

Discontinuous Neural Networks and Discontinuity Learning

Original

Discontinuous Neural Networks and Discontinuity Learning / Della, Santa; Pieraccini, Sandra. - In: JOURNAL OF COMPUTATIONAL AND APPLIED MATHEMATICS. - ISSN 0377-0427. - ELETTRONICO. - 419:(2023), p. 114678. [10.1016/j.cam.2022.114678]

Availability:

This version is available at: 11583/2851043 since: 2022-08-30T12:01:39Z

Publisher:

Elsevier

Published

DOI:10.1016/j.cam.2022.114678

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

Discontinuous Neural Networks and Discontinuity Learning

Francesco Della Santa ^{*†‡§}, Sandra Pieraccini ^{¶†}

September 17, 2021

Abstract

In the framework of discontinuous function approximation and discontinuity interface detection, we consider an approach involving Neural Networks. In particular, we define a novel typology of Neural Network layers endowed with new learnable parameters and discontinuities in the space of the activations. These layers allow to create a new kind of Neural Networks, whose main property is to be discontinuous, able not only to approximate discontinuous functions but also to learn and detect the discontinuity interfaces. A sound theoretical analysis concerning the properties of the new discontinuous layers is performed, and some tests on discontinuous functions are proposed, in order to assess the potential of such instruments.

Keywords: Discontinuous functions; neural networks; deep learning; automatic detection of discontinuity interface.

MSC: 68T07, 65D40, 65D99

1 Introduction

In this work, we introduce a new typology of layers for Neural Networks (NNs); the novelty is given by the introduction of discontinuities in the layer's characterizing function and, consequently, in the Neural Network. In the framework of NNs, discontinuities were involved in the first mathematical models of biological neurons, dating back to the 1940s [19], and the first Neural Networks proposed in the 1950s and 1960s [24, 27]. In these models, the activation functions of the NN units were mainly inspired by the mechanisms of the biological neurons and, therefore, were modeled using the Heaviside step function, or suitable variants. Then, the activation functions evolved into the continuous (and often smooth) ones used nowadays in almost all the deep learning algorithms (see [8, ch. 6.2.2, 6.3] and [2]), thanks to the advantages that they grant in adapting the parameters of a NN. To the best of authors' knowledge, recent literature does not report examples of practical use of discontinuous NN layers or discontinuous NNs; nonetheless, a renewed interest on discontinuous activation functions, at least from the theoretical point of view, is witnessed by very recent works, see e.g. [25], considering the floor function $f(x) = \lfloor x \rfloor$ as activation function, and [21], using the Heaviside function.

We are interested in introducing discontinuities in NN learning models, aiming at using feedforward NNs to approximate discontinuous functions and, at the same time, detect the discontinuity interfaces. This latter problem is quite a challenging task, especially for functions with a high-dimensional domain. Moreover, the information can be quite relevant in several applications. To mention an example, the smoothness of the target function can critically affect the behavior of

*corresponding author (francesco.dellasanta@polito.it)

†Dipartimento di Scienze Matematiche, Politecnico di Torino

‡Member of the INdAM-GNCS group

§SmartData@PoliTO, Politecnico di Torino

¶Dipartimento di Ingegneria Meccanica e Aerospaziale, Politecnico di Torino

numerical methods for stochastic collocation in the framework of uncertainty quantification; thus, knowing the discontinuity interfaces and being able to partition the function domain in several regions in which the function is smooth, can be of paramount importance (see, e.g., [13] and references therein).

In the past decades, feedforward NNs have been used mainly for classification-type tasks [16, 9, 26] but they can perform very well also regression tasks, as guaranteed by the universal approximation theorems [17, 23, 14, 21]. In particular, we recall that the universal approximation theorem of Leshno et al. [17] is guaranteed also for discontinuous activation functions, while Park et al. recently showed (see [21, Th. 3]) that a NN with `relu` and Heaviside activation functions is dense in the space of continuous functions from a compact set $K \subset \mathbb{R}^n$ to \mathbb{R}^m .

Concerning the approximation of discontinuous functions with NNs, interesting results have been recently obtained in [22, 11]. A precursor of the use of NN in this framework can be found in [5]. In such a paper the author, leveraging the interpretation of shallow NNs as the superposition of ridge functions, uses a general construction based on ridgelets to build shallow NNs, thus presenting a tool for approximating target functions with spatial inhomogeneities. Functions with linear discontinuity interfaces are well approximated with such an approach, but the method is yet not satisfying for curvilinear interfaces, as stated by the author. Overall, the existing approximation methods based on NNs are not suitable to effectively tackle general discontinuous functions and simultaneously detecting their discontinuity interfaces; indeed, approximating a discontinuous function with a NN is quite a simple task, but the function represented by the NN will actually be continuous, due to the continuous activation functions of the NN layers.

As far as the discontinuity detection problem is concerned, it is a quite challenging task; the main results have been proposed in the last decades, see e.g. [4, 3, 12, 28]. In [4] a polynomial edge detection method is proposed: the discontinuous interfaces of a piece-wise smooth function $F : \mathbb{R}^n \rightarrow \mathbb{R}$, $n \leq 2$, are identified through the reconstruction of the jump function, given a set of function evaluations; the method proposed in [4] is extended in [3] to higher dimensions by applying the detection method for each input dimension to a generalized polynomial chaos approximation of the target function. Nonetheless, the method in [3] suffers the curse of dimensionality, setting practical restrictions on the dimensionality that can be handled; an improvement of [3] is proposed in [12], that exploits sparse grids to develop an adaptive method that increases the possibilities to be used in higher dimensions. The method proposed in [28] is based on the approximation of the hypersurface representing the discontinuity interface with hyper-spherical coordinates, and it is well suited also for large n , but the method is designed for detecting a single interface, which is assumed to satisfy the star-convexity assumption. In general, a discontinuity interface detection method can be generalized to the case of functions $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with $m > 1$, through the common practice of applying m times the method for functions with codomain of dimension one.

The present work aims at building new discontinuous NNs able to approximate discontinuous functions with other discontinuous functions, whose discontinuity interfaces are relatively easy to be detected. More specifically, the new NNs are endowed with trainable discontinuity jump parameters that allow the model to learn the discontinuity interface of the target function $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m$; then, analyzing the function compositions that define the NN, the discontinuity interfaces and the continuity regions of \mathbf{F} in the domain can be characterized. The main advantages of this method are that the NN both returns an approximation of \mathbf{F} and an approximation of its discontinuity interfaces. Moreover, by a theoretical point of view, there are no restrictions on the applicability of the method proposed, indeed: i) the method does not require assumptions about the regularity of \mathbf{F} (at least, \mathbf{F} must be approximable according to one of the universal approximation theorems); ii) there are no theoretical restrictions on the dimensions (there are only practical restrictions, due to the curse of dimensionality problem suffered by the regression task).

The work is organized as follows. In the next subsection, the main notations used herein are listed. In Section 2, the new discontinuous layer for NNs is presented. In Section 3, all the theoretical results that characterize a discontinuous NN are described. Section 4 illustrates numerical results on some examples assessing the potential of the new discontinuous NNs. We end with some conclusions drawn in Section 5.

1.1 Notation and basic results

In this subsection we introduce some useful notations and simple results for the following sections. First, we introduce the notation adopted for the description of the inner operations taking place in a fully-connected layer of a *Multi Layer Perceptron* (MLP). Then, we introduce some notations that will be useful in Section 3 to analyze discontinuities in NNs.

1.1.1 Notation for Neural Networks

Let N be an H -layers perceptron; i.e., N is an MLP characterized by $H \in \mathbb{N}$ hidden layers. We use the following notation to describe its architecture and the related mathematical entities:

- L_1, \dots, L_H denote the hidden layers; L_0 and L_{H+1} are the *input* and *output* layers, respectively;
- for each $h = 0, \dots, H+1$, $N_h \in \mathbb{N}$ is the number of units of layer L_h ;
- for each $h = 0, \dots, H$, $W^{(h+1)} \in \mathbb{R}^{N_h \times N_{h+1}}$ is the matrix of weights between layers L_h and L_{h+1} and $\mathbf{b}^{(h+1)} \in \mathbb{R}^{N_{h+1}}$ is the vector of biases of layer L_{h+1} ;
- $\mathbf{f}_h : \mathbb{R}^{N_h} \rightarrow \mathbb{R}^{N_h}$ denotes the element-wise application of the activation function $f_h : \mathbb{R} \rightarrow \mathbb{R}$ used in layer L_h ;
- for each $h = 0, \dots, H$, $\mathcal{L}_{h+1} : \mathbb{R}^{N_h} \rightarrow \mathbb{R}^{N_{h+1}}$ denotes the characterizing function of the fully-connected layer L_{h+1} , i.e. the map from the outputs of layer L_h to the ones of layer L_{h+1} , defined as

$$\mathcal{L}_{h+1}(\mathbf{x}^{(h)}) = \mathbf{f}_{h+1} \left(W^{(h+1)T} \mathbf{x}^{(h)} + \mathbf{b}^{(h+1)} \right), \quad \text{with } \mathbf{x}^{(h)} \in \mathbb{R}^{N_h}. \quad (1)$$

- we generalize the notation used for maps between pairs of layers to sequences of layers as follows:

$$\mathcal{L}_{h_1}^{h_2} := \mathcal{L}_{h_2} \circ \mathcal{L}_{h_2-1} \circ \dots \circ \mathcal{L}_{h_1+1} \circ \mathcal{L}_{h_1},$$

for each $1 \leq h_1 < h_2 \leq H+1$. Note that, following this notation, \mathcal{L}_1^{H+1} is the characterizing function $\widehat{\mathbf{F}}$ of the perceptron N and $\mathcal{L}_{h_1}^{h_2}$ is the function characterizing a sub-NN of N given by layers $L_{h_1-1}, L_{h_1}, \dots, L_{h_2-1}, L_{h_2}$.

Remark 1.1. Equation (1) characterizes the action of the so-called fully-connected layers; however it is easy to prove that it can describe also the action of convolutional layers (e.g., see [18]) or the connection of layers that are not fully-connected, setting to zero specific elements of the weight matrix. Then, formula (1) can be used as representative of the general characterizing function of a NN layer. Analogously, almost any feedforward (i.e., non-recurrent) NN can be represented by an equivalent MLP or, at most, by a composition of MLPs; then, in this work, we use the H -layers perceptron N as representative of a generic feedforward NN.

1.1.2 Notation for hyperplanes and corresponding partitions

Let $\Pi = \{\Pi_1, \dots, \Pi_m\}$ be a set of hyperplanes in \mathbb{R}^n , each one characterized by the equation $\mathbf{x}^T \mathbf{w}_j + b_j = 0$, for $j = 1, \dots, m$. Then, we will use the following notation to denote some special subsets of \mathbb{R}^n characterized by Π_1, \dots, Π_m :

- for each pair of disjoint subsets of hyperplanes $\{\Pi_{i_1}, \dots, \Pi_{i_s}\}, \{\Pi_{k_1}, \dots, \Pi_{k_t}\} \subset \Pi$ we denote by $C(\{\Pi_{i_1}, \dots, \Pi_{i_s}\}; \{\Pi_{k_1}, \dots, \Pi_{k_t}\})$ the subset of vectors $\mathbf{x} \in \mathbb{R}^n$ such that

$$\begin{cases} \mathbf{x}^T \mathbf{w}_i + b_i \geq 0 & \forall i = i_1, \dots, i_s \\ \mathbf{x}^T \mathbf{w}_k + b_k < 0 & \forall k = k_1, \dots, k_t \end{cases}.$$

We observe that, if $C(\{\Pi_{i_1}, \dots, \Pi_{i_s}\}; \{\Pi_{k_1}, \dots, \Pi_{k_t}\})$ is not empty, then it is convex.

- the set Π generates a partition $\mathcal{C}(\Pi)$ of convex subsets of \mathbb{R}^n defined as

$$\mathcal{C}(\Pi) = \{C(P; P^C) \mid P \in \mathcal{P}(\Pi)\} \setminus \{\emptyset\}, \quad (2)$$

where P^C is the complement of P in Π and $\mathcal{P}(\Pi)$ is the power set of Π .

Remark 1.2 (Special cases). Let $\Pi = \{\Pi_1, \dots, \Pi_m\}$ be a set of hyperplanes in \mathbb{R}^n and $\mathcal{C}(\Pi)$ the partition (2). Then, the following special cases may occur:

1. Let $\Pi_i, \Pi_j \in \Pi$, $i \neq j$, be such that $\Pi_i = \Pi_j$ and $\mathbf{w}_i = a\mathbf{w}_j$, $b_i = ab_j$, for an $a \in \mathbb{R} \setminus \{0\}$. If $a < 0$, for each $P \in \mathcal{P}(\Pi)$ such that $\Pi_i, \Pi_j \in P$, we have that the set $C(P; P^C)$ lies on the hyperplane $\Pi_i = \Pi_j$, while the set $C(P^C; P)$ is empty.
2. Let us admit as possible elements of Π also the degenerate hyperplanes $\Pi_0 = \mathbb{R}^n$ and $\Pi_\emptyset = \emptyset$ defined by equations $\mathbf{x}^T \mathbf{0} + 0 = 0$ and $\mathbf{x}^T \mathbf{0} + b_\emptyset = 0$, respectively, where $b_\emptyset \neq 0$. Then, if $\Pi_0 \in \Pi$ and/or $\Pi_\emptyset \in \Pi$, we have that $\mathcal{C}(\Pi)$ is still a partition of \mathbb{R}^n in convex subsets and, in particular, $\mathcal{C}(\Pi) = \mathcal{C}(\Pi \setminus \{\Pi_0, \Pi_\emptyset\})$.

