# Mathematics of Neural Networks

J. Vybíral

Czech Technical University
Faculty of Nuclear Sciences and Physical Engineering
https://kmlinux.fjfi.cvut.cz/~vybirja2/

Nečas Center for Mathematical Modeling
https://ncmm.karlin.mff.cuni.cz/

Ostrava, January 2025

## Overview

- **Introduction**: History, notation, activation functions, architectures, interpretability

- **Examples:** Hat function, square function, composition, addition, $\cos(ax)$

- **Function classes**: Upper bounds for function spaces, fixed architecture, variable architecture

- **Barron classes**: Definition, recovery by ANN's

- **Lower bounds**: VC-dimension, continuous widths, lower bounds for shallow and deep networks

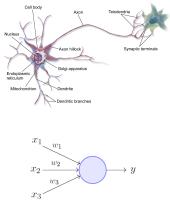- **Riesz basis**: New univariate and multivariate Riesz basis of ANN's, avoiding curse of dimensionality

# Introduction

## Neural Networks

Artificial neural networks (ANN):

- Neurology and biological neural networks & mathematical models: Warren McCulloch and Walter Pitts, 1943

- General idea: Simple building blocks (biological neurons) connected in a large network

- In mathematical biology, the effect of complex behavior arising from simple interactions of agents is well-known

- Frank Rosenblatt (1958), Marvin Minsky and Seymour Papert: Perceptrons (1969) - model of "perceptron", i.e., artificial neuron

# Biological and artificial neurons



Perceptron Model (Minsky-Papert in 1969)

## Artificial neuron - perceptron

- Biological motivation: Neuron combines several inputs (from other neurons), if the inputs are strong enough, the neuron gets "activated" and sends some signal out
- Mathematical model:
  - Neuron gets $n$ real inputs $x_1, \ldots, x_n \in \mathbb{R}$
  - combines them linearly - i.e., computes an inner product $w \cdot x = w_1 x_1 + \cdots + w_n x_n = \langle w, x \rangle$, where $w = (w_1, \ldots, w_n)$ are *weights*
  - if this inner product is large enough (larger than some real bias $b \in \mathbb{R}$), it gives some output, otherwise the output is zero
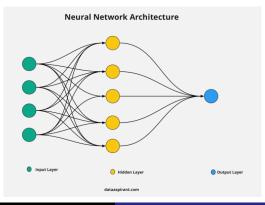- Perceptron is modeled by a function

$$f(x) = \mathbb{1}_{\{\langle w, x \rangle \geq b\}} = \begin{cases} 1, & \text{if } \langle w, x \rangle - b \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

## ANN's - activation functions

- The Heaviside function $\mathbb{1}$ can be replaced by many other choices: monotone, non-linear, bounded or unbounded range, (continuously) differentiable, etc.

- ReLU activation function, K. Fukushima (1969)

- $\text{ReLU}(x) = \max(0, x)$   rectified linear unit

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

- $\tanh(x)$   hyperbolic tangent

- $\ln(1 + e^x)$   Softplus

- $\frac{x}{1+e^{-x}}$   Sigmoid linear unit (SiLU)
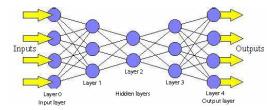
## Networks of neurons

- (Large) number of neurons get connected by a graph and output of some neurons become inputs of the others, based on the connections in the graph
- Example:



Neural Network Architecture

Input Layer     Hidden Layer     Output Layer

dataaspirant.com

## Feedforward neural networks

Feedforward neural networks:

- Artificial neurons are ordered in layers
- outputs of $j$th layer become inputs for $(j+1)$th layer
- Layer0 are the inputs
- Last layer are the outputs

## Feedforward neural networks

- Affine function: $f : \mathbb{R}^{n_1} \to \mathbb{R}^{n_2}$:

  $f(x) = Mx + b$, where $M \in \mathbb{R}^{n_2 \times n_1}$ is a matrix and $b \in \mathbb{R}^{n_2}$.

- Let $d, W, L$ be positive integers

- $\text{ReLU}(x_1, \ldots, x_d) = (\text{ReLU}(x_1), \ldots, \text{ReLU}(x_d))$

- a feed-forward ReLU network $\mathcal{N}$ with width $W$ and depth $L$ is a collection of $L + 1$ affine mappings $A^{(0)}, \ldots, A^{(L)}$, where $A^{(0)} : \mathbb{R}^d \to \mathbb{R}^W$, $A^{(j)} : \mathbb{R}^W \to \mathbb{R}^W$ for $j = 1, \ldots, L - 1$ and $A^{(L)} : \mathbb{R}^W \to \mathbb{R}$. Each such a network $\mathcal{N}$ generates a function of $d$ variables

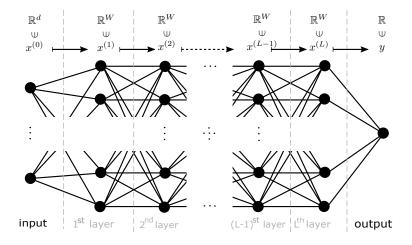$$A^{(L)} \circ \text{ReLU} \circ A^{(L-1)} \circ \cdots \circ \text{ReLU} \circ A^{(0)}.$$

- $\Upsilon^{W,L}$: the set of all functions, which are generated by some feed-forward ReLU network with width $W$ and depth $L$

Figure: Feed-forward $\mathrm{ReLU}$ network with length $L$, width $W$

## Types of networks

Many different architectures:

- Convolutional Neural Networks (CNNs): used in image analysis, different type of layers, apply local filters

- Recurrent Neural Networks (RNNs): sequential data, used in language processing and time series prediction

- Transformers: used in Language Processing

## Training

- We assume some training data $(x_i, g(x_i))$, $i = 1, \ldots, N$ is given

- We fix the architecture and its parameters

- We set the weights to fit as good as possible to the known data (backpropagation algorithm)

- Large networks require large amount of data and large number of parameters (overfitting?)

- Extensive use of data and computational power (GPU's by NVIDIA)

## Interpretability

Difficult interpretability - issue in critical applications



Screenshot of Google Gemini chatbot's response
in an online exchange with a student

## Success stories

- Pre-NN's: Deep Blue - IBM chess playing; 1996 - first matches with Garry Kasparov, winning in 1997

- DeepMind (London-based company, acquired by Google): AlphaGo (2015)

- Other games: poker (incomplete information game), bridge, Starcraft II, DotA 2

- Computer Vision: character recognition, facial recognition, medical imaging, and autonomous vehicles

- Natural Language Processing (NLP): language translation, sentiment analysis, and chatbots

- Speech Recognition: Siri, Alexa, ...

