

---

# Expression of Fractals Through Neural Network Functions

---

**Nadav Dym\***

Department of Mathematics and Rhodes Information Initiative  
Duke University  
Durham, NC 27708  
nadavdym@gmail.com

**Barak Sober\***

Department of Mathematics and Rhodes Information Initiative  
Duke University  
Durham, NC 27708  
barakino@math.duke.edu

**Ingrid Daubechies**

Department of Mathematics and Rhodes Information Initiative  
Department of Electrical and Computer Engineering  
Duke University  
Durham, NC 27708  
ingrid@math.duke.edu

## Abstract

To help understand the underlying mechanisms of neural networks (NNs), several groups have, in recent years, studied the number of linear regions  $\ell$  of piecewise linear functions generated by deep neural networks (DNN). In particular, they showed that  $\ell$  can grow exponentially with the number of network parameters  $p$ , a property often used to explain the advantages of DNNs over shallow NNs in approximating complicated functions. Nonetheless, a simple dimension argument shows that DNNs cannot generate all piecewise linear functions with  $\ell$  linear regions as soon as  $\ell > p$ . It is thus natural to seek to characterize specific families of functions with  $\ell$  linear regions that can be constructed by DNNs. Iterated Function Systems (IFS) generate sequences of piecewise linear functions  $F_k$  with a number of linear regions exponential in  $k$ . We show that, under mild assumptions,  $F_k$  can be generated by a NN using only  $\mathcal{O}(k)$  parameters. IFS are used extensively to generate, at low computational cost, natural-looking landscape textures in artificial images. They have also been proposed for compression of natural images, albeit with less commercial success. The surprisingly good performance of this fractal-based compression suggests that our visual system may lock in, to some extent, on self-similarities in images. The combination of this phenomenon with the capacity, demonstrated here, of DNNs to efficiently approximate IFS may contribute to the success of DNNs, particularly striking for image processing tasks, as well as suggest new algorithms for representing self similarities in images based on the DNN mechanism.

In the past few years Deep Neural Networks (DNN) have been remarkably successful in solving problems in computer vision (e.g., Krizhevsky et al. (2012); CireşAn et al. (2012); Qi et al. (2017)),

Preprint. Under review.

natural language processing (e.g., Manning et al. (2014); Collobert and Weston (2008)), signal processing (e.g., Piczak (2015)), and even art investigation (e.g., Sabetsarvestani et al. (2019)). This success has spawned a vast body of literature focusing on both practical and theoretical aspects of DNNs.

One of the main topics in the theoretical investigation of DNNs is the understanding of their expressiveness and approximation power. The universality theorem shows that even NNs with one hidden layer can approximate any continuous function Cybenko (1989); Hornik et al. (1989). However, these results do not give a reasonable explanation to the power of depth, which had shown to be pivotal to the success of DNNs. In recent years, several works showed that shallow networks are not as efficient as deep ones in terms of their approximation power Poggio et al. (2017); Yarotsky (2017); Liang and Srikant (2016); Cohen et al. (2016); Eldan and Shamir (2016).

In contrast with the study of approximation, the study of expressiveness focuses on understanding the class of functions which DNN can express *exactly*. As modern DNNs are typically composed of affine transformations and linear activation units (ReLU), functions expressed by DNNs must be continuous and piecewise linear (CPwL) by construction. Moreover, it is well established that neural networks with  $p$  parameters can express standard families of CPwL functions with  $\mathcal{O}(p)$  linear regions such as free knot splines and CPwL bases over triangulations (Arora et al. (2018); Yarotsky (2017); He and Zheng (2018); Daubechies et al. (2019)). As these results do not necessarily require the network to be deep this can be considered as an analogue of the universality theorem in terms of expressiveness. However, DNNs can also express CPwL functions with  $\ell \gg p$  linear regions. Several recent works Arora et al. (2018); Serra et al. (2017); Montúfar (2017); Montufar et al. (2014)) have studied the maximal number of linear regions generated by DNNs. Their results show that for DNNs  $\ell$  can in fact grow exponentially faster than  $p$ , if the depth of the network is allowed to grow. This property is often used as a possible explanation of the advantages of DNNs over shallow networks in expressing complex functions.

In this paper our aim is to further our understanding of the class of CPwL functions in the domain  $\ell \gg p$ . Using a simple dimension argument we show that DNNs with  $p$  parameters cannot span all piecewise linear functions with  $\ell \gg p$  linear regions (see Lemma 3 in Appendix A). What then are the CPwL functions which can be expressed by a DNN with a small number of parameters? To the best of our knowledge there is not much known on this topic, a notable exception being the works of Daubechies et al. (2019); Grohs et al. (2019) which construct some examples of such scalar functions.

Our main result is that DNNs are capable of expressing fractal-like functions, such as the (almost) indicator function of the fractals shown in figure 1. This construction is applicable for a wide family of fractals arising from *Iterated Functions Systems* (IFS) satisfying some weak conditions.

IFS generate intricate fractal-type shapes by iteratively applying all possible combinations of a small fixed number of functions. Throughout the history of computer vision research IFS have played an important role Welstead (1999). IFS make it possible to generate, at low computational cost, natural-looking landscape textures Barnsley (2014); they can be used for artificial generation of such images in entertainment applications such as movies and video games Van Pabst and Jense (1996). Self similarity in images has been suggested as a tool for several tasks in computer vision including deblurring, dehazing and super resolution Bahat and Irani (2016); Michaeli and Irani (2013); Irani (2009). IFS can also obtain surprisingly good results in image compression Jacquin (1992), by encoding geometric similarity of small parts of an (arbitrary) image with other, slightly larger parts elsewhere in the same (natural) image. That images compressed and then reconstructed via this coding are perceived as reasonable approximations of the original image may be connected as much (or more) with our visual system as with the properties of natural images; this observation is reminiscent of the role played by similar comparisons between parts of images in Smale et al. (2010). The particular efficiency of NNs at encoding functions or shapes generated by IFS, demonstrated in this paper, is intriguing in light of this possible connection of IFS with the perception of natural images through our own biological neural system, together with the observation that NNs were originally inspired, at least in part, by an (idealized) version of the network formed by biological neurons.