Each element $X \in \mathcal{C}(\Pi)$ can be identified by a unique vector with elements in $\{0, 1\}$, as highlighted by the following definition.

Definition 1.1 (Region Vectors and Region Function). Let $\Pi = \{\Pi_1, \dots, \Pi_m\}$ be a set of m hyperplanes in \mathbb{R}^n , possibly including the degenerate cases Π_0 and Π_\emptyset , each one characterized by the equation $\mathbf{x}^T \mathbf{w}_j + b_j = 0$, for $j = 1, \dots, m$. Let $g: \mathbb{R}^n \rightarrow \{0, 1\}^m$ be the function

$$g(\mathbf{x}) := \mathcal{H}(W^T \mathbf{x} + \mathbf{b}), \quad (3)$$

where \mathcal{H} denotes the component-wise application of the Heaviside function

$$\mathcal{H}(x) = \begin{cases} 1, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

and $W = [\mathbf{w}_1, \dots, \mathbf{w}_m] \in \mathbb{R}^{n \times m}$ and $\mathbf{b} = [b_1, \dots, b_m]^T \in \mathbb{R}^m$. The function g defined in (3) is called region function associated to Π , and $g(\mathbf{x}) \in \{0, 1\}^m$ is called region vector of \mathbf{x} .

The region function g introduced in Definition 1.1 characterizes uniquely the subsets of \mathbb{R}^n of the partition $\mathcal{C}(\Pi)$, as stated in the following Lemma, whose proof is straightforward.

Lemma 1.1. Let Π be a set of m hyperplanes as in Definition 1.1. Then, for each pair of vectors $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$ such that $\mathbf{x}_1 \in X_1, \mathbf{x}_2 \in X_2$ with $X_1, X_2 \in \mathcal{C}(\Pi)$, it holds that $X_1 = X_2$ if and only if $g(\mathbf{x}_1) = g(\mathbf{x}_2)$.

Due to Lemma 1.1, each $X_i \in \mathcal{C}(\Pi)$ is uniquely identified by a vector $\mathbf{k}_i \in \{0, 1\}^m$ such that $\mathbf{k}_i = g(\mathbf{x})$, for each $\mathbf{x} \in X_i$.

Definition 1.2 (Region Vectors of Subsets). Let Π be as in Definition 1.1, and let g be the region function associated to Π . Let $\mathbf{k}_i \in \{0, 1\}^m$ be a vector such that $\mathbf{k}_i = g(\mathbf{x})$ for each $\mathbf{x} \in X_i$, given a fixed $X_i \in \mathcal{C}(\Pi)$. Then \mathbf{k}_i is called Region Vector of X_i with respect to the hyperplanes of Π .

We end this section with an example of a partition $\mathcal{C}(\Pi)$ of \mathbb{R}^2 and the corresponding region vectors; this example is illustrated in Figure 1. Six hyperplanes $\Pi = \{\Pi_1, \dots, \Pi_6\}$ are considered, and some special cases are also included, as we have $\Pi_1 = \Pi_2$ and $\Pi_3 = -\Pi_4$; these situations correspond to the one discussed in Remark 1.2, item 1, with $a > 0$ and $a < 0$, respectively.

In this example, for each $X_i, X_j \in \mathcal{C}(\Pi)$, we observe that the region vectors $\mathbf{k}_i, \mathbf{k}_j$, are strictly related to subsets connection. For example, looking at X_1 and X_3 we observe that their boundaries do intersect in the point given by the intersection $\bigcap_{k=1}^4 \Pi_k$, which are the hyperplanes identified by elements of \mathbf{k}_1 and \mathbf{k}_3 that are different. In other cases, such as X_3 and X_6 or X_1 and X_{10} , we observe that the intersection of the boundaries is contained in the intersection of the hyperplanes

identified by the region vector elements that are different. More in general, still focusing on this example, let \mathbf{k}_i and \mathbf{k}_j differ for $1 \leq t \leq 6$ components of indices ℓ_1, \dots, ℓ_t : if the shared boundary $\partial X_i \cap \partial X_j$ is not empty (the union of the closures $\bar{X}_i \cup \bar{X}_j$ is therefore a connected set), then

$$(\partial X_i \cap \partial X_j) \subseteq (\Pi_{\ell_1} \cap \dots \cap \Pi_{\ell_t});$$

otherwise, the set $\Pi_{\ell_1} \cap \dots \cap \Pi_{\ell_t}$ does not intersect anyone of the boundaries $\partial X_i, \partial X_j$.

The relationship between the elements of the region vectors $\mathbf{k}_i, \mathbf{k}_j$ and connectivity of $\bar{X}_i \cup \bar{X}_j$, observed in the example of Figure 1, is generalizable to partitions in \mathbb{R}^n . However, the example is not intended to be exhaustive of all possible cases which may occur.

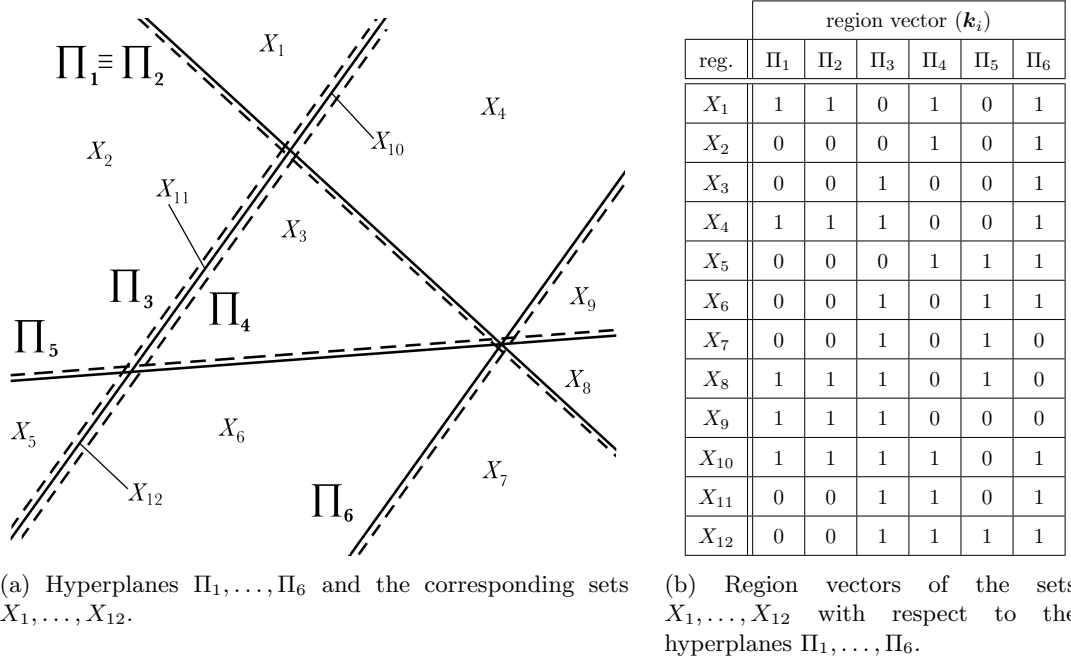


Figure 1: Left: example of partition of \mathbb{R}^2 by six hyperplanes Π_1, \dots, Π_6 . Dotted lines denote the part of the plane with $\mathbf{x}^T \mathbf{w}_j + b_j < 0$, for each $j = 1, \dots, 6$. Right: region vectors $\mathbf{k}_i \in \{0, 1\}^6$ corresponding to each subset $X_i, i = 1, \dots, 12$, of the partition.

2 Discontinuity for Neural Networks

Let N be an H -layers perceptron. We recall that the map \mathcal{L}_{h+1} , characterizing the transformations performed by layer L_{h+1} on the outputs of layer L_h , is defined by (1).

From this formula, it is straightforward to note that the characterizing function $\hat{\mathbf{F}} = \mathcal{L}_1^{H+1} : \mathbb{R}^{N_0} \rightarrow \mathbb{R}^{N_{H+1}}$ of N is a continuous function if $\mathbf{f}_1, \dots, \mathbf{f}_{H+1}$ are all continuous functions. In the earliest works on Neural Networks [19, 24], the first models for artificial neurons did not consider continuous activation functions but the Heaviside function (or suitable variations of it), mainly used to model the “on/off” activation of the neurons.

In the subsequent development of artificial intelligence, the Heaviside function \mathcal{H} has been abandoned, as the subderivative constantly equal to zero prevents the use of the gradient descent during NN training; this phenomenon is equivalent to the asymptotic end of the so-called vanishing gradient problem (see [8, ch. 8.2.5]). In place of \mathcal{H} , many other continuous functions have been

introduced; as a consequence, all recent NNs described in literature are characterized by continuous functions since they are given by the composition of continuous functions.

In this work, we describe a novel approach to reintroduce discontinuities in NNs in such a way that NNs can not only approximate discontinuous functions but also learn discontinuities. In the sequel, we will refer to this property of the NN as the ability of learning discontinuity interfaces, and such discontinuity interfaces will be called “learnable discontinuities”.

2.1 Adding Heaviside to Activation Functions

The main idea behind learnable discontinuities for NNs is to apply the effects of a bias “outside” the activation function f only when the inputs satisfy certain conditions, for example to be not smaller than zero. Thanks to this new bias, the NN has a new trainable parameter that introduces a discontinuity in the function of the NN. The discontinuity introduced depends both on the new parameters and on the weights and biases of the NN; for this reason we will refer to *learnable discontinuities*.

Herein, in order to introduce possible discontinuities in the layers, we add to the activation functions a jump, whose size is expressed by a parameter ε . In details, we add to each element of the right-hand-side of (1) a multiple of the Heaviside function applied component-wise to $W^{(h+1)T}\mathbf{x}^{(h)} + \mathbf{b}^{(h+1)}$, namely we set:

$$\mathbf{x}^{(h+1)} = \mathbf{f}_{h+1} \left(W^{(h+1)T}\mathbf{x}^{(h)} + \mathbf{b}^{(h+1)} \right) + \boldsymbol{\varepsilon}^{(h+1)} \odot \mathcal{H} \left(W^{(h+1)T}\mathbf{x}^{(h)} + \mathbf{b}^{(h+1)} \right), \quad (4)$$

where the symbol \odot denotes the Hadamard (element-wise) product and $\boldsymbol{\varepsilon}^{(h+1)} \in \mathbb{R}^{N_{h+1}}$ is the vector collecting the N_{h+1} jumps introduced (see Figure 2).

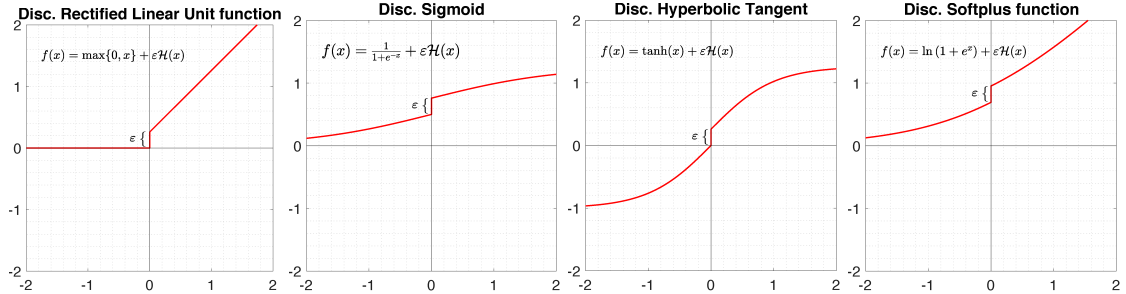


Figure 2: Examples of activation functions plus a multiple of \mathcal{H} .

Definition 2.1 (Discontinuous Layer). *A discontinuous fully-connected layer L with input in \mathbb{R}^c , output in \mathbb{R}^d , and activation function f for the continuous part, is a layer with incoming connections and output signals defined by the characterizing function $\mathcal{L} : \mathbb{R}^c \rightarrow \mathbb{R}^d$ such that:*

$$\mathcal{L}(\mathbf{x}) = \mathbf{f} \left(W^T \mathbf{x} + \mathbf{b} \right) + \boldsymbol{\varepsilon} \odot \mathcal{H} \left(W^T \mathbf{x} + \mathbf{b} \right), \quad (5)$$

where $\boldsymbol{\varepsilon} \in \mathbb{R}^d$ is the vector of trainable discontinuity jumps and where \mathbf{f} , W , and \mathbf{b} are the element-wise application of f , the weight matrix, and the bias vector, respectively.

In the following, a discontinuous layer L of a NN and the corresponding characterizing function \mathcal{L} , defined by (5), will be denoted by δL and $\delta \mathcal{L}$, respectively.

Definition 2.2 (Discontinuous Neural Network). *Let N be a NN with at least one discontinuous layer. Then N is called a discontinuous Neural Network and will be denoted by δNN .*

We remark that the d parameters $\varepsilon_1, \dots, \varepsilon_d$ appearing in (5) are learnable parameters, since the derivatives of $\delta \mathcal{L}(\mathbf{x})$ with respect to them are not constantly equal to zero, as stated in the following Proposition.

