

Best approximation by linear combinations of characteristic functions of half-spaces

Paul C. Kainen
Department of Mathematics
Georgetown University
Washington, D.C. 20057-1233, USA

Věra Kůrková
Institute of Computer Science, Academy of Sciences of the Czech Republic
Pod vodárenskou věží 2, 182 07 Prague 8, Czechia

Andrew Vogt
Department of Mathematics
Georgetown University
Washington, D.C. 20057-1233, USA

Abstract

It is shown that for any positive integer n and any function f in $\mathcal{L}_p([0, 1]^d)$ with $p \in [1, \infty)$ there exist n half-spaces such that f has a best approximation by a linear combination of their characteristic functions. Further, any sequence of linear combinations of n half-space characteristic functions converging in distance to the best approximation distance has a subsequence converging to a best approximation, i.e., the set of such n -fold linear combinations is an approximatively compact set.

Keywords. Best approximation, proximal, approximatively compact, boundedly compact, Heaviside perceptron networks, plane waves.

1 Introduction

An important type of nonlinear approximation is *variable-basis approximation*, where the set of approximating functions is formed by linear combinations of n functions from a given set. This approximation scheme has been widely investigated: it includes splines with free nodes, trigonometric polynomials with free frequencies, sums of wavelets, and feedforward neural networks.

To estimate rates of variable-basis approximation, it is helpful to study properties like existence, uniqueness, and continuity of corresponding approximation operators.

We investigate the existence property for one-hidden-layer Heaviside perceptron networks. Here the approximations are by linear combinations of characteristic functions of closed half-spaces. (The characteristic function of any subset A is the function χ_A with value 1 on the subset and 0 elsewhere.) Such a characteristic function may also be described as a plane wave obtained by composing the Heaviside function with an affine function. We show that for all positive integers n, d in $\mathcal{L}_p([0, 1]^d)$ with $p \in [1, \infty)$ there exists a best approximation mapping to the set of functions computable by Heaviside perceptron networks with n hidden and d input units. Thus for any p -integrable function on $[0, 1]^d$ there is a linear combination of n characteristic functions of closed half-spaces that is nearest in the \mathcal{L}_p -norm. A related proposition is proved by Chui, Li, and Mhaskar in [2], where certain sequences are shown to have subsequences that converge a. e. These authors work in R^d rather than $[0, 1]^d$ and show a. e. convergence rather than \mathcal{L}_p convergence.

2 Heaviside perceptron networks

Feedforward networks compute parametrized sets of functions dependent both on the type of computational units and their interconnections. *Computational units* compute functions of two vector variables: an *input vector* and a *parameter vector*. A standard type of computational unit is the perceptron. A *perceptron* with an *activation function* $\psi : \mathcal{R} \rightarrow \mathcal{R}$ (where \mathcal{R} denotes the set of real numbers) computes real-valued functions on $\mathcal{R}^d \times \mathcal{R}^{d+1}$ of the form $\psi(\mathbf{v} \cdot \mathbf{x} + b)$, where $\mathbf{x} \in \mathcal{R}^d$ is an *input vector*, $\mathbf{v} \in \mathcal{R}^d$ is an *input weight vector*, and $b \in \mathcal{R}$ is a *bias*.

The most common activation functions are sigmoidals, i.e., functions with ess-shaped graph. Both continuous and discontinuous sigmoidals are used. Here we study networks based on the archetypal discontinuous sigmoidal, namely, the *Heaviside function* ϑ defined by $\vartheta(t) = 0$ for $t < 0$ and $\vartheta(t) = 1$ for $t \geq 0$.

Let H_d denote the set of functions on $[0, 1]^d$ computable by Heaviside perceptrons, i.e.,

$$H_d = \{f : [0, 1]^d \rightarrow \mathcal{R} : f(\mathbf{x}) = \vartheta(\mathbf{v} \cdot \mathbf{x} + b), \mathbf{v} \in \mathcal{R}^d, b \in \mathcal{R}\}.$$

H_d is the set of characteristic functions of closed half-spaces of \mathcal{R}^d restricted to $[0, 1]^d$, which is a subset of the set of plane waves (see, e.g., Courant and Hilbert [3, pp.676–681]).