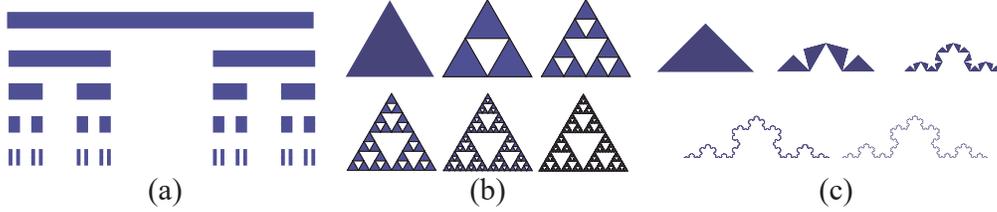


Figure 1: The first few sets in the iterative construction of (a) the Cantor set, (b) The Sierpinski triangle, and (c) the Koch curve.

## 1 Problem setup and main results

Prior to quoting our main result, we wish to introduce some notations pertaining to neural networks, and give a brief introduction to *iterated function systems*.

### 1.1 Neural network notation

The term neural network is being used in the literature to describe both learning algorithms (based upon artificial neural network architecture) and realizations of such architectures; i.e., a given architecture with fixed coefficients to be used as a function. In this paper, we focus on a family of functions that can be expressed through such architecture. Thus, to emphasize this distinction, throughout the paper we refer to such functions as Neural Network Functions (NNF). Furthermore, we limit our discussion to fully connected architectures with ReLU non-linear activation. We also restrict our attention to the behavior of  $f$  on some fixed arbitrarily large compact set  $\mathcal{K}$ .

**Definition 1** (ReLU). *Let  $x = (x_1, \dots, x_N)^T$ , then*

$$\text{ReLU}(x) \stackrel{\text{def}}{=} (\max(x_1, 0), \dots, \max(x_N, 0))^T, \quad (1)$$

**Definition 2** (NNF). *A neural network function is a function of the form*

$$f(x) = \eta_1 \circ \eta_2 \circ \dots \circ \eta_{2d+1}(x), \forall x \in \mathcal{K}, \quad (2)$$

where  $\eta_{2i-1}$  is an affine transformation and  $\eta_{2i} = \text{ReLU}$  for all  $i = 1, \dots, d$ .

The representation of  $f$  in the form (2) is non-unique, and we refer to a specific representation as  $f_\eta$ . Given such a representation we denote  $L(f_\eta) = d$  and refer to it as the *depth* of  $f_\eta$ . The *width* of  $f_\eta$  is the maximal input or output dimension of all  $\eta_i$  comprising  $f_\eta$  and is denoted by  $W(f_\eta)$ . We denote the family of NNFs with depth at most  $d$  and width at most  $w$  by

$$\text{NNF}_{w,d}^{n_1 \rightarrow n_2} \stackrel{\text{def}}{=} \{f : \mathcal{K} \subseteq \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2} \mid \exists f_\eta \text{ s.t. } f = f_\eta, L(f_\eta) \leq d, W(f_\eta) \leq w\}, \quad (3)$$

Note that the number of parameters needed to express  $f \in \text{NNF}_{w,d}^{n_1 \rightarrow n_2}$  is  $\mathcal{O}(w^2 \cdot d)$ .

### 1.2 Iterated Function System

We provide a short review on *Iterated Function System* (IFS). The results stated here can be found in the classical book on IFS Barnsley (2014).

An IFS on  $\mathbb{R}^d$  is a method for constructing very detailed patterns from a small number of building blocks. Mathematically, they are defined as a finite collection  $\mathcal{F} = \{f_1 \dots, f_J\}$  of contractive mappings  $f_j : \mathbb{R}^d \rightarrow \mathbb{R}^d$ . Given an initial bounded subset  $A_0 \subset \mathbb{R}^d$ , an IFS generates a sequence of increasingly detailed sets  $A_k$  by recursively applying the *Hutchinson operator*

$$H(A) = \bigcup_{j=1}^J f_j A,$$

to obtain a sequence of sets

$$A_{k+1} = H(A_k). \quad (4)$$

For any compact set  $A_0$ , the sequence  $A_k = H^k A_0$  converges in the Hausdorff metric to a compact set  $K$  which is called the attractor of the IFS. However, despite the asymptotic independence on the choice of the initial set, in practice from visual as well as computational considerations it is useful to work with special sets which we will call “nice” sets:

**Definition 3** (Nice sets). *We say that a bounded set  $A$  is nice with respect to  $\mathcal{F} = \{f_1, \dots, f_J\}$ , if*

$$f_j A \subseteq A, j = 1, \dots, J \quad (5)$$

and

$$f_i A \cap f_j A = \emptyset, \forall i \neq j, \quad (6)$$

The existence of nice sets is an important property in the study of IFS (Barnsley (2014) page 129):

**Definition 4** (Totally disconnected). *An IFS is totally disconnected if it possesses a nice compact set.*

**Definition 5** (Just touching). *An IFS is just touching if it possesses a nice open set.*

From here on, for the sake of clarity and convenience in notation, we will denote compact sets by the character  $C$  or  $K$ , open sets by  $U$  and general sets by  $A$ .

