previous two sections to construct tanh neural networks that approximate localized Taylor polynomials in Sobolev norm. We prove that the function f can be approximated by a tanh neural network with two hidden layers and we provide explicit bounds on the width and approximation error.

Theorem 5.1. *Let $d, s, k \in \mathbb{N}$ with $s > k$, $R > 0$ as in (53), $\delta > 0$ and $f \in W^{s,\infty}([0,1]^d)$. There exists constants $\mathcal{C}(d, k, s, f), N_0(d) > 0$, such that for every $N \in \mathbb{N}$ with $N > N_0(d)$ there exists a tanh neural network \widehat{f}^N with two hidden layers, one of width at most $3 \lceil \frac{s}{2} \rceil |P_{s-1,d+1}| + d(N-1)$ and another of width at most $3 \lceil \frac{d+2}{2} \rceil |P_{d+1,d+1}| N^d$ (or $3 \lceil \frac{s}{2} \rceil + N - 1$ and $6N$ for $d = 1$), such that,*

$$(73) \quad \left\| f - \widehat{f}^N \right\|_{L^\infty([0,1]^d)} \leq (1 + \delta) \frac{\mathcal{C}(d, 0, s, f)}{N^s},$$

and for $k \geq 1$,

$$(74) \quad \left\| f - \widehat{f}^N \right\|_{W^{k,\infty}([0,1]^d)} \leq 3^d \left(1 + \frac{\delta}{3} \right) (2(k+1))^{3k} \max \left\{ R^k, \ln^k \left(\beta N^{s+d+2} \right) \right\} \frac{\mathcal{C}(d, k, s, f)}{N^{s-k}},$$

where we define

$$(75) \quad \beta = \frac{k^3 2^d \sqrt{d} \max\{1, \|f\|_{W^{k,\infty}([0,1]^d)}^{1/2}\}}{\delta \min\{1, \sqrt{\mathcal{C}(d, k, s, f)}\}}.$$

If $f \in C^s([0,1]^d)$, then it holds that

$$(76) \quad \mathcal{C}(d, k, s, f) = \max_{0 \leq \ell \leq k} \frac{1}{(s-\ell)!} \left(\frac{3d}{2} \right)^{s-\ell} |f|_{W^{s,\infty}([0,1]^d)}, \quad N_0(d) = \frac{3d}{2},$$

and else it holds that

$$(77) \quad \mathcal{C}(d, k, s, f) = \max_{0 \leq \ell \leq k} \frac{\pi^{1/4} \sqrt{s}}{(s-\ell-1)!} \left(5d^2 \right)^{s-\ell} |f|_{W^{s,\infty}([0,1]^d)}, \quad N_0(d) = 5d^2.$$

In addition, the weights of \widehat{f}^N scale as $O\left(\mathcal{C}^{-s/2} N^{s^2/2} (s(s+2))^{3s(s+2)}\right)$.

Proof. We will prove the theorem in the following manner. We divide the unit cube into N^d cubes of edge length $1/N$. On each of these cubes, f can be approximated in Sobolev norm by a polynomial. The global approximation can then be constructed by multiplying each polynomial with the indicator function of the corresponding cubes and summing over all cubes. We then prove that replacing these polynomials, multiplications and indicator functions with the tanh neural networks from the previous sections results in a new approximation that has approximately the same accuracy. In the last step we'll calculate the size of the required neural network.

Step 1: construction of the approximation. Let us denote $J_j^N = \times_{i=1}^d ((j_i - 2)/N, (j_i + 1)/N)$. We calculate that $\text{diam}(J_j^N) = \frac{3\sqrt{d}}{N}$ and that there exists a ball with diameter $\frac{1}{\sqrt{d}} \text{diam}(J_j^N)$ such that J_j^N is star-shaped with respect to every point in this ball. As a consequence, the Bramble-Hilbert lemma (Lemma

A.8) ensures the existence of a polynomial p_j^N of degree at most $s - 1$ such that

$$(78) \quad \begin{aligned} \left\| f - p_j^N \right\|_{W^{\ell, \infty}(J_j^N)} &\leq \frac{\pi^{1/4} \sqrt{s}}{(s - \ell - 1)!} \left(\frac{5d^2}{N} \right)^{s - \ell} |f|_{W^{s, \infty}([0, 1]^d)} \\ &\leq \max_{0 \leq m \leq \ell} \frac{\pi^{1/4} \sqrt{s} (5d^2)^{s - m}}{(s - m - 1)!} \frac{|f|_{W^{s, \infty}([0, 1]^d)}}{N^{s - \ell}} =: \frac{\mathcal{C}(d, \ell, s, f)}{N^{s - \ell}}, \end{aligned}$$

for all $0 \leq \ell \leq k$, under the assumption that $N > 5d^2$, and where we used that $3\sqrt{e} \leq 5$. If moreover $f \in C^s([0, 1]^d)$, then Taylor's theorem (Lemma A.9 with $\delta = \frac{3}{2N}$) ensures the existence of a polynomial p_j^N of degree at most $s - 1$ such that

$$(79) \quad \begin{aligned} \left\| f - p_j^N \right\|_{W^{\ell, \infty}(J_j^N)} &\leq \frac{1}{(s - \ell)!} \left(\frac{3d}{2N} \right)^{s - \ell} |f|_{W^{s, \infty}([0, 1]^d)} \\ &\leq \max_{0 \leq m \leq \ell} \frac{1}{(s - m)!} \left(\frac{3d}{2} \right)^{s - m} \frac{|f|_{W^{s, \infty}([0, 1]^d)}}{N^{s - \ell}} =: \frac{\mathcal{C}(d, \ell, s, f)}{N^{s - \ell}}, \end{aligned}$$

for all $0 \leq \ell \leq k$, under the assumption that $N > 3d/2$. The remainder of the argument will be independent of which polynomial p_j^N and which definition of $\mathcal{C}(d, \ell, s, f)$ is used. To simplify notation, we also define $p^N = \sum_j p_j^N \chi_j$, where χ_j denotes the indicator function on I_j^N . Next, let q_j^N be a tanh neural network as in Section 1 such that

$$(80) \quad \left\| q_j^N - p_j^N \right\|_{W^{k, \infty}([0, 1]^d)} \leq \eta.$$

In addition, we define

$$(81) \quad q_j^N(x) \widehat{\times} \Phi_j^{N, d}(x) := \widehat{\times}_{d+1}^h(q_j^N(x), \phi_{j_1}^{N, d}(x_1), \dots, \phi_{j_d}^{N, d}(x_d)),$$

where $\widehat{\times}_{d+1}^h$ is the network from Corollary 3.7 and $h = h(N)$ will be defined in the remainder of the proof. We then define our approximation as

$$(82) \quad \widehat{f}^N(x) = \sum_{j \in \{1, \dots, N\}^d} q_j^N(x) \widehat{\times} \Phi_j^{N, d}(x).$$

Step 2: estimating the error of the approximation. The triangle inequality gives us

$$(83) \quad \begin{aligned} \left\| f - \widehat{f}^N \right\|_{W^{k, \infty}([0, 1]^d)} &\leq \left\| f - \sum_{j \in \{1, \dots, N\}^d} f \cdot \Phi_j^{N, d} \right\|_{W^{k, \infty}([0, 1]^d)} + \left\| \sum_{j \in \{1, \dots, N\}^d} (f - q_j^N) \cdot \Phi_j^{N, d} \right\|_{W^{k, \infty}([0, 1]^d)} \\ &\quad + \left\| \sum_{j \in \{1, \dots, N\}^d} (q_j^N \cdot \Phi_j^{N, d} - q_j^N \widehat{\times} \Phi_j^{N, d}) \right\|_{W^{k, \infty}([0, 1]^d)} \end{aligned}$$

We proceed by bounding each term of the right hand side separately.

Step 2a: First term of (83). Let $i \in \{0, \dots, N\}^d$ be arbitrary. Recalling that $\mathcal{V}_d = \{v \in \mathbb{Z}^d : \|v\|_\infty \leq 1\}$, we observe for $k \geq 1$,

$$\begin{aligned}
(84) \quad & \left\| f - \sum_{j \in \{1, \dots, N\}^d} f \cdot \Phi_j^{N,d} \right\|_{W^{k,\infty}(I_i^N)} \leq 2^k \|f\|_{W^{k,\infty}(I_i^N)} \left\| 1 - \sum_{v \in \mathcal{V}_d} \Phi_{i+v}^{N,d} \right\|_{W^{k,\infty}(I_i^N)} \\
& \quad + 2^k \|f\|_{W^{k,\infty}(I_i^N)} \left\| \sum_{\substack{j \in \{1, \dots, N\}^d \\ j-i \notin \mathcal{V}_d}} \Phi_j^{N,d} \right\|_{W^{k,\infty}(I_i^N)} \\
& \leq 2^k \|f\|_{W^{k,\infty}(I_i^N)} (2^{kd} d \epsilon + N^d (2k)^{2k} \alpha^k \epsilon) \\
& \leq 2^k \|f\|_{W^{k,\infty}(I_i^N)} 2^{kd} d \epsilon \\
& \quad + 2^k \|f\|_{W^{k,\infty}(I_i^N)} N^d (2k)^{2k} N^k (k+1)^k \max \left\{ R^k, \ln^k \left(\frac{2Nk^2}{\epsilon^{\frac{1}{k+1}} e} \right) \right\} \epsilon \\
& \leq \delta (2(k+1))^{3k} \max \left\{ R^k, \ln^k \left(\frac{2Nk^2}{\epsilon^{\frac{1}{k+1}} e} \right) \right\} \frac{\mathcal{C}(d, k, s, f)}{N^{s-k}},
\end{aligned}$$

where we used Lemma A.6, Lemma 4.1, Lemma 4.2 and Lemma A.5, as well as a suitable definition of ϵ , satisfying

$$(85) \quad \epsilon \leq \frac{\delta \mathcal{C}(d, k, s, f)}{2^{(k+1)d} d N^{s+d} \|f\|_{W^{k,\infty}([0,1]^d)}}.$$

Analogously, for $k = 0$, one can obtain that

$$(86) \quad \left\| f - \sum_{j \in \{1, \dots, N\}^d} f \cdot \Phi_j^{N,d} \right\|_{L^\infty(I_i^N)} \leq \|f\|_{L^\infty(I_i^N)} (d\epsilon + N^d \epsilon) \leq \frac{\delta \mathcal{C}(d, k, s, f)}{3 N^{s-k}}.$$