The simplest type of multilayer feedforward network has one hidden layer and one linear output. Such networks with Heaviside perceptrons in the hidden layer compute functions of the form

$$\sum_{i=1}^n w_i \vartheta(\mathbf{v}_i \cdot \mathbf{x} + b_i),$$

where n is the number of hidden units, $w_i \in \mathcal{R}$ are output weights, and $\mathbf{v}_i \in \mathcal{R}^d$ and $b_i \in \mathcal{R}$ are input weights and biases respectively.

The set of all such functions is the *set of all linear combinations of n elements of H_d* and is denoted by $span_n H_d$.

It is known that for all positive integers d , $\cup_{n \in \mathcal{N}_+} span_n H_d$ (where \mathcal{N}_+ denotes the set of all positive integers) is dense in $(\mathcal{C}([0, 1]^d), \|\cdot\|_C)$, the linear space of all continuous functions on $[0, 1]^d$ with the supremum norm, as well as in $(\mathcal{L}_p([0, 1]^d), \|\cdot\|_p)$ with $p \in [1, \infty]$ (see, e.g., Mhaskar and Micchelli [11] or Leshno et al. [10]). We study best approximation in $span_n H_d$ for a fixed n .

3 Best approximation and approximative compactness

Existence of a best approximation has been formalized in approximation theory by the concept of proximal set (sometimes also called “existence” set). A subset M of a normed linear space $(X, \|\cdot\|)$ is called *proximal* if for every $f \in X$ the distance $\|f - M\| = \inf_{g \in M} \|f - g\|$ is achieved for some element of M , i.e., $\|f - M\| = \min_{g \in M} \|f - g\|$ (Singer [14]). Clearly a proximal subset must be closed.

A sufficient condition for proximality of a subset M of a normed linear space $(X, \|\cdot\|)$ is compactness. Indeed, for each $f \in X$ the functional $e_{\{f\}} : M \rightarrow \mathcal{R}$ defined by $e_{\{f\}}(m) = \|f - m\|$ is continuous [14, p. 391] and hence must achieve its minimum on the compact set M . Two generalizations of compactness also imply proximality. A set M is called *boundedly compact* if the closure of its intersection with any bounded set is compact. A set M is called *approximatively compact* if for each $f \in X$ and any sequence $\{g_i : i \in \mathcal{N}_+\}$ in M such that $\lim_{i \rightarrow \infty} \|f - g_i\| = \|f - M\|$, there exists $g \in M$ such that $\{g_i : i \in \mathcal{N}_+\}$ converges subsequentially to g [14, p.368]. Any closed, boundedly compact set is approximatively compact, and any approximatively compact set is proximal [14, p. 374].

Gurvits and Koiran [6] have shown that for all positive integers d the set of characteristic functions of half-spaces H_d is compact in $(\mathcal{L}_p([0, 1]^d), \|\cdot\|_p)$ with $p \in [1, \infty)$. This can be easily verified once the set H_d is reparametrized by elements of the unit sphere S^d in \mathcal{R}^{d+1} . Indeed, a function $\vartheta(\mathbf{v} \cdot \mathbf{x} + b)$, with the vector $(v_1, \dots, v_d, b) \in \mathcal{R}^{d+1}$ nonzero, is equal to $\vartheta(\hat{\mathbf{v}} \cdot \mathbf{x} + \hat{b})$, where $(\hat{v}_1, \dots, \hat{v}_d, \hat{b}) \in S^d$ is obtained from $(v_1, \dots, v_d, b) \in \mathcal{R}^{d+1}$ by normalization. Strictly speaking, H_d is parametrized by equivalence classes in S^d since different parametrizations may represent the same member of H_d when restricted to $[0, 1]^d$. Since S^d is compact, and the quotient spaces formed by the equivalence classes is likewise, so is H_d .

However, $\text{span}_n H_d$ is not compact for any positive integer n . Nor is it boundedly compact.

The following theorem shows that $\text{span}_n H_d$ is approximatively compact in \mathcal{L}_p -spaces. It extends a result of Kůrková [9], who showed that $\text{span}_n H_d$ is closed in \mathcal{L}_p -spaces with $p \in (1, \infty)$.

Theorem 3.1 *For every n, d positive integers and for every $p \in [1, \infty)$ $\text{span}_n H_d$ is an approximatively compact subset of $(\mathcal{L}_p([0, 1]^d), \|\cdot\|_p)$.*

To prove the theorem we need the following lemma. For a set A , let $\mathcal{P}(A)$ denotes the set of all subsets of A .