Proposition 2.1. Let y_j be the j -th element of $\delta\mathcal{L}(\mathbf{x}) =: \mathbf{y} \in \mathbb{R}^d$ given by (5). Then it holds $\partial y_j / \partial \varepsilon_j \neq 0$ for each $j = 1, \dots, d$.

Proof. The proof is straightforward since, for each $j = 1, \dots, d$, we have that

$$\partial y_j / \partial \varepsilon_j = \mathcal{H}(W_{\cdot,j}^T \mathbf{x} + b_j),$$

that is not constantly equal to zero. \square

From the previous proposition, we can easily deduce also the following one.

Proposition 2.2. Let N be an H -layers perceptron. Assume that for a fixed $h \in \{0, \dots, H\}$, the $(h+1)$ -th layer of N is a discontinuous one (therefore, denoted by δL_{h+1}). For each $\mathbf{x}^{(0)} \in \mathbb{R}^{N_0}$, let y_j be the j -th element of $\mathbf{y} = \mathcal{L}_1^{H+1}(\mathbf{x}^{(0)})$ and let $\mathbf{x}^{(i)}$ denote the vector $\mathcal{L}_1^i(\mathbf{x}^{(0)}) \in \mathbb{R}^{N_i}$, for each $i = 1, \dots, H+1$. Let \mathcal{J} denote the jacobian of \mathcal{L}_{h+2}^{H+1} .

For each input vector $\mathbf{x}^{(0)}$, if the derivatives corresponding to the j -th row $J_{j,\cdot}$ of \mathcal{J} exist at $\mathbf{x}^{(h+1)}$, if the (j,k) -th function $J_{j,k}$ of \mathcal{J} is not constantly equal to zero at $\mathbf{x}^{(h+1)}$, and if $W_{\cdot,k}^{(h+1)T} \mathbf{x}^{(h+1)} + b_k^{(h+1)}$ is not constantly negative, then

$$\frac{\partial y_j}{\partial \varepsilon_k^{(h+1)}} \neq 0,$$

where $j \in \{1, \dots, N_{H+1}\}$, $k \in \{1, \dots, N_{h+1}\}$, and $\varepsilon_k^{(h+1)}$ is the k -th element of $\varepsilon^{(h+1)}$ of $\delta\mathcal{L}_{h+1}$ (see (4)).

Proof. The proof is straightforward since $\mathbf{y} = (\mathcal{L}_{h+2}^{H+1} \circ \mathcal{L}_1^{h+1})(\mathbf{x}^{(0)})$. Indeed, for each $j = 1, \dots, N_{H+1}$ and $k = 1, \dots, N_{h+1}$, we have that

$$\begin{aligned} \frac{\partial y_j}{\partial \varepsilon_k^{(h+1)}} &= \left(\nabla \left(\mathcal{L}_{h+2}^{H+1}(\mathbf{x}^{(h+1)}) \right)_j \right)^T \cdot \nabla_{\varepsilon_k^{(h+1)}} \mathbf{x}^{(h+1)} = \\ &= J_{j,\cdot}^T(\mathbf{x}^{(h+1)}) \cdot \nabla_{\varepsilon_k^{(h+1)}} \mathbf{x}^{(h+1)} = \\ &= J_{j,k}(\mathbf{x}^{(h+1)}) \cdot \frac{\partial x_k^{(h+1)}}{\partial \varepsilon_k^{(h+1)}} = J_{j,k}(\mathbf{x}^{(h+1)}) \cdot \mathcal{H} \left(W_{\cdot,k}^{(h+1)T} \mathbf{x}^{(h+1)} + b_k^{(h+1)} \right), \end{aligned}$$

and from the assumptions the thesis follows. \square

Generalizing Propositions 2.1 and 2.2, we observe that the derivatives with respect to the discontinuity parameters of the loss function can't be constantly equal to zero (excluding special cases) and, therefore, the discontinuity parameters are trainable. To the best of the authors' knowledge, in literature there are no discontinuous NNs characterized by *trainable* discontinuity jumps.

Remark 2.1 (Theoretical justification for discontinuous activation functions). Concerning the universal approximation properties of NNs, we recall the theoretical results of [17, 21], that involve also discontinuous activation functions. These result legitimate the development of discontinuous layers. Indeed, even if there are no clear universal approximation results for H -layers perceptrons with discontinuous activation functions, there are clues suggesting that discontinuities may enhance the approximation (e.g., see [21, Th. 3] or [25]).

Remark 2.2 (Training discontinuous Neural Networks). Introducing discontinuous layers inside a NN implies that the loss function can be piece-wise continuous. Then, during the training, the loss can be non-differentiable and/or discontinuous with respect to the current weights and biases. Actually, non-smoothness is already quite frequent in the working environment of NNs and, due to the discrete representations of functions in a computational domain, the non-differentiability can be "safely disregarded" [8, ch. 6.3]. We can extend the same reasoning when a derivative is needed at a discontinuity point x_0 ; indeed, it is likely that the underlying value is $\tilde{x}_0 = x_0 \pm \epsilon$, with a small

$\epsilon > 0$ (say, less than the machine precision) [8, ch. 6.3]. Furthermore, we observe that derivative-based methods typically perform sufficiently well on piece-wise continuous functions, generally converging to local minima (see the case studies in [29]). Then, we can “safely disregard” the discontinuity of the activation functions during the gradient-based training. Indeed, the training methods for NNs are typically gradient-based and stochastic and, usually, local minima are sought, to avoid overfitting.

However, even if in discontinuous NNs the training ability is preserved (see the results of Section 4), the optimization methods available in the Deep Learning frameworks (e.g., TensorFlow [1]) are indeed not designed to work efficiently with discontinuous functions. The study of more efficient optimization algorithms for the training of discontinuous NNs is certainly of great interest and deserves future investigation.

Remark 2.3 (Discontinuity and increased capacity). We observe that the introduction of discontinuities in NN architectures increases their capacity (see [8, ch. 5.2]), adding a new set of discontinuous functions to the set of functions represented by continuous NNs, the latter being equivalent to formulation (5) with $\epsilon = \mathbf{0}$. Concerning this point, we remark that we refer to the layers characterized by (5) as discontinuous layers independently of the actual values learned for ϵ .

In the next section, we analyze the properties of NNs with at least one discontinuous layer; the aim of this analysis is to understand how a discontinuous layer characterizes the discontinuities of the function $\mathcal{L}_1^{H+1} = \hat{\mathbf{F}}$ of the NN.

3 Properties of Discontinuous Neural Networks

An interesting property of δ NNs is that, in principle, it is possible to exactly find the discontinuities of their characterizing functions; indeed, these NNs are piece-wise continuous functions with known analytical expression (see Corollary 3.3 in the following). This property may be useful to improve the approximations in regression problems but it can be also extremely important for problems in which the discontinuity interfaces of functions are sought, especially in high-dimensional domains (see [6], for an example of application). With this new typology of NN architectures, we propose a novel approach to the discontinuity detection problem, showing that such a kind of NNs can be potentially useful both for discontinuity function approximation and for learning discontinuities.

In this section, we introduce some statements that describe properties related to discontinuity of δ NNs.

3.1 Theoretic Foundations of Discontinuous Neural Networks

The following propositions (Proposition 3.1 and Proposition 3.2) represent some basic results concerning properties of NNs characterized by discontinuous layers. In a nutshell, the propositions state that in a δ NN characterized by a function \mathcal{L}_1^{H+1} , for each discontinuous layer $\delta\mathcal{L}_{h+1}$ of N we have that:

- the discontinuity interfaces of the map $\delta\mathcal{L}_{h+1}$ are affine hyperplanes in \mathbb{R}^{N_h} , characterized by the columns of $W^{(h+1)}$ and the elements of $\mathbf{b}^{(h+1)}$;
- the existence of discontinuity interfaces for $\delta\mathcal{L}_{h+1}$ depends on the nonzero elements of $\epsilon^{(h+1)}$ corresponding to non-null columns of $W^{(h+1)}$.

Moreover, assuming that N has only one discontinuous layer $\delta\mathcal{L}_{h+1}$, a necessary condition for a point $\hat{\mathbf{x}}^{(0)}$ to be a discontinuity point for \mathcal{L}_1^{H+1} is that its image through the first h layers, i.e. $\hat{\mathbf{x}}^{(h)} = \mathcal{L}_1^h(\hat{\mathbf{x}}^{(0)})$, is a discontinuity point for the map $\delta\mathcal{L}_{h+1}$.

In view of the next results, we introduce here the following sets defined for an H -layers perceptron with at least one discontinuous layer δL_{h+1} :

- $\Pi_j^{(h+1)}$ denotes the (possibly degenerate) hyperplane of \mathbb{R}^{N_h} defined by the j -th column of weights and the j -th bias of δL_{h+1} , i.e.:

$$\Pi_j^{(h+1)} := \left\{ \mathbf{x}^{(h)} \in \mathbb{R}^{N_h} \mid W_{:,j}^{(h+1)T} \mathbf{x}^{(h)} + b_j^{(h+1)} = 0 \right\};$$

- $\Pi^{(h+1)}$ denotes the set of all the sets $\Pi_j^{(h+1)}$ defined by the weights and the biases of δL_{h+1} , i.e..

$$\Pi^{(h+1)} := \left\{ \Pi_1^{(h+1)}, \dots, \Pi_{N_{h+1}}^{(h+1)} \right\};$$

- $\Delta^{(h+1)}$ denotes the set of all and only the discontinuity points in \mathbb{R}^{N_h} for $\delta \mathcal{L}_{h+1}$;
- Δ denotes the set of all and only the discontinuity points in \mathbb{R}^{N_0} for \mathcal{L}_1^{H+1} (i.e., for the NN).
- Γ_{h+1} denotes the (possibly empty) counterimage of $\Delta^{(h+1)}$ through \mathcal{L}_1^h , namely:

$$\Gamma_{h+1} := \left\{ \mathbf{x}^{(0)} \in \mathbb{R}^{N_0} \mid \mathcal{L}_1^h(\mathbf{x}^{(0)}) \in \Delta^{(h+1)} \right\} =: (\mathcal{L}_1^h)^{-1}(\Delta^{(h+1)}), \quad (6)$$

where we set \mathcal{L}_1^0 as the identity function, by convention;

- let $\delta L_{h_1+1}, \dots, \delta L_{h_M+1}$ be all and only the discontinuous layers of the H -layers perceptron. Then, we denote by Γ the union of all the counterimages $\Gamma_{h_1+1}, \dots, \Gamma_{h_M+1}$:

$$\Gamma := \bigcup_{m=1}^M \Gamma_{h_m+1}. \quad (7)$$

Proposition 3.1. *Let N be an H -layers perceptron. Assume that for a fixed $h \in \{0, \dots, H\}$, the $(h+1)$ -th layer is a discontinuous one. Let $\mathcal{C}(\Pi^{(h+1)})$ be the partition of \mathbb{R}^{N_h} generated by $\Pi^{(h+1)}$ as in (2) and characterized by $P \in \mathbb{N}$ non-empty subsets such that $\mathcal{C}(\Pi^{(h+1)}) = \{X_1^{(h)}, \dots, X_P^{(h)}\}$. For each $X_p^{(h)} \in \mathcal{C}(\Pi^{(h+1)})$, let $\mathbf{k}_p^{(h)}$ be the region vector of $X_p^{(h)}$ introduced in Definition 1.2. Then, the following assertions are true:*

1. for each $\mathbf{x}^{(h)} \in \mathbb{R}^{N_h}$, equation (4) can be rewritten as

$$\mathbf{x}^{(h+1)} = \mathbf{f}_{h+1} \left(W^{(h+1)T} \mathbf{x}^{(h)} + \mathbf{b}^{(h+1)} \right) + \boldsymbol{\varepsilon}^{(h+1)} \odot \mathbf{k}_i^{(h)}, \quad (8)$$

where $i \in \{1, \dots, P\}$ is such that $\mathbf{x}^{(h)} \in X_i^{(h)}$;

2. $\delta \mathcal{L}_{h+1}$ is discontinuous at $\widehat{\mathbf{x}}^{(h)} \in \mathbb{R}^{N_h}$ if and only if exists $j \in \{1, \dots, N_{h+1}\}$ such that $\widehat{\mathbf{x}}^{(h)} \in \Pi_j^{(h+1)}$, $\varepsilon_j^{(h+1)} \neq 0$, and $W_{:,j}^{(h+1)} \neq \mathbf{0}$. In other words:

$$\Delta^{(h+1)} = \bigcup_{\substack{j=1, \dots, N_{h+1} \\ \varepsilon_j^{(h+1)} \neq 0 \\ W_{:,j}^{(h+1)} \neq \mathbf{0}}} \Pi_j^{(h+1)}. \quad (9)$$

Proof.