Step 2b: Second term of (83) for $k = 0$. In order to bound the second term, we first make some auxiliary calculations. To begin with, we consider the case where $k = 0$. We find that

$$\begin{aligned}
(87) \quad & \left| \sum_{v \in \mathcal{V}_d} (f - q_{i+v}^N) \Phi_{i+v}^{N,d} \right| \leq \max_{v \in \mathcal{V}_d} |f - q_{i+v}^N| \left(\sum_{v \in \mathcal{V}_d} |\Phi_{i+v}^{N,d}| \right) \\
& = \max_{v \in \mathcal{V}_d} |f - q_{i+v}^N| \left| \sum_{v \in \mathcal{V}_d} \Phi_{i+v}^{N,d} \right|
\end{aligned}$$

where all functions are evaluated at some $x \in I_i^N$. We can then use the bounds

$$(88) \quad |f - q_{i+v}^N| \leq \frac{\mathcal{C}(d, 0, s, f)}{N^s} + \eta,$$

which follows from (78) and (80), and,

$$(89) \quad \left| \sum_{v \in \mathcal{V}_d} \Phi_{i+v}^{N,d} \right| \leq 1 + d\epsilon,$$

which follows from Lemma 4.1. As a consequence, we find that

$$(90) \quad \left\| \sum_{v \in \mathcal{V}_d} (f - q_{i+v}^N) \Phi_{i+v}^{N,d} \right\|_{L^\infty(I_i^N)} \leq \left(\frac{\mathcal{C}(d, k, s, f)}{N^{s-k}} + \eta \right) (1 + d\epsilon).$$

Combining this result with the triangle inequality, (78), (80) and Lemma 4.2, we find that

$$(91) \quad \begin{aligned} & \left\| \sum_{j \in \{1, \dots, N\}^d} (f - q_j^N) \cdot \Phi_j^{N,d} \right\|_{L^\infty(I_i^N)} \\ & \leq \left\| \sum_{v \in \mathcal{V}_d} (f - q_{i+v}^N) \Phi_{i+v}^{N,d} \right\|_{L^\infty(I_i^N)} + \sum_{\substack{j \in \{1, \dots, N\}^d \\ j-i \notin \mathcal{V}_d}} \left\| (f - q_j^N) \right\|_{L^\infty(I_i^N)} \left\| \Phi_j^{N,d} \right\|_{L^\infty(I_i^N)} \\ & \leq \left(\frac{\mathcal{C}(d, k, s, f)}{N^{s-k}} + \eta \right) (1 + d\epsilon) + N^d (\mathcal{C}(d, k, s, f) + \eta) \epsilon \\ & \leq \left(1 + \frac{\delta}{3} \right) \frac{\mathcal{C}(d, k, s, f)}{N^{s-k}}. \end{aligned}$$

where we obtain the last inequality by making a suitable choice of η and ϵ .

Step 2c: Second term of (83) for $k \geq 1$. Next we consider the case where $0 < k < s$. Let $\beta \in \mathbb{N}_0^d$ be such that $|\beta| \leq k$. Then as a consequence of the general Leibniz rule we find that

$$(92) \quad D^\beta \left(\sum_{v \in \mathcal{V}_d} (f - q_{i+v}^N) \Phi_{i+v}^{N,d} \right) \leq \sum_{\beta' \leq \beta} \binom{\beta}{\beta'} \sum_{v \in \mathcal{V}_d} \left| D^{\beta'} (f - q_{i+v}^N) \right| \left| D^{\beta - \beta'} \Phi_{i+v}^{N,d} \right|$$

where all functions are evaluated at some $x \in I_i^N$. For every $v \in \mathcal{V}_d$ and $\beta' \leq \beta$ with $\ell := |\beta - \beta'|$, we can then use the bounds

$$(93) \quad \left| D^{\beta'} (f - q_{i+v}^N) \right| \leq \left\| f - q_{i+v}^N \right\|_{W^{k-\ell, \infty}(I_i^N)} \leq \frac{\mathcal{C}(d, k - \ell, s, f)}{N^{s-k+\ell}} + \eta,$$

which follows from (78) and (80), and,

$$(94) \quad \left| D^{\beta - \beta'} \Phi_{i+v}^{N,d} \right| \leq \alpha^\ell (2\ell)^{2\ell} = N^\ell (2\ell)^{2\ell} \max \left\{ R^\ell, \ln^\ell \left(\frac{(2k)^{k+1} (Nk)^k}{e^k \epsilon} \right) \right\},$$

which follows from Lemma A.4 and Lemma A.5. As $\sum_{\beta' \leq \beta} \binom{\beta}{\beta'} \leq 2^k$ (as a consequence of the multi-binomial theorem), we find that

$$(95) \quad \left\| \sum_{v \in \mathcal{V}_d} (f - q_{i+v}^N) \Phi_{i+v}^{N,d} \right\|_{W^{k, \infty}(I_i^N)} \leq 2^k 3^d \left(\frac{\mathcal{C}(d, k, s, f)}{N^{s-k}} + \eta N^k \right) (2k)^{2k} \max \left\{ R^k, \ln^k \left(\frac{(2Nk^2)^{k+1}}{\epsilon e^k} \right) \right\}.$$

Combining this result with the triangle inequality, Lemma A.5, Lemma A.6, (78), (80), Lemma 4.2 and the fact that $\ln(x) \leq \sqrt{x}$ for $x > 0$, we find that

$$\begin{aligned}
(96) \quad & \left\| \sum_{j \in \{1, \dots, N\}^d} (f - q_j^N) \cdot \Phi_j^{N,d} \right\|_{W^{k,\infty}(I_i^N)} \\
& \leq \left\| \sum_{v \in \mathcal{V}_d} (f - q_{i+v}^N) \Phi_{i+v}^{N,d} \right\|_{W^{k,\infty}(I_i^N)} + \sum_{\substack{j \in \{1, \dots, N\}^d \\ j-i \notin \mathcal{V}_d}} \left\| (f - q_j^N) \Phi_j^{N,d} \right\|_{W^{k,\infty}(I_i^N)} \\
& \leq \left\| \sum_{v \in \mathcal{V}_d} (f - q_{i+v}^N) \Phi_{i+v}^{N,d} \right\|_{W^{k,\infty}(I_i^N)} + \sum_{\substack{j \in \{1, \dots, N\}^d \\ j-i \notin \mathcal{V}_d}} 2^k \left\| (f - q_j^N) \right\|_{W^{k,\infty}(I_i^N)} \left\| \Phi_j^{N,d} \right\|_{W^{k,\infty}(I_i^N)} \\
& \leq 2^k 3^d \left(\frac{\mathcal{C}(d, k, s, f)}{N^{s-k}} + \eta N^k \right) (2k)^{2k} \max \left\{ R^k, \ln^k \left(\frac{(2Nk^2)^{k+1}}{\epsilon e^k} \right) \right\} \\
& \quad + N^d 2^k (\mathcal{C}(d, k, s, f) + \eta) (2k)^{2k} N^k (k+1)^k \left(\frac{2Nk^2}{e} \right)^{k/2} \sqrt{\epsilon} \\
& \leq 3^d \left(1 + \frac{\delta}{3} \right) (2(k+1))^{3k} \max \left\{ R^k, \ln^k \left(\frac{2Nk^2}{\epsilon^{\frac{1}{k+1}} e} \right) \right\} \frac{\mathcal{C}(d, k, s, f)}{N^{s-k}},
\end{aligned}$$

where we obtain the last inequality by making a suitable choice of η and ϵ , satisfying

$$(97) \quad \eta \leq \frac{\delta \mathcal{C}}{6N^s} \quad \text{and} \quad \epsilon \leq \frac{\delta^2}{N^{2s+2d+k} k^k},$$

where we assumed that $0 < \delta < 5/6$.

Step 2d: Third term of (83). Finally, using the triangle inequality, Lemma A.7, Lemma 3.5 and Lemma A.4 we obtain that for some $C_k > 0$ depending only on k ,

$$\begin{aligned}
(98) \quad & \left\| \sum_{j \in \{1, \dots, N\}^d} (q_j^N \cdot \Phi_j^{N,d} - q_j^N \widehat{\times} \Phi_j^{N,d}) \right\|_{W^{k,\infty}(I_i^N)} \\
& \leq N^d C_k (d+1)^d d^{2k} \cdot \left\| \widehat{\times}_{d+1}^h - \prod_{i=1}^{d+1} x_i \right\|_{W^{k,\infty}} \left(\|f\|_{W^{k,\infty}([0,1]^d)} + \|f - q_j^N\|_{W^{k,\infty}([0,1]^d)} + \|\rho_i^N\|_{W^{k,\infty}([0,1]^d)} \right) \\
& \leq N^d C_k (d+1)^d d^{2k} \cdot h \left(\|f\|_{W^{k,\infty}([0,1]^d)} + \frac{\mathcal{C}(d, k, s, f)}{N^{s-k}} + \eta + (2\alpha k)^{k+1} \right) \\
& \leq 2^k \frac{\delta \mathcal{C}(d, k, s, f)}{3 N^{s-k}},
\end{aligned}$$

where we obtain the last inequality by making a suitable choice of h and η , satisfying

$$(99) \quad \eta \leq \|f\|_{W^{k,\infty}([0,1]^d)} \quad \text{and} \quad h \leq \frac{2^k \delta \mathcal{C}}{3N^{d+s-k} C_k (d+1)^d d^{2k} (2\|f\|_{W^{k,\infty}([0,1]^d)} + \mathcal{C} + (2\alpha k)^{k+1})}.$$

Step 2e: Final error bound. As i was chosen arbitrary, combining the contributions from the three terms of (83) then proves that

$$(100) \quad \left\| f - \widehat{f}^N \right\|_{L^\infty([0,1]^d)} \leq (1 + \delta) \frac{\mathcal{C}(d, 0, s, f)}{N^s}.$$

Moreover, from (85) and (97) we find that a suitable definition of ϵ is given by

$$(101) \quad \epsilon = \frac{\delta^2 \min\{1, \mathcal{C}(d, k, s, f)\}}{N^{2s+2d+k} k^k 2^{(k+1)d} d \max\{1, \|f\|_{W^{k, \infty}([0,1]^d)}\}},$$

from which it follows that for $k \geq 1$,

$$(102) \quad \epsilon^{\frac{1}{k+1}} \geq \frac{\delta \min\{1, \sqrt{\mathcal{C}(d, k, s, f)}\}}{N^{s+d+1} k 2^d \sqrt{d} \max\{1, \|f\|_{W^{k, \infty}([0,1]^d)}^{1/2}\}}.$$

Combining this observation with all previous steps of the proof then leads to the error bound

$$(103) \quad \|f - \widehat{f}^N\|_{W^{k, \infty}([0,1]^d)} \leq 3^d \left(1 + \frac{\delta}{3}\right) (2(k+1))^{3k} \max\left\{R^k, \ln^k\left(\beta N^{s+d+2}\right)\right\} \frac{\mathcal{C}(d, k, s, f)}{N^{s-k}},$$

for $k \geq 1$ and where we define

$$(104) \quad \beta = \frac{k^3 2^d \sqrt{d} \max\{1, \|f\|_{W^{k, \infty}([0,1]^d)}^{1/2}\}}{\delta \min\{1, \sqrt{\mathcal{C}(d, k, s, f)}\}}.$$

Step 3: Estimating the network and weights sizes. The first hidden layer requires $3 \lceil \frac{s}{2} \rceil |P_{s-1, d+1}|$ neurons for the computation of all multivariate monomials (cf. Lemma 3.5). For $d = 1$, the result follows from Lemma 3.2 instead of Lemma 3.5. For the computation of all $\rho_j^N(x_i)$ another $d(N-1)$ neurons are needed in the first hidden layer. The second hidden layer needs at most $3 \lceil \frac{d+2}{2} \rceil |P_{d+1, d+1}|$ neurons for realizing $\widehat{\times}_{d+1}^h$, which needs to be performed N^d times. For $d = 1$, six neurons are sufficient to approximate the multiplication (see (51)).