Lemma 3.2 *Let m be a positive integer, $\{a_{jk} : k \in \mathcal{N}_+, j = 1, \dots, m\}$ be m sequences of real numbers, and $\mathcal{S} \subseteq \mathcal{P}(\{1, \dots, m\})$ be such that for each $S \in \mathcal{S}$, $\lim_{k \rightarrow \infty} \sum_{j \in S} a_{jk} = c_S$ for some $c_S \in \mathcal{R}$. Then there exist real numbers $\{a_j : j = 1, \dots, m\}$ such that for each $S \in \mathcal{S}$, $\sum_{j \in S} a_j = c_S$.*

Proof. Let $p = \text{card } \mathcal{S}$ and let $\mathcal{S} = \{S_1, \dots, S_p\}$. Define $T : \mathcal{R}^m \rightarrow \mathcal{R}^p$ by $T(x_1, \dots, x_m) = (\sum_{j \in S_1} x_j, \dots, \sum_{j \in S_p} x_j)$. Then T is linear, and hence its range is a subspace of \mathcal{R}^p and so is a closed set. Since $(c_{S_1}, \dots, c_{S_p}) \in \text{cl } T(\mathcal{R}^m) = T(\mathcal{R}^m)$, there exists $(a_1, \dots, a_m) \in \mathcal{R}^m$ with $(c_{S_1}, \dots, c_{S_p}) = T(a_1, \dots, a_m)$. \square

Proof of Theorem 3.1

Let $f \in \mathcal{L}_p([0, 1]^d)$ and let $\{\sum_{j=1}^n a_{jk} g_{jk} : k \in \mathcal{N}_+\}$ be a sequence of elements of $\text{span}_n H_d$

such that $\lim_{k \rightarrow \infty} \|f - \sum_{j=1}^n a_{jk} g_{jk}\|_p = \|f - \text{span}_n H_d\|_p$. Since H_d is compact, by passing to suitable subsequences we can assume that for all $j = 1, \dots, n$, there exist $g_j \in H_d$ such that $\lim_{k \rightarrow \infty} g_{jk} = g_j$ (here and in the sequel, we use the notation $\lim_{k \rightarrow \infty}$ to mean a limit of a suitable subsequence).

We shall show that there exist real numbers a_1, \dots, a_n such that

$$\|f - \text{span}_n H_d\|_p = \|f - \sum_{j=1}^n a_j g_j\|_p. \quad (1)$$

Then using (1) we shall show even that $\{\sum_{j=1}^n a_{jk} g_{jk} : k \in \mathcal{N}_+\}$ converges to $\sum_{j=1}^n a_j g_j$ in $\|\cdot\|_p$ subsequentially.

Decompose $\{1, \dots, n\}$ into two disjoint subsets I and J such that I consists of those j for which the sequences $\{a_{jk} : k \in \mathcal{N}_+\}$ have convergent subsequences, and J of those j for which the sequences $\{|a_{jk}| : k \in \mathcal{N}_+\}$ diverge. Again, by passing to suitable subsequences we can assume that for all $j \in I$, $\lim_{k \rightarrow \infty} a_{jk} = a_j$. Thus $\{\sum_{j \in I} a_{jk} g_{jk} : k \in \mathcal{N}_+\}$ converges subsequentially to $\sum_{j \in I} a_j g_j$.

Set $h = f - \sum_{j \in I} a_j g_j$. Since for all $j \in I$, the chosen subsequences $\{a_{jk} : k \in \mathcal{N}_+\}$ and $\{g_{jk} : k \in \mathcal{N}_+\}$ are bounded, we have $\|f - \text{span}_n H_d\|_p = \lim_{k \rightarrow \infty} \|f - \sum_{j=1}^n a_{jk} g_{jk}\|_p = \lim_{k \rightarrow \infty} \|h - \sum_{j \in J} a_{jk} g_{jk}\|_p$.

