# Approximation of functions with one-bit neural networks

C. Sinan Güntürk[*]        Weilin Li[†]

March 17, 2023

**Abstract**

The celebrated universal approximation theorems for neural networks roughly state that any reasonable function can be arbitrarily well-approximated by a network whose parameters are appropriately chosen real numbers. This paper examines the approximation capabilities of one-bit neural networks – those whose nonzero parameters are $\pm a$ for some fixed $a \neq 0$. One of our main theorems shows that for any $f \in C^s([0,1]^d)$ with $\|f\|_\infty < 1$ and error $\varepsilon$, there is a $f_{NN}$ such that $|f(\boldsymbol{x}) - f_{NN}(\boldsymbol{x})| \leq \varepsilon$ for all $\boldsymbol{x}$ away from the boundary of $[0,1]^d$, and $f_{NN}$ is either implementable by a $\{\pm 1\}$ quadratic network with $O(\varepsilon^{-2d/s})$ parameters or a $\{\pm\frac{1}{2}\}$ ReLU network with $O(\varepsilon^{-2d/s}\log(1/\varepsilon))$ parameters, as $\varepsilon \to 0$. We establish new approximation results for iterated multivariate Bernstein operators, error estimates for noise-shaping quantization on the Bernstein basis, and novel implementation of the Bernstein polynomials by one-bit quadratic and ReLU neural networks.

**Keywords:** neural network, quantization, one-bit, Bernstein, polynomial approximation

**MSC2020:** 41A10, 41A25, 41A36, 41A63, 42C15, 68T07

## 1 Introduction

### 1.1 Motivation

In this paper, we address the following question regarding the approximation capabilities of quantized neural networks: what classes of functions can be approximated by neural networks of a given size such that their weights and biases are constrained to be in a small set of allowable values, especially with regards to the extreme one-bit case? While this is an interesting mathematical question by itself and warrants special attention given the growing importance of machine learning across numerous scientific disciplines, here we illuminate two particularly important broader questions that motivate this theoretical study.

Our first motivation comes from a practical problem. State-of-the-art neural networks oftentimes contain a massive number of parameters and are trained on enormous computational machines. It appears that in regards to performance of neural networks, the "bigger is better" philosophy largely holds true [7]. As higher resolution audio, image, and video data become increasingly more common, the size of high performance networks will only continue to grow and require more computational resources to utilize. This conflicts with the desire to use them on portable and low power devices such as smartphones. Quantization is the process of replacing high resolution floating point numbers with coarser ones. It is a natural solution to this issue, since simpler binary

---

[*]Courant Institute, New York University. Email: gunturk@cims.nyu.edu

[†]City University of New York, City College. Email: wli6@ccny.cuny.edu

operations can help alleviate the costly computational burden that comes with using expensive floating point operations.

Our second motivation is related to the over-parameterization phenomena. Since state-of-the-art neural networks often contain many more parameters than both the number of training samples and data dimensionality [5], it is widely believed that there is not a unique set of parameters that can be used to represent a given function. If we imagine that the set of all parameters generating a prescribed function is a manifold embedded in a high dimensional parameter space, then our main question is closely related to whether this manifold is sufficiently close to some lattice point. Since certain parameter choices may be more desirable than others, this is a central question in not only quantization, but also in neural network compression and model reduction.

## 1.2 Quantized neural networks

In this paper, we exclusively examine strict neural networks. We use the adjective "strict" to emphasize that such networks do not have any skip connections and apply the same activation function to each node except for the final affine layer. More specifically, fix a function $\beta\colon \mathbb{R} \to \mathbb{R}$, and slightly abusing notation, for each $m \geq 1$, we extend it to a map $\beta\colon \mathbb{R}^m \to \mathbb{R}^m$ defined as $\beta(\boldsymbol{x}) := (\beta(x_1), \ldots, \beta(x_m))$.

**Definition 1.1.** A strict neural network with activation $\beta$ is any function $f\colon \mathbb{R}^d \to \mathbb{R}^m$ of the form,

$$f(\boldsymbol{x}) := W_L\beta(W_{L-1}\cdots\beta(W_1(\boldsymbol{x}))), \quad W_\ell(\boldsymbol{u}) := A_\ell\boldsymbol{u} + \boldsymbol{b}_\ell \quad \text{for} \quad \ell = 1, \ldots, L.$$

In this definition, each $A_\ell \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$ is referred to as a weight matrix, $\boldsymbol{b}_\ell \in \mathbb{R}^{N_\ell}$ is called a bias vector, $L$ is the number of layers, and $N_0 = d$ and $N_L = m$. For each $1 \leq \ell \leq L-1$, we refer to $\boldsymbol{u}_\ell := \beta(W_\ell\beta(\cdots\beta(W_1(\boldsymbol{x}))))$ as the $\ell$-th layer's output. This network has $L$ layers, and applies the same activation function to each node except for the final linear layer. It has $N_\ell$ nodes in layer $\ell$, hence has $\sum_{\ell=1}^{L} N_\ell$ nodes in total, and we define the number of parameters as however many nonzero entries in $\{A_\ell\}_{\ell=1}^{L}$ and $\{\boldsymbol{b}_\ell\}_{\ell=1}^{L}$ there are. Since we place no restrictions on the weight matrices' structures, our framework allows for fully connected and convolutional networks. We say a function $g$ can be implemented by neural network provided it can be written in the above form with appropriate weights.

Other common definitions of a neural network allow for additional operations and flexibility. For instance, some conventions allow for the use of skip connections whereby $\boldsymbol{u}_{\ell+1}$ is allowed to depend on any previous layer's outputs $\boldsymbol{u}_\ell, \ldots, \boldsymbol{u}_1$. Another example is the use of different activation functions per node and possibly in each layer, including the identity function which effectively bypasses an activation. Following the terminology from [19], to differentiate between the strict notion we use versus less stringent notions, we refer to the latter as generalized neural networks.

We are interested in strict quantized neural networks. While it is convenient to treat the entries of $\{A_\ell\}_{\ell=1}^{L}$ and $\{\boldsymbol{b}_\ell\}_{\ell=1}^{L}$ as real numbers for theoretical analysis, when used for computations, each entry is traditionally stored as a 32-bit float in memory. To reduce the number of bits, we consider a fixed finite $\mathcal{A} \subseteq \mathbb{R}$ called the alphabet.

**Definition 1.2.** A strict $\mathcal{A}$-quantized neural network is a strict neural network where all nonzero entries of the weight matrices $\{A_\ell\}_{\ell=1}^{L}$ and bias vectors $\{\boldsymbol{b}_\ell\}_{\ell=1}^{L}$ belong to $\mathcal{A}$.

An alphabet $\mathcal{A}$ that only consists of a small set of allowable values, such as the extreme one-bit case, is theoretically interesting as it poses stringent constraints and is computationally relevant as one can take advantage of special hardware and software for one-bit floating point operations. Since the weights and biases are selected from the same alphabet $\mathcal{A}$, this definition imposes a particular scaling on the associated network, which can be altered by dilating the domain.

## 1.3 The challenge: approximation by coarsely quantized networks

The problem of neural network quantization can be studied from an approximation theory perspective. Let $\mathcal{N} := \mathcal{N}(\mathcal{A}, \beta, L, N, P)$ be the set of functions that can be expressed as a strict $\mathcal{A}$-quantized neural network with activation $\beta$ and has at most $L$ layers, $N$ nodes, and $P$ parameters. For a prescribed function class $\mathcal{F}$, and distortion measure $\mathcal{E} \colon \mathcal{F} \times \mathcal{N} \to \mathbb{R}$, we study the

$$\text{approximation error} \quad := \quad \sup_{f \in \mathcal{F}} \ \inf_{g \in \mathcal{N}} \ \mathcal{E}(f, g).$$

The approximation error achieved by generalized neural networks has been extensively studied in the traditional case where there is no quantization. Well known classical universal approximation theorems [10, 3] are qualitative statements for shallow networks. Modern versions [39, 34, 30, 12, 19] provide quantitative approximation rates in terms of the function class, number of layers, parameters, etc. This is only a partial list of references, and additional ones can be found in the bibliography of a recent comprehensive survey [13]. Perhaps this is a suitable place to mention that the set of strict neural networks is a subset of their generalized counterparts, so any function class that can be approximated by strict quantized neural networks can also be well-approximated by generalized unquantized ones.

It is natural to wonder if these results or their proof strategies can be adapted to strict quantized networks. If $\mathcal{A}$ is of sufficiently high resolution and covers a wide range of numbers, such as $\mathcal{A} = \delta \mathbb{Z} \cap [-M, M]$, for sufficiently large $M$ and small $\delta > 0$, then the aforementioned approximation results for unquantized generalized networks can be suitably adapted, such as in [6, Lemma 3.7]. However, this approach requires using an increasingly higher resolution alphabet (reducing $\delta$) to achieve smaller errors. Major difficulties come into play once we fix an alphabet $\mathcal{A}$ and a function class $\mathcal{F}$, and ask to approximate any $f \in \mathcal{F}$ up to any error $\varepsilon$ by a $\mathcal{A}$-quantized strict neural network.

Many of the aforementioned papers employ the following ubiquitous strategy. For a prescribed $f \in \mathcal{F}$, we approximate $f$ by a particular linear combination, $\sum_{k=1}^{N} a_k \phi_k$, where the coefficients $\{a_k\}_{k=1}^{N}$ depend on $f$ and the span of $\{\phi_k\}_{k=1}^{\infty}$ is dense in $\mathcal{F}$. Examples of $\{\phi_k\}_{k=1}^{\infty}$ include local polynomials, ridge functions, and wavelets. While it is possible that $\phi_k$ is not implementable as a generalized neural network, it is enough to find a $\psi_k$ that is a close approximation of $\phi_k$ which is implementable. After $\{\psi_k\}_{k=1}^{N}$ are implemented, the summation $\sum_{k=1}^{N} a_k \psi_k$ is carried out by a linear last layer whose weights are $\{a_k\}_{k=1}^{N}$.

However, this approximation and implementation strategy becomes problematic when $\mathcal{A}$ is a small set, especially for the one-bit case. For example, if we were to closely follow the same strategy, we need each $\psi_k$ to be implementable by a $\mathcal{A}$-quantized strict neural network, and require $a_k \in \mathcal{A}$ to perform the linear combination $\sum_{k=1}^{N} a_k \psi_k$. From this point of view, it is natural to desire the following three properties:

$(P_1)$ Approximation. Given a large function class $\mathcal{F}$, finite linear combinations of $\{\phi_k\}_{k=1}^{\infty}$ with real coefficients can efficiently approximate any $f \in \mathcal{F}$.

$(P_2)$ Quantization. For any $f \in \mathcal{F}$ and its approximation $\sum_{k=1}^{N} a_k \phi_k$, the real coefficients $\{a_k\}_{k=1}^{N}$ can be replaced with suitable ones from just $\mathcal{A}$ without incurring to much additional error.

$(P_3)$ Implementation. For each $\phi_k$, there is a good approximant $\psi_k$ that can be implemented by a strict $\mathcal{A}$-quantized neural network.

While there many satisfactory choices for $(P_1)$, it is not immediate if any of those can be made compatible with $(P_2)$ and $(P_3)$. On the other hand, there are a few known choices of $\{\phi_k\}_{k=1}^{\infty}$ that satisfy $(P_2)$. Some are for very restrictive function classes: when $f$ is bandlimited and $\{\phi_k\}_{k=1}^{\infty}$

contain shifts of a sinc-kernel [11, 20] or when $f$ is a power series of a single complex variable and $\{\phi_k\}_{k=1}^{\infty}$ is the standard polynomial power basis [21].

The approach taken in this paper builds upon our recent publication [22], and there, we showed that any continuous function on $[0, 1]$ can be approximated by a $\pm 1$ linear combination of Bernstein polynomials. A Bernstein polynomial of order $n$ and index $\boldsymbol{k} = (k_1, \ldots, k_d)$ with $0 \leq k_\ell \leq n$ is the function $p_{n,\boldsymbol{k}} \colon \mathbb{R}^d \to [0, 1]$ defined as

$$p_{n,\boldsymbol{k}}(\boldsymbol{x}) = \binom{n}{\boldsymbol{k}} \boldsymbol{x}^{\boldsymbol{k}} (1 - \boldsymbol{x})^{n-\boldsymbol{k}} = \prod_{\ell=1}^{d} \binom{n}{k_\ell} x_\ell^{k_\ell} (1 - x_\ell)^{n-k_\ell}.$$

Continuing with this line of research, we investigate the multivariate Bernstein system as a potential candidate that satisfies all three of our desired properties.

## 1.4 Main contributions

The main theorems are proved by decomposing the total approximation error into three Bernstein related terms,

$$f - f_{NN} \quad = \quad \underbrace{f - f_B}_{\text{Bern. approx. error}} \quad + \quad \underbrace{f_B - f_Q}_{\text{Bern. quan. error}} \quad + \quad \underbrace{f_Q - f_{NN}}_{\text{Bern. implementation error}} \quad .$$

Here, $f_B$ is a linear combination of multivariate Bernstein polynomials whose coefficients are real and appropriately bounded, $f_Q$ is a $\pm a$ linear combination of multivariate Bernstein polynomials for appropriate $a$, and $f_{NN}$ is a function implementable by a strict $\{\pm a\}$-quantized neural network.

Our first main theorem concerns $f - f_Q = (f - f_B) + (f_B - f_Q)$. It shows that any smooth multivariate $f$ can be approximated by a $\pm 1$ linear combination of Bernstein polynomials with a quantitative rate that exploits smoothness of $f$. In the following, $\| \cdot \|_{C^s}$ is a norm on the space of $s$-times continuously differentiable functions and $\| \cdot \|_{C^1 \text{Lip}}$ is a norm on the space of continuously differentiable functions whose first order partial derivatives are Lipschitz.

**Theorem A.** Let $s, d, n \geq 1$, $\mu \in (0, 1)$, and $f \in C^s([0, 1]^d)$ with $\|f\|_\infty \leq \mu$. If $s \geq 3$, also assume that $n \geq \dfrac{\sqrt{2^{s+1}} d}{4(1 - \mu)} \|f\|_{C^1 \text{Lip}}$. Then for any $1 \leq \ell \leq d$, there exists a sequence $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ such that $\sigma_{\boldsymbol{k}} \in \{\pm 1\}$ for each $\boldsymbol{k}$ and for all $\boldsymbol{x} \in [0, 1]^d$,

$$\left| f(\boldsymbol{x}) - \sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| \lesssim_{s,d,\mu} \|f\|_{C^s} \min \left( 1, n^{-s/2} x_\ell^{-s} (1 - x_\ell)^{-s} \right).$$

Theorem A is proved constructively by first approximating $f$ with a suitable $f_B$ of the form $\sum_{0 \leq \boldsymbol{k} \leq n} a_{\boldsymbol{k}} p_{n,\boldsymbol{k}}$. Then secondly, the coefficients $\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ are fed into an algorithm called $\Sigma\Delta$ quantization to produce the desired one-bit sequence $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ from which we get $f_Q$ of the form $\sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}$. The approximant $f_Q$ is a $\pm 1$ linear combination of Bernstein polynomials and can be numerically computed, provided that we have point samples of $f$ on a sufficiently dense grid, as summarized in Algorithm 1. The $\Sigma\Delta$ algorithm falls under a general class of "noise-shaping" methods, where the main idea is to compute the signs $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ so that the "noise" $\{a_{\boldsymbol{k}} - \sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$, when fed into the synthesis operator $c \mapsto \sum_{0 \leq \boldsymbol{k} \leq n} c_{\boldsymbol{k}} p_{n,\boldsymbol{k}}$, is small in a suitable sense.

Our second main result deals with the implementation error $f_Q - f_{NN}$. It shows that any one-bit linear combination of Bernstein polynomials is implementable by strict one-bit neural networks, with either the quadratic $\rho(t) = \frac{1}{2} t^2$ or ReLU $\sigma(t) = \max(t, 0)$ activation. The proof is constructive and schematic diagrams for the constructed networks are shown in Figures 1a and 1b.
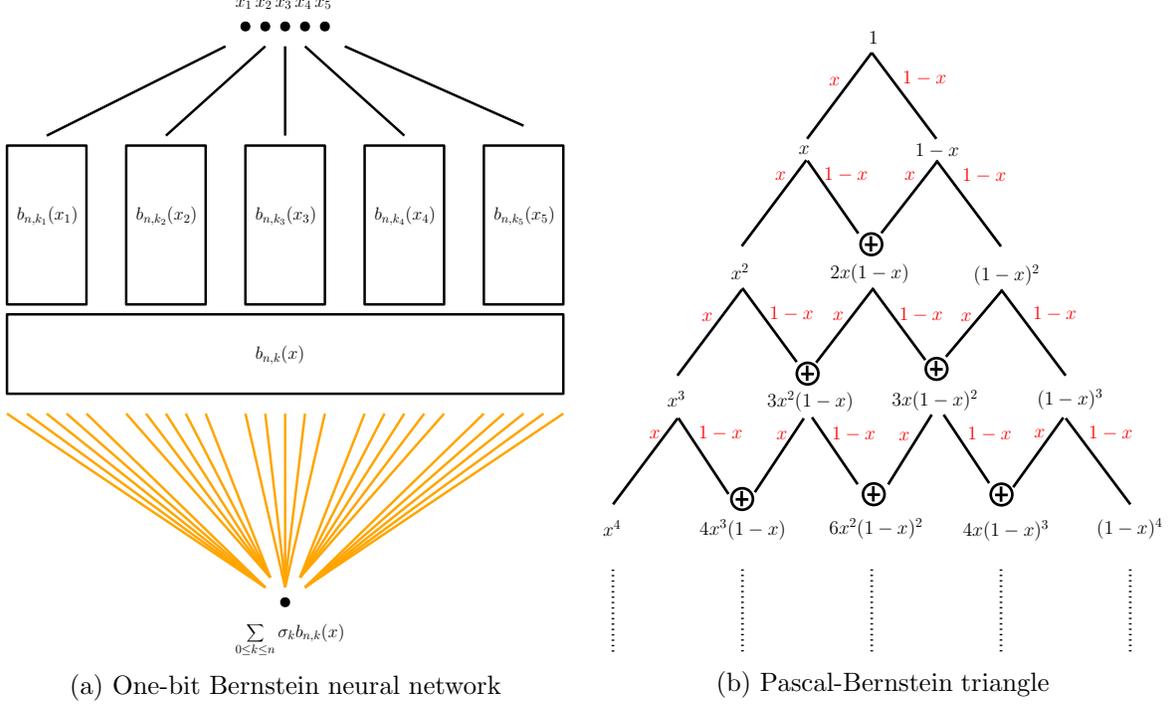
4

(a) One-bit Bernstein neural network      (b) Pascal-Bernstein triangle

Figure 1: Schematic diagrams of the constructed one-bit neural networks.

**Theorem B.** For any integers $d, n \geq 1$ and function $f = \sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}$ where $\sigma_{\boldsymbol{k}} \in \{\pm 1\}$ for each $\boldsymbol{k}$, the following hold.

- There is a $\{\pm 1\}$-quantized quadratic neural network $f_{NN,\rho}$ that has $O(n)$ layers and $O(n^d)$ nodes and parameters, as $n \to \infty$, such that $f_{NN,\rho} = f$.

- For each $\varepsilon$, there exists a $\{\pm\frac{1}{2}\}$-quantized ReLU neural network $f_{NN,\sigma}$ with $O(n \log(n/\varepsilon))$ layers and $O((n^2 + n^d) \log(n/\varepsilon))$ nodes and parameters, as $n \to \infty$ and $\varepsilon \to 0$, such that $\|f - f_{NN,\sigma}\|_\infty \leq \varepsilon$.

Combining the previous theorems, we obtain the final main theorem of this paper.

**Theorem C.** Let $s, d, n \geq 1$, $\mu \in (0, 1)$, and $f \in C^s([0, 1]^d)$ with $\|f\|_\infty \leq \mu$. If $s \geq 3$, also assume that $n \geq \dfrac{\sqrt{2^{s+1}} d}{4(1 - \mu)} \|f\|_{C^1 \mathrm{Lip}}$. For any $1 \leq \ell \leq d$, there exist functions $f_{NN,\rho}$ and $f_{NN,\sigma}$ such that:

- for $f_{NN} \in \{f_{NN,\rho}, f_{NN,\sigma}\}$ and any $\boldsymbol{x} \in [0, 1]^d$,

$$|f(\boldsymbol{x}) - f_{NN}(\boldsymbol{x})| \lesssim_{s,d,\mu} \|f\|_{C^s} \min\left(1, n^{-s/2} x_\ell^{-s} (1 - x_\ell)^{-s}\right).$$

- $f_{NN,\rho}$ is implementable by a strict $\{\pm 1\}$-quantized quadratic neural network that has $O(n)$ layers and $O(n^d)$ nodes and parameters, as $n \to \infty$.

- $f_{NN,\sigma}$ is implementable by a strict $\{\pm\frac{1}{2}\}$-quantized ReLU neural network that has $O(n \log n)$ layers and $O((n^2 + n^d) \log n)$ nodes and parameters, as $n \to \infty$.

5

## 1.5 Why Bernstein? Other contributions

In this paper, we show that the multivariate Bernstein polynomials $\{p_{n,\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n}$ of order $n$ satisfy the three important properties $(P_1)$, $(P_2)$, and $(P_3)$, which are then used to prove the main theorems.

$(P_1)$ Approximation. In Theorem 2.4, we use iterated Bernstein operators to approximate any $f \in C^s([0,1]^d)$, with a rate of approximation that exploits smoothness of $f$, which generalizes the one-dimensional results of Micchelli [32] and Felbecker [16]. This allows us to avoid a known saturation result [14, Chapter 10, Theorem 3.1], which says that the usual Bernstein operator is unable to exploit higher order derivatives of the target function. In Theorem 2.6, we convert the resulting iterated approximant into a linear combination of Bernstein polynomials without much amplification in the resulting coefficients.

$(P_2)$ Quantization. In Theorem 3.5, we show that any linear combination of Bernstein polynomials with real coefficients can be replaced with a suitable sequence of $\pm 1$ linear combination, without significant error. This step is done constructively through a directional $\Sigma\Delta$ algorithm. This is perhaps surprising because previous applications of noise-shaping quantization utilize some notion of redundancy in the generating system, whereas the Bernstein polynomials are linearly independent. One explanation is that the Bernstein system of order $n$ and of a single variable span a subspace whose numerical rank is approximately $\sqrt{n}$.

$(P_3)$ Implementation. For both the ReLU and quadratic activation functions, our implementation strategy exploits a natural Pascal triangle interpretation of the univariate Bernstein polynomials, as shown in Figure 1b. This connection is vital in being able to implement the Bernstein polynomials in a stable and efficient way. Since we exclusively employ strict neural networks, our constructions are more constrained than those found in papers that employ generalized neural networks. The constructions are found in A.

## 1.6 Additional related work

This paper addresses the universal approximation capabilities of coarsely quantized neural networks: whether they can arbitrarily well-approximate large function classes. Related work on approximation by unquantized neural networks or with variable high resolution alphabets were discussed in Section 1.3. There we also explain why the problem of universal approximation with coarsely quantized networks is significantly different and requires novel technical developments.

On the other hand, several algorithms for neural network quantization have been developed and their performances have been evaluated empirically, see the survey article [25]. In essence, existing algorithms take a pre-trained network and replace each layer with a quantized approximant, and/or directly train the network by quantizing the back-propagation vectors. While it has been empirically observed that they can compress networks without sacrificing substantial accuracy, there is little theory explaining why. Several recent works [2, 31] address this gap by proposing gradient-based quantization methods with provable guarantees.

This paper's material builds upon our previous work [23]. There we analyzed the approximation properties of one-bit linear combinations of univariate Bernstein polynomials, and Theorem C is a multivariate generalization of [23, Theorem 7]. While [23, Appendix B] briefly touches upon implementation by quadratic neural networks, it only pertained to a single dimension and did not provide quantitative bounds for the size of such networks. On the other hand, the results in this paper deal with arbitrary dimensions, and for both the quadratic and ReLU activation functions.

This version improves upon the first draft of this paper in two ways: we implement the approximation strategy given in Theorem A with strict neural networks, and the ReLU case only necessitates a one-bit alphabet.

## 1.7 Organization

The organization of subsequent sections follow the ordering $(P_1)$, $(P_2)$, and $(P_3)$. The Bernstein approximation error $f - f_B$ is treated in Section 2. The quantization error in the Bernstein basis $f_B - f_Q$ is dealt with in Section 3. The implementation error of these one-bit approximations $f_Q - f_{NN}$ for both the quadratic and ReLU activation functions is studied in Section 4, while the constructions of $f_{NN}$ are carried out in Appendix A. Proofs of Theorems A, B, and C are provided in Sections 3.4, 4.3, and 4.4, respectively. Final remarks and other aspects of this paper, including the computational method and entropy considerations, are contained in Section 5.

## 1.8 Basic notation

We let $\mathbb{R}$ be the reals and $\mathbb{N}$ be the natural numbers including zero. For reasons that will become evident, we let log denote the base 2 logarithm. The ceiling and floor functions are denoted $\lceil \cdot \rceil$ and $\lfloor \cdot \rfloor$, respectively. We use the notation $A \lesssim B$ to mean that there is a universal constant $C$ such that $A \leq CB$. When we write $A \lesssim_{a,b,c} B$, we mean that there is a $C > 0$ that possibly depends on $a, b, c$ for which $A \leq CB$.

