

# Analysis of Neural Network

Jinchao Xu 许进超

Penn State University

[xu@math.psu.edu](mailto:xu@math.psu.edu) <http://www.math.psu.edu/xu/>

CUHK, June 2021

Acknowledgement:

NSF: DMS-1819157

## 1 Introduction

## 2 Finite Element and Neural Networks

- Neural Network Functions
- Connection of ReLU DNN and linear FEM

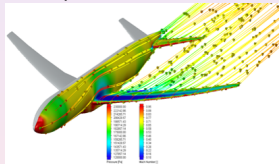
## 3 Neural Network Approximation Class

- Density of Shallow Neural Network
- Approximation Properties
- Approximation Spaces

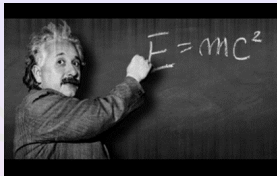
# Four major methods of scientific research



Experimental Science



Computational Science



Theoretical Science

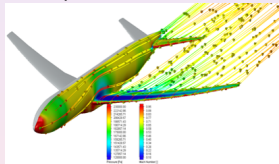


Data Science

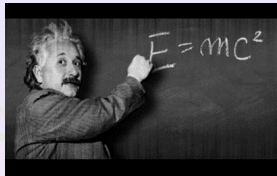
# Four major methods of scientific research



Experimental Science



Computational Science



Theoretical Science



Data Science

Question:

*Is "Data Science" really a science?*

# A basic AI problem: classification

- Can a machine (function) tell the difference ?



# Supervised learning

- Function interpolation (data fitting)
  - ▶ Each image = a big vector of pixel values
    - ★  $d = 1280 \times 720 \times 3$  (width  $\times$  height  $\times$  RGB channel)  $\approx 3\text{M}$ .

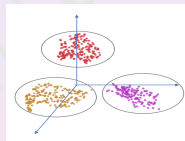
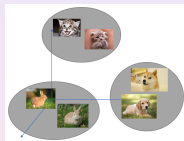
# Supervised learning

- Function interpolation (data fitting)

- ▶ Each image = a big vector of pixel values

- ★  $d = 1280 \times 720 \times 3$  (width  $\times$  height  $\times$  RGB channel)  $\approx 3\text{M}$ .

- ▶ 3 different sets of points in  $\mathbb{R}^d$ , are they separable?



# Supervised learning

- Function interpolation (data fitting)

- ▶ Each image = a big vector of pixel values
  - ★  $d = 1280 \times 720 \times 3$  (width  $\times$  height  $\times$  RGB channel)  $\approx 3\text{M}$ .
- ▶ 3 different sets of points in  $\mathbb{R}^d$ , are they separable?



- ▶ Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:



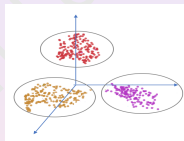
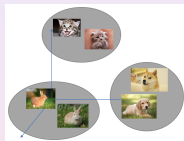
# Supervised learning

- Function interpolation (data fitting)

- ▶ Each image = a big vector of pixel values

- ★  $d = 1280 \times 720 \times 3$  (width  $\times$  height  $\times$  RGB channel)  $\approx 3\text{M}$ .

- ▶ 3 different sets of points in  $\mathbb{R}^d$ , are they separable?



- ▶ Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

$$f(\text{cat image}; \Theta) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

# Supervised learning

- Function interpolation (data fitting)

- ▶ Each image = a big vector of pixel values

- ★  $d = 1280 \times 720 \times 3$  (width  $\times$  height  $\times$  RGB channel)  $\approx 3\text{M}$ .

- ▶ 3 different sets of points in  $\mathbb{R}^d$ , are they separable?



- ▶ Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

$$f(\text{cat image}; \Theta) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad f(\text{dog image}; \Theta) \approx \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

# Supervised learning

- Function interpolation (data fitting)

- ▶ Each image = a big vector of pixel values

- ★  $d = 1280 \times 720 \times 3$  (width  $\times$  height  $\times$  RGB channel)  $\approx 3M$ .

- ▶ 3 different sets of points in  $\mathbb{R}^d$ , are they separable?



- ▶ Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

$$f\left(\begin{array}{c} \text{cat} \end{array}; \Theta\right) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad f\left(\begin{array}{c} \text{dog} \end{array}; \Theta\right) \approx \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad f\left(\begin{array}{c} \text{cat} \end{array}; \Theta\right) \approx \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

# Supervised learning

- Function interpolation (data fitting)

- ▶ Each image = a big vector of pixel values

- ★  $d = 1280 \times 720 \times 3$  (width  $\times$  height  $\times$  RGB channel)  $\approx 3M$ .

- ▶ 3 different sets of points in  $\mathbb{R}^d$ , are they separable?



- ▶ Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

$$f(\text{cat image}; \Theta) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad f(\text{dog image}; \Theta) \approx \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad f(\text{cat image}; \Theta) \approx \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

# How to formulate “learning”?

Jinchao Xu  
xu@multigrid.org

# How to formulate “learning”?

- Data:  $\{x_i, y_i\}_{i=1}^m$

Jinchao Xu  
xu@multigrid.org

# How to formulate “learning”?

- Data:  $\{x_i, y_i\}_{i=1}^m$
- Find  $f^*$  in some function class such that  $f^*(x_i) \approx y_i$ .

# How to formulate “learning”?

- Data:  $\{x_i, y_i\}_{i=1}^m$
- Find  $f^*$  in some function class such that  $f^*(x_i) \approx y_i$ .
- Mathematically: solve the **optimization problem** by parameterizing the abstract function class

$$\min_{\Theta} L(x, y, \Theta) \quad (1)$$

where

$$L(x, y, \Theta) = \mathbb{E}_{(x, y) \sim \mathcal{D}} [L(f(x; \Theta), y)] \approx \frac{1}{N} \sum_{i=1}^N \|f(x_i; \Theta) - y_i\|^2$$

- ▶ Or combine the feature map  $f$  with the logistic regression model to obtain **cross-entropy** loss function.



# How to formulate “learning”?

- Data:  $\{x_i, y_i\}_{i=1}^m$
- Find  $f^*$  in some function class such that  $f^*(x_i) \approx y_i$ .
- Mathematically: solve the **optimization problem** by parameterizing the abstract function class

$$\min_{\Theta} L(x, y, \Theta) \quad (1)$$

where

$$L(x, y, \Theta) = \mathbb{E}_{(x, y) \sim \mathcal{D}} [L(f(x; \Theta), y)] \approx \frac{1}{N} \sum_{i=1}^N \|f(x_i; \Theta) - y_i\|^2$$

- ▶ Or combine the feature map  $f$  with the logistic regression model to obtain **cross-entropy** loss function.
- Application: image classification:

# How to formulate “learning”?

- Data:  $\{x_i, y_i\}_{i=1}^m$
- Find  $f^*$  in some function class such that  $f^*(x_i) \approx y_i$ .
- Mathematically: solve the **optimization problem** by parameterizing the abstract function class

$$\min_{\Theta} L(x, y, \Theta) \quad (1)$$

where

$$L(x, y, \Theta) = \mathbb{E}_{(x, y) \sim \mathcal{D}} [L(f(x; \Theta), y)] \approx \frac{1}{N} \sum_{i=1}^N \|f(x_i; \Theta) - y_i\|^2$$

- ▶ Or combine the feature map  $f$  with the logistic regression model to obtain **cross-entropy** loss function.

- Application: image classification:

$$f(\boxed{?}; \Theta) = \begin{pmatrix} 0.7 \\ 0.2 \\ 0.1 \end{pmatrix}$$

# How to formulate “learning”?

- Data:  $\{x_i, y_i\}_{i=1}^m$
- Find  $f^*$  in some function class such that  $f^*(x_i) \approx y_i$ .
- Mathematically: solve the **optimization problem** by parameterizing the abstract function class

$$\min_{\Theta} L(x, y, \Theta) \quad (1)$$

where

$$L(x, y, \Theta) = \mathbb{E}_{(x, y) \sim \mathcal{D}} [L(f(x; \Theta), y)] \approx \frac{1}{N} \sum_{i=1}^N \|f(x_i; \Theta) - y_i\|^2$$

- ▶ Or combine the feature map  $f$  with the logistic regression model to obtain **cross-entropy** loss function.

- Application: image classification:

$$f(\boxed{?}; \Theta) = \begin{pmatrix} 0.7 \\ 0.2 \\ 0.1 \end{pmatrix} \implies \boxed{?} =$$

# How to formulate “learning”?

- Data:  $\{x_i, y_i\}_{i=1}^m$
- Find  $f^*$  in some function class such that  $f^*(x_i) \approx y_i$ .
- Mathematically: solve the **optimization problem** by parameterizing the abstract function class

$$\min_{\Theta} L(x, y, \Theta) \quad (1)$$

where

$$L(x, y, \Theta) = \mathbb{E}_{(x, y) \sim \mathcal{D}} [L(f(x; \Theta), y)] \approx \frac{1}{N} \sum_{i=1}^N \|f(x_i; \Theta) - y_i\|^2$$

- ▶ Or combine the feature map  $f$  with the logistic regression model to obtain **cross-entropy** loss function.

- Application: image classification:

$$f(\boxed{?}; \Theta) = \begin{pmatrix} 0.7 \\ 0.2 \\ 0.1 \end{pmatrix} \implies \boxed{?} = \text{cat}$$

# Deep Learning

*Machine learning using a special function class:  
deep neural networks!*

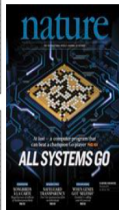
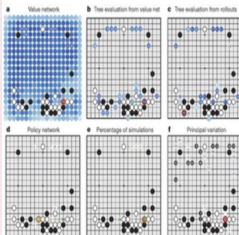
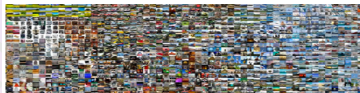
# Deep learning and its great successes

CNN has been successfully used in:

- Computer vision
  - ▶ Classification, detection, segmentation...
  - ▶ Medical image processing,
  - ▶ Face recognition,
- Reinforcement learning
  - ▶ AlphaGo,
  - ▶ Automated driving,
- Natural language processing
  - ▶ Speech recognition,
  - ▶ Machine translation,
- ...