Let \mathcal{S} denotes the set of all subsets of J . Decompose \mathcal{S} into two disjoint subsets \mathcal{S}_1 and \mathcal{S}_2 such that \mathcal{S}_1 consists of those $S \in \mathcal{S}$ for which by passage to suitable subsequences $\lim_{k \rightarrow \infty} \sum_{j \in S} a_{jk} = c_S$ for some $c_S \in \mathcal{R}$, and \mathcal{S}_2 consists of those $S \in \mathcal{S}$ for which $\lim_{k \rightarrow \infty} |\sum_{j \in S} a_{jk}| = \infty$. Note that the empty set is in \mathcal{S}_1 with the convention $\sum_{j \in \emptyset} = 0$.

Using Lemma 3.2, for all $j \in \cup \mathcal{S}_1$, we get $a_j \in \mathcal{R}$ such that for all $S \in \mathcal{S}_1$, $\sum_{j \in S} a_j = c_S$. For $j \in J - \cup \mathcal{S}_1$, set $a_j = 0$.

Since $\sum_{j=1}^n a_j g_j \in \text{span}_n H_d$, we have $\|f - \text{span}_n H_d\|_p \leq \|f - \sum_{j=1}^n a_j g_j\|_p$ and thus to prove (1), it is sufficient to show that $\|f - \text{span}_n H_d\|_p \geq \|f - \sum_{j=1}^n a_j g_j\|_p$ or equivalently

$$\lim_{k \rightarrow \infty} \int_{[0,1]^d} \left| h - \sum_{j \in J} a_{jk} g_{jk} \right|^p d\mu \geq \int_{[0,1]^d} \left| h - \sum_{j \in J} a_j g_j \right|^p d\mu \quad (2)$$

where μ is Lebesgue measure on $[0, 1]^d$.

To verify (2), for each $k \in \mathcal{N}_+$ we shall decompose the integration over $[0, 1]^d$ into a sum of integrals over convex regions where the functions $\sum_{j \in J} a_{jk} g_{jk}$ are constant. To describe such regions, we shall define partitions of $[0, 1]^d$ determined by families of characteristic functions $\{g_{jk} : j \in J, k \in \mathcal{N}_+\}$, and $\{g_j : j \in J\}$. The partitions are indexed by the elements of the set \mathcal{S} of all subsets of J . For $k \in \mathcal{N}_+$, a partition $\{T_k(S) : S \in \mathcal{S}\}$ is defined by $T_k(S) = \{x \in [0, 1]^d : (g_{jk}(x) = 1 \Leftrightarrow j \in S)\}$, and similarly a partition $\{T(S) : S \in \mathcal{S}\}$ is defined by $T(S) = \{x \in [0, 1]^d : g_j(x) = 1 \Leftrightarrow j \in S\}$. Notice that since for all $j = 1, \dots, n$, $\lim_{k \rightarrow \infty} g_{jk} = g_j$ in $\mathcal{L}_p([0, 1]^d)$, we have $\lim_{k \rightarrow \infty} \mu(T_k(S)) = \mu(T(S))$ for all $S \in \mathcal{S}$. Indeed, the characteristic function of $T_k(S)$, $\chi_{T_k(S)}$, equals the product $\prod_{j \in S} g_{jk} \prod_{j \notin S} (1 - g_{jk})$ and converges in $\mathcal{L}_p([0, 1]^d)$ to the characteristic function of $T(S)$, $\chi_{T(S)}$, the latter equal to $\prod_{j \in S} g_j \prod_{j \notin S} (1 - g_j)$.

Using the definition of $T_k(S)$ (in particular its property guaranteeing that for all $S \in \mathcal{S}$, $T_k(S)$ is just the region where for all $j \in S$ and no other $j \in J$, g_{jk} is equal to 1), we get

$$\begin{aligned}
& \lim_{k \rightarrow \infty} \int_{[0,1]^d} \left| h - \sum_{j \in J} a_{jk} g_{jk} \right|^p d\mu = \lim_{k \rightarrow \infty} \sum_{S \in \mathcal{S}} \int_{T_k(S)} \left| h - \sum_{j \in S} a_{jk} \right|^p d\mu = \\
& \lim_{k \rightarrow \infty} \left(\sum_{S \in \mathcal{S}_1} \int_{T_k(S)} \left| h - \sum_{j \in S} a_{jk} \right|^p d\mu + \sum_{S \in \mathcal{S}_2} \int_{T_k(S)} \left| h - \sum_{j \in S} a_{jk} \right|^p d\mu \right) \geq \\
& \lim_{k \rightarrow \infty} \sum_{S \in \mathcal{S}_1} \int_{T_k(S)} \left| h - \sum_{j \in S} a_{jk} \right|^p d\mu. \tag{3}
\end{aligned}$$