- ChatGPT (Generative pre-trained transformer)

## ChatGPT vs. DeepSeek

Forty percent of families in a city have one child, thirty percent of families have two children, twenty percent of families have three children, and ten percent of families have four children.

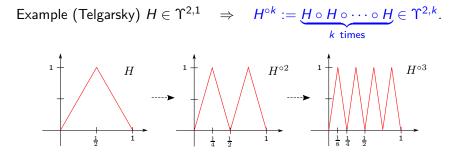What is the probability that a randomly selected child has a sister?

# Examples

## Example: the hat function

The **hat function** $H : [0, 1] \rightarrow [0, 1]$,

$$
H(x) \;=\; \begin{cases} 2x, & 0 \leq x \leq \frac{1}{2} \\ 2(1 - x), & \frac{1}{2} < x \leq 1 \end{cases}
$$

$$
=\; 2(x - 0)_+ - 4 \left( x - \frac{1}{2} \right)_+
$$



$$
=\; \begin{bmatrix} 2 & -4 \end{bmatrix} \operatorname{ReLU} \left\{ \begin{bmatrix} 1 \\ 1 \end{bmatrix} x + \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix} \right\}
$$

belongs to $\Upsilon^{2,1}$.

## Compositions and deep networks

Example (Telgarsky) $H \in \Upsilon^{2,1}$  $\Rightarrow$  $H^{\circ k} := \underbrace{H \circ H \circ \cdots \circ H}_{k \text{ times}} \in \Upsilon^{2,k}.$

## Hat function revisited

Hat function $H : \mathbb{R} \to [0, 1]$ on the whole $\mathbb{R}$:

$$H(x) = \begin{cases} 2x, & \text{if } 0 \leq x < 1/2, \\ 2(1-x), & \text{if } 1/2 \leq x \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

$$H(x) = 2\,\text{ReLU}(x) - 4\,\text{ReLU}(x - 1/2) + 2\,\text{ReLU}(x - 1)$$
$$= W_2 \circ \text{ReLU} \circ W_1$$
$$W_1(x) = (x, x - 1/2, x - 1)^T$$
$$W_2(y) = 2y_1 - 4y_2 + 2y_3.$$

. . . can be again iteratively composed with itself.

## Quadratic function

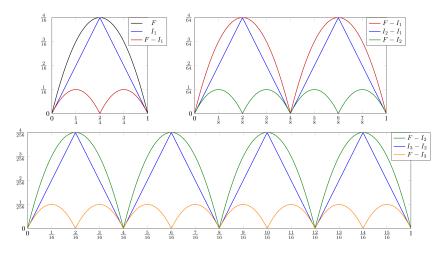We would like to reproduce the product function $(x, y) \rightarrow x \cdot y$

With the ReLU activation function - it is surprisingly difficult!

The ANN's with ReLU generate only piecewise linear functions

Due to

$$x \cdot y = (x + y)^2/4 - (x - y)^2/4$$

it is enough to approximate the univariate function $x \rightarrow x^2$

We approximate $F(x) = x - x^2$ by iterated $H$ functions!

From: Deep Neural Network Approximation Theory (by Elbrächter, Perekrestenko, Grohs, and Bölcskei, 2021)

## Composition and addition

- If $f_1 : \mathbb{R}^{d_1} \to \mathbb{R}^{d_2}$ and $f_2 : \mathbb{R}^{d_2} \to \mathbb{R}^{d_3}$ are well representable by ANN, then so is $f_2 \circ f_1$

  - We first calculate $f_1(x)$ by an ANN, the output serves as an input for the ANN of $f_2$.

  - We "glue" the networks after each other, we add the lengths

- If $f_1 : \mathbb{R}^d \to \mathbb{R}$ and $f_2 : \mathbb{R}^d \to \mathbb{R}$ can be represented by neural networks, then $f_1 + f_2$ can be also well represented.

  - We "glue" the networks on the top of each other: length is the same, the widths add together

  - Or we glue the networks after each other, the inputs are passed as well as the already calculated outputs: lengths are added

  - Or we do something in between

## Composition and addition

### Proposition (Properties of $\Upsilon^{W,L}$, $d = 1$)

Let $W \geq 2$. For any $\mathcal{Y}_1 \in \Upsilon^{W,L_1}, \ldots, \mathcal{Y}_k \in \Upsilon^{W,L_k}$ the following holds:

(i) The composition of the $\mathcal{Y}_i$ satisfies

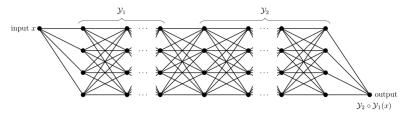$$\mathcal{Y}_k \circ \cdots \circ \mathcal{Y}_1 \in \Upsilon^{W,L}, \qquad L = L_1 + \cdots + L_k.$$

(ii) The sum of the $\mathcal{Y}_i$ satisfies

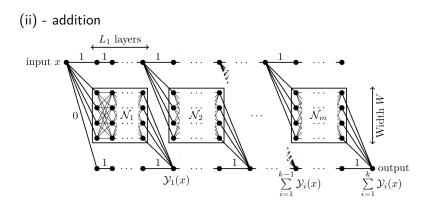$$\mathcal{Y}_1 + \cdots + \mathcal{Y}_k \in \Upsilon^{W+2,L}, \qquad L = L_1 + \cdots + L_k.$$

## Composition and addition

Proof: (i) - composition



The network computing $\mathcal{Y}_2 \circ \mathcal{Y}_1$.

## Composition and addition

(ii) - addition

# Approximation of $\cos(ax)$

Now we can

- Approximate $x \to x^k$ for $k \geq 2$
- Approximate $x \to \sum_{i=0}^{m} a_i x^i$
- Approximate $\cos(ax)$ by a partial sum of its Taylor series
- Put all these pieces together!
- Theorem (EPGB): For $a, D > 0$ there is an ANN, which approximates $\cos(ax)$ on $[-D, D]$ uniformly up to $\varepsilon > 0$ error with length $L \leq C(\log^2(\varepsilon^{-1}) + \log(\lceil aD \rceil))$ and width $W \leq 9$

# Function classes

# Universal approximation theorem(s)

- Appeared around 1990

- Show the density of functions generated by ANN's in some function classes, typically $C(\Omega)$

- Limit (aka asymptotic) theorems - usually no guarantee on $W$, $L$

- Cybenko (1989): Sums $\sum_{j=1}^{N} \alpha_j \sigma(y_j^T x + \theta_j)$ are dense in $C([0,1]^d)$ - the width is arbitrary, $\sigma$ is a rather general function

- Hornik, Stinchcombe, White (1989): Universal approximation of ANN's with one hidden layer.