### 1.2.1 Examples

**The cantor set:** The cantor set  $K_{\text{cantor}}$  is the attractor of an IFS on  $\mathbb{R}$  defined by  $\mathcal{F} = \{f_1, f_2\}$ , where

$$f_1(x) = 1/3x \text{ and } f_2(x) = 1/3x + 2/3. \quad (7)$$

The IFS  $\mathcal{F}$  is totally disconnected since the set  $C_0 = [0, 1]$  is a nice compact set. Figure 1(a) shows the first few sets  $C_i$  obtained from this choice of  $C_0$ .

**Just touching** Figure 1(b)-(c) shows two examples of *just touching* IFS which generate the Sierpinski triangle and the Koch curve. A full description of these IFS is given in the Appendix for completeness.

### 1.3 Main results

Our goal is to use NNF to efficiently construct functions which represent sets  $A_k$  of the form (4). The classical choice of a function which represents a set  $A$  is the standard indicator function

$$\chi_A(x) = \begin{cases} 1, & x \in A \\ 0, & x \notin A \end{cases}. \quad (8)$$

However, this function cannot be realized as an NNF as it is not continuous. Instead, we define the notion of a CPwL indicator function

**Definition 6.** *We say that  $\varphi_C(x)$  is a CPwL indicator function of the compact set  $C$ , if*

$$\begin{cases} \varphi_C(x) \geq 0, & x \in C \\ \varphi_C(x) < 0, & x \notin C \end{cases}.$$

Note that a compact set  $C$  may admit many CPwL indicator functions  $\varphi_C$ . CPwL indicators can be used to approximate indicator functions to arbitrary precision due to the following simple lemma.

**Lemma 1.** *If  $\varphi_C \in \text{NNF}_{W,L}^{d \rightarrow 1}$  is a CPwL indicator function of  $C$ , then there exist  $\psi_n \in \text{NNF}_{\max\{W,2\},L+1}^{d \rightarrow 1}$  which converge to  $\chi_C$  pointwise and in  $L^p(\mathcal{K})$  for all  $1 \leq p < \infty$ .*

The proof of this lemma is given in Appendix A.

We now state our main result:

**Theorem 1.** *Let  $\mathcal{F} = \{f_1, \dots, f_J\}$  be an IFS on  $\mathbb{R}^d$ . Let  $C_0$  be a nice compact subset of  $\mathbb{R}^d$ , and assume*

1.  $C_0$  is a finite union of convex polytopes.
2.  $f_j$  are invertible affine transformations.

*Then, there exist constants  $W_0, L_0$ , such that for all  $k$ , there exists a CPwL indicator function  $\varphi_k \in \text{NNF}_{W_0, L_0 k}^{d \rightarrow 1}$  for the set  $C_k = H^k C_0$ .*

The constants  $W_0, L_0$  are given in Equation (31). They depend on  $d, J$ , the number  $p_0$  of convex polytopes which compose  $C_0$ , and the maximal number of hyperplanes  $m$  needed to express each of the polytopes.

One example for which Theorem 1 applies is the cantor set we discussed earlier, with the initial compact nice set  $C_0 = [0, 1]$ . The set  $C_0$  is an interval- a 1D polytope, and the functions  $f_1, f_2$  defined in (7) are affine and invertible. In this case  $d = 1, p_0 = 1, J = 2, m = 2$ , and so the theorem together with the parameter count in (31) shows that we can represent a CPwL indicator function for  $C_k$  as an NNF with width 10 and depth  $4k$ .

We note that by combining the theorem with Lemma 1 we also obtain that the indicator  $\chi_{C_k}$  can be approximated to machine precision using a NNF with almost identical complexity. Finally we note that as can be seen in (30) the DNN construction we propose is composed of an identical network applied  $k$  times to obtain a CPwL indicator of  $C_k$ , and so by using a recurrent NN architecture instead of the fully connected architecture we described here the number of parameters necessary to represent  $C_k$  can be completely independent of  $k$ .

We now turn to discuss *just touching* IFS; that is, IFS which possesses a nice open set  $U$ . In this case, we can use the mechanism of the proof of Theorem 1 to approximate  $\chi_{U_k}$  to arbitrary precision, though we cannot construct an exact CPwL indicator in this case:

**Theorem 2.** *Let  $\mathcal{F} = \{f_1, \dots, f_J\}$  be an IFS on  $\mathbb{R}^d$ . Let  $U_0$  be a nice open subset of  $\mathbb{R}^d$ , and assume*

1.  $U_0$  is a finite union of disjoint open convex polytopes.
2.  $f_j$  are invertible affine transformations.

*Then, there exist constants  $W_0, L_0$ , such that for all  $k$ , there exist a sequence of functions  $\psi_n \in \text{NNF}_{W_0, L_0 k+1}^{d \rightarrow 1}$  which converge to the indicator function of  $U_k = H^k U_0$  pointwise and in  $L^p(\mathcal{K})$  for all compact  $\mathcal{K}$  and  $1 \leq p < \infty$ .*

Finally, we note that any IFS on  $\mathbb{R}^d$  can be lifted to a totally disconnected IFS on  $\mathbb{R}^{d+1}$ . We discuss this simple construction in Appendix C.

The remainder of the paper is organized as follows: In Section 2 we discuss some theoretical preliminaries, which will be useful for the proof of Theorem 1, presented in Section 3. An outline of the proof of Theorem 2, the proof of Lemma 1, and the dimension argument mentioned in the introduction (i.e., the fact that DNNs with  $p$  parameters cannot span all CPwL with  $\ell \gg p$  linear regions), can be found in Appendix A.

