

On studying neural network expressiveness using topological data analysis and knot theory

Alexandre Louvet

November 9, 2020

Abstract

In this paper we summarize the state of the art on the question of neural network expressiveness both on the theoretical approach to the problem with the study of universal approximators and some practical approaches using topological data analysis and trajectories. We then propose an analysis of the question from a knot theory perspective and share results using studied methods for datasets in dimension 3 and 4.

Contents

1	Neural network expressiveness	3
1.1	Definition	3
1.2	Universal Approximator	3
1.2.1	Sigmoidal functions	4
1.2.2	ReLU functions	5
2	Homology	7
2.1	The idea of homology	7
2.2	Δ -complexes	8
2.3	Simplicial homology	8
2.4	Singular homology	8
2.5	Homology in our research	8
3	Topological data Analysis	9
3.1	Fundamental concepts	9
3.2	Persistent homology	9
4	Knot theory	9
4.1	Definition	9
4.2	Knot determinant	9
4.2.1	Definition	9
4.2.2	Algorithms	9
5	Measuring neural network expressiveness	9
5.1	Using topological data analysis	9
5.2	Using trajectories	9
6	The study of trajectories from a knot theory perspective	9
6.1	Methodology	9
6.2	Algorithms	9
6.3	Results	9
7	Extending the study of expressiveness with topological data analysis	9
7.1	Methodology	9
7.2	Algorithms	9
7.3	Results	9

1 Neural network expressiveness

1.1 Definition

Let I_n denote the n -dimensional unit cube $[0, 1]^n$ and $\mathcal{F}(I_n, \mathbb{R})$ be the space of functions from I_n to \mathbb{R} . We want to study the density of the subsets S_f of $\mathcal{F}(I_n, \mathbb{R})$ that can be written as follows:

$$S_f = \{G_N(x) \in \mathcal{F}(I_n, \mathbb{R}) \mid G(x) = \sum_{i=1}^N \alpha_i f(y_i^T x + \theta_i)\}, N \in \mathbb{N}$$

depending on the choice of $f \in \mathcal{F}(\mathbb{R}, \mathbb{R})$. In the previous equation $y_j \in \mathbb{R}^n$ and $\alpha_j, \theta \in \mathbb{R}$, y^T is the transpose of y and $y^T x$ is the inner product of y and x .

The study of neural network expressiveness consists of the problem described above when f is a function used as an activation function for neural network. The study of density can be on the whole set $\mathcal{F}(I_n, \mathbb{R})$ or on subsets of it such as $\mathbb{C}(I_n, \mathbb{R})$ the set of continuous functions from I_n to \mathbb{R} .

In particular if S_f is dense in a subset $A \subseteq \mathcal{F}(I_n, \mathbb{R})$ we will say that a single-layer feed-forward neural network (Fig. 1) with f as its activation function is a *universal approximator* of A . Considering a neural network has a finite number of nodes neural network expressiveness also consists of the study of the rate of approach of the approximation, i.e. the study of

$$\lim_{N \rightarrow \infty} H(N) = \max_{h \in A} \left(\min_{G_N \in S_f} (\|G_N - h\|) \right) \text{ with } \|\cdot\| \text{ the canonical norm on } \mathcal{F}(I_n, \mathbb{R})$$

The study of that limit and especially of its asymptotic approximation gives an idea of the efficiency of the approximator, i.e. the amount of node to add to the network to improve the approximation.

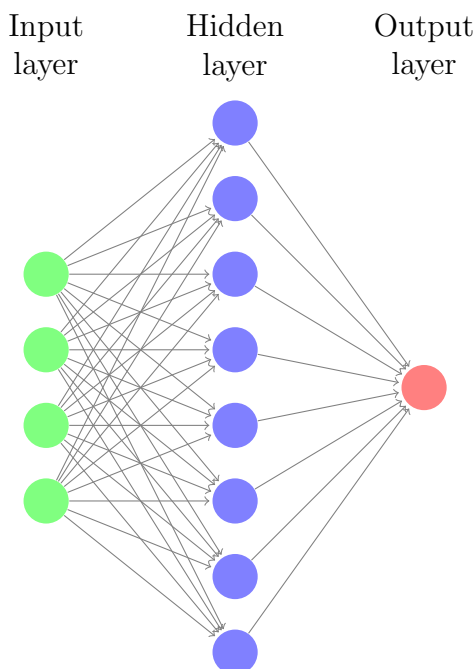


Fig 1: A single-layer feed-forward neural network with $n = 4$ and $N = 8$

1.2 Universal Approximator

In this section we will study the different subsets on which the logistic and ReLU functions acts as universal approximators.