In the proof we achieved the wanted accuracy by making suitable choices of η, ϵ, h . From equation (100) and Lemma A.5, it follows that

$$(105) \quad \alpha = O(Ns \ln(\mathcal{C}N)).$$

For the approximate multiplication, (99) requires that $h^{-1} = O(N^{2s+d} s^{k+1})$. Corollary 3.7 then proves that the weights of $\widehat{\times}_{d+1}^h$ grow as $O(N^{d(s+d/2)} s^{d(k+1)/2})$. Finally, the condition $\eta^{-1} = O(\mathcal{C}^{-1} N^s)$ from (97) corresponds to weights growing as

$$(106) \quad O\left(\mathcal{C}^{-s/2} N^{s^2/2} (s(s+2))^{3s(s+2)}\right)$$

as a consequence of Corollary 3.6. This concludes the proof. \square

Remark 5.2. *The result of Theorem 5.1 can be generalized to functions $f \in W^{k,p}(\Omega)$ for $p < \infty$. For this, a slightly more general version of Lemmas A.6 and A.8 is needed. The convergence rate will still be as in Theorem 5.1, only the constant \mathcal{C} will be different.*

Remark 5.3. *Recently, it has been shown that the curse of dimensionality can be lessened for functions in so-called Korobov spaces [49, 35, 4]. In particular, in [35, Theorem 4.2], this framework is used to show how ReQU neural networks can approximate a $C^k([-1,1]^d)$ -function to an accuracy of $\epsilon > 0$ in supremum norm with at most $O\left(\frac{d}{k} \ln\left(\frac{1}{\epsilon}\right) + d\right)$ hidden layers and at most $O\left(\epsilon^{-\frac{1+\delta}{k}} \left(\frac{1+\delta}{k} \ln\left(\frac{1}{\epsilon}\right)\right)^{d-1}\right)$ neurons and non-zero weights. As their proof builds upon the mimicking of polynomials, it is clear from our results that similar approximation rates can be obtained using tanh neural networks.*

One particularly useful consequence of Theorem 5.1 is that it provides an explicit error bound on the approximation of Lipschitz functions using tanh neural networks.

Corollary 5.4. *Let $d \in \mathbb{N}$ and let $f : [0, 1]^d \rightarrow \mathbb{R}$ be a Lipschitz continuous function with Lipschitz constant $L > 0$. For every $N \in \mathbb{N}$ with $N > 5d^2$ there exists a tanh neural network \widehat{f}^N with two hidden layers of widths at most $d(N - 1)$ and $3 \left\lceil \frac{d+1}{2} \right\rceil |P_{d,d}| N^d$ (or $N - 1$ and $6N$ for $d = 1$), such that*

$$(107) \quad \left\| f - \widehat{f}^N \right\|_{L^\infty([0,1]^d)} \leq \frac{7d^2 L}{N}.$$

Moreover, the weights of \widehat{f}^N scale as $O(\sqrt{N})$.

Proof. The corollary follows directly from Theorem 5.1 by setting $k = 0$, $s = 1$, choosing $\delta > 0$ in such a way that $5(1 + \delta)\pi^{1/4} \leq 7$ and observing that $|f|_{W^{1,\infty}} \leq L$ because of the Lipschitz continuity of f . The constructed network in Theorem 5.1 is based on localized $(s - 1)$ -th order polynomials. For $s = 1$ this corresponds to constant functions, thereby removing the need to mimic monomials. As a consequence, the network width can be simplified to the widths stated in the corollary. \square

5.2. Approximation of analytic functions. We now investigate how we can apply Theorem 5.1 to analytic functions. As the class of analytic functions coincides with the Gevrey class G^1 , it follows that for every analytic function there exists a constant $C_f > 0$ such that

$$(108) \quad |f|_{W^{s,\infty}([0,1]^d)} \leq C_f^{s+1} s! \text{ for all } s \in \mathbb{N}_0.$$

A related concept is the class of (Q, R) -analytic functions [15, 6, 7], where $Q, R > 0$, consisting of analytic functions for which the following smoothness condition holds,

$$(109) \quad |f|_{W^{s,\infty}([0,1]^d)} \leq QR^{-s} s! \text{ for all } s \in \mathbb{N}_0.$$

Note that any analytic function is (C_f, C_f^{-1}) -analytic by the previous characterization of analyticity. Hence, a function is analytic on some compact interval if and only if it is (Q, R) -analytic for some $Q, R > 0$ on that interval. The following corollaries discuss multiple ways to approximate analytic functions using tanh neural networks. All constants mentioned in the statements can be easily calculated from the proof.

We start with the basic consequence of Theorem 5.1 for (Q, R) -analytic functions. It provides explicit estimates on both the approximation error in supremum norm and the network size. It can easily be generalized to Sobolev norm using Theorem 5.1.

Corollary 5.5. *Let $d \in \mathbb{N}$, $\delta, Q, R > 0$, $\Omega \subset \mathbb{R}^d$ open with $[0, 1]^d \subset \Omega$ and let f be (Q, R) -analytic on Ω . Then for every $s \in \mathbb{N}_0, N \in \mathbb{N}$ with $N > 3d/2$, there is a tanh neural network $\widehat{f}^{N,s}$ with two hidden layers of widths at most $3 \left\lceil \frac{s}{2} \right\rceil |P_{s-1,d+1}| + d(N - 1)$ and $3 \left\lceil \frac{d+2}{2} \right\rceil |P_{d+1,d+1}| N^d$ (or $3 \left\lceil \frac{s}{2} \right\rceil + N - 1$ and $6N^d$ for $d = 1$) such that*

$$(110) \quad \left\| f - \widehat{f}^{N,s} \right\|_{L^\infty([0,1]^d)} \leq (1 + \delta) Q \left(\frac{3d}{2RN} \right)^s.$$

Moreover, if $R > d/2$ then for every $s \in \mathbb{N}_0$ there is a shallow tanh neural network \widehat{f}^s with at most $\frac{3s}{2} |P_{s-1,d+1}|$ (or $3 \left\lceil \frac{s}{2} \right\rceil$ for $d=1$) in its hidden layer such that

$$(111) \quad \left\| f - \widehat{f}^s \right\|_{W^{k,\infty}([0,1]^d)} \leq (1 + \delta) Q \left(\frac{d}{2R} \right)^s.$$

Proof. The first part of the statement follows directly from Theorem 5.1 by taking $\delta = 1/3$. The second part follows from taking $N = 1$ in Theorem 5.1 and observing that the proof can be simplified in this case. Indeed, one can then directly use

Taylor's theorem (Lemma A.9) with $\delta = \frac{1}{2}$ instead of $\delta = \frac{3}{2N}$ and there is no more need for an approximate partition of unity, thereby also removing the need for a second hidden layer. \square

The following corollary enables a consistent comparison with the available literature, as it bounds the approximation error in terms of one single parameter. Whereas other papers focus on the number of non-zero weights and biases as complexity measure, we opted for the network width. This is useful in practice as the network can be directly chosen, whereas it is very challenging to exactly control the sparsity of the neural network (i.e. the number of non-zero weights and biases). Moreover, many bounds on the generalization error require an estimate of the network width [3, 25].

Corollary 5.6. *Let $d \in \mathbb{N}$, $k \in \mathbb{N}_0$, $\delta, Q, R > 0$, $\Omega \subset \mathbb{R}^d$ open with $[0, 1]^d \subset \Omega$ and let f be a (Q, R) -analytic on Ω . Then there exists a constant $c_{d,k,\alpha,f} > 0$ such that for every $\mathcal{N} \in \mathbb{N}$ there exists a tanh neural network $\hat{f}^{\mathcal{N}}$ with two hidden layers of width at most $O(\mathcal{N})$ for $\mathcal{N} \rightarrow \infty$ such that*

$$(112) \quad \left\| f - \hat{f}^{\mathcal{N}} \right\|_{W^{k,\infty}([0,1]^d)} \leq c_{d,k,\alpha,f} \mathcal{N}^{\frac{k}{d+1}} \exp\left(-\alpha \mathcal{N}^{\frac{1}{d+1}} \log(\mathcal{N})\right) \leq \frac{c_{d,k,\alpha,f}}{\mathcal{N}^{\alpha-k/(d+1)}}.$$

In particular, for $k = 0$ it holds that

$$(113) \quad \left\| f - \hat{f}^{\mathcal{N}} \right\|_{L^\infty([0,1]^d)} \leq (1 + \delta)Q \cdot \exp\left(-\alpha \mathcal{N}^{\frac{1}{d+1}} \log(\mathcal{N})\right) \leq \frac{(1 + \delta)Q}{\mathcal{N}^\alpha}.$$

Moreover, the weights $\hat{f}^{\mathcal{N}}$ of grow as $O\left(\mathcal{N}^{\frac{4(k+\alpha d \mathcal{N}^{1/(d+1)})}{d+1}}\right)$ for $\mathcal{N} \rightarrow \infty$.

Proof. First we observe that for every $\gamma > 0$ it holds that

$$(114) \quad \ln^k\left(\beta N^{s+d+2}\right) = O\left(\left(\frac{2NR}{3d}\right)^{\gamma s}\right)$$

for large s and N . From Theorem 5.1 we then find that for every N and s there is a network $\hat{f}^{N,s}$ such that

$$(115) \quad \left\| f - \hat{f}^{N,s} \right\|_{W^{k,\infty}([0,1]^d)} = O\left(\left(\frac{s}{R}\right)^k \left(\frac{3d}{2NR}\right)^{s(1-\gamma)-k}\right)$$

From Theorem 5.1 with the choices $s = k + \alpha(1-\gamma)^{-1}(d+1)\mathcal{N}^{\frac{1}{d+1}}$ and $N = \frac{3d}{2R}\mathcal{N}^{\frac{1}{d+1}}$ for some $\mathcal{N} \in \mathbb{N}$, gives that there exists a constant $c_{d,k,\alpha,f} > 0$ such that

$$(116) \quad \begin{aligned} \left\| f - \hat{f}^{\mathcal{N}} \right\|_{W^{k,\infty}([0,1]^d)} &\leq c_{d,k,\alpha,f} \mathcal{N}^{\frac{k}{d+1}} \left(\frac{1}{\mathcal{N}^{\frac{1}{d+1}}}\right)^{\alpha(d+1)\mathcal{N}^{\frac{1}{d+1}}} \\ &= c_{d,k,\alpha,f} \mathcal{N}^{\frac{k}{d+1}} \exp\left(-\alpha \mathcal{N}^{\frac{1}{d+1}} \log(\mathcal{N})\right) \\ &\leq \frac{c_{d,k,\alpha,f}}{\mathcal{N}^{\alpha-k/(d+1)}}. \end{aligned}$$

In particular, for $k = 0$ we find that

$$(117) \quad \left\| f - \hat{f}^{\mathcal{N}} \right\|_{L^\infty([0,1]^d)} \leq (1 + \delta)Q \cdot \exp\left(-\alpha \mathcal{N}^{\frac{1}{d+1}} \log(\mathcal{N})\right) \leq \frac{(1 + \delta)Q}{\mathcal{N}^\alpha}.$$

Using Lemma 2.1, we find that the network widths are respectively $O((e\alpha d)^{d+1}\mathcal{N})$ and $O(d5^d \mathcal{N}^{\frac{d}{d+1}})$ for large \mathcal{N} and d (the exact sizes can be easily calculated). Finally, from Theorem 5.1 follows as well that the weights grow as $O\left(\mathcal{N}^{\frac{4(k+\alpha d \mathcal{N}^{1/(d+1)})}{d+1}}\right)$. \square

We thus find that tanh neural networks with two hidden layers result in an exponential convergence rate. Moreover, the above corollary shows that a convergence rate that is independent of the dimension can be obtained, thereby lessening the curse of dimensionality. The proof however shows that even though the rate is free of the curse of dimensionality, the constant implied in the Landau notation still depends (super)exponentially on the dimension. Similar papers observe the same phenomena [54], or do not discuss this.