Since for all $S \in \mathcal{S}$, $\lim_{k \rightarrow \infty} \mu(T_k(S)) = \mu(T(S))$ and for all $S \in \mathcal{S}_1$, $\lim_{k \rightarrow \infty} \sum_{j \in S} a_{jk} = c_S = \sum_{j \in S} a_j$, we have

$$\begin{aligned}
& \lim_{k \rightarrow \infty} \sum_{S \in \mathcal{S}_1} \int_{T_k(S)} \left| h - \sum_{j \in S} a_{jk} \right|^p d\mu = \lim_{k \rightarrow \infty} \sum_{S \in \mathcal{S}_1} \int_{T_k(S)} \left| h - \sum_{j \in S} a_j \right|^p d\mu = \\
& \sum_{S \in \mathcal{S}_1} \int_{T(S)} \left| h - \sum_{j \in S} a_j \right|^p d\mu.
\end{aligned}$$

For all $S \in \mathcal{S}$, by the triangle inequality in $\mathcal{L}_p(T_k(S))$

$$\begin{aligned}
& \lim_{k \rightarrow \infty} \left(\int_{T_k(S)} \left| \sum_{j \in S} a_{jk} \right|^p d\mu \right)^{1/p} \leq \\
& \lim_{k \rightarrow \infty} \left(\left(\int_{T_k(S)} \left| h - \sum_{j \in S} a_{jk} \right|^p d\mu \right)^{1/p} + \left(\int_{T_k(S)} |h|^p d\mu \right)^{1/p} \right) \leq \\
& \lim_{k \rightarrow \infty} \left(\int_{[0,1]^d} \left| h - \sum_{j \in J} a_{jk} g_{jk} \right|^p d\mu \right)^{1/p} + \left(\int_{[0,1]^d} |h|^p d\mu \right)^{1/p} = \|f - \text{span}_n H_d\|_p + \|h\|_p.
\end{aligned}$$

Thus for all $S \in \mathcal{S}$, $\lim_{k \rightarrow \infty} \int_{T_k(S)} \left| \sum_{j \in S} a_{jk} \right|^p d\mu$ is finite. In particular this is true when $S \in \mathcal{S}_2$, for which $\lim_{k \rightarrow \infty} \left| \sum_{j \in S} a_{jk} \right|^p = \infty$, and so $\lim_{k \rightarrow \infty} \mu(T_k(S)) = 0 = \mu(T(S))$ for $S \in \mathcal{S}_2$. Hence, we can replace integration over $\cup_{S \in \mathcal{S}_1} T(S)$ by integration over the whole of $[0,1]^d$, obtaining

$$\sum_{S \in \mathcal{S}_1} \int_{T(S)} \left| h - \sum_{j \in S} a_j \right|^p d\mu = \int_{[0,1]^d} \left| h - \sum_{j \in J} a_j g_j \right|^p d\mu,$$

which proves (2). Moreover, as a byproduct we even get that

$$\lim_{k \rightarrow \infty} \sum_{S \in \mathcal{S}_2} \int_{T_k(S)} \left| h - \sum_{j \in S} a_{jk} \right|^p d\mu = 0, \quad (4)$$

since in (3) the first expression is equal to the last (both are equal to $\|f - \text{span}_n H_d\|_p^p$).

So we have shown that $\text{span}_n H_d$ is proximal. Now we shall verify that it is even approximatively compact by showing that $\{\sum_{j \in J} a_{jk} g_{jk} : k \in \mathcal{N}_+\}$ converges subsequentially to $\sum_{j \in J} a_j g_j$, or equivalently

$$\lim_{k \rightarrow \infty} \int_{[0,1]^d} \left| \sum_{j \in J} (a_{jk} g_{jk} - a_j g_j) \right|^p d\mu = 0. \quad (5)$$

As above, we start by decomposing the integration into a sum of integrals over convex regions. The left hand side of (5) is equal to

$$\lim_{k \rightarrow \infty} \sum_{S \in \mathcal{S}_1} \int_{T_k(S)} \left| \sum_{j \in S} (a_{jk} - a_j g_j) \right|^p d\mu + \sum_{S \in \mathcal{S}_2} \int_{T_k(S)} \left| \sum_{j \in S} (a_{jk} - a_j g_j) \right|^p d\mu.$$