Let $\boldsymbol{e}_1, \ldots, \boldsymbol{e}_d$ denote the canonical orthonormal basis for $\mathbb{R}^d$. For any $p \in [1, \infty]$, let $\| \cdot \|_p$ be the usual $p$-th norm on vectors in $\mathbb{R}^d$. We also use the same notation for the $\ell^p$-norm of a function defined on a countable set, and the $L^p$ norm of a function.

We denote multi-indices and tuples with boldface letters. For $\boldsymbol{k}, \boldsymbol{\ell} \in \mathbb{N}^d$ and $m \in \mathbb{N}$, we define the following. We write $\boldsymbol{k} \leq \boldsymbol{\ell}$ as shorthand for $k_j \leq \ell_j$ for each $j$. Likewise, we let $\boldsymbol{k} \leq m$ (or $m \leq \boldsymbol{k}$) to mean that $k_j \leq m$ (or $m \leq k_j$) for each $j$. For any $\boldsymbol{x} \in \mathbb{R}^d$, let $\boldsymbol{x}^{\boldsymbol{k}} := x_1^{k_1} \cdots x_d^{k_d}$, and $|\boldsymbol{k}| = k_1 + \cdots + k_d$. The degree of $\boldsymbol{x}^{\boldsymbol{k}}$ is defined to be $|\boldsymbol{k}|$. For any $t \in \mathbb{R}$, we let $t\boldsymbol{k} = (tk_1, \ldots, tk_d)$. We follow standard convention for dealing with combinatorial factors:

$$\binom{\boldsymbol{k}}{\boldsymbol{\ell}} := \prod_{j=1}^{d} \binom{k_j}{\ell_j}, \quad \text{and} \quad \binom{m}{\boldsymbol{k}} := \prod_{j=1}^{d} \binom{m}{k_j}.$$

## 2 Approximation error

In this section, we concentrate on the error incurred by approximating a smooth $f$ with $f_B = \sum_{0 \leq \boldsymbol{k} \leq n} a_{\boldsymbol{k}} p_{n,\boldsymbol{k}}$, a linear combination of Bernstein polynomials of order $n$, where each coefficient $a_{\boldsymbol{k}}$ can be suitably controlled in terms of $f$. For reasons that will become apparent later on when we discuss the quantization error, we will construct a $f_B$ whose coefficients in the Bernstein basis are not much larger than $\|f\|_\infty$. More specifically, the results of this section will provide us with upper bounds on the approximation error by Bernstein polynomials,

$$\mathcal{E}_{n,c}^{\text{approx}}(f) := \inf_{\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}, |a_{\boldsymbol{k}}| \leq c} \left\| f - \sum_{0 \leq \boldsymbol{k} \leq n} a_{\boldsymbol{k}} p_{n,\boldsymbol{k}} \right\|_\infty.$$

A famous theorem of Bernstein showed that $B_n(f)$ converges uniformly to $f$ on $[0,1]$ provided that $f$ is continuous, which generalizes to higher dimensions, see [27] for the two-dimensional case. While $B_n(f)$ may appear to be a natural candidate for $f_B$, there is a well known saturation result

which says that, even for infinitely differentiable $f$, the fastest rate of decay is $1/n$, see [14, Chapter 10, Theorem 3.1].

To circumvent this saturation, we follow an approach of Micchelli [32], and show in Theorem 2.4, that appropriate iterates of multivariate Bernstein operators achieve improved convergence rates that exploit smoothness of $f$. Theorem 2.5 allows us to convert these iterated approximations to a linear combination of Bernstein polynomials and Theorem 2.6 provides us with the existence of a suitable $f_B$.

## 2.1   Background on Bernstein polynomials

For a fixed integer $n \geq 1$, we denote the set of univariate Bernstein polynomials by $\mathcal{B}_n := \{p_{n,k}\}_{k=0}^n$. Each $p_{n,k} \colon [0,1] \to [0,1]$ is a polynomial of degree $n$ and

$$p_{n,k}(x) := \binom{n}{k} x^k (1-x)^{n-k}.$$

The Bernstein polynomials are nonnegative, form a partition of unity for $[0,1]$, and form a basis for the vector space of $n$ degree algebraic polynomials. They can be used to give a constructive proof of the classical Weierstrass theorem, by showing that any continuous $f$ on $[0,1]$ can be uniformly approximated by the Bernstein polynomial of $f$,

$$B_n(f)(x) := \sum_{k=0}^n f\left(\frac{k}{n}\right) p_{n,k}(x).$$

For reasons that will be apparent later, it will be convenient for us to extend the index $k$ beyond $n$ in the definition of $p_{n,k}$. If $k < 0$ or $k > n$, then we define $p_{n,k} := 0$.

The Bernstein polynomials can also be viewed from a probabilistic perspective. For any $x \in [0,1]$, the quantity $p_{n,k}(x)$ is the probability that $k$ successes occur in $n$ independent Bernoulli trials each with probability of success $x$. Since the expected value is $nx$, for each integer $s \geq 0$, the $s$-th central moment is

$$T_{n,s}(x) := \sum_{k=0}^n (k - nx)^s p_{n,k}(x).$$

It is known that $T_{n,s}$ is a polynomial in $x$ of degree at most $s$ and in $n$ of degree at most $\lfloor s/2 \rfloor$. We will employ a few specific formulas for small $s$. For each $x \in [0,1]$, with the short hand notation $X := x(1-x)$, we have

$$
\begin{aligned}
&T_{n,0}(x) = 1, \quad T_{n,1}(x) = 0, \quad T_{n,2}(x) = nX, \\
&T_{n,3}(x) = n(1-2x)X, \quad T_{n,4}(x) = 3n^2 X^2 + n(X - 6X^2).
\end{aligned}
\tag{2.1}
$$

For each integer $s \geq 0$, there is a constant $A_s$ such that for all $n \geq 1$ and $x \in [0,1]$, we have

$$0 \leq T_{n,2s}(x) \leq A_s n^s. \tag{2.2}$$

See [14, Chapter 10] for proofs of the above results. Although we will not need explicit upper bounds for $A_s$, some can be found in [1, 33].

For dimension $d \geq 1$ and integer $n \geq 1$, Bernstein polynomials are defined to be tensor products of single variable Bernstein polynomials.

**Definition 2.1.** The multivariate Bernstein polynomials of order $n$ is $\mathcal{B}_n := \{p_{n,\boldsymbol{k}}\}_{0 \le \boldsymbol{k} \le n}$, where each $p_{n,\boldsymbol{k}} \colon [0,1]^d \to [0,1]$ is defined as,

$$p_{n,\boldsymbol{k}}(\boldsymbol{x}) := p_{n,k_1}(x_1) \cdots p_{n,k_d}(x_d) = \binom{n}{\boldsymbol{k}} \boldsymbol{x}^{\boldsymbol{k}} (1 - \boldsymbol{x})^{n-\boldsymbol{k}}.$$

Whenever there is a $1 \le \ell \le d$ such that either $k_\ell < 0$ or $k_\ell > n$, then we define $p_{n,\boldsymbol{k}} = 0$. It is important to mention that we exclusively use multivariate Bernstein polynomials that are formed as tensor products. They are significantly different from Bernstein polynomials on the canonical $d$-dimensional simplex.

Several properties of multivariate Bernstein polynomials can be deduced by exploiting their tensor product structure. They form a partition of unity for $[0,1]^d$, and in particular,

$$\sum_{k_\ell=0}^{n} p_{n,k_\ell}(x_\ell) = 1 \quad \text{for all} \quad 1 \le \ell \le d \quad \text{and} \quad \boldsymbol{x} \in [0,1]^d. \tag{2.3}$$

The Bernstein polynomial of a multivariate $f$ is defined similar to before,

$$B_n(f)(\boldsymbol{x}) := \sum_{0 \le \boldsymbol{k} \le n} f\left(\frac{\boldsymbol{k}}{n}\right) p_{n,\boldsymbol{k}}(\boldsymbol{x}) = \sum_{k_1=0}^{n} \cdots \sum_{k_d=0}^{n} f\left(\frac{k_1}{n}, \cdots, \frac{k_d}{n}\right) p_{n,k_1}(x_1) \cdots p_{n,k_d}(x_d).$$

Central moments of multivariate Bernstein polynomials can be readily extracted. For any $\boldsymbol{\alpha} \in \mathbb{N}^d$, we define

$$T_{n,\boldsymbol{\alpha}}(\boldsymbol{x}) := \sum_{0 \le \boldsymbol{k} \le n} (\boldsymbol{k} - n\boldsymbol{x})^{\boldsymbol{\alpha}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) = T_{n,\alpha_1}(x_1) \cdots T_{n,\alpha_d}(x_d).$$

Letting $A_{\boldsymbol{\alpha}} := A_{\alpha_1} \cdots A_{\alpha_d}$, it follows from (2.2) that

$$0 \le T_{n,2\boldsymbol{\alpha}}(\boldsymbol{x}) \le A_{\boldsymbol{\alpha}} n^{|\boldsymbol{\alpha}|}. \tag{2.4}$$

## 2.2 Approximation by the Bernstein operator

We start by investigating the approximation properties of the Bernstein operator. We say a function $f \colon [0,1]^d \to \mathbb{R}$ is Lipschitz continuous if there is a $L \ge 0$ for which $|f(\boldsymbol{x}) - f(\boldsymbol{y})| \le L\|\boldsymbol{x} - \boldsymbol{y}\|_2$ for all $\boldsymbol{x}, \boldsymbol{y} \in [0,1]^d$. We let $|f|_{\text{Lip}}$ be the smallest such $L$ for which this inequality holds. The following is an elementary observation regarding the approximation of Lipschitz functions by the Bernstein operator, see also [26].

**Proposition 2.2.** *For any Lipschitz continuous $f$ on $[0,1]^d$, we have*

$$\|f - B_n(f)\|_\infty \le \frac{|f|_{\text{Lip}}}{2} \sqrt{\frac{d}{n}}.$$

*Proof.* Fix any $\boldsymbol{x} \in [0,1]^d$. By the partition of unity property (2.3), we have

$$\left| f(\boldsymbol{x}) - \sum_{0 \le \boldsymbol{k} \le n} f\left(\frac{\boldsymbol{k}}{n}\right) p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| = \left| \sum_{0 \le \boldsymbol{k} \le n} \left( f(\boldsymbol{x}) - f\left(\frac{\boldsymbol{k}}{n}\right) \right) p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right|.$$

We use the Lipschitz condition on $f$ to see that

$$\left| \sum_{0 \le \boldsymbol{k} \le n} \left( f(\boldsymbol{x}) - f\left(\frac{\boldsymbol{k}}{n}\right) \right) p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| \le \frac{|f|_{\text{Lip}}}{n} \sum_{0 \le \boldsymbol{k} \le n} \|n\boldsymbol{x} - \boldsymbol{k}\|_2 \, p_{n,\boldsymbol{k}}(\boldsymbol{x}).$$

9

By Cauchy-Schwarz, partition of unity property (2.3), and moments identity (2.1), we have

$$\sum_{0 \le \boldsymbol{k} \le n} \|n\boldsymbol{x} - \boldsymbol{k}\|_2 \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \le \Big( \sum_{0 \le \boldsymbol{k} \le n} \|nx - \boldsymbol{k}\|_2^2 \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \Big)^{1/2} \Big( \sum_{0 \le \boldsymbol{k} \le n} p_{n,\boldsymbol{k}}(\boldsymbol{x}) \Big)^{1/2}$$

$$= \Big( \sum_{\ell=1}^{d} \sum_{0 \le \boldsymbol{k} \le n} (nx_\ell - k_\ell)^2 \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \Big)^{1/2} = \Big( \sum_{\ell=1}^{d} T_{n,2}(x_\ell) \Big)^{1/2} \le \frac{\sqrt{nd}}{2}.$$

$\square$

For the one dimensional case, if the target function is twice differentiable, then it is possible to obtain a faster rate of decay in $n$, see [14, Chapter 10]. We proceed to derive an analogous result for the multivariate case. Before we state it, for concreteness, let us introduce some notation.

For a $\boldsymbol{\alpha} \in \mathbb{N}^d$, we use the shorthand notation $\partial_{\boldsymbol{x}}^{\boldsymbol{\alpha}} := \partial^{\boldsymbol{\alpha}} := \partial_{x_1}^{\alpha_1} \cdots \partial_{x_d}^{\alpha_d}$. For any integer $s \ge 1$, we let $C^s([0,1]^d)$ be the set of continuous real functions $f$ defined on a $[0,1]^d$ such that $\partial^{\boldsymbol{\alpha}} f$ is continuous for all $|\boldsymbol{\alpha}| \le s$. We define the semi-norm $|f|_{\dot{C}^s} := \sup_{|\boldsymbol{\alpha}|=s} \|\partial^{\boldsymbol{\alpha}} f\|_\infty$ and equip $C^s(U)$ with the norm,

$$\|f\|_{C^s} = \max \Big\{ \|f\|_\infty, \max_{1 \le k \le s} |f|_{\dot{C}^k} \Big\}.$$

We also define $C^s\mathrm{Lip}$ to be the set of functions that are $s$ times differentiable and whose $s$-th order partial derivatives are Lipschitz continuous, and we equip $C^s\mathrm{Lip}$ with the norm

$$\|f\|_{C^s\mathrm{Lip}} = \max \Big\{ \|f\|_{C^s}, \max_{|\boldsymbol{\alpha}|=s} |\partial^{\boldsymbol{\alpha}} f|_{\mathrm{Lip}} \Big\}.$$

**Proposition 2.3.** *For any $f \in C^1\mathrm{Lip}([0,1]^d)$, we have*

$$\|f - B_n(f)\|_\infty \le \frac{d}{8n} \|f\|_{C^1\mathrm{Lip}}.$$

*Proof.* Fix any $\boldsymbol{x} \in [0,1]^d$. For each $0 \le \boldsymbol{k} \le n$, by Taylor's theorem,

$$f\Big(\frac{\boldsymbol{k}}{n}\Big) = f(\boldsymbol{x}) + \sum_{\ell=1}^{d} \partial_\ell f(\boldsymbol{x}) \Big(\frac{k_\ell}{n} - x_\ell\Big)$$

$$+ \sum_{\ell=1}^{d} \Big(\frac{k_\ell}{n} - x_\ell\Big) \int_0^1 \Big[ \partial_\ell f\Big(\boldsymbol{x} + t\Big(\frac{k_\ell}{n} - x_\ell\Big) \boldsymbol{e}_\ell\Big) - \partial_\ell f(\boldsymbol{x}) \Big] \, dt.$$

This equation, the partition of unity property (2.3), and the moments identity (2.1) (in particular that $T_{n,1}(x_\ell) = 0$ for each $\ell$), we deduce

$$B_n(f)(\boldsymbol{x}) - f(\boldsymbol{x}) = \sum_{0 \le \boldsymbol{k} \le n} \Big( f\Big(\frac{\boldsymbol{k}}{n}\Big) - f(\boldsymbol{x}) \Big) p_{n,\boldsymbol{k}}(\boldsymbol{x})$$

$$= \sum_{0 \le \boldsymbol{k} \le n} \sum_{\ell=1}^{d} \Big(\frac{k_\ell}{n} - x_\ell\Big) p_{n,\boldsymbol{k}}(x) \int_0^1 \Big[ \partial_\ell f\Big(\boldsymbol{x} + t\Big(\frac{k_\ell}{n} - x_\ell\Big) \boldsymbol{e}_\ell\Big) - \partial_\ell f(\boldsymbol{x}) \Big] \, dt.$$

From this equation, that the Bernstein polynomials are nonnegative, the Lipschitz assumption on the partial derivatives of $f$, and partition of unity (2.3), we see that

$$|f(\boldsymbol{x}) - B_n(f)(\boldsymbol{x})| \le \frac{\|f\|_{C^1\mathrm{Lip}}}{2n^2} \sum_{0 \le \boldsymbol{k} \le n} \sum_{\ell=1}^{d} (k_\ell - nx_\ell)^2 p_{n,\boldsymbol{k}}(\boldsymbol{x}) = \frac{\|f\|_{C^1\mathrm{Lip}}}{2n^2} \sum_{\ell=1}^{d} T_{n,2}(x_\ell).$$

Applying the moments identity (2.1) completes the proof. $\square$

10

There are other properties of Bernstein operators, such as convergence of partial derivatives [17] that we will not use in this paper.

## 2.3   Approximation of smooth functions with iterated Bernstein

As mentioned earlier, due to the saturation phenomenon for the Bernstein operator, there is no hope of improving the decay rate of $1/n$ in Proposition 2.3 even under additional regularity assumptions on the target function. To overcome this, iterated univariate Bernstein operators were developed by Micchelli [32] and Felbecker [16] as alternative means of approximation. By viewing $B_n$ as an operator on $C([0,1]^d)$, for any integer $r \geq 1$, we define the following iterated Bernstein operator

$$U_{n,r} := I - (I - B_n)^r = \sum_{j=1}^{r} (-1)^{j-1} \binom{r}{j} B_n^j.$$

Note that $U_{n,1} = B_n$ coincides with the usual Bernstein operator. The following theorem generalizes the core results of Micchelli-Felbecker to higher dimensions.

**Theorem 2.4.** *For any integers $s, d \geq 1$ and any $f \in C^s([0,1]^d)$, it holds that*

$$\|f - U_{n,\lceil s/2 \rceil}(f)\|_\infty \lesssim_{s,d} \|f\|_{C^s} n^{-s/2}.$$

*Proof.* It suffices to prove that, for all $s \geq 1$ and $f \in C^s([0,1]^d)$, we have

$$\|f - U_{n,\lceil s/2 \rceil}(f)\|_\infty = \|(I - B_n)^{\lceil s/2 \rceil}(f)\|_\infty \lesssim_{s,d} \|f\|_{C^s} n^{-s/2}. \tag{2.5}$$

We proceed by strong induction. The $s = 1, 2$ cases hold by Propositions 2.2 and 2.3. Hence, assume that inequality (2.5) holds for $s = 1, 2, \ldots, r$ for some $r \geq 2$. Now we will prove that the statement holds for $s = r + 1$.

To this end, fix a $f \in C^{r+1}([0,1]^d)$ and any $\boldsymbol{x} \in [0,1]^d$. For each $0 \leq \boldsymbol{k} \leq n$, there is a $\xi_{\boldsymbol{k},\boldsymbol{x}}$ such that

$$f\left(\frac{\boldsymbol{k}}{n}\right) = f(\boldsymbol{x}) + \sum_{0 < |\boldsymbol{\alpha}| \leq r} \frac{\partial^{\boldsymbol{\alpha}} f(\boldsymbol{x})}{\boldsymbol{\alpha}!} \left(\frac{\boldsymbol{k}}{n} - \boldsymbol{x}\right)^{\boldsymbol{\alpha}} + \sum_{|\boldsymbol{\alpha}| = r+1} \frac{\partial^{\boldsymbol{\alpha}} f(\xi_{\boldsymbol{k},\boldsymbol{x}})}{\boldsymbol{\alpha}!} \left(\frac{\boldsymbol{k}}{n} - \boldsymbol{x}\right)^{\boldsymbol{\alpha}}.$$

From here, we see that

$$(B_n - I)(f)(\boldsymbol{x})$$
$$= \sum_{0 < |\boldsymbol{\alpha}| \leq r} \sum_{0 \leq \boldsymbol{k} \leq n} \frac{\partial^{\boldsymbol{\alpha}} f(\boldsymbol{x})}{\boldsymbol{\alpha}!} \left(\frac{\boldsymbol{k}}{n} - \boldsymbol{x}\right)^{\boldsymbol{\alpha}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) + \underbrace{\sum_{|\boldsymbol{\alpha}| = r+1} \sum_{0 \leq \boldsymbol{k} \leq n} \frac{\partial^{\boldsymbol{\alpha}} f(\xi_{\boldsymbol{k},\boldsymbol{x}})}{\boldsymbol{\alpha}!} \left(\frac{\boldsymbol{k}}{n} - \boldsymbol{x}\right)^{\boldsymbol{\alpha}} p_{n,\boldsymbol{k}}(\boldsymbol{x})}_{:=G_{n,r+1}(\boldsymbol{x})}.$$

To simplify the first term, notice that since the first central moment of Bernstein polynomials is identically zero, $T_{n,1} = 0$. Hence all terms for which $|\boldsymbol{\alpha}| = 1$ disappear. For each $2 \leq |\boldsymbol{\alpha}| \leq r$ and $2 \leq m \leq r$, we define the functions,

$$F_{n,\boldsymbol{\alpha}} := \frac{\partial^{\boldsymbol{\alpha}} f \, T_{n,\boldsymbol{\alpha}}}{n^{|\boldsymbol{\alpha}|/2} \boldsymbol{\alpha}!}, \quad \text{and} \quad F_{n,m} := \sum_{|\boldsymbol{\alpha}| = m} F_{n,\boldsymbol{\alpha}}.$$

It follows that

$$(B_n - I)(f)(\boldsymbol{x}) = \sum_{2 \leq |\boldsymbol{\alpha}| \leq r} \frac{F_{n,\boldsymbol{\alpha}}(\boldsymbol{x})}{n^{|\boldsymbol{\alpha}|/2}} + G_{n,r+1}(\boldsymbol{x}) = \sum_{m=2}^{r} \frac{F_{n,m}(\boldsymbol{x})}{n^{m/2}} + G_{n,r+1}(\boldsymbol{x}).$$

11

This identity holds for each $\boldsymbol{x}$. Noting that $\lceil (r+1)/2 \rceil \geq 1$ since $r \geq 2$, We have

$$(B_n - I)^{\lceil (r+1)/2 \rceil}(f) = \sum_{m=2}^{r} \frac{(B_n - I)^{\lceil (r+1)/2 \rceil - 1} F_{n,m}}{n^{m/2}} + (B_n - I)^{\lceil (r+1)/2 \rceil - 1} G_{n,r+1}. \tag{2.6}$$

We concentrate on the primary term in (2.6) first. To apply the inductive hypothesis, we claim that for each $2 \leq m \leq r$, we have $F_{n,m} \in C^{r+1-m}$ and

$$\|F_{n,m}\|_{C^{r+1-m}} \leq \|f\|_{C^{r+1}} 2^{r+1-m} \sum_{|\boldsymbol{\alpha}|=m} \sqrt{A_{\boldsymbol{\alpha}}}. \tag{2.7}$$

The key part of this assertion is that the upper bound for $\|F_{n,m}\|_{C^{r+1-m}}$ does not depend on $n$. Note that for each $|\boldsymbol{\alpha}| = m$, we have $\partial^{\boldsymbol{\alpha}} f \in C^{r+1-m}$ due to the initial assumption that $f \in C^{r+1}$. Also $T_{n,\boldsymbol{\alpha}}$ is infinitely differentiable since it is a multinomial, so we see that $F_{n,m} \in C^{r+1-m}$. By Leibniz, for each $|\boldsymbol{\beta}| \leq r+1-m$, we have

$$\partial^{\boldsymbol{\beta}} F_{n,m} = \sum_{|\boldsymbol{\alpha}|=m} \frac{1}{n^{m/2} \boldsymbol{\alpha}!} \sum_{0 \leq \boldsymbol{\gamma} \leq \boldsymbol{\beta}} \binom{\boldsymbol{\beta}}{\boldsymbol{\gamma}} \partial^{\boldsymbol{\alpha}+\boldsymbol{\beta}-\boldsymbol{\gamma}} f \, \partial^{\boldsymbol{\gamma}} T_{n,\boldsymbol{\alpha}}.$$

Since $T_{n,\boldsymbol{\alpha}}(\boldsymbol{x})$ is a polynomial in $x_\ell$ of degree at most $\alpha_\ell$, the inside summation can be taken over $0 \leq \boldsymbol{\gamma} \leq \min(\boldsymbol{\alpha}, \boldsymbol{\beta})$. Additionally, Bernstein's inequality for algebraic polynomials and central moment bounds (2.2) yield

$$\|\partial^{\boldsymbol{\gamma}} T_{n,\boldsymbol{\alpha}}\|_\infty \leq \prod_{\ell=1}^{d} \|\partial^{\gamma_\ell} T_{n,\alpha_\ell}\|_\infty \leq \prod_{\ell=1}^{d} \frac{\alpha_\ell!}{(\alpha_\ell - \gamma_\ell)!} \|T_{n,\alpha_\ell}\|_\infty$$

$$\leq \prod_{\ell=1}^{d} \frac{\alpha_\ell!}{(\alpha_\ell - \gamma_\ell)!} \sqrt{\|T_{n,2\alpha_\ell}\|_\infty} \leq \frac{\boldsymbol{\alpha}!}{(\boldsymbol{\alpha}-\boldsymbol{\gamma})!} \sqrt{A_{\boldsymbol{\alpha}}} \, n^{|\boldsymbol{\alpha}|/2}.$$

From this and that $|\boldsymbol{\alpha} + \boldsymbol{\beta} - \boldsymbol{\gamma}| \leq r+1$, it follows that

$$\|\partial^{\boldsymbol{\beta}} F_{n,m}\|_\infty \leq \sum_{|\boldsymbol{\alpha}|=m} \sqrt{A_{\boldsymbol{\alpha}}} \sum_{0 \leq \boldsymbol{\gamma} \leq \min(\boldsymbol{\alpha}, \boldsymbol{\beta})} \binom{\boldsymbol{\beta}}{\boldsymbol{\gamma}} \frac{1}{(\boldsymbol{\alpha}-\boldsymbol{\gamma})!} \|\partial^{\boldsymbol{\alpha}+\boldsymbol{\beta}-\boldsymbol{\gamma}} f\|_\infty$$

$$\leq \|f\|_{C^{r+1}} \sum_{|\boldsymbol{\alpha}|=m} \sqrt{A_{\boldsymbol{\alpha}}} \sum_{0 \leq \boldsymbol{\gamma} \leq \boldsymbol{\beta}} \binom{\boldsymbol{\beta}}{\boldsymbol{\gamma}} = \|f\|_{C^{r+1}} 2^{|\boldsymbol{\beta}|} \sum_{|\boldsymbol{\alpha}|=m} \sqrt{A_{\boldsymbol{\alpha}}}.$$

This completes the proof of (2.7).