Remark 5.7. *One can also restate the previous corollary by saying that an approximation rate of $O(\mathcal{N}^k \exp(-\mathcal{N}))$ can be obtained using a tanh neural network with two hidden layers of widths $O\left(\mathcal{N}^{\binom{\mathcal{N}+d}{\mathcal{N}}}\right)$ and $O(1)$ for $\mathcal{N} \rightarrow \infty$. Since $O\left(\mathcal{N}^{\binom{\mathcal{N}+d}{\mathcal{N}}}\right)$ grows asymptotically slower than $O(\mathcal{N}^d)$, another (very modest) lessening of the curse of dimensionality is revealed.*

Next, we show that, under an additional assumption, shallow tanh neural networks can also approximate analytic functions at an exponential rate. Moreover, in contrast to Corollary 5.6, there are no hidden constants that grow as $O(d^d)$. For simplicity, we restrict ourselves to approximation in supremum norm (i.e. $k = 0$).

Corollary 5.8. *Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$ open with $[0, 1]^d \subset \Omega$ and let f be analytic on Ω . If f satisfies for some $C > 0$ that $|f|_{W^{s,\infty}([0,1]^d)} \leq C^s$ for all $s \in \mathbb{N}$, then for every $\mathcal{N} \in \mathbb{N}$ there exists a shallow tanh neural network $\hat{f}^{\mathcal{N}}$ of width $3 \left\lceil \frac{\mathcal{N}+5Cd}{2} \right\rceil \binom{\mathcal{N}+(5C+1)d}{\mathcal{N}+5Cd}$ (or $3 \left\lceil \frac{\mathcal{N}}{2} \right\rceil$ for $d = 1$) such that*

$$(118) \quad \left\| f - \hat{f}^{\mathcal{N}} \right\|_{L^\infty([0,1]^d)} \leq \exp(-\mathcal{N}).$$

Proof. Assume that f satisfies for some $C > 0$ that $|f|_{W^{s,\infty}([0,1]^d)} \leq C^s$ for all $s \in \mathbb{N}$. We calculate that for $\rho > 1$,

$$(119) \quad \frac{C^s}{s!} \left(\frac{3d\rho}{2} \right)^s = \frac{1}{s!} \left(\frac{3Cd\rho}{2} \right)^s \leq \frac{1}{\sqrt{2\pi}} \left(\frac{3Cde\rho}{2s} \right)^s \leq \frac{1}{\sqrt{2\pi}} e^{3Cd\rho/2},$$

where we used Stirling's approximation and maximized over all s . This proves that f is (Q, R) -analytic with $Q = \frac{1}{\sqrt{2\pi}} e^{3Cd\rho/2}$ and $R = \frac{3d\rho}{2}$. Using Corollary 5.5 with $s = \mathcal{N}$ and $N = 1$ gives us that

$$(120) \quad \left\| f - \hat{f}^{\mathcal{N}} \right\|_{L^\infty([0,1]^d)} \leq 2 \frac{1}{\sqrt{2\pi}} e^{3Cd\rho/2} \rho^{-\mathcal{N}}.$$

If we set $\rho = e$, then $e^{3d\rho/2} \leq e^{5d}$. Therefore it holds that

$$(121) \quad \left\| f - \hat{f}^{\mathcal{N}} \right\|_{L^\infty([0,1]^d)} \leq \exp(-\mathcal{N} + 5Cd).$$

Note that since now $N = 1$, the network architecture is even simpler: there is no need to construct a partition of unity, nor does there need to be a second hidden layer in order to approximately multiply the results of subnetworks. Therefore, a shallow tanh neural network with $3 \left\lceil \frac{\mathcal{N}}{2} \right\rceil |P_{\mathcal{N}-1, d+1}| \leq 3 \left\lceil \frac{\mathcal{N}}{2} \right\rceil \binom{\mathcal{N}+d}{\mathcal{N}}$ neurons in its hidden layer suffices. The statement from the theorem is obtained by making the substitution $\mathcal{N} \leftarrow \mathcal{N} + 5Cd$. \square

Finally, we discuss how dimension-independent convergence rates can be obtained for a class of countably-parametric, holomorphic maps $f : U := [-1, 1]^{\mathbb{N}} \rightarrow \mathbb{R}$, which arise in e.g. elliptic PDEs with uncertain coefficients. This was first discussed in [60] for deep ReLU neural networks and we will show that their results

can be adapted to hold for shallow tanh neural networks. More precisely, their results hold for functions u that admit a representation as a sparse Taylor generalized polynomial chaos expansion

$$(122) \quad f(y) = \sum_{\nu \in \mathcal{F}} \frac{D^\nu f(0)}{\nu!} y^\nu,$$

which is unconditionally convergent for $y \in U$ and where \mathcal{F} is defined by

$$(123) \quad \mathcal{F} = \{\nu \in \mathbb{N}_0^{\mathbb{N}} \mid \nu_j \neq 0 \text{ for only finitely many } j\}.$$

For a multi-index $\nu \in \mathbb{N}_0^{\mathbb{N}}$, we denote by $\text{supp}(\nu) = \{j \in \mathbb{N} \mid \nu_j \neq 0\}$ the support of ν , and we denote by $|\nu| = \sum_{j \in \mathbb{N}} |\nu_j|$ the ℓ^1 -norm of ν .

It is shown in [60, Section 2] that f admits such a representation if f is (b, ϵ) -holomorphic for $b \in \ell^p(\mathbb{N})$, $p \in (0, 1]$ and $\epsilon > 0$. The notion of (b, ϵ) -holomorphy is defined as follows.

Definition 5.9 (Def. 2.1 in [60]). *Let V be a Banach space. Let $b \in \ell^p(\mathbb{N})$, $p \in (0, 1]$ be a monotonically decreasing sequence. A poly-radius $\rho \in [1, \infty)^{\mathbb{N}}$ is called (b, ϵ) -admissible for some $\epsilon > 0$ if*

$$(124) \quad \sum_{j \in \mathbb{N}} b_j (\rho_j - 1) \leq \epsilon.$$

A continuous function $f : U \rightarrow V$ is called (b, ϵ) -holomorphic if there exists a constant $C_f < \infty$ such that the following holds: For every (b, ϵ) -admissible ρ , there exists an extension $\tilde{f} : B_\rho \rightarrow V_{\mathbb{C}}$ of f , i.e. we have $\tilde{f}(y) = f(y)$ for all $y \in U \subset B_\rho$, \tilde{f} is holomorphic in each component and such that $\sup_{z \in B_\rho} \|\tilde{f}(z)\|_{V_{\mathbb{C}}} \leq C_f$. Here, $B_\rho \subset \mathbb{C}^{\mathbb{N}}$ denotes the ball of polyradius ρ :

$$B_\rho = \{z \in \mathbb{C}^{\mathbb{N}} \mid |z_j| < \rho_j, \forall j \in \mathbb{N}\},$$

and $V_{\mathbb{C}} \simeq V + iV$ is the complexification of V .

In [60], it is shown that for a (b, ϵ) -holomorphic function f , an approximation rate of $O(\mathcal{N}^{1-1/p})$ can be obtained using a ReLU neural network of depth $O(\log(\mathcal{N}) \log \log(\mathcal{N}))$ or using a neural network with a smoother activation function of depth $O(\log \log(\mathcal{N}))$. We show that the same approximation rate can be obtained using a shallow tanh neural network.

Corollary 5.10. *Let $f : U = [-1, 1]^{\mathbb{N}} \rightarrow \mathbb{R}$ be (b, ϵ) -holomorphic for $b \in \ell^p(\mathbb{N})$, $p \in (0, 1)$ and $\epsilon > 0$. Then there exists a constant $C > 0$ such that for every $\mathcal{N} \in \mathbb{N}$ there exists a shallow tanh neural network $\hat{f}^{\mathcal{N}}$ of width at most $O(\mathcal{N}(C \log(\mathcal{N}))^{C \log(\mathcal{N})})$ such that*

$$(125) \quad \left\| f - \hat{f}^{\mathcal{N}} \right\|_{L^\infty(U)} = O(\mathcal{N}^{1-1/p}) \quad \text{for } \mathcal{N} \rightarrow \infty.$$

Proof. There exists a constant $C > 0$ and a sequence of index sets $(\Lambda_{\mathcal{N}})_{\mathcal{N} \in \mathbb{N}} \subset \mathcal{F}$ for which it holds that (cf. [60, proof of Thm. 3.9])

$$(126) \quad \sup_{y \in U} \left| f(y) - \sum_{\nu \in \Lambda_{\mathcal{N}}} \frac{D^\nu f(0)}{\nu!} y^\nu \right| = O(\mathcal{N}^{1-1/p}).$$

and such that $|\Lambda_{\mathcal{N}}| = \mathcal{N}$, $\text{supp}(\nu) \subseteq \{1, \dots, \mathcal{N}\}$ for all $\nu \in \Lambda_{\mathcal{N}}$ and for all \mathcal{N} [60, proof of Thm. 3.9], and $\sup_{\mathcal{N} \in \mathbb{N}} |\nu| \leq C(1 + \log(\mathcal{N}))$, where $|\nu| := \sum_{j \in \mathbb{N}} |\nu_j|$ [60, Thm. 2.7]. The latter implies in particular that $\sup_{\mathcal{N} \in \mathbb{N}} |\text{supp}(\nu)| \leq C(1 + \log(\mathcal{N}))$.