1. The proof is immediate as it directly follows from Definition 1.1.

2. The function f_{h+1} is continuous and \mathcal{H} is discontinuous at zero. Then, for each $j = 1, \dots, N_{h+1}$, the function $\mathcal{H}(W_{\cdot, j}^{(h+1)T} \mathbf{x}^{(h)} + b_j^{(h+1)})$ is discontinuous at $\hat{\mathbf{x}}^{(h)}$ if and only if $\hat{\mathbf{x}}^{(h)} \in \Pi_j^{(h+1)}$ and $W_{\cdot, j}^{(h+1)} \neq \mathbf{0}$. Therefore, $\delta\mathcal{L}_{h+1}$ is discontinuous at $\hat{\mathbf{x}}^{(h)}$ if and only if there exists $j \in \{1, \dots, N_{h+1}\}$ such that $\hat{\mathbf{x}}^{(h)} \in \Pi_j^{(h+1)}$, $W_{\cdot, j}^{(h+1)} \neq \mathbf{0}$, and $\varepsilon_j^{(h+1)} \neq 0$.

□

In a nutshell, according to item 2 of the previous Proposition, discontinuity interfaces of $\delta\mathcal{L}_{h+1}$ are affine hyperplanes in \mathbb{R}^{N_h} , whose equations depend on the columns of $W^{(h+1)}$ and the elements of $\mathbf{b}^{(h+1)}$.

While Proposition 3.1 characterizes the discontinuity points of a layer $\delta\mathcal{L}_{h+1}$, in Proposition 3.2 we characterize the discontinuity points of \mathcal{L}_1^{H+1} , assuming that $\delta\mathcal{L}_{h+1}$ is the only discontinuous layer of the δNN .

Proposition 3.2. *Under the assumptions of Proposition 3.1, assume that N has a unique discontinuous layer $\delta\mathcal{L}_{h+1}$ for a fixed $h \in \{0, \dots, H\}$. Let $\hat{\mathbf{x}}^{(0)}$ be a given vector in \mathbb{R}^{N_0} . Then:*

1. \mathcal{L}_1^{h+1} is discontinuous at $\hat{\mathbf{x}}^{(0)}$ if and only if $\hat{\mathbf{x}}^{(0)} \in \Gamma_{h+1}$, i.e. $\mathcal{L}_1^h(\hat{\mathbf{x}}^{(0)}) \in \Delta^{(h+1)}$;
2. if $\hat{\mathbf{x}}^{(0)}$ is a discontinuity point for \mathcal{L}_1^{H+1} then $\hat{\mathbf{x}}^{(0)}$ is a discontinuity point for \mathcal{L}_1^{h+1} ; i.e.:

$$\Delta \subseteq \Gamma_{h+1}.$$

Proof.

1. For each $\hat{\mathbf{x}}^{(0)} \in \Gamma_{h+1}$, the proof that $\hat{\mathbf{x}}^{(0)}$ is a discontinuity point for \mathcal{L}_1^{h+1} is straightforward, as $\delta\mathcal{L}_{h+1}$ is the only discontinuous layer of N . On the other hand, let $\hat{\mathbf{x}}^{(0)} \in \mathbb{R}^{N_0}$ be a discontinuity point for \mathcal{L}_1^{h+1} such that $\mathcal{L}_1^h(\hat{\mathbf{x}}^{(0)}) \notin \Delta^{(h+1)}$; then, $\delta\mathcal{L}_{h+1}$ is continuous at $\mathcal{L}_1^h(\hat{\mathbf{x}}^{(0)})$. But \mathcal{L}_1^h is continuous and $\mathcal{L}_1^{h+1} = \delta\mathcal{L}_{h+1} \circ \mathcal{L}_1^h$; then, \mathcal{L}_1^{h+1} is continuous at $\hat{\mathbf{x}}^{(0)}$, which is a contradiction of the hypothesis.
2. The result is straightforward, as $\delta\mathcal{L}_{h+1}$ is the only discontinuous layer of N .

□

Proposition 3.2 can be generalized to NNs that take into account more discontinuous layers. This generalization is summarized in Theorem 3.1 in the next section.

3.2 Main Results about Discontinuous Neural Networks

The results presented in Section 3.1 describe the discontinuity behavior in a δNN and give a general idea about the potential of such a kind of instruments. Indeed, the discontinuities of a δNN can be quite well characterized.

However, the detection of all the discontinuity interfaces $\Delta \subset \mathbb{R}^{N_0}$ of a map \mathcal{L}_1^{H+1} representing a δNN can be quite an hard task; even just considering a δNN with one discontinuous layer (see Proposition 3.2), the search for points in Γ_{h+1} would require to solve the nonlinear system

$$\begin{cases} W_{\cdot, j_1}^{(h+1)T} \mathcal{L}_1^h(\mathbf{x}) + b_{j_1} = 0 \\ \vdots \\ W_{\cdot, j_k}^{(h+1)T} \mathcal{L}_1^h(\mathbf{x}) + b_{j_k} = 0 \end{cases},$$

where $\varepsilon_j^{(h+1)} \neq 0$ and $W_{\cdot, j}^{(h+1)} \neq \mathbf{0}$ for all and only the $j \in \{j_1, \dots, j_k\} \subseteq \{1, \dots, N_{h+1}\}$.

An alternative idea, useful to avoid the difficulties related to the direct detection of the discontinuity interfaces, is to solve its complementary problem, i.e., find the continuity regions of the domain. Theorem 3.1 and Corollary 3.3 state that δNNs are, actually, piecewise continuous

functions and that given a pair of points in the function domain, it can be easily detected if they belongs or not to a region described by the same continuous piece of function.

Before illustrating these results, we generalize Definition 1.1 with respect to all the discontinuous layers of a δ NN.

Definition 3.1 (Generalized Region Vectors for δ NN). *Let N be an H -layers perceptron with M discontinuous layers. Let δL_{h_m+1} , for $m = 1, \dots, M$, and $0 \leq h_1 < \dots < h_M \leq H$, be the discontinuous layers. For each $m = 1, \dots, M$, let $g_m : \mathbb{R}^{N_{h_m}} \rightarrow \{0, 1\}^{N_{h_m+1}}$ denote the region function corresponding to the weights and biases of δL_{h_m+1} , i.e.*

$$g_m(\mathbf{x}^{(h_m)}) := \mathcal{H} \left(W^{(h_m+1)T} \mathbf{x}^{(h_m)} + \mathbf{b}^{(h_m+1)} \right),$$

for each $\mathbf{x}^{(h_m)} \in \mathbb{R}^{N_{h_m}}$; then, we denote by K_m the image of g_m , i.e. the set

$$K_m = \left\{ \mathbf{k}^{(h_m)} \in \{0, 1\}^{N_{h_m+1}} \mid \exists \mathbf{x}^{(h_m)} \in \mathbb{R}^{N_{h_m}} \text{ s.t. } g_m(\mathbf{x}^{(h_m)}) = \mathbf{k}^{(h_m)} \right\}, \quad (10)$$

representing all the region vectors characterizing the sets of $\mathcal{C}(\Pi^{(h_m+1)})$ in $\mathbb{R}^{N_{h_m}}$ identified by the weights and biases of layer δL_{h_m+1} .

Setting $\delta N = \sum_{m=1}^M N_{h_m+1}$, let $\mathcal{G} : \mathbb{R}^{N_0} \rightarrow \{0, 1\}^{\delta N}$ be defined as

$$\mathcal{G}(\mathbf{x}^{(0)}) = \begin{bmatrix} g_1 \circ \mathcal{L}_1^{h_1}(\mathbf{x}^{(0)}) \\ \vdots \\ g_M \circ \mathcal{L}_1^{h_M}(\mathbf{x}^{(0)}) \end{bmatrix}. \quad (11)$$

Then, \mathcal{G} is called generalized region function of N and we call generalized region vectors of N all the vectors $\mathbf{k} \in K$, where K is the image of \mathcal{G} :

$$K := \left\{ \mathbf{k} \in \{0, 1\}^{\delta N} \mid \exists \mathbf{x}^{(0)} \in \mathbb{R}^{N_0} \text{ s.t. } \mathcal{G}(\mathbf{x}^{(0)}) = \mathbf{k} \right\}. \quad (12)$$

in particular, for each $m = 1, \dots, M$, we denote by $\mathbf{k}^{(h_m)}$ the subvectors of \mathbf{k} belonging to K_m and related to δL_{h_m+1} (see (10)).

Thanks to the generalized region vectors of a given δ NN, it is possible to identify a partition of the domain \mathbb{R}^{N_0} such that the discontinuities of δ NN are contained in the union of all the boundaries of the partition sets, while the map \mathcal{L}_1^{H+1} is continuous in the interior of these sets. These results are better described in the following theorem and corollaries.

Theorem 3.1. *Let N be an H -layers perceptron with M discontinuous layers. Let δL_{h_m+1} , for $m = 1, \dots, M$, and $0 \leq h_1 < \dots < h_M \leq H$, be the discontinuous layers. For each $i = 1, \dots, |K|$, consider the region vector $\mathbf{k}_i \in K$ and let \mathcal{K}_i be defined as*

$$\mathcal{K}_i = \{ \mathbf{x}^{(0)} \in \mathbb{R}^{N_0} \mid \mathcal{G}(\mathbf{x}^{(0)}) = \mathbf{k}_i \}. \quad (13)$$

Then, the following assertions are true:

1. The set Δ of discontinuity points of \mathcal{L}_1^{H+1} is contained in Γ (see (7)).
2. $\{\mathcal{K}_1, \dots, \mathcal{K}_{|K|}\}$ is a partition of \mathbb{R}^{N_0} ;
3. \mathcal{L}_1^{H+1} is continuous in the interior of \mathcal{K}_i (denoted by $\overset{\circ}{\mathcal{K}}_i$), for each $i = 1, \dots, |K|$;
4. Let $\tilde{\Gamma}$ denote the union of all the counterimages, through $\mathcal{L}_1^{h_m}$, of the hyperplanes $\Pi_j^{(h_m+1)}$, i.e.:

$$\tilde{\Gamma} := \bigcup_{m=1}^M \bigcup_{\substack{j=1, \dots, N_{h_m+1} \\ W_{\cdot, j}^{(h_m+1)} \neq \mathbf{0}}} \left(\mathcal{L}_1^{h_m} \right)^{-1} \left(\Pi_j^{(h_m+1)} \right) \quad (14)$$

Then,

$$\bigcup_{i=1}^{|K|} \partial \mathcal{K}_i = \tilde{\Gamma} \quad (15)$$

Proof.

1. Let $\hat{\mathbf{x}}^{(0)}$ be a discontinuity point for \mathcal{L}_1^{H+1} such that $\hat{\mathbf{x}}^{(0)} \notin \Gamma$. Then, for each $m = 1, \dots, M$, we have that $\mathcal{L}_1^{h_m+1}(\hat{\mathbf{x}}^{(0)}) \notin \Delta^{(h_m+1)}$ and $\mathcal{L}_1^{h_m+1}$ is continuous at $\hat{\mathbf{x}}^{(0)}$. Therefore, all the layers are continuous at $\hat{\mathbf{x}}^{(0)}$ and \mathcal{L}_1^{H+1} is continuous at $\hat{\mathbf{x}}^{(0)}$, which is a contradiction of the hypothesis.
2. The result is a direct consequence of the definition of K and $\mathcal{K}_1 \dots, \mathcal{K}_{|K|}$.
3. Let $\hat{\mathbf{x}}^{(0)}$ be an arbitrary point of $\overset{\circ}{\mathcal{K}}_i$, for a fixed $i \in \{1 \dots, |K|\}$; then, there exists an $\hat{\epsilon} > 0$ such that, for each $0 < \epsilon \leq \hat{\epsilon}$, the ball $B_\epsilon(\hat{\mathbf{x}}^{(0)})$ contains only points $\mathbf{x}^{(0)}$ belonging to \mathcal{K}_i and, therefore, for each $m = 1, \dots, M$ it holds

$$\begin{aligned} \lim_{\mathbf{x}^{(0)} \rightarrow \hat{\mathbf{x}}^{(0)}} \mathcal{L}_1^{h_m+1}(\mathbf{x}^{(0)}) &= \lim_{\substack{\mathbf{x}^{(0)} \rightarrow \hat{\mathbf{x}}^{(0)} \\ B_\epsilon(\hat{\mathbf{x}}^{(0)})}} \mathbf{f}_{h_m+1} \left(W^{(h_m+1)T} \mathcal{L}_1^{h_m}(\mathbf{x}^{(0)}) + \mathbf{b}^{(h_m+1)} \right) + \\ &\quad \lim_{\substack{\mathbf{x}^{(0)} \rightarrow \hat{\mathbf{x}}^{(0)} \\ B_\epsilon(\hat{\mathbf{x}}^{(0)})}} \boldsymbol{\varepsilon}^{(h_m+1)} \odot \mathcal{H} \left(W^{(h_m+1)T} \mathcal{L}_1^{h_m}(\mathbf{x}^{(0)}) + \mathbf{b}^{(h_m+1)} \right) = \\ &= \mathbf{f}_{h_m+1} \left(W^{(h_m+1)T} \mathcal{L}_1^{h_m}(\hat{\mathbf{x}}^{(0)}) + \mathbf{b}^{(h_m+1)} \right) + \\ &\quad \boldsymbol{\varepsilon}^{(h_m+1)} \odot \mathbf{k}_i^{(h_m)} = \\ &= \mathcal{L}_1^{h_m+1}(\hat{\mathbf{x}}^{(0)}), \end{aligned}$$

where $[\mathbf{k}_i^{(h_1)T}, \dots, \mathbf{k}_i^{(h_M)T}]^T = \mathbf{k}_i$ is the generalized region vector of \mathcal{K}_i . Then \mathcal{L}_1^{H+1} (and $\mathcal{L}_1^{h_m+1}$, for each $m = 1, \dots, M$) is continuous at $\hat{\mathbf{x}}^{(0)}$; since $\hat{\mathbf{x}}^{(0)}$ is an arbitrary point, \mathcal{L}_1^{H+1} is continuous in $\overset{\circ}{\mathcal{K}}_i$, for each $i = 1, \dots, |K|$.