## 2 Theoretical preliminaries

In this section we wish to define formally some basic actions on NNF as well as two NNF *building blocks* that will be used in our construction later on. The actions we wish to address are *composition*, *splitting*, and *summation*; and the building blocks are the commonly used *min* and *max* NNFs.

### 2.1 Basic actions

**Composition:** Using the aforementioned notation, the composition of two neural network functions is

$$f \in \text{NNF}_{w_1, d_1}^{n_1 \rightarrow n_2}, g \in \text{NNF}_{w_2, d_2}^{n_2 \rightarrow n_3} \Rightarrow g \circ f \in \text{NNF}_{\max(w_1, w_2), d_1 + d_2}^{n_1 \rightarrow n_3}. \quad (9)$$

Accordingly, the number of parameters in  $g \circ f$  is  $\mathcal{O}(\max(w_1, w_2)^2 \cdot (d_1 + d_2))$ .

**Splitting:** It is sometimes desirable to split a function into two separate operations

$$\begin{array}{ccc} & & f(x) \\ & \nearrow & \\ x & & \\ \cap & \searrow & \\ \mathbb{R}^n & & g(x) \end{array} \quad (10)$$

where  $f \in \text{NNF}_{w_1, d_1}^{n \rightarrow n_1}$ ,  $g \in \text{NNF}_{w_2, d_2}^{n \rightarrow n_2}$ . The NNF resulting from the action of splitting (10) is denoted by  $x_g^f$  and

$$x_g^f \in \text{NNF}_{w_1 + w_2, \max(d_1, d_2)}^{n \rightarrow (n_1 + n_2)}. \quad (11)$$

In the case where  $d_1 = d_2$  the definition of  $x_g^f$  is straightforward. The case  $d_1 < d_2$  can be reduced to the case  $d_1 = d_2$  by representing  $f$  using a network of depth  $d_2$ . This is done by shifting  $f$  by such a large value  $b$  so that  $f(\mathcal{K}) + b \geq 0$ , and so  $f + b$ , will be unaffected by ReLU, then applying identity mappings, and finally reapplying translation by  $-b$  to achieve the original function  $f$ .

**Summation:** Finally, given  $f \in \text{NNF}_{w_1, d_1}^{n_1 \rightarrow n_2}$  and  $g \in \text{NNF}_{w_2, d_2}^{n_1 \rightarrow n_2}$ , the function  $f + g$  is well defined and

$$f + g = \begin{array}{ccc} & & f(x) \\ & \nearrow & \searrow \\ x & & \\ \cap & & \\ \mathbb{R}^{n_1} & & g(x) \\ & \searrow & \nearrow \\ & & f(x) + g(x) \end{array} . \quad (12)$$

Thus,  $f + g$  is merely  $x_g^f$  with an added linear layer as output, and so

$$f + g \in \text{NNF}_{(w_1+w_2), \max(d_1, d_2)+1}^{n_1 \rightarrow n_2}. \quad (13)$$

## 2.2 min and max as NNF

Although common neural network architectures apply min and max operations (e.g., max-pooling), in order to maintain consistency with the standard definition of fully connected networks we define them here as NNF using ReLU. We can pronounce the max and min operation as

$$\max\{x, y\} = \text{ReLU}(x - y) + y, \text{ and } \min\{x, y\} = x - \text{ReLU}(x - y), \quad (14)$$

Thus, for  $x, y \in \mathbb{R}$

$$\min\{x, y\}, \max\{x, y\} \in \text{NNF}_{2,1}^{2 \rightarrow 1}. \quad (15)$$

For a vector  $x \in \mathbb{R}^\ell$  we define

$$\min(x) \stackrel{\text{def}}{=} \min\{x_1, \dots, x_\ell\}. \quad (16)$$

This minimum operation can be implemented via a *Divide and Conquer* type approach (e.g., see He and Zheng (2018), Theorem 3.1) with  $\log_2 \ell$  depth, and therefore

$$\min(x) \in \text{NNF}_{\ell, \lceil \log_2(\ell) \rceil}^{\ell \rightarrow 1}. \quad (17)$$

Furthermore, for  $x \in \mathbb{R}^\ell$  and  $y \in \mathbb{R}$  we denote the vector-scalar minimum operation by

$$m\{x, y\} \stackrel{\text{def}}{=} (\min\{x_1, y\}, \dots, \min\{x_\ell, y\}), \quad (18)$$

and it follows that

$$m\{x, y\} \in \text{NNF}_{\ell+1,1}^{\ell+1 \rightarrow \ell}. \quad (19)$$

## 3 Proof of Theorem 1

The outline of the proof is as follows:

1. Since  $C_0$  is a nice set with respect to the IFS, there exists a function  $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , which satisfies  $T \circ f_j(x) = x$ .
2. Then, we can establish the following relation (see Lemma 2)

$$x \in C_k \iff x, Tx, \dots, T^{k-1}x \in C_1. \quad (20)$$

3. Both  $T$  and  $\varphi_1$  (a CPwL indicator of  $C_1$ ) can be constructed as NNF (see Sections 3.1 and 3.2).
4. Then, from the construction of vector *min* as NNF and the fact that composition of NNFs is NNF (see Sections 2.1 and 2.2), we get that a CPwL indicator of  $C_k$  can be pronounced as the following NNF

$$\varphi_k(x) = \min(\varphi_1(x), \varphi_1(Tx), \dots, \varphi_1(T^{k-1}x)). \quad (21)$$

5. Using the basic actions' calculations for width and depth (Section 2.1) we come to the conclusion that

$$\varphi_k \in \text{NNF}_{W_0, L_0 k}^{d \rightarrow 1}, \quad (22)$$

where  $W_0, L_0$  are some constants independent of  $k$  (see Section 3.3).

Let us now show that if there exists a function  $T$ , as described in item 1 in the outline, then the equivalence relation of item 2 holds.