IMAGENET



# Deep Concern

## Quotes (Researchers from U Wash, Princeton, MIT, . . .)

- “Deep learning is killing image processing, natural language processing...”
- “Graduate students only work on deep learning”
- “One time extinction event – graduate students won’t know the fundamental tools”



Source: D. Donoho/ H. Monajemi/ V. Papayan “Stats 38” at Stanford and B. Dong at PKU

# Mathematical understanding and Analysis?

Question:

*Why and how do these neural network machine learning models and relevant algorithms work?*



# Mathematical understanding and Analysis?

Question:

*Why and how do these neural network machine learning models and relevant algorithms work?*

This series of talks:

- 1 Finite Element and Deep Neural Network
  - ▶ ReLU neural networks = linear finite elements
  - ▶ Largest function class that a stable neural network can approximate
  - ▶ Optimal approximation rates for popular neural networks
- 2 Multigrid and Image Classification
  - ▶ Linear separable sets and logistic regression
  - ▶ A model for feature extractions
  - ▶ Image classification by multigrid method
- 3 Neural Network and Numerical PDEs
  - ▶ Error analysis of neural network for numerical PDEs
  - ▶ Numerical quadrature and Rademacher complexity analysis
  - ▶ Training algorithm that achieves the best asymptotic convergence rate

## 1 Introduction

## 2 Finite Element and Neural Networks

- Neural Network Functions
- Connection of ReLU DNN and linear FEM

## 3 Neural Network Approximation Class

- Density of Shallow Neural Network
- Approximation Properties
- Approximation Spaces

# Recall: supervised learning

- Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

# Recall: supervised learning

- Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

$$f(\text{cat}; \Theta) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

# Recall: supervised learning

- Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

$$f(\text{cat image}; \Theta) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad f(\text{dog image}; \Theta) \approx \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

# Recall: supervised learning

- Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

$$f(\text{cat}; \Theta) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad f(\text{dog}; \Theta) \approx \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad f(\text{rabbit}; \Theta) \approx \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

# Recall: supervised learning

- Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

$$f(\text{cat}; \Theta) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad f(\text{dog}; \Theta) \approx \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad f(\text{rabbit}; \Theta) \approx \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

- Image classification:

$$f(\text{?}; \Theta) = \begin{pmatrix} 0.7 \\ 0.2 \\ 0.1 \end{pmatrix}$$

# Recall: supervised learning

- Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

$$f(\text{cat image}; \Theta) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad f(\text{dog image}; \Theta) \approx \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad f(\text{rabbit image}; \Theta) \approx \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

- Image classification:

$$f(\text{?}; \Theta) = \begin{pmatrix} 0.7 \\ 0.2 \\ 0.1 \end{pmatrix} \Rightarrow \text{?} = \text{cat!}$$



# Recall: supervised learning

- Mathematical problem: Find  $f(\cdot; \Theta) : \mathbb{R}^d \rightarrow \mathbb{R}^3$  such that:

$$f(\text{cat image}; \Theta) \approx \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad f(\text{dog image}; \Theta) \approx \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad f(\text{rabbit image}; \Theta) \approx \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

- Image classification:

$$f(\boxed{?}; \Theta) = \begin{pmatrix} 0.7 \\ 0.2 \\ 0.1 \end{pmatrix} \Rightarrow \boxed{?} = \text{cat!}$$

Question:

*What is a good function class for  $f$ ?*

# What is a good function class for $f$ ?

Most favorable function class:

Jinchao Xu  
xu@multigrid.org

# What is a good function class for $f$ ?

Most favorable function class:

- Polynomials!

$$\sum_{\alpha_1 + \alpha_2 + \dots + \alpha_d \leq n} a_{\alpha} x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}$$

# What is a good function class for $f$ ?

Most favorable function class:

- Polynomials!

$$\sum_{\alpha_1 + \alpha_2 + \dots + \alpha_d \leq n} a_{\alpha} x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}$$

Important property:

# What is a good function class for $f$ ?

Most favorable function class:

- Polynomials!

$$\sum_{\alpha_1 + \alpha_2 + \dots + \alpha_d \leq n} a_{\alpha} x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}$$

**Important property:** polynomials can approximate any reasonable function!

# What is a good function class for $f$ ?

Most favorable function class:

- Polynomials!

$$\sum_{\alpha_1 + \alpha_2 + \dots + \alpha_d \leq n} a_{\alpha} x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}$$

**Important property:** polynomials can approximate any reasonable function!

- dense in  $C(\Omega)$  [Weierstrass theorem]

# What is a good function class for $f$ ?

Most favorable function class:

- Polynomials!

$$\sum_{\alpha_1 + \alpha_2 + \dots + \alpha_d \leq n} a_{\alpha} x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}$$

**Important property:** polynomials can approximate any reasonable function!

- dense in  $C(\Omega)$  [Weierstrass theorem]
- dense in all Sobolev spaces:  $L^2(\Omega)$ ,  $W^{m,p}(\Omega)$ ,  $\dots$

# What is a good function class for $f$ ?

Most favorable function class:

- Polynomials!

$$\sum_{\alpha_1 + \alpha_2 + \dots + \alpha_d \leq n} a_{\alpha} x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}$$

**Important property:** polynomials can approximate any reasonable function!

- dense in  $C(\Omega)$  [Weierstrass theorem]
- dense in all Sobolev spaces:  $L^2(\Omega)$ ,  $W^{m,p}(\Omega)$ ,  $\dots$

**Curse of dimensionality:** Number of coefficients for polynomials of degrees  $n$  in  $\mathbb{R}^d$ :

$$N = \binom{d+n}{n} = \frac{(n+d)!}{d!n!}.$$



# What is a good function class for $f$ ?

Most favorable function class:

- Polynomials!

$$\sum_{\alpha_1 + \alpha_2 + \dots + \alpha_d \leq n} a_{\alpha} x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}$$

**Important property:** polynomials can approximate any reasonable function!

- dense in  $C(\Omega)$  [Weierstrass theorem]
- dense in all Sobolev spaces:  $L^2(\Omega)$ ,  $W^{m,p}(\Omega)$ ,  $\dots$

**Curse of dimensionality:** Number of coefficients for polynomials of degrees  $n$  in  $\mathbb{R}^d$ :

$$N = \binom{d+n}{n} = \frac{(n+d)!}{d!n!}.$$

# What is a good function class for $f$ ?

Most favorable function class:

- Polynomials!

$$\sum_{\alpha_1 + \alpha_2 + \dots + \alpha_d \leq n} a_{\alpha} x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}$$

**Important property:** polynomials can approximate any reasonable function!

- dense in  $C(\Omega)$  [Weierstrass theorem]
- dense in all Sobolev spaces:  $L^2(\Omega)$ ,  $W^{m,p}(\Omega)$ ,  $\dots$

**Curse of dimensionality:** Number of coefficients for polynomials of degrees  $n$  in  $\mathbb{R}^d$ :

$$N = \binom{d+n}{n} = \frac{(n+d)!}{d!n!}.$$

For example  $n = 100$ :

$d =$	2	4	8
$N =$	$5 \times 10^3$	$4.6 \times 10^6$	$3.5 \times 10^{11}$

## 1 Introduction

## 2 Finite Element and Neural Networks

- Neural Network Functions
- Connection of ReLU DNN and linear FEM

## 3 Neural Network Approximation Class

- Density of Shallow Neural Network
- Approximation Properties
- Approximation Spaces

# $\sigma$ -DNN: Linears, activation and composition

xu@multigrid.org

# $\sigma$ -DNN: Linears, activation and composition

- 1 Start from a linear function

$$W^0 x + b^0$$

xu@multigrid.org

# $\sigma$ -DNN: Linears, activation and composition

- 1 Start from a linear function

$$W^0 x + b^0$$

- 2 Compose with the activation function:

$$x^{(1)} = \sigma(W^0 x + b^0)$$

# $\sigma$ -DNN: Linears, activation and composition

- 1 Start from a linear function

$$W^0 x + b^0$$

- 2 Compose with the activation function:

$$x^{(1)} = \sigma(W^0 x + b^0)$$

- 3 Compose with another linear function:

$$W^1 x^{(1)} + b^1$$

# $\sigma$ -DNN: Linears, activation and composition

- 1 Start from a linear function

$$W^0 x + b^0$$

- 2 Compose with the activation function:

$$x^{(1)} = \sigma(W^0 x + b^0)$$

- 3 Compose with another linear function:

$$W^1 x^{(1)} + b^1$$

- 4 Compose with the activation function:

$$x^{(2)} = \sigma(W^1 x^{(1)} + b^1)$$



# $\sigma$ -DNN: Linears, activation and composition

- 1 Start from a linear function

$$W^0 x + b^0$$

- 2 Compose with the activation function:

$$x^{(1)} = \sigma(W^0 x + b^0)$$

- 3 Compose with another linear function:

$$W^1 x^{(1)} + b^1$$

- 4 Compose with the activation function:

$$x^{(2)} = \sigma(W^1 x^{(1)} + b^1)$$

- 5 Compose with another linear function

$$f(x; \Theta) = W^2 x^{(2)} + b^2$$

# $\sigma$ -DNN: Linears, activation and composition

- 1 Start from a linear function

$$W^0 x + b^0$$

- 2 Compose with the activation function:

$$x^{(1)} = \sigma(W^0 x + b^0)$$

- 3 Compose with another linear function:

$$W^1 x^{(1)} + b^1$$

- 4 Compose with the activation function:

$$x^{(2)} = \sigma(W^1 x^{(1)} + b^1)$$

- 5 Compose with another linear function

$$f(x; \Theta) = W^2 x^{(2)} + b^2$$

- 6 ...

# $\sigma$ -DNN: Linears, activation and composition

- 1 Start from a linear function

$$W^0 x + b^0$$

- 2 Compose with the activation function:

$$x^{(1)} = \sigma(W^0 x + b^0)$$

- 3 Compose with another linear function:

$$W^1 x^{(1)} + b^1$$

- 4 Compose with the activation function:

$$x^{(2)} = \sigma(W^1 x^{(1)} + b^1)$$

- 5 Compose with another linear function

$$f(x; \Theta) = W^2 x^{(2)} + b^2$$

- 6 ...

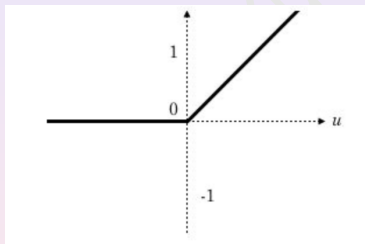
Deep neural network functions with  $\ell$ -hidden layers

$$\Sigma_{n_1:\ell}^{\sigma} = \{W^{\ell} x^{(\ell)} + b^{\ell}, W^i \in \mathbb{R}^{n_i}, b_i \in \mathbb{R}\}$$

# A popular activation function

Ramp or ReLU function

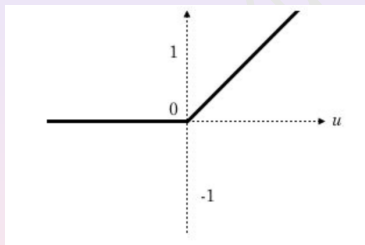
$$\sigma = \text{ReLU}(x) = \max(0, x)$$



# A popular activation function

Ramp or ReLU function

$$\sigma = \text{ReLU}(x) = \max(0, x)$$



Notation:

$$\Sigma_{n_{1:\ell}}^k = \Sigma_{n_{1:\ell}}^{\text{ReLU}^k}.$$

# What does a function in $\Sigma_{n_1:\ell}^1$ look like?

Obviously:

$\Sigma_{n_1:\ell}^1$  = a space of continuous piecewise linear functions!<sup>[1]</sup>

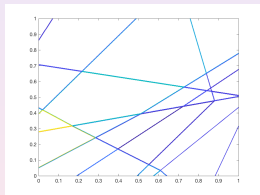
---

[1] [he2018relu](#).

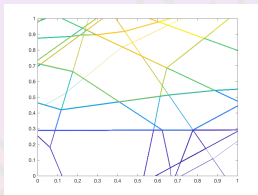
# What does a function in $\Sigma_{n_1:\ell}^1$ look like?

Obviously:

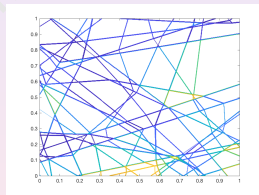
$\Sigma_{n_1:\ell}^1$  = a space of continuous piecewise linear functions!<sup>[1]</sup>



(10, 10)



(20, 20)



(40, 40)

[1] [he2018relu](#).

