# Approximation by Combinations of ReLU and Squared ReLU Ridge Functions With $\ell^1$ and $\ell^0$ Controls

Jason M. Klusowski<sup>(D)</sup>, Student Member, IEEE, and Andrew R. Barron, Fellow, IEEE

Abstract—We establish  $L^{\infty}$  and  $L^2$  error bounds for functions of many variables that are approximated by linear combinations of rectified linear unit (ReLU) and squared ReLU ridge functions with  $\ell^1$  and  $\ell^0$  controls on their inner and outer parameters. With the squared ReLU ridge function, we show that the  $L^2$  approximation error is inversely proportional to the inner layer  $\ell^0$  sparsity and it need only be sublinear in the outer layer  $\ell^0$  sparsity. Our constructions are obtained using a variant of the Maurey-Jones-Barron probabilistic method, which can be interpreted as either stratified sampling with proportionate allocation or two-stage cluster sampling. We also provide companion error lower bounds that reveal near optimality of our constructions. Despite the sparsity assumptions, we showcase the richness and flexibility of these ridge combinations by defining a large family of functions, in terms of certain spectral conditions, that are particularly well approximated by them.

*Index Terms*—Ridge combinations, rectified linear unit, approximation error, spline, stratified sampling, sparse models.

#### I. INTRODUCTION

**F**UNCTIONS of many variables are approximated using linear combinations of ridge functions with one layer of nonlinearities, viz.,

$$f_m(x) = \sum_{k=1}^{m} b_k \phi(a_k \cdot x - t_k),$$
 (1)

where  $b_k \in \mathbb{R}$  are the outer layer parameters and  $a_k \in \mathbb{R}^d$ are the vectors of inner parameters for the single-hidden layer of functions  $\phi(a_k \cdot x - t_k)$ . The activation function  $\phi$  is allowed to be quite general. For example, it can be bounded and Lipschitz, polynomials with certain controls on their degrees, or bounded with jump discontinuities. When the ridge activation function is a sigmoid, (1) is single-hidden layer artificial neural network.

One goal in a statistical setting is to estimate a regression function, i.e., conditional mean response,  $f(x) = \mathbb{E}[Y | X = x]$  with domain  $D \triangleq [-1, 1]^d$  from noisy observations  $\{(X_i, Y_i)\}_{i=1}^n$ , where  $Y = f(X) + \varepsilon$ . In classical literature [1],

Manuscript received July 28, 2016; revised February 24, 2018; accepted September 13, 2018. Date of publication October 8, 2018; date of current version November 20, 2018.

J. M. Klusowski is with the Department of Statistics and Biostatistics, Rutgers University-New Brunswick, Piscataway, NJ 08854-8019 USA (e-mail: jason.klusowski@rutgers.edu).

A. R. Barron is with the Department of Statistics and Data Science, Yale University, New Haven, CT 06511 USA (e-mail: andrew.barron@yale.edu).

Communicated by C. Caramanis, Associate Editor for Machine Learning. Digital Object Identifier 10.1109/TIT.2018.2874447  $L^2(P)$  mean squared prediction error of order  $(d/n)^{1/2}$ , achieved by  $\ell^1$  penalized least squares estimators<sup>1</sup> over the class of models (1), are obtained by optimizing the tradeoff between *approximation error* and *descriptive complexity relative to sample size*. Bounds on the approximation error are obtained by first showing how models of the form (1) with  $\phi(z) = \mathbf{1}\{z > 0\}$  can be used to approximate f satisfying  $\int_{\mathbb{R}^d} \|\omega\|_1 |\mathcal{F}(f)(\omega)| d\omega < +\infty$ , provided f admits a Fourier representation  $f(x) = \int_{\mathbb{R}^d} e^{ix \cdot \omega} \mathcal{F}(f)(\omega) d\omega$  on  $[-1, 1]^d$ . Because it is often difficult to work with discontinuous  $\phi$ (i.e., vanishing or exploding gradient issues), these step functions are replaced with smooth  $\phi$  such that  $\phi(\tau z) \wedge 1 \rightarrow$  $\mathbf{1}\{z > 0\}$  as  $\tau \to +\infty$ . Thus, this setup allows one to work with approximants of the form (1) with smooth  $\phi$ , but at the expense of *unbounded*  $\ell^1$  norm  $||a_k||_1$ .

Like high-dimensional linear regression [2], many applications of statistical inference and estimation require a setting where  $d \gg n$ . In contrast to the aforementioned mean square prediction error of  $(d/n)^{1/2}$ , it has been shown [3] how models of the form (1) with Lipschitz<sup>2</sup>  $\phi$  (reps. Lipschitz derivative  $\phi'$ ) and *bounded* inner parameters  $||a_k||_0$  and  $||a_k||_1$  can be used to give desirable  $L^2(D)$  mean squared prediction error of order  $((\log d)/n)^{1/3}$  (resp.  $((\log d)/n)^{2/5}$ ), also achieved by penalized estimators.<sup>3</sup> In fact, [4] shows that these rates are nearly optimal. A few natural questions arise from restricting the  $\ell^0$  and  $\ell^1$  norms of the inner parameters in the model:

- To what degree do the sparsity assumptions limit the flexibility of the model (1)?
- What condition can be imposed on f so that it can be approximated by  $f_m$  with Lipschitz  $\phi$  (or Lipschitz derivative  $\phi'$ ) and bounded  $||a_k||_0$  and / or  $||a_k||_1$ ?
- How well can f be approximated by  $f_m$ , given these sparsity constraints?

According to classic approximation results [5], [6], if the domain of f is contained in  $[-1, 1]^d$  and f admits a Fourier representation  $f(x) = \int_{\mathbb{R}^d} e^{ix \cdot \omega} \mathcal{F}(f)(\omega) d\omega$ , then the spectral condition  $v_{f,1} < +\infty$ , where  $v_{f,s} \triangleq \int_{\mathbb{R}^d} ||\omega||_1^s |\mathcal{F}(f)(\omega)| d\omega$ , is enough to ensure that f - f(0) can be approximated

<sup>1</sup>That is, the fit minimizes  $(1/n) \sum_{i=1}^{n} (f_m(X_i) - Y_i)^2 + \lambda \sum_{k=1}^{m} |b_k|$  for some appropriately chosen  $\lambda > 0$ .

 $^2 {\rm Henceforth},$  when we say a function is Lipschitz, we assume it has bounded Lipschitz parameter.

<sup>3</sup>With additional  $\ell^0$  inner sparsity, we might also consider an estimator that minimizes  $(1/n) \sum_{i=1}^n (f_m(X_i) - Y_i)^2 + \lambda_0 \psi \left( \sum_{k=1}^m |b_k| ||a_k||_0 \right)$  for some convex function  $\psi$  and appropriately chosen  $\lambda_0 > 0$ .

0018-9448 © 2018 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications\_standards/publications/rights/index.html for more information.

in  $L^{\infty}(D)$  by equally weighted, i.e,  $|b_1| = \cdots = |b_m|$ , linear combinations of functions of the form (1) with  $\phi(z) = \mathbf{1}\{z > 0\}$ . Typical  $L^{\infty}$  error rates  $||f - f_m||_{\infty}$  of an *m*-term approximation (1) are at most  $cv_{f,1}\sqrt{d} \ m^{-1/2}$ , where *c* is a universal constant [5], [7], [8]. A rate of  $c(p)v_{f,1}m^{-1/2-1/(pd)}$  was given in [9, Th. 3] for  $L^p(D)$ for nonnegative even integer *p*. Again, all these bounds are valid when the step activation function is replaced by a smooth approximant  $\phi$  (in particular, *any* sigmoid satisfying  $\lim_{z\to\pm\infty} \phi(z) = \pm 1$ ), but at the expense of unbounded  $||a_k||_1$ .