**Lemma 2.** *Let  $\{f_1, \dots, f_J\}$  be an IFS, and assume  $C_0$  is a nice set, and  $f_j C_0 \subseteq C_0, \forall j$ . Assume  $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$  satisfies  $T \circ f_j(x) = x$  for all  $x \in C_0$ , then*

$$x \in C_k \iff x, Tx, \dots, T^{k-1}x \in C_1.$$

*Proof.* Assume  $x \in C_k$ , then there exist  $j_1, j_2, \dots, j_k \in \{1, \dots, J\}$ , and  $y \in C_0$ , such that

$$x = f_{j_1} \circ f_{j_2} \dots \circ f_{j_k}(y)$$

Thus for all  $0 \leq r \leq k-1$

$$T^r x = f_{j_{r+1}} \circ f_{j_{r+2}} \dots \circ f_{j_k}(y) \in C_{k-r} \subseteq C_1$$

as required.

We prove the opposite direction by induction on  $k$ . For  $k=1$  the claim is obvious. Now assume the claim holds for  $k-1$ , we want to show that it holds for  $k$  as well- we assume  $x, Tx, \dots, T^{k-1}x \in C_1$  and we need to show that  $x \in C_k$ .

Since  $x \in C_1$ , there exists  $y_1 \in C_0$  and  $j_1 \in \{1, \dots, J\}$ , such that

$$x = f_{j_1}(y_1). \quad (23)$$

By applying  $T$  to this equation we obtain

$$Tx = y_1. \quad (24)$$

On the other hand by the induction hypothesis we have that  $Tx \in C_{k-1}$  so that there exists  $y_2 \in C_0$  and  $j_2, \dots, j_k \in \{1, \dots, J\}$  such that

$$Tx = f_{j_2} \circ f_{j_3} \dots \circ f_{j_k}(y_2) \quad (25)$$

and so it follows that  $x \in C_k$  since

$$x \stackrel{(23)}{=} f_{j_1}(y_1) \stackrel{(24)}{=} f_{j_1}(Tx) \stackrel{(25)}{=} f_{j_1} \circ f_{j_2} \circ f_{j_3} \dots \circ f_{j_k}(y_2)$$

□

### 3.1 Construction of $\varphi_1$

We now show how to construct a CPwL indicator function  $\varphi$  for the set  $C_1$ , using the same methodology as in Kovalsky et al. (2019). By assumption,  $C_0$  is a union of some  $p_0$  convex polytopes, each defined as an intersection of at most  $m$  half-spaces. Accordingly  $C_1$  is a finite union of  $p = Jp_0$  such polytopes.

A half-space  $H$  is defined by a linear function  $\varphi_H : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\varphi_H \stackrel{\text{def}}{=} \langle a, x \rangle + b,$$

and

$$H = \{x \mid \varphi_H(x) \geq 0\}.$$

Thus, by definition  $\varphi_H$  is a CPwL indicator of  $H$  and is a member of  $\text{NNF}_{1,0}^{d \rightarrow 1}$ . If  $\mathcal{P}$  is a convex polytope which is an intersection of at most  $m$  half-spaces  $H_i$ , then we obtain a CPwL indicator for  $\mathcal{P}$  via

$$\varphi_{\mathcal{P}}(x) = \min(\varphi_{H_1}(x), \varphi_{H_2}(x), \dots, \varphi_{H_m}(x)).$$

Note that  $\varphi_{\mathcal{P}} \in \text{NNF}_{m, \lceil \log_2 m \rceil}^{d \rightarrow 1}$  due to the splitting, composition and minimum rules from Section 2.

Now, if a set  $C$  is the union of  $p$  polytopes  $\mathcal{P}_i$ , each defined by at most  $m$  half-spaces, then a CPwL indicator  $\varphi_C$  for  $C$  is given by

$$\varphi_C(x) = \max(\varphi_{\mathcal{P}_1}(x), \varphi_{\mathcal{P}_2}(x), \dots, \varphi_{\mathcal{P}_p}(x)),$$

which is in  $\text{NNF}_{mp, \lceil \log_2 m \rceil + \lceil \log_2 p \rceil}^{d \rightarrow 1}$  due to the splitting, composition and maximum rules from Section 2. This is the class of function  $\varphi_1 = \varphi_{C_1}$  belongs to, when we set  $p = Jp_0$ .

### 3.2 Construction of $T$ as NNF

We now construct  $T$  as an NNF.  $T$  will be of the form  $T = \sum_{j=1}^J T_j$  where each  $T_i$  will be an NNF satisfying

$$T_i(x) = \delta_{ij} f_i^{-1}(x), \quad \forall x \in \bigcup_{j=1}^J f_j C_0 \quad (26)$$

Let,

$$M = \|f_i^{-1}\|_{L^\infty(C_0)}, \text{ and } -\delta = \max_{x \in \bigcup_{j \neq i} f_j C_0} \varphi_{f_i C_0}(x) < 0,$$

and  $\mathbf{1}_d$  denotes the constant vector with unit value in all  $d$  entries. Recall that  $f_i^{-1}$  is defined for all points in  $\mathbb{R}$ , as it is the inverse of an affine transformation. Then, we construct  $T_i$  as NNF according to the following scheme

$$\begin{array}{ccccccc} & \nearrow & & & \searrow & & \\ x & \rightarrow & \boxed{\begin{array}{c} f_i^{-1}(x) \\ \underbrace{\varphi_{f_i C_0}(x)} \\ y_1 \in \mathbb{R} \end{array}} & \rightarrow & \boxed{\begin{array}{c} f_i^{-1}(x) \\ \frac{4M}{\delta}(y_1 + \delta/2) \\ y_2 \in \mathbb{R} \end{array}} & \rightarrow & \boxed{\begin{array}{c} m\{f_i^{-1}(x), y_2\} \\ y_3 \in \mathbb{R}^d \end{array}} & \rightarrow \dots \\ & & & & & & \\ \dots & \rightarrow & \boxed{\text{ReLU}(2(y_3 + M\mathbf{1}_d)) - \text{ReLU}(y_3 + 2M\mathbf{1}_d)} & & & & \end{array} \quad (27)$$