# How is $\Sigma_{n_{1:\ell}}^1$ compared with (adaptive) linear FEM?

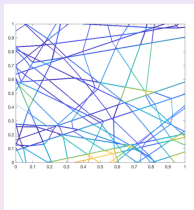


Figure: (40, 40)

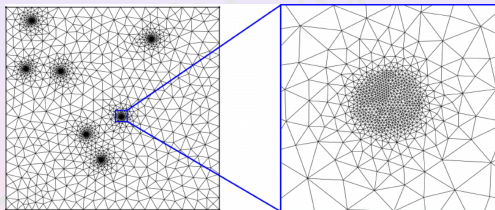


Figure: Adaptive Grid



## 1 Introduction

## 2 Finite Element and Neural Networks

- Neural Network Functions
- Connection of ReLU DNN and linear FEM

## 3 Neural Network Approximation Class


- Density of Shallow Neural Network
- Approximation Properties
- Approximation Spaces

# Finite element: piecewise linear functions

Jinchao Xu  
xu@multigrid.org


# Finite element: piecewise linear functions

- Uniform grid  $\mathcal{T}_h$

$$0 = x_0 < x_1 < \cdots < x_{N+1} = 1, \quad x_j = \frac{j}{N+1} \quad (j = 0 : N+1).$$


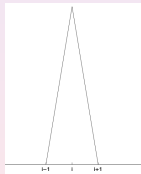
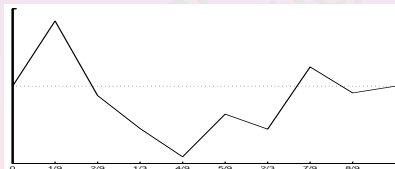
# Finite element: piecewise linear functions

- Uniform grid  $\mathcal{T}_h$

$$0 = x_0 < x_1 < \cdots < x_{N+1} = 1, \quad x_j = \frac{j}{N+1} \quad (j = 0 : N+1).$$


- Linear finite element space

$$V_h = \{v : v \text{ is continuous and piecewise linear w.r.t. } \mathcal{T}_h\}.$$



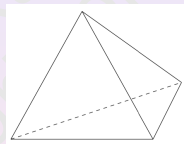
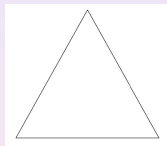
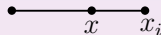
$$v_h(x) = \sum_{i=1}^n v_h(x_i) \phi_i(x).$$

# Linear finite element in multi-dimensions

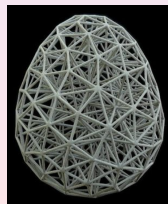
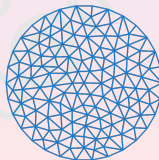
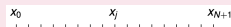
$$w_1 x + b$$

$$w_1 x_1 + w_2 x_2 + b$$

$$w_1 x_1 + w_2 x_2 + w_3 x_3 + b \quad \dots$$



...

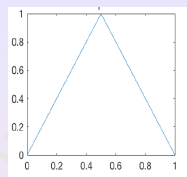


...

# FEM basis function in 1D

- Denote the basis function in  $\mathcal{T}_1$

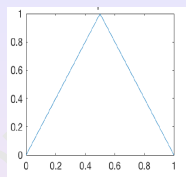
$$\varphi(x) = \begin{cases} 2x & x \in [0, \frac{1}{2}] \\ 2(1-x) & x \in [\frac{1}{2}, 1] \\ 0, & \text{others} \end{cases} \quad (2)$$



# FEM basis function in 1D

- Denote the basis function in  $\mathcal{T}_1$

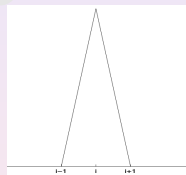
$$\varphi(x) = \begin{cases} 2x & x \in [0, \frac{1}{2}] \\ 2(1-x) & x \in [\frac{1}{2}, 1] \\ 0, & \text{others} \end{cases} \quad (2)$$



- All basis functions  $\varphi_i$  can be written as

$$\varphi_i = \varphi\left(\frac{x - x_{i-1}}{2h}\right) = \varphi(w_h x + b_i). \quad (3)$$

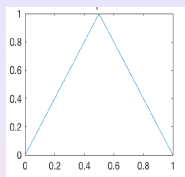
$$\text{with } w_h = \frac{1}{2h}, \quad b_i = \frac{-(i-1)}{2}.$$



# FEM basis function in 1D

- Denote the basis function in  $\mathcal{T}_1$

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



- All basis functions  $\varphi_i$  can be written as

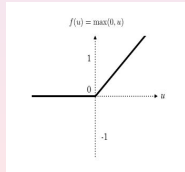
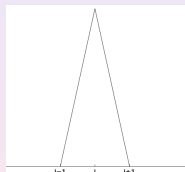
$$\varphi_i = \varphi\left(\frac{x - x_{i-1}}{2h}\right) = \varphi(w_h x + b_i). \quad (3)$$

with  $w_h = \frac{1}{2h}$ ,  $b_i = \frac{-(i-1)}{2}$ .

- Let  $x_+ = \max(0, x) = \text{ReLU}(x)$ ,

$$\varphi(x) = 2x_+ - 4(x - 1/2)_+ + 2(x - 1)_+.$$

- $\varphi_i \in \text{span} \{ (wx + b)_+, w, b \in \mathbb{R}^1 \}$





$$\text{FEM} \implies \Sigma_n^1 (d = 1)$$

- FEM  $\implies \Sigma_n^1$  ( make  $w_h$  and  $b_i$  arbitrary)

$$\text{FEM} \implies \Sigma_n^1 (d = 1)$$

- $\text{FEM} \implies \Sigma_n^1$  ( make  $w_h$  and  $b_i$  arbitrary)

$$\text{FEM} \subset \left\{ \sum_{i=1}^n a_i (\omega_i x + b_i)_+, a_i, \omega_i, b_i \in \mathbb{R}^1 \right\} = \Sigma_n^1.$$

$$\text{FEM} \implies \Sigma_n^1 (d = 1)$$

- $\text{FEM} \implies \Sigma_n^1$  ( make  $w_h$  and  $b_i$  arbitrary)

$$\text{FEM} \subset \left\{ \sum_{i=1}^n a_i (\omega_i x + b_i)_+, a_i, \omega_i, b_i \in \mathbb{R}^1 \right\} = \Sigma_n^1.$$

$$\text{FEM} \implies \Sigma_n^1 (d = 1)$$

- $\text{FEM} \implies \Sigma_n^1$  ( make  $w_h$  and  $b_i$  arbitrary)

$$\text{FEM} \subset \left\{ \sum_{i=1}^n a_i (\omega_i x + b_i)_+, a_i, \omega_i, b_i \in \mathbb{R}^1 \right\} = \Sigma_n^1.$$

$\Sigma_n^1$  is one hidden layer “shallow” neural network with activation function ReLU,  $n$  neurons.

# Generalization to multi-dimension:

Higher dimension  $d \geq 1$

$$\Sigma_n^1 = \left\{ \sum_{i=1}^n a_i (w_i x + b_i)_+ : w_i \in \mathbb{R}^{1 \times d}, b_i \in \mathbb{R} \right\} \quad (4)$$

where  $w_i x = \sum_{j=1}^d w_{ij} x_j$ .

# Connection of ReLU-DNN and Linear FEM

①  $d = 1,$

$$\text{FE} \subset \Sigma_n^1.$$

---

[2] [he2018relu](#).

[3] [arora2016understanding](#).

# Connection of ReLU-DNN and Linear FEM

1  $d = 1,$

$$\text{FE} \subset \Sigma_n^1.$$

2  $d \geq 2^{[2]},$

$$\text{FE} \not\subset \Sigma_n^1, \quad \forall n \geq 1.$$

---

[2] [he2018relu](#).

[3] [arora2016understanding](#).

# Connection of ReLU-DNN and Linear FEM

1  $d = 1,$

$$\text{FE} \subset \Sigma_n^1.$$

2  $d \geq 2^{[2]},$

$$\text{FE} \not\subset \Sigma_n^1, \quad \forall n \geq 1.$$

3  $d \geq 2^{[3]},$

$$\text{FE} \subset \Sigma_{n_1:\ell}^1 \quad \text{for some } \ell > 1,$$

where  $\Sigma_{n_1:\ell}^1$  is ReLU-DNN with  $\ell$  layers.

---

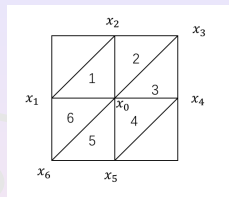
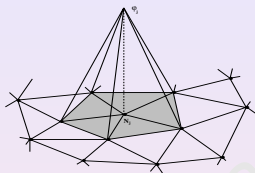
[2] [he2018relu](#).

[3] [arora2016understanding](#).



# A 2D example: FE basis function

Consider a 2D FE basis function,  $\phi(x)$ :

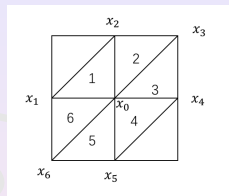
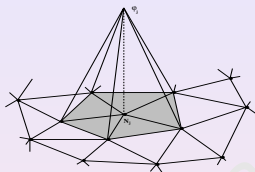


[4] [he2018relu](#).

[5] [arora2016understanding](#).

# A 2D example: FE basis function

Consider a 2D FE basis function,  $\phi(x)$ :



Here  $g_i$  is linear in Domain  $i$ , and  $x_7 = x_1$ , satisfying

$$g_i(x_0) = 1 \quad g_i(x_i) = 0 \quad g_i(x_{i+1}) = 0$$

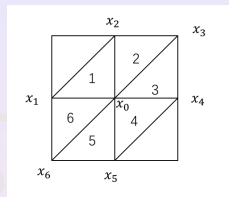
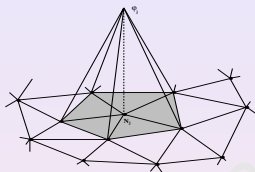
$$\phi(x) = \begin{cases} g_i(x), & x \in \text{Domain } i \\ 0, & x \in \mathbb{R}^2 - \overline{x_1 x_2 x_3 x_4 x_5 x_6} \end{cases} \quad (5)$$

[4] [he2018relu](#).

[5] [arora2016understanding](#).

# A 2D example: FE basis function

Consider a 2D FE basis function,  $\phi(x)$ :



Here  $g_i$  is linear in Domain  $i$ , and  $x_7 = x_1$ , satisfying

$$g_i(x_0) = 1 \quad g_i(x_i) = 0 \quad g_i(x_{i+1}) = 0$$

$$\phi(x) = \begin{cases} g_i(x), & x \in \text{Domain } i \\ 0, & x \in \mathbb{R}^2 - \overline{x_1 x_2 x_3 x_4 x_5 x_6} \end{cases} \quad (5)$$

We have<sup>[4][5]</sup>

$$\phi \in \Sigma_{n_{1:4}}^1. \quad (6)$$

[4] [he2018relu](#).

[5] [arora2016understanding](#).