Towards giving partial answers to the questions we posed, in Section II, we show how functions of the form (1) with ReLU (also known as a ramp or first order spline)  $\phi(z) = (z)_+^2 = 0 \lor z$  (which is Lipchitz)<sup>4</sup> or squared ReLU  $\phi(z) = (z)_+^2$  (which has Lipschitz derivative) activation function can be used to give desirable  $L^{\infty}(D)$  approximation error bounds, even when  $||a_k||_1 = 1$ ,  $0 \le t_k \le 1$ , and  $|b_1| = \cdots = |b_m|$ . Because of the widespread popularity of the ReLU activation function and its variants, these simpler forms may also be of independent interest for computational and algorithmic reasons as in [10]–[14], to name a few.

Unlike the case with step activation functions, our analysis makes no use of the combinatorial properties of half-spaces as in Vapnik-Chervonenkis theory [15], [16]. The  $L^2(D)$  case for ReLU ridge functions (also known as hinging hyperplanes) with  $\ell^1$ -bounded inner parameters was considered in [17, Th. 3] and our  $L^{\infty}(D)$  bounds improve upon that line of work and, in addition, increase the exponent from 1/2 to 1/2 + O(1/d). Our proof techniques are substantively different than [17] and, importantly, are more amenable to empirical process theory, which is the key to showing our error bounds.

These tighter rates of approximation, with ReLU and squared ReLU activation functions, are possible under two different conditions – finite  $v_{f,2}$  or  $v_{f,3}$ , respectively. The main idea we use originates from [9] and [18] and can be seen as stratified sampling with proportionate allocation. This technique is widely applied in survey sampling as a means of variance reduction [19].

At the end of Section II, we will also discuss the degree to which these bounds can be improved by providing companion lower bounds on the minimax rates of approximation.

Section III will focus on how accurate estimation can be achieved even when  $||a_k||_0$  is also bounded. In particular, we show how an *m*-term linear combination (1) with  $||a_k||_0 \le \sqrt{m}$  and  $||a_k||_1 = 1$  can approximate *f* satisfying  $v_{f,3} < +\infty$ in  $L^2(D)$  with error at most  $\sqrt{2}v_{f,3}m^{-1/2}$ . In other words, the  $L^2(D)$  approximation error is inversely proportional to the inner layer sparsity and it need only be sublinear in the outer layer sparsity. The constructions that achieve these error bounds are obtained using a variant of the Maurey-Jones-Barron probabilistic method, which can be interpreted as two-stage cluster sampling.

Throughout this paper, we will state explicitly how our bounds depend on d so that the reader can fully appreciate the complexity of approximation. If a is a vector in

<sup>4</sup>It is perhaps more conventional to write  $(z)^+$  for  $0 \lor z$ , however, to avoid clutter in the exponent, we use the current notation.

Euclidean space, we use the notation a(k) to denote its k-th component.

## II. $L^{\infty}$ Approximation With Bounded $\ell^1$ Norm

### A. Positive Results

In this section, we provide the statements and proofs of the existence results for  $f_m$  with bounded  $\ell^1$  norm of inner parameters. We would like to point out that the results of Theorem 1 hold when all occurrences of the ReLU or squared ReLU activation functions are replaced by general  $\phi$  which is Lipschitz or has Lipschitz derivative  $\phi'$ , respectively.

Theorem 1: Suppose f admits an integral representation

$$f(x) = v \int_{[0,1] \times \{a: \|a\|_1 = 1\}} \eta(t,a) (a \cdot x - t)_+^{s-1} dP(t,a), \quad (2)$$

for x in  $D = [-1, 1]^d$  and  $s \in \{2, 3\}$ , where P is a probability measure on  $[0, 1] \times \{a \in \mathbb{R}^d : ||a||_1 = 1\}$  and  $\eta(t, a)$  is either -1 or +1. There exists a linear combination of ridge functions of the form

$$f_m(x) = \frac{v}{m} \sum_{k=1}^m b_k (a_k \cdot x - t_k)_+^{s-1},$$
 (3)

with  $b_k \in [-1, 1]$ ,  $||a_k||_1 = 1$ ,  $0 \le t_k \le 1$  such that

$$\sup_{x \in D} |f(x) - f_m(x)| \le c\sqrt{d + \log m} \ m^{-1/2 - 1/d}, \ s = 2$$

and

S

$$\sup_{e \in D} |f(x) - f_m(x)| \le c\sqrt{d} \ m^{-1/2 - 1/d}, \quad s = 3,$$

for some universal constant c > 0. Furthermore, if the  $b_k$  are restricted to  $\{-1, 1\}$ , the upper bound is of order

 $\sqrt{d + \log m} m^{-1/2 - 1/(d+2)}, \quad s = 2$ 

and

$$\sqrt{d} m^{-1/2-1/(d+2)}, s = 3.$$

Theorem 2: Let  $D = [-1, 1]^d$ . Suppose f admits a Fourier representation  $f(x) = \int_{\mathbb{R}^d} e^{ix \cdot \omega} \mathcal{F}(f)(\omega) d\omega$  and

$$v_{f,2} = \int_{\mathbb{R}^d} \|\omega\|_1^2 |\mathcal{F}(f)(\omega)| d\omega < +\infty.$$

There exists a linear combination of ReLU ridge functions of the form

$$f_m(x) = b_0 + a_0 \cdot x + \frac{v}{m} \sum_{k=1}^m b_k (a_k \cdot x - t_k)_+$$
(4)

with  $b_k \in [-1, 1]$ ,  $||a_k||_1 = 1$ ,  $0 \le t_k \le 1$ ,  $b_0 = f(0)$ ,  $a_0 = \nabla f(0)$ , and  $v \le 2v_{f,2}$  such that

$$\sup_{x \in D} |f(x) - f_m(x)| \le c v_{f,2} \sqrt{d + \log m} \ m^{-1/2 - 1/d},$$

for some universal constant c > 0. Furthermore, if the  $b_k$  are restricted to  $\{-1, 1\}$ , the upper bound is of order

$$v_{f,2}\sqrt{d+\log m} \ m^{-1/2-1/(d+2)}$$

Authorized licensed use limited to: New York University. Downloaded on April 03,2024 at 00:18:36 UTC from IEEE Xplore. Restrictions apply.

Theorem 3: Under the setup of Theorem 2, suppose

$$v_{f,3} = \int_{\mathbb{R}^d} \|\omega\|_1^3 |\mathcal{F}(f)(\omega)| d\omega < +\infty$$

There exists a linear combination of squared ReLU ridge functions of the form

$$f_m(x) = b_0 + a_0 \cdot x + x^T A_0 x + \frac{v}{2m} \sum_{k=1}^m b_k (a_k \cdot x - t_k)_+^2 \quad (5)$$

with  $b_k \in [-1, 1]$ ,  $||a_k||_1 = 1$ ,  $0 \le t_k \le 1$ ,  $b_0 = f(0)$ ,  $a_0 = \nabla f(0)$ ,  $A_0 = \nabla \nabla^T f(0)$ , and  $v \le 2v_{f,3}$  such that

$$\sup_{x \in D} |f(x) - f_m(x)| \le c v_{f,3} \sqrt{d} \ m^{-1/2 - 1/d}$$

for some universal constant c > 0. Furthermore, if the  $b_k$  are restricted to  $\{-1, 1\}$ , the upper bound is of order

$$v_{f,3}\sqrt{d} m^{-1/2-1/(d+2)}$$
.

The key observation for proving Theorem 2 and Theorem 3 is that f modulo linear or quadratic terms with finite  $v_{f,s}$ can be written in the integral form (2). Unlike in [17, Th. 3] where an interpolation argument is used, our technique of writing f as the mean of a random variable allows for more straightforward use of empirical process theory to bound the expected sup-error of the empirical average of *m* independent draws from its population mean. Our argument is also more flexible than [17] and can be readily adapted to the case of squared ReLU activation function. We should also point out that our  $L^{\infty}(D)$  error bounds immediately imply  $L^{p}(D)$  error bounds for all p > 1. In fact, using nearly exactly the same techniques, it can be shown that the results in Theorem 1, Theorem 2, and Theorem 3 hold verbatim in  $L^2(D)$ , sans the  $\sqrt{d + \log m}$  or  $\sqrt{d}$  factors, corresponding to the ReLu or squared ReLU cases, respectively.

