

Jinchao Xu

Deep Learning Algorithms and Analysis

Summer 2020

Contents

1	FEM and DNN	5
1.1	Motivation: from finite element to neural network	5
1.2	Why we need deep neural networks via composition	7
1.2.1	FEM ans DNN_1 in 1D	7
1.2.2	Linear finite element cannot be recovered by DNN_1 for $d \geq 2$	8
1.3	Definition of neural network space	11

1.1 Motivation: from finite element to neural network

In this chapter, we will introduce the so-called shallow neural network (deep neural network with one hidden layer) from the viewpoint of finite element method.

Let us first consider the linear finite element functions on the unit interval $\bar{\mathcal{Q}} = [0, 1]$ in 1D. We then consider a set of equidistant grids \mathcal{Q}_ℓ of level ℓ on the unit interval $\bar{\mathcal{Q}} = [0, 1]$ and mesh length $h_\ell = 2^{-\ell}$. The grid points $x_{\ell,i}$ are given by

$$(1.1) \quad x_{\ell,i} := ih_\ell, \quad 0 \leq i \leq 2^\ell.$$

For $\ell = 1$, we denote the special hat function:

$$(1.2) \quad \varphi(x) = \begin{cases} 2x & x \in [0, \frac{1}{2}] \\ 2(1-x) & x \in [\frac{1}{2}, 1] \\ 0, & \text{others} \end{cases}.$$

The next diagram shows this basis function:

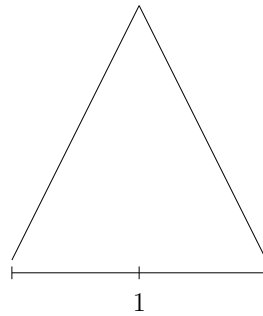


Fig. 1.1. Diagram of $\varphi(x)$

Then, for any nodal basis function $\varphi_{\ell,i}$ as below:

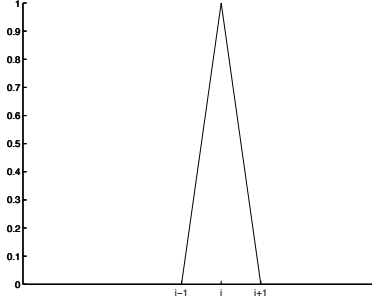


Fig. 1.2. Diagram of $\varphi_{\ell,i}(x)$

in a fine grid \mathcal{T}_ℓ can be written as

$$(1.3) \quad \varphi_{\ell,i} = \varphi\left(\frac{x - x_{\ell,i-1}}{2h_\ell}\right) = \varphi(w_\ell x + b_{\ell,i}).$$

That is to say, any $\varphi_{\ell,i}(x)$ can be obtained from $\varphi(x)$ by scaling (dilation) and translation with

$$(1.4) \quad w_\ell = 2^{\ell-1}, \quad b_{\ell,i} = \frac{-(i-1)}{2},$$

in $\varphi_{\ell,i} = \varphi(w_\ell x + b_{\ell,i})$.

Let recall the finite element interpolation as

$$(1.5) \quad u(x) \approx u_\ell(x) := \sum_{0 \leq i \leq 2^\ell} u(x_{\ell,i}) \varphi_{\ell,i}(x),$$

for any smooth function $u(x)$ on $(0, 1)$. The above interpolation will converge as $\ell \rightarrow \infty$, which show that

$$(1.6) \quad \text{span} \{ \varphi(w_\ell x + b_{\ell,i}) \} \quad \text{is dense in} \quad H^1(0, 1).$$

Thus, we may have the next concise relation:

$$(1.7) \quad \text{FE space} = \text{span} \{ \varphi(w_\ell x + b_{\ell,i}) \mid 0 \leq i \leq 2^\ell, \ell = 1, 2, \dots \} \subset \text{span} \{ \varphi(wx + b) \mid w, b \in \mathbb{R} \}.$$

In other words, the finite element space can be understood as the linear combination of $\varphi(wx + b)$ with certain special choice of w and b .