Using the triangle inequality, (4), and $\lim_{k \rightarrow \infty} \mu(T_k(S)) = 0$ for all $S \in \mathcal{S}_2$, we get

$$\begin{aligned} & \lim_{k \rightarrow \infty} \sum_{S \in \mathcal{S}_2} \int_{T_k(S)} \left| \sum_{j \in S} (a_{jk} - a_j g_j) \right|^p d\mu \leq \\ & \lim_{k \rightarrow \infty} \left(\sum_{S \in \mathcal{S}_2} \int_{T_k(S)} \left| h - \sum_{j \in S} a_{jk} \right|^p d\mu + \sum_{S \in \mathcal{S}_2} \int_{T_k(S)} \left| h - \sum_{j \in S} a_j g_j \right|^p d\mu \right) = \\ & \lim_{k \rightarrow \infty} \sum_{S \in \mathcal{S}_2} \int_{T_k(S)} \left| h - \sum_{j \in S} a_j g_j \right|^p d\mu = \sum_{S \in \mathcal{S}_2} \int_{T(S)} \left| h - \sum_{j \in S} a_j g_j \right|^p d\mu = 0 \end{aligned}$$

since $\mu(T(S)) = 0$ for $S \in \mathcal{S}_2$.

Thus $\lim_{k \rightarrow \infty} \sum_{S \in \mathcal{S}_2} \int_{T_k(S)} \left| \sum_{j \in S} (a_{jk} - a_j g_j) \right|^p d\mu = 0$, which implies that the left hand side of (5) is equal to

$$\lim_{k \rightarrow \infty} \sum_{S \in \mathcal{S}_1} \int_{T_k(S)} \left| \sum_{j \in S} (a_{jk} - a_j g_j) \right|^p d\mu = \sum_{S \in \mathcal{S}_1} \int_{T(S)} \left| \sum_{j \in S} a_j g_j - a_j g_j \right|^p d\mu = 0$$

because $(\sum_{j \in S} a_{jk} g_{jk}) \chi_{T_k(S)} = (\sum_{j \in S} a_{jk}) \chi_{T_k(S)}$ converges to $c_S \chi_{T(S)} = (\sum_{j \in S} a_j g_j) \chi_{T(S)}$ in $\mathcal{L}_p([0,1]^d)$.

So $\lim_{k \rightarrow \infty} \sum_{j \in J} a_{jk} g_{jk} = \sum_{j \in J} a_j g_j$, the same is already known to be true when J is replaced by I , and hence also $\lim_{k \rightarrow \infty} \sum_{j=1}^n a_{jk} g_{jk} = \sum_{j=1}^n a_j g_j$ subsequentially in $\mathcal{L}_p([0,1]^d)$. \square

Theorem 3.1 shows that a function in $\mathcal{L}_p([0, 1]^d)$ has a best approximation among functions computable by one-hidden-layer networks with a single linear output unit and n Heaviside perceptrons in the hidden layer. In other words, in the space of parameters of networks of this type, there exists a global minimum of the error functional defined as \mathcal{L}_p -distance from the function to be approximated.

Combining Theorem 3.1 with Theorem 2.2 of [8] (see also [7]), we note that while such best approximation operators exist from $\mathcal{L}_p([0, 1]^d)$ to $\text{span}_n H_d$, they cannot be continuous for $p \in (1, \infty)$.

4 Discussion

In Proposition 3.3 of [2] the authors show that any sequence $\{P_k\}$ in $\text{span}_n H_d$ (domain taken to be R^d here), with the property that $\limsup_k \|P_k\|_{\mathcal{L}_1(K)} \leq 1$ for every compact set K in R^d , has a subsequence converging a. e. in R^d to a member of $\text{span}_n H_d$. Although the proof techniques in [2] do have some overlap with those used here, the results there are different. A. e. convergence need not imply \mathcal{L}_p convergence for $p \in [1, \infty)$: the sequence $P_k = (k)^{\frac{1}{p}} \chi_A$, where $A = [0, \frac{1}{k}] \times [0, 1]^{d-1}$, converges a. e. in $\mathcal{L}_p(R^d)$ but has no convergent subsequence in the \mathcal{L}_p -norm.