1.2.1 Sigmoidal functions

We say that a function $\sigma \in \mathcal{F}(I_n, \mathbb{R})$ is a sigmoidal function if:

$$\sigma(x) \rightarrow \begin{cases} 0 & \text{as } t \rightarrow +\infty \\ 1 & \text{as } t \rightarrow -\infty \end{cases}$$

The sigmoidal functions include the logistic function defined as:

$$f(x) = \frac{1}{1+e^{-x}}$$

widely used as an activation function for neural networks.

The first study of neural network expressiveness with sigmoidal functions date back to by G.Cybenko in 1989 [1]. He proves that S_σ for σ a sigmoidal function is dense in regards of the supremum norm in $C(I_n, \mathbb{R})$. The demonstration goes as follows.

We denote $M(I_n)$ the space of signed regular Borel measures on I_n

Definition 1 σ is discriminatory if $\mu \in M(I_n)$ and

$$\forall y \in \mathbb{R}^n, \theta \in \mathbb{R} \int_{I_n} \sigma(y^T x + \theta) d\mu(x) = 0 \implies \mu = 0$$

Theorem 1 Let σ be a continuous discriminatory function. Then finite sums of the form

$$G(x) = \sum_{i=1}^N \alpha_i \sigma(y^T x + \theta_i)$$

are dense in $C(I_n, \mathbb{R})$

PROOF: Let $S \subset C(I_n)$ be the set of the function of the form $G(x)$. S is a linear subset of $C(I_n)$. Let us show that the closure of S is $C(I_n)$.

Assume it is not the case. Then the closure of S , denoted R , is a proper subspace of $C(I_n)$. Using the Hahn-Banach theorem, there exists L a bounded linear functional on $C(I_n)$ with $L \neq 0$ and $L(R) = L(S) = 0$

Using the Riesz Representation Theorem, we obtain:

$$L(h) = \int_{I_n} h(x) d\mu(x)$$

for some $\mu \in M(I_n)$, for all $h \in C(I_n)$. Since $\sigma(y^T x + \theta_i) \in R$, we have

$$\forall y, \theta \int_{I_n} \sigma(y^T x + \theta) d\mu(x) = 0$$

Since σ is discriminatory, we have $\mu = 0$ and $L = 0$ follows! Hence the closure of S is $C(I_n)$ and by definition S is dense in $C(I_n)$ \square

Now let us show that sigmoidal functions are discriminatory.

Lemma 1 Any bounded, measurable sigmoidal function, a , is discriminatory. In particular, any continuous sigmoidal function is discriminatory.

PROOF: First note $\forall x, y, \theta, \phi$

$$\sigma_\lambda(x) = \sigma(\lambda(y^T x + \theta) + \phi) \begin{cases} \rightarrow 1 & \text{for } y^T x + \theta > 0 \text{ as } \lambda \rightarrow +\infty \\ \rightarrow 0 & \text{for } y^T x + \theta < 0 \text{ as } \lambda \rightarrow +\infty \\ = \sigma(\phi) & \text{for } y^T x + \theta = 0 \end{cases}$$

Thus $\sigma_\lambda(x)$ converges pointwise and boundedly to:

$$\gamma(x) \begin{cases} = 1 & \text{for } y^T x + \theta > 0 \\ = 0 & \text{for } y^T x + \theta < 0 \\ = \sigma(\phi) & \text{for } y^T x + \theta = 0 \end{cases}$$

as $\lambda \rightarrow +\infty$

Let $\Pi_{y,\theta} = \{x \mid y^T x + \theta = 0\}$ and let $H_{y,\theta} = \{x \mid y^T x + \theta > 0\}$. Lebesgue bounded convergence theorem gives us:

$$\begin{aligned} 0 &= \int \sigma_\lambda(x) d\mu(x) \\ &= \int_{I_n} \gamma(x) d\mu(x) \\ &= \sigma(\phi) \mu(\Pi_{y,\theta}) + \mu(H_{y,\theta}) \end{aligned}$$

for all ϕ, θ, y

Fix y , we write

$$F(h) = \int_{I_n} h(y^T x) d\mu(x)$$

Note that F is a bounded function on $L^\infty(\mathbb{R})$ since μ is a signed measure. By choosing h as the indicator function on $[\theta, \infty[$, we have:

$$F(h) = \int_{I_n} h(y^T x) d\mu(x) = \mu(\Pi_{y,-\theta}) + \mu(H_{y,-\theta}) = 0$$