Remark 1: In [18], it was shown that the standard order  $m^{-1/2} L^{\infty}(D)$  error bound alluded to earlier could be improved to be of order  $\sqrt{\log m} m^{-1/2-1/(2d)}$  under an alternate condition of finite  $v_{f,1}^* \triangleq \sup_{u \in \mathbb{S}^{d-1}} \int_0^\infty r^d |\mathcal{F}(f)(ru)| dr$ , but with the requirement that  $||a_k||_1$  be unbounded. In general, our assumptions are neither stronger nor weaker than this since the function f with Fourier transform  $\mathcal{F}(f)(\omega) = e^{-||\omega-\omega_0||}/||\omega-\omega_0||$  for  $\omega_0 \neq 0$  and  $d \geq 2$  has infinite  $v_{f,1}^*$  but finite  $v_{f,s}$  for  $s \geq 0$ , while the function f with Fourier transform  $\mathcal{F}(f)(\omega) = 1/(1+||\omega||)^{d+2}$  has finite  $v_{f,1}^*$  but infinite  $v_{f,s}$  for  $s \geq 2$ .

Proof of Theorem 1:

**Case I:** s = 2. Let  $\mathcal{B}_1, \ldots, \mathcal{B}_M$  be a partition of the space  $\Omega = \{(\eta, t, a)' : \eta \in \{-1, +1\}, 0 \le t \le 1, \|a\|_1 = 1\}$  such that

$$\inf_{(\tilde{\eta},\tilde{t},\tilde{a})'\in\mathcal{B}_{k},\ k=1,\dots,M} \sup_{(\eta,t,a)'\in\Omega} \|h(\tilde{\eta},\tilde{t},\tilde{a}) - h(\eta,t,a)\|_{\infty} < \epsilon,$$
(6)

where  $h(\eta, t, a)(x) = h(x) = \eta(a \cdot x - t)_{+}^{s-1}$ . It is not hard to show that  $M \simeq \epsilon^{-d}$ . For k = 1, ..., M define

$$dP_k(t,a) = dP(t,a)\mathbf{1}\{(\eta(t,a),t,a)' \in \mathcal{B}_k\}/L_k,$$

where  $L_k$  is chosen to make  $P_k$  a probability measure. A very important property we will use is that  $\operatorname{Var}_{P_k}[h] \leq \epsilon$ , which follows from (6). Let *m* be a positive integer and define a sequence of *M* independent random variables  $\{m_k\}_{1 \le k \le M}$  as follows: let  $m_k$  equal  $\lfloor mL_k \rfloor$  and  $\lceil mL_k \rceil$  with probabilities chosen to make its mean equal to  $mL_k$ . Given,  $\underline{m} = \{m_k\}_{1 \le k \le M}$ , take a random sample  $\underline{a} = \{(t_{j,k}, a_{j,k})'\}_{1 \le j \le n_k, 1 \le k \le M}$  of size  $n_k = m_k + 1\{m_k = 0\}$  from  $P_k$ . Thus, we split the population  $\Omega$  into *M* "strata"  $\mathcal{B}_1, \ldots, \mathcal{B}_M$  and allocate the number of within-stratum samples to be proportional to the "size" of the stratum  $m_1, \ldots, m_M$  (i.e., proportionate allocation). The within-stratum variability of *h* (i.e.,  $\operatorname{Var}_{P_k}[h]$ ) is now smaller than the population level variability (i.e.,  $\operatorname{Var}_P[h]$ ) by a factor of  $\epsilon$  as evidenced by (6). Note that the  $n_k$  sum to be at most m + M because

$$\sum_{k=1}^{M} n_k = \sum_{k=1}^{M} m_k \mathbf{1}\{m_k > 0\} + \sum_{k=1}^{M} \mathbf{1}\{m_k = 0\}$$
  

$$\leq \sum_{k=1}^{M} (mL_k + 1)\mathbf{1}\{m_k > 0\} + \sum_{k=1}^{M} \mathbf{1}\{m_k = 0\}$$
  

$$= m \sum_{k=1}^{M} L_k \mathbf{1}\{m_k > 0\} + M$$
  

$$\leq m + M, \qquad (7)$$

where the last inequality follows from  $\sum_{k=1}^{M} L_k \leq 1$ . For  $j = 1, \ldots, m_k$ , let  $h_{j,k} = h(\eta(t_{j,k}, a_{j,k}), t_{j,k}, a_{j,k})$  and  $f_k = \frac{vm_k}{mn_k} \sum_{j=1}^{n_k} h_{j,k}$ . Also, let  $\overline{f}_m = \sum_{k=1}^{M} f_k$ . A simple calculation shows that the mean of  $\overline{f}_m$  is f. Write  $\sum_{k=1}^{M} (f_k(x) - \mathbb{E}f_k(x)) = \frac{v}{m} \left( \sum_{k=1}^{M} (m_k - L_km)\mathbb{E}_{P_k}h(x) \right) + \frac{v}{m} \left( \sum_{k=1}^{M} \sum_{j=1}^{n_k} \frac{m_k}{n_k} (h_{j,k}(x) - \mathbb{E}_{P_k}h(x)) \right)$ . By the triangle inequality, we upper bound

$$\mathbb{E}\sup_{x\in D} |\overline{f}_m(x) - f(x)| = \mathbb{E}\sup_{x\in D} |\sum_{k=1}^M (f_k(x) - \mathbb{E}f_k(x))|$$

by

$$\frac{v}{m} \mathbb{E}_{\underline{m}} \sup_{x \in D} |\sum_{k=1}^{M} (m_k - L_k m) \mathbb{E}_{P_k} h(x)| + \frac{v}{m} \mathbb{E}_{\underline{m}} \mathbb{E}_{\underline{a}|\underline{m}} \sup_{x \in D} |\sum_{k=1}^{M} \sum_{j=1}^{n_k} \frac{m_k}{n_k} (h_{j,k}(x) - \mathbb{E}_{P_k} h(x))|.$$
(8)

Now

$$\mathbb{E}_{\underline{a}|\underline{m}} \sup_{x \in D} |\sum_{k=1}^{M} \sum_{j=1}^{n_{k}} \frac{m_{k}}{n_{k}} (h_{j,k}(x) - \mathbb{E}_{P_{k}} h(x))| \\ \leq 2\mathbb{E}_{\underline{a}|\underline{m}} \sup_{x \in D} |\sum_{k=1}^{M} \sum_{j=1}^{n_{k}} \sigma_{j,k} \frac{m_{k}}{n_{k}} [h_{j,k}(x) - \mu_{j,k}(x)]|, \quad (9)$$

where  $\{\sigma_{j,k}\}$  is a sequence of independent identically distributed Rademacher variables and  $\{x \mapsto \mu_{j,k}(x)\}$  is any sequence of functions defined on *D* [see for example Lemma 2.3.6 in [20]]. For notational brevity, we define  $\tilde{h}_{j,k}(x) = \frac{m_k}{n_k} [h_{j,k}(x) - \mu_{j,k}(x)]$ . By Dudley's entropy integral method (see [21, Corollary 13.2])], the quantity in (9) can be bounded by

7652

$$24\int_0^{\delta/2}\sqrt{N(u,D)}du,\tag{10}$$

where N(u, D) is the *u*-metric entropy of *D* with respect to the norm  $\kappa(x, x')$  (i.e., the logarithm of the smallest *u*-net that covers *D* with respect to  $\kappa$ ) defined by

$$\kappa^{2}(x, x') \triangleq \sum_{k=1}^{M} \sum_{j=1}^{n_{k}} (\widetilde{h}_{j,k}(x) - \widetilde{h}_{j,k}(x'))^{2} \\\leq (m+M) \|x - x'\|_{\infty}^{2},$$
(11)

and  $\delta^2 = \sup_{x \in D} \sum_{k=1}^{M} \sum_{j=1}^{n_k} |\tilde{h}_{j,k}(x)|^2$ . If we set  $\mu_{j,k}$  to equal  $\frac{m_k}{n_k} h(\eta(t_k, a_k), t_k, a_k)$ , where  $(\eta_k, t_k, a_k)'$  is any fixed point in  $\mathcal{B}_k$ , it follows from (6) and (7) that  $\delta \leq \sqrt{m + M\epsilon}$  and from (11) that  $N(u, D) \leq d \log(3\sqrt{m + M}/u)$ . By evaluating the integral in (10), we can bound the second term in (8) by

$$24v\sqrt{d} \ m^{-1/2}\epsilon\sqrt{-\log\epsilon+1}\sqrt{1+M/m}.$$
 (12)