Since the sequence $\{P_k\}$ above is bounded and lies in $\text{span}_1 H_d$ (with respect to $\mathcal{L}_p([0, 1]^d)$), it also illustrates that $\text{span}_1 H_d$, and hence $\text{span}_n H_d$, are not boundedly compact. Another example of an approximatively compact set that is not boundedly compact is any closed infinite-dimensional subspace of a uniformly convex Banach space [14, pp. 368-9].

Theorem 3.1 cannot be extended to perceptron networks with differentiable activation functions, e.g., the logistic sigmoid or hyperbolic tangent. For such functions, sets $\text{span}_n P_d(\psi)$ (where $P_d(\psi) = \{f : [0, 1]^d \rightarrow \mathcal{R} : f(\mathbf{x}) = \psi(\mathbf{v} \cdot \mathbf{x} + b), \mathbf{v} \in \mathcal{R}^d, b \in \mathcal{R}\}$) are not closed and hence cannot be proximal. This was first observed by Girosi and Poggio [5] and later exploited by Leshno et al. [10] for a proof of the universal approximation property.

Theorem 3.1 does not offer any information on the error of the best approximation. Estimates in the literature (DeVore, Howard, and Micchelli [4], Pinkus [12], Pinkus [13]) that give lower bounds on such errors and depend on continuity of best approximation operators are not applicable by the remarks at the end of section 3.

Cheang and Barron [1] show that linear combinations of characteristic functions of closed half-spaces with relatively few terms can yield good approximations of such functions as the characteristic function χ_B of a ball. However, χ_B is not approximated by the linear combination itself but rather by the characteristic function of the set where the linear combination exceeds a certain threshold. This amounts to replacing a linear output in the corresponding neural network by a threshold unit.

Acknowledgment

V. Kůrková was partially supported by GA ČR grants 201/99/0092 and 201/02/0428. Collaboration of P. Kainen, V. Kůrková, and A. Vogt was supported by NRC COBASE grants and by a Faculty Development grant from Georgetown University.

References

- [1] G. H. L. Cheang and A. R. Barron, A better approximation for balls. *J. Approx. Theory* **104** (2000), 183–203.
- [2] C. K. Chui, X. Li, and H. N. Mhaskar, Neural networks for localized approximation. *Math. of Computation* **63** (1994), 607–623.
- [3] R. Courant and D. Hilbert, “Methods of Mathematical Physics”, vol. II, Wiley, New York, 1962.
- [4] R. DeVore, R. Howard, and C. Micchelli, Optimal nonlinear approximation. *Manuscripta Math.* **63** (1989), 469–478.
- [5] F. Girosi and T. Poggio, Networks and the best approximation property, *Biological Cybernetics* **63** (1990), 169–176.
- [6] L. Gurvits and P. Koiran, Approximation and learning of convex superpositions. *J. of Computer and System Sciences* **55**, (1997), 161–170.
- [7] P. C. Kainen, V. Kůrková, and A. Vogt, Geometry and topology of continuous best and near best approximations. *J. Approx. Theory* **105** (2000), 252–262.
- [8] P. C. Kainen, V. Kůrková, and A. Vogt, Continuity of approximation by neural networks in L^p spaces. *Annals of Operations Research* **101** (2001), 143–147.
- [9] V. Kůrková, Approximation of functions by perceptron networks with bounded number of hidden units. *Neural Networks* **8** (1995), 745–750.
- [10] M. Leshno, V. Y. Lin, A. Pinkus, and S. Schocken, Multilayer feedforward networks with a nonpolynomial activation can approximate any function. *Neural Networks* **6** (1993), 861–867.
- [11] H. N. Mhaskar and C. Micchelli, Approximation by superposition of sigmoidal and radial basis functions. *Advances in Applied Math.* **13** (1992), 350–373.
- [12] A. Pinkus, “n-Width in Approximation Theory”, Springer-Verlag, Berlin, 1989.
- [13] A. Pinkus, Approximation theory of the MLP model in neural networks. *Acta Numerica* **8** (1999), 143–195.
- [14] I. Singer, “Best Approximation in Normed Linear Spaces by Elements of Linear Subspaces”, Springer-Verlag, Berlin, 1970.