Based on these results from [60], it therefore suffices to show that we can accurately approximate all monomials $y \mapsto y^\nu$ for $\nu \in \Lambda_{\mathcal{N}}$ with shallow tanh neural

networks. For a fixed $\nu \in \Lambda_{\mathcal{N}}$, the monomial $y \mapsto y^\nu$ can be approximated (to arbitrary accuracy) using Corollary 3.6 with $d = n \leftarrow C(1 + \log(\mathcal{N}))$, resulting in a shallow tanh neural network of width $O\left((Ce(1 + \log(\mathcal{N})))^{C(1 + \log(\mathcal{N}) + 1)}\right)$. The network $\hat{f}^{\mathcal{N}}$ from the statement can then be constructed by parallelizing all the networks that approximate the individual monomials, yielding an approximation

$$\|f - \hat{f}^{\mathcal{N}}\|_{L^\infty(U)} = O(\mathcal{N}^{1-1/p}).$$

To be precise, we take the input of this network to be $(y_1, \dots, y_{\mathcal{N}})$ instead of y , which is possible since $\text{supp}(\nu) \subseteq \{1, \dots, \mathcal{N}\}$ for all $\nu \in \Lambda_{\mathcal{N}}$. As $|\Lambda_{\mathcal{N}}| = \mathcal{N}$, the resulting width of $\hat{f}^{\mathcal{N}}$ is $O\left(\mathcal{N}(Ce(1 + \log(\mathcal{N})))^{C(1 + \log(\mathcal{N}) + 1)}\right)$, which is asymptotically equivalent to the width from the statement for $\mathcal{N} \rightarrow \infty$. \square

This result implies in particular, that linear functionals of parametric solutions of PDEs can be approximated by shallow tanh neural networks [60]. Following [60, Section 4], the result can also be extended to directly approximate the parametric solution manifold, e.g. to approximate (b, ϵ) -holomorphic operators of the form $f : [-1, 1]^{\mathbb{N}} \rightarrow H_0^1([0, 1])$.

5.3. Examples. In this section, we illustrate the bounds, derived in Theorem 5.1 with several prototypical examples. In particular, we will investigate the width, weights and sparsity of the networks from the proof of Theorem 5.1.

First, we demonstrate how large the networks of Theorem 5.1 are for a simple function approximation example with $d = 1$ and $k = 0$. We consider the functions

$$(127) \quad f_a : [0, 1] \rightarrow [-1, 1] : x \mapsto \sin(ax), \quad a > 0.$$

For a given error tolerance $\varepsilon > 0$, we look for a three-layer tanh neural network $\hat{f}^{N,s}$, as given by Theorem 5.1, such that provably

$$(128) \quad \|f_a - \hat{f}^{N,s}\|_{\infty} \leq \varepsilon.$$

From all the networks that satisfy this condition, we take the one with the minimal width. More rigorously, we select

$$(129) \quad \operatorname{argmin}_{s, N \in \mathbb{N} : (3a/2N)^s / s! < \varepsilon} \max \left\{ 3 \left\lceil \frac{s}{2} \right\rceil + N - 1, 6N \right\},$$

where we used that $|f_a|_{W^{s,\infty}} \leq a^s$ for $s \geq 1$. Alternatively, one can also set $N = 1$ in Theorem 5.1, which makes the bound more efficient as no more partition of unity is needed, thereby reducing the need for a second hidden layer. This is similar to the proof of Corollary 5.5. In this case, we select

$$(130) \quad \operatorname{argmin}_{s \in \mathbb{N} : (a/2)^s / s! < \varepsilon} \left\{ 3 \left\lceil \frac{s}{2} \right\rceil \right\}.$$

We present the result in Figure 2. For the chosen examples, a shallow (i.e. two-layer) tanh neural network achieves a similar level of error as a three-layer network of the same width. This can be explained by the fact that $|f_a|_{W^{s,\infty}}$ grows as $O(a^s)$ and not as $O(a^s s!)$, such that setting $N > 1$ is not required for the bound of Theorem 5.1 to be non-vacuous. Moreover, the networks suggested by Theorem 5.1 are not unreasonably large for this simple example. Yet, they still remain overestimates: we found that e.g. $f_{2\pi}$ can already be approximated to an error of 1% by a shallow tanh network of width four. Finally, the exponential convergence is evident as a small increase in the network width already leads to a very large improvement in the accuracy.

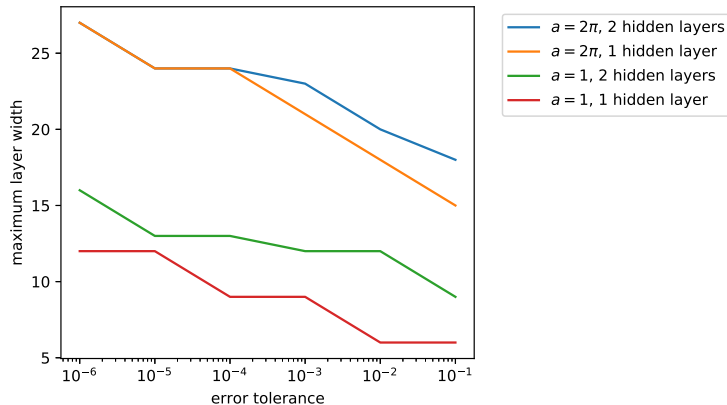


FIGURE 2. Needed layer width according to Theorem 5.1 to approximate the function f_a to a given error tolerance.

power	1	3	5	7	9
MSE	$4.91 \cdot 10^{-7}$	$1.54 \cdot 10^{-6}$	$2.87 \cdot 10^{-5}$	$1.13 \cdot 10^{-4}$	$6.13 \cdot 10^{-5}$
largest weight	3.13	2.14	2.27	4.55	4.41

TABLE 1. MSE and largest weight (in absolute value) of shallow tanh neural networks that approximate univariate monomials with odd powers on $[0, 1]$.

Next, we investigate whether the blow-up of the network weights from the theoretical results is observed in practice. We approximate univariate monomials of odd power in supremum norm on the interval $[0, 1]$ using shallow neural networks whose sizes are determined by Lemma 3.1. We generate a training set using 2000 randomly generated points, based on the uniform distribution on $[0, 1]$ and minimize the training loss for 2000 epochs using the Adam optimizer [28]. The results can be found in Table 1 and show that the weights do not blow up in practice. Rather, the weights remain small for this example. This is possibly a consequence of the phenomenon of implicit regularization in deep learning, e.g. [51].

Finally, we show that the neural networks constructed in the proof of Theorem 5.1 are not very sparse i.e., the fraction of non-zero weights of the network, compared to the total number of weights, is not small. Figure 3 shows that the fraction of non-zero weights of the network increases with increasing s and decreasing N (for $d = 1$). For analytic functions, it is (asymptotically) more efficient to increase s than N , as the convergence rate is $O(N^{-s})$. This lets us conclude that the constructed networks corresponding to sensible choices of s and N are, in general, quite dense. This is in agreement with what one observes in practice. This is in contrast to the theoretical results for deep ReLU (and other) neural networks, where the sparsity of the constructed networks generally increases with increasing accuracy [65, 54, 20].

6. DISCUSSION

The main aim of this paper was to provide explicit bounds on the error (in high-order Sobolev) norms with which a neural network with a tanh activation function approximates Sobolev regular functions, C^k functions and analytic functions. To this end, we prove such explicit bounds on the approximation error for Sobolev functions in Theorem 5.1 and for analytic functions in Corollary 5.5. In both

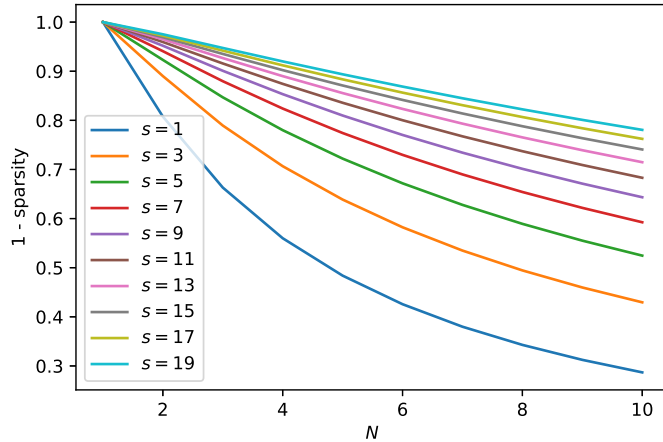


FIGURE 3. Fraction of non-zero weights (i.e. $1 - \text{sparsity}$) of the networks of Theorem 5.1 for different values of s and N in the case where $d = 1$.

cases, we prove these bounds for a tanh neural network with just 2 hidden layers. Our proofs are constructive and the construction relies on three key ideas: (1) the approximation of monomials by finite differences of a smooth activation function (Lemma 3.5), (2) the approximation of the multiplication operator (Lemma 3.7) and (3) the approximation of a partition of unity (Section 4).

At this point, it is instructive to compare and contrast our approach and results with the large body of literature on approximation of functions with artificial neural networks.

First, we compare our approach, as stated above, with other related works. The simple, yet very effective trick of approximating monomials by finite differences of smooth activation function has been around for decades [55], but is still a building block in the constructive proofs of many recent papers on neural network expressivity, e.g. [59, 52, 62, 20]. To the best of the authors' knowledge, all available results build upon the observation that there is a $x \in \mathbb{R}$ such that $\sigma^{(n)}(x) \neq 0$ for all $n \in \mathbb{N}$. As such, this construction does not allow for explicit estimates on the approximation error and the network weights (see Remark 3.4). *Our key novelty in this paper is to circumvent this issue by first approximating univariate monomials of odd powers and then expanding to even powers and multivariate monomials.* This allows us to obtain *uniform* explicit bounds for the error in approximating multivariate monomials of a varying degree and paves the way for explicit bounds on the approximation error.

The approximation of the multiplication operator in d dimensions by a shallow neural network (Corollary 3.7) was discussed in [37, Appendix A] for activation functions σ that satisfy that $\sigma(x) = \sum_k \sigma_k x^k$ where $\sigma_k \neq 0$ for $0 \leq k \leq d$. In particular, they prove that 2^d neurons are both sufficient and necessary. However, the hyperbolic tangent (and the sigmoid) activation function does not satisfy this assumption. Here, we propose a *novel construction* of the multiplication operator with a shallow tanh neural network.

As for the partition of unity, which serves as an essential ingredient in our proofs, an exact partition of unity for ReLU neural networks can be readily constructed [65]. Approximations of partitions of unity with neural networks with sigmoidal

activation function can be found in [10, 11, 52] and a general framework for approximations of partitions of unity was proposed in [20]. In this paper, we have constructed novel approximations of partitions of unity with shallow tanh neural networks, that were motivated by localized polynomials which arise in the Bramble-Hilbert Lemma and the Taylor's theorem. Compared to the other works mentioned above, our results on partitions of unity stand out for the explicit bounds on the approximation error and the weights.