# Properties of ReLU

We can write the basis function  $\phi(x)$  defined on hexagon  $\overline{x_1 x_2 x_3 x_4 x_5 x_6}$  which is satisfied that important identities:

# Properties of ReLU

We can write the basis function  $\phi(x)$  defined on hexagon  $\overline{x_1 x_2 x_3 x_4 x_5 x_6}$  which is satisfied that important identities:

$$x = \text{ReLU}(x) - \text{ReLU}(-x), |x| = \text{ReLU}(x) + \text{ReLU}(-x)$$

# Properties of ReLU

We can write the basis function  $\phi(x)$  defined on hexagon  $\overline{x_1 x_2 x_3 x_4 x_5 x_6}$  which is satisfied that important identities:

$$x = \text{ReLU}(x) - \text{ReLU}(-x), |x| = \text{ReLU}(x) + \text{ReLU}(-x)$$

and

$$\min(a, b) = \frac{a + b}{2} - \frac{|a - b|}{2} = v \cdot \text{ReLU}(W \cdot [a, b]^T) \quad (7)$$

# Properties of ReLU

We can write the basis function  $\phi(x)$  defined on hexagon  $\overline{x_1 x_2 x_3 x_4 x_5 x_6}$  which is satisfied that important identities:

$$x = \text{ReLU}(x) - \text{ReLU}(-x), |x| = \text{ReLU}(x) + \text{ReLU}(-x)$$

and

$$\min(a, b) = \frac{a+b}{2} - \frac{|a-b|}{2} = v \cdot \text{ReLU}(W \cdot [a, b]^T) \quad (7)$$

where

$$v = \frac{1}{2}[1, -1, -1, -1] \quad W = \begin{bmatrix} 1 & 1 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \end{bmatrix}$$

# Example of DNN and FEM: 2D-FEM basis function

$$\min(a, b, c) = \min(\min(a, b), c)$$

Jinchao Xu  
xu@multigrid.org



## Example of DNN and FEM: 2D-FEM basis function

$$\begin{aligned}\min(a, b, c) &= \min(\min(a, b), c) \\ &= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} \min(a, b) \\ c \end{pmatrix} \right)\end{aligned}$$

# Example of DNN and FEM: 2D-FEM basis function

$$\begin{aligned}\min(a, b, c) &= \min(\min(a, b), c) \\ &= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} \min(a, b) \\ c \end{pmatrix} \right) \\ &= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} v \cdot \text{ReLU}(W \cdot [a, b]^T) \\ \text{ReLU}(c) - \text{ReLU}(-c) \end{pmatrix} \right)\end{aligned}$$

## Example of DNN and FEM: 2D-FEM basis function

$$\begin{aligned}\min(a, b, c) &= \min(\min(a, b), c) \\&= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} \min(a, b) \\ c \end{pmatrix} \right) \\&= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} v \cdot \text{ReLU}(W \cdot [a, b]^T) \\ \text{ReLU}(c) - \text{ReLU}(-c) \end{pmatrix} \right) \\&= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} v \cdot \text{ReLU}(W \cdot [a, b]^T) \\ [1, -1] \cdot \text{ReLU}([1, -1]^T c) \end{pmatrix} \right)\end{aligned}$$

## Example of DNN and FEM: 2D-FEM basis function

$$\begin{aligned}\min(a, b, c) &= \min(\min(a, b), c) \\&= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} \min(a, b) \\ c \end{pmatrix} \right) \\&= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} v \cdot \text{ReLU}(W \cdot [a, b]^T) \\ \text{ReLU}(c) - \text{ReLU}(-c) \end{pmatrix} \right) \\&= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} v \cdot \text{ReLU}(W \cdot [a, b]^T) \\ [1, -1] \cdot \text{ReLU}([1, -1]^T c) \end{pmatrix} \right) \\&= v \cdot \text{ReLU} \left( W_2 \cdot \text{ReLU}(W_1 \cdot [a, b, c]^T) \right)\end{aligned}$$

# Example of DNN and FEM: 2D-FEM basis function

$$\begin{aligned}\min(a, b, c) &= \min(\min(a, b), c) \\&= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} \min(a, b) \\ c \end{pmatrix} \right) \\&= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} v \cdot \text{ReLU}(W \cdot [a, b]^T) \\ \text{ReLU}(c) - \text{ReLU}(-c) \end{pmatrix} \right) \\&= v \cdot \text{ReLU} \left( W \cdot \begin{pmatrix} v \cdot \text{ReLU}(W \cdot [a, b]^T) \\ [1, -1] \cdot \text{ReLU}([1, -1]^T c) \end{pmatrix} \right) \\&= v \cdot \text{ReLU} \left( W_2 \cdot \text{ReLU}(W_1 \cdot [a, b, c]^T) \right)\end{aligned}$$

where

$$W_2 = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & 1 & -1 \\ -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & -1 & 1 \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & -1 & 1 \\ -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & 1 & -1 \end{bmatrix}, \quad W_1 = \begin{bmatrix} 1 & 1 & 0 \\ -1 & -1 & 0 \\ 1 & -1 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & -1 \end{bmatrix}.$$

# Example of DNN and FEM: 2D-FEM basis function

It follows that

$$\begin{aligned} g &\equiv \min(g_1, g_2, g_3, g_4, g_5, g_6) = \min(\min(g_1, g_2, g_3), \min(g_4, g_5, g_6)) \\ &= v \cdot \text{ReLU}\left(W \cdot \begin{bmatrix} \min(g_1, g_2, g_3) \\ \min(g_4, g_5, g_6) \end{bmatrix}\right) \\ &= v \cdot \text{ReLU}\left(W \cdot \begin{bmatrix} v \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot [g_1, g_2, g_3]^T)) \\ v \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot [g_4, g_5, g_6]^T)) \end{bmatrix}\right) \end{aligned}$$

# Example of DNN and FEM: 2D-FEM basis function

It follows that

$$\begin{aligned} g &\equiv \min(g_1, g_2, g_3, g_4, g_5, g_6) = \min(\min(g_1, g_2, g_3), \min(g_4, g_5, g_6)) \\ &= v \cdot \text{ReLU}\left(W \cdot \begin{bmatrix} \min(g_1, g_2, g_3) \\ \min(g_4, g_5, g_6) \end{bmatrix}\right) \\ &= v \cdot \text{ReLU}\left(W \cdot \begin{bmatrix} v \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot [g_1, g_2, g_3]^T)) \\ v \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot [g_4, g_5, g_6]^T)) \end{bmatrix}\right) \end{aligned}$$

Thus

$$\phi = \text{ReLU}\left(v \cdot \text{ReLU}\left(W \cdot \begin{bmatrix} v \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot [g_1, g_2, g_3]^T)) \\ v \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot [g_4, g_5, g_6]^T)) \end{bmatrix}\right)\right)$$

# Example of DNN and FEM: 2D-FEM basis function

It follows that

$$\begin{aligned} g &\equiv \min(g_1, g_2, g_3, g_4, g_5, g_6) = \min(\min(g_1, g_2, g_3), \min(g_4, g_5, g_6)) \\ &= v \cdot \text{ReLU}\left(W \cdot \begin{bmatrix} \min(g_1, g_2, g_3) \\ \min(g_4, g_5, g_6) \end{bmatrix}\right) \\ &= v \cdot \text{ReLU}\left(W \cdot \begin{bmatrix} v \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot [g_1, g_2, g_3]^T)) \\ v \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot [g_4, g_5, g_6]^T)) \end{bmatrix}\right) \end{aligned}$$

Thus

$$\phi = \text{ReLU}\left(v \cdot \text{ReLU}\left(W \cdot \begin{bmatrix} v \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot [g_1, g_2, g_3]^T)) \\ v \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot [g_4, g_5, g_6]^T)) \end{bmatrix}\right)\right) \in \Sigma_{n_1:4}^1.$$



# ReLU-DNN and Linear FEM for $H^1$

$$\text{ReLU-DNN} = \Sigma_{n_1:\ell}^1$$

# ReLU-DNN and Linear FEM for $H^1$

$$\text{ReLU-DNN} = \Sigma_{n_{1:\ell}}^1 = \text{Linear FEM} \subset H^1(\Omega)$$

## 1 Introduction

## 2 Finite Element and Neural Networks

- Neural Network Functions
- Connection of ReLU DNN and linear FEM

## 3 Neural Network Approximation Class

- Density of Shallow Neural Network
- Approximation Properties
- Approximation Spaces

## 1 Introduction

## 2 Finite Element and Neural Networks

- Neural Network Functions
- Connection of ReLU DNN and linear FEM

## 3 Neural Network Approximation Class

- **Density of Shallow Neural Network**
- Approximation Properties
- Approximation Spaces

# Shallow Neural Networks

- Superpositions of *ridge functions*

$$\Sigma_n^\sigma := \left\{ \sum_{i=1}^n a_i \sigma(\omega_i \cdot x + b_i), a_i \in \mathbb{R}, \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R} \right\}$$

# Shallow Neural Networks

- Superpositions of *ridge functions*

$$\Sigma_n^\sigma := \left\{ \sum_{i=1}^n a_i \sigma(\omega_i \cdot x + b_i), a_i \in \mathbb{R}, \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R} \right\}$$

- Shallow networks means one hidden layer, i.e.

$$\Sigma_{[1]}^\sigma = \bigcup_{n=1}^{\infty} \Sigma_n^\sigma \quad (8)$$

# Shallow Neural Networks

- Superpositions of *ridge functions*

$$\Sigma_n^\sigma := \left\{ \sum_{i=1}^n a_i \sigma(\omega_i \cdot x + b_i), a_i \in \mathbb{R}, \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R} \right\}$$

- Shallow networks means one hidden layer, i.e.

$$\Sigma_{[1]}^\sigma = \bigcup_{n=1}^{\infty} \Sigma_n^\sigma \quad (8)$$

- Let  $\Omega \subset \mathbb{R}^d$  be a bounded domain

# Shallow Neural Networks

- Superpositions of *ridge functions*

$$\Sigma_n^\sigma := \left\{ \sum_{i=1}^n a_i \sigma(\omega_i \cdot x + b_i), a_i \in \mathbb{R}, \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R} \right\}$$

- Shallow networks means one hidden layer, i.e.

$$\Sigma_{[1]}^\sigma = \bigcup_{n=1}^{\infty} \Sigma_n^\sigma \quad (8)$$

- Let  $\Omega \subset \mathbb{R}^d$  be a bounded domain
- Want to use shallow networks to approximate functions  $f$  on  $\Omega$



# Shallow Neural Networks

- Superpositions of *ridge functions*

$$\Sigma_n^\sigma := \left\{ \sum_{i=1}^n a_i \sigma(\omega_i \cdot x + b_i), a_i \in \mathbb{R}, \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R} \right\}$$

- Shallow networks means one hidden layer, i.e.

$$\Sigma_{[1]}^\sigma = \bigcup_{n=1}^{\infty} \Sigma_n^\sigma \quad (8)$$

- Let  $\Omega \subset \mathbb{R}^d$  be a bounded domain
- Want to use shallow networks to approximate functions  $f$  on  $\Omega$
- Consider two problems:
  - ▶ Density of  $\Sigma_{[1]}^\sigma$  in  $L^2(\Omega)$
  - ▶ Approximation rates in  $L^2(\Omega)$

# Are Shallow Networks Dense?

Can shallow networks approximate arbitrary functions?

- Is  $\Sigma_{[1]}^\sigma$  dense in  $L^2(\Omega)$ ?

---

[6] [hornik1989multilayer](#), [hornik1991approximation](#).

# Are Shallow Networks Dense?

Can shallow networks approximate arbitrary functions?

- Is  $\Sigma_{[1]}^\sigma$  dense in  $L^2(\Omega)$ ?
- It depends upon  $\sigma$
- If  $\sigma$  is a polynomial,

---

[6] [hornik1989multilayer](#), [hornik1991approximation](#).

# Are Shallow Networks Dense?

Can shallow networks approximate arbitrary functions?

- Is  $\Sigma_{[1]}^\sigma$  dense in  $L^2(\Omega)$ ?
- It depends upon  $\sigma$
- If  $\sigma$  is a polynomial, No!

---

[6] hornik1989multilayer, hornik1991approximation.

# Are Shallow Networks Dense?

Can shallow networks approximate arbitrary functions?

- Is  $\Sigma_{[1]}^\sigma$  dense in  $L^2(\Omega)$ ?
- It depends upon  $\sigma$
- If  $\sigma$  is a polynomial, No!
- Yes! As long as  $\sigma \in C^k$  is not a polynomial<sup>[6]</sup>.

---

[6] hornik1989multilayer, hornik1991approximation.