Returning back to the proof at hand, notice that for each $2 \leq m \leq r$, if we define

$$q(m) := \left\lceil \frac{r+1}{2} \right\rceil - 1 - \left\lceil \frac{r+1-m}{2} \right\rceil,$$

then we have $q(m) \geq 0$ since $m \geq 2$, and $q(m) \leq (m-1)/2$. It follows from the inductive hypothesis that there exist constants $M_{r+1-m,d} > 0$ for each $2 \leq m \leq r$ such that

$$\left\| (B_n - I)^{\lceil (r+1)/2 \rceil - 1} F_{n,m} \right\|_\infty \leq \|B_n - I\|_\infty^{q(m)} \left\| (B_n - I)^{\lceil (r+1-m)/2 \rceil} F_{n,m} \right\|_\infty$$

$$\leq 2^{(m-1)/2} M_{r+1-m,d} \|F_{n,m}\|_{C^{r+1-q}} n^{-(r+1-m)/2}. \tag{2.8}$$

We next control the remainder term involving $G_{n,r+1}$. We have the following upper bound for $\|G_{n,r+1}\|_\infty$. By Cauchy-Schwarz, the partition of unity property (2.3), and central moment bounds (2.4), we see that for each $\boldsymbol{x}$,

$$
\begin{aligned}
|G_{n,r+1}(\boldsymbol{x})| &\leq \frac{|f|_{\dot{C}^{r+1}}}{n^{r+1}} \sum_{0 \leq \boldsymbol{k} \leq n} \sum_{|\boldsymbol{\alpha}|=r+1} \frac{1}{\boldsymbol{\alpha}!} \left|(\boldsymbol{k}-n\boldsymbol{x})^{\boldsymbol{\alpha}}\right| p_{n,\boldsymbol{k}}(\boldsymbol{x}) \\
&\leq \frac{|f|_{\dot{C}^{r+1}}}{n^{r+1}} \sum_{|\boldsymbol{\alpha}|=r+1} \frac{1}{\boldsymbol{\alpha}!} \left( \sum_{0 \leq \boldsymbol{k} \leq n} (\boldsymbol{k}-n\boldsymbol{x})^{2\boldsymbol{\alpha}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right)^{1/2} \leq \frac{|f|_{\dot{C}^{r+1}}}{n^{(r+1)/2}} \sum_{|\boldsymbol{\alpha}|=r+1} \frac{\sqrt{A_{\boldsymbol{\alpha}}}}{\boldsymbol{\alpha}!}.
\end{aligned} \tag{2.9}
$$

Moreover, $\|B_n - I\|_\infty \leq 2$, and $\lceil (r+1)/2 \rceil - 1 \leq r/2$. Hence we have

$$
\|(B_n - I)^{\lceil (r+1)/2 \rceil - 1} G_{n,r+1}\|_\infty \leq 2^{\lceil (r+1)/2 \rceil - 1} \|G_{n,r+1}\|_\infty \leq 2^{r/2} \|G_{n,r+1}\|_\infty. \tag{2.10}
$$

Now we are ready to complete the proof. Combining (2.6), (2.7), (2.8), (2.9), and (2.10), we see that

$$
\begin{aligned}
\|(B_n &- I)^{\lceil (r+1)/2 \rceil}(f)\|_\infty \\
&\leq \sum_{m=2}^{r} \frac{\|(B_n - I)^{\lceil (r+1)/2 \rceil - 1} F_{n,m}\|_\infty}{n^{m/2}} + \|(B_n - I)^{\lceil (r+1)/2 \rceil - 1} G_{n,r+1}\|_\infty \\
&\leq \sum_{m=2}^{r} \frac{2^{(m-1)/2} M_{r+1-m,d} \|F_{n,m}\|_{C^{r+1-m}}}{n^{(r+1)/2}} + \frac{2^{r/2} |f|_{\dot{C}^{r+1}}}{n^{(r+1)/2}} \sum_{|\boldsymbol{\alpha}|=r+1} \frac{\sqrt{A_{\boldsymbol{\alpha}}}}{\boldsymbol{\alpha}!} \\
&\lesssim_{r+1,d} \frac{\|f\|_{C^{r+1}}}{n^{(r+1)/2}}.
\end{aligned}
$$

This completes the proof by induction.

$\square$

Let us briefly discuss our result in the context of classical approximation results of Micchelli [32] and Felbecker [16]. The main distinction is that our result holds for iterates of multivariate Bernstein polynomials (formed as tensor products), while the classical papers only treat the univariate case. One small improvement we made deals the parity of $s$. Micchelli only treated the even case, which loses a $n^{-1/2}$ factor for the odd cases, whereas the odd case was satisfactorily derived by Felbecker. In our proof, we combined the even and odd cases together seamlessly. The extension of these classical results to higher dimensions was perhaps known by experts, though we were unable to find a reference.

On the other hand, there are alternative generalizations of Bernstein polynomials to the canonical simplex in $\mathbb{R}^d$, as opposed to the unit cube. Such polynomials are significantly different from the tensor product ones employed in this paper. For Bernstein polynomials on the simplex, approximation rates of Micchelli-Felbecker iterations have been studied, see [18, 15] and references therein.

## 2.4  From iterated Bernstein to linear combinations

In the next section on quantization, we will show to quantize the coefficients of a function written in the Bernstein basis. For this reason, we show how to relate the iterated Bernstein approximation $U_{n,r}(f)$ to the Bernstein polynomial $B_n(f_{n,r})$ of a possibly different function $f_{n,r}$, which can be found constructively via the formula,

$$
f_{n,r} := \left( I + \sum_{m=1}^{r-1} (I - B_n)^m \right)(f). \tag{2.11}
$$

13

**Theorem 2.5.** *For any integers $r, d \geq 1$, any $f \in C([0,1]^d)$, and any $n \geq 1$, we have*

$$U_{n,r}(f) = B_n(f_{n,r}),$$

*where $f_{n,r}$ is defined in (2.11). Further, it holds that*

$$\|f_{n,r}\|_\infty \leq \|f\|_\infty + (2^{r-1} - 1)\|f - B_n(f)\|_\infty.$$

*Proof.* [22, Theorem 5] proved this for $d = 1$ case, but the same argument extends to the Bernstein operator on $C([0,1]^d)$ without any modifications to the proof. □

This theorem shows that not only $U_{n,r}(f) = B_n(f_{n,r})$, but it also implies that the coefficients $\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ in the following theorem are not much larger than $\|f\|_\infty$. This will be important in the next section, where the coefficients will be fed into a particular quantization algorithm called $\Sigma\Delta$, which requires a $\ell^\infty$ assumption on the coefficients.

**Theorem 2.6.** *Let $d, s, n \geq 1$, $\mu, \delta \in (0,1)$, and $f \in C^s([0,1]^d)$ with $\|f\|_\infty \leq \mu$. If $s \geq 3$, also assume that $n \geq \dfrac{\sqrt{2^{s+1}}d}{8\delta}\|f\|_{C^1\mathrm{Lip}}$. There exist $\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ such that $\|a\|_\infty \leq \mu + \delta$ and*

$$\left\|f - \sum_{0 \leq \boldsymbol{k} \leq n} a_{\boldsymbol{k}} p_{n,\boldsymbol{k}}\right\|_\infty \lesssim_{s,d} \|f\|_{C^s} n^{-s/2}. \tag{2.12}$$

*Proof.* We consider different cases depending on $s$.

If $s = 1$, then we let $a_{\boldsymbol{k}} = f(\boldsymbol{k}/n)$, so $\|a\|_\infty \leq \|f\|_\infty \leq \mu$. The conclusion follows from Proposition 2.2. In this case, the implicit constant in (2.12) is $\sqrt{d}/2$.

If $s = 2$, then we let $a_{\boldsymbol{k}} = f(\boldsymbol{k}/n)$, so $\|a\|_\infty \leq \|f\|_\infty \leq \mu$. The conclusion follows from Proposition 2.3. In this case, the implicit constant in (2.12) is $d^2/8$.

Suppose $s \geq 3$. Consider the function $h := U_{n,\lceil s/2\rceil}(f)$. From Theorem 2.4, we have

$$\|f - h\|_\infty \lesssim_{s,d} \|f\|_{C^s} n^{-s/2}.$$

From Theorem 2.5 and Proposition 2.3, we have $h = B_n(f_{n,\lceil s/2\rceil})$, where $f_{n,\lceil s/2\rceil}$ is defined in (2.11), and

$$\|f_{n,\lceil s/2\rceil}\|_\infty \leq \|f\|_\infty + (2^{\lceil s/2\rceil} - 1)\|f - B_n(f)\|_\infty \leq \mu + \frac{\sqrt{2^{s+1}}d}{8n}\|f\|_{C^1\mathrm{Lip}}.$$

Pick any $n$ sufficiently large so that the right hand side is bounded above by $\mu + \delta$. We set $a_{\boldsymbol{k}} := f_{n,\lceil s/2\rceil}(\boldsymbol{k}/n)$ so that

$$h = B_n(f_{n,\lceil s/2\rceil}) = \sum_{0 \leq \boldsymbol{k} \leq n} a_{\boldsymbol{k}} p_{n,\boldsymbol{k}},$$

which completes the proof. □

## 3 Quantization error

The goal of this section is to show that, given a polynomial whose coefficients in the Bernstein basis is $\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ where $\|a\|_\infty \leq \mu$, for a prescribed $\mu \in (0,1)$, we can find a sequence $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n} \subseteq \{\pm 1\}$ such that the quantization error induced on the Bernstein basis,

$$\mathcal{E}_{n,a}^{\mathrm{quan}}(\boldsymbol{x}) := \sum_{0 \leq \boldsymbol{k} \leq n} (a_{\boldsymbol{k}} - \sigma_{\boldsymbol{k}}) \, p_{n,\boldsymbol{k}}(\boldsymbol{x}),$$

is small in a suitable sense. Theorem 3.5 will provide us with an explicit algorithm that will convert $\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ into an appropriate $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$.

## 3.1 Background on $\Sigma\Delta$ quantization

$\Sigma\Delta$ quantization (or modulation) refers to a large family of algorithms designed to convert any given sequence $y := \{y_k\}_{k\in\mathbb{N}}$ of real numbers in a given set $\mathcal{Y}$ to another sequence $q := \{q_k\}_{k\in\mathbb{N}}$ taking values in a prescribed discrete set $\mathcal{A}$, typically selected as an arithmetic progression. This is done in such a way that the quantization error is a "high-pass" sequence, in the sense that inner products of $q$ with slowly varying sequences are small.

A canonical way of ensuring this is to ask that for any $y \in \mathcal{Y}$, there exist $q$ and a "state" sequence $\{u_k\}_{k\in\mathbb{Z}}$ that satisfy the $r$-th order difference equation

$$y - q = \Delta^r u, \quad \text{and} \quad (\Delta u)_k := u_k - u_{k-1}. \tag{3.1}$$

If this is possible, we then say that $q$ is an $r$-th order noise-shaped quantization of $y$. When (3.1) is implemented recursively, it means that each $q_k$ is found by means of a quantization rule of the form

$$q_k = F(u_{k-1}, u_{k-2}, \dots, y_k, y_{k-1}, \dots),$$

and $u_k$ is updated via

$$u_k = \sum_{j=1}^{r} (-1)^{j-1} \binom{r}{j} u_{k-j} + y_k - q_k.$$

**Definition 3.1.** A quantization rule is stable for $\mathcal{Y}$ if for each $y \in \mathcal{Y}$, there is a sequence $u$ satisfying (3.1) such that $\|u\|_\infty$ is bounded uniformly in $y$.

Stability is a desirable property for a quantization algorithm as it allows one to control the error $y - q$ uniformly in $y \in \mathcal{Y}$. Establishing the existence of a stable $r$-th order scheme for fixed and finite $\mathcal{A}$, especially in the extreme one-bit case where $\mathcal{A} = \{\pm a\}$ for some $a \neq 0$, is difficult. The first breakthrough on this problem was made in the seminal paper of Daubechies and DeVore [11]. There it was shown for $\mathcal{A} = \{\pm 1\}$, any $r \geq 1$, and any $\mu \in (0,1)$, there exists a stable $r$-th order $\Sigma\Delta$ quantizer such that whenever $\|y\|_\infty \leq \mu$, it holds that $\|u\|_\infty \leq C_{r,\mu}$. The constant $C_{r,\mu}$ depends only on $r$ and $\mu$, and blows up as $r \to \infty$ or $\mu \to 1$, except for when $r = 1$. For $r = 1$, it suffices to take $C_{r,\mu} = 1$ and $\mu \in (0,1]$. Another family of stable $\Sigma\Delta$ quantizers, but with more favorable $C_{r,\mu}$, was subsequently proposed in [20].

In this paper, $r$-th order $\Sigma\Delta$ quantization refers to either of the stable rules in [11] and [20] for when $\mathcal{A} = \{\pm 1\}$. Except for Section 5, we will use $\Sigma\Delta$ solely as a method of approximation, and will not need explicit descriptions of these rules or precise estimates for $C_{r,\mu}$. If new stable $r$-th order $\Sigma\Delta$ quantization methods are developed, those can be used instead of the ones mentioned here, without changes to the remaining parts of this paper.

We will need to extend $\Sigma\Delta$ to sequences indexed by $\mathbb{N}^d$. Let $u$ be a function on $\mathbb{N}^d$. For each integer $1 \leq \ell \leq d$, we denote the finite difference operator in the $\ell$-th coordinate by $\Delta_\ell$, which acts on $u$ by the formula

$$(\Delta_\ell u)_{\boldsymbol{k}} := u_{k_1,\dots,k_{\ell-1},k_\ell,k_{\ell+1},\dots,k_d} - u_{k_1,\dots,k_{\ell-1},k_\ell-1,k_{\ell+1},\dots,k_d}$$

The following shows that there exists a stable directional $r$-th order $\Sigma\Delta$ quantizer.

**Proposition 3.2.** *For any $\mu \in (0,1)$, and integers $r, d \geq 1$, there exists a $C_{r,\mu} > 0$ that depends only on $r$ and $\mu$ such that the following hold. For any function $y$ defined on $\mathbb{N}^d$ such that $\|y\|_\infty \leq \mu$ and any integer $1 \leq \ell \leq d$, there exist a function $q$ on $\mathbb{N}^d$ where $q_{\boldsymbol{k}} \in \{\pm 1\}$ for all $\boldsymbol{k} \in \mathbb{N}^d$ and a function $u$ on $\mathbb{Z}^d$ supported in $\mathbb{N}^d$, such that $\|u\|_\infty \leq C_{r,\mu}$ and*

$$y - q = \Delta_\ell^r u.$$

*For $r = 1$, we only require that $\mu \in (0,1]$ and the statement holds for $C_{1,\mu} = 1$.*

*Proof.* Fix an integer $1 \leq \ell \leq d$, and for each $\boldsymbol{k} \in \mathbb{N}^d$, we define

$$\boldsymbol{k}' = (k_1, \ldots, k_{\ell-1}, k_{\ell+1}, \ldots, k_d) \in \mathbb{N}^{d-1}.$$

Now we form a sequence from $y$ by setting all indices except for the $\ell$-th one to be $\boldsymbol{k}'$,

$$y_{\boldsymbol{k}'} := \left\{ y_{k_1, \ldots, k_{\ell-1}, j, k_{\ell+1}, \ldots, k_d} \right\}_{j \in \mathbb{N}}.$$

Using either one of the two quantization rules pointed out in Section 3.1, we define a $q$ on $\mathbb{N}^d$ such that $q_{\boldsymbol{k}'} \in \{\pm 1\}$ and is defined as a solution to the $r$-th order difference equation

$$y_{\boldsymbol{k}'} - q_{\boldsymbol{k}'} = \Delta^r(u_{\boldsymbol{k}'}).$$

This scheme is stable, in the sense that there exists $C_{r,\mu} > 0$ such that for each $\boldsymbol{k}'$, we have $\|u_{\boldsymbol{k}'}\|_\infty \leq C_{r,\mu}$. Hence $\|u\|_\infty \leq C_{r,\mu}$. $\qquad\square$

It follows from this proposition that any stable $r$-th order $\Sigma\Delta$ can be used to generate a stable $r$-th order scheme in arbitrary dimensions.

**Definition 3.3.** A $r$-th order $\Sigma\Delta$ applied to the $\ell$-th direction is a map that satisfies the conclusions of Proposition 3.2.

## 3.2 Directional $\Sigma\Delta$ quantization on Bernstein polynomials

Going back to our original goal of quantizing the coefficients $a = \{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ of a polynomial in the Bernstein basis, we process the coefficients with a $r$-th order $\Sigma\Delta$ applied to the $\ell$-th direction, as in Proposition 3.2. Letting $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ be the $\pm 1$ sequence produced by this algorithm and $u$ denote the state, the quantization error induced on the Bernstein basis is

$$\mathcal{E}_{n,a}^{\text{quan}}(\boldsymbol{x}) = \sum_{0 \leq \boldsymbol{k} \leq n} (\Delta_\ell^r u)_{\boldsymbol{k}} \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) = \sum_{0 \leq \boldsymbol{k} \leq n} u_{\boldsymbol{k}} \left( (\Delta_\ell^*)^r p_{n,\cdot}(\boldsymbol{x}) \right)_{\boldsymbol{k}}. \tag{3.2}$$

Here, we let $\Delta_\ell^*$ be the adjoint of $\Delta_\ell$, and

$$(\Delta_\ell^* p_{n,\cdot}(\boldsymbol{x}))_{\boldsymbol{k}} := \left( p_{n,k_\ell}(x_\ell) - p_{n,k_\ell+1}(x_\ell) \right) \prod_{j \neq \ell} p_{n,k_j}(x_j).$$

Notice that with our convention that $p_{n,\boldsymbol{k}} = 0$ if there is a $\ell$ such that $k_\ell > n$, all the boundary terms are correctly included in (3.2). We are now ready to state the following result for first order $\Sigma\Delta$ on the Bernstein basis.

**Proposition 3.4.** *For any integers* $n, d \geq 1$, $\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ *with* $\|a\|_\infty \leq 1$, *and* $1 \leq \ell \leq d$, *if* $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ *is the output of a stable first order* $\Sigma\Delta$ *quantizer applied in the* $\ell$-th *direction with input* $a$, *then*

$$\left| \sum_{0 \leq \boldsymbol{k} \leq n} (a_{\boldsymbol{k}} - \sigma_{\boldsymbol{k}}) \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| \leq \min \left( 2, n^{-1/2} x_\ell^{1/2} (1 - x_\ell)^{-1/2} \right).$$

*Proof.* It follows from the inequality $\|u\|_\infty \leq C_{1,\mu} = 1$ and identity (3.2) for $r = 1$, we have

$$\left| \sum_{0 \leq \boldsymbol{k} \leq n} (a_{\boldsymbol{k}} - \sigma_{\boldsymbol{k}}) \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| \leq \sum_{0 \leq \boldsymbol{k} \leq n} \left| (\Delta_\ell^* p_{n,\cdot}(\boldsymbol{x}))_{\boldsymbol{k}} \right|.$$

16

The consecutive differences of the $p_{n,k}$, for the one dimensional case, satisfies the following identity: for all $x \in [0,1]$ and $0 \leq k \leq n$,

$$p_{n,k}(x) - p_{n,k+1}(x) = \frac{(k+1)-(n+1)x}{(n+1)x(1-x)} \, p_{n+1,k+1}(x). \tag{3.3}$$

See [14, Chapter 10] and [29] for this identity. When interpreting the right hand side of this equation for $x = 0$ or $x = 1$, it should be observed that the polynomial $((k+1)-(n+1)x)p_{n+1,k+1}(x)$ is divisible by $x(1-x)$ for each $k$. We proceed to extend (3.3) to the multivariate case. For each integer $1 \leq \ell \leq d$, we have

$$\left( \Delta_\ell^* p_{n,\cdot}(\boldsymbol{x}) \right)_{\boldsymbol{k}} = \left( \frac{(k_\ell + 1) - (n+1)x_\ell}{(n+1)x_\ell(1-x_\ell)} \, p_{n+1,k_\ell+1}(x_\ell) \right) \prod_{j \neq \ell} p_{n,k_j}(x_j). \tag{3.4}$$

By identity (3.4), the partition of unity property (2.3), Cauchy-Schwarz, and central moment bounds (2.2), we have

$$\sum_{0 \leq \boldsymbol{k} \leq n} \left| (\Delta_\ell^* p_{n,\cdot}(\boldsymbol{x}))_{\boldsymbol{k}} \right| = \sum_{0 \leq \boldsymbol{k} \leq n} \left| \frac{(k_\ell + 1) - (n+1)x_\ell}{(n+1)x_\ell(1-x_\ell)} \right| p_{n+1,k_\ell+1}(x_\ell) \prod_{j \neq \ell} p_{n,k_j}(x_j)$$

$$= \frac{1}{(n+1)x_\ell(1-x_\ell)} \sum_{k_\ell=0}^{n} \left| (k_\ell + 1) - (n+1)x_\ell \right| p_{n+1,k_\ell+1}(x_\ell)$$

$$\leq \frac{\sqrt{T_{n+1,2}(x_\ell)}}{(n+1)x_\ell(1-x_\ell)} \left( \sum_{k_\ell=0}^{n} p_{n+1,k_\ell+1}(x_\ell) \right)^{1/2}$$

$$= \frac{1}{\sqrt{(n+1)x_\ell(1-x_\ell)}}.$$

On the other hand, we have the trivial bound for the quantization error,

$$\left| \sum_{0 \leq \boldsymbol{k} \leq n} (a_{\boldsymbol{k}} - \sigma_{\boldsymbol{k}}) \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| \leq \|a - \sigma\|_\infty \sum_{0 \leq \boldsymbol{k} \leq n} p_{n,\boldsymbol{k}}(\boldsymbol{x}) = \|\Delta_\ell u\|_\infty \leq 2.$$

$\square$

For larger values of $r$, due to an increasing complexity in the formulas for $(\Delta_\ell^*)^r p_{n,\cdot}(\boldsymbol{x})$, we do not provide explicit upper bounds. Our strategy for deriving an upper bound for the quantization error builds upon a corresponding one-dimensional result in [22, Theorem 6].

**Theorem 3.5.** *For any integers* $n, d, r \geq 1$, $\mu \in (0,1)$, $\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ *with* $\|a\|_\infty \leq \mu$, *and* $1 \leq \ell \leq d$, *if* $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ *is the output of a stable $r$-th order $\Sigma\Delta$ quantization applied in the $\ell$-th direction with input* $\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$, *then*

$$\left| \sum_{0 \leq \boldsymbol{k} \leq n} (a_{\boldsymbol{k}} - \sigma_{\boldsymbol{k}}) \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| \lesssim_{r,\mu} \min \left( 1, n^{-r/2} x_\ell^{-r}(1 - x_\ell)^{-r} \right).$$

*Proof.* Since $\Delta_\ell^*$ is applied to the $\ell$-th direction and the multivariate Bernstein polynomials are tensor products, we have

$$\left( (\Delta_\ell^*)^r p_{n,\cdot}(\boldsymbol{x}) \right)_{\boldsymbol{k}} = \left( (\Delta^*)^r p_{n,\cdot}(x_\ell) \right)_{k_\ell} \left( \prod_{j \neq \ell} p_{n,k_j}(x_j) \right).$$

17

Employing (3.2), non-negativity of the Bernstein polynomials, and partition of unity (2.3),

$$\Big| \sum_{0 \le \boldsymbol{k} \le n} (a_{\boldsymbol{k}} - \sigma_{\boldsymbol{k}}) \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \Big| \le \|u\|_\infty \sum_{0 \le \boldsymbol{k} \le n} \big|((\Delta_\ell^*)^r p_{n,\cdot}(\boldsymbol{x}))_{\boldsymbol{k}}\big| \le C_{\mu,r} \sum_{k_\ell=0}^{n} \big|((\Delta^*)^r p_{n,\cdot}(x_\ell))_{k_\ell}\big|.$$

This quantity was upper bounded in [22, Theorem 6], and by the referenced bound, we have

$$\sum_{k_\ell=0}^{n} \big|((\Delta^*)^r p_{n,\cdot}(x_\ell))_{k_\ell}\big| \lesssim_r n^{-r/2} x_\ell^{-r} (1 - x_\ell)^{-r}.$$

On the other hand, we have the trivial upper bound for the quantization error,

$$\Big| \sum_{0 \le \boldsymbol{k} \le n} (a_{\boldsymbol{k}} - \sigma_{\boldsymbol{k}}) \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \Big| \le \|a - \sigma\|_\infty \sum_{0 \le \boldsymbol{k} \le n} p_{n,k}(\boldsymbol{x}) = \|\Delta_\ell^r u\|_\infty \le 2^r \|u\|_\infty \le 2^r C_{r,\mu}.$$

$\square$

## 3.3 Comments on the quantization results

In contrast to Theorem 3.5 for the Bernstein basis, we explain why $\{\pm 1\}$ linear combinations of the power basis $\mathcal{P}_n := \{x \mapsto x^k\}_{k=0}^{n}$, by which we mean any function of the form $p(x) = \sum_{k=0}^{n} \sigma_k x^k$ where $\sigma_k \in \{\pm 1\}$ for each $k$, cannot accurately approximate continuous functions on $[0, 1]$. We provide two different set of explanations.