Given the afore-mentioned novel ideas, we were able to obtain explicit bounds on the approximation error. A suitable avenue to compare our results with results obtained in related works lies in the approximation error bounds for analytic functions. We recall that in Corollary 5.5, we prove approximation rates for analytic functions in the $W^{k,\infty}$ -norm. Although approximation rates in this norm were proved for Sobolev functions in the very recent paper [20], it is unclear if their results can be extended for analytic functions. A key reason for this lies in the fact that the widths of their constructed networks are not explicitly stated and the depth increases with maximal degree of monomials, inhibiting uniform control that is necessary for approximating analytic functions.

Exponential convergence (in terms of network size) of neural networks to analytic functions in the L^∞ -norm was first proven in [43] for neural networks with smooth activation functions and in [64] for ReLU neural networks. In [54, 22], the authors prove exponential convergence in $W^{1,\infty}$ -norm for ReLU neural networks. We compare our results for approximation of analytic functions with these papers in Table 2. For [43], the parameter ρ is related to the polyradius of the ellipse to which the function needs to be holomorphically extendable. In [64], the additional assumption is made that the analytic function admits a Taylor expansion on $[-1, 1]^d$ that converges absolutely and uniformly, which does not hold for general analytic functions. For [54], the parameter β is at least inversely proportional to the dimension d and also depends on the radius of the Bernstein ellipse to which the analytic function can be holomorphically extended. From Table 2, one can clearly observe that our results yield an asymptotically faster convergence in terms of network width than the other related works. In addition, our results hold in stronger norms and we provide explicit bounds on the approximation error and weights, in contrast to other papers. For instance in [43, 54], the convergence rate even depends on the (unknown) polyradius of the ellipse to which the function can be holomorphically extended.

All of our approximation results, including the approximation bounds on analytic functions hold for a tanh neural network with only two hidden layers. This result, see also [43], runs contrary to the prevailing view that depth of neural networks is essential for function approximation and establishes that shallow but wide neural networks can be very expressive when it comes to function approximation and might provide some justification for the use of very shallow and wide neural networks in scientific computing [40, 39, 46].

Finally, it is essential to mention that although we highlight our contribution in terms of the tanh activation function as it is the most commonly used of the smooth activation functions. Our results apply verbatim to the logistic or sigmoid activation function as it is a shifted and scaled tanh. However, our constructions also apply to a much larger class of smooth activation functions as elaborated in Section 3.

We conclude by pointing out some limitations of the presented results. The most important limitation is the fact that the amplitude of the weights in our constructive network can grow very fast (Theorem 5.1). In practice, implicit and explicit regularization mechanisms during training will ensure that such growth of

source	norm	activation	depth	width	error bound
[43, Thm 2.3]	$L^\infty([-1, 1]^d)$	C^∞	2	\mathcal{N}	$O(\rho^{-\mathcal{N}/d})$
[64, Thm. 6]	$L^\infty([-1 + \delta, 1 - \delta]^d)$	ReLU	$O(\mathcal{N})$	$d + 4$	$O\left(\exp\left(-d\delta\mathcal{N}^{1/2d}\right)\right)$
[54, Thm. 3.6]	$W^{1,\infty}([-1, 1]^d)$	ReLU	$O(\mathcal{N}^{\frac{1}{\alpha+1}} \log(\mathcal{N}))$	$O(\mathcal{N})$	$O\left(\exp\left(-\beta\mathcal{N}^{\frac{1}{\alpha+1}}\right)\right)$
this work	$W^{k,\infty}([0, 1]^d)$	tanh	3	$O(\mathcal{N})$	$O\left(\mathcal{N}^{\frac{k}{\alpha+1}} \exp\left(-\mathcal{N}^{\frac{1}{\alpha+1}} \ln(\mathcal{N})\right)\right)$

TABLE 2. Comparison of upper bounds on the approximation error for analytic functions by neural networks.

weights will not happen. In fact, we present examples to empirically show that the gradient-descent based training procedure manages to find rather small weights and biases that still provide a very high accuracy. We therefore believe that our bounds are useful in practice, more as upper bounds for setting the network size.

Another limitation, which we share with other published results on approximation with neural networks, is that our results suffer from the curse of dimensionality. We could however prove that it is possible to obtain a dimension-independent convergence rate to analytic functions in corollaries 5.6 and 5.8. Another possible mitigation of the curse of dimensionality for the approximation rate is when the underlying map is (b, ϵ) -holomorphic, see Corollary 5.10 for the precise result. However in these cases, the constants (and hence the network size) can still grow exponentially in the input dimension. Fortunately, one can argue that a large number of high-dimensional functions are in fact compositions of low-dimensional functions, which might explain the success of deep learning in high dimensions [56]. For instance, the Kolmogorov-Arnold superposition theorem [29] even states that all d -variate functions are in fact compositions of univariate functions and the sum of d numbers, which also can be used to lessen the curse of dimensionality [50].

Finally, the weights in our constructed networks are continuous with respect to the function of interest f . It has been proven that the best neural approximation cannot be achieved using continuous weight selection [26]. An example of how discontinuous weight selection can improve the approximation rate can be found in [66].

REFERENCES

- [1] A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information theory*, 39(3):930–945, 1993.
- [2] A. R. Barron. Approximation and estimation bounds for artificial neural networks. *Machine learning*, 14(1):115–133, 1994.
- [3] C. Beck, A. Jentzen, and B. Kuckuck. Full error analysis for the training of deep neural networks. *arXiv preprint arXiv:1910.00121*, 2019.
- [4] M. Blanchard and M. Bennouna. The representation power of neural networks: Breaking the curse of dimensionality. *arXiv preprint arXiv:2012.05451*, 2020.
- [5] K. N. Boyadzhiev. Derivative polynomials for tanh, tan, sech and sec in explicit form. *arXiv preprint arXiv:0903.0117*, 2009.
- [6] E. Candes, L. Demanet, and L. Ying. Fast computation of Fourier integral operators. *SIAM Journal on Scientific Computing*, 29(6):2464–2493, 2007.
- [7] E. Candes, L. Demanet, and L. Ying. A fast butterfly algorithm for the computation of Fourier integral operators. *Multiscale Modeling & Simulation*, 7(4):1727–1750, 2009.
- [8] K. Cho, B. van Merriënboer, C. Gulcehre, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP 2014)*, 2014.
- [9] G. Constantine and T. Savits. A multivariate Faà di Bruno formula with applications. *Transactions of the American Mathematical Society*, 348(2):503–520, 1996.
- [10] D. Costarelli and R. Spigler. Approximation results for neural network operators activated by sigmoidal functions. *Neural Networks*, 44:101–106, 2013.

- [11] D. Costarelli and R. Spigler. Multivariate neural network operators with sigmoidal activation functions. *Neural Networks*, 48:72–77, 2013.
- [12] F. Cucker and S. Smale. On the mathematical foundations of learning. *Bulletin of the American Mathematical Society*, 39(1):1–49, 2002.
- [13] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2(4):303–314, Dec 1989.
- [14] K. R. Davidson and A. P. Donsig. *Real analysis and applications: theory in practice*. Springer Science & Business Media, 2009.
- [15] L. Demanet and L. Ying. On Chebyshev interpolation of analytic functions. *preprint*, 2010.
- [16] R. G. Durán. On polynomial approximation in Sobolev spaces. *SIAM journal on numerical analysis*, 20(5):985–988, 1983.
- [17] W. E, J. Han, and A. Jentzen. Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations. *Communications in Mathematics and Statistics*, 5(4):349–380, 2017.
- [18] P. Grohs and F. Voigtlaender. Proof of the theory-to-practice gap in deep learning via sampling complexity bounds for neural network approximation spaces. *arXiv preprint arXiv:2104.02746*, 2021.
- [19] I. Gühring, G. Kutyniok, and P. Petersen. Error bounds for approximations with deep ReLU neural networks in $W^{s,p}$ norms. *Analysis and Applications*, 18(05):803–859, 2020.
- [20] I. Gühring and M. Raslan. Approximation rates for neural networks with encodable weights in smoothness spaces. *Neural Networks*, 134:107–130, 2021.
- [21] J. Han, A. Jentzen, and W. E. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [22] L. Herrmann, J. Opschoor, and C. Schwab. Constructive deep ReLU neural network approximation. *SAM research report 2021-04, ETH Zürich*, 2021.
- [23] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [24] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359 – 366, 1989.
- [25] A. Jentzen and T. Welti. Overall error analysis for the training of deep neural networks via stochastic gradient descent with random initialisation. *arXiv preprint arXiv:2003.01291*, 2020.
- [26] P. C. Kainen, V. Kůrková, and A. Vogt. Approximation by neural networks is not continuous. *Neurocomputing*, 29(1-3):47–56, 1999.
- [27] H. Katsuura. Summations involving binomial coefficients. *The College Mathematics Journal*, 40(4):275–278, 2009.
- [28] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [29] A. N. Kolmogorov. On the representation of continuous functions of many variables by superposition of continuous functions of one variable and addition. In *Doklady Akademii Nauk*, volume 114, pages 953–956. Russian Academy of Sciences, 1957.
- [30] G. Kutyniok, P. Petersen, M. Raslan, and R. Schneider. A theoretical analysis of deep neural networks and parametric PDEs. *arXiv preprint arXiv:1904.00377*, 2019.
- [31] F. Laakmann and P. Petersen. Efficient approximation of solutions of parametric linear transport equations by ReLU DNNs. *arXiv preprint arXiv:2001.11441*, 2020.
- [32] I. E. Lagaris, A. Likas, and D. I. Fotiadis. Artificial neural networks for solving ordinary and partial differential equations. *IEEE Transactions on Neural Networks*, 9(5):987–1000, 2000.
- [33] S. Lanthaler, S. Mishra, and G. E. Karniadakis. Error estimates for DeepOnets: A deep learning framework in infinite dimensions. *arXiv preprint arXiv:2102.09618*, 2021.
- [34] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- [35] B. Li, S. Tang, and H. Yu. Better approximations of high dimensional smooth functions by deep neural networks with rectified power units. *arXiv preprint arXiv:1903.05858*, 2019.
- [36] Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, and A. Anandkumar. Fourier neural operator for parametric partial differential equations. *arXiv preprint arXiv:2010.08895v1*, 2020.
- [37] H. W. Lin, M. Tegmark, and D. Rolnick. Why does deep and cheap learning work so well? *Journal of Statistical Physics*, 168(6):1223–1247, 2017.
- [38] M. Longo, S. Mishra, C. Schwab, and T. K. Rusch. Higher-order Quasi-Monte Carlo training of deep neural networks. *arXiv preprint arXiv:2009.02713*, 2021.
- [39] L. Lu, P. Jin, and G. E. Karniadakis. DeepOnet: Learning nonlinear operators for identifying differential equations based on the universal approximation theorem of operators. *arXiv preprint arXiv:1910.03193*, 2019.
- [40] L. Lu, Y. Su, and G. E. Karniadakis. Collapse of deep and narrow neural nets. *arXiv preprint arXiv:1808.04947*, 2018.