Here, we need to point that this span $\{ \varphi(wx + b) \mid w, b \in \mathbb{R} \}$ is exactly the deep neural networks with one hidden layer (shallow neural networks) with activation function $\varphi(x)$. More precisely,

$$(1.8) \quad f \in \text{span} \{ \varphi(wx + b) \mid w, b \in \mathbb{R} \},$$

means there exist positive integer N and $w_j, b_j \in \mathbb{R}$ such that

$$(1.9) \quad f = \sum_{j=1}^N a_j \varphi(w_j x + b_j).$$

The above function is also called one hidden neural network function with N neurons.

Remark 1. 1. By making w_ℓ and $b_{\ell,i}$ arbitrary, we get a much larger class of function which is exact a special neural network with activation function $\varphi(x)$.

2. Generalizations:

- a) φ can be different, such as $\text{ReLU}(x) = \max\{0, x\}$.
- b) There is a natural extension for high dimension d as

$$(1.10) \quad \{\varphi(w \cdot x + b)\},$$

where $w \in \mathbb{R}^d$, $b \in \mathbb{R}$ and $w \cdot x = \sum_{i=1}^d w_i x_i$. This is called “deep” neural network with one hidden layer.

1.2 Why we need deep neural networks via composition

1.2.1 FEM and DNN₁ in 1D

Thanks to the connection between $\varphi(x)$ and $\text{ReLU}(x) = \{0, x\}$

$$(1.11) \quad \varphi(x) = 2\text{ReLU}(x) - 4\text{ReLU}(x - \frac{1}{2}) + 2\text{ReLU}(x - 1).$$

It suffices to show that each basis function $\varphi_{\ell,i}$ can be represented by a ReLU DNN. We first note that the basis function φ_i has the support in $[x_{i-1}, x_{i+1}]$ can be easily written as

$$(1.12) \quad \varphi_{\ell,i}(x) = \frac{1}{h_\ell} \text{ReLU}(x - x_{\ell,i-1}) - (\frac{2}{h_\ell}) \text{ReLU}(x - x_{\ell,i}) + \frac{1}{h_\ell} \text{ReLU}(x - x_{\ell,i+1}).$$

More generally, if function φ_i is not on the uniform grid but has support in $[x_{i-1}, x_{i+1}]$ can be easily written as

$$(1.13) \quad \varphi_i(x) = \frac{1}{h_{i-1}} \text{ReLU}(x - x_{i-1}) - (\frac{1}{h_{i-1}} + \frac{1}{h_i}) \text{ReLU}(x - x_i) + \frac{1}{h_i} \text{ReLU}(x - x_{i+1}),$$

where $h_i = x_{i+1} - x_i$.

Thus is to say, we have the next theorem.

Theorem 1. For $d = 1$, and $\Omega \subset \mathbb{R}^d$ is a bounded interval, then DNN₁ can be used to cover all linear finite element function in on Ω .

1.2.2 Linear finite element cannot be recovered by DNN_1 for $d \geq 2$

In view of Theorem 1 and the fact that $\text{DNN}_J \subseteq \text{DNN}_{J+1}$, it is natural to ask that how many layers are needed at least to recover all linear finite element functions in \mathbb{R}^d for $d \geq 2$. In this section, we will show that

$$(1.14) \quad J_d \geq 2, \quad \text{if } d \geq 2,$$

where J_d is the minimal J such that all linear finite element functions in \mathbb{R}^d can be recovered by DNN_J .

In particular, we will show the following theorem.

Theorem 2. *If $\Omega \subset \mathbb{R}^d$ is either a bounded domain or $\Omega = \mathbb{R}^d$, DNN_1 can not be used to recover all linear finite element functions on Ω .*

Proof. We prove it by contradiction. Let us assume that for any continuous piecewise linear function $f : \Omega \rightarrow \mathbb{R}$, we can find finite $N \in \mathbb{N}$, $w_i \in \mathbb{R}^{1,d}$ as row vector and $\alpha_i, b_i, \beta \in \mathbb{R}$ such that

$$f = \sum_{i=1}^N \alpha_i \text{ReLU}(w_i x + b_i) + \beta,$$

with $f_i = \alpha_i \text{ReLU}(w_i x + b_i)$, $\alpha_i \neq 0$ and $w_i \neq 0$. Consider the finite element functions, if this one hidden layer ReLU DNN can recover any basis function of FEM, then it can recover the finite element space. Thus let us assume f is a locally supported basis function for FEM. Furthermore, if Ω is a bounded domain, we assume that

$$(1.15) \quad d(\text{supp}(f), \partial\Omega) > 0,$$

with

$$d(A, B) = \inf_{x \in A, y \in B} \|x - y\|,$$

as the distance of two closed sets.