- Maiorov and Pinkus (1999): There is an activation function $\sigma$, such that for any $f \in C([0,1]^d)$ there is a two-layer network with width bounded by $c \cdot d$, which approximates $f$ up to $\varepsilon$ in the uniform norm.

## What we do?

We are interested in approximation properties of ANN's on whole function classes, not only for one particular function (like $\cos(ax)$)

Questions we address:

- (mainly) ReLU activation function
- Different function classes (Sobolev, Besov, $BV$, Lipschitz, etc.)
- Dependence of the (worst-case) error of approximation on width $W$ and length $L$
- Decay rates - upper bounds
- Lower bounds

## What we do not do?

(Important) questions we do *not* address

- The role of different activation functions

- Relation to backpropagation

- - and to trained networks

- Lack of overfitting?!

- Typical good local minima?!

- . . . and many others

## Re-using existing decompositions

- For many function spaces $X(\Omega)$ we have suitable (and optimal) decompositions into Fourier basis, splines, wavelets, curvelets and other -lets

- If we can represent/approximate these building blocks as ANN, then we first decompose $f \in X(\Omega)$ into these blocks, then represent these blocks as ANN's and finally, we sum it all up!

- *Transference principle*

Yarotsky (2017):

Error bounds for approximations with deep ReLU networks

- 
$$\|f\|_{W^{n,\infty}([0,1]^d)} := \max_{\alpha:|\alpha|\le n} \sup_{x\in[0,1]^d} \|D^\alpha f(x)\|$$

- $|\alpha| = \alpha_1 + \cdots + \alpha_d$

- "Error" is always the error in the uniform norm

- Theorem: For $n, d \in \mathbb{N}$ and $1 > \varepsilon > 0$, there is an ReLU ANN
  architecture, such that (with appropriate weights) approximates
  every $f$ from the unit ball of $W^{n,\infty}([0,1]^d)$ up to error $\varepsilon$
  - fixed architecture for different $f$'s
  - $L \le c(\ln(1/\varepsilon) + 1)$, number of weights $\le c\varepsilon^{-d/n}(\ln(1/\varepsilon) + 1)$
  - $c = c(d, n)$

Yarotsky (2017):

Proof idea:

- Decomposition of unity $\phi_m(x)$, $m = (m_1, \ldots, m_d) \in \{0, \ldots, N\}^d$

- $P_m(x)$ - local polynomial approximation of $f$ around $m/N$

- $f_1(x) = \sum_m \phi_m P_m(x)$

- $\|f - f_1\|_\infty \leq \varepsilon/2$

- We approximate terms $\phi_m P_m(x)$ by ANN's

- And sum them up!

## Speed up through adaptivity - Yarotsky (2017)

**Choosing different architecture
for every $f$ can speed up the approximation!**

Consider $d = 1$ and $f \in W^{1,\infty}([0,1])$

How many weights do we need for uniform $\varepsilon$-error?

**Previous approach:** number of weights $O(\varepsilon^{-1} \cdot \ln(1/\varepsilon))$
**Piece-wise linear approximation:** number of weights $O(\varepsilon^{-1})$
**Choosing adaptive network architecture:** $O(\varepsilon^{-1}/\ln(1/\varepsilon))$

# Speed up through adaptivity - Yarotsky (2017)

Proof idea:

- We first approximate $f$ by piece-wise linear function $\tilde{f}_1$, but on a scale $m\varepsilon, m \approx \ln(1/\varepsilon) > 1$

- We add to $\tilde{f}_1$ a correction term $\tilde{f}_2$, which bridges the grid of scale $m\varepsilon$ to the scale $\varepsilon$

- The function $f - \tilde{f}_1$ vanishes at 0 and $m\varepsilon$. We split $[0, m\varepsilon]$ into $m$ intervals and collect functions $\gamma$ with

    $$\gamma(0) = \gamma(m\varepsilon) = 0, \qquad \gamma(j\varepsilon) - \gamma((j-1)\varepsilon) \in \{-2\varepsilon, 0, 2\varepsilon\}$$

- Finitely many pre-cached functions; adding the correct terms reduces the error of approximation

## General function classes

Previous approach can be vastly generalized!

- We can consider $\Omega = [0,1]^d$, or a general domain $\Omega \subset \mathbb{R}^d$
- $f$ can be from the Sobolev space $W^s(L_q(\Omega))$ or Besov space $B_r^s(L_q(\Omega))$
- The error of approximation can be measured in $L_p(\Omega)$
- Siegel (2023): For $W$ fixed ($W = 25d + 31$) and $L \to \infty$, the error decay rate is $L^{-2s/d}$
- The factor 2 is due to the "bit-extraction technique" of Bartlett, Maiorov, and Meir (1998)

## Shallow vs. deep neural networks

- Approximation-theoretic properties of deep neural networks:
  - R. DeVore, B. Hanin, and G. Petrova, *Neural network approximation*, Acta Numerica 30 (2021): 327–444
  - R. Gribonval, G. Kutyniok, M. Nielsen, and F. Voigtlaender, *Approximation spaces of deep neural networks*, Constr. Appr. 55(1) (2022), 259–367
  - D. Elbrächter, D. Perekrestenko, P. Grohs, and H. Bölcskei, *Deep neural network approximation theory*, IEEE Trans. Inf. Theory 67(5) (2021), 2581–2623.

- Existence of functions which, can be expressed through a small three-layer network, but can only be represented through a very large two-layer network: R. Eldan and O. Shamir (2016)

- Further examples: H. N. Mhaskar and T. Poggio (2016)

# Barron classes

## Barron classes - definition

**Barron 1994**:

Let $f : \mathbb{R}^d \to \mathbb{R}$ (or $\mathcal{C}$), let $\hat{f}$ be its Fourier transform.

**Barron classes on** $\mathbb{R}^d$: We define

$$\|f\|_B := \int_{\mathbb{R}^d} \|\xi\|_2 \cdot |\hat{f}(\xi)| d\xi,$$

where $\|\xi\|_2^2 = \sum_j |\xi_j|^2$

**Barron classes on** $\Omega \subset \mathbb{R}^d$: By restriction

## Convex combinations and random sampling

We need "empirical method of B. Maurey",

a.k.a. "Approximate Caratheodory's theorem"

**Theorem:** Let $T$ be a subset of the unit ball of $\mathbb{R}^D$. Then, for every $n$ and every $x \in \text{conv}(T)$, there are $x_1, \ldots, x_n \in T$ with

$$\left\| x - \frac{1}{n} \sum_{j=1}^{n} x_j \right\|_2 \leq \frac{1}{\sqrt{n}}.$$

**Remarks:**

- $D$ does not matter, Hilbert space $H$ possible
- Probabilistic proof!?