By linearity $F(h) = 0$ for indicator function on any interval and hence for any simple function (sum of indicator functions) and since simple functions are dense in $L^\infty(\mathbb{R})$, $F = 0$. In particular it is true for the bounded function $s(u) = \sin(m.u)$ and $c(u) = \cos(m.u)$. It gives:

$$F(s + ic) = \int_{I_n} \cos(m^T x) + i \sin(m^T x) d\mu(x) = \int_{I_n} \exp(im^T x) d\mu(x) = 0$$

for all m . Therefore the fourier transform of μ is 0 and μ must be 0. Hence σ is discriminatory. \square

This proves that any function of $C(I_n, \mathbb{R})$ can be approximated by a single-layer network with sigmoidal functions as activation function.

In 1991, Hornik extended in [2] gave the proof for bounded non-constant functions of $F(I_n, \mathbb{R})$. The proof is very similar to the original one.

1.2.2 ReLU functions

ReLU functions stands for Rectified Linear Units, there are formed of two pieces of linear functions. They gained interest recently by showing convincing results in a lot of different applications.

In 2016 Arora et al. [3] showed a version of the universal approximation theorem for piecewise linear function that includes ReLU functions.

Definition 2 We say a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous piecewise linear (PWL) if there exists a finite set of polyhedra whose union is \mathbb{R}^n , and f is affine linear over each polyhedron (implies continuity because affine regions are closed and cover \mathbb{R}^n)

Proposition 1 Every function in $L^q(\mathbb{R}^n)$, ($1 \leq q \leq \infty$) can be arbitrarily approximated in the L^q norm (which for a function f is given by $\|f\|_q = (\int |f|^q)^{1/q}$) by a ReLU Deep Neural Network.

PROOF:

We know that any ReLU DNN represents a PWL function, let's prove the converse.

By theorem 1 in [4] any piecewise linear function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, can be written:

$$f = \sum_{j=1}^p s_j (\max_{i \in S_j} l_i)$$

with l_i ($i \leq i \leq k$) linear functions, S_i ($1 \leq i \leq p$) $\subseteq \{1, \dots, k\}$ with $\forall i, |S_i| \leq n + 1$ and $\forall j \in \{1, \dots, p\}, s_j \in \{-1, +1\}$

It means that any PWL convex function can be represented as a linear combination of at most $n + 1$ affine pieces. That's to say a ReLU DNN with size $n + 1$.

Let $p \in \mathbb{N}^*$, let $f \in L^p(\mathbb{R}^n)$, consider the function sequence:

$$f_n(x) = (x - \frac{k}{n})f(\frac{k}{n}) + (1 - x + \frac{k}{n})f(\frac{k+1}{n}) \text{ with } \frac{k}{n} \leq x < \frac{k+1}{n}, n \geq 1$$

$f_n \xrightarrow[n \rightarrow \infty]{} f$, and $\forall n, f_n$ is a PWL continuous function. Therefore the continuous PWL functions are dense in $L^p(\mathbb{R}^n)$.

Since any PWL function $\mathbb{R}^n \rightarrow \mathbb{R}$ is representable by a ReLU DNN, universal approximation theorem follows from the fact that the family of continuous PWL functions is dense in any $L^p(\mathbb{R}^n)$ space for $1 \leq p \leq \infty$.

□

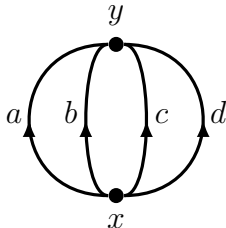
In this article, Arora et al. also give a higher bound for the depth of the network. they show that the network required for any function $f \in L^q(\mathbb{R}^n)$ is at most $\lceil \log_2(n + 1) \rceil$.

In 2018, Hannin & Selke [5] showed that the minimal width a ReLU neural network must have in order to approximate any function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is $n + m$. Giving a lower bound to the width of the ReLU network.