For convenience, we enumerate the sequential operations (marked by the boxes) by  $\ell_1, \dots, \ell_4$ . We claim that  $T_i$  constructed in this form fulfills (26). Note that the magnitudes of the coordinates of  $f_i^{-1}(x)$  are smaller than  $M$  for all  $x \in C_0$ . By construction,  $y_2 > 2M$  if  $x \in f_i(C_0)$  and is smaller than  $-2M$  if  $x \in f_j(C_0)$  for  $j \neq i$ . It then follows that if  $x \in f_i(C_0)$  then  $y_3 = f_i^{-1}(x)$  and if  $x \in f_j(C_0)$  for  $j \neq i$  then  $y_3 \leq -2M$  elementwise. Finally by applying the function

$$\text{ReLU}(2(t + M)) - \text{ReLU}(t + 2M) = \begin{cases} 0 & \text{if } t \leq -2M \\ -t - 2M & \text{if } -2M \leq t \leq -M \\ t & \text{if } t \geq -M \end{cases} \quad (28)$$

elementwise to  $y_3$ , we obtain  $T_i(x)$  which satisfies (26). Using the parameter counting rules from Section 2 we have

$$\begin{aligned} \ell_1 &= x_{\varphi_{f_i C_0}}^{f_i^{-1}} \in \text{NNF}_{d+mp_0, \lceil \log_2 m \rceil + \lceil \log_2 p_0 \rceil}^{d \rightarrow d+1} \\ \ell_2 &\in \text{NNF}_{d+1,0}^{d+1 \rightarrow d+1}; \quad \ell_3 \in \text{NNF}_{2d,1}^{d+1 \rightarrow d}; \quad \ell_4 \in \text{NNF}_{2d,1}^{d \rightarrow d} \end{aligned}$$

and so using the composition we have that  $T_i \in \text{NNF}_{1/JW_T, L_T}^{d \rightarrow d}$  for

$$W_T = J \max\{d + mp_0, 2d\}, \quad L_T = 2 + \lceil \log_2 m \rceil + \lceil \log_2 p_0 \rceil, \quad (29)$$

and by using the summation rule to count the parameters of  $T$  we obtain that  $T \in \text{NNF}_{W_T, L_T}^{d \rightarrow d}$ .

### 3.3 Construction of $C_k$ through NNF

To conclude the proof of the theorem, we show how to construct a CPwL indicator of  $C_k$  as an NNF based on (21) and our construction of  $T$  and  $\varphi$ . We set

$$m_0(x) = \varphi_1(x), \quad m_j(x) = m_j(m_{j-1}, T^j x) = \min\{m_{j-1}(x), \varphi_1(T^j x)\}, \quad 1 \leq j \leq k-1.$$

and note that  $m_j \in \text{NNF}_{mp, \lceil \log_2 m \rceil + \lceil \log_2 p \rceil + 1}^{d+1 \rightarrow 1}$  while  $m_0$  has one layer less. Note that the indicator  $\varphi_k$  from (21) is equal to  $m_{k-1}$ , and so  $\varphi_k$  can be computed by the following NNF:

$$\begin{array}{ccccccc} & & Tx & \rightarrow & T^2x & \rightarrow & \dots & T^{k-1}x \\ & \nearrow & & & \searrow & & & \\ x & \rightarrow & m_0(x) & \rightarrow & m_1(x) & \rightarrow & \dots & m_{k-2}(x) & \rightarrow & m_{k-1}(x) \end{array} \quad (30)$$

The width of  $m_j$  is the same as the width of  $\varphi_1$  and it is one layer deeper, therefore  $\varphi_k(x) = m_{k-1}(x)$  can be constructed as an NNF with constant width  $W_0$  and depth  $kL_0$ , where these constants are defined by

$$W_0 = J \max\{d + mp_0, 2d\} + Jmp_0, \quad L_0 = \lceil \log_2 m \rceil + \lceil \log_2(Jp_0) \rceil + 2. \quad (31)$$

## 4 Future work

Our results can be interpreted as showing that black and white images of 2D fractals- the indicator functions of  $C_k$ - can be compressed efficiently by DNN. DNN are also known to provide efficient construction of wavelets Grohs et al. (2019). As both IFS and wavelets have been used successfully to compress images we believe it may be possible to devise improved compression algorithms based on DNNs. We intend to investigate this in the near future.