For the first expectation in (8), we follow a similar approach. As before,

$$\mathbb{E}_{\underline{m}} \sup_{x \in D} |\sum_{k=1}^{M} (m_k - L_k m) \mathbb{E}_{P_k} h(x)|$$

$$\leq 2 \mathbb{E}_{\underline{m}} \sup_{x \in D} |\sum_{k=1}^{M} \sigma_k (m_k - L_k m) \mathbb{E}_{P_k} h(x)|, \quad (13)$$

where  $\{\sigma_k\}$  is a sequence of independent identically distributed Rademacher variables. For notational brevity, we write  $\tilde{h}_k(x) = (m_k - L_k m) \mathbb{E}_{P_k} h(x)$ . We can also bound (13) by (10), except this time N(u, D) is the *u*-metric entropy of *D* with respect to the norm  $\rho(x, x')$  defined by

$$\rho^{2}(x, x') \triangleq \sum_{k=1}^{M} (\widetilde{h}_{k}(x) - \widetilde{h}_{k}(x'))^{2}$$
$$\leq M \|x - x'\|_{\infty}^{2}, \qquad (14)$$

where the last line follows from  $|m_k - L_k m| \le 1$  and  $|\mathbb{E}_{P_k}h(x) - \mathbb{E}_{P_k}h(x')| \le ||x - x'||_{\infty}$ . The quantity  $\delta$  is also less than  $\sqrt{M}$ , since  $\sup_{x \in D} |\tilde{h}_k(x)| \le 1$  and moreover  $N(u, D) \le d \log(3\sqrt{M}/u)$ . Evaluating the integral in (10) with these specifications yields a bound on the first term in (8) of

$$\frac{48v\sqrt{d}\sqrt{M}}{m}.$$
 (15)

Adding (15) and (12) together yields a bound on  $\mathbb{E} \sup_{x \in D} |\overline{f}_m(x) - f(x)|$  of

$$48v\sqrt{d}m^{-1/2}(\sqrt{M/m} + \epsilon\sqrt{1 + M/m}\sqrt{-\log\epsilon + 1}).$$
 (16)

Choose

$$M = m \frac{\epsilon^2 (-\log \epsilon + 1)}{1 - \epsilon^2 (-\log \epsilon + 1)}.$$
(17)

Consequently,  $\mathbb{E} \sup_{x \in D} |\overline{f}_m(x) - f(x)|$  is at most

$$96v\sqrt{d}m^{-1/2}\frac{\epsilon\sqrt{-\log\epsilon+1}}{\sqrt{1-\epsilon^2(-\log\epsilon+1)}}.$$
(18)

We stated earlier that  $M \simeq \epsilon^{-d}$ . Thus (17) determines  $\epsilon$  to be at most of order  $m^{-1/(d+2)}$ . Since the inequality (17) holds on average, there is a realization of  $\overline{f}_m$  for which  $\sup_{x\in D} |\overline{f}_m(x) - f(x)|$  has the same bound. Note that  $\overline{f}_m$  has the desired equally weighted form.

For the second conclusion, we set  $m_k = mL_k$  and  $n_k = \lceil m_k \rceil$ . In this case, the first term in (8) is zero and hence  $\mathbb{E} \sup_{x \in D} |\overline{f_m}(x) - f(x)|$  is not greater than (12). The conclusion follows with M = m and  $\epsilon$  of order  $m^{-1/d}$ .

**Case II:** s = 3. The metric  $\kappa(x, x')$  is in fact bounded by a constant multiple of  $\sqrt{m + M\epsilon} ||x - x'||_{\infty}$ . To see this, we note that the function  $\tilde{h}_{j,k}(x)$  has the form

$$\pm \frac{m_k}{n_k} [(a \cdot x - t)_+^2 - (a_k \cdot x - t_k)_+^2],$$

with  $||a - a_k||_1 + |t - t_k| < \epsilon$ . Thus, the gradient of  $\tilde{h}_{j,k}(x)$  with respect to x has the form

$$\nabla \widetilde{h}_{j,k}(x) = \pm \frac{2m_k}{n_k} [(a(a \cdot x - t)_+ - a_k(a_k \cdot x - t_k)_+].$$

Adding and subtracting  $\frac{2m_k}{n_k}a(a_k \cdot x - t_k)_+$  to the above expression yields the bound of order  $\epsilon$  for  $\sup_{x \in D} \|\nabla \tilde{h}_{j,k}(x)\|_1$ . Taylor's theorem yields the desired bound on  $\kappa(x, x')$ . Again using Dudley's entropy integral, we can bound  $\mathbb{E} \sup_{x \in D} |\overline{f_m}(x) - f(x)|$  by a universal constant multiple of either  $v\sqrt{dm^{-1/2}}(\sqrt{M/m} + \epsilon\sqrt{1+M/m})$  or  $v\sqrt{dm^{-1/2}}\epsilon\sqrt{1+M/m}$  corresponding to the equally weighted or non-equally weighted cases, respectively. The corresponding results follow with  $M = m\epsilon^2/(1-\epsilon^2)$  and  $\epsilon$  of order  $m^{-1/(d+2)}$  or M = m and  $\epsilon$  of order  $m^{-1/d}$ . Note that here the additional smoothness afforded by the stronger assumption  $v_{f,3} < +\infty$  allows one to remove the  $\sqrt{-\log \epsilon + 1}$  factor that appeared in the final bound in the proof of Theorem 2. This rate is the same as what was achieved in Theorem 2. Without a  $\sqrt{(\log m)/d + 1}$  factor.  $\Box$ 

*Proof of Theorem 2:* If  $|z| \le c$ , we note the identity

$$-\int_{0}^{c} [(z-u)_{+}e^{iu} + (-z-u)_{+}e^{-iu}]du = e^{iz} - iz - 1.$$
(19)

If  $c = \|\omega\|_1$ ,  $z = \omega \cdot x$ ,  $a = a(\omega) = \omega/\|\omega\|_1$ , and  $u = \|\omega\|_1 t$ ,  $0 \le t \le 1$ , we find that

$$-\|\omega\|_{1}^{2} \int_{0}^{1} [(a \cdot x - t)_{+} e^{i\|\omega\|_{1}t} + (-a \cdot x - t)_{+} e^{-i\|\omega\|_{1}t}]dt$$
$$= e^{i\omega \cdot x} - i\omega \cdot x - 1.$$

Multiplying the above by  $\mathcal{F}(f)(\omega) = e^{ib(\omega)} |\mathcal{F}(f)(\omega)|$ , integrating over  $\mathbb{R}^d$ , and applying Fubini's theorem yields

$$f(x) - x \cdot \nabla f(0) - f(0) = \int_{\mathbb{R}^d} \int_0^1 g(t, \omega) dt d\omega,$$

where

$$g(t,\omega) = -[(a \cdot x - t)_{+} \cos(\|\omega\|_{1}t + b(\omega)) + (-a \cdot x - t)_{+} \cos(\|\omega\|_{1}t - b(\omega))]\|\omega\|_{1}^{2}|\mathcal{F}(f)(\omega)|$$

Consider the probability measure on  $\{-1, 1\} \times [0, 1] \times \mathbb{R}^d$  defined by

$$dP(z,t,\omega) = \frac{1}{v} |\cos(z\|\omega\|_1 t + b(\omega))| \|\omega\|_1^2 |\mathcal{F}(f)(\omega)| dt d\omega,$$
(20)

Authorized licensed use limited to: New York University. Downloaded on April 03,2024 at 00:18:36 UTC from IEEE Xplore. Restrictions apply.