## Convex combinations and random sampling - proof

**Proof:**

- Let $x \in \text{conv}(T)$, i.e., $x = \sum_k \gamma_k y_k$; $\gamma_k \geq 0, \sum_k \gamma_k = 1, y_k \in T$
- Let $g$ be a random variable: $\mathbb{P}(g = y_k) = \gamma_k$
- Let $g_1, \ldots, g_n$ be $n$ independent copies of $g$; put $\tilde{x} = \frac{1}{n} \sum_{j=1}^{n} g_j$
- $\mathbb{E}g = \sum_k y_k \cdot \mathbb{P}(g = y_k) = x$; $\mathbb{E}\tilde{x} = x$
- Finally

$$
\begin{aligned}
\mathbb{E}\|x - \tilde{x}\|^2 = \mathbb{E}\Big\|\frac{1}{n}\sum_{j=1}^{n}(x - g_j)\Big\|^2 &= \frac{1}{n^2}\sum_{j,k=1}^{n}\mathbb{E}\langle x - g_j, x - g_k\rangle \\
&= \frac{1}{n^2}\sum_{j=1}^{n}\mathbb{E}\|x - g_j\|^2 = \frac{1}{n}\mathbb{E}\|x - g\|^2 = \frac{1}{n}\mathbb{E}\|g - \mathbb{E}g\|^2 \\
&= \frac{1}{n}(\mathbb{E}\|g\|^2 - \|\mathbb{E}g\|^2) \leq \frac{1}{n}.
\end{aligned}
$$

- There is some realization of $g_1, \ldots, g_n$, which is at least so good!

## Barron class - approximation by ANN's

Let $B = \{x \in \mathbb{R}^d : \|x\|_2 \leq 1\}$ be the unit ball, $\mu$ - prob. measure on $B$

**Theorem (Barron - 1993):**
For every $f : \mathbb{R}^d \to \mathbb{R}$ with $\int \|\xi\|_2 \cdot |\hat{f}(\xi)| d\xi = C$
there exists

$$f_n(x) = \sum_{k=1}^{n} c_k \phi(a_k \cdot x + b_k) + c_0$$

with

$$\|f - f_n\|_{L_2(B,\mu)}^2 \leq \frac{(2C)^2}{n}.$$

## Barron class - approximation by ANN's

**Proof idea:** $\hat{f}(\xi) = e^{i\theta(\xi)}|\hat{f}(\xi)|$

$$f(x) = Re \int_{\mathbb{R}^d} e^{i\xi \cdot x}\hat{f}(\xi)d\xi = Re \int_{\mathbb{R}^d} e^{i\xi \cdot x}e^{i\theta(\xi)}|\hat{f}(\xi)|d\xi$$
$$= \int_{\mathbb{R}^d} \cos(\xi \cdot x + \theta(\xi))|\hat{f}(\xi)|d\xi$$

- $f(x)$ is a linear combination of functions $x \to c \cos(\xi \cdot x + b)$

- come over to convex combination

- $\cos(\langle \cdot, x \rangle + b)$ is in the convex hull of $\phi(\langle \cdot, x \rangle + b)$

- apply the "empirical method of B. Maurey"

# Barron class - extensions

- For general domain, replace $\|\xi\|_2$ with $\|\xi\|_B := \sup_{x \in B} |\xi \cdot x|$

- There are no implicit constant depending on $d$

- Use of these spaces avoids the curse of dimensionality

- The construction is only implicit, randomized

# Lower bounds

## Lower bounds - general setting

- By definition, we want to show (for example) that

  $$\sup_{f \in X} \inf_{\mathcal{N} \in \Upsilon^{W,L}} \|f - \mathcal{N}\|_Y$$

  is large.

- Formally, it is enough to find one $f \in X$, which does not look like ANN's from $\Upsilon^{W,L}$

- Instead, we want to show that the class $X$ of all possible functions $f \in X$, which we want to recover/approximate is too large, to be approximated by ANN's with a small number of parameters.

- We introduce some complexity measure of a class of functions and show that the complexity of $X$ is larger than that of $\Upsilon^{W,L}$, if $W$ and $L$ are small.

## VC dimension - definition

**Vapnik–Chervonenkis dimension**

Let $X$ be a set and $\mathcal{H}$ a system of its subsets. We say that $C \subset X$ is shattered by $\mathcal{H}$ if $H \cap C$ takes all $2^{|C|}$ different values when $H$ runs over $\mathcal{H}$, i.e., if

- $|\{H \cap C : H \in \mathcal{H}\}| = 2^{|C|}$, or
- for every $B \subset C$ there exists $H \in \mathcal{H}$ such that $H \cap C = B$

The Vapnik–Chervonenkis dimension of $\mathcal{H}$ (the VC dimension of $\mathcal{H}$) is the size of the largest $C$ shattered by $\mathcal{H}$

## VC dimension - examples

**Examples**

- If $C$ is shattered by $\mathcal{H}$, then $2^{|C|} = |\{H \cap C : H \in \mathcal{H}\}| \leq |\mathcal{H}|$
- ... hence VC-dim$(\mathcal{H}) \leq \log_2(|\mathcal{H}|)$
- $X = \{1, \ldots, n\}$ and $\mathcal{H}$ are all subsets of $X$ with at most $k$ elements, then VC-dim$(\mathcal{H}) = k$
- Sauer–Shelah lemma: VC-dim$(\mathcal{H}) \leq k$ implies $|\mathcal{H}| \leq \sum_{j=0}^{k} \binom{n}{j}$
- Rademacher complexity of $H \subset \mathbb{R}^d$

$$\mathsf{Rad}(H) := \frac{1}{d} \, \mathbb{E}_\sigma \left[ \sup_{\sigma \in H} \sum_{j=1}^{d} \sigma_j a_j \right]$$