## References

- Arora, R., Basu, A., Mianjy, P., and Mukherjee, A. (2018). Understanding deep neural networks with rectified linear units. In *International Conference on Learning Representations*.
- Bahat, Y. and Irani, M. (2016). Blind dehazing using internal patch recurrence. In *2016 IEEE International Conference on Computational Photography (ICCP)*, pages 1–9. IEEE.
- Barnsley, M. F. (2014). *Fractals everywhere*. Academic press.
- CireşAn, D., Meier, U., Masci, J., and Schmidhuber, J. (2012). Multi-column deep neural network for traffic sign classification. *Neural networks*, 32:333–338.
- Cohen, N., Sharir, O., and Shashua, A. (2016). On the expressive power of deep learning: A tensor analysis. In *Conference on Learning Theory*, pages 698–728.
- Collobert, R. and Weston, J. (2008). A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM.
- Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314.
- Daubechies, I., DeVore, R., Foucart, S., Hanin, B., and Petrova, G. (2019). Nonlinear approximation and (deep) relu networks.
- Eldan, R. and Shamir, O. (2016). The power of depth for feedforward neural networks. In *Conference on learning theory*, pages 907–940.
- Falconer, K. (2004). *Fractal geometry: mathematical foundations and applications*. John Wiley & Sons.
- Grohs, P., Perekrestenko, D., Elbrächter, D., and Bölcskei, H. (2019). Deep neural network approximation theory. *CoRR*, abs/1901.02220.
- He, J., L. L. X. J. and Zheng, C. (2018). Relu deep neural networks and linear finite elements. *arXiv preprint arXiv:1807.03973*.
- Hornik, K., Stinchcombe, M., and White, H. (1989). Multilayer feedforward networks are universal approximators. *Neural Netw.*, 2(5):359–366.
- Irani, D. G. S. B. M. (2009). Super-resolution from a single image. In *Proceedings of the IEEE International Conference on Computer Vision, Kyoto, Japan*, pages 349–356.
- Jacquin, A. E. (1992). Image coding based on a fractal theory of iterated contractive image transformations. *IEEE transactions on image processing*, 1(1):18–30.
- Kovalsky, S., Dov, D., and Dym, N. (2019). Poisson reconstruction over neural network function spaces. *Technical report*.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105.
- Liang, S. and Srikant, R. (2016). Why deep neural networks for function approximation? *arXiv preprint arXiv:1610.04161*.

- Manning, C., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S., and McClosky, D. (2014). The stanford corenlp natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations*, pages 55–60.
- Michaeli, T. and Irani, M. (2013). Nonparametric blind super-resolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 945–952.
- Montúfar, G. (2017). Notes on the number of linear regions of deep neural networks.
- Montufar, G. F., Pascanu, R., Cho, K., and Bengio, Y. (2014). On the number of linear regions of deep neural networks. In *Advances in neural information processing systems*, pages 2924–2932.
- Piczak, K. J. (2015). Environmental sound classification with convolutional neural networks. In *2015 IEEE 25th International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6.
- Poggio, T., Mhaskar, H., Rosasco, L., Miranda, B., and Liao, Q. (2017). Why and when can deep-but not shallow-networks avoid the curse of dimensionality: a review. *International Journal of Automation and Computing*, 14(5):503–519.
- Qi, C. R., Su, H., Mo, K., and Guibas, L. J. (2017). Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 652–660.
- Sabetsarvestani, Z., Sober, B., Higgitt, C., Daubechies, I., and Rodrigues, M. R. D. (2019). Artificial intelligence for art investigation: Meeting the challenge of separating x-ray images of the ghent altarpiece. *Science Advances*.
- Serra, T., Tjandraatmadja, C., and Ramalingam, S. (2017). Bounding and counting linear regions of deep neural networks. *arXiv preprint arXiv:1711.02114*.
- Smale, S., Rosasco, L., Bouvrie, J., Caponnetto, A., and Poggio, T. (2010). Mathematics of the neural response. *Foundations of Computational Mathematics*, 10(1):67–91.
- Van Pabst, J. V. L. and Jense, H. (1996). Dynamic terrain generation based on multifractal techniques. In *High Performance Computing for Computer Graphics and Visualisation*, pages 186–203. Springer.
- Welstead, S. T. (1999). *Fractal and wavelet image compression techniques*. SPIE Optical Engineering Press Bellingham, Washington.
- Yarotsky, D. (2017). Error bounds for approximations with deep relu networks. *Neural Networks*, 94:103–114.

## A Additional proofs

**Proof of Theorem 2.** Looking at the outline of the proof of Theorem 1, the only item that needs to be adapted to the case of open sets is the construction of  $T$ , as it is not clear how we should define it on the boundary of the domains  $f_j U_0$ . Assume w.l.o.g. that  $0 \notin \bar{U}_1$ . We define  $T$  to be a (non-continuous) function which is equal to  $f_j^{-1}(x)$  if  $x \in f_j U_0$  for some  $j$ , and is equal to zero otherwise. Note that Lemma 2 does not assume compactness of  $C_0$ , and so the sufficient and necessary conditions of item 2 of the outline holds for the non-continuous  $T$ .

Let  $\varphi$  be the CPwL indicator function of  $\bar{U}_1$ , which is constructed in Subsection 3.1, and note that  $\varphi$  satisfies that  $x \in U_1$  if and only if  $\varphi(x) > 0$ . Thus we can use (21) to define a (non-continuous) function  $\varphi_k$  which satisfies  $\varphi_k(x) > 0$  if and only if  $x \in U_k$ .

Next we can use the same construction as in Subsection 3.2 to define for  $\delta > 0$  CPwL maps  $T_i^\delta(x)$  which are zero if  $\varphi_i(x) \leq 0$ , and are equal to  $f_i^{-1}$  if  $\varphi_i(x) \geq \delta$ . We take the sum of these maps to obtain a CPwL map  $T^\delta$  which converge pointwise to  $T$  as  $\delta \rightarrow 0$ . Next we replace  $T$  with  $T^\delta$  in the definition of  $\varphi_k$  in (21) to obtain CPwL functions  $\varphi_k^\delta$  which converge pointwise to  $\varphi_k$  as  $\delta \rightarrow 0$ .