2 Homology

The study of homology is part of the field of algebraic topology. We will first introduce the concept of homology, then take a look at the theoretical concept behind it by introducing Δ -complex: a primitive topological structure that allow the study of homology. Then with the introduction of simplicial homology and finally singular homology we will complete the presentation of the fundamental objects used in the study of homology. We will then reconsider these tools in the context of our research. This part heavily relies on the book *Algebraic Topology* by A. Hatcher [6] and the class given by P. Albin at University of Illinois [7] with the same book as class material.

2.1 The idea of homology

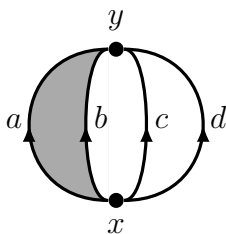


Consider the space X_1 in the figure drawn above. We can see two points x and y as well as four edges a , b , c and d . Consider the loops formed by travelling along the edges, for instance the path ab^{-1} is a loop with x as the basepoint. Something to consider is that the loop $b^{-1}a$ is basically the same loop but starting from a different basepoint, y in that case. By abelianizing, we can consider cycles instead of loops without any basepoint.

The abelian groups having only one operation, we will then switch to additive notations to remain in accordance with Hatcher's notations. With that new notation we can write equalities such as $(a - c) + (b - d) = (a - d) + (b - c)$. This is justified by the fact that from an algebraic point of view, there is no difference between these two cycles.

Let us call any linear combination of edges, chains. Then the condition for a chain $ka + lb + mc + nd$ with $k, l, m, n \in \mathbb{N}$ to be a cycle is $k + l + m + n = 0$.

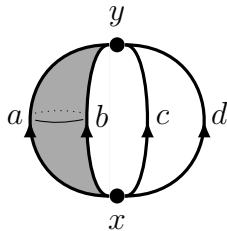
Now let us start to formalize these concepts. Consider C_0 the free abelian group spanned by the vertices x, y and let C_1 be the free abelian group spanned by the edges a, b, c, d . One can define an homomorphism $\delta : C_1 \rightarrow C_0$ that sends any basis element of C_1 to $x - y$. Thus we have $\delta(ka + lb + mc + nd) = (k + l + m + n)a - (k + l + m + n)y$ and $\ker(\delta) = \{ \text{cycles of } X_1 \}$. $\{a - b, b - c, c - d\}$ forms a basis for this kernel, i.e. any cycle is a linear combination of the most obvious cycles of X_1 , that conveys the information that X_1 has three visible "holes" formed by the space between its four edges.



Now fill the hole between edges a and b with a 2-cell, namely A to create a new space X_2 . Considering it as being oriented clockwise, its boundary is $a - b$. The cycle $a - b$ is homotopic to a point since it can now be contracted by sliding over A , it no longer encloses a hole in G . That suggests that homology should consider only the quotient group of what we found previously by $a - b$. In this quotient group the cycles $a - c$ and $b - c$ are equivalent. Which is consistent with the fact that they are homotopic in X_2 .

Algebraically we can now consider a pair of homomorphism δ_1, δ_2 such that $C_2 \xrightarrow{\delta_2} C_1 \xrightarrow{\delta_1} C_0$ where C_2 is the cyclic group spanned by A and $\delta_2(A) = a - b$. The quotient group we are interested in is $\ker(\delta_1)/\text{Im}(\delta_2)$, i.e. the 1-dimensional cycles modulo those that are

boundaries. We will call it the homology group $H_1(X_2)$. It can also be computed in X_1 by considering C_2 to be zero since there are no two cells in X_2 . $H_1(X_2)$ now only admits two generators $b - c$ and cd expressing the geometric fact that the number of holes reduced from 3 to 2.