$\sigma_1, \ldots, \sigma_d$ - independent Rademacher variables

## VC dimension of sets of functions

**VC-dimension of a set of binary functions**

Let $X$ be a set and let $\mathcal{F}$ be a set of functions $f : X \to \{0, 1\}$. We define $H_f := \{x : f(x) = 1\} \subset X$ and

$$\text{VC-dim}(\mathcal{F}) := \text{VC-dim}(\{H_f : f \in \mathcal{F}\})$$

**VC-dimension of a set of functions**

Let $X$ be a set and let $\mathcal{F}$ be a set of functions $f : X \to \mathbb{R}$. We define

$$\text{VC-dim}(\mathcal{F}) := \text{VC-dim}(\{\text{sgn}(f) : f \in \mathcal{F}\})$$

## VC dimension of ANN's

Bartlett, Maiorov, Meir (1998):

- Fix an architecture of an ANN with $\omega$ parameters and $L$ layers and piecewise polynomial activation function. The set $\mathcal{F}$ of all functions, which we obtain by varying the $\omega$ parameters has

$$\text{VC-dim}(\text{sgn}(\mathcal{F})) \leq c(\omega L \log \omega + \omega L^2)$$

- For fixed $L$: $O(\omega \log \omega)$
- Lower bound: There is an architecture, where
$$\text{VC-dim}(\text{sgn}(\mathcal{F})) \geq c\,\omega L$$

# VC dimension of ANN's - proof idea

- Upper bound: We count, layer-by-layer the number of all possible vectors

$$(\mathrm{sgn}(f(x_1)), \ldots, \mathrm{sgn}(f(x_n))) \in \{0, 1\}^n$$

  if $f$ goes over $\mathcal{F}$

- Lower bound:

  Construction of a net with $O(\nu)$ weights and $O(\mu)$ layers

  If $a_i = \sum_{j=1}^{\mu} 2^{-j} a_{i,j}$ is a binary representation of the parameters $a_1, \ldots, a_\nu$, then for the input $x = (e_l, e_m) \in \{0, 1\}^{\nu + \mu}$ the net outputs $a_{l,m}$

  The set $\{(e_l, e_m) : 1 \leq l \leq \nu, 1 \leq m \leq \mu\}$ with $\nu \cdot \mu$ elements is shattered by the set of ANN's when $a \in \{0, 1\}^{\mu \times \nu}$

## Lower bounds for function recovery - widths

**Lower bounds by *continuous nonlinear widths***

**Theorem (DeVore, Howard, Micchelli, 1989)**:
Let $\mathcal{M} : W^{n,\infty}([0,1]^d) \to \mathbb{R}^\omega$ be continuous and let
$\eta : \mathbb{R}^\omega \to C([0,1]^d)$ be arbitrary. If $\|f - \eta(\mathcal{M}(f))\|_\infty \leq \varepsilon$ for all $f$ in
the unit ball of $W^{n,\infty}$, then $\omega \geq c\varepsilon^{-d/n}$.

Hence:

If the architecture is fixed, the parameters of the network depend
continuously on $f$, and it approximates every $f \in W^{n,\infty}$ up to uniform
$\varepsilon$-error, then the network must have at least $\varepsilon^{-d/n}$ parameters.

- The upper bound of Yarotsky (fixed architecture) is nearly optimal
- Kainen, Kůrková, Vogt (1999): The optimal weights do not
  depend continuously on $f$

## Lower bounds for function recovery - VC dimension

Yarotsky (2017):

- Consider points $x_1, \ldots, x_{N^d} \in [0,1]^d$, mutual distance $\geq 1/N$
- Construct $f \in W^{n,\infty}$, which takes $\pm 2\varepsilon$ values at $x_1, \ldots, x_{N^d}$
- The VC-dimension of $W^{n,\infty}$ is large
- Use upper bounds on VC-dimension of ANN's

**Theorem:**

- Fixed architecture for approximating $W^{n,\infty}([0,1]^d)$ up to $\varepsilon$ must have at least $\omega \geq \varepsilon^{-d/(2n)}$ weights.
- If, moreover, $L \leq c \ln^p(1/\varepsilon)$, then $\omega \geq \varepsilon^{-d/n} \ln^{-2p-1}(1/\varepsilon)$.

# Riesz basis of neural networks

# Multivariate Riesz basis of ReLU neural networks

**Question:** **Can we find a (nearly) orthonormal system, which is easily reproducible by** ReLU-**neural networks?**

- Studied for $d = 1$ in [DDFHP22]:

  📄 I. Daubechies, R. DeVore, S. Foucart, B. Hanin, and G. Petrova.
  *Nonlinear Approximation and (Deep) ReLU Networks.*
  Constr. Appr. 55:127–172 (2022).

  --→ They produced a system of piece-wise linear functions, which resembles the trigonometric Fourier system

## $d = 1$: Univariate Riesz basis

### Definition

- For $x \in [0,1]$, we define
$$\mathcal{C}(x) = 4\left|x - \frac{1}{2}\right| - 1 = \begin{cases} 1 - 4x, & x \in [0, 1/2), \\ 4x - 3, & x \in [1/2, 1] \end{cases}$$

  and
$$\mathcal{S}(x) = \left|2 - 4\left|x - \frac{1}{4}\right|\right| - 1 = \begin{cases} 4x, & x \in [0, 1/4), \\ 2 - 4x, & x \in [1/4, 3/4), \\ 4x - 4, & x \in [3/4, 1]. \end{cases}$$

- For $x \in \mathbb{R}$, we extend this definition periodically, i.e.
$$\mathcal{C}(x) = \mathcal{C}(x - \lfloor x \rfloor) \quad \text{and} \quad \mathcal{S}(x) = \mathcal{S}(x - \lfloor x \rfloor).$$

- If $k \geq 1$ and $x \in \mathbb{R}$, we put $\mathcal{C}_k(x) = \mathcal{C}(kx)$ and $\mathcal{S}_k(x) = \mathcal{S}(kx)$.

## Univariate Riesz basis



$\dashrightarrow$ $\mathcal{C}$, $\mathcal{S}$ emulate/resemble trigonometric functions

## Univariate Riesz basis

$\mathcal{C}_k(x)$ and $\mathcal{S}_k(x)$ can be easily reproduced by the ReLU ANN's:

### Lemma

For $k \in \mathbb{N}$ it holds $\mathcal{C}_k \in \Upsilon^{2,\lceil \log_2 k \rceil + 1}$ and $\mathcal{S}_k \in \Upsilon^{2,\lceil \log_2 k \rceil + 2}$

Proof :
$$H \in \Upsilon^{2,1} \implies \mathcal{C} = 1 - 2H \in \Upsilon^{2,1}$$
$$\implies H^{\circ m} := H \circ H \circ \cdots \circ H \in \Upsilon^{2,m}$$

which implies $\mathcal{C}_{2^m}(x) = \mathcal{C}(H^{\circ m}(x)) \in \Upsilon^{2,m+1}$

For arbitrary $k \leq 2^m$ with $m := \lceil \log_2 k \rceil$ we derive

$$\mathcal{C}_k(x) = \mathcal{C}_{2^m}(k2^{-m}x) \in \Upsilon^{2,m+1} = \Upsilon^{2,\lceil \log_2 k \rceil + 1}.$$

For $\mathcal{S}$ use identity $\mathcal{S}(x) = \mathcal{C}_2(\frac{x}{2} + \frac{3}{8}) \in \Upsilon^{2,2}$.