where

$$v = \int_{\mathbb{R}^d} \int_0^1 [|\cos(\|\omega\|_1 t + b(\omega))| + |\cos(\|\omega\|_1 t - b(\omega))|] \|\omega\|_1^2 |\mathcal{F}(f)(\omega)| dt d\omega \le 2v_{f,2}.$$

Define a function h(z, t, a)(x) that equals

$$(za \cdot x - t)_+ \eta(z, t, \omega),$$

where  $\eta(z, t, \omega) = -\operatorname{sgn} \cos(\|\omega\|_1 zt + b(\omega))$ . Note that h(z, t, a)(x) has the form  $\pm(\pm a \cdot x - t)_+$ . Thus, we see that

$$f(x) - x \cdot \nabla f(0) - f(0) = v \int_{\{-1,1\} \times [0,1] \times \mathbb{R}^d} h(z,t,a)(x) dP(z,t,\omega).$$
(21)

The result follows from an application of Theorem 1.  $\Box$ 

*Proof of Theorem 3:* For the result in Theorem 3, we will use exactly the same techniques. The function  $f(x) - x^T \nabla \nabla^T f(0) x/2 - x \cdot \nabla f(0) - f(0)$  can be written as the real part of

$$\int_{\mathbb{R}^d} (e^{i\omega \cdot x} + (\omega \cdot x)^2/2 - i\omega \cdot x - 1)\mathcal{F}(f)(\omega)d\omega.$$
 (22)

As before, the integrand in (22) admits an integral representation given by

$$(i/2)\|\omega\|_1^3 \int_0^1 [(-a \cdot x - t)_+^2 e^{-i\|\omega\|_1 t} - (a \cdot x - t)_+^2 e^{i\|\omega\|_1 t}] dt,$$

which can be used to show that  $f(x) - x^T \nabla \nabla^T f(0) x/2 - x \cdot \nabla f(0) - f(0)$  equals

$$\frac{v}{2} \int_{\{-1,1\}\times[0,1]\times\mathbb{R}^d} h(z,t,a)(x) dP(z,t,\omega),$$
(23)

where

$$h(z, t, a) = \operatorname{sgn} \sin(z \|\omega\|_1 t + b(\omega)) (za \cdot x - t)_+^2$$

and

$$dP(z,t,\omega) = \frac{1}{\nu} |\sin(z||\omega||_1 t + b(\omega))|||\omega||_1^3 |\mathcal{F}(f)(\omega)| dt d\omega,$$

$$v = \int_{\mathbb{R}^d} \int_0^1 [|\sin(\|\omega\|_1 t + b(\omega))| + |\sin(\|\omega\|_1 t - b(\omega))|] \|\omega\|_1^3 |\mathcal{F}(f)(\omega)| dt d\omega \le 2v_{f,3}.$$

The result follows from an application of Theorem 1.  $\Box$ 

Remark 2: By slightly modifying the definition of h from the proofs of Theorem 2 and Theorem 3 (in particular, multiplying it by a sinusoidal function of  $\omega$  and t), it suffices to sample instead from the density  $dP(t, \omega) = \frac{\|\omega\|_1^s |\mathcal{F}(f)(\omega)|}{v_{f,s}} dt d\omega$  on  $[0, 1] \times \mathbb{R}^d$ .

Remark 3: For unit bounded x, the expression  $e^{i\omega \cdot x} - i\omega \cdot x - 1$  is bounded in magnitude by  $\|\omega\|_1^2$ , so one only needs Fourier representation of  $f(x) - x \cdot \nabla f(0) - f(0)$  when using the integrability with the  $\|\omega\|_1^2$  factor. Similarly,  $e^{i\omega \cdot x} + (\omega \cdot x)^2/2 - i\omega \cdot x - 1$  is bounded in magnitude by  $\|\omega\|_3^3$ , so one only needs Fourier representation of  $f(x) - x^T \nabla \nabla^T f(0)x - x \cdot \nabla f(0) - 1$  when using the integrability with the  $\|\omega\|_3^3$  factor. Remark 4: Note that in Theorem 2 and Theorem 3, we work with integrals with respect to the absolutely continuous measure  $d\mathcal{F}(f)(\omega)$ . In general, a (complex) Fourier measure  $d\mathcal{F}(f)(\omega)$  does not need to be absolutely continuous. For instance, it can be discrete on a lattice of values of  $\omega$ , associated with a multivariate Fourier series representation for bounded domains x (and periodic extensions thereof). Indeed, for bounded domains, one might have access to both Fourier series and Fourier transforms of extensions of f to  $\mathbb{R}^d$ . The best extension is one that gives the smallest Fourier norm  $\int_{\mathbb{R}^d} \|\omega\|_1^s |d\mathcal{F}(f)(\omega)|$ . For further discussion along these lines, see [6].

Next, we investigate the optimality of the rates from Section II.

## B. Lower Bounds

Let  $\mathcal{H}_s = \{x \mapsto \eta(a \cdot x - t)_+^{s-1} : \|a\|_1 \leq 1, 0 \leq t \leq 1, \eta \in \{-1, +1\}\}$  and for  $p \in [2, +\infty]$  let  $\mathcal{F}_p^s$  denote the closure of the convex hull of  $\mathcal{H}_s$  with respect to the  $\|\cdot\|_p$  norm on  $L^p(D, P)$  for p finite, where P is the uniform probability measure on D, and  $\|\cdot\|_{\infty}$  (the supremum norm over D) for  $p = +\infty$ . We let  $\mathcal{C}_m^s$  denote the collection of all convex combinations of m terms from  $\mathcal{H}_s$ . By Theorem 2 and Theorem 3, after possibly subtracting a linear or quadratic term,  $f/(2v_{f,2})$  and  $f/v_{f,3}$  belongs to  $\mathcal{F}_p^2$  and  $\mathcal{F}_p^3$ , respectively. For  $p \in [2, +\infty]$  and  $\epsilon > 0$ , we define the  $\epsilon$ -covering number  $N_p(\epsilon)$  by

$$\min\{n: \exists \mathcal{F} \subset \mathcal{F}_p^s, |\mathcal{F}| = n, \text{ s.t. } \inf_{f' \in \mathcal{F}} \sup_{f \in \mathcal{F}_p^s} ||f - f'||_p \le \epsilon\}.$$

and the  $\epsilon$ -packing number  $M_p(\epsilon)$  by

$$\max\{n: \exists \mathcal{F} \subset \mathcal{F}_p^s, |\mathcal{F}| = n, \text{ s.t.} \inf_{f, f' \in \mathcal{F}} \|f - f'\|_p > \epsilon\}.$$

Theorem 1 implies that  $\inf_{f_m \in C_m^s} \sup_{f \in \mathcal{F}_\infty^s} ||f - f_m||_{\infty}$  achieves the bounds as stated therein.

*Theorem 4: For*  $p \in [2, +\infty]$  *and*  $s \in \{2, 3\}$ *,* 

$$\inf_{f_m \in \mathcal{C}_m^s} \sup_{f \in \mathcal{F}_p^s} \|f - f_m\|_p \ge (Amd^{2s+1}\log(md))^{-1/2-s/d},$$

for some universal positive constant A.

Ignoring the dependence on d and logarithmic factors in m, this result coupled with Theorem 1 implies that  $\inf_{f_m \in C_m^2} \sup_{f \in \mathcal{F}_p^2} ||f - f_m||_p$  is between  $m^{-1/2-2/d}$  and  $m^{-1/2-1/d}$ ; for large d, the rates are essentially the same. Compare this with [9, Th. 4] or [5, Th. 3], where a lower bound of  $c(\delta, d) m^{-1/2-1/d-\delta}$ ,  $\delta > 0$  arbitrary, was obtained for approximants of the form (1) with Lipschitz  $\phi$ , but with inner parameter vectors of unbounded  $\ell^1$  norm.

We only give the proof of Theorem 4 for s = 2, since the other case s = 3 is handled similarly. First, we provide a few ancillary results that will be used later on. The next result is contained in [22, Lemma 4.2] and is useful for giving a lower bound on  $M_p(\epsilon)$ .