# Non-polynomial activation function $\Leftrightarrow$ approximation

## Lemma

$\Sigma_{[1]}^\sigma$  is dense in  $C^0(\Omega) \iff \sigma$  is not a polynomial.

# Non-polynomial activation function $\Leftrightarrow$ approximation

## Lemma

$\Sigma_{[1]}^\sigma$  is dense in  $C^0(\Omega) \iff \sigma$  is not a polynomial.

## Proof.

- $[\sigma((\omega + he_j) \cdot x + b) - \sigma(\omega \cdot x + b)]/h \in \Sigma_{[1]}^\sigma,$

# Non-polynomial activation function $\Leftrightarrow$ approximation

## Lemma

$\Sigma_{[1]}^\sigma$  is dense in  $C^0(\Omega) \iff \sigma$  is not a polynomial.

## Proof.

- $[\sigma((\omega + h e_j) \cdot x + b) - \sigma(\omega \cdot x + b)]/h \in \Sigma_{[1]}^\sigma,$
- $\frac{\partial}{\partial \omega_j} \sigma(\omega \cdot x + b)|_{\omega=0} = x_j \sigma'(b) \in \overline{\Sigma_{[1]}^\sigma}$



# Non-polynomial activation function $\Leftrightarrow$ approximation

## Lemma

$\Sigma_{[1]}^\sigma$  is dense in  $C^0(\Omega) \iff \sigma$  is not a polynomial.

## Proof.

- $[\sigma((\omega + he_j) \cdot x + b) - \sigma(\omega \cdot x + b)]/h \in \Sigma_{[1]}^\sigma,$
- $\frac{\partial}{\partial \omega_j} \sigma(\omega \cdot x + b)|_{\omega=0} = x_j \sigma'(b) \in \overline{\Sigma_{[1]}^\sigma}$   
 $\Rightarrow x_j \in \overline{\Sigma_{[1]}^\sigma}$  if  $\sigma'(b) \neq 0$  for some  $b$ .

# Non-polynomial activation function $\Leftrightarrow$ approximation

## Lemma

$\Sigma_{[1]}^\sigma$  is dense in  $C^0(\Omega) \iff \sigma$  is not a polynomial.

## Proof.

- $[\sigma((\omega + h e_j) \cdot x + b) - \sigma(\omega \cdot x + b)]/h \in \Sigma_{[1]}^\sigma$ ,
- $\frac{\partial}{\partial \omega_j} \sigma(\omega \cdot x + b)|_{\omega=0} = x_j \sigma'(b) \in \overline{\Sigma_{[1]}^\sigma}$   
 $\Rightarrow x_j \in \overline{\Sigma_{[1]}^\sigma}$  if  $\sigma'(b) \neq 0$  for some  $b$ .
- Similarly  $x^\alpha = x_1^{\alpha_1} \cdots x_d^{\alpha_d} \in \overline{\Sigma_{[1]}^\sigma}$  if  $\sigma^{(|\alpha|)}(b) \neq 0$  for some  $b$ .

# Non-polynomial activation function $\Leftrightarrow$ approximation

## Lemma

$\Sigma_{[1]}^\sigma$  is dense in  $C^0(\Omega) \iff \sigma$  is not a polynomial.

## Proof.

- $[\sigma((\omega + he_j) \cdot x + b) - \sigma(\omega \cdot x + b)]/h \in \Sigma_{[1]}^\sigma$ ,
- $\frac{\partial}{\partial \omega_j} \sigma(\omega \cdot x + b)|_{\omega=0} = x_j \sigma'(b) \in \overline{\Sigma_{[1]}^\sigma}$   
 $\Rightarrow x_j \in \overline{\Sigma_{[1]}^\sigma}$  if  $\sigma'(b) \neq 0$  for some  $b$ .
- Similarly  $x^\alpha = x_1^{\alpha_1} \cdots x_d^{\alpha_d} \in \overline{\Sigma_{[1]}^\sigma}$  if  $\sigma^{(|\alpha|)}(b) \neq 0$  for some  $b$ .
- $\overline{\Sigma_{[1]}^\sigma}$  contains all polynomials.



# Rich function class generated by activation function

We note that, for any  $w_i, b_i$

$\left\{ w_i x + b_i \right\}_{i=1}^m$  are linearly **dependent** if  $m > d + 1$ !

# Rich function class generated by activation function

We note that, for any  $w_i, b_i$

$\left\{ w_i x + b_i \right\}_{i=1}^m$  are linearly **dependent** if  $m > d + 1$ !

Question:

Given  $w_i$  and  $b_i$ , when are  $\left\{ \sigma(w_i x + b_i) \right\}_{i=1}^n$  linearly independent?

# Rich function class generated by activation function

We note that, for any  $w_i, b_i$

$\left\{ w_i x + b_i \right\}_{i=1}^m$  are linearly **dependent** if  $m > d + 1$ !

Question:

Given  $w_i$  and  $b_i$ , when are  $\left\{ \sigma(w_i x + b_i) \right\}_{i=1}^n$  linearly independent?

But

**Theorem (He, Lin, Xu and Zheng 2018, Siegel and Xu 2019)**

*Assume that  $\sigma$  is NOT a polynomial, then*

$\left\{ \sigma(w_i x + b_i) \right\}_{i=1}^m$  are linearly **independent** for any  $m \geq 1$

*if*

# Rich function class generated by activation function

We note that, for any  $w_i, b_i$

$\left\{ w_i x + b_i \right\}_{i=1}^m$  are linearly **dependent** if  $m > d + 1$ !

Question:

Given  $w_i$  and  $b_i$ , when are  $\left\{ \sigma(w_i x + b_i) \right\}_{i=1}^n$  linearly independent?

But

**Theorem (He, Lin, Xu and Zheng 2018, Siegel and Xu 2019)**

*Assume that  $\sigma$  is NOT a polynomial, then*

$\left\{ \sigma(w_i x + b_i) \right\}_{i=1}^m$  are linearly **independent** for any  $m \geq 1$

*if  $\{w_i\}$  are pairwise linearly independent.*

# A keyword in deep learning

Jinchao Xu  
xu@multigrid.org



# A keyword in deep learning

- nonlinearity

# A keyword in deep learning

- nonlinearity  $\Rightarrow$  Non-polynomial

# Implications

- Shallow networks are universal function approximators!

# Implications

- Shallow networks are universal function approximators!
- What about deeper networks?

# Implications

- Shallow networks are universal function approximators!
- What about deeper networks?
  - ▶ Each layer can approximate the identity to arbitrary accuracy

# Implications

- Shallow networks are universal function approximators!
- What about deeper networks?
  - ▶ Each layer can approximate the identity to arbitrary accuracy
  - ▶ So single hidden layer result extends to deeper networks as well

# Implications

- Shallow networks are universal function approximators!
- What about deeper networks?
  - ▶ Each layer can approximate the identity to arbitrary accuracy
  - ▶ So single hidden layer result extends to deeper networks as well
  - ▶ One way of writing this is

$$\Sigma_{[1:k]}^{\sigma} \subset \overline{\Sigma_{[1:k+1]}^{\sigma}} \quad (9)$$

if  $\sigma$  is not a polynomial!

- ▶ Note that we don't always have  $\Sigma_{[1:k]}^{\sigma} \subset \Sigma_{[1:k+1]}^{\sigma}$

# Implications

- Shallow networks are universal function approximators!
- What about deeper networks?
  - ▶ Each layer can approximate the identity to arbitrary accuracy
  - ▶ So single hidden layer result extends to deeper networks as well
  - ▶ One way of writing this is

$$\Sigma_{[1:k]}^{\sigma} \subset \overline{\Sigma_{[1:k+1]}^{\sigma}} \quad (9)$$

if  $\sigma$  is not a polynomial!

- ▶ Note that we don't always have  $\Sigma_{[1:k]}^{\sigma} \subset \Sigma_{[1:k+1]}^{\sigma}$
- How efficiently can functions be approximated?
  - ▶ No control of the size of the network or parameters
  - ▶ Want approximation rates for networks with controlled weights



## 1 Introduction

## 2 Finite Element and Neural Networks

- Neural Network Functions
- Connection of ReLU DNN and linear FEM

## 3 Neural Network Approximation Class

- Density of Shallow Neural Network
- **Approximation Properties**
- Approximation Spaces

# Sampling Argument

xu@multigrid.org

# Sampling Argument

Define

$$\mathbb{E}g := \int_G g(\omega, x) \lambda(\omega) d\omega. \quad (10)$$

# Sampling Argument

Define

$$\mathbb{E}g := \int_G g(\omega, x) \lambda(\omega) d\omega. \quad (10)$$

$$\mathbb{E}_n g := \int_{G^n} g(\omega_1, \dots, \omega_n) \lambda(\omega_1) \lambda(\omega_2) \dots \lambda(\omega_n) d\omega_1 d\omega_2 \dots d\omega_n. \quad (11)$$

# Sampling Argument

Define

$$\mathbb{E}g := \int_G g(\omega, x) \lambda(\omega) d\omega. \quad (10)$$

$$\mathbb{E}_n g := \int_{G^n} g(\omega_1, \dots, \omega_n) \lambda(\omega_1) \lambda(\omega_2) \dots \lambda(\omega_n) d\omega_1 d\omega_2 \dots d\omega_n. \quad (11)$$

## Lemma

For any  $g \in L^\infty(G)$ , we have

$$\mathbb{E}_n \left( \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i) \right)^2 = \frac{1}{n} \left( \mathbb{E}(g^2) - (\mathbb{E}(g))^2 \right) \leq \frac{1}{n} \mathbb{E}(g^2). \quad (12)$$

# Sampling Argument

1 Let  $u = \int_G g(\omega, x) \lambda(x) dx$  and the sampling  $u_n = \frac{1}{n} \sum_{i=1}^n g(\omega_i, x)$

$$u(x) - u_n(x) = \mathbb{E}g(x) - \frac{1}{n} \sum_{i=1}^n g(\omega_i, x); \quad (13)$$

# Sampling Argument

1 Let  $u = \int_G g(\omega, x) \lambda(x) dx$  and the sampling  $u_n = \frac{1}{n} \sum_{i=1}^n g(\omega_i, x)$

$$u(x) - u_n(x) = \mathbb{E}g(x) - \frac{1}{n} \sum_{i=1}^n g(\omega_i, x); \quad (13)$$

2 It holds that

$$\mathbb{E}_n \left( \left\| \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i) \right\|^2 \right) \leq \frac{1}{n} \mathbb{E}(\|g\|^2) \leq \frac{1}{n}; \quad (14)$$

# Sampling Argument

1 Let  $u = \int_G g(\omega, x) \lambda(x) dx$  and the sampling  $u_n = \frac{1}{n} \sum_{i=1}^n g(\omega_i, x)$

$$u(x) - u_n(x) = \mathbb{E}g(x) - \frac{1}{n} \sum_{i=1}^n g(\omega_i, x); \quad (13)$$

2 It holds that

$$\mathbb{E}_n \left( \left\| \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i) \right\|^2 \right) \leq \frac{1}{n} \mathbb{E}(\|g\|^2) \leq \frac{1}{n}; \quad (14)$$

3 There exist  $\{\omega_i^*\}_{i=1}^n$  such that  $u_n = \frac{1}{n} \sum_{i=1}^n g(\omega_i^*, x)$  satisfies

$$\|u - u_n\|_0 \leq n^{-\frac{1}{2}} \quad \text{or} \quad \left\| \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i, x) \right\|_0 \leq n^{-\frac{1}{2}}. \quad (15)$$



# Approximation Rates for Cosine Networks

# Approximation Rates for Cosine Networks

- Cosine networks

$$\Sigma_n^{\cos} = \left\{ u_n : u_n = \sum_{i=1}^n a_i \cos(w_i x + b_i), \quad \forall a_i, w_i, b_i \right\}$$