## Univariate Riesz basis

### Definition

Let $H$ be a real Hilbert space. $(x_n)_n \subset H$ is a **Riesz sequence** if $\exists A, B > 0$:

$$A \sum_n \alpha_n^2 \leq \left\| \sum_n \alpha_n x_n \right\|^2 \leq B \sum_n \alpha_n^2 \qquad (\blacktriangle)$$

for every $(\alpha_n)_n \in \ell_2$.

$(x_n)_n$ is called a **Riesz basis** if the closed span of $(x_n)_n$ is $H$.

## Univariate Riesz basis

### Definition

Let $H$ be a real Hilbert space. $(x_n)_n \subset H$ is a **Riesz sequence** if $\exists A, B > 0$:

$$A \sum_n \alpha_n^2 \leq \left\| \sum_n \alpha_n x_n \right\|^2 \leq B \sum_n \alpha_n^2 \qquad (\blacktriangle)$$

for every $(\alpha_n)_n \in \ell_2$.

$(x_n)_n$ is called a **Riesz basis** if the closed span of $(x_n)_n$ is $H$.

### Theorem (DDFHP22)

*The system*
$$\mathcal{R}_1 := \{1\} \cup \{\sqrt{3}\mathcal{C}_k, \sqrt{3}\mathcal{S}_k : k \in \mathbb{N}\}$$

*is a Riesz basis of* $L_2([0,1])$.

*A and B can be chosen as* $A = 1/2$ *and* $B = 3/2$.

## Univariate Riesz basis

We provide an alternative proof based on:

- Properties of the Gram matrix $(\langle \mathcal{C}_k, \mathcal{C}_l \rangle)_{k,l}$

- Gershgorin's theorem

- Further Tools:
  - Trigonometric series
  - Euler products
  - Ramanujan's formula

📄 C. Schneider and J.Vybíral.

Multivariate Riesz basis of ReLU neural networks.

Appl. Comput. Harmonic Anal. 68 (2024), 101605

## Sketch of the proof:

**Step 1: Reformulate the definition of a (finite) Riesz sequence as an eigenvalue problem of its Gram matrix**

$H$ - a real Hilbert space; $\{x_i\}_{i=1}^{N} \subset H$; $\alpha = (\alpha_1, \ldots, \alpha_N)^T \in \mathbb{R}^N$

$$\left\| \sum_{i=1}^{N} \alpha_i x_i \right\|^2 = \sum_{i,j=1}^{N} \alpha_i \alpha_j \langle x_i, x_j \rangle = \alpha^T G \alpha,$$

where $G = (g_{i,j})_{i,j=1}^{N}$ with $g_{i,j} = \langle x_i, x_j \rangle$ is the Gram matrix of $\{x_i\}_{i=1}^{N}$.

## Sketch of the proof:

**Step 1: Reformulate the definition of a (finite) Riesz sequence as an eigenvalue problem of its Gram matrix**

$H$ - a real Hilbert space; $\{x_i\}_{i=1}^{N} \subset H$; $\alpha = (\alpha_1, \ldots, \alpha_N)^T \in \mathbb{R}^N$

$$\left\| \sum_{i=1}^{N} \alpha_i x_i \right\|^2 = \sum_{i,j=1}^{N} \alpha_i \alpha_j \langle x_i, x_j \rangle = \alpha^T G \alpha,$$

where $G = (g_{i,j})_{i,j=1}^{N}$ with $g_{i,j} = \langle x_i, x_j \rangle$ is the Gram matrix of $\{x_i\}_{i=1}^{N}$.

Therefore, (▲) is equivalent to

$$A \alpha^T \alpha \le \alpha^T G \alpha \le B \alpha^T \alpha \quad \text{for every} \quad \alpha \in \mathbb{R}^N,$$

## Sketch of the proof:

**Step 2: Gershgorin circle theorem**

$$\sigma(G) \subset \bigcup_{i=1}^{N}\Big[g_{i,i} - \sum_{j \neq i}|g_{i,j}|, g_{i,i} + \sum_{j \neq i}|g_{i,j}|\Big].$$

## Sketch of the proof:

**Step 2: Gershgorin circle theorem**
$$\sigma(G) \subset \bigcup_{i=1}^{N} \Big[ g_{i,i} - \sum_{j \neq i} |g_{i,j}|, g_{i,i} + \sum_{j \neq i} |g_{i,j}| \Big].$$

Lemma: Let $i, j \in \mathbb{N}$. Then

- $\langle \mathcal{C}_i, \mathcal{S}_j \rangle = 0$;

- $\langle \mathcal{C}_i, \mathcal{C}_j \rangle = \langle \mathcal{S}_i, \mathcal{S}_j \rangle = 0$ if $i/\gcd(i,j)$ is odd and $j/\gcd(i,j)$ is even (or vice versa), i.e., if the prime factorizations of $i$ and $j$ contain a different power of 2;

- If $i/\gcd(i,j)$ and $j/\gcd(i,j)$ are both odd, then
$$3 \cdot \langle \mathcal{C}_i, \mathcal{C}_j \rangle = 3 \cdot |\langle \mathcal{S}_i, \mathcal{S}_j \rangle| = \frac{\gcd(i,j)^4}{i^2 \cdot j^2}.$$
The sign of $\langle \mathcal{S}_i, \mathcal{S}_j \rangle$ is negative if, $(i+j)/(2\gcd(i,j))$ is even.

- In particular, we get $\langle \mathcal{C}_i, \mathcal{C}_i \rangle = \langle \mathcal{S}_i, \mathcal{S}_i \rangle = 1/3$ for all $i \in \mathbb{N}$.