To obtain pointwise convergence to the indicator of  $U_k$  we need to compose  $\varphi_k^\delta$  with a simple scalar function: For  $a < b$  we define the function

$$f_{a,b}(x) = \frac{1}{b-a} (\text{ReLU}(x-a) - \text{ReLU}(x-b)) = \begin{cases} 0 & \text{if } x \leq a \\ \frac{x-a}{b-a} & \text{if } a \leq x \leq b \\ 1 & \text{if } x \geq b \end{cases} \quad (32)$$

Now we set  $a = \varphi_k(0)$  which is negative since  $0 \notin \bar{U}_1$ , and  $b = 0$ , and obtain that  $\psi_\delta = f_{a,b} \circ \varphi_k^\delta$  converges pointwise to the indicator of  $U_k$  as  $\delta \rightarrow 0$ . The  $L^p(\mathcal{K})$  convergence argument is identical to the argument in Lemma 1  $\square$

**Proof of Lemma 1.** By choosing  $f = f_{a,b}$  as defined in (32) with  $a = -2, b = -1$ , we obtain that  $\psi_t(x) = f \circ \varphi_C(tx)$  converges elementwise to  $\chi_C$  as  $t \rightarrow \infty$ . For a fixed compact set  $\mathcal{K}$  and fixed  $1 \leq p < \infty$ , we show  $L^p(\mathcal{K})$  convergence by showing that the integral of the functions

$$g_t(x) = |\psi_t(x) - \chi_C(x)|^p$$

converge to zero. This follows easily from the dominated convergence theorem since  $g_t$  converge pointwise to zero and are bounded uniformly for all  $t$  by the constant function 1 which is integrable in  $\mathcal{K}$ .  $\square$

### Dimension argument

**Lemma 3.** *If  $\ell$  is larger than the number of parameters  $p$  defining  $\text{NNF}_{W,L}^{1 \rightarrow 1}$ , then not all CPwL functions with  $\ell$  knots are contained in  $\text{NNF}_{W,L}^{1 \rightarrow 1}$ .*

*Proof.* Denote the set of parameters by  $\theta$ , and denote by  $f_\theta$  the NNF defined through  $\theta$ . Choose any  $t_1, t_2, \dots, t_\ell \in \mathbb{R}$ . It is sufficient to show that the function  $F : \mathbb{R}^p \rightarrow \mathbb{R}^\ell$  defined by

$$F_i(\theta) = f_\theta(t_i), i = 1, \dots, \ell$$

is not onto. Indeed,  $F$  is continuous as the parameter space  $\mathbb{R}^p$  can be partitioned into a finite number of sets on which  $F$  is a multivariate polynomial whose degree is determined by the degree of the network (for example, the composition of two affine transformations results with a second order polynomial in the coefficients). It follows that the restriction of  $F$  to a closed ball  $\bar{B}_N$  of radius  $N$  is Lipschitz, and therefore it is known (see Chapter 2 in Falconer (2004)) that its image has Hausdorff dimension  $\leq p$ . The image of  $F$  is the countable union of  $F(B_N)$ , hence it has Hausdorff dimension  $\leq p < \ell$  as well and so  $F$  is not onto.  $\square$

## B IFS examples

**The Sierpinski triangle:** The Sierpinski triangle  $K_{\text{Sierpinski}}$  is the attractor of an IFS on  $\mathbb{R}^2$ . It is defined using an equilateral triangle  $T$  with vertices  $v_1, v_2, v_3 \in \mathbb{R}^2$ , where  $\mathcal{F} = \{f_1, f_2, f_3\}$  are defined to be

$$f_i(x) = 1/2x + v_i, \quad i = 1, 2, 3.$$

The open triangle  $U_0 = T^\circ$  is a nice open set. Figure 1(b) shows the first few sets  $U_i$  obtained from this choice of  $U_0$ .

**The Koch curve:** The Koch curve  $K_{\text{Koch}}$  is the attractor of an IFS on  $\mathbb{R}^2$ . It is defined via four maps  $f_i(x) = A_i(x) + b_i$  where

$$A_1 = \frac{1}{3} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, b_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad A_2 = \begin{bmatrix} 1/6 & -\sqrt{3}/6 \\ \sqrt{3}/6 & 1/6 \end{bmatrix}, b_2 = \begin{bmatrix} 1/3 \\ 0 \end{bmatrix}$$

$$A_3 = A_2^T, b_3 = \begin{bmatrix} 1/2 \\ \sqrt{3}/6 \end{bmatrix}, \quad A_4 = A_1, b_4 = \begin{bmatrix} 2/3 \\ 0 \end{bmatrix}.$$

The open triangle with vertices

$$v_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, v_2 = \begin{bmatrix} 1/2 \\ \sqrt{3}/2 \end{bmatrix}, v_3 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

is a nice open set w.r.t. this IFS. Figure 1(c) shows the first few sets  $U_i$  obtained from this choice of  $U_0$ .

## C Lifting

Any IFS on  $\mathbb{R}^d$  composed of affine and invertible functions can be lifted to an IFS on  $\mathbb{R}^{d+1}$  fulfilling the conditions of 1. Assume we are given such an IFS  $\{f_1, \dots, f_J\}$  on  $\mathbb{R}^d$  and an initial convex polytope  $C_0$  fulfilling (5) (such a set always exists). We can then define a new IFS  $\{\hat{f}_1, \dots, \hat{f}_J\}$  on  $\mathbb{R}^{d+1}$  by

$$\hat{f}_j(x, t) = \left( f_j(x), \frac{1}{3(J+1)}t + \frac{j}{J+1} \right).$$

The convex polytope  $\hat{C}_0 = C_0 \times [0, 1]$  is a nice compact subset of  $\mathbb{R}^{d+1}$ , and so we can use Theorem 1 to construct the indicator function of  $\hat{C}_k$ , which is related to the indicator function of  $C_k$  via

$$\chi_{C_k}(x) = \max_{t \in [0, 1]} \chi_{\hat{C}_k}(x, t)$$