- Integral representation of  $u$  in terms of cosine functions

# Approximation Rates for Cosine Networks

- Cosine networks

$$\Sigma_n^{\cos} = \left\{ u_n : u_n = \sum_{i=1}^n a_i \cos(w_i x + b_i), \quad \forall a_i, w_i, b_i \right\}$$

- Integral representation of  $u$  in terms of cosine functions

$$u(x) = \operatorname{Re} \int_{\mathbb{R}^d} e^{2\pi i \omega \cdot x} \hat{u}(\omega) d\omega$$

# Approximation Rates for Cosine Networks

- Cosine networks

$$\Sigma_n^{\cos} = \left\{ u_n : u_n = \sum_{i=1}^n a_i \cos(w_i x + b_i), \quad \forall a_i, w_i, b_i \right\}$$

- Integral representation of  $u$  in terms of cosine functions

$$\begin{aligned} u(x) &= \operatorname{Re} \int_{\mathbb{R}^d} e^{2\pi i \omega \cdot x} \hat{u}(\omega) d\omega \\ &= \int_{\mathbb{R}^d} \cos(2\pi \omega \cdot x + b(\omega)) |\hat{u}(\omega)| d\omega \quad (\text{by } \hat{u}(\omega) = |\hat{u}(\omega)| e^{ib(\omega)}) \end{aligned}$$

# Approximation Rates for Cosine Networks

- Cosine networks

$$\Sigma_n^{\cos} = \left\{ u_n : u_n = \sum_{i=1}^n a_i \cos(w_i x + b_i), \quad \forall a_i, w_i, b_i \right\}$$

- Integral representation of  $u$  in terms of cosine functions

$$\begin{aligned} u(x) &= \operatorname{Re} \int_{\mathbb{R}^d} e^{2\pi i \omega \cdot x} \hat{u}(\omega) d\omega \\ &= \int_{\mathbb{R}^d} \cos(2\pi \omega \cdot x + b(\omega)) |\hat{u}(\omega)| d\omega \quad (\text{by } \hat{u}(\omega) = |\hat{u}(\omega)| e^{ib(\omega)}) \\ &= \|\hat{u}\|_{L^1} \int_{\mathbb{R}^d} \cos(2\pi \omega \cdot x + b(\omega)) \lambda(\omega) d\omega \quad (\lambda(\omega) = \frac{|\hat{u}(\omega)|}{\|\hat{u}\|_{L^1}}) \end{aligned}$$

# Approximation Rates for Cosine Networks

- Cosine networks

$$\Sigma_n^{\cos} = \left\{ u_n : u_n = \sum_{i=1}^n a_i \cos(w_i x + b_i), \quad \forall a_i, w_i, b_i \right\}$$

- Integral representation of  $u$  in terms of cosine functions

$$\begin{aligned} u(x) &= \operatorname{Re} \int_{\mathbb{R}^d} e^{2\pi i \omega \cdot x} \hat{u}(\omega) d\omega \\ &= \int_{\mathbb{R}^d} \cos(2\pi \omega \cdot x + b(\omega)) |\hat{u}(\omega)| d\omega \quad (\text{by } \hat{u}(\omega) = |\hat{u}(\omega)| e^{ib(\omega)}) \\ &= \|\hat{u}\|_{L^1} \int_{\mathbb{R}^d} \cos(2\pi \omega \cdot x + b(\omega)) \lambda(\omega) d\omega \quad (\lambda(\omega) = \frac{|\hat{u}(\omega)|}{\|\hat{u}\|_{L^1}}) \\ &= \|\hat{u}\|_{L^1} \mathbb{E}(g(\omega, x)) \quad (g(\omega, x) = \cos(2\pi \omega \cdot x + b(\omega))) \end{aligned}$$

# Approximation Rates for Cosine Networks

- Cosine networks

$$\Sigma_n^{\cos} = \left\{ u_n : u_n = \sum_{i=1}^n a_i \cos(w_i x + b_i), \quad \forall a_i, w_i, b_i \right\}$$

- Integral representation of  $u$  in terms of cosine functions

$$\begin{aligned} u(x) &= \operatorname{Re} \int_{\mathbb{R}^d} e^{2\pi i \omega \cdot x} \hat{u}(\omega) d\omega \\ &= \int_{\mathbb{R}^d} \cos(2\pi \omega \cdot x + b(\omega)) |\hat{u}(\omega)| d\omega \quad (\text{by } \hat{u}(\omega) = |\hat{u}(\omega)| e^{ib(\omega)}) \\ &= \|\hat{u}\|_{L^1} \int_{\mathbb{R}^d} \cos(2\pi \omega \cdot x + b(\omega)) \lambda(\omega) d\omega \quad (\lambda(\omega) = \frac{|\hat{u}(\omega)|}{\|\hat{u}\|_{L^1}}) \\ &= \|\hat{u}\|_{L^1} \mathbb{E}(g(\omega, x)) \quad (g(\omega, x) = \cos(2\pi \omega \cdot x + b(\omega))) \end{aligned}$$

# Approximation Rates for Cosine Networks<sup>[7]</sup>

xu@multigrid.org

---

<sup>[7]</sup> jones1992simple.



# Approximation Rates for Cosine Networks<sup>[7]</sup>

1 The preceding *sampling argument* gives the *approximation rate*:

## Theorem

There exists  $u_n \in \Sigma_{n,M}^{\cos} = \left\{ \sum_{i=1}^n a_i \cos(\omega_i \cdot x + b_i), \sum_{i=1}^n |a_i| \leq M \right\}$  such that

$$\|u - u_n\| \lesssim n^{-\frac{1}{2}} \|\hat{u}\|_{L^1(\mathbb{R}^d)}.$$

---

[7] jones1992simple.

# Spectral Barron Norm

Generalization:

$$\inf_{u_n \in \Sigma_{n,M}^{\text{cos}}} \|u - u_n\|_{H^m(\Omega)} \lesssim n^{-\frac{1}{2}} \|u\|_{B^m(\Omega)}$$

where  $B^m(\Omega)$  is the spectral Barron space.

$$\|u\|_{B^m(\Omega)} = \inf_{u_e|_{\Omega}=u} \int_{\mathbb{R}^d} (1 + |\omega|)^m |\hat{u}_e(\omega)| d\omega$$

# Approximation Rates for the $\text{ReLU}^k$ Network

# Approximation Rates for the $\text{ReLU}^k$ Network

Note that

$$u(x) - \sum_{|\alpha| \leq k} \frac{1}{\alpha!} D^\alpha u(0) x^\alpha = \|\rho\|_{L^1(G)} \int_G g(x, \theta) \lambda(\theta) d\theta = \|\rho\|_{L^1(G)} \mathbb{E}(g)$$

with  $\theta = (z, t, \omega) \in G = \{-1, 1\} \times [0, T] \times \mathbb{R}^d$ ,

$$\begin{aligned} \rho(\theta) &= \frac{1}{(2\pi)^d} |s(zt, \omega)| |\hat{u}(\omega)| \|\omega\|^{k+1} \\ s(zt, \omega) &= \begin{cases} (-1)^{\frac{k+1}{2}} \cos(z\|\omega\|t + b(\omega)) & k \text{ is odd,} \\ (-1)^{\frac{k+2}{2}} \sin(z\|\omega\|t + b(\omega)) & k \text{ is even.} \end{cases} \end{aligned} \quad (16)$$

$$g(x, \theta) = (z\bar{\omega} \cdot x - t)_+^k \text{sgn} s(zt, \omega), \quad \lambda(\theta) = \frac{\rho(\theta)}{\|\rho\|_{L^1(G)}}. \quad (17)$$

# Approximation Rates for the ReLU<sup>k</sup>[8][9]

## Lemma (Sampling Analysis)

There exist  $\omega_i, t_i$  such that

$$\|u - u_n\|_{H^m} \lesssim n^{-\frac{1}{2}} \|u\|_{B^{m+k+1}} \quad (18)$$

$$\text{with } u_n = \sum_{|\alpha| \leq k} \frac{1}{\alpha!} D^\alpha u(0) x^\alpha + \frac{c}{k!n} \sum_{i=1}^n \beta_i (\omega_i \cdot x - t_i)_+^k.$$

---

[8] [klusowski2016risk](#).

[9] [xu2020finite](#).

# Approximation Rates for the ReLU<sup>k</sup>[8][9]

## Lemma (Sampling Analysis)

There exist  $\omega_i, t_i$  such that

$$\|u - u_n\|_{H^m} \lesssim n^{-\frac{1}{2}} \|u\|_{B^{m+k+1}} \quad (18)$$

$$\text{with } u_n = \sum_{|\alpha| \leq k} \frac{1}{\alpha!} D^\alpha u(0) x^\alpha + \frac{c}{k!n} \sum_{i=1}^n \beta_i (\omega_i \cdot x - t_i)_+^k.$$

## Lemma (More Refined Analysis: Stratified Method)

There exist  $\omega_i, t_i$  such that

$$\|u - u_n\|_{H^m} \lesssim n^{-\frac{1}{2} - \frac{1}{d}} \|u\|_{B^{m+k+1}} \quad (19)$$

$$\text{with } \bar{\omega}_i = \frac{\omega_i}{\|\omega_i\|} \text{ and } u_n = \sum_{|\alpha| \leq k} \frac{1}{\alpha!} D^\alpha u(0) x^\alpha + \frac{c}{k!n} \sum_{i=1}^n \beta_i (\bar{\omega}_i \cdot x - t_i)_+^k.$$

[8] klusowski2016risk.

[9] xu2020finite.

# More General Activation Functions

- Using similar techniques, we can extend these rates to more general activation functions as well<sup>[10]</sup>

## Theorem

Suppose that  $\sigma \in L^\infty$  and satisfies the decay condition

$$|\sigma(t)| \lesssim (1 + |t|)^p \quad (20)$$

for some  $p > 1$ . Let  $u \in L^2(\Omega)$ . Then we have

$$\inf_{u_n \in \Sigma_n^\sigma} \|u - u_n\|_{L^2(\Omega)} \lesssim \|u\|_{B^1(\Omega)} n^{-\frac{1}{2}}. \quad (21)$$

---

[10] siegel2020approximation.

# More General Activation Functions (cont.)

- At the cost of a worse rate, we can even drop almost all assumptions on<sup>[11]</sup>  $\sigma$

## Theorem

Suppose that  $\sigma \in L^\infty$  and  $\hat{\sigma}$  is continuous and non-zero at a single point. Let  $u \in L^2(\Omega)$ . Then we have

$$\inf_{u_n \in \Sigma_n^\sigma} \|u - u_n\|_{L^2(\Omega)} \lesssim \|u\|_{B^1(\Omega)} n^{-\frac{1}{4}}. \quad (22)$$

- In particular:
  - ▶ Holds for all  $\sigma \in L^1 \cap L^\infty$
  - ▶ Holds for all  $\sigma \in BV$

---

[11] siegel2020approximation.



## 1 Introduction

## 2 Finite Element and Neural Networks

- Neural Network Functions
- Connection of ReLU DNN and linear FEM

## 3 Neural Network Approximation Class

- Density of Shallow Neural Network
- Approximation Properties
- Approximation Spaces

# Stable Neural Network Approximation

- Recall: approximation of cosine neural network from the function class (by sampling argument):

$$\Sigma_{n,M}^{\cos} = \left\{ \sum_{i=1}^n a_i \cos(\omega_i \cdot x + b_i), \sum_{i=1}^n |a_i| \leq M, \quad M = \|\hat{u}\|_{L^1} \right\}$$

# Stable Neural Network Approximation

- Recall: approximation of cosine neural network from the function class (by sampling argument):

$$\Sigma_{n,M}^{\cos} = \left\{ \sum_{i=1}^n a_i \cos(\omega_i \cdot x + b_i), \sum_{i=1}^n |a_i| \leq M, \quad M = \|\hat{u}\|_{L^1} \right\}$$

- How far can the sampling argument go?