## Sketch of the proof:

**Step 3: Proof of the Lemma:**

By Fourier series: $c_k(x) = \sqrt{2}\cos(2\pi k x)$

$$\sqrt{3}\,\mathcal{C}_k = \mu \sum_{m \geq 0} \frac{1}{(2m+1)^2} c_{(2m+1)k}, \quad \mu^2 \frac{\pi^4}{96} = 1$$

Then

$$3\langle \mathcal{C}_i, \mathcal{C}_j \rangle = \sum_{m,n=0}^{\infty} \frac{\mu^2}{(2m+1)^2(2n+1)^2} \delta_{(2m+1)i,(2n+1)j},$$

Solve $(2m+1)i = (2n+1)j$: $g = \gcd(i,j)$ and

$$2m+1 = \frac{j}{g} \cdot (2l+1), \quad 2n+1 = \frac{i}{g} \cdot (2l+1), \quad l \in \mathbb{N}_0.$$

## Sketch of the proof:

**Step 4:** For $i$ odd, estimate $\sum_{j \text{ odd}} \langle \mathcal{C}_i, \mathcal{C}_j \rangle$

$i = q_1^{\alpha_1} \ldots q_n^{\alpha_n}$, primes $q_1, \ldots, q_n \geq 3$, $\alpha_1, \ldots, \alpha_n \geq 1$

$j = q_1^{\beta_1} \ldots q_n^{\beta_n} \cdot J$, with $\beta_1, \ldots, \beta_n \geq 0$ and $J$ odd with $\gcd(J, i) = 1$

$\gcd(i, j) = \prod_{u=1}^{n} q_u^{\min(\alpha_u, \beta_u)}$

Then

$$\sum_{j \in \mathbb{N}, \, j \text{ odd}} 3 \cdot \langle \mathcal{C}_i, \mathcal{C}_j \rangle = \cdots \leq \prod_{p \geq 3 : p \text{ prime}} \frac{1 + 1/p^2}{1 - 1/p^2} =$$

## Sketch of the proof:

**Step 4:** For $i$ odd, estimate $\sum_{j \text{ odd}} \langle \mathcal{C}_i, \mathcal{C}_j \rangle$

$i = q_1^{\alpha_1} \ldots q_n^{\alpha_n}$, primes $q_1, \ldots, q_n \geq 3$, $\alpha_1, \ldots, \alpha_n \geq 1$

$j = q_1^{\beta_1} \ldots q_n^{\beta_n} \cdot J$, with $\beta_1, \ldots, \beta_n \geq 0$ and $J$ odd with $\gcd(J, i) = 1$

$\gcd(i, j) = \prod_{u=1}^{n} q_u^{\min(\alpha_u, \beta_u)}$

Then

$$\sum_{j \in \mathbb{N}, \, j \text{ odd}} 3 \cdot \langle \mathcal{C}_i, \mathcal{C}_j \rangle = \cdots \leq \prod_{p \geq 3 : p \text{ prime}} \frac{1 + 1/p^2}{1 - 1/p^2} = \frac{3}{2}$$

$\implies \sigma(G) \subset [1/2, 3/2]$.

## $d = 1$: ReLU neural networks

One can reproduce linear combinations of $\mathcal{C}_k$ and $\mathcal{S}_k$ via ReLU
networks with good control on the depth $L$:

### Theorem (DDFHP22, Thm. 6.2)

*Let $W \geq 6$. For every $k \geq 1$, and set of indices $\Lambda \subset \mathbb{N}$ with $|\Lambda| = k$,
the set*

$$\mathcal{F}_\Lambda := \left\{ \sum_{j \in \Lambda} (a_j \mathcal{C}_j + b_j \mathcal{S}_j), \ a_j, b_j \in \mathbb{R}, \ j \in \Lambda, \ |\Lambda| = k \right\} \subset \Upsilon^{W,L}$$

$$L = 2 \left\lceil \frac{k}{\lfloor \frac{W-2}{4} \rfloor} \right\rceil (\lceil \log_2(\lambda) \rceil + 2) \qquad \text{with} \quad \lambda := \max\{j : \ j \in \Lambda\}.$$

$\dashrightarrow$ deep NN approximation ($W$ is fixed)

# $d > 1$: Multivariate Riesz basis

**Good news:** A tensor product of two Riesz sequences is a Riesz sequence

## $d > 1$: Multivariate Riesz basis

**Good news:** A tensor product of two Riesz sequences is a Riesz sequence

**Bad news:** The Riesz constants get multiplied!  ... $A^d, B^d$

$\dashrightarrow$ bad dependence on $d$

# $d > 1$: Multivariate Riesz basis

**Good news:** A tensor product of two Riesz sequences is a Riesz sequence

**Bad news:** The Riesz constants get multiplied!   ... $A^d, B^d$

$--\rightarrow$ bad dependence on $d$

**Another bad news:** For ReLU neural networks it is rather complicated to calculate $x \cdot y$

## $d > 1$: Multivariate Riesz basis

**Good news:** A tensor product of two Riesz sequences is a Riesz sequence

**Bad news:** The Riesz constants get multiplied! ... $A^d, B^d$

$\dashrightarrow$ bad dependence on $d$

**Another bad news:** For ReLU neural networks it is rather complicated to calculate $x \cdot y$

**Way out!** If $\alpha = (\alpha_1, \ldots, \alpha_d) \in \mathbb{Z}^d$, we say that $\alpha \rhd 0$ if the first non-zero index of $\alpha$ is positive.

## $d > 1$: Multivariate Riesz basis

**Good news:** A tensor product of two Riesz sequences is a Riesz sequence

**Bad news:** The Riesz constants get multiplied!   ... $A^d, B^d$

$\dashrightarrow$ bad dependence on $d$

**Another bad news:** For ReLU neural networks it is rather complicated to calculate $x \cdot y$

**Way out!** If $\alpha = (\alpha_1, \ldots, \alpha_d) \in \mathbb{Z}^d$, we say that $\alpha \rhd 0$ if the first non-zero index of $\alpha$ is positive.

Theorem (SV24): The system

$$\mathcal{R}_d := \{1\} \cup \{\sqrt{3}\mathcal{C}(\alpha \cdot x) : \alpha \rhd 0\} \cup \{\sqrt{3}\mathcal{S}(\alpha \cdot x) : \alpha \rhd 0\}$$

is a Riesz basis of $L_2([0,1]^d)$. The constants can be chosen as $A = 1/2$ and $B = 3/2$, independently of $d$!

## $d > 1$: ReLU neural networks

One can reproduce the multivariate Riesz basis via ReLU networks:

Theorem (SV24):

- For $j \in \mathbb{N}$ it holds on $[0, 1]$,

$$\mathcal{C}_j \in \Upsilon^{2, \lceil \log_2 j \rceil + 1} \quad \text{and} \quad \mathcal{S}_j \in \Upsilon^{2, \lceil \log_2 j \rceil + 2}$$

  Entries of weight matrices and the bias vectors are bounded by 8.