Lemma 1: Let H be a Hilbert space equipped with a norm  $\|\cdot\|$  and containing a finite set  $\mathcal{H}$  with the following properties. (i)  $|\mathcal{H}| \geq 3$ ,

- (*ii*)  $\sum_{h,h'\in\mathcal{H}, h\neq h'} |\langle h, h' \rangle| \le \delta^2$ (*iii*)  $\delta^2 \le \min_{h\in\mathcal{H}} \|h\|^2$

Then there exists a collection  $\Omega \subset \{0,1\}^{|\mathcal{H}|}$  with car-dinality at least  $2^{(1-H(1/4))|\mathcal{H}|-1}$ , where H(1/4) is the entropy of a Bernoulli random variable with success probability 1/4, such that each pair of elements in the set  $\mathcal{F} = \left\{ \frac{1}{|\mathcal{H}|} \sum_{h \in \mathcal{H}} \omega_h h : (\omega_h : h \in \mathcal{H}) \in \Omega \right\} \text{ is separated by at}$ least  $\frac{1}{2}\sqrt{\frac{\min_{h\in\mathcal{H}}\|h\|^2-\delta^2}{|\mathcal{H}|}}$  in  $\|\cdot\|$ .

Lemma 2: If  $\theta$  belongs to  $[R]^d = \{1, 2, \dots, R\}^d$ ,  $R \in \mathbb{Z}^+$ , then the collection of functions

$$\mathcal{H} = \{x \mapsto \sin(\pi \theta \cdot x) / (4\pi \|\theta\|_1^2) : \theta \in [R]^d\}$$

satisfies the assumption of Lemma 1 with  $H = L^2(D, P)$ , where P is the uniform probability measure on D. Moreover,  $|\mathcal{H}| = R^d$ ,  $\delta = 0$ ,  $\min_{h \in \mathcal{H}} ||h|| = 1/(4\sqrt{2\pi}d^2R^2)$ , and  $\mathcal{F} \subset \mathcal{F}_p^1$  for all  $p \in [2, +\infty]$ . Consequently, if  $\epsilon = 1/(8\sqrt{2\pi} d^2 R^{2+d/2})$ , then

$$\log M_p(\epsilon) \ge (\log 2)(1 - H(1/4)) \left(8\epsilon \sqrt{2\pi} d^2\right)^{-\frac{2d}{4+d}} - 1$$
$$\ge \left(c\epsilon d^2\right)^{-\frac{2d}{4+d}},\tag{24}$$

for some universal constant c > 0.

*Proof:* We first observe the identity

$$\begin{aligned} \sin(\pi \theta \cdot x) / (4\pi \|\theta\|_1^2) \\ &= \theta \cdot x / (4\pi \|\theta\|_1^2) \\ &+ \frac{\pi}{4} \int_0^1 [(-a \cdot x - t)_+ - (a \cdot x - t)_+] \sin(\pi \|\theta\|_1 t) dt, \end{aligned}$$

where  $a = a(\theta) = \theta / \|\theta\|_1$ . Note that above integral can also be written as an expectation of

$$-z \operatorname{sgn}(\sin(\pi \|\theta\|_1 t)) (za \cdot x - t)_+ \in \mathcal{H}_2$$

with respect to the density

$$p_{\theta}(z,t) = \frac{\pi}{4} |\sin(\pi \|\theta\|_1 t)|,$$

on  $\{-1, 1\} \times [0, 1]$ . The fact that  $p_{\theta}$  integrates to one is a consequence of the identity

$$\int_0^1 |\sin(\pi \, \|\theta\|_1 t)| dt = 2/\pi$$

Since  $\int_D |\sin(\pi\theta \cdot x)|^2 dP(x) = 1/2$ , each member of  $\mathcal{H}$ has norm equal to  $1/(4\sqrt{2\pi} \|\theta\|_1^2)$  and each pair of elements is orthogonal so that  $\delta = 0$ . Integrations over D involving  $\sin(\pi\theta \cdot x)$  are easiest to see using an instance of Euler's formula, viz.,  $\sin(\alpha \cdot x) = \frac{1}{2i} (\prod_{k=1}^{d} e^{i\alpha(k)x(k)} - \prod_{k=1}^{d} e^{-i\alpha(k)x(k)}).$ 

*Proof of Theorem 4:* Let A > 0 be arbitrary. Suppose contrary to the hypothesis,

$$\inf_{f_m \in \mathcal{C}_m^2} \sup_{f \in \mathcal{F}_p^2} \|f - f_m\|_p < (Amd^5 \log(md))^{-1/2 - 2/d}$$
$$\triangleq \epsilon_0/3.$$

Note that each element of  $C_m^2$  has the form  $\sum_{k=1}^m \lambda_k h_k$ , where  $\sum_{k=1}^m \lambda_k = 1$  and  $h_k \in \mathcal{H}_s$ . Next, consider the subcollection  $\tilde{C}_m^2$  with elements of the form  $\sum_{k=1}^m \tilde{\lambda}_k \tilde{h}_k$ , where  $\tilde{\lambda}_k$   $in D = [-1, 1]^d$  and  $s \in \{2, 3\}$ , where P is a probability measure on  $[0, 1] \times \{a \in \mathbb{R}^d : \|a\|_1 = 1\}$  and  $\eta(t, a)$  is either

belongs to an  $\epsilon_0/3$ -net  $\widetilde{\mathcal{P}}$  of the m-1 dimensional probability simplex  $\mathcal{P}_m$  and  $\tilde{h}_k$  belongs to an  $\epsilon_0/3$ -net  $\mathcal{H}$  of  $\mathcal{H}_s$ . By a stars and bars argument, there are at most  $|\mathcal{P}|\binom{m+|\mathcal{H}|-1}{m}$  such functions. Furthermore, since  $\sup_{h \in \mathcal{H}_{\epsilon}} \|h\|_{\infty} \leq 1$ , we have

$$\inf_{f_m \in \widetilde{\mathcal{C}}_m^2} \sup_{f \in \mathcal{F}_p^2} \|f - f_m\|_2 \leq \inf_{f_m \in \mathcal{C}_m^2} \sup_{f \in \mathcal{F}_p^2} \|f - f_m\|_2 + \inf_{\widetilde{h} \in \widetilde{\mathcal{H}}} \sup_{h \in \mathcal{H}_s} \|h - \widetilde{h}\|_2 + \inf_{\widetilde{\lambda} \in \widetilde{\mathcal{P}}} \sup_{\lambda \in \mathcal{P}_m} \|\lambda - \widetilde{\lambda}\|_1 < \epsilon_0/3 + \epsilon_0/3 + \epsilon_0/3 = \epsilon_0.$$

Since  $|\widetilde{\mathcal{H}}| \simeq \epsilon_0^{-d-1}$  and  $|\widetilde{\mathcal{P}}| \simeq \epsilon_0^{-m+1}$ , it follows that

$$\log N_{p}(\epsilon_{0}) \leq \log |\tilde{C}_{m}^{2}|$$

$$\leq c_{0} \log \left[ \epsilon_{0}^{-m-1} \binom{m+c_{1}\epsilon_{0}^{-d-1}-1}{m} \right]$$

$$\leq c_{2}dm \log(1/\epsilon_{0})$$

$$\leq c_{3}dm \log(Adm), \qquad (25)$$

for some positive universal constants  $c_0 > 0$ ,  $c_1 > 0$ ,  $c_2 > 0$ , and  $c_3 > 0$ .