4. The inclusion $\tilde{\Gamma} \subseteq \bigcup_{i=1}^{|K|} \partial \mathcal{K}_i$ is straightforward, by the definition of the sets $\mathcal{K}_1, \dots, \mathcal{K}_{|K|}$, so we are only left to prove the opposite inclusion.

Let $\hat{\mathbf{x}}^{(0)}$ be a boundary point for a given \mathcal{K}_i , $i \in \{1 \dots, |K|\}$, i.e. $\hat{\mathbf{x}}^{(0)} \in \partial \mathcal{K}_i$; then, for each $\epsilon > 0$, the ball $B_\epsilon(\hat{\mathbf{x}}^{(0)})$ contains both an internal point $\mathbf{x}_{\text{in}}^{(0)} \in \mathcal{K}_i$ and an external point $\mathbf{x}_{\text{out}}^{(0)} \in \mathcal{K}_j$, $j \neq i$. Since these two vectors belong to two different continuity regions, they have different generalized region vectors, i.e. $\mathbf{k}_i := \mathcal{G}(\mathbf{x}_{\text{in}}^{(0)}) \neq \mathcal{G}(\mathbf{x}_{\text{out}}^{(0)}) =: \mathbf{k}_j$.

Let \mathcal{K}_j , $j \neq i$, be a region that shares the boundary with \mathcal{K}_i through $\hat{\mathbf{x}}^{(0)} \in \partial \mathcal{K}_i$, and let $m \in \{1, \dots, M\}$ be the first index such that $\mathbf{k}_i^{(h_m)} \neq \mathbf{k}_j^{(h_m)}$, for each pair of internal and external points $\mathbf{x}_{\text{in}}^{(0)} \in \mathcal{K}_i$ and $\mathbf{x}_{\text{out}}^{(0)} \in \mathcal{K}_j$, respectively, in the ball $B_\epsilon(\hat{\mathbf{x}}^{(0)})$, for each $\epsilon > 0$.

For item 3, the sub-NN characterized by $\mathcal{L}_1^{h_m}$ is continuous on $B_\epsilon(\hat{\mathbf{x}}^{(0)})$; then, the image $B_\epsilon^{(h_m)}(\hat{\mathbf{x}}^{(h_m)}) := \mathcal{L}_1^{h_m}(B_\epsilon(\hat{\mathbf{x}}^{(0)}))$ is a connected neighborhood of $\hat{\mathbf{x}}^{(h_m)}$ in $\mathbb{R}^{N_{h_m}}$, where we denoted by $\hat{\mathbf{x}}^{(h_m)}$ the image of $\hat{\mathbf{x}}^{(0)}$ through $\mathcal{L}_1^{h_m}$.

Let us denote by $\mathbf{x}_{\text{in}}^{(h_m)}, \mathbf{x}_{\text{out}}^{(h_m)} \in B_\epsilon^{(h_m)}(\hat{\mathbf{x}}^{(h_m)})$ the images of $\mathbf{x}_{\text{in}}^{(0)}, \mathbf{x}_{\text{out}}^{(0)} \in B_\epsilon(\hat{\mathbf{x}}^{(0)})$, respectively, through $\mathcal{L}_1^{h_m}$. Since we have $\mathbf{k}_i^{(h_m)} \neq \mathbf{k}_j^{(h_m)}$, it holds that $\mathbf{x}_{\text{in}}^{(h_m)}, \mathbf{x}_{\text{out}}^{(h_m)}$ belong to two distinct sets $X_i^{(h_m)}, X_j^{(h_m)} \in \mathcal{C}(\Pi^{(h_m+1)})$, respectively, where $\partial X_i^{(h_m)} \cap \partial X_j^{(h_m)} \neq \emptyset$.

Then, for each $\epsilon > 0$, we have that $\hat{\mathbf{x}}^{(h_m)} \in \partial X_i^{(h_m)} \cap \partial X_j^{(h_m)}$ and, therefore, $\hat{\mathbf{x}}^{(0)} \in \tilde{\Gamma}$ because $\hat{\mathbf{x}}^{(h_m)}$ belongs to one of the hyperplanes in $\Pi^{(h_m+1)}$. For the generality of the choice of \mathcal{K}_i and $\hat{\mathbf{x}}^{(0)} \in \partial \mathcal{K}_i$, we have that $\bigcup_{i=1}^{|K|} \partial \mathcal{K}_i \subseteq \tilde{\Gamma}$.

□

Corollary 3.2. *Let the hypotheses of Theorem 3.1 be satisfied. Then, the set Δ of all the discontinuity points of \mathcal{L}_1^{H+1} is contained in $\bigcup_{i=1}^{|K|} \partial\mathcal{K}_i$ and, in particular, we have that*

$$\Delta \subseteq \Gamma \subseteq \bigcup_{i=1}^{|K|} \partial\mathcal{K}_i. \quad (16)$$

Proof. We observe that item 3 of Theorem 3.1 implies $\Delta \subseteq \bigcup_{i=1}^{|K|} \partial\mathcal{K}_i$ and we recall that item 1 of Theorem 3.1 implies $\Delta \subseteq \Gamma$. By construction, $\Gamma \subseteq \tilde{\Gamma}$ and for item 4 we have $\tilde{\Gamma} = \bigcup_{i=1}^{|K|} \partial\mathcal{K}_i$; then, (16) is proved. □

The Theorem and the Corollary above highlight the importance of the partition given by sets (13). Due to the properties illustrated in the theorem, we denote the sets of the partition by the names of *continuity regions*, through the following definition.

Definition 3.2 (Continuity Regions of a δ NN). *Let N be an H -layers perceptron with M discontinuous layers, and let $\mathcal{K}_1 \dots, \mathcal{K}_{|K|} \subseteq \mathbb{R}^{N_0}$ be defined as in (13). The sets $\mathcal{K}_1 \dots, \mathcal{K}_{|K|}$ are called continuity regions of N .*

Definition 3.2 makes even more sense looking at the results illustrated in the following Corollary, immediate consequence of Theorem 3.1, which states that δ NNs are piece-wise continuous functions.

Corollary 3.3. *Let the hypotheses of Theorem 3.1 be satisfied. Let $\mathcal{M}_{h_m+1}^{(i)}$ be the function*

$$\mathcal{M}_{h_m+1}^{(i)}(\mathbf{x}^{(h_m)}) = \mathbf{f}_{h+1} \left(W^{(h_m+1)T} \mathbf{x}^{(h_m)} + \mathbf{b}^{(h+1)} \right) + \boldsymbol{\varepsilon}^{(h_m+1)} \odot \mathbf{k}_i^{(h_m)},$$

for $m = 1, \dots, M$, where $\mathbf{k}_i^{(h_m)}$ is the subvector of \mathbf{k}_i (see Definition 3.1).

Then \mathcal{L}_1^{H+1} is a piecewise continuous function such that

$$\mathcal{L}_1^{H+1}(\mathbf{x}^{(0)}) = \begin{cases} \mathcal{F}_1(\mathbf{x}^{(0)}), & \text{if } \mathbf{x}^{(0)} \in \mathcal{K}_1 \\ \vdots \\ \mathcal{F}_{|K|}(\mathbf{x}^{(0)}), & \text{if } \mathbf{x}^{(0)} \in \mathcal{K}_{|K|} \end{cases}$$

where, for each $i = 1, \dots, |K|$, \mathcal{F}_i is the continuous function defined by

$$\mathcal{F}_i = \mathcal{L}_{h_M+2}^{H+1} \circ \mathcal{M}_{h_M+1}^{(i)} \circ \mathcal{L}_{h_{M-1}+2}^{h_M} \circ \mathcal{M}_{h_{M-1}+1}^{(i)} \circ \dots \circ \mathcal{M}_{h_1+1}^{(i)} \circ \mathcal{L}_1^{h_1}.$$

We conclude this section making a resume and some observations concerning the outcomes of Theorem 3.1 and Corollary 3.3 that are useful for practical applications. First, the results state that it is possible to identify the points where the map \mathcal{L}_1^{H+1} of the δ NN is certainly continuous; these regions of the domain are the open sets $(\mathcal{K}_i \setminus \partial\mathcal{K}_i)$, i.e. the interior of the continuity regions \mathcal{K}_i . As a consequence, we have the precise indication that the discontinuity interfaces of \mathcal{L}_1^{H+1} are contained in the union of the boundaries $\partial\mathcal{K}_i$. Furthermore, since the continuity regions are characterized by the generalized region vectors and the set K of these vectors is contained in $K_1 \times \dots \times K_M$ (see (12)), we can control the number of continuity regions through the following proposition.

Proposition 3.3 (Maximum number of continuity regions). *Let the hypotheses of Theorem 3.1 be satisfied. Then, the number $|K|$ of continuity regions for N is such that*

$$|K| \leq 2^{\sum_{m=1}^M N_{h_m+1}}.$$

The continuity regions \mathcal{K}_i of a δ NN are a bit more complex than their simpler counterparts in the discontinuous layers, which are sets $X_i^{(h)}$ identified by hyperplanes and region vectors (see Section 1.1.2). Indeed, sets \mathcal{K}_i are not necessarily convex nor connected sets in \mathbb{R}^{N_0} , due to the function compositions that take place inside the δ NN; this is also the reason for the difficulties in directly computing the discontinuity interfaces, even if theoretically possible. On the other hand, given a set $X^{(0)} \subseteq \mathbb{R}^{N_0}$, the computation of the generalized region vector $\mathcal{G}(\mathbf{x}^{(0)})$ for each $\mathbf{x}^{(0)} \in X^{(0)}$ is extremely easy and fast since, by definition, \mathcal{G} is characterized by sub-networks of the δ NN considered (see (11)). Therefore, in practice, the regions \mathcal{K}_i can be deduced by sampling a large enough set of points in the domain and then computing the generalized region vector for all such points.

Since the continuity regions can be non-convex and/or disconnected, the actual significance of Proposition 3.3 is to define an upper bound to the number of continuous functions that define the equation of \mathcal{L}_1^{H+1} through (3.3).

4 Learning Discontinuities: Numerical Experiments

In the previous sections we have focused on the analysis of the function \mathcal{L}_1^{H+1} characterizing a (trained) δ NN, highlighting properties related to the discontinuous layers present in the network. In this section we show experimental evidence about the ability of a δ NN to learn discontinuity interfaces, given a target function to approximate. To this aim, we consider different test cases, with increasing complexity, given by discontinuous functions with several kinds of discontinuity interfaces. Several δ NN architectures have been tested; since the main target of this section is to show the viability of δ NNs for discontinuous function regression and discontinuity interfaces detection, we do not focus on the problem of finding the best hyper-parameters or architectures for regression performance, but we focus on the parameters and characteristics related to the discontinuous layers of the δ NNs, aiming to highlight the actual sensitivity of the δ NN approximation performances with respect to the new discontinuous layers.

4.1 Test cases

In all the test cases here considered, the underlying function is a scalar function with domain in \mathbb{R}^2 . In all the cases, we consider the functions restricted to the region $D = [-2, 2] \times [-2, 2]$. The test functions used are the following.

1. **Test 1.** We consider the function $g_\ell : D \rightarrow \mathbb{R}$ defined as

$$g_\ell(\mathbf{x}) = \begin{cases} 2 \sin(1.25\pi \|\mathbf{x}\|) + 4 & \text{if } x_2 \leq 2x_1 \\ 2 \sin(0.75\pi \|\mathbf{x}\|) & \text{otherwise} \end{cases}, \quad (17)$$

The discontinuity interface is the line $\ell : x_2 = 2x_1$ that halves the square domain D (see Figure 3-(a)).

2. **Test 2.** As a second test function, we consider $g_s : D \rightarrow \mathbb{R}$ defined as

$$g_s(\mathbf{x}) = \begin{cases} -2x_1^2 + 6, & \text{if } \mathbf{x} \in [-1, 1] \times [0, +\infty) \\ 4e^{(1-x_1^2)/2}, & \text{otherwise} \end{cases}, \quad (18)$$

The discontinuity interface is the segment $s := \{\mathbf{x}(\lambda) = \lambda[-1, 0]^T + (1-\lambda)[1, 0]^T, 0 \leq \lambda \leq 1\}$ (see Figure 3-(b)). Then, in this case, the discontinuity is a sort of straight “rip” for the graph of the function.

3. **Test 3.** The third test function we consider is $g_\eta : D \rightarrow \mathbb{R}$, defined as

$$g_\eta(\mathbf{x}) = \begin{cases} \sin(0.4\pi(x_1 + x_2)) , & \text{if } x_2 \geq e^{x_1} \\ \sin(0.7\pi(x_1 + x_2)) - 4 , & \text{if } x_2 < e^{x_1} - 1 \\ \sin(\pi(x_1 + x_2)) + 4 , & \text{otherwise} \end{cases} , \quad (19)$$

The discontinuity interfaces are two curves, η_1 and η_2 , that split the domain in three regions (see Figure 3-(c)). In particular, $\eta_1 : x_2 = e^{x_1}$ and $\eta_2 : x_2 = e^{x_1} - 1$.