A more important observation is that $\nabla f : \Omega \rightarrow \mathbb{R}^d$ is a piecewise constant vector function. The key point is to consider the discontinuous points for $g := \nabla f = \sum_{i=1}^N \nabla f_i$.

For more general case, we can define the set of discontinuous points of a function by

$$D_g := \{x \in \Omega \mid x \text{ is a discontinuous point of } g\}.$$

Because of the property that

$$(1.16) \quad D_{f+g} \supseteq D_f \cup D_g \setminus (D_f \cap D_g),$$

we have

$$(1.17) \quad D_{\sum_{i=1}^N g_i} \supseteq \bigcup_{i=1}^N D_{g_i} \setminus \bigcup_{i \neq j} (D_{g_i} \cap D_{g_j}).$$

Note that

$$(1.18) \quad g_i = \nabla f_i(x) = \nabla (\alpha_i \text{ReLU}(w_i x + b_i)) = (\alpha_i H(w_i x + b_i)) w_i \in \mathbb{R}^d,$$

for $i = 1 : N$ with H be the Heaviside function defined as:

$$H(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ 1 & \text{if } x > 0. \end{cases}$$

This means that

$$(1.19) \quad D_{g_i} = \{x \mid w_i x + b_i = 0\}$$

is a $d - 1$ dimensional affine space in \mathbb{R}^d .

Without loss of generality, we can assume that

$$(1.20) \quad D_{g_i} \neq D_{g_j}.$$

When the other case occurs, i.e. $D_{g_{\ell_1}} = D_{g_{\ell_2}} = \dots = D_{g_{\ell_k}}$, by the definition of g_i in (1.18) and D_{g_i} in (1.19), this happens if and only if there is a row vector (w, b) such that

$$(1.21) \quad c_{\ell_i} (w \ b) = (w_{\ell_i} \ b_{\ell_i}),$$

with some $c_{\ell_i} \neq 0$ for $i = 1 : k$. We combine those g_{ℓ_i} as

$$\begin{aligned} \tilde{g}_\ell &= \sum_{i=1}^k g_{\ell_i} = \sum_{i=1}^k \alpha_{\ell_i} H(w_{\ell_i} x + b_{\ell_i}) w_{\ell_i}, \\ &= \sum_{i=1}^k (c_{\ell_i} \alpha_{\ell_i} H(c_{\ell_i} (wx + b))) w, \\ &= \begin{cases} \left(\sum_{i=1}^k c_{\ell_i} \alpha_{\ell_i} H(c_{\ell_i}) \right) w & \text{if } wx + b > 0, \\ \left(\sum_{i=1}^k c_{\ell_i} \alpha_{\ell_i} H(-c_{\ell_i}) \right) w & \text{if } wx + b \leq 0. \end{cases} \end{aligned}$$

Thus, if

$$\left(\sum_{i=1}^k c_{\ell_i} \alpha_{\ell_i} H(c_{\ell_i}) \right) = \left(\sum_{i=1}^k c_{\ell_i} \alpha_{\ell_i} H(-c_{\ell_i}) \right),$$

\tilde{g}_ℓ is a constant vector function, that is to say $D_{\sum_{i=1}^k g_{\ell_i}} = D_{\tilde{g}_\ell} = \emptyset$. Otherwise, \tilde{g}_ℓ is a piecewise constant vector function with the property that

$$D_{\sum_{i=1}^k g_{\ell_i}} = D_{\tilde{g}_\ell} = D_{g_{\ell_i}} = \{x \mid wx + b = 0\}.$$

This means that we can use condition (1.21) as an equivalence relation and split $\{g_i\}_{i=1}^N$ into some groups, and we can combine those g_{ℓ_i} in each group as what we do above. After that, we have

$$\sum_{i=1}^N g_i = \sum_{\ell=1}^{\tilde{N}} \tilde{g}_\ell,$$

with $D_{\tilde{g}_s} \neq D_{\tilde{g}_t}$. Finally, we can have that $D_{\tilde{g}_s} \cap D_{\tilde{g}_t}$ is an empty set or a $d - 2$ dimensional affine space in \mathbb{R}^d . Since $\tilde{N} \leq N$ is a finite number,

$$D := \bigcup_{i=1}^N D_{\tilde{g}_i} \setminus \bigcup_{s \neq t} (D_{\tilde{g}_s} \cap D_{\tilde{g}_t})$$

is an unbounded set.