On the other hand, using (24) from Lemma 2 coupled with the fact that  $N_p(\epsilon_0) \ge M_p(2\epsilon_0)$ , we have

$$\log N_{p}(\epsilon_{0}) \geq \log M_{p}(2\epsilon_{0})$$
$$\geq \left(2c\epsilon_{0}d^{2}\right)^{-\frac{2d}{4+d}}$$
$$\geq c_{4}Adm\log(dm), \qquad (26)$$

for some universal constant  $c_4 > 0$ . Combining (25) and (26), we find that

 $c_4Adm\log(dm) \le c_3dm\log(Adm).$ 

If A is large enough (independent of m or d), we reach a contradiction. This proves the lower bound.  $\square$ 

# III. $L^2$ Approximation With Bounded $\ell^0$ and $\ell^1$ Norm

In Section II, we explored conditions for which good approximation in  $L^{\infty}(D)$  could be achieved even with  $\ell^1$  controls on the inner parameter vectors. In this section, we show how similar statements can be made in  $L^{2}(D)$ , but with control on the  $\ell^0$  norm as well. Note that unlike Theorem 1, we see in Theorem 5 how the smoothness of the activation function directly affects the rate of approximation. The proof is obtained by applying the Maurey-Jones-Barron probabilistic method in two stages (similar to two-stage cluster sampling), first on the outer layer coefficients, and then on the inner layer coefficients.

Theorem 5: Suppose f admits an integral representation

$$f(x) = v \int_{[0,1] \times \{a: \|a\|_1 = 1\}} \eta(t,a) \ (a \cdot x - t)_+^{s-1} dP(t,a),$$

-1 or +1. There exists a linear combination of ridge functions of the form

$$f_{m,m_0}(x) = \frac{v}{m} \sum_{k=1}^m b_k \left( a_k \cdot x - t_k \right)_+^{s-1},$$

where  $||a_k||_0 \le m_0$ ,  $||a_k||_1 = 1$ , and  $b_k \in \{-1, +1\}$  such that

$$||f - f_{m,m_0}||_2 \le v \sqrt{\frac{1}{m} + \frac{1}{m_0^{s-1}}}.$$

Furthermore, the same rates for s = 2 or s = 3 are achieved for general f adjusted by a linear or quadratic term with  $v = 2v_{f,2} < +\infty$  or  $v = v_{f,3} < +\infty$ , respectively.

Remark 5: In particular, taking  $m_0 = \sqrt{m}$ , it follows that there exists an m-term linear combination of squared ReLU ridge functions, with  $\sqrt{m}$ -sparse inner parameter vectors, that approximates f with  $L^2(D)$  error at most  $\sqrt{2} \text{cm}^{-1/2}$ . In other words, the  $L^2(D)$  approximation error is inversely proportional to the inner layer sparsity and it need only be sublinear in the outer layer sparsity.

*Proof:* Take a random sample  $\underline{a} = \{(t_k, a_k)'\}_{1 \le k \le m}$ from *P*. Given  $\underline{a}$ , take a random sample  $\underline{\widetilde{a}} = \{\widetilde{a}_{\ell,k}\}_{1 \le \ell \le m_0, \ 1 \le k \le m}$ , where  $\mathbb{P}[\widetilde{a}_{\ell,k} = \operatorname{sgn}(a_k(j))e_j] = |a_k(j)|$ for  $j = 1, \ldots, d$ ,  $a_k = (a_k(1), \ldots, a_k(d))'$ , and  $e_j$  is the *j*-th standard basis vector for  $\mathbb{R}^d$ . Note that

$$\mathbb{E}_{\underline{\tilde{a}}|\underline{a}}[\overline{\tilde{a}}_{\ell,k}] = a_k \tag{27}$$

and

$$\operatorname{Var}_{\underline{\widetilde{a}}|\underline{a}}[\widetilde{a}_{\ell,k} \cdot x] \leq \mathbb{E}_{\underline{\widetilde{a}}|\underline{a}}[\widetilde{a}_{\ell,k} \cdot x]^2 = \sum_{j=1}^d |a_k(j)| |x(j)|^2$$
$$\leq ||a_k||_1 ||x||_{\infty}^2 \leq 1.$$
(28)

Define

$$\overline{f}_{m,m_0}(x) = \frac{v}{m} \sum_{k=1}^m \eta(t_k, a_k) \left( \frac{1}{m_0} \sum_{\ell=1}^{m_0} \widetilde{a}_{\ell,k} \cdot x - t_k \right)_+^{s-1}.$$
 (29)

By the bias-variance decomposition,

$$\mathbb{E}\|f-\overline{f}_{m,m_0}\|_2^2 = \mathbb{E}\|\overline{f}_{m,m_0} - \mathbb{E}\overline{f}_{m,m_0}\|_2^2 + \|f-\mathbb{E}\overline{f}_{m,m_0}\|_2^2.$$

Note that  $\mathbb{E} \|\overline{f}_{m,m_0} - \mathbb{E}\overline{f}_{m,m_0}\|_2^2 \leq \frac{v^2}{m}$ . Next, observe that

$$f(x) - \mathbb{E}f_{m,m_0}(x)$$

$$= \frac{v}{m} \sum_{k=1}^{m} \mathbb{E}_{\underline{a}} \bigg[ \eta(t_k, a_k) \times \mathbb{E}_{\underline{\widetilde{a}}|\underline{a}} \left( (a_k \cdot x - t_k)_+^{s-1} - \left( \frac{1}{m_0} \sum_{\ell=1}^{m_0} \widetilde{a}_{\ell,k} \cdot x - t_k \right)_+^{s-1} \right) \bigg],$$

which, by an application of the triangle inequality, implies that

$$|f(x) - \mathbb{E}f_{m,m_0}(x)| \leq \frac{v}{m} \sum_{k=1}^m \mathbb{E}_{\underline{a}} \left| (a_k \cdot x - t_k)_+^{s-1} - \mathbb{E}_{\underline{\widetilde{a}}|\underline{a}} \left( \frac{1}{m_0} \sum_{\ell=1}^{m_0} \widetilde{a}_{\ell,k} \cdot x - t_k \right)_+^{s-1} \right|.$$

Next, we use the following two properties of  $(z)_+^{s-1}$ : for all z and z' in  $\mathbb{R}$ ,

$$|(z)_{+} - (z')_{+}| \le |z - z'|, \qquad (30)$$
  
$$|(z)_{+}^{2} - (z')_{+}^{2} - 2(z - z')(z')_{+}| \le |z - z'|^{2}. \qquad (31)$$

If s = 2, we have by (30), (27), and (28) that

$$\begin{split} \mathbb{E}_{\underline{a}} \left| (a_k \cdot x - t_k)_+ - \mathbb{E}_{\underline{\widetilde{a}}|\underline{a}} \left( \frac{1}{m_0} \sum_{\ell=1}^{m_0} \widetilde{a}_{\ell,k} \cdot x - t_k \right)_+ \right| \\ &\leq \mathbb{E}_{\underline{a}} \mathbb{E}_{\underline{\widetilde{a}}|\underline{a}} \left| a_k \cdot x - \frac{1}{m_0} \sum_{\ell=1}^{m_0} \widetilde{a}_{\ell,k} \cdot x \right| \\ &\leq \mathbb{E}_{\underline{a}} \sqrt{\mathbb{E}_{\underline{\widetilde{a}}|\underline{a}}} \left| a_k \cdot x - \frac{1}{m_0} \sum_{\ell=1}^{m_0} \widetilde{a}_{\ell,k} \cdot x \right|^2 \\ &= \mathbb{E}_{\underline{a}} \sqrt{\frac{\operatorname{Var}_{\underline{\widetilde{a}}|\underline{a}}[\widetilde{a}_{\ell,k} \cdot x]}{m_0}} \leq \frac{1}{\sqrt{m_0}}. \end{split}$$

This shows that  $\|f - \mathbb{E}\overline{f}_{m,m_0}\|_2^2 \leq \frac{v^2}{m_0}$ . If s = 3, we have from (31), (27), and (28) that

$$\mathbb{E}_{\underline{a}} \left| (a_k \cdot x - t_k)_+^2 - \mathbb{E}_{\underline{\widetilde{a}}|\underline{a}} \left( \frac{1}{m_0} \sum_{\ell=1}^{m_0} \widetilde{a}_{\ell,k} \cdot x - t_k \right)_+^2 \right|$$
  
$$\leq \mathbb{E}_{\underline{a}} \mathbb{E}_{\underline{\widetilde{a}}|\underline{a}} \left| a_k \cdot x - \frac{1}{m_0} \sum_{\ell=1}^{m_0} \widetilde{a}_{\ell,k} \cdot x \right|^2$$
  
$$= \mathbb{E}_{\underline{a}} \left[ \frac{\mathsf{Var}_{\underline{\widetilde{a}}|\underline{a}}[\widetilde{a}_{\ell,k} \cdot x]}{m_0} \right] \leq \frac{1}{m_0}.$$

This shows that  $||f - \mathbb{E}\overline{f}_{m,m_0}||_2^2 \leq \frac{v^2}{m_0^2}$ . Since these bounds hold on average, there exists a realization of (29) for which the bounds are also valid. Note that the vector  $\frac{1}{m_0} \sum_{\ell=1}^{m_0} \widetilde{a}_{\ell,k}$  has  $\ell^0$  norm at most  $m_0$  and unit  $\ell^1$  norm.