4. **Test 4.** As a last example we consider a function $g_\gamma : D \rightarrow \mathbb{R}$ characterized by a discontinuity interface that is a closed curve, the circumference $\gamma : \|\mathbf{x}\|^2 = 1$. The function g_γ (see Figure 3-(d)) is defined as

$$g_\gamma(\mathbf{x}) = \begin{cases} \sin(\pi(x_1 + x_2)) + 4 , & \text{if } \|\mathbf{x}\|^2 \leq 1 \\ \sin(0.4\pi(x_1 + x_2)) , & \text{otherwise} \end{cases} . \quad (20)$$

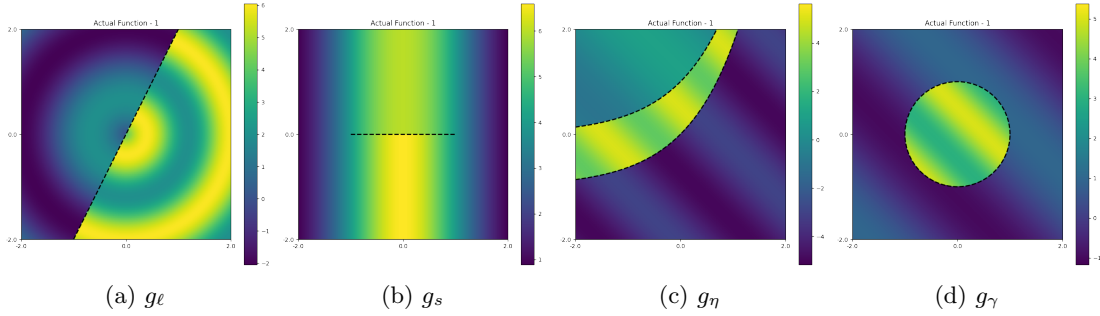


Figure 3: Top view of the functions g_ℓ , g_s , g_η , and g_γ (see equations (17)-(20), respectively). The dotted curves correspond to the discontinuity interfaces.

4.2 Architectures and Performance Measures

In the numerical experiments, we use discontinuous H -layers perceptrons. We consider three archetypes of δ NN architectures, varying the depth, the number of discontinuous layers, and the size and position of these layers. The number of units in the fully connected (i.e., non-discontinuous) hidden layers is fixed to 128, while we let d denote the number of units in each discontinuous layer.

The three architecture archetypes we consider are:

1. *Architectures with one discontinuous layer only.* We consider a unique discontinuous layer with d units, which is the h -th inner layer, with $h \in \{1, \dots, H\}$ (see Figure 4-(a)). We let $\mathcal{A}_{h,H}^d$ denote such an architecture.
2. *Architectures with two discontinuous layers.*
 - 2.1. We consider two consecutive discontinuous layers with d units each, which are the h -th and $(h+1)$ -th inner layers, for $h \in \{1, \dots, H-1\}$ (see Figure 4-(b)). We let $\mathcal{B}_{h,H}^d$ denote such an architecture.
 - 2.2. We consider two discontinuous layers with d units each, separated by a fully-connected layer; they are the h -th and $(h+2)$ -th inner layers, with $h \in \{1, \dots, H-2\}$ (see Figure 4-(c)). We let $\mathcal{C}_{h,H}^d$ denote such an architecture.

In our tests we considered: architecture $\mathcal{A}_{h,H}^d$ with $H = 5$, $h = 1, \dots, 5$ and $d = 2, 4, 8$; architecture $\mathcal{B}_{h,H}^d$ with $H = 5, 7$, $h = 1, \dots, H-1$ and $d = 2, 4, 8$; architecture $\mathcal{C}_{h,H}^d$ with $H = 5, 7$, $h = 1, \dots, H-2$ and $d = 2, 4, 8$. As a whole, we tested 69 architectures.

The activation function for all the hidden layers of the δ NNs is the `elu` activation function [7], chosen after a preliminary investigation, while in the output layer the linear activation function is used. The depth H for the NNs has been chosen to analyze the effects of the discontinuous layers with respect to their position in the network, while guaranteeing a good approximation of the target function, and trying to avoid the so-called degradation problem [10]. The size of the discontinuous layers d has been chosen to keep limited the maximum number of continuity regions (see Proposition 3.3) to ease the analysis.

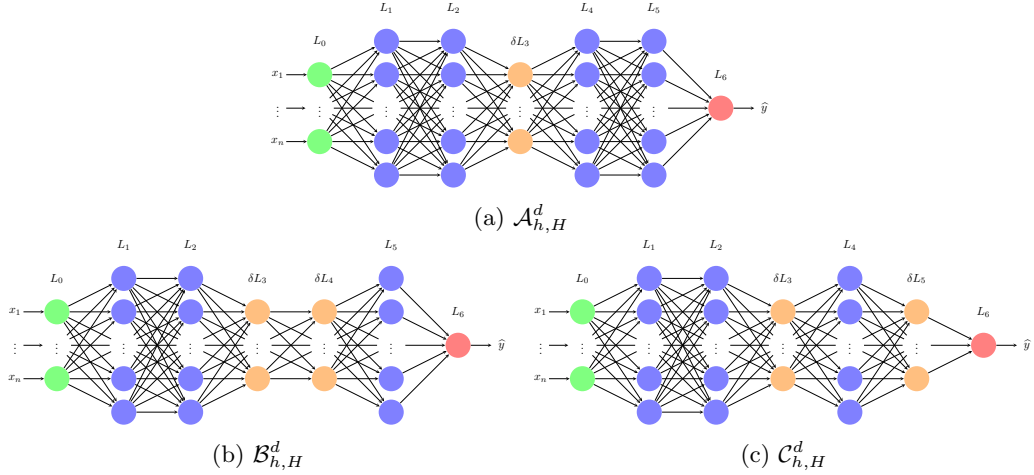


Figure 4: The three archetypes $\mathcal{A}_{h,H}^d$, $\mathcal{B}_{h,H}^d$, $\mathcal{C}_{h,H}^d$ (example with $h = 3$ and $H = 5$). Discontinuous layers are represented by orange units, fully-connected hidden layers by purple units, input layers by green units, and output layers by red units.

All the networks are trained with the same training options and configurations:

- **Dataset:** for each testcase $g = g_\ell, g_s, g_\eta, g_\gamma$, the dataset \mathcal{D} is made of 10 000 pairs (\mathbf{x}_i, y_i) , with \mathbf{x}_i randomly sampled with uniform distribution from the domain of the target function and $y_i = g(\mathbf{x}_i)$. Then, \mathcal{D} is randomly split into the training set (5 600 pairs), validation set (1 400 pairs) and test set (3 000 pairs);
- **Data preprocessing:** z-normalization (see [20]) of the input data with respect to the training set;
- **Training options:** Mean Square Error (MSE) loss function, Adam optimizer (learning rate $\epsilon = 10^{-3}$, decay rates $\beta_1 = 0.9$ and $\beta_2 = 0.999$; see [15]) with learning rate reduction on plateaus (factor 0.75, patience 50 epochs), mini-batch size of 64 samples, 5 000 maximum number of epochs, early stopping with best-weights restoration (patience of 250 epochs) as regularizer.

The performance measure that we adopt to evaluate the results on the test set is the Mean Absolute Error (MAE), considering that the values of all the test functions are approximately between -2 and 6 :

$$\text{MAE}(\delta\text{NN}, \mathcal{P}) := \frac{1}{|\mathcal{P}|} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{P}} |y_i - \hat{y}_i|,$$

where \mathcal{P} is the test set and \hat{y}_i denotes the output prediction of the δ NN for \mathbf{x}_i .

4.3 Performance Analysis for the Test Functions

We start our analysis verifying the performance of the δ NN in the regression task. Upon training all the NNs, with respect to all the test cases, we analyze the MAE of the models on the test set. The errors, reported in Table 1, prove that δ NNs behave quite well in the regression task, and discontinuous layers do not hinder the regression abilities of Neural Networks.

MAE	g_ℓ	g_s	g_η	g_γ
mean	0.0906	0.0218	0.1868	0.0605
std	0.1597	0.1363	0.2706	0.0557
median	0.0274	0.0029	0.0695	0.0427

Table 1: Statistics of the MAEs over all the δ NNs on the test sets of the functions $g_\ell, g_s, g_\eta, g_\gamma$.

Once the good approximation abilities of the δ NNs are verified, we focus the analysis on the ability to identify the continuity regions. We recall that the computation of the generalized region vectors for an arbitrary set of points in the domain is extremely easy and fast since, by definition, \mathcal{G} is characterized by sub-networks of the δ NN (see (11)). In a nutshell, when we compute the prediction $\mathcal{L}_1^{H+1}(\mathbf{x}^{(0)})$ for a generic vector $\mathbf{x}^{(0)} \in \mathbb{R}^{N_0}$, we can easily obtain from the NN, at no additional cost, also the intermediate values $\mathcal{L}_1^h(\mathbf{x}^{(0)})$ needed in (11), returned by any arbitrary hidden layer L_h ; the computation of \mathcal{G} is therefore an extremely easy and fast task, for an arbitrary batch of vectors.

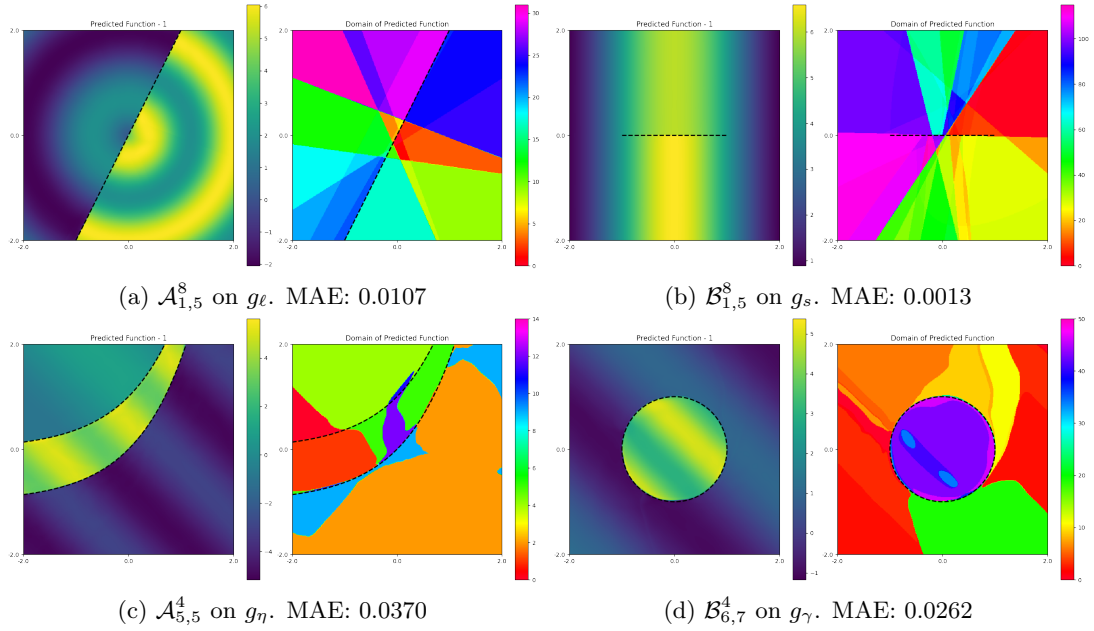


Figure 5: Approximation of $g_\ell, g_s, g_\eta, g_\gamma$ with the δ NNs $\mathcal{A}_{1,5}^8, \mathcal{B}_{1,5}^8, \mathcal{A}_{5,5}^4, \mathcal{B}_{6,7}^4$, respectively. For each subfigure we have the δ NN function values (left) and the continuity regions (right). In the plots of the continuity regions, each region is identified by a different color, according to the colorbar on the right.

In Figure 5 we report some examples of results of the regression and interface detection tasks obtained on all the tests considered. For each sub-figure the left plot reports the approximation of the corresponding function, and the right panel reports the continuity regions identified by the δ NN. These last figures are obtained by picking a large enough number of points in the region D ,

and computing the corresponding region vector. Points with the same region vector are labeled with the same color. Note that according to Proposition 3.3 the number of continuity regions can be up to $2^{\delta N}$, being δN the total number of units in the discontinuous layers.

In general, we observe a very good approximation of both the test function and the actual continuity regions by the δ NNs. In Figure 5, we show the results obtained for each test function on a selected δ NN. Looking at the continuity regions of the δ NNs, we observe the following phenomena.

Boundaries of the continuity regions. If the first hidden layer is discontinuous, then there are necessarily some continuity regions with straight boundaries. This is due to the fact that the boundaries of the continuity regions correspond to the counterimages through \mathcal{L}_1^h of the hyperplanes $\Pi_j^{(h+1)}$ introduced by the discontinuous layers (see Theorem 3.1, item 4). When the first layer is a discontinuous one, since by convention \mathcal{L}_1^0 is the identity function (see (6)), the continuity region boundaries identified by the hyperplanes $\Pi_j^{(1)}$ are the hyperplanes themselves. In particular, if δL_1 is the only discontinuous layer, then the continuity regions have *only* linear boundaries. This phenomenon is depicted in Figure 5: figure top, left is obtained with a unique discontinuous layer which is the first hidden layer, whereas figure top, right is obtained with two discontinuous layers, which are the first and second hidden layers. Note that in the second case the boundaries separating the continuity regions identified are either lines or mildly curvilinear lines.