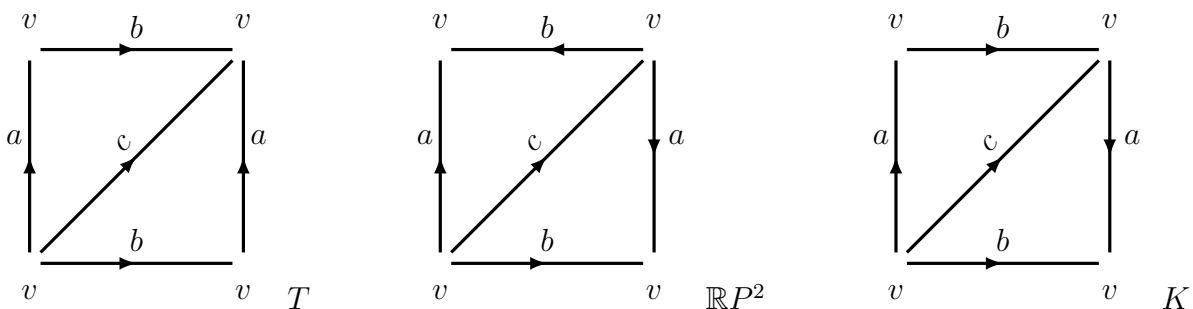
Suppose we attach another 2-cell between a and b to create X_3 , namely B . C_2 now consists of linear combinations of A and B and $\delta_2(A) = \delta_2(B) = a - b$. In the one hand, $H_1(X_3) = H_1(X_2) \approx \mathbb{Z} \times \mathbb{Z}$, but this time $\ker(\delta_2) \neq 0$ and $A - B$ is its generator, hence we have $H_2(X_3) = \ker(\delta_2) \approx \mathbb{Z}$. Topologically the cycle $A - B$ is equivalent to a sphere and it detects the presence of a "hole" enclosed by this sphere rather than a circle.

The pattern we can see appear with these examples is rather clear. The n -cell complexes of X forms free abelian groups $C_n(X)$ and one can define homomorphisms $\delta_n : C_n \rightarrow C_{n-1}$ to define homology groups $H_n(X) = \ker(\delta_n) / \text{Im}(\delta_{n+1})$. The only problem now is to define δ_n for any n . If it is simple for small n (head minus tail for a vertex), it becomes rather complicated when n grows and more complex polyhedral cells appear in X . The most efficient approach is to decompose polyhedra into simplices which allows simple orientation and boundary computation. For that purpose we will define nice cellular complexes that will be the fundamental element for homology computation in the homology theories presented after.

2.2 Δ -complexes

The most important theory of algebraic topology is called simplicial homology, since it is technically complicated we will first introduce a simpler version of it called simplicial homology. The natural definition spaces of simplicial homology is called Δ -complexes that we will introduce in this part.

The projective plane, the klein bottle and the torus can be obtained from squares by identifying edges and giving them orientation as draw below.



2.3 Simplicial homology

2.4 Singular homology

2.5 Homology in our research

Homology exists as a computable alternative homotopy.

3 Topological data Analysis

3.1 Fundamental concepts

3.2 Persistent homology

4 Knot theory

4.1 Definition

4.2 Knot determinant

4.2.1 Definition

4.2.2 Algorithms

5 Measuring neural network expressiveness

5.1 Using topological data analysis

5.2 Using trajectories

6 The study of trajectories from a knot theory perspective

6.1 Methodology

6.2 Algorithms

6.3 Results

7 Extending the study of expressiveness with topological data analysis

7.1 Methodology

7.2 Algorithms

7.3 Results

References

- [1] G. Cybenko, “Approximation by superpositions of a sigmoidal function,” *Mathematics of Control, Signals and Systems*, vol. 2, pp. 303–314, Dec. 1989.
- [2] K. Hornik, “Approximation capabilities of multilayer feedforward networks,” *Neural Networks*, vol. 4, pp. 251–257, Jan. 1991.
- [3] R. Arora, A. Basu, P. Mianjy, and A. Mukherjee, “Understanding Deep Neural Networks with Rectified Linear Units,” *arXiv:1611.01491 [cond-mat, stat]*, Feb. 2018. arXiv: 1611.01491.
- [4] S. Wang and X. Sun, “Generalization of hinging hyperplanes,” *IEEE Transactions on Information Theory*, vol. 51, pp. 4425–4431, Dec. 2005.
- [5] B. Hanin and M. Sellke, “Approximating Continuous Functions by ReLU Nets of Minimal Width,” *arXiv:1710.11278 [cs, math, stat]*, Mar. 2018. arXiv: 1710.11278.
- [6] A. Hatcher, *Algebraic topology*. Cambridge ; New York: Cambridge University Press, 2002.
- [7] P. Albin, “1. History of Algebraic Topology; Homotopy Equivalence - Pierre Albin,” 2018.