The fact that the bounds also hold for f adjusted by a linear or quadratic term (under an assumption of finite  $v_{f,2}$  or  $v_{f,3}$ ) follows from (21) and (23).

#### ACKNOWLEDGMENT

The authors would like to thank the Associate Editor and two anonymous reviewers whose detailed feedback led to dramatic improvements to this paper. They also thank Joowon Kim for helpful comments on earlier drafts of this manuscript.

#### REFERENCES

- [1] A. R. Barron, "Approximation and estimation bounds for artificial neural networks," *Mach. Learn.*, vol. 14, no. 1, pp. 115–133, 1994.
- [2] G. Raskutti, M. J. Wainwright, and B. Yu, "Minimax rates of estimation for high-dimensional linear regression over *l<sub>q</sub>*-balls," *IEEE Trans. Inf. Theory*, vol. 57, no. 10, pp. 6976–6994, Oct. 2011, doi: 10.1109/TIT. 2011.2165799.
- [3] J. M. Klusowski and A. R. Barron, "Risk bounds for high-dimensional ridge function combinations including neural networks," Work. Paper, 2018.
- [4] J. M. Klusowski and A. R. Barron, "Minimax lower bounds for ridge combinations including neural nets," in *Proc. IEEE Int. Symp. Inf. Theory*, Aachen, Germany, Jun. 2017, pp. 1377–1380.
- [5] A. R. Barron, "Neural net approximation," in Proc. Yale Workshop Adapt. Learn. Syst. New Haven, CT, USA: Yale Univ. Press, 1992.

- [6] A. R. Barron, "Universal approximation bounds for superpositions of a sigmoidal function," *IEEE Trans. Inf. Theory*, vol. 39, no. 3, pp. 930–945, May 1993, doi: 10.1109/18.256500.
- [7] G. H. L. Cheang and A. R. Barron, "A better approximation for balls," J. Approximation Theory, vol. 104, no. 2, pp. 183–203, 2000, doi: 10.1006/ jath.1999.3441.
- [8] J. E. Yukich, M. B. Stinchcombe, and H. White, "Sup-norm approximation bounds for networks through probabilistic methods," *IEEE Trans. Inf. Theory*, vol. 41, no. 4, pp. 1021–1027, Jul. 1995, doi: 10.1109/18. 391247.
- [9] Y. Makovoz, "Random approximants and neural networks," J. Approximation Theory, vol. 85, pp. 98–109, Apr. 1996, doi: 10.1006/jath. 1996.0031.
- [10] P. L. Bartlett, "The sample complexity of pattern classification with neural networks: The size of the weights is more important than the size of the network," *IEEE Trans. Inf. Theory*, vol. 44, no. 2, pp. 525–536, Mar. 1998, doi: 10.1109/18.661502.
- [11] S. Ioannidis and A. Montanari. (2017). "Learning combinations of sigmoids through gradient estimation." [Online]. Available: https://arxiv. org/abs/1708.06678
- [12] M. Janzamin, H. Sedghi, and A. Anandkumar. (2015). "Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods." [Online]. Available: https://arxiv.org/abs/1506.08473
- [13] Y. Zhang, J. D. Lee, M. J. Wainwright, and M. I. Jordan. (2015). "Learning halfspaces and neural networks with random initialization." [Online]. Available: https://arxiv.org/abs/1511.07948
- [14] A. Brutzkus and A. Globerson. (Feb. 2017). "Globally optimal gradient descent for a ConvNet with Gaussian inputs." [Online]. Available: https://arxiv.org/abs/1702.07966
- [15] V. N. Vapnik and A. J. Červonenkis, "The uniform convergence of frequencies of the appearance of events to their probabilities," *Teor. Verojatnost. Primenen.*, vol. 16, pp. 264–279, 1971.
- [16] D. Haussler, "Sphere packing numbers for subsets of the Boolean *n*-cube with bounded Vapnik-Chervonenkis dimension," *J. Combinat. Theory, A*, vol. 69, no. 2, pp. 217–232, 1995, doi: 10.1016/0097-3165(95)90052-7.
- [17] L. Breiman, "Hinging hyperplanes for regression, classification, and function approximation," *IEEE Trans. Inf. Theory*, vol. 39, no. 3, pp. 999–1013, May 1993, doi: 10.1109/18.256506.
- [18] Y. Makovoz, "Uniform approximation by neural networks," J. Approximation Theory, vol. 95, no. 2, pp. 215–228, 1998, doi: 10.1006/jath. 1997.3217.
- [19] J. Neyman, "On the two different aspects of the representative method: The method of stratified sampling and the method of purposive selection," J. Roy. Stat. Soc., vol. 97, no. 4, pp. 558–625, 1934.

- [20] A. W. van der Vaart and J. A. Wellner, Weak Convergence and Empirical Processes: With Applications to Statistics (Springer Series in Statistics). New York, NY, USA: Springer-Verlag, 1996, doi: 10.1007/978-1-4757-2545-2.
- [21] S. Boucheron, G. Lugosi, and P. Massart, "A nonasymptotic theory of independence, with a foreword by Michel Ledoux," in *Concentration Inequalities*. Oxford, U.K.: Oxford Univ. Press, 2013, doi: 10.1093/ acprof:oso/9780199535255.001.0001.
- [22] V. Kůrková and M. Sanguineti, "Estimates of covering numbers of convex sets with slowly decaying orthogonal subsets," *Discrete Appl. Math.*, vol. 155, no. 15, pp. 1930–1942, 2007, doi: 10.1016/j.dam.2007.04.007.

**Jason M. Klusowski** (S'16) received a B.S. (Hons.) in Statistics and Mathematics from the University of Manitoba in 2013, M.A. in Statistics from Yale University in 2017, and Ph.D. in Statistics and Data Science from Yale University in 2018. He is currently Assistant Professor in the Department of Statistics and Biostatistics at Rutgers University–New Brunswick. His research interests are in high-dimensional statistics and network analysis.

Andrew R. Barron (S'84-M'85-SM'00-F'12) has research interests in the areas of statistical information theory, model selection, the minimum description length principle, probability limit theorems, asymptotics of Bayes procedures, estimation of functions of many variables, artificial neural networks, approximation theory, investment theory, universal data compression, sparse regression codes for practical capacity-achieving communications, and trajectory control optimization. Barron is a fellow of the IEEE, Medallion Prize winner of the Institute of Mathematical Statistics and a winner along with Bertrand Clarke of the IEEE Thompson Prize (for Best Paper in all IEEE Journals for authors under 30 at time of submission). He has served as secretary of the Board of Governors and subsequently has served multiple terms as a member of the Board of Governors of the IEEE Information Theory Society. He has chaired the Thomas M. Cover Dissertation Prize Committee. Received Ph.D., Electrical Engineering, Stanford University; M.S., Electrical Engineering, Stanford University; B.S. E.E. and Math Science, Rice University. From 1985 - 1992 Andrew was Assistant and then Associate Professor of Statistics and Electrical & Computer Engineering, University of Illinois. From 1992 to present, Andrew is a Professor of Statistics and Data Science at Yale and has served terms as department chair, director of graduate studies, director of undergraduate studies in Statistics, director of undergraduate studies in Applied Mathematics, and courtesy appointments as Professor of Electrical Engineering at Yale.