All possible real numbers that can be realized as an output of a $\{\pm 1\}$ linear combination in the power basis up to degree $n$ with input $x$ is the set

$$P_n(x) = \Big\{ \sum_{k=0}^{n} \sigma_k x^k \colon \sigma_k \in \{\pm 1\} \text{ for all } k \Big\}.$$

A plot of $P_{10}$ is shown in Section 3.3. It is straightforward to see that

$$P_n(x) \subseteq [1 - r_n(x), \, 1 + r_n(x)] \cup [-1 - r_n(x), \, -1 + r_n(x)], \quad \text{where} \quad r_n(x) = x \frac{1 - x^n}{1 - x}.$$

Hence, $P_n(x)$ is a strict subset of $[-2, 2]$ whenever $x \in (0, 1/2)$ and it becomes an increasingly smaller subset as $x \to 0$. For any $c \in (0, 1/2)$, all sufficiently small $\varepsilon > 0$, and all $n \ge 1$, we have

$$\sup_{f \in C([0,1])} \, \inf_{\sigma_0, \ldots, \sigma_n \in \{\pm 1\}} \, \sup_{x \in [0,c]} \Big| f(x) - \sum_{k=0}^{n} \sigma_k x^k \Big| > \varepsilon.$$

The key component of this statement is that $c$ is fixed independent of $n$. In contrast to Theorem 3.5, the regions for which the approximation error is large shrinks to zero as $n \to \infty$. For all sufficiently large $n$, this region consists of two $d$-dimensional rectangles whose measures are $O(1/\sqrt{n})$.

Related to this discussion is the observation that $P(x) := \bigcup_{n=0}^{\infty} P_n(x)$ for $x \in (0, 1/2)$ has a fractal structure and is totally disconnected. Indeed, $P(x)$ is the set of all real numbers that can be written in base $x$ with $-1, 0, +1$ digits and of length $n$. for any $a, b \in P(x)$ with $a < b$, we can express them as a sequence in base $x$, and let $m$ be the first digit for which their base $x$ expansions disagree. Call their first $m - 1$ digits $\sigma_0, \ldots, \sigma_{m-1}$. Any $c \in P(x)$ with $a < c < b$ can be written in base $x$ as $\sigma_0, \ldots, \sigma_{m-1}, \varepsilon_m, \ldots$, where $\varepsilon_k \in \{-1, 0, 1\}$. We recall the basic observation that
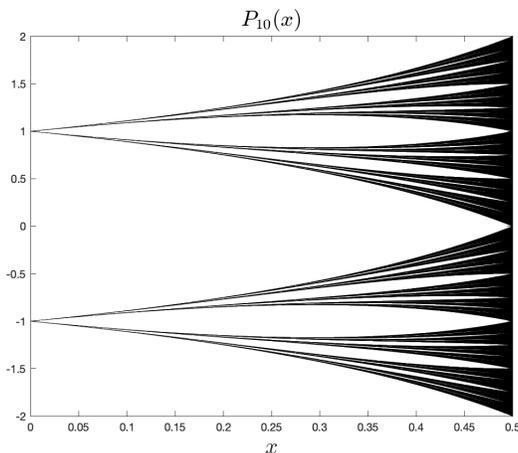
Figure 2: Plot of $P_{10}$ on $[0, \frac{1}{2}]$.

$\sum_{k=m}^{\infty} \varepsilon_k x^k \in (-x^m/(1-x), x^m/(1-x)) \subseteq (-x^{m-1}, x^{m-1})$ since $x \in (0, 1/2)$. Hence we can always find a $c$ such that $a < c < b$ and $c \notin P(x)$.

Of course, the previous discussion only pertains to the power basis' inability to approximate continuous functions near the origin. It was shown in [21] that $\{\pm 1\}$ linear combinations in the power basis are able to approximate certain power series in the complex plane near the point $z = 1$, so approximation is possible in other regions.

From a quantization perspective, Theorem 3.5 is perhaps surprising because previous applications of noise-shaping quantization (including $\Sigma\Delta$) [11, 20, 21, 8, 24], utilize some notion of redundancy in the system. However, the multivariate Bernstein polynomials of order $n$ forms a basis for the space of multivariate polynomials whose degree in each variable is at most $n$, so it does not exhibit redundancy in the traditional sense.

Instead, the Bernstein system exhibits a different type of redundancy. To make this notion more precise, we define the synthesis operator $S_n : \mathbb{R}^{n+1} \to \mathcal{P}_n$ by $S_n u := \sum_{k=0}^{n} u_k p_{n,k}$. Using the usual inner products on both $\mathbb{R}^{n+1}$ and $\mathcal{P}_n$, a direct calculation shows that the frame operator $S_n^* S_n : \mathbb{R}^{n+1} \to \mathbb{R}^{n+1}$ is represented as the matrix $B$ such that $B_{j,k} := \int_0^1 p_{n,j}(x) p_{n,k}(x) \, dx$. The $\varepsilon$ numerical rank of a $(n+1) \times (n+1)$ matrix $A$ by

$$d_\varepsilon(A) := \max\{0 \leq m \leq n \colon \sigma_m(A) \geq \varepsilon_0(A)\}.$$

Then [22, Appendix B] showed that for fixed $\varepsilon > 0$, we have

$$d_\varepsilon(B) = \sqrt{2 \ln(1/\varepsilon)} \sqrt{n}(1 + o(1)) \quad \text{as} \quad n \to \infty.$$

This partially explains why noise-shaping quantization in the Bernstein basis is possible and why Theorems 2.4 and 3.5 exhibit decay rates of $(\sqrt{n})^{-s}$ and $(\sqrt{n})^{-r}$ respectively.

Continuing this line of discussion and to connect it with the prior comparison between the Bernstein and power basis, notice that the synthesis operator for the power basis is $T_n : \mathbb{R}^{n+1} \to \mathcal{P}_n$ where $T_n u = \sum_{k=0}^{n} u_k x^k$. Then the frame operator $T_n^* T_n$ can be identified with the matrix $H$ which has entries $H_{j,k} := \int_0^1 x^j x^k \, dx = 1/(j + k + 1)$. This is precisely the Hilbert matrix, and it follows from [4, Corollary 4.2] that

$$d_\varepsilon(H) \leq \left\lceil \frac{\log(8n - 4) \log(4/\varepsilon)}{\pi^2} \right\rceil.$$

Hence, the numerical rank for the power basis is roughly $\log(n)$ compared to that of $\sqrt{n}$ for the Bernstein basis.

## 3.4 Proof of Theorem A

*Proof.* According to Theorem 2.6, using $(1-\mu)/2$ as $\delta$ in the referenced theorem and the assumption that $n \geq \dfrac{\sqrt{2^{s+1}d}}{4(1-\mu)}\|f\|_{C^1\mathrm{Lip}}$, if $s \geq 3$, there exist $\{a_{\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n}$ such that $\|a\|_\infty \leq (\mu+1)/2 < 1$ and

$$\left\| f - \sum_{0\leq\boldsymbol{k}\leq n} a_{\boldsymbol{k}} p_{n,\boldsymbol{k}} \right\|_\infty \lesssim_{s,d} \|f\|_{C^s} n^{-s/2}.$$

For any $1 \leq \ell \leq d$, applying $s$-th order $\Sigma\Delta$ in the $\ell$-th direction on $\{a_{\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n}$, noting that $\|a\|_\infty \leq (\mu+1)/2 < 1$, we obtain $\{\sigma_{\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n} \subseteq \{\pm 1\}$. According to Theorem 3.5, the induced quantization error satisfies for all $\boldsymbol{x} \in [0,1]^d$,

$$\left| \sum_{0\leq\boldsymbol{k}\leq n} a_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) - \sum_{0\leq\boldsymbol{k}\leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| \lesssim_{s,\mu} \min\left(1, n^{-s/2} x_\ell^{-s}(1-x_\ell)^{-s}\right).$$

Combining these inequalities completes the proof. $\qquad\square$

## 4 Implementation error

In this section, we concentrate on the implementation error by one-bit quantized neural networks. That is, suppose for some $\{\sigma_{\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n}$ such that $\sigma_{\boldsymbol{k}} \in \mathcal{A}$, we would like find functions $\{b_{n,\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n}$ such that $\sum_{0\leq\boldsymbol{k}\leq n} \sigma_{\boldsymbol{k}} b_{n,\boldsymbol{k}}$ is implementable by a strict $\mathcal{A}$-quantized neural network with activation $\beta$ for which the implementation error

$$\mathcal{E}^{\mathrm{imp}}_{n,\sigma,\mathcal{A},\beta}(\boldsymbol{x}) := \sum_{0\leq\boldsymbol{k}\leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) - \sum_{0\leq\boldsymbol{k}\leq n} \sigma_{\boldsymbol{k}} b_{n,\boldsymbol{k}}(\boldsymbol{x}),$$

is suitably controlled. We address two pairs of activation functions and one-bit alphabets.

(a) The first pair is the quadratic activation $\rho\colon \mathbb{R} \to \mathbb{R}$ defined as $\rho(t) = \frac{1}{2}t^2$, together with the one-bit alphabet $\mathcal{A}_1 := \{\pm 1\}$.

(b) The second pair is the ReLU activation $\sigma\colon \mathbb{R} \to \mathbb{R}$ defined as $\sigma(t) = \max(t,0)$, together with the one-bit alphabet $\mathcal{A}_{1/2} := \{\pm\frac{1}{2}\}$.

For both activation functions, the main mechanism behind our construction hinges on a Pascal triangle interpretation of the Bernstein polynomials, which we proceed to explain in the univariate case first. Each Bernstein polynomial of degree $m+1$ can be made by multiplying at most two pairs of Bernstein polynomials, with degrees $m$ and 1, and adding them together. The key formulas are, for each $m \geq 1$,

$$p_{m+1,k}(x) = \begin{cases} (1-x)p_{m,0}(x) & \text{if } k=0, \\ x p_{m,k-1}(x) + (1-x)p_{m,k}(x) & \text{if } 0 < k < m, \\ x p_{m,m}(x) & \text{if } k = m+1. \end{cases} \tag{4.1}$$

These recurrence relations are summarized in a Pascal-like network as shown in Figure 1b. Due to the combinatorial factors that appear in the Bernstein polynomials, the Pascal triangle interpretation is crucial when attempting to implement approximations of the Bernstein polynomials in a stable way when using only parameters from a small set.

20

Definitions of strict (unquantized and quantized) neural networks were provided in Definitions 1.1 and 1.2. As discussed earlier, strict neural networks do not use skip connections and all intermediate layers use the same activation. Here, we introduce additional definitions and some basic concepts that will help facilitate the subsequent proofs.

One cumbersome issue is that if $F \colon \mathbb{R}^d \to \mathbb{R}^m$ and $G \colon \mathbb{R}^m \to \mathbb{R}^n$ are both implementable by strict neural networks, then their composition $G \circ F$ is not necessarily implementable because by using the output of $F$ as the input of $G$ results in network with a layer without activation. For this reason, we introduce the following definition.

**Definition 4.1.** An activated neural network with activation $\beta$ is any function $F \colon \mathbb{R}^d \to \mathbb{R}^m$ of the form,

$$F(\boldsymbol{x}) := \beta(W_L \beta(W_{L-1} \cdots \beta(W_1(\boldsymbol{x})))), \quad W_\ell(\boldsymbol{u}) := A_\ell \boldsymbol{u} + \boldsymbol{b}_\ell, \quad \text{for} \quad \ell = 1, \ldots, L.$$

Hence, an activated neural network also does not have skip connections, but each layer, including the output layer, uses the same activation function $\beta$. We say the network is unquantized if the weights and biases are allowed to use any real number, while it is $\mathcal{A}$-quantized if they are selected from $\mathcal{A}$ only. The number of layers, nodes, and parameters of an activated neural network are defined in the same way as for strict neural networks. With this terminology in place, we make the following basic observations.

(a) If $F \colon \mathbb{R}^d \to \mathbb{R}^m$ is implementable by an activated network of size $(L_1, N_1, P_1)$ and $G \colon \mathbb{R}^m \to \mathbb{R}^n$ is implementable by a strict (resp., activated) neural network of size $(L_2, N_2, P_2)$, then $G \circ F$ is implementable by a strict (resp., activated) network of size $(L_1 + L_2, N_1 + N_2, P_1 + P_2)$. Naturally, we say that the second network is appended to the first or that the networks are composed.

(b) If $F \colon \mathbb{R}^d \to \mathbb{R}^m$ and $G \colon \mathbb{R}^d \to \mathbb{R}^n$ are both implementable by strict (resp., activated) networks of size $(L, N_1, P_1)$ and $(L, N_2, P_2)$, then the function $F \oplus G \colon \mathbb{R}^d \to \mathbb{R}^{m+n}$ is implementable by a strict (resp., activated) network of size $(L, N_1 + N_2, P_1 + P_2)$. We say that these networks are placed in parallel.

## 4.1 Overview for $\{\pm 1\}$-quantized quadratic networks

In this subsection we provide a high level discussion of our constructions for $\mathcal{A}_1$-quantized neural networks with quadratic activation $\rho$, and postpone the details for Appendix A.1. In the expository portions, we just use the term "network" since the alphabet and activation do not change in this subsection.

To turn the schematic diagrams shown in Figures 1a and 1b into a proper network, it suffices to convert multiplications by $1 - x$ and $x$ into neural network operations. Our methodology is inspired by the basic identity,

$$ab = \rho(a + b) - \rho(a) - \rho(b). \tag{4.2}$$

It is important to remark that this identity cannot be (directly) used for intermediate layers, since strict neural networks require every node to be activated and hence such a node can only produce $\rho(ab)$, not $ab$ itself. Hence, we cannot realize the product function as an activated neural network.

This technical issue can be circumvented under many situations in light of the following observation. If an intermediate layer's nodes output $\rho(a + b)$, $\rho(a)$, $\rho(b)$, and $c$, then a subsequent intermediate layer can implement $\rho(ab + c)$ in view the identity

$$\rho(ab + c) = \rho(\rho(a + b) - \rho(a) - \rho(b) + c).$$

21

Hence, although $ab$ cannot be directly implemented in an intermediate layer, it can be used as an input into a subsequent intermediate or final layer. This type of reasoning can be adapted to more general situations beyond implementing $\rho(ab + c)$.

To see how this observation is relevant to Bernstein, we define the functions

$$U(x) := \rho(1 - x), \quad V(x) := \rho(x), \quad X_{0,0}(x) := \rho(1),$$
$$Y_{0,0}(x) := \rho(1 - x + 1), \quad Z_{0,0}(x) := \rho(x + 1).$$

For each integer $m \geq 1$ and $k = 0, \ldots, m$, we define

$$X_{m,k}(x) = \rho(p_{m,k}(x)), \quad Y_{m,k}(x) = \rho(1 - x + p_{m,k}(x)), \quad Z_{m,k}(x) = \rho(x + p_{m,k}(x)).$$

It follows from the recurrence relation (4.1) and identity (4.2) that

$$p_{m,k} = \begin{cases} Y_{m-1,0} - U - X_{m-1,0} & \text{if } k = 0, \\ Z_{m-1,k-1} - V - X_{m-1,k-1} + Y_{m-1,k} - U - X_{m-1,k} & \text{if } 0 < k < m, \\ Z_{m-1,m-1} - V - X_{m-1,m-1} & \text{if } k = m. \end{cases} \quad (4.3)$$

These formulas show that each $U, V, X, Y, Z$ can be realized as an intermediate layer's output by composing networks, e.g., $X_{m+1,k}$ can be produced given $X_{m,j}, Y_{m,j}, Z_{m,j}, U, V$ for each $0 \leq j \leq m$. While each Bernstein polynomial does not actually correspond to the output of a network used in our final construction, it will be a node's preactivation. This implies that any linear combination $\sum_{k=0}^{n} \sigma_k p_{n,k}$ with $\sigma_k \in \{\pm 1\}$ for each $k$, is implementable by a network.

For the multivariate case, in view of (4.3), each $p_{n,\boldsymbol{k}}$ is a sum of tensor products consisting of various combinations from $U, V, X, Y, Z$. To implement tensor products, we need multiplication. Similar to the bivariate case (4.2), a multivariate product cannot be realized as an activated network, and it will instead be made implicitly in a subsequently layer. For example, to use the product $abc$ as an argument in a future intermediate layer, if one layer outputs

$$\rho(a), \ \rho(a + b), \ \rho(b), \ \rho(c), \ \rho(c + 1), \ \rho(1),$$

then it possible for the next layer to output $\rho(ab + c), \ \rho(ab), \ \rho(c)$, due to the identities

$$\rho(ab + c) = \rho\big(\rho(a + b) - \rho(a) - \rho(b) + \rho(c + 1) - \rho(c) - \rho(1)\big),$$
$$\rho(ab) = \rho(a + b) - \rho(a) - \rho(b),$$
$$\rho(c) = \rho(c + 1) - \rho(c) - \rho(1).$$

The significance here is that a $\{\pm 1\}$ linear combination of $\rho(ab + c), \ \rho(ab), \ \rho(c)$, is $abc$ which can be implicitly formed in the layer afterwards. By continuing this process, it is possible to create the terms necessary to produce multivariate multiplication. After these considerations, it is possible to show that for any $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ with $\sigma_{\boldsymbol{k}} \in \{\pm 1\}$, the sum $\sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}$ is implementable.

## 4.2  Overview for $\{\pm\frac{1}{2}\}$-quantized ReLU networks

In this subsection we provide a high level discussion of our constructions for $\mathcal{A}_{1/2}$-quantized neural networks with ReLU activation $\sigma$, and postpone the details for Appendix A.2. In the expository portions, we just use the term "network" since the alphabet and activation do not change in this subsection.

Since ReLU is the conventional activation used in practice, the approximation properties of unquantized ReLU networks have been thoroughly studied. We use some standard ideas popularized

by [39]. The starting point is the tent function $\phi\colon [0,1] \to [0,1]$, where $\phi(x) := 2x$ for $x \in [0, \frac{1}{2}]$ and $\phi(x) := 2 - 2x$ for $x \in [\frac{1}{2}, 1]$, and the identity,

$$x(1-x) = \sum_{k=1}^{\infty} \frac{\phi^{\circ k}(x)}{4^k}, \tag{4.4}$$

where the series converges uniformly in $x$ and $\phi^{\circ k}$ refers to the composition of $\phi$ with itself $k$-times, with the convention that $\phi^{\circ 1} := \phi$. Since $\phi$ can be implemented, this identity provides a natural method for implementation of the squaring function and hence multiplication via the basic identity

$$ab = \frac{1}{2}(a+b)^2 - \frac{1}{2}a^2 - \frac{1}{2}b^2. \tag{4.5}$$

Let us make some comments about how this strategy can be modified in order to account for the lack of skip connections and one-bit quantization. For technical reasons, we observed that approximation of the squaring function is insufficient as it leads to requiring additional bits. We instead approximate $x^2$ from both above and below by nonnegative $S_\varepsilon^+$ and $S_\varepsilon^-$ respectively, with error $\varepsilon$ uniformly in $x$. The implementations of $S_\varepsilon^+$ and $S_\varepsilon^-$ are also further complicated by the constraint of not being able to use skip connections. Although we use "duplication" networks to pass values down a specified number of layers, exaggerated use results in bloated networks and care is taken to use them sparingly. By using $S_\varepsilon^\pm$ and mimicking (4.5), we define the approximate multiplication function

$$P_\varepsilon(x,y) = \sigma\left( 2S_{\varepsilon/6}^-\left(\frac{x+y}{2}\right) - 2S_{\varepsilon/6}^+\left(\frac{x}{2}\right) - 2S_{\varepsilon/6}^+\left(\frac{y}{2}\right) \right).$$

Not only does $P_\varepsilon$ approximate the product function uniformly with error $\varepsilon$, it also satisfies the inequalities $0 \leq P_\varepsilon(x,y) \leq xy$. This is important since we use strict neural networks and the ReLU function is the identity on nonnegative real numbers.

With approximate bivariate multiplication at hand, approximations of univariate Bernstein polynomials can be constructed mimicking recurrence (4.1). That is, we define $\{b_{m,k}\}_{0 \leq k \leq n}$ recursively starting with $b_{1,0}(x) = 1 - x$ and $b_{1,1}(x) = x$ and for $m \geq 1$,

$$b_{m+1,k}(x) = \begin{cases} P_\varepsilon(1-x, b_{m,0}(x)) & \text{if } k = 0, \\ P_\varepsilon(x, b_{m,k-1}(x)) + P_\varepsilon(1-x, b_{m,k}(x)) & \text{if } 0 < k < m, \\ P_\varepsilon(x, b_{m,m}(x)) & \text{if } k = m+1. \end{cases}$$

These can be approximately multiplied together to yield a collection $\{b_{n,\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ that approximate the multivariate Bernstein polynomials $\{p_{n,\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ uniformly on $[0,1]^d$. Finally, we approximate any $f = \sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}$ by $\sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{n,\boldsymbol{k}} b_{n,\boldsymbol{k}}$, which is implementable under the assumption that $\sigma_{\boldsymbol{k}} \in \{\pm \frac{1}{2}\}$.

## 4.3   Proof of Theorem B

*Proof.* Quadratic case and $\{\pm 1\}$ alphabet. From identity (4.3), we see that each univariate Bernstein polynomial $p_{n,k}$ is a $\pm 1$ linear combination of at most six functions in the set

$$\{U, V, X_{n-1,k}, Y_{n-1,k}, Z_{n-1,k} \colon k = 0, \ldots, n-1\}.$$

To simplify the notation, denote this set of $3n + 2$ functions by $\{\psi_j\}_{j=1}^{3n+2}$. Then for each $0 \leq k \leq n$, there exists $\{\varepsilon_{k,j}\}_{j=1}^{3n+2}$ such that $\varepsilon_{k,j} \in \{0,1\}$ with at most six that are nonzero, and $p_{n,k} = \sum_{j=1}^{3n+2} \varepsilon_{k,j} \psi_j$.

23

To handle the multivariate case, we set $\psi_{\boldsymbol{j}}(\boldsymbol{x}) := \psi_{j_1}(x_1) \cdots \psi_{j_d}(x_d)$ and $\varepsilon_{\boldsymbol{k},\boldsymbol{j}} := \varepsilon_{k_1,j_1} \cdots \varepsilon_{k_d,j_d}$. Then we have

$$\sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}} = \sum_{0 \leq \boldsymbol{k} \leq n} \sum_{1 \leq \boldsymbol{j} \leq 3n+2} \sigma_{\boldsymbol{k}} \varepsilon_{\boldsymbol{k},\boldsymbol{j}} \psi_{\boldsymbol{j}} = \sum_{1 \leq \boldsymbol{j} \leq 3n+2} \Big( \sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} \varepsilon_{\boldsymbol{k},\boldsymbol{j}} \Big) \psi_{\boldsymbol{j}}. \tag{4.6}$$

It is important to remark that this is not necessarily a $\pm 1$ combination of $\psi_{\boldsymbol{j}}$'s because it is possible for a $\psi_{\boldsymbol{j}}$ to repeat. In fact, there will be many repetitions as $U(x_1)U(x_2) \cdots U(x_d)$ appears $(n-1)^d$ times. However, it is possible to express the right side of (4.6) as a $\pm 1$ sum of $\psi_{\boldsymbol{j}}$'s if we allow for repetitions, and that this $\pm 1$ sum only requires at most $6^d(n+1)^d$ terms, in view of the observation that $|\{\boldsymbol{j} \colon \varepsilon_{\boldsymbol{k},\boldsymbol{j}} \neq 0\}| \leq 6^d$ uniformly in $\boldsymbol{k}$. Hence, there is a finite sequence $I$ with $|I| \leq 6^d(n+1)^d$ and $\{\widetilde{\sigma}_{\boldsymbol{j}}\}_{\boldsymbol{j} \in I}$ with $\widetilde{\sigma}_{\boldsymbol{j}} \in \{\pm 1\}$ such that

$$\sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) = \sum_{\boldsymbol{j} \in I} \widetilde{\sigma}_{\boldsymbol{j}} \psi_{\boldsymbol{j}}(\boldsymbol{x}).$$

We are now ready to implement this approximation strategy as a neural network. If $d = 1$, notice that

$$\sum_{k=0}^{n} \sigma_k p_{n,k} = \sigma_0 Y_{n-1,0} - \sigma_0 U - \sigma_0 X_{n-1,0}$$

$$+ \sum_{k=1}^{n-1} (\sigma_k Z_{n-1,k-1} - \sigma_k V - \sigma_k X_{n-1,k-1} + \sigma_k Y_{n-1,k} - \sigma_k U - \sigma_k X_{n-1,k})$$

$$+ \sigma_{n-1} Z_{n-1,n-1} - \sigma_{n-1} U - \sigma_{n-1} X_{n-1,n-1}.$$

Notice that each $X_{n-1,k}$ appears twice. We use two networks described in Lemma A.1 placed parallel with each other, so that they outputs all terms on the right hand side, and a linear layer produces the final summation. This network has size $O(n)$.