On the other hand, the more the discontinuous layers are located toward the output layer, the more the counter-images of the corresponding hyperplanes are obtained from the application of nonlinear functions, and the more the boundaries of the continuity regions can have a curved shape. This behaviour is still depicted in Figure 5: the bottom figures are obtained with discontinuous layers which are either the last one (bottom left case) or the last two (bottom right case): in both cases, highly curved boundaries are obtained.

Trade-off between approximation and discontinuity detection. The more the discontinuous layers are located toward the output layer, the more they are forced to focus on learning the discontinuity jumps of the target function. Indeed, a discontinuous layer near to the input can mainly focus on learning the discontinuity interface while spending not much effort on the precise values of the jump parameters ε , since the following layers can enlarge/shorten these jumps to reach appropriate values for the approximation task. Then, a discontinuous layer followed by few layers is forced also to learn the jumps, having less help available from the following layers.

In few cases, we observe that the efforts of the δ NN is mostly spent in learning the discontinuity interfaces, harming the function approximation. However, this mainly happens with discontinuous layers with few units (i.e., $d = 2$) that are not near to the input layer. This rare problems clearly depend on the dimensionality reduction inside the NN (from \mathbb{R}^{128} to \mathbb{R}^2) that occurs toward the end of the network.

Since the δ NNs can learn much more discontinuity interfaces than the ones actually characterizing the test functions, we observe that the models use the discontinuous layers also trying to adjust and improve the regression. In particular, we observe that in many cases the boundaries of the continuity regions partially “follow” the level curves of the target function but discontinuities are almost imperceptible.

From all the previous observations, the following indications can be deduced. If we are looking for discontinuities with an almost linear discontinuity interface, it is preferred to introduce discontinuous layers near to the input layer, whereas if we are looking for discontinuities with highly curvilinear interfaces, it is preferred to introduce discontinuous layers near to the output layer. Moreover, too small discontinuous layers should be avoided, and also the use of a discontinuous layer as first hidden layer should be avoided, unless it is *a priori* known that the discontinuity interface is a straight line or segment.

We remark that in case of an incorrect choice of the discontinuous layers, we observe that the approximation abilities very rarely are compromised even if the discontinuity interfaces are not learned. Typical examples are the ones illustrated in Figure 6, where two δ NN approximating the function g_γ (see Figure 3-(d)) have continuity regions with boundaries that are “not-enough” curvilinear to approximate the circumference γ . Following these indications, the architectures used in Figure 5 are the best δ NNs among the ones with $h < \lceil H/2 \rceil$ for g_ℓ and g_s and among the ones with $h \geq \lceil H/2 \rceil$ for g_η and g_γ , where we recall that h is the index of the first discontinuous layer in the NN architecture and H is the depth of the NN.

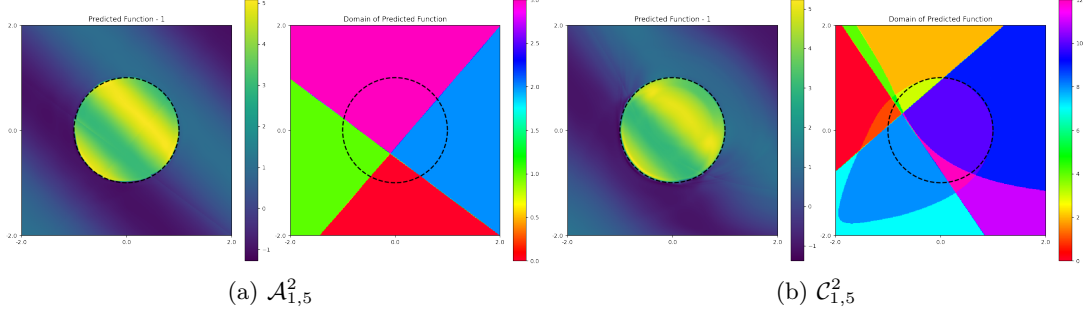


Figure 6: Approximation of g_γ with the δ NNs $\mathcal{A}_{1,5}^2$ and $\mathcal{C}_{1,5}^2$. For each subfigure we have the δ NN function values (left) and the continuity regions (right). In the plots of the continuity regions, each region is identified by a different color, according to the colorbar on the right.

4.4 Continuity Region Clustering

In the previous section, we observed that most of the trained δ NNs approximate with a good quality level both the functions and their continuity regions of the test functions. Nonetheless, in most of the cases, the continuity regions of the test functions are approximated by the δ NNs with many continuity regions that are separated by negligible/inexistent discontinuities; in practice the method in its basic form largely overestimates the number of continuity region. For example, in Figure 5-(b), we see more than 100 continuity regions but, clearly, the regions in the bottom half of the square can be considered as one single continuity region, as no discontinuities are perceived in the left plot; the same applies to the top half of the square. Then, to identify the actual continuity regions of the test functions, we introduce a method to group the continuity regions of the δ NNs if the discontinuities on the boundaries are negligible/inexistent.

The development of such a kind of method is by far not trivial. From Proposition 3.1 we know that if a discontinuity jump parameter is zero, the δ NN does not introduce a discontinuity on the corresponding continuity region boundary. Nonetheless, small values of the discontinuity jump parameters not necessarily correspond to negligible discontinuity jumps; indeed, as previously observed, a jump introduced by a discontinuous layer can be enlarged/shortened by the following layers. Then, we developed a clusterization method based both on the values of the discontinuity jump parameters and on the action they play inside the δ NN.

Let N be a discontinuous H -layers perceptron defined as in Theorem 3.1, i.e. with M discontinuous layers $\delta L_{h_1+1}, \dots, \delta L_{h_M+1}$ such that $0 \leq h_1 < \dots < h_M \leq H$. We recall that N_{h_m+1} is the number of units of the discontinuous layer δL_{h_m+1} and that δN denote the total number of units in discontinuous layers, namely $\delta N := \sum_{m=1}^M N_{h_m+1}$; we also recall that the generalized vector function is $\mathcal{G} : \mathbb{R}^{N_0} \rightarrow \{0, 1\}^{\delta N}$. In view of the formal description of the clusterization method, we introduce here the following further notation:

- for each $\{i_1, \dots, i_k\} \subseteq \{1, \dots, \delta N\}$, we denote by $\mathcal{G}|_{i_1, \dots, i_k}$ the vector valued function whose elements are elements i_1, \dots, i_k of \mathcal{G} (see Definition 3.1);
- after introducing a global indexing for the discontinuity jump parameters, now labeled ε_i

for $i = 1, \dots, \delta N$, we denote by $\mathcal{L}_1^{H+1}|_{\varepsilon_i=0}$ the characterizing function of the δ NN obtained from N by setting to zero the discontinuity jump parameter ε_i ;

- for any finite set of vectors $X = \{\mathbf{x}_1, \dots, \mathbf{x}_q\} \subset \mathbb{R}^{N_0}$, we denote by $\mathcal{L}_1^{H+1}(X)$ the matrix in $\mathbb{R}^{N_{H+1} \times q}$ defined as

$$\mathcal{L}_1^{H+1}(X) := \left[\mathcal{L}_1^{H+1}(\mathbf{x}_1) \dots \mathcal{L}_1^{H+1}(\mathbf{x}_q) \right];$$

we adopt the same convention for $\mathcal{L}_1^{H+1}|_{\varepsilon_i=0}(X)$.

We can now describe the clusterization method for the continuity regions of N . The proposed method defines new regions by suitably merging some continuity regions of the δ NN, leveraging the function $\mathcal{G}|_{i_1, \dots, i_k}$. To understand the merging procedure, consider for example the function $\mathcal{G}|_i$, restricted to the index $i \in \{1, \dots, \delta N\}$ only: this restriction is blind to all the interfaces not related to ε_i and therefore only two regions in \mathbb{R}^{N_0} are retained, which are, respectively: i) the union of all the continuity regions \widehat{C} such that $(\mathcal{G}(\mathbf{x}^{(0)}))_i = \mathcal{G}|_i(\mathbf{x}^{(0)}) = 0$, for all $\mathbf{x}^{(0)} \in \widehat{C}$; ii) the union of all the continuity regions \widetilde{C} such that $(\mathcal{G}(\mathbf{x}^{(0)}))_i = \mathcal{G}|_i(\mathbf{x}^{(0)}) = 1$ for all $\mathbf{x}^{(0)} \in \widetilde{C}$.

The number $k \in \mathbb{N}$ of discontinuity interfaces can be arbitrarily fixed. The clustering is based both on the values of the discontinuity jump parameters ε_i and on the action they play inside the δ NN. Indeed, we introduce for each ε_i a rank value which depends not only on the size of ε_i itself, but also on its effect on the δ NN; the latter dependence is obtained measuring the difference between \mathcal{L}_1^{H+1} and $\mathcal{L}_1^{H+1}|_{\varepsilon_i=0}$, i.e. switching-off ε_i . Formally, let $X = \{\mathbf{x}_1, \dots, \mathbf{x}_q\} \subset \mathbb{R}^{N_0}$ be a finite set of vectors. Let $\|\cdot\|$ be a norm on $\mathbb{R}^{N_{H+1} \times q}$ and let $k \in \mathbb{N}$, $k \leq \delta N$, be fixed. For each $i = 1, \dots, \delta N$, we compute the rank value

$$\rho_i(X) := |\varepsilon_i| \cdot \|\mathcal{L}_1^{H+1}(X) - \mathcal{L}_1^{H+1}|_{\varepsilon_i=0}(X)\|. \quad (21)$$

In a nutshell, $\rho_i(X)$ measures how much ε_i contributes to the outputs of the δ NN, weighted by the value of ε_i itself. The higher $\rho_i(X)$ is, the more likely the discontinuity interface corresponding to ε_i approximates a real discontinuity interface of the target function. Then, retaining only the discontinuity interfaces corresponding to the k highest rank values only, we can merge the continuity regions that are unlikely to be separated by a discontinuity interface of the target function. Following this idea, the rank values are sorted in descending order, and continuity regions separated by hyperplanes corresponding to parameters ε_i with the smallest rank values are merged, in such a way that we end with a fixed number of discontinuity interfaces.

The above procedure can be sketched in the following algorithm.

Algorithm 4.1 (Clusterization Method for Continuity Regions of a δ NN). *Let N be a discontinuous H -layers perceptron defined as in Theorem 3.1 and let $k \in \mathbb{N}$, $k \leq \delta N$, be the number of indices with respect to which I want to perform the continuity region clustering of N . Then:*

1. For each $i = 1, \dots, \delta N$, compute the rank value $\rho_i(X)$ as in (21);
2. sort the rank values in descending order: $\rho_{i_1}(X) \geq \dots \geq \rho_{i_{\delta N}}(X)$. Let i_1, \dots, i_k be the indices corresponding to the largest rank values;
3. Compute the new regions with respect to $\mathcal{G}|_{i_1, \dots, i_k}$.

Algorithm 4.1 represents a first attempt to build an effective method to identify the continuity regions of a target unknown function using a δ NN. Testing it on the best δ NNs selected for the test functions (see Figure 5), we observe extremely good results. In this tests we use the infinity norm for the rank values computations (see (21)), as a preliminary analysis showed better clustering performance with respect to the ℓ_1 and ℓ_2 norms.

As far as functions g_ℓ , g_s , and g_η are concerned, regions returned by Algorithm 4.1 follow very well the actual discontinuity interfaces of the test function, for each $1 \leq k \leq \delta N$ (see Figures 7-9). The only exception is the case of g_γ , in which the method is not able to catch the circumference for $k = 1$ (see Figure 10), and $k = 3$ is needed to reach the target. However, we observe that the

method is able to detect the circumference γ with other non-optimal δ NNs (e.g., see Figure 11); this phenomenon is still under investigation and may suggest that a basic error-based criterion not necessarily select the δ NN that best identify the discontinuity interfaces of the target function.

Nevertheless, δ NNs proved to have the potential for being a new useful tool for the discontinuity detection problem.