- [41] K. O. Lye, S. Mishra, and D. Ray. Deep learning observables in computational fluid dynamics. *Journal of Computational Physics*, 410:109339, 2020.
- [42] K. O. Lye, S. Mishra, D. Ray, and P. Chandrashekar. Iterative surrogate model optimization (ISMO): An active learning algorithm for PDE constrained optimization with deep neural networks. *Computer Methods in Applied Mechanics and Engineering*, 374:113575, 2021.
- [43] H. N. Mhaskar. Neural networks for optimal approximation of smooth and analytic functions. *Neural computation*, 8(1):164–177, 1996.
- [44] S. Mishra and R. Molinaro. Estimates on the generalization error of physics informed neural networks (PINNs) for approximating PDEs. *arXiv preprint <https://arxiv.org/pdf/2006.16144.pdf>*, 2020.
- [45] S. Mishra and R. Molinaro. Estimates on the generalization error of physics-informed neural networks (PINNs) for approximating PDEs II: A class of inverse problems. *arXiv preprint [arXiv:2007.01138](https://arxiv.org/pdf/2007.01138.pdf)*, 2020.
- [46] S. Mishra and R. Molinaro. Physics-informed neural networks for simulating radiative transfer. *arXiv preprint [arXiv:2009.13291](https://arxiv.org/pdf/2009.13291.pdf)*, 2020.
- [47] S. Mishra and T. K. Rusch. Enhancing accuracy of deep learning algorithms by training on low-discrepancy sequences. *arXiv preprint [arXiv:2005.12564](https://arxiv.org/pdf/2005.12564.pdf)*, 2021.
- [48] D. S. Moak. Combinatorial multinomial matrices and multinomial Stirling numbers. *Proceedings of the American Mathematical Society*, pages 1–8, 1990.
- [49] H. Montanelli and Q. Du. New error bounds for deep ReLU networks using sparse grids. *SIAM Journal on Mathematics of Data Science*, 1(1):78–92, 2019.
- [50] H. Montanelli and H. Yang. Error bounds for deep ReLU networks using the Kolmogorov-Arnold superposition theorem. *Neural Networks*, 129:1–6, 2020.
- [51] B. Neyshabur, R. Tomioka, and N. Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. In *ICLR (Workshop)*, 2015.
- [52] I. Ohn and Y. Kim. Smooth function approximation by deep neural networks with general activation functions. *Entropy*, 21(7):627, 2019.
- [53] J. A. Opschoor, P. C. Petersen, and C. Schwab. Deep ReLU networks and high-order finite element methods. *Analysis and Applications*, 18(05):715–770, 2020.
- [54] J. A. Opschoor, C. Schwab, and J. Zech. Exponential ReLU DNN expression of holomorphic maps in high dimension. *SAM Research Report*, 2019, 2019.
- [55] A. Pinkus. Approximation theory of the MLP model in neural networks. *Acta numerica*, 8(1):143–195, 1999.
- [56] T. Poggio, H. Mhaskar, L. Rosasco, B. Miranda, and Q. Liao. Why and when can deep-but not shallow-networks avoid the curse of dimensionality: a review. *International Journal of Automation and Computing*, 14(5):503–519, 2017.
- [57] M. Raissi and G. E. Karniadakis. Hidden physics models: Machine learning of nonlinear partial differential equations. *Journal of Computational Physics*, 357:125–141, 2018.
- [58] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [59] D. Rolnick and M. Tegmark. The power of deeper networks for expressing natural functions. In *International Conference on Learning Representations*, 2018.
- [60] C. Schwab and J. Zech. Deep learning in high dimension: Neural network expression rates for generalized polynomial chaos expansions in UQ. *Analysis and Applications*, 17(01):19–55, 2019.
- [61] S. Shalev-Shwartz and S. Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [62] J. W. Siegel and J. Xu. Approximation rates for neural networks with general activation functions. *Neural Networks*, 128:313–321, 2020.
- [63] K. Weierstrass. Über die analytische Darstellbarkeit sogenannter willkürlicher Functionen einer reellen veränderlichen. *Sitzungsberichte der Königlich Preussischen Akademie der Wissenschaften zu Berlin*, 2:633–639, 1885.
- [64] E. Weinan and Q. Wang. Exponential convergence of the deep neural network approximation for analytic functions. *Science China Mathematics*, 61(10):1733–1740, 2018.
- [65] D. Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Networks*, 94:103–114, 2017.
- [66] D. Yarotsky. Optimal approximation of continuous functions by very deep ReLU networks. In *Conference on Learning Theory*, pages 639–649. PMLR, 2018.

APPENDIX A. AUXILIARY RESULTS

Lemma A.1. *It holds for every $n \in \mathbb{N}$ that $\left| \tanh^{(2n-1)}(0) \right| \geq 1$.*

Proof. For $|x| < \frac{\pi}{2}$, the power series expansion of \tanh at x is given by

$$(131) \quad \tanh(x) = \sum_{n=1}^{\infty} \frac{2^{2n}(2^{2n}-1)B_{2n}}{(2n)!} x^{2n-1},$$

where B_n is the n -th Bernoulli number. One can then calculate that for every $n \in \mathbb{N}$,

$$(132) \quad \left| \tanh^{(2n-1)}(0) \right| = \left| \frac{2^{2n}(2^{2n}-1)B_{2n}}{2n} \right| \geq 1.$$

This concludes the proof of the statement. \square

Lemma A.2. *Let $s \in 2\mathbb{N} - 1$. It holds that*

$$(133) \quad \inf_{\alpha > 0} \frac{2^{s/2}(1+\alpha)^{(s^2+s)/2}}{\alpha^{s/2}} \leq \sqrt{e}(2es)^{s/2},$$

where the infimum is reached at $\alpha = 1/s$.

Proof. Some elementary calculations show that

$$(134) \quad \inf_{\alpha > 0} \frac{2^{s/2}(1+\alpha)^{(s^2+s)/2}}{\alpha^{s/2}} = 2^{s/2} \left(1 + \frac{1}{s}\right)^{(s^2+s)/2} s^{s/2}.$$

Moreover, it holds that

$$(135) \quad \left(1 + \frac{1}{s}\right)^{(s^2+s)/2} = \left(\left(1 + \frac{1}{s}\right)^s\right)^{(s+1)/2} \leq e^{(s+1)/2}.$$

The statement follows immediately from these inequalities. \square

Lemma A.3. *Let $n, q \in \mathbb{N}$ and $P_{n,q} = \{\alpha \in \mathbb{N}_0^q : |\alpha| = n\}$. If $D = (D_{\alpha,\beta})_{\alpha,\beta \in P_{n,q}}$ is defined as in (39), then D is invertible and*

$$(136) \quad \left\| D^{-1} \right\|_{\infty} \leq (n!)^3 |P_{n,q}|^2 2^n.$$

Proof. Following [48], let $P_{n,q} = \{\alpha \in \mathbb{N}_0^q : |\alpha| = n\}$ and $I_{n,q} = \{\alpha' \in \mathbb{N}_0^{q-1} : |\alpha'| \leq n\}$. Let $s(k, m)$ be Stirling numbers of the first kind, for $k \leq m$, defined by

$$(137) \quad x(x-1) \cdots (x-m+1) = \sum_{k=0}^m s(k, m) x^k.$$

For $\alpha', \beta' \in I_{n,q}$, define $S_{\alpha',\beta'} = \prod_{i=1}^{q-1} s(\alpha'_i, \beta'_i)$, where $S(\alpha', \beta') = 0$ unless $\alpha'_i \leq \beta'_i$ for all i . Denote by S the corresponding matrix, where the order of rows and columns reflects the lexicographic order on $I_{n,q}$. Next, define B by

$$(138) \quad B_{\alpha',\beta'} = \binom{\alpha'}{\beta'} (-1)^{|\beta'|} := \prod_{i=1}^{q-1} \binom{\alpha'_i}{\beta'_i} (-1)^{\beta'_i} \quad \text{for } \alpha', \beta' \in I_{n,q}.$$

Finally, let L and Λ be diagonal matrices defined by $L_{\alpha',\alpha'} = (-1)^{|\alpha'|}/\alpha'!$ and $\Lambda_{\alpha',\alpha'} = n(n-1) \cdots (n-|\alpha'|+1)/\alpha'!$ for $\alpha' \in I_{n,q}$. It then holds that [48, Corollary 2],

$$(139) \quad D^{-1} = BL^{-1}\Lambda^{-1}SLB.$$

To prove an upper bound on the supremum norm of D^{-1} , we first note that $|s(k, m)| \leq \sum_{k=0}^m |s(k, m)| = m!$ (which can be seen by setting $x = -1$ in (137)) and thus $S(\alpha', \beta') \leq \beta'!$. In addition, it holds for any $\alpha' \in I_{n,q}$ that

$$(140) \quad 1 \leq \frac{n!}{\alpha'!(n-|\alpha'|)!} \leq \frac{n!}{\alpha'!},$$

which gives us that $\max_{\alpha' \in I_{n,q}} (\alpha'!) \leq n!$. This gives us consequently

(141)

$$\left| L^{-1} \Lambda^{-1} S \right|_{\alpha', \beta'} \leq (\alpha'!)^2 \beta'! \leq (n!)^2 \beta'!,$$

(142)

$$\left| L^{-1} \Lambda^{-1} S L \right|_{\alpha', \beta'} \leq (n!)^2,$$

(143)

$$\left| L^{-1} \Lambda^{-1} S L B \right|_{\alpha', \beta'} \leq (n!)^2 \sum_{\gamma' \in I_{n,q}} \binom{\gamma'}{\beta'} \leq (n!)^3 |I_{n,q}|,$$

(144)

$$\left| B L^{-1} \Lambda^{-1} S L B \right|_{\alpha', \beta'} \leq (n!)^3 |I_{n,q}| \sum_{\gamma' \in I_{n,q}} \binom{\alpha'}{\gamma'} = (n!)^3 |I_{n,q}| 2^{|\gamma'|} \leq (n!)^3 |I_{n,q}| 2^n.$$

This and (139) let us conclude that $\|D^{-1}\|_{\infty} \leq (n!)^3 |I_{n,q}|^2 2^n$. The lemma then follows from the existence of a one-to-one correspondence of elements in $I_{n,q}$ and $P_{n,q}$. \square

Lemma A.4. *Let $m \in \mathbb{N}$. Then it holds that*

$$(145) \quad \left| \sigma^{(m)}(x) \right| \leq (2m)^{m+1} \min\{\exp(-2x), \exp(2x)\} \quad \text{for all } x \in \mathbb{R}.$$