For $d \geq 2$, we use Lemma A.1 to generate $d$ networks in parallel, so that layer $n+1$ produces, for each $1 \leq \ell \leq d$, the outputs

$$U(x_\ell), V(x_\ell), \{X_{n-1,k}(x_\ell)\}_{k=0}^{n-1}, \{Y_{n-1,k}(x_\ell)\}_{k=0}^{n-1}, \{Z_{n-1,k}(x_\ell)\}_{k=0}^{n-1}. \tag{4.7}$$

Doing so requires a network of size $O(n)$. Let us momentarily fix a $\boldsymbol{j} \in I$. Since $\psi_{\boldsymbol{j}}(\boldsymbol{x}) = \psi_{j_1}(x_1) \cdots \psi_{j_d}(x_d)$ and each $\psi_{j_\ell}(x_\ell)$ is a $\pm 1$ summation of terms in (4.7) due to identity (4.3), we use the network constructed in Lemma A.2 that outputs the quantities

$$\rho\Big( \prod_{\ell \leq d_*} \psi_{\boldsymbol{j}}(x_\ell) + \prod_{\ell > d_*} \psi_{\boldsymbol{j}}(x_\ell) \Big), \ \rho\Big( \prod_{\ell \leq d_*} \psi_{\boldsymbol{j}}(x_\ell) \Big), \ \rho\Big( \prod_{\ell > d_*} \psi_{\boldsymbol{j}}(x_\ell) \Big).$$

We do this for each $\boldsymbol{j} \in I$ and place these networks in parallel, hence further requiring a network with $O(1)$ layers and $O(n^d)$ nodes and parameters. The final linear layer produces $\sum_{\boldsymbol{j} \in I} \widetilde{\sigma}_{\boldsymbol{j}} \psi_{\boldsymbol{j}}(\boldsymbol{x})$ since $\widetilde{\sigma}_{\boldsymbol{j}} \in \{\pm 1\}$ and $\psi_{\boldsymbol{j}}(\boldsymbol{x})$ is a $\pm 1$ linear combination of the above terms.

ReLU case and $\{\pm\frac{1}{2}\}$ alphabet. Let $\varepsilon > 0$. As shown in Lemma A.7, there is an activated $\{\pm\frac{1}{2}\}$-quantized ReLU neural network that implements a set of functions $\{b_{n,\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ such that $\|p_{n,\boldsymbol{k}} - b_{n,\boldsymbol{k}}\|_\infty \leq \varepsilon$ for each $0 \leq \boldsymbol{k} \leq n$. This network has $O(n \log(n/\varepsilon))$ layers and $O(n^2 \log(n/\varepsilon) + n^d \log(1/\varepsilon))$ nodes and parameters, as $n \to \infty$ and $\varepsilon \to 0$. We use two copies of these networks placed in parallel, so that the function

$$f_{NN,\sigma} := \sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} b_{n,\boldsymbol{k}} = \sum_{0 \leq \boldsymbol{k} \leq n} \frac{\sigma_{\boldsymbol{k}}}{2} b_{n,\boldsymbol{k}} + \sum_{0 \leq \boldsymbol{k} \leq n} \frac{\sigma_{\boldsymbol{k}}}{2} b_{n,\boldsymbol{k}}$$

24

is still implementable by a $\{\pm\frac{1}{2}\}$-quantized neural network after placing a linear layer with weights given by two copies of $\{\sigma_{\boldsymbol{k}}/2\}_{0\leq\boldsymbol{k}\leq n}$. This last layer has size $(1,1,2(n+1)^d)$. In total, $f_{NN,\sigma}$ is implementable by a $\{\pm\frac{1}{2}\}$-quantized neural network with $O(n\log(n/\varepsilon))$ layers and $O(n^2\log(n/\varepsilon) + n^d\log(1/\varepsilon))$ nodes and parameters. Since $|\sigma_{\boldsymbol{k}}| = 1$, we have

$$\Big\| \sum_{0\leq\boldsymbol{k}\leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}} - \sum_{0\leq\boldsymbol{k}\leq n} \sigma_{\boldsymbol{k}} b_{n,\boldsymbol{k}} \Big\|_\infty \leq \sum_{0\leq\boldsymbol{k}\leq n} \|p_{n,\boldsymbol{k}} - b_{n,\boldsymbol{k}}\|_\infty \leq (n+1)^d \varepsilon.$$

$\square$

### 4.4 Proof of Theorem C

*Proof.* Quadratic case and $\{\pm 1\}$ alphabet. It follows from Theorem A that for any $f \in C^s([0,1]^d)$ with $\|f\|_\infty \leq \mu$, there exist $\{\sigma_{\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n} \subseteq \{\pm 1\}$ such that for all $\boldsymbol{x} \in [0,1]^d$,

$$\Big| f(\boldsymbol{x}) - \sum_{0\leq\boldsymbol{k}\leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) \Big| \lesssim_{s,\mu} \|f\|_{C^s} \min\big(1, n^{-s/2} x_\ell^{-s}(1-x_\ell)^{-s}\big).$$

Using Theorem B for the quadratic case completes the proof.

ReLU case and $\{\pm\frac{1}{2}\}$ alphabet. We apply Theorem A to obtain there exist $\{\sigma_{\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n} \subseteq \{\pm 1\}$ such that for all $\boldsymbol{x} \in [0,1]^d$,

$$\Big| f(\boldsymbol{x}) - \sum_{0\leq\boldsymbol{k}\leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) \Big| \lesssim_{s,\mu} \|f\|_{C^s} \min\big(1, n^{-s/2} x_\ell^{-s}(1-x_\ell)^{-s}\big).$$

Using Theorem B for the ReLU case with $\varepsilon = \|f\|_{C^s} n^d n^{-s/2}$ completes the proof. $\square$

## 5 Final Remarks

### 5.1 Algorithm for computing the one-bit coefficients and stability to noise

The binary sequence $\{\sigma_{\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n}$ that appears in Theorem A, which is also used in the neural network constructions in Theorem C, can be numerically computed from samples of $f$ on the lattice $\{\boldsymbol{k}/n\}_{0\leq\boldsymbol{k}\leq n}$ without any other additional information about $f$, and is summarized in Algorithm 1.

---
**Algorithm 1** Binary Bernstein algorithm

---
**Require:** smoothness of the target function $s$, direction $\ell$ for $\Sigma\Delta$, samples $\{f(\frac{\boldsymbol{k}}{n})\}_{0\leq\boldsymbol{k}\leq n}$.
   1. Calculate $\{a_{\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n}$ defined to be $a_{\boldsymbol{k}} := f_{n,\lceil s/2\rceil}(\frac{\boldsymbol{k}}{n})$.
   2. Abort if $\|a\|_\infty \geq 1$.
   3. Apply $s$-th order $\Sigma\Delta$ in direction $\ell$ on $\{a_{\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n}$.
**Ensure:** One-bit coefficients $\{\sigma_{\boldsymbol{k}}\}_{0\leq\boldsymbol{k}\leq n}$ and approximant $\sum_{0\leq\boldsymbol{k}\leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(\boldsymbol{x})$.

---

Tracing through the proof of Theorem A, the first step of our approximation scheme is to compute the real coefficients

$$a_{\boldsymbol{k}} := f_{n,\lceil s/2\rceil}\Big(\frac{\boldsymbol{k}}{n}\Big) := \Big( \sum_{m=0}^{\lceil s/2\rceil-1} (I - B_n)^m(f) \Big)\Big(\frac{\boldsymbol{k}}{n}\Big).$$

We explain how to compute these from the samples $\{f(\boldsymbol{k}/n)\}_{0 \leq \boldsymbol{k} \leq n}$. Fix any bijection from $\{\boldsymbol{k} \colon 0 \leq \boldsymbol{k} \leq n\}$ to $\{1, \ldots, (n+1)^d\}$, such as the lexicographic ordering. Let $\boldsymbol{f} \in \mathbb{R}^{(n+1)^d}$ be the vector such that $\boldsymbol{f_k} = f(\boldsymbol{k}/n)$, where the subscript on $\boldsymbol{f}$ should be understood as the image of $\boldsymbol{k}$ under whichever ordering was selected. Let $P$ be the $(n+1)^d \times (n+1)^d$ matrix whose $(\boldsymbol{j}, \boldsymbol{k})$ entry is $p_{n,\boldsymbol{k}}(\boldsymbol{j}/n)$. It follows from a direct calculation that $B_n(f)(\boldsymbol{k}/n) = (P\boldsymbol{f})_{\boldsymbol{k}}$. Furthermore, it holds that for each $m \geq 1$,

$$B_n^m(f)\Big(\frac{\boldsymbol{k}}{n}\Big) = (P^m \boldsymbol{f})_{\boldsymbol{k}}.$$

This can be shown by induction since by linearity of the Bernstein operator,

$$B_n^{m+1}(f)\Big(\frac{\boldsymbol{k}}{n}\Big) = \sum_{0 \leq \boldsymbol{\ell} \leq n} B_n^m(f)\Big(\frac{\boldsymbol{\ell}}{n}\Big) p_{n,\boldsymbol{\ell}}\Big(\frac{\boldsymbol{k}}{n}\Big) = \sum_{0 \leq \boldsymbol{\ell} \leq n} (P^m \boldsymbol{f})_{\boldsymbol{\ell}} P_{\boldsymbol{k},\boldsymbol{\ell}} = (P^{m+1}\boldsymbol{f})_{\boldsymbol{k}}.$$

It follows from the above that

$$a_{\boldsymbol{k}} = f_{n,\lceil s/2 \rceil}\Big(\frac{\boldsymbol{k}}{n}\Big) = \Big( \sum_{m=0}^{\lceil s/2 \rceil - 1} (I - P)^m \boldsymbol{f} \Big)_{\boldsymbol{k}}.$$

Hence computation of $\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ amounts to matrix vector operations.

It follows from Theorem 2.5 that for sufficiently large $n$, we can guarantee that $\|a\|_\infty < 1$. From here, calculation of $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ just requires feeding them into a stable $s$-th order $\Sigma\Delta$ quantization scheme applied to the $\ell$-th direction, for any $\ell$ chosen beforehand. Since directional $\Sigma\Delta$ follows directly from its corresponding one-dimensional version, we drop the dependence on $\ell$ in this expository portion.

For the reader's convenience, let us fully describe the scheme introduced in [20]. Instead of solving (3.1) directly, one considers the equation (where $y, q, h, v$ are sequences indexed by $\mathbb{Z}$),

$$y_k - q_k = v_k - (h * v)_k, \quad \text{where} \quad (h * v)_k = \sum_{j=1}^{k} h_j v_{k-j}. \tag{5.1}$$

Let $\mu \in (0, 1)$, which will serve as an upper bound for $\|y\|_\infty$, and fix an integer $r \geq 1$. Pick any natural number $\gamma > 6$ such that $\mu \leq 2 - \cosh(\pi\gamma^{-1/2})$. Define the integers $z_k := \gamma(k-1)^2 + 1$ for $k = 1, \ldots, r$, and let $h$ be a sequence supported in $\{z_1, \ldots, z_r\}$ with $h_{z_k} = d_k$, where $\{d_1, \ldots, d_r\}$ are found as solutions to the Vandermonde system,

$$\begin{bmatrix} 1 & 1 & \cdots & 1 \\ z_1 & z_2 & \cdots & z_r \\ \vdots & \vdots & & \vdots \\ z_1^{r-1} & z_2^{r-1} & \cdots & z_r^{r-1} \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_r \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

Hence $h$ is readily computed by solving a linear system. Given input $y$, we compute the $\{\pm 1\}$ sequence $q$ recursively by

$$q_k := \text{sign}((h * v)_k + y_k), \quad \text{and} \quad v_k := y_k - q_k + (h * v)_k.$$

Here, we use the convention that $\text{sign}(0) = 1$. As for computations, this is enough since we are only interested in computing $q$, which serves as the $\{\sigma_{\boldsymbol{k}}\}_{k_\ell=0}^{n}$ for fixed $\boldsymbol{k}' = (k_1, \ldots, k_{\ell-1}, k_{\ell+1}, k_d)$.

For theoretical purposes, in order to control the quantization error, it was shown in the referenced paper that there exists an auxiliary sequence $g$ for which $v_k - (h * v)_k = \Delta^r(g * v)$, so (5.1)

is in fact a $\Sigma\Delta$ scheme as in (3.1) with $u = g * v$ instead. With this transformation in place, it was shown that

$$\|u\|_\infty \leq \frac{3}{\sqrt{2\pi r}}(\gamma e)^r r^r.$$

Hence, the right hand side term serves as the implicit constant $C_{\mu,r}$ that controls the stability of this particular quantization scheme.

Next, we examine the stability of our approximation scheme to perturbations of the input function $f$. Since the approximation method only depends on the samples $\{f(\boldsymbol{k}/n)\}_{0 \leq \boldsymbol{k} \leq n}$, we consider noisy samples of the form,

$$\widetilde{y}_{\boldsymbol{k}} = f\left(\frac{\boldsymbol{k}}{n}\right) + \eta_{\boldsymbol{k}},$$

where $\{\eta_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ represents any unknown perturbation. We measure the noise level in via the $\ell^\infty$ norm $\|\eta\|_\infty$. The following theorem evaluates the resulting approximation error if we use the perturbed samples of $f$ in Algorithm 1.

**Theorem 5.1.** *Let $s, d, n \geq 1$, $\varepsilon > 0$, and $f \in C^s([0,1]^d)$ such that*

$$\beta := \|f\|_\infty + \sqrt{2^{s+1}}\varepsilon + \frac{\sqrt{2^{s-1}}d}{8n}\|f\|_{C^1\mathrm{Lip}} < 1. \tag{5.2}$$

*For any $1 \leq \ell \leq d$ and $\{\eta_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ such that $\|\eta\|_\infty \leq \varepsilon$, let $\{\widetilde{\sigma}_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n} \subseteq \{\pm 1\}$ be the output of Algorithm 1 given inputs $\widetilde{y}_{\boldsymbol{k}} = f(\boldsymbol{k}/n) + \eta_{\boldsymbol{k}}$ for each $0 \leq \boldsymbol{k} \leq n$. Then we have*

$$\left| f(x) - \sum_{0 \leq \boldsymbol{k} \leq n} \widetilde{\sigma}_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(x) \right| \lesssim_{s,d,\beta} \|f\|_{C^s} \min\left(1, n^{-s/2} x_\ell^{-s}(1 - x_\ell)^{-s}\right) + \varepsilon.$$

*Proof.* Let $\{\widetilde{a}_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ be the real coefficients produced by the first step of Algorithm 1 given input $\{\widetilde{y}_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$. It will be helpful to produce a function $\widetilde{f}$ for which $\widetilde{f}(\boldsymbol{k}/n) = y_{\boldsymbol{k}}$. Let $\varphi$ be any $C^\infty$ function compactly supported in $[-\frac{1}{4}, \frac{1}{4}]^d$ such that $\varphi(0) = 1$ and $0 \leq \varphi \leq 1$. Define

$$\widetilde{f}(x) := f(x) + \sum_{0 \leq \boldsymbol{k} \leq n} \eta_{\boldsymbol{k}}\varphi(nx - \boldsymbol{k}),$$

so that $\widetilde{f} \in C^s([0,1]^d)$ and $\widetilde{f}(\boldsymbol{k}/n) = \widetilde{y}_{\boldsymbol{k}}$. Note that the set of functions $\{\varphi(n \cdot - \boldsymbol{k})\}_{0 \leq \boldsymbol{k} \leq n}$ have disjoint supports, and consequently,

$$\|\widetilde{f} - f\|_\infty = \left\| \sum_{0 \leq \boldsymbol{k} \leq n} \eta_{\boldsymbol{k}}\varphi(n \cdot - \boldsymbol{k}) \right\|_\infty \leq \|\eta\|_\infty \leq \varepsilon.$$

To simplify the following notation, let $r = \lceil s/2 \rceil$. Define the function

$$\widetilde{f}_{n,r} = \sum_{m=0}^{r-1} (I - B_n)^m(\widetilde{f}).$$

By Theorem 2.5, we have that $B_n(\widetilde{f}_{n,\lceil s/2 \rceil}) = U_{n,\lceil s/2 \rceil}(f)$ and

$$\begin{aligned}
\|\widetilde{f}_{n,r}\|_\infty &\leq \|\widetilde{f}\|_\infty + (2^{r-1} - 1)\|\widetilde{f} - B_n(\widetilde{f})\|_\infty \\
&\leq \|f\|_\infty + \|\eta\|_\infty + (2^{r-1} - 1)\big(\|\widetilde{f} - f\|_\infty + \|f - B_n(f)\|_\infty + \|B_n(f) - B_n(\widetilde{f})\|_\infty\big) \\
&\leq \|f\|_\infty + (2^r - 1)\|\eta\|_\infty + \frac{(2^{r-1} - 1)d}{8n}\|f\|_{C^1\mathrm{Lip}},
\end{aligned}$$

27

where the final inequality follows from Proposition 2.3.

Note that $\widetilde{a}_{\boldsymbol{k}} = \widetilde{f}_{n,r}(\boldsymbol{k}/n)$. Hence, under assumption (5.2), we see that $\|\widetilde{a}\|_\infty \leq \|\widetilde{f}_{n,\lceil s/2\rceil}\|_\infty \leq \beta < 1$. This permits us to employ a stable $s$-th order $\Sigma\Delta$ scheme in direction $\ell$ to generate signs $\{\widetilde{\sigma}_{\boldsymbol{k}}\}_{0\leq \boldsymbol{k}\leq n}$ from $\{\widetilde{a}_{\boldsymbol{k}}\}_{0\leq \boldsymbol{k}\leq n}$, which satisfy the difference equation $\widetilde{a} - \widetilde{\sigma} = (\Delta_\ell)^r \widetilde{u}$, for a bounded $\widetilde{u}$. We proceed to compare $\{\widetilde{a}_{\boldsymbol{k}}\}$ with $\{a_{\boldsymbol{k}}\}_{0\leq \boldsymbol{k}\leq n}$, which are coefficients generated by Algorithm 1 given noiseless samples. Note that $a_{\boldsymbol{k}} := f_{n,r}(\boldsymbol{k}/n)$ and

$$\left| a_{\boldsymbol{k}} - \widetilde{a}_{\boldsymbol{k}} \right| = \left| \widetilde{f}_{n,r}\left(\frac{\boldsymbol{k}}{n}\right) - f_{n,r}\left(\frac{\boldsymbol{k}}{n}\right) \right| \leq \Big\| \sum_{m=0}^{r-1} (I - B_n)^m (\widetilde{f} - f) \Big\|_\infty \leq 2^r \|\widetilde{f} - f\|_\infty \leq 2^r \varepsilon.$$

Using this bound on $\|a - \widetilde{a}\|_\infty$ and that the Bernstein polynomials are nonnegative and form a partition of unity, we have

$$\left| f(x) - \sum_{0\leq \boldsymbol{k}\leq n} \widetilde{\sigma}_{\boldsymbol{k}} p_{n,\boldsymbol{k}} \right| \leq \left| f(x) - \sum_{0\leq \boldsymbol{k}\leq n} a_{\boldsymbol{k}} p_{n,\boldsymbol{k}} \right| + \Big| \sum_{0\leq \boldsymbol{k}\leq n} \left((\Delta_\ell)^r \widetilde{u}\right)_{\boldsymbol{k}} p_{n,\boldsymbol{k}} \Big| + \Big| \sum_{0\leq \boldsymbol{k}\leq n} \left(a_{\boldsymbol{k}} - \widetilde{a}_{\boldsymbol{k}}\right) p_{n,\boldsymbol{k}} \Big|$$

$$\leq \left| f(x) - \sum_{0\leq \boldsymbol{k}\leq n} a_{\boldsymbol{k}} p_{n,\boldsymbol{k}} \right| + \Big| \sum_{0\leq \boldsymbol{k}\leq n} \left((\Delta_\ell)^r \widetilde{u}\right)_{\boldsymbol{k}} p_{n,\boldsymbol{k}} \Big| + 2^r \varepsilon.$$

Notice that the first term is the approximation error by iterated Bernstein operators and can be controlled using Theorem 2.6. The second term is the quantization error on the Bernstein basis and is upper bounded using Theorem 3.5. Doing so completes the proof. $\square$

Several comments about this theorem are in order. First, notice that the error bound in Theorem 5.1 is that of Theorem A plus a contribution from the noise. This theorem holds for arbitrary (hence deterministic and adversarial) perturbations and perhaps could be improved if additional assumptions are made, such as a statistical model. Moreover, the upper bound does not increase in $n$, which is not obvious since the noise energy $\|\eta\|_2$ may grown in $n$ without additional assumptions on the noise. Second, condition (5.2) ensures that the samples $\{\widetilde{f}_{n,r}(\boldsymbol{k}/n)\}$ have absolute value strictly less than 1, which is only used to guarantee that a $s$-th order directional $\Sigma\Delta$ scheme is stable. Such a condition is not always necessarily since it may be plausible that $\Sigma\Delta$ is stable for a larger class of sequences, especially those generated from uniformly sampling smooth functions which have significant correlation between samples. Third, notice that condition (5.2) becomes easier, not more difficult, to satisfy for fixed $s$ and increasing $n$. Again, it is not obvious that such a behavior is possible since the noise energy $\|\eta\|_2$ increases in $n$.

## 5.2 Bits and minimax code length

In this subsection, we discuss our main results on approximation by quantized neural networks in the context of codes. Any function that can be implemented by a strict quantized neural network can be identified by its parameters and topology. Due to Theorem C, every smooth function can be encoded with small loss of information, and then approximately reconstructed by mapping back to its associated network. We must first explain what we mean by the number of bits.

We define the minimax code length in an abstract setting before returning back to neural networks. Let $(X, d_X)$ be a metric space and $\mathcal{F} \subseteq X$ a class of functions. Suppose that for any $\varepsilon > 0$, there is an integer $B_\varepsilon \geq 0$ and maps $E\colon \mathcal{F} \to \{\pm 1\}^{B_\varepsilon}$ and $D\colon \{\pm 1\}^{B_\varepsilon} \to X$, for which

$$\sup_{f\in\mathcal{F}} \; d_X\big(f, D(E(f))\big) \leq \varepsilon.$$

28

In which case, we call $(E, D)$ as an $\varepsilon$-approximate encoder-decoder pair and $B_\varepsilon$ the number of bits required for this $\varepsilon$-approximate encoder-decoder pair. Usually the quantity of interest is the growth rate of $B_\varepsilon$ as $\varepsilon \to 0$. To examine the optimality of a given pair, the smallest $B_\varepsilon$ for which the above holds is defined to be the minimax code length,

$$B_\varepsilon^* := \min \left\{ B : \exists\, (E, D) \text{ for which } \sup_{f \in \mathcal{F}} d\big(f, D(E(f))\big) \leq \varepsilon \right\}.$$

Of course $B_\varepsilon \geq B_\varepsilon^*$, but to determine $B_\varepsilon^*$, it is well known that the minimax code length is precisely the smallest natural number that upper bounds the Kolmogorov $\varepsilon$-entropy of $\mathcal{F}$.

Let us now discuss how this is connected to quantized neural networks, in an abstract setting. Let $\mathcal{A}$ be finite, and for each $\varepsilon > 0$, suppose $\mathcal{N}_\varepsilon \subseteq X$ is a finite set of $\mathcal{A}$-quantized neural networks, where $\mathcal{A}$ is assumed to be finite, for which we have the approximation property that: for each $f \in \mathcal{F}$, there is a $g \in \mathcal{N}_\varepsilon$ such that $d_X(f, g) \leq \varepsilon$. Since $\mathcal{A}$ is finite and any $g \in \mathcal{N}_\varepsilon$ can be identified by its network parameters and topology, there is a $B_\varepsilon$ and injective $E \colon N_\varepsilon \to \{\pm 1\}^{B_\varepsilon}$. By definition, there is a decoder $D$ such that $D \circ E$ is the identity map on $\mathcal{N}_\varepsilon$. We extend $E$ to $\mathcal{F}$ by first mapping $f$ to an $\varepsilon$ approximation $g$ and then using the bit representation of $g$. Hence $(E, D)$ is an $\varepsilon$-approximate encoder-decoder pair, and we call $B_\varepsilon$ the number of bits used by the neural network.

When one refers to the "number of bits", perhaps an immediate inclination is to envision the number of nonzero bits required to store the weights of a single network into memory. This is usually the perspective taken if the goal is to quantize a particular network to reduce its memory cost, see [9]. The number of bits to specify a single neural network is (often significantly) smaller than the minimax code length, which we proceed to explain below.