# Stable Neural Network Approximation

- Recall: approximation of cosine neural network from the function class (by sampling argument):

$$\Sigma_{n,M}^{\cos} = \left\{ \sum_{i=1}^n a_i \cos(\omega_i \cdot x + b_i), \sum_{i=1}^n |a_i| \leq M, \quad M = \|\hat{u}\|_{L^1} \right\}$$

- How far can the sampling argument go?
- Consider approximation from the class

$$\Sigma_{n,M}^{\sigma} := \left\{ \sum_{i=1}^n a_i \sigma(\omega_i \cdot x + b_i), \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R}, \sum_{i=1}^n |a_i| \leq M \right\} \quad (23)$$

of neural networks with  $\ell^1$ -bounded outer coefficients.

# Stable Neural Network Approximation

- Recall: approximation of cosine neural network from the function class (by sampling argument):

$$\Sigma_{n,M}^{\cos} = \left\{ \sum_{i=1}^n a_i \cos(\omega_i \cdot x + b_i), \sum_{i=1}^n |a_i| \leq M, \quad M = \|\hat{u}\|_{L^1} \right\}$$

- How far can the sampling argument go?
- Consider approximation from the class

$$\Sigma_{n,M}^{\sigma} := \left\{ \sum_{i=1}^n a_i \sigma(\omega_i \cdot x + b_i), \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R}, \sum_{i=1}^n |a_i| \leq M \right\} \quad (23)$$

of neural networks with  $\ell^1$ -bounded outer coefficients.

- More generally for a dictionary  $\mathbb{D} \subset H = L^2(\Omega)$ , consider

$$\Sigma_{n,M}(\mathbb{D}) = \left\{ \sum_{i=1}^n a_i h_i, h_i \in \mathbb{D}, \sum_{i=1}^n |a_i| \leq M \right\} \quad (24)$$

# Stable Neural Network Approximation

- Recall: approximation of cosine neural network from the function class (by sampling argument):

$$\Sigma_{n,M}^{\cos} = \left\{ \sum_{i=1}^n a_i \cos(\omega_i \cdot x + b_i), \sum_{i=1}^n |a_i| \leq M, \quad M = \|\hat{u}\|_{L^1} \right\}$$

- How far can the sampling argument go?
- Consider approximation from the class

$$\Sigma_{n,M}^{\sigma} := \left\{ \sum_{i=1}^n a_i \sigma(\omega_i \cdot x + b_i), \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R}, \sum_{i=1}^n |a_i| \leq M \right\} \quad (23)$$

of neural networks with  $\ell^1$ -bounded outer coefficients.

- More generally for a dictionary  $\mathbb{D} \subset H = L^2(\Omega)$ , consider

$$\Sigma_{n,M}(\mathbb{D}) = \left\{ \sum_{i=1}^n a_i h_i, h_i \in \mathbb{D}, \sum_{i=1}^n |a_i| \leq M \right\} \quad (24)$$

- Let  $M < \infty$  be fixed and consider approximation as  $n \rightarrow \infty$ .

# Stable Dictionary Approximation Space

Siegel & Xu, 2021<sup>[12]</sup>:

- Define a closed convex hull of  $\pm\mathbb{D}$ :

$$B_1(\mathbb{D}) = \overline{\bigcup_{n=1}^{\infty} \Sigma_{n,1}^{\sigma}}, \quad (25)$$

---

[12] [siegel2021optimal](#).

# Stable Dictionary Approximation Space

Siegel & Xu, 2021<sup>[12]</sup>:

- Define a closed convex hull of  $\pm\mathbb{D}$ :

$$B_1(\mathbb{D}) = \overline{\bigcup_{n=1}^{\infty} \Sigma_{n,1}^{\sigma}}, \quad (25)$$

- Define a norm

$$\|f\|_{\mathcal{K}_1(\mathbb{D})} = \inf\{r > 0 : f \in rB_1(\mathbb{D})\}, \quad (26)$$

as the guage of the set  $B_1(\mathbb{D})$ .

---

[12] siegel2021optimal.



# Stable Dictionary Approximation Space

Siegel & Xu, 2021<sup>[12]</sup>:

- Define a closed convex hull of  $\pm\mathbb{D}$ :

$$B_1(\mathbb{D}) = \overline{\bigcup_{n=1}^{\infty} \Sigma_{n,1}^{\sigma}}, \quad (25)$$

- Define a norm

$$\|f\|_{\mathcal{K}_1(\mathbb{D})} = \inf\{r > 0 : f \in rB_1(\mathbb{D})\}, \quad (26)$$

as the gauge of the set  $B_1(\mathbb{D})$ .

- The unit ball is

$$\{f \in H : \|f\|_{\mathcal{K}_1(\mathbb{D})} \leq 1\} = B_1(\mathbb{D}). \quad (27)$$

---

[12] siegel2021optimal.

# Stable Dictionary Approximation Space

Siegel & Xu, 2021<sup>[12]</sup>:

- Define a closed convex hull of  $\pm\mathbb{D}$ :

$$B_1(\mathbb{D}) = \overline{\cup_{n=1}^{\infty} \Sigma_{n,1}^{\sigma}}, \quad (25)$$

- Define a norm

$$\|f\|_{\mathcal{K}_1(\mathbb{D})} = \inf\{r > 0 : f \in rB_1(\mathbb{D})\}, \quad (26)$$

as the gauge of the set  $B_1(\mathbb{D})$ .

- The unit ball is

$$\{f \in H : \|f\|_{\mathcal{K}_1(\mathbb{D})} \leq 1\} = B_1(\mathbb{D}). \quad (27)$$

- We have

$$\{f \in H : \|f\|_{\mathcal{K}_1(\mathbb{D})} < \infty\} = \cup_{M>0} \overline{\cup_{n=1}^{\infty} \Sigma_{n,M}^{\sigma}} \quad (28)$$

is a Banach space.

---

[12] siegel2021optimal.

## Example: $H = \ell^2$

- Let  $H = \ell^2$ ,  $\mathbb{D} = \{\mathbf{e}_1, \mathbf{e}_2, \dots\}$ .
- What is  $B_1(\mathbb{D})$ ?

## Example: $H = \ell^2$

- Let  $H = \ell^2$ ,  $\mathbb{D} = \{e_1, e_2, \dots\}$ .
- What is  $B_1(\mathbb{D})$ ?
- The convex hull of  $\pm\mathbb{D}$  is

$$B_1(\mathbb{D}) = \{(a_1, a_2, \dots) \in \ell^2 : \sum_{i=1}^{\infty} |a_i| \leq 1\} \quad (29)$$

## Example: $H = \ell^2$

- Let  $H = \ell^2$ ,  $\mathbb{D} = \{\mathbf{e}_1, \mathbf{e}_2, \dots\}$ .
- What is  $B_1(\mathbb{D})$ ?
- The convex hull of  $\pm\mathbb{D}$  is

$$B_1(\mathbb{D}) = \{(a_1, a_2, \dots) \in \ell^2 : \sum_{i=1}^{\infty} |a_i| \leq 1\} \quad (29)$$

- Thus the norm is given by

$$\mathcal{K}_1(\mathbb{D}) = \ell^1 \subset \ell^2. \quad (30)$$

# Stable Dictionary Approximation Space

## Theorem (Siegel & Xu 2021)

A function  $f \in H = L^2(\Omega)$  can be approximated at all, i.e.

$$\lim_{n \rightarrow \infty} \inf_{f_n \in \Sigma_{n,M}(\mathbb{D})} \|f - f_n\|_H = 0, \quad (31)$$

for a fixed  $M < \infty$  if and only if

$$f \in MB_1(\mathbb{D}) \subset \mathcal{K}_1(\mathbb{D}).$$

# Stable Dictionary Approximation Space

## Theorem (Siegel & Xu 2021)

A function  $f \in H = L^2(\Omega)$  can be approximated at all, i.e.

$$\lim_{n \rightarrow \infty} \inf_{f_n \in \Sigma_{n,M}(\mathbb{D})} \|f - f_n\|_H = 0, \quad (31)$$

for a fixed  $M < \infty$  if and only if

$$f \in MB_1(\mathbb{D}) \subset \mathcal{K}_1(\mathbb{D}).$$

Furthermore, if

$$\|\mathbb{D}\| \equiv \sup_{h \in \mathbb{D}} \|h\|_H < \infty$$

we have

$$\inf_{f_n \in \Sigma_{n,M}(\mathbb{D})} \|f - f_n\|_H \leq n^{-\frac{1}{2}} \|\mathbb{D}\| \|f\|_{\mathcal{K}_1(\mathbb{D})}. \quad (32)$$

# The Spectral Barron Space

- Let  $f \in B_1(\mathbb{D})$ ,  $H = L^2(\Omega)$ ,  $\Omega = B_1^d = \{x \in \mathbb{R}^d : |x| \leq 1\}$ , and

$$\mathbb{D} = \mathbb{F}_s^d := \{(1 + |\omega|)^{-s} e^{2\pi i \omega \cdot x} : \omega \in \mathbb{R}^d\} \quad (33)$$

---

[13] [siegel2021optimal](#).



# The Spectral Barron Space

- Let  $f \in B_1(\mathbb{D})$ ,  $H = L^2(\Omega)$ ,  $\Omega = B_1^d = \{x \in \mathbb{R}^d : |x| \leq 1\}$ , and

$$\mathbb{D} = \mathbb{F}_s^d := \{(1 + |\omega|)^{-s} e^{2\pi i \omega \cdot x} : \omega \in \mathbb{R}^d\} \quad (33)$$

- In this case the norm is characterized by<sup>[13]</sup>

$$\|f\|_{\mathcal{K}_1(\mathbb{F}_s^d)} = \inf_{f_e|_{B_1^d} = f} \int_{\mathbb{R}^d} (1 + |\xi|)^s |\hat{f}_e(\xi)| d\xi, \quad (34)$$

where the infimum is taken over all extensions  $f_e \in L^1(\mathbb{R}^d)$ .

---

[13] siegel2021optimal.

# The Spectral Barron Space

- Let  $f \in B_1(\mathbb{D})$ ,  $H = L^2(\Omega)$ ,  $\Omega = B_1^d = \{x \in \mathbb{R}^d : |x| \leq 1\}$ , and

$$\mathbb{D} = \mathbb{F}_s^d := \{(1 + |\omega|)^{-s} e^{2\pi i \omega \cdot x} : \omega \in \mathbb{R}^d\} \quad (33)$$

- In this case the norm is characterized by<sup>[13]</sup>

$$\|f\|_{\mathcal{K}_1(\mathbb{F}_s^d)} = \inf_{f_e|_{B_1^d} = f} \int_{\mathbb{R}^d} (1 + |\xi|)^s |\hat{f}_e(\xi)| d\xi, \quad (34)$$

where the infimum is taken over all extensions  $f_e \in L^1(\mathbb{R}^d)$ .

- Property:

$$H^{s+\frac{d}{2}+\varepsilon}(\Omega) \hookrightarrow B^s(\Omega) \hookrightarrow W^{s,\infty}(\Omega). \quad (35)$$

---

[13] siegel2021optimal.