Proof. In [5], the following formula for the derivative of the hyperbolic tangent is proven,

$$(146) \quad \sigma^{(m)}(x) = (-2)^m (\sigma(x) + 1) \sum_{k=0}^m \frac{k!}{2^k} \left\{ \begin{matrix} m \\ k \end{matrix} \right\} (\sigma(x) - 1)^k,$$

where $\left\{ \begin{matrix} m \\ k \end{matrix} \right\}$ denote Stirling numbers of the second kind, for which it holds that $\left\{ \begin{matrix} m \\ k \end{matrix} \right\} \leq \frac{k^m}{k!}$. This then gives us

$$(147) \quad \left| \sigma^{(m)}(x) \right| \leq 2^m |1 + \sigma(x)| \sum_{k=0}^m k^m \leq 2^m m^{m+1} |1 + \sigma(x)| \leq (2m)^{m+1} \exp(2x),$$

as $\sum_{k=0}^m k^m \leq m \cdot m^m \leq m^{m+1}$. Furthermore one can note that $\sigma^{(m)}(-x) = -\sigma^{(m)}(x)$, which gives us

$$(148) \quad \left| \sigma^{(m)}(x) \right| \leq 2^m m^{m+1} |1 - \sigma(x)| \leq (2m)^{m+1} \exp(-2x).$$

The statement follows easily. \square

Lemma A.5. *The conditions stated in (53) for $k > 0$ are satisfied if*

$$(149) \quad \alpha = N \max \left\{ R, \ln \left(\frac{(2k)^{k+1} (Nk)^k}{e^k \epsilon} \right) \right\}.$$

Proof. The first condition of (53) is trivially satisfied when α is chosen as in the statement. From Lemma A.4, it follows that

$$(150) \quad \alpha^k (2k)^{k+1} \exp(-2\alpha/N) \leq \epsilon$$

is a sufficient condition that implies the other conditions of (53). Using $\max_{\alpha > 0} \alpha^k \exp(-\alpha/N) \leq (Nk)^k \exp(-k)$ we find that

$$(151) \quad \alpha^k \exp(-2\alpha/N) = \alpha^k \exp(-\alpha/N) \exp(-\alpha/N) \leq (Nk)^k \exp(-k) \exp(-\alpha/N).$$

The statement follows directly. \square

Lemma A.6. *Let $d \in \mathbb{N}$, $k \in \mathbb{N}_0$, $\Omega \subset \mathbb{R}^d$ and $f, g \in C^k(\Omega)$. Then it holds that*

$$(152) \quad \|fg\|_{W^{k,\infty}} \leq 2^k \|f\|_{W^{k,\infty}} \|g\|_{W^{k,\infty}}.$$

Proof. The statement follows directly from the general Leibniz rule. \square

Lemma A.7. *Let $d, m, n \in \mathbb{N}$, $\Omega_1 \subset \mathbb{R}^d$, $\Omega_2 \subset \mathbb{R}^m$, $f \in C^n(\Omega_1; \Omega_2)$ and $g \in C^n(\Omega_2; \mathbb{R})$. Then it holds that*

$$(153) \quad \|g \circ f\|_{W^{n,\infty}} \leq 16(e^2 n^4 m d^2)^n \|g\|_{W^{n,\infty}} \max_{1 \leq i \leq m} \|(f)_i\|_{W^{n,\infty}}^n.$$

Proof. Let $\nu \in \mathbb{N}^d$ with $|\nu| = n$. We use the multivariate Faà di Bruno formula [9],

$$(154) \quad D^\nu(g \circ f) = \sum_{1 \leq |\lambda| \leq n} D^\lambda g \sum_{p(\nu, \lambda)} (\nu!) \prod_{j=1}^n \frac{(f_{l_j})^{k_j}}{k_j! (l_j!)^{|k_j|}},$$

where $(f_\mu)_i = D^\mu(f)_i$ for $1 \leq i \leq m$ and the set $p(\nu, \lambda)$ is defined as

$$(155) \quad \begin{aligned} p(\nu, \lambda) = \{(\kappa, \ell) := (k_1, \dots, k_n; l_1, \dots, l_n) : & \text{for some } 1 \leq s \leq n, \\ & k_i = 0 \text{ and } l_i = 0 \text{ for } 1 \leq i \leq n-s; |k_i| > 0 \text{ for } n-s+1 \leq i \leq n; \\ & \text{and } 0 \prec l_{n-s+1} \prec \dots \prec l_n \text{ are such that} \\ & \sum_{i=1}^n k_i = \lambda, \sum_{i=1}^n |k_i| l_i = \nu\}, \end{aligned}$$

where $a \prec b$ either means that $|a| < |b|$ or $a < b$ according to lexicographic ordering; furthermore the vectors k_i are m -dimensional and the l_i are d -dimensional. From the stated conditions, it follows directly that $\sum_{i=1}^n |k_i| \leq n$ and $\sum_{i=1}^n |l_i| \leq n$. Next, we bound the complexity of $p(\nu, \lambda)$. From $\sum_{i=1}^n |k_i| \leq n$, it follows that the number of κ is bounded above by $|P_{n, (m+1)n}|$, which can in turn be bounded by $\sqrt{\pi} e^n (mn)^n$ by Lemma 2.1. Similarly, it follows that the number of ℓ is bounded above by $|P_{n, (d+1)n}|$, which can in turn be bounded by $\sqrt{\pi} e^n (dn)^n$ by Lemma 2.1. Therefore, $|p(\nu, \lambda)| \leq \pi (e^2 n^2 m d)^n$. Finally, we can make the estimates that $|\{\lambda : 1 \leq |\lambda| \leq n\}| \leq |P_{n, d+1}| \leq \sqrt{\pi} e^n d^n$, $D^\lambda g \leq \|g\|_{W^{n,\infty}}$, $\nu! \leq n!$ and $\prod_{j=1}^n (f_{l_j})^{k_j} \leq \max_{1 \leq i \leq m} \|(f)_i\|_{W^{n,\infty}}^n$. Together with Stirling's approximation, this yields

$$(156) \quad \begin{aligned} \|D^\nu(g \circ f)\|_\infty &\leq \sqrt{\pi} e^n d^n \|g\|_{W^{n,\infty}} \cdot \pi (e^2 n^2 m d)^n \cdot n! \cdot \max_{1 \leq i \leq m} \|(f)_i\|_{W^{n,\infty}}^n \\ &\leq 16(e^2 n^4 m d^2)^n \|g\|_{W^{n,\infty}} \max_{1 \leq i \leq m} \|(f)_i\|_{W^{n,\infty}}^n. \end{aligned}$$

\square

Lemma A.8 (Bramble-Hilbert). *Let $\Omega \subset \mathbb{R}^d$ be an open and bounded set of diameter $0 < h < e^{-1/2} d^{-3/2}$ which is star-shaped with respect to every point in an open ball $B \subset \Omega$ with diameter ph . Then for every $f \in W^{s,\infty}(\Omega)$ there exists a polynomial \hat{f} of degree at most $s-1$ such that for any $k \in \mathbb{N}_0$ with $k < s$ it holds that,*

$$(157) \quad \left\| f - \hat{f} \right\|_{W^{k,\infty}(\Omega)} \leq \frac{\sqrt{s} \pi^{1/4} (d\sqrt{deh})^{s-k}}{(s-k-1)!} |f|_{W^{s,\infty}(\Omega)}.$$

Proof. By setting $p = q = \infty$ in the penultimate equation in the proof of the main theorem in [16] (note that this reference uses a different definition of Sobolev norm),

it follows that there exists a polynomial \hat{f} of degree at most $s - 1$ such that for $0 \leq m < s$ it holds that,

$$(158) \quad \left| f - \hat{f} \right|_{W^{m, \infty}(\Omega)} \leq (s-m) \left(\sum_{\beta \in P_{s-m, d}} (\beta!)^{-2} \right)^{1/2} h^{s-m} \sqrt{|P_{s-m, d}|} \|f\|_{W^{s, \infty}(\Omega)}.$$

Using Lemma 2.1, we find that $\sqrt{|P_{s-m, d}|} \leq \pi^{1/4} (ed)^{(s-m)/2}$ and from the multinomial theorem it follows that

$$(159) \quad \sum_{\beta \in P_{s-m, d}} (\beta!)^{-2} \leq \sum_{\beta' \in P_{2(s-m), 2d}} (\beta'!)^{-1} = \frac{(2d)^{2(s-m)}}{(2(s-m))!}.$$

One can also calculate that $(2(s-m))! \geq 4^{s-m} (s-m)((s-m-1)!)^2$. Combining the previous observations, we find

$$(160) \quad \left| f - \hat{f} \right|_{W^{m, \infty}(\Omega)} \leq \frac{\sqrt{s-m} (d\sqrt{deh})^{s-m}}{(s-m-1)!} \|f\|_{W^{s, \infty}(\Omega)}.$$

Majorizing over $0 \leq m \leq k$ then gives the upper bound from the statement. \square

Lemma A.9 (Taylor's theorem). *Let $d, s \in \mathbb{N}$, $0 < \delta < 1/d$. Then for every $f \in C^s([- \delta, \delta]^d)$ there exists a polynomial \hat{f} of degree at most $s - 1$ such that for any $k \in \mathbb{N}_0$ with $k < s$ it holds that,*

$$(161) \quad \left\| f - \hat{f} \right\|_{W^{k, \infty}([- \delta, \delta]^d)} \leq \frac{(d\delta)^{s-k}}{(s-k)!} \|f\|_{W^{s, \infty}([- \delta, \delta]^d)}$$

Proof. We give a constructive proof. For $f \in C^s([- \delta, \delta]^d)$, we define the polynomial \hat{f} as

$$(162) \quad \hat{f}(x) = \sum_{|\alpha| \leq s-1} \frac{D^\alpha f(0)}{\alpha!} x^\alpha.$$

Then take $\beta \in \mathbb{N}_0^d$ with $|\beta| \leq k$. It then holds that

$$(163) \quad D^\beta \hat{f}(x) = \sum_{\substack{|\alpha| \leq s-1 \\ \alpha \geq \beta}} \frac{D^\alpha f(0)}{\alpha!} \frac{\alpha!}{(\alpha - \beta)!} x^{\alpha - \beta} = \sum_{|\gamma| \leq s-1 - |\beta|} \frac{D^\gamma D^\beta f(0)}{\gamma!} x^\gamma.$$

For $x \in [- \delta, \delta]^d$, Taylor's theorem guarantees the existence of a constant $c \in (0, 1)$ such that

$$(164) \quad D^\beta f(x) = \sum_{|\gamma| \leq s-1 - |\beta|} \frac{D^\gamma D^\beta f(0)}{\gamma!} x^\gamma + \sum_{|\gamma| = s - |\beta|} \frac{D^\gamma D^\beta f(cx)}{\gamma!} x^\gamma.$$

The previous equalities, together with the multinomial theorem, then prove that

$$(165) \quad \left\| D^\beta f - D^\beta \hat{f} \right\|_{L^\infty([- \delta, \delta]^d)} \leq C_s \sum_{|\gamma| = s - |\beta|} \frac{\delta^{|\gamma|}}{\gamma!} = \frac{C_s (d\delta)^{s - |\beta|}}{(s - |\beta|)!},$$

where $C_s := \|f\|_{W^{s, \infty}([- \delta, \delta]^d)}$. Under the assumption that $\delta < 1/d$ we then can conclude that

$$(166) \quad \left\| f - \hat{f} \right\|_{W^{k, \infty}([- \delta, \delta]^d)} \leq \frac{C_s (d\delta)^{s-k}}{(s-k)!}.$$

\square