For each $\mathcal{A}$-quantized neural network $g \in \mathcal{N}_\varepsilon$, we can list out the parameters and topology of the network as a finite sequence $y(q)$ in $\mathcal{A} \cup \{0\}$, where '0' is used to specify that a weight or bias is not being used from one node to another. Since $\mathcal{A}$ is a finite set, it can be encoded with $\lceil \log |\mathcal{A}| \rceil$ bits. Discarding each 0 in $y$ and replacing each nonzero term in $y$ with its corresponding $\pm 1$ representation, we obtain a finite sequence $q(g) \in \{\pm 1\}^{B_\varepsilon(g)}$. This quantity $B_\varepsilon(g)$ is sometimes referred to as the number of bits required to store the network $g$. However, if we compare $B_\varepsilon^\circ := \sup_{g \in \mathcal{N}_\varepsilon} B_\varepsilon(g)$ to the number of bits $B_\varepsilon$, it is possible that $B_\varepsilon^\circ = o(B_\varepsilon)$ as $\varepsilon \to 0$. The short explanation is that the map $\mathcal{N}_\varepsilon \mapsto \{\pm 1\}^{B_\varepsilon^\circ}$ given as $g \mapsto q(g)$ is not necessarily injective because '0' is required to specify the network's topology, which has been omitted from the quantity $B_\varepsilon(g)$. In this case, if $\mathcal{N}_\varepsilon$ contains networks with vastly different topologies, then additional bits may be required to distinguish between the networks' topologies and it may not be feasible to simply discard all zeros.

Going back to the content of this paper, our strategy is to approximate $f \in C^s([0,1]^d)$ by a one-bit linear combination of (possibly approximate) Bernstein polynomials. The latter set only depends on the function class and prescribed error, and not on a particular $f$. Hence the networks used to approximate $C^s([0,1]^d)$ all have the same topology, and for this reason, we do not need to include '0' as a bit. The following theorem quantifies the number of bits.

**Theorem 5.2.** *For any $s, d \geq 1$ and sufficiently small $\mu > 0$, the following hold. There exist an encoder $E \colon C^s([0,1]^d) \to \{\pm 1\}^B$ and decoder $D \colon \{\pm 1\}^B \to \mathcal{N}(\{\pm 1\}, \rho, L, N, P)$ such that $L = O(\varepsilon^{-2/s})$ and $\max(B, N, P) = O(\varepsilon^{-2d/s})$ as $\varepsilon \to 0$, and*

$$\sup_{\|f\|_\infty \leq \mu,\, \|f\|_{C^s} \leq 1} \|f - D(E(f))\|_\infty \leq \varepsilon.$$

*Proof.* For any $f \in C^s([0,1]^d)$, Whitney's extension theorem provides us with a $F \in C^s(\mathbb{R}^d)$ such

that $F = f$ on $[0,1]^d$ and $\|F\|_{C^s(\mathbb{R}^d)} \lesssim_{d,s} \|f\|_{C^s([0,1]^d)}$, where importantly, the implicit constant that appears in this inequality does not depend on $f$, see [38] and also [36, Chapter 6, Theorem 4].

Consider the set $U := [-\frac{1}{2}, \frac{3}{2}] \times [0,1]^{d-1}$. We will only consider Whitney extensions $f$ of that are compactly supported in $U$ – this is possible since given any Whitney extension of $f$, which might not necessarily be compactly supported, we can multiply it by a smooth function that is compactly supported in $U$ and identically equal to 1 on $[0,1]^d$.

Hence, let $W_0$ and $W_2$ be the best possible constants for which $\|F\|_{L^\infty(E)} \leq W_0\|f\|_\infty$ and $\|F\|_{C^2(E)} \leq W_2\|f\|_{C^2}$, where $F$ is a Whitney extension of $f$ that is compactly supported in $U$. Note that $W_0$ and $W_2$ only depend on $d, s, U$.

For now, fix a $f \in C^s([0,1]^d)$ such that $\|f\|_\infty \leq \mu$ and $\|f\|_{C^s} \leq 1$. Let $F \in C^s(\mathbb{R}^d)$ be a Whitney extension of $f$ that is compactly supported in $U$. Consider the function $\widetilde{F} \in C^s([0,1]^d)$ defined as $\widetilde{F}(\boldsymbol{x}) := F(2x_1 - \frac{1}{2}, x_2, \ldots, x_d)$. Notice that $\|\widetilde{F}\|_\infty = \|F\|_\infty \leq W_0\|f\|_\infty \leq W_0\mu$ and that $\|\widetilde{F}\|_{C^2([0,1]^d)} \leq 4\|F\|_{C^2} \leq 4W_2\|f\|_{C^2}$.

From here onward, assume that $\mu$ is sufficiently small so that $W_0\mu \leq 1/2$, and note that $\mu$ only needs to be sufficiently small depending on $s, d$ since $W_0$ only depends on these quantities. For any integer $n \geq 2sd^2W_2$, so that $n \geq sd^2W_2/(1 - W_0\mu)$, we apply Theorem A to $\widetilde{F}$ with $\ell = 1$. Hence, there exists $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n} \subseteq \{\pm 1\}$ such that for all $\boldsymbol{x} \in [0,1]^d$ with $\frac{1}{4} \leq x_1 \leq \frac{3}{4}$, we have

$$\left|\widetilde{F}(\boldsymbol{x}) - \widetilde{H}(\boldsymbol{x})\right| \lesssim_{s,d} n^{-s/2}, \quad \text{where} \quad \widetilde{H}(\boldsymbol{x}) := \sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(\boldsymbol{x}). \tag{5.3}$$

Define $H \colon [0,1]^d \to \mathbb{R}$ by $H(\boldsymbol{x}) := \widetilde{H}(\frac{1}{2}x_1 + \frac{1}{4}, x_2, \ldots, x_d)$. For each $\boldsymbol{x} \in [0,1]^d$,

$$|f(\boldsymbol{x}) - H(\boldsymbol{x})| = \left|\widetilde{F}\left(\frac{1}{2}x_1 + \frac{1}{4}, x_2, \ldots, x_d\right) - \widetilde{H}\left(\frac{1}{2}x_1 + \frac{1}{4}, x_2, \ldots, x_d\right)\right| \lesssim_{s,d} n^{-s/2}, \tag{5.4}$$

where the final inequality follows from (5.3). Hence, for any $\varepsilon$ sufficiently small, by making $n$ large enough depending only on $\varepsilon$, $d$, and $s$, the right hand side of (5.4) can be made smaller than $\varepsilon$. We see that $n = O(\varepsilon^{-2/s})$ as $\varepsilon \to 0$ suffices.

We next show that $H$ can be implemented by a neural network and we determine the number of bits used in this encoding. Recall that

$$H(\boldsymbol{x}) = \sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}\left(\frac{1}{2}x_1 + \frac{1}{4}, x_2, \ldots, x_d\right).$$

The function $x_1 \mapsto \frac{1}{2}x_1 + \frac{1}{4}$ can be made by a neural of size $(4, 4, 9)$ because we first make $\frac{1}{2} = \rho(0 \cdot x_1 + 1)$, then $\frac{1}{4} = \rho(\frac{1}{2})$, then $\rho(x_1 + \frac{1}{4})$ and $\rho(x_1 - \frac{1}{4})$, and finally $\frac{1}{2}x_1 + \frac{1}{4} = \rho(x_1 + \frac{1}{4}) - \rho(x_1 - \frac{1}{4}) + \frac{1}{4}$. As shown in Theorem B, the multivariate Bernstein polynomials $\{p_{n,\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ are implementable by a $\{\pm 1\}$-quantized quadratic neural network that does not depend on the target function $f$, and it has $O(n)$ layers and $O(n^d)$ nodes and parameters as $n \to \infty$. Hence, we can implement the functions

$$\left\{\boldsymbol{x} \mapsto p_{n,\boldsymbol{k}}\left(\frac{1}{2}x_1 + \frac{1}{4}, x_2, \ldots, x_d\right)\right\}_{0 \leq \boldsymbol{k} \leq n}$$

using a $\{\pm 1\}$-quantized quadratic network that only depends on $n$ (which is selected in terms of $s, d$) and not on $f$. Hence the number of bits to encode this part is $O(\frac{2d}{s} \log(1/\varepsilon))$. By incorporating a linear last layer of size $(1, 1, (n+1)^d)$, whose weights are specified by $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$, we see that $H$ is implementable. The signs $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ depend on $f$ and $\sigma_{\boldsymbol{k}} \in \{\pm 1\}$, so the total number of bits used to encode this part is $(n+1)^d = O(\varepsilon^{-2d/s})$ as $\varepsilon \to 0$.

$\square$

To put Theorem 5.2 in context, Kolmogorov entropy tells us that the minimum number of bits necessary to encode the unit ball of $C^s([0,1]^d)$ with error at most $\varepsilon$ measured in the uniform norm is $O(\varepsilon^{-d/s})$, regardless of which encode-decoder pair is used, see [35, Chapter 7.5] and [28] for the one-dimensional version and [37, page 86] for a full proof the multidimensional case. Although our approximation strategy via Bernstein polynomials does not attain the entropy rate, it has several other desirable features that are not captured by entropy considerations. Our method only uses a one-bit alphabet (which is not enforced by bit counting considerations), can be readily computed from queries, and is stable to perturbations of the unknown function.

## 5.3   Beyond directional $\Sigma\Delta$?

One unsatisfying aspect of the quantization schemes employed and analyzed in Section 3 is that they are inherently one-dimensional. It is due to this that the quantization error bound in Theorem 3.5 contains the $x_\ell^{-r}(1 - x_\ell)^{-r}$ term which blows up at the boundaries $x_\ell = 0$ and $x_\ell = 1$. This term is then propagated to Theorems A and C.

On one hand, the inability to approximate the target function at the endpoints is unavoidable. To see why, $p_{n,k}(0) = \delta_{0,k}$ and $p_{n,k}(1) = \delta_{n,k}$, and consequently, if $\sigma_k \in \{\pm 1\}$ for each $k$, then $\sum_{k=0}^n \sigma_k p_{n,k}$ can only be equal to $\pm 1$ at the endpoints. This means that a $\{\pm 1\}$ linear combination of Bernstein polynomials cannot possibly approximate any continuous function uniformly on $[0,1]$, without additional assumptions on the target function near the endpoints.

On the other hand, one may wonder if it is possible to spread the error out to other faces of the $d$-dimensional cube, instead of concentrating all it to the faces where $x_\ell = 0$ and $x_\ell = 1$. More precisely, perhaps one could replace the term $x_\ell^{-r}(1 - x_\ell)^{-r}$ with

$$\boldsymbol{x}^{-\boldsymbol{r}}(1 - \boldsymbol{x})^{-\boldsymbol{r}} = \prod_{\ell=1}^d x_\ell^{-r_\ell}(1 - x_\ell)^{-r_\ell}, \quad \text{for any} \quad |\boldsymbol{r}| = r.$$

The main bottleneck of carrying out the more general case stems from difficulties with $\Sigma\Delta$ quantization. Currently, there is no such available stable one-bit multidimensional higher order $\Sigma\Delta$ scheme, beyond the directional one used in this paper. More precisely, we define a stable one-bit $\boldsymbol{r}$-th order $\Sigma\Delta$ quantizer with alphabet $\mathcal{A}$ to be a map that takes any sequence $y$ with $\|y\|_\infty \leq \mu < 1$ and outputs $q$, such that $q_{\boldsymbol{k}} \in \mathcal{A}$ and there is an associated $u$ with $\|u\|_\infty \leq C_{\mu,\boldsymbol{r},\mathcal{A}}$ (that does not depend on $y$) satisfying the equation

$$y - q = \Delta^{\boldsymbol{r}} u, \quad \text{where} \quad \Delta^{\boldsymbol{r}} = \Delta_d^{r_d} \Delta_{d-1}^{r_{d-1}} \cdots \Delta_1^{r_1}.$$

We say this scheme is truly multidimensional if there are $j \neq k$ such that $r_j, r_k \geq 1$.

If a stable $\boldsymbol{r}$-th order $\Sigma\Delta$ scheme exists, then it is straightforward to modify the analysis given in the proof of Theorem 3.5. Indeed, let $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ with $\sigma_{\boldsymbol{k}} \in \mathcal{A}$ be the output given input $\{a_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$. Since $a - \sigma = \Delta^{\boldsymbol{r}} u$ for some $\|u\|_\infty \leq C_{\mu,\boldsymbol{r},\mathcal{A}}$, a summation by parts yields

$$\left| \sum_{0 \leq \boldsymbol{k} \leq n} (a_{\boldsymbol{k}} - \sigma_{\boldsymbol{k}}) \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| = \left| \sum_{0 \leq \boldsymbol{k} \leq n} (\Delta^{\boldsymbol{r}} u)_{\boldsymbol{k}} \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| \lesssim_{\mu,\boldsymbol{r},\mathcal{A}} \sum_{0 \leq \boldsymbol{k} \leq n} \left| ((\Delta^{\boldsymbol{r}})^* p_{n,\cdot}(\boldsymbol{x}))_{\boldsymbol{k}} \right|.$$

Note that $(\Delta^{\boldsymbol{r}})^* = (\Delta_d^*)^{r_d} (\Delta_{d-1}^*)^{r_{d-1}} (\Delta_1^*)^{r_1}$. The Bernstein polynomials are tensor products, so

$$\sum_{0 \leq \boldsymbol{k} \leq n} \left| ((\Delta^{\boldsymbol{r}})^* p_{n,\cdot}(\boldsymbol{x}))_{\boldsymbol{k}} \right| = \prod_{\ell=1}^d \sum_{k_\ell=0}^n \left| ((\Delta^*)^{r_\ell} p_{n,\cdot}(x_\ell))_{k_\ell} \right| \lesssim_{r_1,\ldots,r_d} \prod_{\ell=1}^d n^{-r_\ell} x_\ell^{-r_\ell}(1 - x_\ell)^{-r_\ell},$$

31

where we used the one dimensional result in [22, Theorem 6]. Together with the trivial estimate that the quantization error is bounded by 2, we have proved that for all $\boldsymbol{x} \in [0,1]^d$,

$$\left| \sum_{0 \leq \boldsymbol{k} \leq n} (a_{\boldsymbol{k}} - \sigma_{\boldsymbol{k}}) \, p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| \lesssim_{\mu,\boldsymbol{r},\mathcal{A}} \min\left(1, n^{-|\boldsymbol{r}|/2} \boldsymbol{x}^{-\boldsymbol{r}}(1-\boldsymbol{x})^{-\boldsymbol{r}}\right).$$

With this at hand, Theorem Theorem A can be modified as follows. Under the same assumptions on $d, s, n, \mu, f$, for any $\boldsymbol{s}$ with $|\boldsymbol{s}| = s$, there exist $\{\sigma_{\boldsymbol{k}}\}_{0 \leq \boldsymbol{k} \leq n}$ with $\sigma_{\boldsymbol{k}} \in \mathcal{A}$ such that

$$\left| f(\boldsymbol{x}) - \sum_{0 \leq \boldsymbol{k} \leq n} \sigma_{\boldsymbol{k}} p_{n,\boldsymbol{k}}(\boldsymbol{x}) \right| \lesssim_{\boldsymbol{s},d,\mu,\mathcal{A}} \|f\|_{C^s} \min\left(1, n^{-s/2} \boldsymbol{x}^{-\boldsymbol{s}}(1-\boldsymbol{x})^{-\boldsymbol{s}}\right).$$

In the second step of Algorithm 1, one can use a $\boldsymbol{s}$-th order $\Sigma\Delta$ scheme instead of a directional one. It is important to emphasize again that there is no known stable one-bit $\Sigma\Delta$ quantizer, so this discussion cannot be extended to the one-bit case.

# A    Neural network constructions

This appendix constructs the desired one-bit neural networks. In the subsequent proofs, we just use the term "network" since the alphabet and activation do not vary in each subsection. To visualize the topology of network, we typically draw a schematic diagram. Nodes belonging to the same layer are placed horizontally, layers are stacked vertically, with the input nodes on top and output nodes on bottom. We say the $k$-th node in layer $\ell - 1$ is connected to the $j$-th node in layer $\ell$ if either $(A_\ell)_{j,k}$ or $(\boldsymbol{b}_\ell)_j$ is nonzero. Two connected nodes are represented in a schematic diagram by a line segment adjoining them.

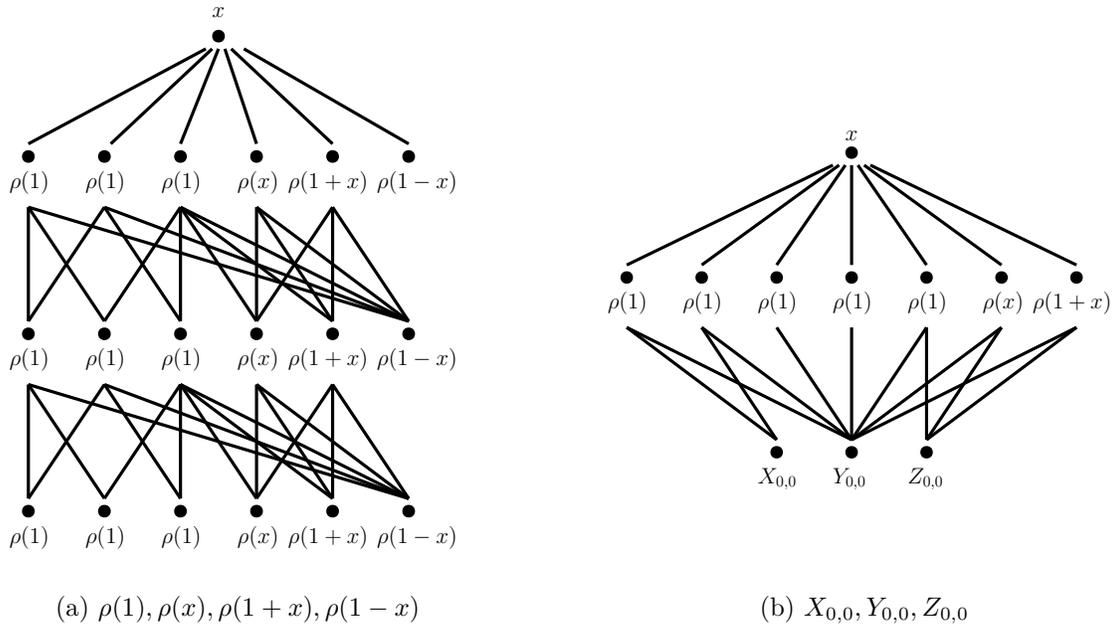## A.1    Implementation of $\{\pm 1\}$-quantized quadratic networks



(a) $\rho(1), \rho(x), \rho(1+x), \rho(1-x)$

(b) $X_{0,0}, Y_{0,0}, Z_{0,0}$

Figure 3: Helper activated one-bit quadratic networks

32

**Lemma A.1.** *For any integer $n \geq 1$, there exists an activated $\mathcal{A}_1$-quantized quadratic neural network with $n+1$ layers and $O(n)$ nodes and parameters that implements the function*

$$x \mapsto \Big( U(x), V(x), \{X_{n-1,k}(x)\}_{k=0}^{n-1}, \{Y_{n-1,k}(x)\}_{k=0}^{n-1}, \{Z_{n-1,k}(x)\}_{k=0}^{n-1} \Big).$$

*Proof.* We will create an activated network such that $\{X_{m,k}\}_{k=0}^{m}$, $\{Y_{m,k}\}_{k=0}^{m}$ and $\{Z_{m,k}\}_{k=0}^{m}$ are outputs in layer $m+2$. Throughout this proof, we will implicitly use the identity $x = \rho(x+1) - \rho(x) - \rho(1)$ without explicit mention. To carry out this strategy, notice that (4.1) and (4.3) imply

$$X_{m,k} = \begin{cases} \rho(Y_{m-1,0} - U - X_{m-1,0}) & \text{if } k = 0, \\ \rho(Z_{m-1,k-1} - V - X_{m-1,k-1} + Y_{m-1,k} - U - X_{m-1,k}) & \text{if } 0 < k < m, \\ \rho(Z_{m-1,m-1} - V - X_{m-1,m-1}) & \text{if } k = m. \end{cases}$$

$$Y_{m,k} = \begin{cases} \rho(1 - x + Y_{m-1,0} - U - X_{m-1,0}) & \text{if } k = 0, \\ \rho(1 - x + Z_{m-1,k-1} - V - X_{m-1,k-1} + Y_{m-1,k} - U - X_{m-1,k}) & \text{if } 0 < k < m, \quad \text{(A.1)} \\ \rho(1 - x + Z_{m-1,m-1} - V - X_{m-1,m-1}) & \text{if } k = m. \end{cases}$$

$$Z_{m,k} = \begin{cases} \rho(x + Y_{m-1,0} - U - X_{m-1,0}) & \text{if } k = 0, \\ \rho(x + Z_{m-1,k-1} - V - X_{m-1,k-1} + Y_{m-1,k} - U - X_{m-1,k}) & \text{if } 0 < k < m, \\ \rho(x + Z_{m-1,m-1} - V - X_{m-1,m-1}) & \text{if } k = m. \end{cases}$$

To employ the recurrence (A.1), we first construct a network with $n+1$ layers such that each layer has six nodes that output $\rho(1)$, $\rho(1)$, $\rho(1)$, $\rho(x)$, $\rho(1+x)$, and $\rho(1-x)$. For layer 1, these six outputs are readily made from the input $x$. The same six outputs can be generated by $\rho$ applied to $\pm 1$ linear combinations of the same terms, since

$$\rho(1) = \rho(\rho(1) + \rho(1)), \quad \rho(1-x) = \rho(\rho(1) + \rho(1) + \rho(1) - \rho(1+x) + \rho(x)),$$
$$\rho(x) = \rho(\rho(1+x) - \rho(x) - \rho(1)), \quad \rho(1+x) = \rho(\rho(1+x) - \rho(x) + \rho(1)).$$

This is shown in Figure 3a. Since a constant number of nodes and parameters are added to increase the depth by one, terminating at layer $n+1$, this network has size $O(n)$. This establishes that $U, V$ are outputs of layer $n+1$.

We can produce $X_{0,0}, Y_{0,0}, Z_{0,0}$ with a network of size $(2, 10, 20)$, as shown in Figure 3b, since

$$X_{0,0} = \rho(1) = \rho(\rho(1) + \rho(1)),$$
$$Y_{0,0} = \rho(1 - x + 1) = \rho(\rho(1) + \rho(1) - \rho(x+1) + \rho(x) + \rho(1) + \rho(1) + \rho(1)),$$
$$Z_{0,0} = \rho(x + 1) = \rho(\rho(x+1) - \rho(x) + \rho(1)).$$

It follows from recurrence (A.1) that if layer $m$ outputs the quantities $\{X_{m-1,k}\}_{k=0}^{m-1}$, $\{Y_{m-1,k}\}_{k=0}^{m-1}$, $\{Z_{m-1,k}\}_{k=0}^{m-1}$ and $\rho(1), \rho(1), \rho(1), \rho(x), \rho(1+x)$, then it is possible to generate the quantities $\{X_{m,k}\}_{k=0}^{m}$, $\{Y_{m,k}\}_{k=0}^{m}$, $\{Z_{m,k}\}_{k=0}^{m}$ in the subsequent layer $m+1$ by adding only a constant number of nodes and parameters. $\square$

**Lemma A.2.** *For any $d \geq 2$, there exist an integer $d_*$ and an activated $\mathcal{A}_1$-quantized quadratic neural network with $O(\log d)$ layers and $O(d)$ nodes and parameters, as $d \to \infty$, that implements the function*

$$\boldsymbol{x} = (x_1, \ldots, x_d) \mapsto \Big( \rho\Big( \prod_{j \leq d_*} x_j + \prod_{j > d_*} x_j \Big), \rho\Big( \prod_{j \leq d_*} x_j \Big), \rho\Big( \prod_{j > d_*} x_j \Big) \Big).$$

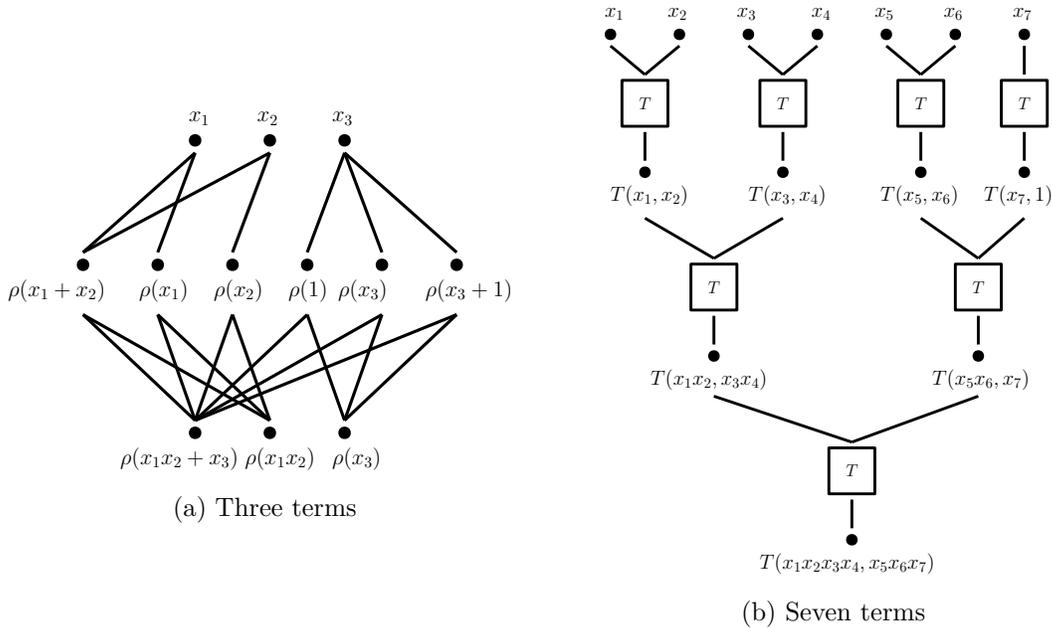33

(a) Three terms

(b) Seven terms

Figure 4: Pseudo-multiplication networks

*Proof.* The $d = 2$ case is handled by (4.2) whereby $d_* = 1$. For $d = 3$, we create a two layer network of size $(2,9,20)$ whose output is $(\rho(x_1 x_2 + x_3), \rho(x_1 x_2), \rho(x_3))$, as shown in Figure 4a. This is possible because $x_1 x_2 = \rho(x_1 + x_2) - \rho(x_1) - \rho(x_2)$ and $x_3 = \rho(x_3 + 1) - \rho(x_3) - \rho(1)$, which are $\pm 1$ linear combinations of outputs from the first layer. Hence, $d_* = 2$ when $d = 3$.