# The Barron Space

The results are proved in Siegel and Xu 2021<sup>[14]</sup>

- Let  $H = L^2(\Omega)$ ,  $\Omega = B_1^d = \{x \in \mathbb{R}^d : |x| \leq 1\}$ , and

$$\mathbb{D} = \mathbb{P}_k^d := \{\sigma_k(\omega \cdot x + b) : \omega \in S^{d-1}, b \in [-2, 2]\}, \quad (36)$$

where  $\sigma_k = [\max(0, x)]^k$ .

---

[14] **siegel2021optimal**.

[15] **ma2019barron**.

[16] **barron1993universal**.

# The Barron Space

The results are proved in Siegel and Xu 2021<sup>[14]</sup>

- Let  $H = L^2(\Omega)$ ,  $\Omega = B_1^d = \{x \in \mathbb{R}^d : |x| \leq 1\}$ , and

$$\mathbb{D} = \mathbb{P}_k^d := \{\sigma_k(\omega \cdot x + b) : \omega \in S^{d-1}, b \in [-2, 2]\}, \quad (36)$$

where  $\sigma_k = [\max(0, x)]^k$ .

- When  $k = 1$ ,  $\mathcal{K}_1(\mathbb{P}_k^d)$  is equivalent to the Barron space (introduced in<sup>[15]</sup>).

---

[14] **siegel2021optimal**.

[15] **ma2019barron**.

[16] **barron1993universal**.

# The Barron Space

The results are proved in Siegel and Xu 2021<sup>[14]</sup>

- Let  $H = L^2(\Omega)$ ,  $\Omega = B_1^d = \{x \in \mathbb{R}^d : |x| \leq 1\}$ , and

$$\mathbb{D} = \mathbb{P}_k^d := \{\sigma_k(\omega \cdot x + b) : \omega \in S^{d-1}, b \in [-2, 2]\}, \quad (36)$$

where  $\sigma_k = [\max(0, x)]^k$ .

- When  $k = 1$ ,  $\mathcal{K}_1(\mathbb{P}_k^d)$  is equivalent to the Barron space (introduced in<sup>[15]</sup>).
- When  $k = 0$ ,  $d = 1$ ,  $\mathcal{K}_1(\mathbb{P}_k^d) = BV([-1, 1])$ .

---

[14] **siegel2021optimal**.

[15] **ma2019barron**.

[16] **barron1993universal**.

# The Barron Space

The results are proved in Siegel and Xu 2021<sup>[14]</sup>

- Let  $H = L^2(\Omega)$ ,  $\Omega = B_1^d = \{x \in \mathbb{R}^d : |x| \leq 1\}$ , and

$$\mathbb{D} = \mathbb{P}_k^d := \{\sigma_k(\omega \cdot x + b) : \omega \in S^{d-1}, b \in [-2, 2]\}, \quad (36)$$

where  $\sigma_k = [\max(0, x)]^k$ .

- When  $k = 1$ ,  $\mathcal{K}_1(\mathbb{P}_k^d)$  is equivalent to the Barron space (introduced in<sup>[15]</sup>).
- When  $k = 0$ ,  $d = 1$ ,  $\mathcal{K}_1(\mathbb{P}_k^d) = BV([-1, 1])$ .
- We have  $\mathcal{K}_1(\mathbb{P}_k^d) \supset \mathcal{K}_1(\mathbb{F}_{k+1}^d)$  (for  $k = 0$ , Barron 1993<sup>[16]</sup>)

---

[14] **siegel2021optimal**.

[15] **ma2019barron**.

[16] **barron1993universal**.

# Previous Best Results

For some dictionaries  $\mathbb{D}$ , the  $n^{-\frac{1}{2}}$  approximation rate can be improved!

- For  $\mathbb{D} = \mathbb{P}_0^d$ , we have<sup>[17]</sup>

$$\sup_{f \in B_1(\mathbb{D})} \inf_{f_n \in \Sigma_{n,M}} \|f - f_n\|_{L^2(B_1^d)} \lesssim n^{-\frac{1}{2} - \frac{1}{2d}}. \quad (37)$$

- For  $\mathbb{D} = \mathbb{P}_k^d$  for  $k \geq 1$ , we have<sup>[18],[19]</sup>, if  $f$  is in some spectral Barron space:

$$\inf_{f_n \in \Sigma_{n,M}} \|f - f_n\|_{L^2(B_1^d)} \lesssim n^{-\frac{1}{2} - \frac{1}{d}}. \quad (38)$$

---

[17] makovoz1996random.

[18] klusowski2018approximation.

[19] CiCP-28-1707.

# Previous Best Results

For some dictionaries  $\mathbb{D}$ , the  $n^{-\frac{1}{2}}$  approximation rate can be improved!

- For  $\mathbb{D} = \mathbb{P}_0^d$ , we have<sup>[17]</sup>

$$\sup_{f \in B_1(\mathbb{D})} \inf_{f_n \in \Sigma_{n,M}} \|f - f_n\|_{L^2(B_1^d)} \lesssim n^{-\frac{1}{2} - \frac{1}{2d}}. \quad (37)$$

- For  $\mathbb{D} = \mathbb{P}_k^d$  for  $k \geq 1$ , we have<sup>[18],[19]</sup>, if  $f$  is in some spectral Barron space:

$$\inf_{f_n \in \Sigma_{n,M}} \|f - f_n\|_{L^2(B_1^d)} \lesssim n^{-\frac{1}{2} - \frac{1}{d}}. \quad (38)$$

- What are the optimal approximation rates?

---

[17] makovoz1996random.

[18] klusowski2018approximation.

[19] CiCP-28-1707.



# New Optimal Bounds

We have the optimal approximation rates<sup>[20]</sup>

## Theorem

For  $\mathbb{D} = \mathbb{P}_k^d$  for  $k \geq 1$ , we have

$$n^{-\frac{1}{2} - \frac{2k+1}{2d}} \lesssim \sup_{f \in B_1(\mathbb{D})} \inf_{f_n \in \Sigma_{n,M}} \|f - f_n\|_{L^2(\Omega)} \lesssim n^{-\frac{1}{2} - \frac{2k+1}{2d}} \quad (39)$$

---

[20] siegel2021optimal.

[21] lin2014lower.

# New Optimal Bounds

We have the optimal approximation rates<sup>[20]</sup>

## Theorem

For  $\mathbb{D} = \mathbb{P}_k^d$  for  $k \geq 1$ , we have

$$n^{-\frac{1}{2} - \frac{2k+1}{2d}} \lesssim \sup_{f \in B_1(\mathbb{D})} \inf_{f_n \in \Sigma_{n,M}} \|f - f_n\|_{L^2(\Omega)} \lesssim n^{-\frac{1}{2} - \frac{2k+1}{2d}} \quad (39)$$

In comparison: optimal bound for finite elements<sup>[21]</sup>

## Theorem

Assume that  $V_h^k$  is a finite element of degree  $k$  on quasi-uniform mesh  $\{\mathcal{T}_h\}$  of  $\mathcal{O}(N)$  elements. Assume  $u$  is sufficiently smooth and not piecewise polynomials, then we have

$$c(u)n^{-\frac{k}{d}} \leq \inf_{v_h \in V_h^k} \|u - v_h\|_{L^2(\Omega)} \leq C(u)n^{-\frac{k}{d}} = \mathcal{O}(h^k). \quad (40)$$

---

[20] siegel2021optimal.

[21] lin2014lower.

# Removing the constraint that $\sum_{i=1}^n |a_i| \leq M$

Define

$$\Sigma_n^k := \left\{ \sum_{i=1}^n a_i \sigma_k(\omega_i \cdot x + b_i), \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R}, a_i \in \mathbb{R} \right\}. \quad (41)$$

Then  $\Sigma_n^k$  has the following approximation property<sup>[22]</sup>

## Theorem (Siegel and Xu)

$$\inf_{f_n \in \Sigma_n^k} \|f - f_n\| \lesssim \begin{cases} n^{-\frac{1}{2}} & \|f\|_{\mathcal{K}_1(\mathbb{R}_s^d)} \\ n^{-(k+1)} \log n & \|f\|_{\mathcal{K}_1(\mathbb{R}_s^d)} \end{cases} \quad \begin{matrix} \text{if } s = \frac{1}{2} \\ \text{for some } s > 1 \end{matrix} \quad (42)$$

- Improves result of Barron<sup>[23]</sup> by relaxing condition on  $f$

<sup>[22]</sup> siegel2020high.

<sup>[23]</sup> barron1993universal.

# Removing the constraint that $\sum_{i=1}^n |a_i| \leq M$

Define

$$\Sigma_n^k := \left\{ \sum_{i=1}^n a_i \sigma_k(\omega_i \cdot x + b_i), \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R}, a_i \in \mathbb{R} \right\}. \quad (41)$$

Then  $\Sigma_n^k$  has the following approximation property<sup>[22]</sup>

## Theorem (Siegel and Xu)

$$\inf_{f_n \in \Sigma_n^k} \|f - f_n\| \lesssim \begin{cases} n^{-\frac{1}{2}} & \|f\|_{\mathcal{K}_1(\mathbb{R}_s^d)} \\ n^{-(k+1)} \log n & \|f\|_{\mathcal{K}_1(\mathbb{R}_s^d)} \end{cases} \quad \begin{matrix} \text{if } s = \frac{1}{2} \\ \text{for some } s > 1 \end{matrix} \quad (42)$$

- Improves result of Barron<sup>[23]</sup> by relaxing condition on  $f$
- Shows that very high order approximation rates can be attained with sufficient smoothness

---

[22] siegel2020high.

[23] barron1993universal.

# Removing the constraint that $\sum_{i=1}^n |a_i| \leq M$

Define

$$\Sigma_n^k := \left\{ \sum_{i=1}^n a_i \sigma_k(\omega_i \cdot x + b_i), \omega_i \in \mathbb{R}^d, b_i \in \mathbb{R}, a_i \in \mathbb{R} \right\}. \quad (41)$$

Then  $\Sigma_n^k$  has the following approximation property<sup>[22]</sup>

## Theorem (Siegel and Xu)

$$\inf_{f_n \in \Sigma_n^k} \|f - f_n\| \lesssim \begin{cases} n^{-\frac{1}{2}} & \|f\|_{\mathcal{K}_1(\mathbb{R}_s^d)} \\ n^{-(k+1)} \log n & \|f\|_{\mathcal{K}_1(\mathbb{R}_s^d)} \end{cases} \quad \begin{matrix} \text{if } s = \frac{1}{2} \\ \text{for some } s > 1 \end{matrix} \quad (42)$$

- Improves result of Barron<sup>[23]</sup> by relaxing condition on  $f$
- Shows that very high order approximation rates can be attained with sufficient smoothness
- Comparison with FEM:

$$\inf_{w \in \Sigma_n^k} \|u - w\| \approx \left\{ \inf_{v \in V_n^k} \|u - v\| \right\}^d.$$

[22] siegel2020high.

[23] barron1993universal.

# References

- J. He, L. Li, J. Xu and C. Zheng, ReLU Deep Neural Networks and Linear Finite Elements. Journal of Computational Mathematics 38.3 (2020), pp. 502527.
- J. Xu, The Finite Neuron Method and Convergence Analysis, Commun. Comput. Phys., 28, pp. 1707-1745, (2020).
- J. W. Siegel and J. Xu, "High-Order Approximation Rates for Neural Networks with  $\text{ReLU}^k$  Activation Functions." arXiv preprint arXiv:2012.07205 (2020).
- J. W. Siegel and J. Xu, "Approximation Rates for Neural Networks with General Activation Functions." Neural Networks (2020).
- J. W. Siegel and J. Xu, "Optimal Approximation Rates and Entropy Bounds for  $\text{ReLU}^k$  Networks." arXiv preprint arXiv:2101.12365 (2021)

**Thank you!**