- If $\Omega = \mathbb{R}^d$,

$$\text{supp}(f) \supseteq D_g = D_{\sum_{i=1}^N g_i} = D_{\sum_{\ell=1}^{\tilde{N}} \tilde{g}_\ell} \supseteq D,$$

is contradictory to the assumption that f is locally supported.

- If Ω is a bounded domain,

$$d(D, \partial\Omega) = \begin{cases} s > 0 & \text{if } D_{\tilde{g}_i} \cap \Omega = \emptyset, \forall i \\ 0 & \text{otherwise.} \end{cases}$$

Note again that all $D_{\tilde{g}_i}$'s are $d - 1$ dimensional affine spaces, while $D_{\tilde{g}_i} \cap D_{\tilde{g}_j}$ is either an empty set or a $d - 2$ dimensional affine space. If $d(D, \partial\Omega) > 0$, this implies that ∇f is continuous in Ω , which contradicts the assumption that f is a basis function in FEM. If $d(D, \partial\Omega) = 0$, this contradicts the previous assumption in (1.15).

Hence DNN_1 cannot recover any piecewise linear function in Ω for $d \geq 2$. \square

Following the proof above, we have the following theorem.

Theorem 3. $\{\text{ReLU}(w_i x + b_i)\}_{i=1}^m$ are linearly independent if (w_i, b_i) and (w_j, b_j) are linearly independent in $\mathbb{R}^{1 \times (d+1)}$ for any $i \neq j$.

1.3 Definition of neural network space

1. Primary variables $n_0 = d$

$$x^0 = x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

2. n_1 hyperplanes

$$W^1 x + b^1 = \begin{pmatrix} w_1^1 x + b_1^1 \\ w_2^1 x + b_2^1 \\ \vdots \\ w_n^1 x + b_n^1 \end{pmatrix}$$

3. n_1 -neurons:

$$x^1 = \sigma(W^1 x + b^1) = \begin{pmatrix} \sigma(w_1^1 x + b_1^1) \\ \sigma(w_2^1 x + b_2^1) \\ \vdots \\ \sigma(w_n^1 x + b_n^1) \end{pmatrix}$$

4. n_2 -hyperplanes

$$W^2 x^1 + b^2 = \begin{pmatrix} w_1^2 x^1 + b_1^2 \\ w_2^2 x^1 + b_2^2 \\ \vdots \\ w_n^2 x^1 + b_n^2 \end{pmatrix}$$

Shallow neural network functions:

(1.22)

$${}_n\mathbf{N}(n_1, n_2) = {}_n\mathbf{N}(\sigma; n_1, n_2) = \left\{ W^2 x^1 + b^2, x^1 = \sigma(W^1 x + b^1) \text{ with } W^\ell \in \mathbb{R}^{n_\ell \times n_{\ell-1}}, b^\ell \in \mathbb{R}^{n_\ell}, \ell = 1, 2, n_0 = d \right\}$$

(1.23)

$${}_n\mathbf{N}(\sigma; n_1, n_2, \dots, n_L) = \left\{ W^L x^{L-1} + b^L, x^\ell = \sigma(W^\ell x^{\ell-1} + b^\ell) \text{ with } W^\ell \in \mathbb{R}^{n_\ell \times n_{\ell-1}}, b^\ell \in \mathbb{R}^{n_\ell}, \ell = 1 : L, n_0 = d, x^0 = x \right\}$$

First, let us first define the so-called 1-hidden layer (shallow) neural network.

The 1-hidden layer (shallow) neural network is defined as:

$$(1.24) \quad {}_n\mathbf{N} = {}_n\mathbf{N}(\sigma) = {}_n\mathbf{N}^1(\sigma) = \bigcup_{n_1 \geq 1} {}_n\mathbf{N}(\sigma; n_1, 1)$$

The 2-hidden layer (shallow) neural network is defined as:

$$(1.25) \quad {}_n\mathbf{N}^2(\sigma) = \bigcup_{n_1, n_2 \geq 1} {}_n\mathbf{N}(\sigma; n_1, n_2, 1)$$