From now on, assume that $d \geq 4$. The main idea is that each layer, we generate the three terms on the right hand side (4.2) for $\lfloor d/2 \rfloor$ pairs of products, and then iterate the identity, which finally terminates after $O(\log d)$ iterations. Each step of the iteration depends on the number theoretic properties of $d$. The basic observation that allows for termination of this argument is that: if $abcd = \prod_{j=1}^{d} x_j$ and there is a layer that outputs

$$\rho(a), \rho(a + b), \rho(b), \rho(c), \rho(c + d), \rho(d),$$

then by appending another network of size $(1, 3, 12)$, we can implement $\rho(ab)$, $\rho(ab + cd)$, and $\rho(cd)$, which would complete the proof. To simplify the resulting argument, we introduce the shorthand notation

$$T(a, b) := (\rho(a), \rho(a + b), \rho(b)).$$

Each step of the reduction depends on the factorization of $d$. For layer 1, if $d$ is even, we create $3d/2$ nodes whose outputs are

$$T(x_1, x_2), T(x_3, x_4), \ldots, T(x_{d-1}, x_d).$$

If $d$ is odd, in layer 1, we create $3\lceil d/2 \rceil + 1$ nodes whose outputs are

$$T(x_1, x_2), T(x_3, x_4), \ldots, T(x_{d-2}, x_{d-1}), T(x_d, 1), \rho(1).$$

Let $d_1 := \lceil d/2 \rceil$, and we can assume that $d_1 \geq 4$, otherwise we are finished by the above basic operation. The construction of layer 2 depends on both the parities of $d_1$ and $d$. If both $d$ and $d_1$ are even, we use (4.2) again to create $3d_1/2$ nodes whose outputs are

$$T(x_1 x_2, x_3 x_4), \ldots, T(x_{d-3} x_{d-2}, x_{d-1} x_d).$$

34

The remaining three cases are fairly similar except for when both $d$ and $d_1$ are odd. In which case, we can simply reproduce $T(x_d, 1), \rho(1)$ in the next layer from the $T(x_d, 1), \rho(1)$ nodes in the previous layer. Indeed, in layer 2, we have nodes that output

$$T(x_1 x_2, x_3 x_4), \ldots, T(x_{d-4} x_{d-3}, x_{d-2} x_{d-1}), T(x_d, 1), \rho(1).$$

We continue this strategy until we end with the situation expressed in the basic observation. Notice that there are at most $O(\log d)$ iterations, and each iteration requires appending a network with only a single layer. Furthermore, the number of nodes and parameters used in layer $\ell$ is $O(d/2^\ell)$, so there are at most $O(d)$ many nodes and parameters in total. □

Figure 4b displays the corresponding network constructed in Lemma A.2 for $d = 7$ whereby $d_* = 4$. The next theorem shows that any $\pm 1$ linear combination of multivariate Bernstein polynomials is implementable by a neural network. We remark that the following does not claim that the Bernstein polynomials are implementable by an activated neural network.

## A.2 Implementation of $\{\pm\frac{1}{2}\}$-quantized ReLU networks

We begin with some basic properties of the ReLU function, which we will use without explicit reference. For any $t \geq 0$, $s \in \mathbb{R}$, and $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n$, it holds that

$$\sigma(t\boldsymbol{x}) = t\sigma(\boldsymbol{x}), \quad \sigma(\sigma(\boldsymbol{x})) = \sigma(\boldsymbol{x}), \quad \sigma(\boldsymbol{x} + \boldsymbol{y}) \leq \sigma(\boldsymbol{x}) + \sigma(\boldsymbol{y}), \quad \text{and} \quad |\sigma(s) - t| \leq |s - t|.$$

We next describe a few functions that are implementable by an activated $\mathcal{A}_{1/2}$-quantized ReLU network. For any input $x \in \mathbb{R}$, the number $1/2$ can be implemented by a network of size $(1, 1, 1)$ since $1/2 = \sigma(0x + 1/2)$. Implementation of the map $x \mapsto 2^{-m}\sigma(x)$ for any integer $m \geq 1$ can be achieved by a network of size $(m, m, m)$ by

$$x \mapsto \sigma\left(\frac{1}{2}\sigma\left(\frac{1}{2}\sigma\left(\cdots\sigma\left(\frac{1}{2}x\right)\right)\right)\right) = \frac{1}{2^m}\sigma(x).$$

Summation of nonnegative numbers can be carried out by a $\mathcal{A}_{1/2}$-quantized network of size $(2, 9, 16)$ because for nonnegative $a, b$, we have

$$a + b = \sigma\left(\frac{1}{2}\sigma\left(\frac{1}{2}a\right) + \cdots + \frac{1}{2}\sigma\left(\frac{1}{2}a\right) + \frac{1}{2}\sigma\left(\frac{1}{2}b\right) + \cdots + \frac{1}{2}\sigma\left(\frac{1}{2}b\right)\right).$$

Likewise, for nonnegative $a, b$, the quantity $\sigma(a - b)$ is implementable as well.

For any integer $L \geq 2$, there is a $L$ layer activated $\mathcal{A}_{1/2}$-quantized ReLU network $\zeta_L$ with $O(L)$ nodes and parameters that implements the map $x \mapsto \sigma(x)$. For $L = 2$ and $L = 3$, this can be done via the two networks shown in Figure 5a, which have size $(2, 5, 8)$ and $(3, 7, 14)$ respectively. For any $L \geq 4$, we can compose these maps so that the resulting $\zeta_L$ network has $L$ layers and $O(L)$ nodes and parameters. We call $\zeta_L$ a $L$ layer duplication network, and will be used to propagate nonnegative numbers down an arbitrary number of layers without the use of skip connections.

The tent function $\phi$ can be implemented by a ReLU network with real parameters since $\phi(x) = \sigma(2x) - 2\sigma(2x - 1)$. By making some adjustments, it is not difficult to see that $\phi$ can be implemented

35

(a) Duplication networks

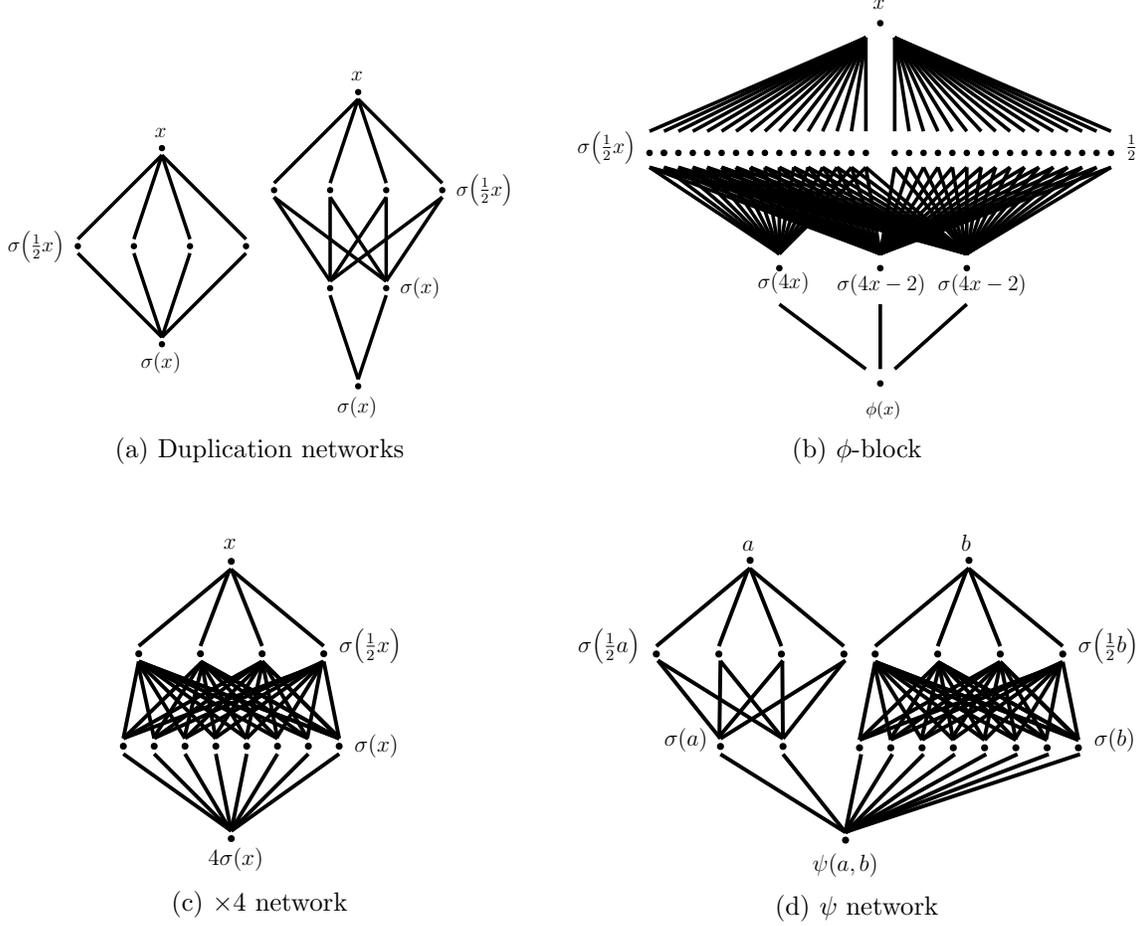(b) $\phi$-block



(c) $\times 4$ network

(d) $\psi$ network

Figure 5: Implementations of basic functions

with an activated $\mathcal{A}_{1/2}$-quantized network of size $(3, 36, 107)$ because of the identities

$$\phi(x) = \sigma\Big(\frac{1}{2}\sigma(4x) - \frac{1}{2}\sigma(4x-2) - \frac{1}{2}\sigma(4x-2)\Big),$$

$$\sigma(4x) = \sigma\Big(\underbrace{\frac{1}{2}\sigma\Big(\frac{1}{2}x\Big) + \cdots + \frac{1}{2}\sigma\Big(\frac{1}{2}x\Big)}_{16 \text{ times}}\Big),$$

$$\sigma(4x-2) = \sigma\Big(\underbrace{\frac{1}{2}\sigma\Big(\frac{1}{2}x\Big) + \cdots + \frac{1}{2}\sigma\Big(\frac{1}{2}x\Big)}_{16 \text{ times}} - \underbrace{\frac{1}{2}\cdot\frac{1}{2} - \cdots - \frac{1}{2}\cdot\frac{1}{2}}_{8 \text{ times}}\Big).$$

For convenience, we call this network a $\phi$-block and it is shown in Fig. 5b.

**Lemma A.3.** *For any $\varepsilon > 0$, there exist nonnegative functions $S_\varepsilon^+, S_\varepsilon^- : [0,1] \to \mathbb{R}$ such that*

$$\text{for all } x \in [0,1], \quad S_\varepsilon^-(x) \le x^2 \le S_\varepsilon^+(x), \quad |S_\varepsilon^+(x) - x^2| \le \varepsilon, \quad |S_\varepsilon^-(x) - x^2| \le \varepsilon.$$

*Furthermore, there exist two activated $\mathcal{A}_{1/2}$-quantized ReLU neural network both with the same number of layers and each of size $O(\log(1/\varepsilon))$ that implement $S_\varepsilon^+$ and $S_\varepsilon^-$.*

*Proof.* Fix $\varepsilon > 0$ and $m \geq 1$ will be an integer chosen later depending only on $\varepsilon$. We define

$$S_\varepsilon^+(x) := \sigma\left(x - \sum_{k=1}^m \frac{\phi^{\circ k}(x)}{4^k}\right) \quad \text{and} \quad S_\varepsilon^-(x) := \sigma\left(S_\varepsilon^+(x) - \frac{1}{2}\frac{1}{4^m}\right).$$

Since $\phi^{\circ k}$ is nonnegative for all $k \geq 1$, it follows from (4.4) that $S_\varepsilon^+(x) \geq x^2$ and

$$|S_\varepsilon^+(x) - x^2| = S_\varepsilon^+(x) - x^2 = \sum_{k=m+1}^\infty \frac{\phi^{\circ k}(x)}{4^k} \leq \sum_{k=m+1}^\infty \frac{1}{4^k} = \frac{1}{3}\frac{1}{4^m}.$$

We first show that $S_\varepsilon^-(x) \leq x^2$. This trivially holds if $S_\varepsilon^+(x) \leq 1/(2 \cdot 4^m)$. Otherwise, for $x \in [0,1]$ such that $S_\varepsilon^+(x) \geq 1/(2 \cdot 4^m)$, we have

$$S_\varepsilon^-(x) = S_\varepsilon^+(x) - \frac{1}{2}\frac{1}{4^m} = x^2 + \sum_{k=m+1}^\infty \frac{\phi^{\circ k}(x)}{4^k} - \frac{1}{2}\frac{1}{4^m} \leq x^2,$$

where the last inequality follows from the observation that

$$\sum_{k=m+1}^\infty \frac{\phi^{\circ k}(x)}{4^k} \leq \sum_{k=m+1}^\infty \frac{\|\phi^{\circ k}\|_\infty}{4^k} \leq \sum_{k=m+1}^\infty \frac{1}{4^k} = \frac{1}{3}\frac{1}{4^m} < \frac{1}{2}\frac{1}{4^m}.$$

Moreover, since $x^2 \geq 0$, we have that

$$|S_\varepsilon^-(x) - x^2| = \left|\sigma\left(S_\varepsilon^+(x) - \frac{1}{2}\frac{1}{4^m}\right) - x^2\right| \leq |S_\varepsilon^+(x) - x^2| + \frac{1}{2}\frac{1}{4^m} \leq \frac{1}{4^m}.$$

Thus, we pick any integer $m \geq \log(1/\varepsilon)/2$.

We proceed to discuss implementations of $S_\varepsilon^+$ and $S_\varepsilon^-$ by activated $\mathcal{A}_{1/2}$-quantized ReLU neural networks, as shown in Figure 6. The network for $S_\varepsilon^+$ consists of only the black nodes and weights, and has a two column structure. The network for $S_\varepsilon^-$ consists of all the nodes and weights, and has a three column structure.

We focus on $S_\varepsilon^+$ first. The left column nodes implement the functions $\{\phi^{\circ k}\}_{k=1}^m$, which can be done by composing $m$ many $\phi$-blocks. Each $\phi^{\circ k}$ is produced in the $3k$-th layer, so the left column of the network has $3m$ layers and $O(m)$ nodes and parameters. The layer 3 right node outputs $4x$, which can be done by using the top network shown in Figure 5c, which has size $(3, 12, 44)$. Recall that the layer 3 left node outputs $\phi(x)$. We can implement the function $\psi(a,b) := \sigma(4b - a)$ on nonnegative $a, b$ via a network of size $(3, 19, 58)$, which is the bottom network in Figure 5d. If $4b - a \geq 0$, then $\psi(a,b) = 4b - a$. Using this $\psi$ network, and that $16x - \phi(x) \geq 0$ by identity (4.4), we see that the layer 6 right node produces $16x - \phi(x)$. Now we proceed a similar fashion. For $k = 2, \ldots, m-1$, the left and right nodes in layer $3k$ output $\phi^{\circ k}(x)$ and $4^{k-1}x - \sum_{j=1}^{k-1} 4^{k-1-j}\phi^{\circ j}(x)$, respectively. Using the $\psi$ network and identity (4.4) again, the $3k + 3$ layer right node outputs $4^k x - \sum_{j=1}^k 4^{k-j}\phi^{\circ j}(x)$. Thus, generating $4^m x - \sum_{j=1}^m 4^{m-j}\phi^{\circ j}(x)$ requires a network with $3m + 3$ layers and $O(m)$ nodes and parameters. Finally, we divide by $4^m$, which can be done with a network of size $(2m, 2m, 2m)$ to produce $S_\varepsilon^+(x)$ in layer $5m + 3$. We use a duplication network with two layers to produce $S_\varepsilon^+(x)$ in layer $5m + 5$.

For $S_\varepsilon^-$, we use the same network as for $S_\varepsilon^+$ and include a third column. The layer 1 rightmost node outputs $1/2$, which is implementable by a network of size $(1,1,1)$. We carry down $1/2$ another $3m+2$ layers by a network of size $O(m)$. Then we divide by $1/4^m$ via a network of size $(2m, 2m, 2m)$. In layer $5m + 3$, the middle and right nodes output $S_\varepsilon^+(x)$ and $1/2^{2m+1}$. From here, we can easily generate $S_\varepsilon^-(x)$ using another two layer network that implements $(a,b) \mapsto \sigma(a - b)$. $\qquad\square$
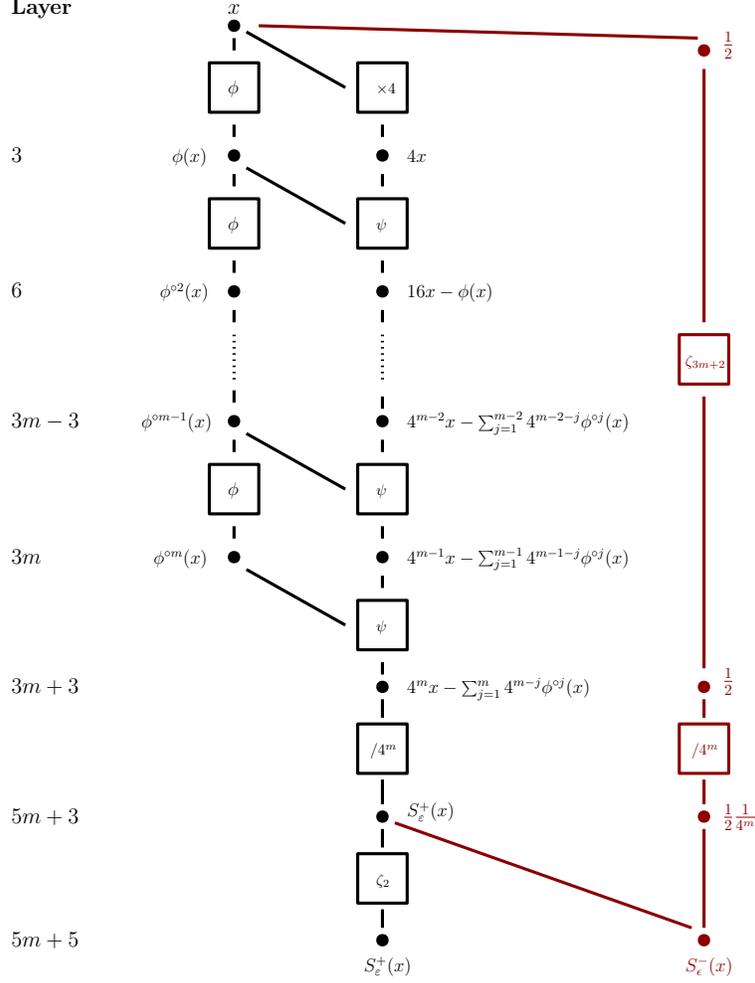
Figure 6: Implementations of $S_\varepsilon^+$ and $S_\varepsilon^-$

**Lemma A.4.** *For any $\varepsilon > 0$, there exists an activated $\mathcal{A}_{1/2}$-quantized ReLU neural network of size $O(\log(1/\varepsilon))$ that implements a nonnegative function $P_\varepsilon \colon [0,1]^2 \to \mathbb{R}$ such that*

$$\text{for any } x, y \in [0,1], \quad P_\varepsilon(x,y) \leq xy, \quad \text{and} \quad |P_\varepsilon(x,y) - xy| \leq \varepsilon.$$

*Proof.* Let $\varepsilon \in (0,1)$, and $S_\delta^+$ and $S_\delta^-$ be the approximate squaring function from Lemma A.3, where $\delta = \varepsilon/6$. We define the functions

$$P_{\varepsilon,0}(x,y) := 2\left(S_\delta^-\left(\frac{x+y}{2}\right) - S_\delta^+\left(\frac{x}{2}\right) - S_\delta^+\left(\frac{y}{2}\right)\right) \quad \text{and} \quad P_\varepsilon(x,y) := \sigma(P_{\varepsilon,0}(x,y)).$$

Since $xy = 2(\frac{x+y}{2})^2 - 2(\frac{x}{2})^2 - 2(\frac{y}{2})^2$ and $xy \geq 0$, it follows from Lemma A.3 that

$$\left|P_\varepsilon(x,y) - xy\right| = \left|\sigma(P_{\varepsilon,0}(x,y)) - xy\right| \leq \left|P_{\varepsilon,0}(x,y) - xy\right| \leq 6\delta = \varepsilon.$$

Additionally, using that $S_\varepsilon^+(x) \geq x^2$ and $S_\varepsilon^-(x) \leq x^2$, we see that

$$P_{\varepsilon,0}(x,y) \leq 2\left(\left(\frac{x+y}{2}\right)^2 - \left(\frac{x}{2}\right)^2 - \left(\frac{y}{2}\right)^2\right) = xy.$$
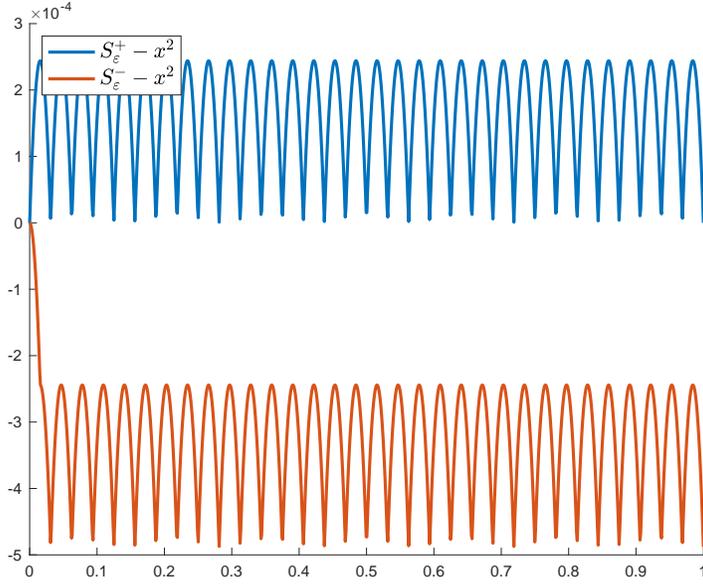
38

Figure 7: $S_\varepsilon^+ - x^2$ and $S_\varepsilon^- - x^2$ for $\varepsilon = 10^{-3}$

This inequality together implies $P_\varepsilon(x, y) \leq xy$ as well, because either $P_{\varepsilon,0}(x, y) \leq 0$ in which case $P_\varepsilon(x, y) = 0 \leq xy$, or $P_{\varepsilon,0}(x, y) > 0$ and $P_\varepsilon(x, y) = P_{\varepsilon,0}(x, y) \leq xy$.

Now we count the number of parameters. For inputs $x, y \in [0, 1]$, we first generate $x/2 = \sigma(x/2)$, $y/2 = \sigma(y/2)$, and $(x+y)/2 = \sigma(x/2+y/2)$. This can be done with a $\mathcal{A}_{1/2}$-quantized ReLU network of size $(1, 3, 4)$. Next, we place three networks in parallel, one $S_\delta^-$ network that takes $(x + y)/2$ as input, and two $S_\delta^+$ networks that take $x/2$ and $y/2$ as input. Note the implementations of $S_\delta^+$ and $S_\delta^-$ each has size $O(\log(1/\varepsilon))$ and the same number of layers. Then we generate

$$\sigma\left(\frac{1}{2}S_\delta^-\left(\frac{x+y}{2}\right) - \frac{1}{2}S_\delta^+\left(\frac{x}{2}\right) - \frac{1}{2}S_\delta^+\left(\frac{y}{2}\right)\right)$$

with a network of size $(1, 1, 3)$. Finally, we need a multiplication by 4, which can be implemented by the top network shown in Figure 5c. □

While $S_\varepsilon^+$ approximates the quadratic function from above, we also constructed $S_\varepsilon^-$, which approximates the quadratic function from below. An example of their approximation errors is shown in Figure 7. We are not aware of other neural network papers that have used $S_\varepsilon^-$, but it is necessary to ensure that $P_\varepsilon$ is nonnegative, which is an important property in subsequent results.

Now we move onto the product of several real numbers. In many neural network approximation papers, approximate multiplication is done sequentially, such as $P_\delta(a, P_\delta(b, P_\delta(c, d)))$ for four inputs $a, b, c, d$. This strategy is inefficient without use of skip connections since $a, b$ would need to be brought down a considerable number of layers via duplication networks. We perform multiplication of $d$ numbers in a dyadic manner for improved efficiency.