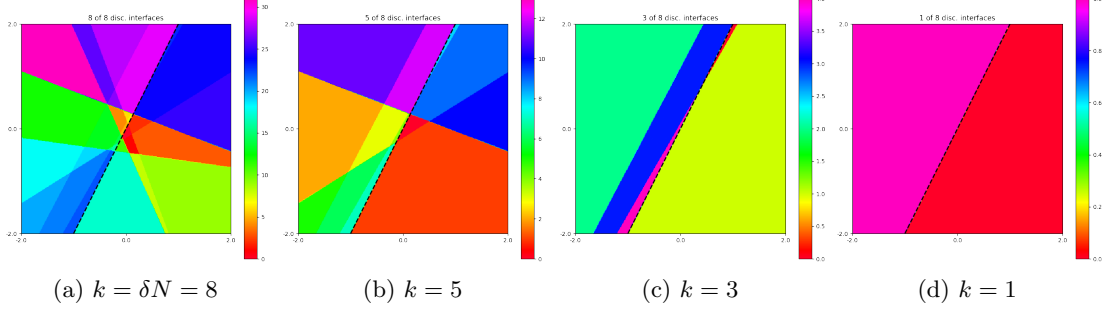


Figure 7: Test 1. Regions returned by Algorithm 4.1 and $\mathcal{A}_{1,5}^8$

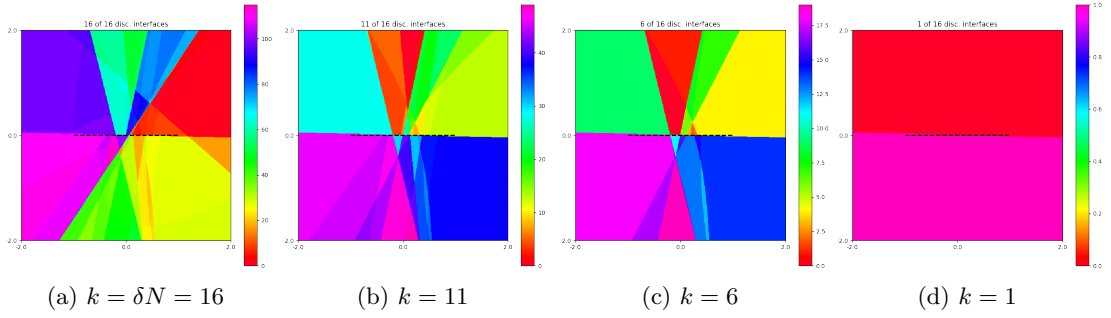


Figure 8: Test 2. Regions returned by Algorithm 4.1 and $\mathcal{B}_{1,5}^8$

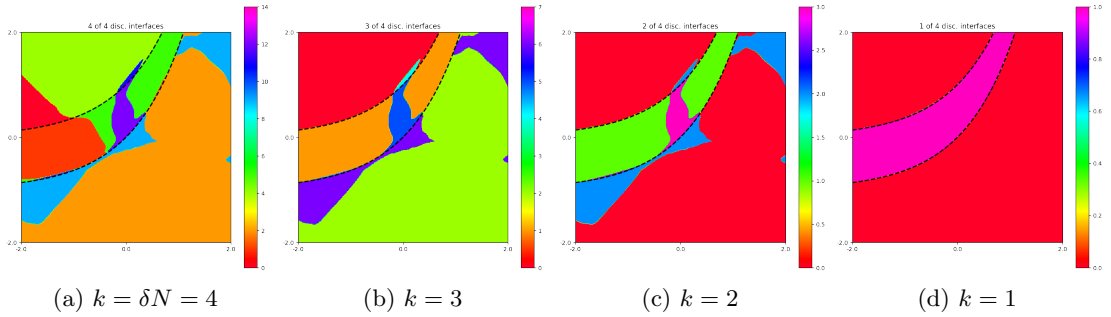


Figure 9: Test 3. Regions returned by Algorithm 4.1 and $\mathcal{A}_{5,5}^4$

5 Conclusions

We have presented a novel typology of layers for Neural Network models, characterized by a discontinuous map where the discontinuity action is obtained adding a vector of multiples of the

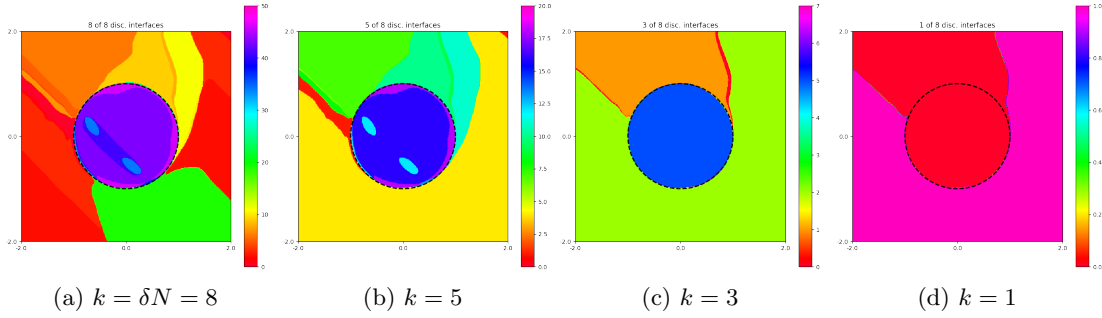


Figure 10: Test 4. Regions returned by Algorithm 4.1 and $\mathcal{B}_{6,7}^4$

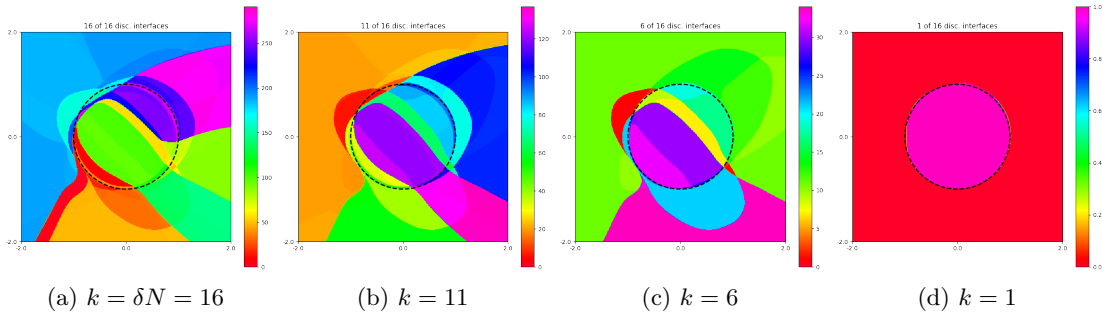


Figure 11: Test 4. Regions returned by Algorithm 4.1 and $\mathcal{B}_{4,7}^8$

Heaviside function applied to input signals of the layer (see (4)); consequently, the function \mathcal{L}_1^{H+1} of a NN with at least one of these discontinuous layers could be discontinuous, too. Denoting by δNN such a NN, we analyzed and studied the theoretical properties that characterize their maps \mathcal{L}_1^{H+1} . Some useful results have been proven (see Section 3), concerning the discontinuities introduced in the δNN s.

We have also illustrated some possible applications of δNN s to discontinuous functions. We considered different examples with increasing complexity and we analyzed the sensitivity of the new NN models in both approximating the discontinuous functions and detecting their discontinuity interfaces using the continuity regions. Extremely interesting results have been obtained, showing a real deal of potential for the new δNN s; indeed, the δNN s proved to have a remarkable ability in detecting the actual discontinuity interfaces of the approximated function, without compromising the function approximation ability typical of the NNs. Since in its basic form the method proposed overestimate the number of continuity regions, we also propose a method for clustering the continuity regions of the δNN in order to have a more precise identification of the actual continuity regions of the original function.

Acknowledgement

Research performed in the framework of the Italian MIUR Award “Dipartimento di Eccellenza 2018-2022” to the Department of Mathematical Sciences, Politecnico di Torino, CUP: E11G1800-0350001. The research has also been partially supported by INdAM-GNCS, by the Italian MIUR PRIN project 201752HKKH8.003, and by the SmartData@PoliTO center for Big Data and Machine Learning technologies.

References

- [1] M. ABADI, A. AGARWAL, P. BARHAM, E. BREVDO, Z. CHEN, C. CITRO, G. S. CORRADO, A. DAVIS, J. DEAN, M. DEVIN, S. GHEMAWAT, I. GOODFELLOW, A. HARP, G. IRVING, M. ISARD, Y. JIA, R. JOZEFOWICZ, L. KAISER, M. KUDLUR, J. LEVENBERG, D. MANÉ, R. MONGA, S. MOORE, D. MURRAY, C. OLAH, M. SCHUSTER, J. SHLENS, B. STEINER, I. SUTSKEVER, K. TALWAR, P. TUCKER, V. VANHOUCKE, V. VASUDEVAN, F. VIÉGAS, O. VINYALS, P. WARDEN, M. WATTENBERG, M. WICKE, Y. YU, AND X. ZHENG, *TensorFlow: Large-scale machine learning on heterogeneous systems*, 2015. Software available from tensorflow.org.
- [2] A. APICELLA, F. DONNARUMMA, F. ISGRÒ, AND R. PREVETE, *A survey on modern trainable activation functions*, arXiv, (2020).
- [3] R. ARCHIBALD, A. GELB, R. SAXENA, AND D. XIU, *Discontinuity detection in multivariate space for stochastic simulations*, Journal of Computational Physics, 228 (2009), pp. 2676–2689.
- [4] R. ARCHIBALD, A. GELB, AND J. YOON, *Polynomial fitting for edge detection in irregularly sampled signals and images*, SIAM Journal on Numerical Analysis, 43 (2005), pp. 259–279.
- [5] E. J. CANDÈS, *Ridgelets: Estimating with ridge functions*, Annals of Statistics, 31 (2003), pp. 1561–1599.
- [6] C. CANUTO, S. PIERACCINI, AND D. XIU, *Uncertainty quantification of discontinuous outputs via a non-intrusive bifidelity strategy*, Journal of Computational Physics, 398 (2019).
- [7] D. CLEVERT, T. UNTERTHINER, AND S. HOCHREITER, *Fast and accurate deep network learning by exponential linear units (elus)*, in 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings, Y. Bengio and Y. LeCun, eds., 2016.
- [8] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep Learning*, MIT Press, 2016. www.deeplearningbook.org.
- [9] I. J. GOODFELLOW, M. MIRZA, A. COURVILLE, AND Y. BENGIO, *Multi-prediction deep Boltzmann machines*, Advances in Neural Information Processing Systems, (2013), pp. 1–9.
- [10] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2016-Decem (2016), pp. 770–778.
- [11] M. IMAIZUMI, *Deep Neural Networks Learn Non-Smooth Functions Effectively*, Proceedings of Machine Learning Research, 89 (2019), pp. 869–878.
- [12] J. D. JAKEMAN, R. ARCHIBALD, AND D. XIU, *Characterization of discontinuities in high-dimensional stochastic problems on adaptive sparse grids*, Journal of Computational Physics, 230 (2011), pp. 3977–3997.
- [13] J. D. JAKEMAN, A. NARAYAN, AND D. XIU, *Minimal multi-element stochastic collocation for uncertainty quantification of discontinuous functions*, Journal of Computational Physics, 242 (2013), pp. 790 – 808.
- [14] P. KIDGER AND T. LYONS, *Universal Approximation with Deep Narrow Networks*, in Proceedings of Thirty Third Conference on Learning Theory, J. Abernethy and S. Agarwal, eds., vol. 125 of Proceedings of Machine Learning Research, PMLR, 09–12 Jul 2020, pp. 2306–2327.

- [15] D. P. KINGMA AND J. L. BA, *Adam: A method for stochastic optimization*, 3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings, (2015), pp. 1–15.
- [16] A. KRIZHEVSKY, I. SUTSKEVER, AND G. HINTON, *ImageNet Classification with Deep Convolutional Neural Networks*, Advances in Neural Information Processing Systems, (2012).
- [17] M. LESHNO, V. Y. LIN, A. PINKUS, AND S. SCHOCKEN, *Multilayer feedforward networks with a nonpolynomial activation function can approximate any function*, Neural Networks, 6 (1993), pp. 861–867.
- [18] W. MA AND J. LU, *An Equivalence of Fully Connected Layer and Convolutional Layer*, Tech. Rep. 3, 2017.
- [19] W. S. MCCULLOCH AND W. H. PITTS, *A Logical Calculus of the Ideas Immanent in Nervous Activity*, Bulletin of Mathematical Biophysics, 5 (1943), pp. 115–133.
- [20] N. M. NAWI, W. H. ATOMI, AND M. REHMAN, *The Effect of Data Pre-processing on Optimized Training of Artificial Neural Networks*, Procedia Technology, 11 (2013), pp. 32–39.
- [21] S. PARK, C. YUN, J. LEE, AND J. SHIN, *Minimum width for universal approximation*, in International Conference on Learning Representations, 2021.
- [22] P. PETERSEN AND F. VOIGTLAENDER, *Optimal approximation of piecewise smooth functions using deep ReLU neural networks*, Neural Networks, 108 (2018), pp. 296–330.
- [23] A. PINKUS, *Approximation theory of the MLP model in neural networks*, Acta Numerica, 8 (1999), pp. 143–195.
- [24] F. ROSENBLATT, *The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain*, Psychological Review, 65 (1958), pp. 386–408.
- [25] Z. SHEN, H. YANG, AND S. ZHANG, *Deep Network With Approximation Error Being Reciprocal of Width to Power of Square Root of Depth*, Neural Computation, 33 (2021), pp. 1005–1036.
- [26] Y. TAIGMAN, M. YANG, M. RANZATO, AND L. WOLF, *DeepFace: Closing the Gap to Human-Level Performance in Face Verification*, Computer Vision Foundation, (2014).
- [27] B. WIDROW AND M. E. HOFF, *Adaptive switching circuits*, 1960 IRE WESCON Convention Record, 4 (1960), pp. 96–104.
- [28] G. ZHANG, C. G. WEBSTER, M. GUNZBUERGER, AND J. BURKARDT, *Hyperspherical sparse approximation techniques for high-dimensional discontinuity detection*, SIAM Review, 58 (2016), pp. 517–551.
- [29] S. ZHANG, X. ZOU, J. AHLQUIST, I. M. NAVON, AND J. G. SELA, *Use of differentiable and nondifferentiable optimization algorithms for variational data assimilation with discontinuous cost functions*, Monthly Weather Review, 128 (2000), pp. 4031–4044.