**Lemma A.5.** *For any integer $d \geq 2$ and $\varepsilon > 0$, there exists an activated $\mathcal{A}_{1/2}$-quantized ReLU neural network of size $O(\log(1/\varepsilon))$ that implements a nonnegative function $P_{\varepsilon,d} \colon [0, 1]^d \to \mathbb{R}$ such that*

$$\text{for all } \boldsymbol{x} \in [0, 1]^d, \quad 0 \leq P_{\varepsilon,d}(\boldsymbol{x}) \leq x_1 \cdots x_d, \quad \text{and} \quad |P_{\varepsilon,d}(\boldsymbol{x}) - x_1 \cdots x_d| \leq \varepsilon.$$

39

*Proof.* Let $\varepsilon \in (0,1)$, set $m = \lceil \log d \rceil$, and we will pick $\delta_1, \ldots, \delta_m$ in terms of $\varepsilon$ and $d$ later. We let $P_{\delta_k}$ be both the function and its network implementation from Proposition A.4, and $L_k = O(\log(1/\delta_k))$ be the number of layers in $P_{\delta_k}$. Fix an arbitrary $\boldsymbol{x} \in [0,1]^d$ and for reasons that will become apparent, let $u_{0,\ell} = x_\ell$.

Every consecutive pair of input nodes $x_1, \ldots, x_d$ are fed into $\lfloor d/2 \rfloor$ approximate multiplication networks $P_{\delta_1}$ placed in parallel. If $d$ is odd, then we simply carry $x_d$ down $L_1$ layers with a copy network of size $O(L_1)$. In layer $L_1$, there are $n_1 = \lceil d/2 \rceil$ nodes whose outputs we call $u_{1,1}, \ldots, u_{1,n_1}$. For example, $u_{1,1} = P_{\delta_1}(x_1, x_2)$.

Each consecutive pair of $u_{1,1}, \ldots, u_{1,n_1}$ are fed into at most $\lfloor n_1/2 \rfloor$ networks $P_{\delta_2}$ placed in parallel, and if $n_1$ is odd, then $u_{1,n_1}$ is carried down $L_2$ layers by a copy network of size $O(L_2)$. In layer $L_1 + L_2$, we call the outputs $u_{2,1}, \ldots, u_{2,n_2}$, where $n_2 = \lceil n_1/2 \rceil$. We continue this process with $P_{\delta_3}, \ldots, P_{\delta_m}$ and denote the outputs following each $P_{\delta_k}$ or copy network by $u_{k,1}, \ldots, u_{k,n_k}$. The final output is $u_m := P_{\varepsilon,d}(\boldsymbol{x})$.

We first show how to pick the $\delta_1, \ldots, \delta_m$ properly and quantify the approximation error. First note that from Proposition A.4, for any $a, b \in [0,1]$ and $\delta > 0$, we have $0 \le P_\delta(a,b) \le ab \le 1$, so by induction, we have that $0 \le u_{k,\ell} \le 1$. The quantities $U := \{u_{k,\ell}\}_{k=0,\ldots,m, \; \ell=1,\ldots,n_k}$ that are created have a tree structure with leaves $x_1, \ldots, x_d$. For a $u \in U$, we let $p(u)$ be product of all $x_1, \ldots, x_d$ that are connected to $u$. We have that

$$u_{1,\ell} \le p(u_{1,\ell}) \le 1 \quad \text{and} \quad |u_{1,\ell} - p(u_{1,\ell})| \le \delta_1 \quad \text{for} \quad \ell = 1, \ldots, n_1.$$

Fix any $u_{k,\ell}$. If $u_{k,\ell}$ is not generated as an output of a $P_{\delta_k}$ network, then it is created by copying a $u_{k-1,\ell'}$, in which case, $|u_{k,\ell} - p(u_{k,\ell})| = 0$. If not, $u_{k,\ell} = P_{\delta_k}(u_{k-1,a}, u_{k-1,b})$ for distinct $1 \le a, b \le n_{k-1}$. Then

$$\begin{aligned}
|u_{k,\ell} - p(u_{k,\ell})| &= |P_{\delta_k}(u_{k-1,a}, u_{k-1,b}) - p(u_{k,\ell})| \\
&= |P_{\delta_k}(u_{k-1,a}, u_{k-1,b}) - u_{k-1,a}u_{k-1,b}| + |u_{k-1,a}u_{k-1,b} - u_{k-1,a}p(u_{k-1,b})| \\
&\quad + |u_{k-1,a}p(u_{k-1,b}) - p(u_{k-1,a})p(u_{k-1,b})| \\
&\le \delta_k + |u_{k-1,b} - p(u_{k-1,b})| + |u_{k-1,a} - p(u_{k-1,a})|.
\end{aligned}$$

Hence, the error made by $u_{k,\ell}$ is bounded above by $\delta_k$ plus the errors made in the previous two $u_{k-1,a}$ and $u_{k-1,b}$. Inducting on $k$, we see that

$$|u_m - p(u_m)| \le 2^{m-1}\delta_1 + 2^{m-2}\delta_2 + \cdots + \delta_m \le d\big(\delta_1 + 2^{-1}\delta_2 + \cdots + 2^{-m+1}\delta_m\big),$$

where we noted that $2^{m-1} = 2^{\lceil \log d \rceil - 1} \le d$. To make the final error bounded by $\varepsilon$, we pick $\delta_k = \varepsilon/(2d)$ for each $k = 1, \ldots, m$. This proves that $P_{\varepsilon,d}(\boldsymbol{x}) - x_1 \cdots x_d| \le \varepsilon$ for all $\boldsymbol{x} \in [0,1]^d$.

It remains to count the size of this network that implements $P_{\varepsilon,d}$. Generating $\{u_{k,\ell}\}_{\ell=1,\ldots,n_k}$ from $\{u_{k-1,\ell}\}_{\ell=1,\ldots,n_{k-1}}$ requires $n_k$ networks each of size $O(\log(1/\delta_k)) = O(\log(1/\varepsilon))$. Hence, the number of layers in $P_{\varepsilon,d}$ is $O(m \log(1/\varepsilon)) = O(\log(1/\varepsilon))$. Since $n_k \le 2d/2^k$ and

$$\sum_{k=1}^m n_k \log(1/\delta_k) \le \sum_{k=1}^m \frac{2d}{2^k} \log(2d/\varepsilon) \le 2d \log(2d/\varepsilon),$$

the resulting network that implements $P_{\varepsilon,d}$ has $O(\log(1/\varepsilon))$ nodes and parameters. $\qquad\square$

**Lemma A.6.** *For any $\varepsilon > 0$ and integer $n \ge 1$, there exists an activated $\mathcal{A}_{1/2}$-quantized ReLU neural network with $O(n \log(n/\varepsilon))$ layers and $O(n^2 \log(n/\varepsilon))$ nodes and parameters that implements a function $b_n \colon [0,1] \to \mathbb{R}^{n+1}$ such that for each $0 \le k \le n$,*

$$0 \le b_{n,k} \le p_{n,k}, \quad \text{and} \quad \|b_{n,k} - p_{n,k}\|_\infty \le \varepsilon.$$

*Proof.* This construction will be done recursively. Fix an $\varepsilon > 0$, and we will pick an appropriate $\delta > 0$ depending only on $\varepsilon$ and $n$ later. Let $P_\delta$ denote both the approximate binary multiplication function and its network implementation from Proposition A.4 and $L_\delta = O(\log(1/\delta))$ be the number of layers.

Let $x \in [0,1]$ denote the input. We generate $1 - x$ via the formula,

$$1 - x = \sigma\Big(\frac{1}{2}\sigma\Big(-\frac{1}{2}x + \frac{1}{2}\Big) + \frac{1}{2}\sigma\Big(-\frac{1}{2}x + \frac{1}{2}\Big) + \frac{1}{2}\sigma\Big(-\frac{1}{2}x + \frac{1}{2}\Big) + \frac{1}{2}\sigma\Big(-\frac{1}{2}x + \frac{1}{2}\Big)\Big).$$

This is carried out by a network of size $(2, 5, 12)$. We next use a copy network to produce $x$ in the second layer. Hence, the first degree Bernstein polynomials $p_{1,0}(x) = 1 - x$ and $p_{1,1}(x) = x$ are exactly implementable, and are outputs of nodes in the second layer. We define

$$b_{1,0}(x) := \sigma(x) \quad \text{and} \quad b_{1,1}(x) := \sigma(1 - x).$$

Since we will need to use $x, 1 - x$ repeatedly throughout, we will use copy networks to bring them down however many layers we need. The size of these networks will be proportional to the total number of layers in $b_n$. We will see that the size of these copy networks will be dominated by the other portions of the final network.

Set $L_1 = 2$. We recursively define the following. For each $1 \le m \le n - 1$ and $1 \le k \le m$, let

$$b_{m+1,k}(x) := \mathrm{sum}\Big(P_\delta\big(x, b_{m,k-1}(x)\big), P_\delta\big(1 - x, b_{m,k}(x)\big)\Big),$$

where $\mathrm{sum}(\cdot, \cdot)$ refers to the two layer summation network. This shows that $b_{m+1,k}(x)$ can be generated provided that $b_{m,k-1}(x), b_{m,k}(x), x, 1 - x$ are all outputs of nodes that appear in layer $L_m$. If so, $b_{m+1,1}(x), \ldots, b_{m+1,m}(x)$ as outputs of nodes in layer $L_{m+1} := L_m + L_\delta + 2$. For the remaining two endpoint cases, let $\zeta_2$ be a two layer duplication network. We define

$$b_{m+1,0}(x) := \zeta_2(P_\delta(1 - x, b_{m,0}(x))) \quad \text{and} \quad b_{m+1,m+1}(x) := \zeta_2(P_\delta(x, b_{m,k-1}(x))),$$

which are also outputs in layer $L_{m+1}$. Finally, we copy $x, 1 - x$ from layer $L_m$ down to layer $L_{m+1}$ which requires a network of size $O(L_\delta)$.

We still need to show that these functions are well defined, because $P_\delta$ takes inputs in $[0,1]^2$. To establish this, we prove the stronger statement that $0 \le b_{m,k} \le p_{m,k}$ for each $1 \le m \le n$ and $0 \le k \le m$. We proceed by induction on $m$. For the base case, we have $b_{1,0} = p_{1,0}$, and $b_{1,1} = p_{1,1}$. Assume that for some $m$, we have that $b_{m,k} \le p_{m,k}$ for each $0 \le k \le m$. Now, consider any $1 \le k \le m + 1$. By Lemma A.4, for all $x \in [0,1]$, we have

$$\begin{aligned}
b_{m+1,k}(x) &= P_\delta\big(x, b_{m,k-1}(x)\big) + P_\delta\big(1 - x, b_{m,k}(x)\big) \\
&\le x b_{m,k-1}(x) + (1 - x) b_{m,k}(x) \\
&\le x p_{m,k-1}(x) + (1 - x) p_{m,k}(x) = p_{m+1,k}(x).
\end{aligned}$$

The remaining two cases $k = 0$ and $k = m + 1$ follow from an analogous argument. To summarize, we have shown that for $1 \le m \le n - 1$, layer $L_m$ has $m + 3$ nodes whose outputs are $b_{m,0}(x), \ldots, b_{m,m}(x), x, 1 - x$.

We proceed to examine the approximation error. For convenience, let

$$\alpha_{m,k} := \|b_{m,k} - p_{m,k}\|_\infty, \quad \text{and} \quad \alpha_m := \max_{0 \le k \le m} \alpha_{m,k}.$$

Hence $\alpha_1 = 0$. Using the recurrence relation (4.1), triangle inequality, and Lemma A.4, we have

$$
\begin{aligned}
\alpha_{m+1,k} &\leq \sup_{x \in [0,1]} \left( \left| x p_{m,k-1}(x) - x b_{m,k-1}(x) + (1-x) p_{m,k}(x) - (1-x) b_{m,k}(x) \right| \right. \\
&\qquad\qquad \left. + \left| P_\delta(x, b_{m,k-1}(x)) - x b_{m,k-1}(x) + P_\delta(1-x, b_{m,k}(x)) - (1-x) b_{m,k}(x) \right| \right) \\
&\leq \sup_{x \in [0,1]} \left( \delta + \delta + |x| \alpha_{m,k-1} + |1-x| \alpha_{m,k} \right) \\
&\leq \sup_{x \in [0,1]} \left( 2\delta + x \alpha_m + (1-x) \alpha_m \right) = 2\delta + \alpha_m.
\end{aligned}
\tag{A.2}
$$

Repeating the same argument for $k = 0$ and $k = m + 1$ provides us with

$$
\max(\alpha_{m+1,0}, \alpha_{m+1,m+1}) \leq \delta + \alpha_m.
\tag{A.3}
$$

Combining equations (A.2) and (A.3), we see that $\alpha_{m+1} \leq 2\delta + \alpha_m$, for all $1 \leq m \leq n$. A telescoping argument shows that, $\alpha_m = \alpha_m - \alpha_0 \leq 2(m-1)\delta$. Thus, we pick $\delta = \varepsilon/(2n)$ to see that

$$
\alpha_m = \max_{0 \leq k \leq m} \|b_{m,k} - p_{m,k}\|_\infty \leq \varepsilon \quad \text{for all } m \leq n.
$$

Now, we proceed to count the number of parameters. For the first row of this Pascal triangle, $b_{1,0}$ and $b_{1,1}$ can be made with a network of constant size. Computing each $b_{m+1,k}$ from the previous row $\{b_{m,k}\}_{k=0}^m$ requires at most two approximate multiplication networks $P_\delta$ with $\delta = \varepsilon/(2n)$ and a summation, which requires a network of size $O(\log(n/\varepsilon))$. Hence, computing $\{b_{m=1,k}\}_{k=0}^{m+1}$ from $\{b_{m,k}\}_{k=0}^m$ requires $O(\log(n/\varepsilon))$ layers and $O(m \log(n/\varepsilon))$ nodes and parameters. We do this from $m = 1$ to $m = n - 1$. $\qquad\square$

**Lemma A.7.** *For any $\varepsilon > 0$ and integers $n, d \geq 1$, there exists an activated $\mathcal{A}_{1/2}$-quantized ReLU neural network with $O(n \log(n/\varepsilon))$ layers and $O(n^2 \log(n/\varepsilon) + n^d \log(1/\varepsilon))$ nodes and parameters, as $n \to \infty$ and $\varepsilon \to 0$, that implements a function $b_n \colon [0,1]^d \to \mathbb{R}^{(n+1)^d}$ such that*

$$
\text{for each } 0 \leq \boldsymbol{k} \leq n, \quad b_{n,\boldsymbol{k}} \geq 0 \quad \text{and} \quad \|b_{n,\boldsymbol{k}} - p_{n,\boldsymbol{k}}\|_\infty \leq \varepsilon.
$$

*Proof.* Fix $\varepsilon \in (0,1)$, and we will pick appropriate $\delta, \gamma \in (0,1)$ later. Let $\boldsymbol{x} \in [0,1]^d$ be the input. For each $1 \leq \ell \leq d$, we use Lemma A.6 to provide us with an activated network with $O(n \log(n/\gamma))$ layers and $O(n^2 \log(n/\gamma))$ nodes and parameters that produces $\{b_{n,k_\ell}(x_\ell)\}_{k_\ell=0}^n$. Each of these $d$ networks have exactly the same number of layers, so placing all $d$ of them in parallel, we obtain a network that outputs $\{b_{n,k_\ell}(x_\ell)\}_{1 \leq \ell \leq d,\, 0 \leq k_\ell \leq n}$ in the same layer. We also have

$$
b_{n,k_\ell} \leq p_{n,k_\ell}, \quad \text{and} \quad \|p_{n,k_\ell} - b_{n,k_\ell}\|_\infty \leq \gamma.
$$

For each $0 \leq \boldsymbol{k} \leq n$, we use a $d$-term approximate product ReLU neural network $P_{\delta,d}$ as in Lemma A.5 and define

$$
b_{n,\boldsymbol{k}}(\boldsymbol{x}) := P_{\delta,d}\big(b_{n,k_1}(x_1), \ldots, b_{n,k_d}(x_d)\big).
$$

This is well-defined since $b_{n,k_\ell} \leq p_{n,k_\ell} \leq 1$.

We need $(n+1)^d$ many such $P_{\delta,d}$ networks placed in parallel, and each one has size $O(\log(1/\delta))$. The entire implementation of $\{b_{n,\boldsymbol{k}}(\boldsymbol{x})\}_{0 \leq \boldsymbol{k} \leq n}$ can be done by a network with

$$
\begin{aligned}
O\big(n \log(n/\gamma) + \log(1/\delta)\big) &\quad \text{layers,} \\
O\big(n^2 \log(n/\gamma) + n^d \log(1/\delta)\big) &\quad \text{nodes and parameters.}
\end{aligned}
$$

Next we compute the errors between $p_{n,\boldsymbol{k}}$ and $b_{n,\boldsymbol{k}}$, and then optimize over the parameters. For each $0 \leq \boldsymbol{k} \leq n$, we first apply Lemma A.5 to get

$$
\begin{aligned}
|p_{n,\boldsymbol{k}}(\boldsymbol{x}) - b_{n,\boldsymbol{k}}(\boldsymbol{x})| &\leq \left| P_{\delta,d}\Big(b_{n,k_1}(x_1), \ldots, b_{n,k_d}(x_d)\Big) - b_{n,k_1}(x_1) \cdots b_{n,k_d}(x_d) \right| \\
&\quad + |p_{n,\boldsymbol{k}}(\boldsymbol{x}) - b_{n,k_1}(x_1) \cdots b_{n,k_d}(x_d)| \\
&\leq \delta + |p_{n,\boldsymbol{k}}(\boldsymbol{x}) - b_{n,k_1}(x_1) \cdots b_{n,k_d}(x_d)|.
\end{aligned}
$$

To control the right hand side, we use that $p_{n,\boldsymbol{k}}(\boldsymbol{x})$ is a tensor product and peel off one term at a time. Then

$$
\begin{aligned}
|p_{n,\boldsymbol{k}}(\boldsymbol{x}) &- b_{n,k_1}(x_1) \cdots b_{n,k_d}(x_d)| \\
&\leq |p_{n,k_1}(x_1) - b_{n,k_1}(x_1)||p_{n,k_2}(x_2) \cdots p_{n,k_d}(x_d)| \\
&\quad + |b_{n,k_1}(x_1)||p_{n,k_2}(x_2) \cdots p_{n,k_d}(x_d) - b_{n,k_2}(x_2) \cdots b_{n,k_d}(x_d)| \\
&\leq \gamma + |p_{n,k_2}(x_2) \cdots p_{n,k_d}(x_d) - b_{n,k_2}(x_2) \cdots b_{n,k_d}(x_d)|.
\end{aligned}
$$

Continuing in this manner, we obtain the inequality

$$
|p_{n,\boldsymbol{k}}(\boldsymbol{x}) - b_{n,\boldsymbol{k}}(\boldsymbol{x})| \leq \delta + |p_{n,\boldsymbol{k}}(\boldsymbol{x}) - b_{n,k_1}(x_1) \cdots b_{n,k_d}(x_d)| \leq \delta + d\gamma.
$$

We select $\gamma = \varepsilon/(2d)$ and $\delta = \varepsilon/2$ to complete the proof. $\qquad\square$

## Acknowledgments

## References

[1] José A Adell, Jorge Bustamante, and José M Quesada. Estimates for the moments of Bernstein polynomials. *Journal of Mathematical Analysis and Applications*, 432(1):114–128, 2015.

[2] Jonathan Ashbrock and Alexander M. Powell. Stochastic Markov gradient descent and training low-bit neural networks. *Sampling Theory, Signal Processing, and Data Analysis*, 19(2):1–23, 2021.

[3] Andrew R. Barron. Approximation and estimation bounds for artificial neural networks. *Machine Learning*, 14(1):115–133, 1994.

[4] Bernhard Beckermann and Alex Townsend. On the singular values of matrices with displacement structure. *SIAM Journal on Matrix Analysis and Applications*, 38(4):1227–1248, 2017.

[5] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.

[6] Helmut Bolcskei, Philipp Grohs, Gitta Kutyniok, and Philipp Petersen. Optimal approximation with sparsely connected deep neural networks. *SIAM Journal on Mathematics of Data Science*, 1(1):8–45, 2019.

[7] Alfredo Canziani, Adam Paszke, and Eugenio Culurciello. An analysis of deep neural network models for practical applications. *arXiv preprint arXiv:1605.07678*, 2016.

[8] Evan Chou, C. Sinan Güntürk, Felix Krahmer, Rayan Saab, and Özgür Yılmaz. Noise-shaping quantization methods for frame-based and compressive sampling systems. *Sampling Theory, A Renaissance*, pages 157–184, 2015.

[9] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. *Advances in neural information processing systems*, 28, 2015.

[10] George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2(4):303–314, 1989.

[11] Ingrid Daubechies and Ronald DeVore. Approximating a bandlimited function using very coarsely quantized data: A family of stable sigma-delta modulators of arbitrary order. *Annals of Mathematics*, 158(2):679–710, 2003.

[12] Ingrid Daubechies, Ronald DeVore, Simon Foucart, Boris Hanin, and Guergana Petrova. Non-linear approximation and (deep) ReLU networks. *Constructive Approximation*, pages 1–46, 2021.

[13] Ronald DeVore, Boris Hanin, and Guergana Petrova. Neural network approximation. *Acta Numerica*, 30:327–444, 2021.

[14] Ronald A. DeVore and George G. Lorentz. *Constructive Approximation*, volume 303. Springer Science & Business Media, 1993.

[15] Chunmei Ding and Feilong Cao. K-functionals and multivariate Bernstein polynomials. *Journal of Approximation Theory*, 155(2):125–135, 2008.

[16] Günter Felbecker. Linearkombinationen von iterierten Bernsteinoperatoren. *Manuscripta Mathematica*, 29(2):229–248, 1979.

[17] Adrian Fellhauer. Approximation of smooth functions using bernstein polynomials in multiple variables. *arXiv preprint arXiv:1609.01940*, 2016.

[18] Heinz H Gonska and Xin-long Zhou. Approximation theorems for the iterated boolean sums of bernstein operators. *Journal of Computational and applied mathematics*, 53(1):21–31, 1994.

[19] Rémi Gribonval, Gitta Kutyniok, Morten Nielsen, and Felix Voigtlaender. Approximation spaces of deep neural networks. *Constructive approximation*, 55(1):259–367, 2022.

[20] C. Sinan Güntürk. One-bit sigma-delta quantization with exponential accuracy. *Communications on Pure and Applied Mathematics*, 56(11):1608–1630, 2003.

[21] C. Sinan Güntürk. Approximation by power series with $\pm 1$ coefficients. *International Mathematics Research Notices*, 2005(26):1601–1610, 2005.

[22] C. Sinan Güntürk and Weilin Li. Approximation with one-bit polynomials in Bernstein form. *Submitted*, 2021.

[23] C. Sinan Güntürk and Weilin Li. Approximation with one-bit polynomials in Bernstein form. *Constructive Approximation*, pages 1–30, 2022.

[24] C. Sinan Güntürk and Weilin Li. Quantization for spectral super-resolution. *Constructive Approximation*, 56(3):619–648, 2022.

[25] Yunhui Guo. A survey on methods and theories of quantized neural networks. *arXiv preprint arXiv:1808.04752*, 2018.

[26] Clemens Heitzinger. *Simulation and Inverse Modeling of Semiconductor Manufacturing Processes*. PhD thesis, Technische Universität Wien: Vienna, Austria, 2002.

[27] Theophil H. Hildebrandt and Issac J. Schoenberg. On linear functional operations and the moment problem for a finite interval in one or several dimensions. *Annals of Mathematics*, pages 317–328, 1933.

[28] Andrei Nikolaevich Kolmogorov and Vladimir Mikhailovich Tikhomirov. $\varepsilon$-entropy and $\varepsilon$-capacity of sets in function spaces. *Uspekhi Matematicheskikh Nauk*, 14(2):3–86, 1959.

[29] George G. Lorentz. *Bernstein Polynomials*. Chelsea, 2nd edition, 1986.

[30] Jianfeng Lu, Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network approximation for smooth functions. *SIAM Journal on Mathematical Analysis*, 53(5):5465–5506, 2021.

[31] Eric Lybrand and Rayan Saab. A greedy algorithm for quantizing neural networks. *Journal of Machine Learning Research*, 22(156):1–38, 2021.

[32] Charles Micchelli. The saturation class and iterates of the Bernstein polynomials. *Journal of Approximation Theory*, 8(1):1–18, 1973.

[33] Giuseppe Molteni. Explicit bounds for even moments of bernstein's polynomials. *Journal of Approximation Theory*, 273:105658, 2022.

[34] Uri Shaham, Alexander Cloninger, and Ronald R. Coifman. Provable approximation properties for deep neural networks. *Applied and Computational Harmonic Analysis*, 44(3):537–557, 2018.

[35] Albert N. Shiryayev. *Selected Works of AN Kolmogorov: Volume III: Information Theory and the Theory of Algorithms*, volume 27. Springer, 1993.

[36] Elias M. Stein. *Singular Integrals and Differentiability Properties of Functions*, volume 2. Princeton University Press, 1970.

[37] Anatoliĭ Georgievich Vitushkin. *Theory of the Transmission and Processing of Information*. Pergamon Press, 1961.

[38] Hassler Whitney. Analytic extensions of differentiable functions defined in closed sets. *Transactions of the American Mathematical Society*, 36(1):63–89, 1934.

[39] Dmitry Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Networks*, 94:103–114, 2017.