## $d > 1$: ReLU neural networks

One can reproduce the multivariate Riesz basis via ReLU networks:

Theorem (SV24):

- For $j \in \mathbb{N}$ it holds on $[0, 1]$,

$$\mathcal{C}_j \in \Upsilon^{2, \lceil \log_2 j \rceil + 1} \quad \text{and} \quad \mathcal{S}_j \in \Upsilon^{2, \lceil \log_2 j \rceil + 2}$$

  Entries of weight matrices and the bias vectors are bounded by 8.

- Let $d > 1$, $\alpha \in \mathbb{Z}^d \setminus \{0\}$. Then on $[0, 1]^d$

$$\mathcal{C}(\alpha \cdot x) \in \Upsilon^{2, \lceil \log_2 \|\alpha\|_1 \rceil + 2} \quad \text{and} \quad \mathcal{S}(\alpha \cdot x) \in \Upsilon^{2, \lceil \log_2 \|\alpha\|_1 \rceil + 3},$$

  ($\|\alpha\|_1 = |\alpha_1| + \ldots + |\alpha_d|$, weights and biases are bounded by 8)

## $d > 1$: ReLU neural networks

One can reproduce the multivariate Riesz basis via ReLU networks:

Theorem (SV24):

- For $j \in \mathbb{N}$ it holds on $[0,1]$,

$$\mathcal{C}_j \in \Upsilon^{2,\lceil \log_2 j \rceil + 1} \quad \text{and} \quad \mathcal{S}_j \in \Upsilon^{2,\lceil \log_2 j \rceil + 2}$$

  Entries of weight matrices and the bias vectors are bounded by 8.

- Let $d > 1$, $\alpha \in \mathbb{Z}^d \setminus \{0\}$. Then on $[0,1]^d$

$$\mathcal{C}(\alpha \cdot x) \in \Upsilon^{2,\lceil \log_2 \|\alpha\|_1 \rceil + 2} \quad \text{and} \quad \mathcal{S}(\alpha \cdot x) \in \Upsilon^{2,\lceil \log_2 \|\alpha\|_1 \rceil + 3},$$

  ($\|\alpha\|_1 = |\alpha_1| + \ldots + |\alpha_d|$, weights and biases are bounded by 8)

- Let $d \geq 1$, $\{\alpha_1, \ldots, \alpha_k, \beta_1, \ldots, \beta_l\} \subset \mathbb{Z}^d \setminus \{0\}$. Then on $[0,1]^d$

$$f(x) = \sum_{i=1}^{k} a_i \mathcal{C}(\alpha_i \cdot x) + \sum_{j=1}^{l} b_j \mathcal{S}(\beta_j \cdot x) \in \Upsilon^{W,L}, \quad x \in [0,1]^d$$

  $W = 2(k+l)$ and $L = 2 + \max_{\substack{i=1,\ldots,k; \\ j=1,\ldots,l}} \{\lceil \log_2(\|\alpha_i\|_1) \rceil, \lceil \log_2(\|\beta_j\|_1) \rceil + 1\}$

  (weights and biases bounded by $\max_{\substack{i=1,\ldots,k; \\ j=1,\ldots,l}} \{8|a_i|, 8|b_j|, 8\}$)

# Open Problems

- Open Problem 1

  Are the constants $A = \frac{1}{2}$, $B = \frac{3}{2}$ optimal?     *(probably not)*

# Open Problems

- Open Problem 1

   Are the constants $A = \frac{1}{2}$, $B = \frac{3}{2}$ optimal?   *(probably not)*

- Open Problem 2

   Is it even possible to obtain a tight frame in the sense that $A = B$?

# Open Problems

- Open Problem 1

  Are the constants $A = \frac{1}{2}$, $B = \frac{3}{2}$ optimal?     *(probably not)*

- Open Problem 2

  Is it even possible to obtain a tight frame in the sense that
  $A = B$?

- Open Problem 3

  Does a frame exist which is better than the Riesz basis?

## Applications to approximations of functions

- Previous results on function recovery (often) used localization - decomposition of unity on a tensor grid

- Implicit constants (often) depend exponentially on dimension $d$

- We apply the (non-local) Riesz basis!

- We decompose $f$ into this basis, keep largest elements and recover them as ANN's

## Univariate setting

Let $s \geq 0$. Then $W^s([0,1])$ is the set of 1-periodic functions

$$f(x) = a_0 + \sum_{m=1}^{\infty} a_m \cos(2\pi m x) + b_m \sin(2\pi m x), \quad x \in [0,1]$$

with $\|f\|_{W^s}^2 := a_0^2 + \sum_{m=1}^{\infty} m^{2s}(a_m^2 + b_m^2) < \infty$.

## Univariate setting

Let $s \geq 0$. Then $W^s([0,1])$ is the set of 1-periodic functions

$$f(x) = a_0 + \sum_{m=1}^{\infty} a_m \cos(2\pi m x) + b_m \sin(2\pi m x), \quad x \in [0,1]$$

with $\|f\|_{W^s}^2 := a_0^2 + \sum_{m=1}^{\infty} m^{2s}(a_m^2 + b_m^2) < \infty$.

And $\mathcal{F}^s([0,1])$ is the set of 1-periodic functions

$$f(x) = \alpha_0 + \sum_{k=1}^{\infty} \alpha_k \mathcal{C}_k(x) + \beta_k \mathcal{S}_k(x), \quad x \in [0,1],$$

with $\|f\|_{\mathcal{F}^s}^2 := \alpha_0^2 + \sum_{m=1}^{\infty} m^{2s}(\alpha_m^2 + \beta_m^2) < \infty$.

**Theorem:** Let $0 \leq s < 1$. Then, $W^s([0,1]) = \mathcal{F}^s([0,1])$.

## Multivariate version

For $d \geq 2$ and $s \geq 0$, we proceed in the same way!

For

$$f(x) = \sum_{m \in \mathbb{Z}^d} a_m e^{2\pi i m \cdot x}, \quad x \in [0,1]^d$$

we put $\|f\|_{W^s}^2 = |a_0|^2 + \sum_{m \in \mathbb{Z}^d \setminus \{0\}} \|m\|_2^{2s} \cdot |a_m|^2 < \infty.$

## Multivariate version

For $d \geq 2$ and $s \geq 0$, we proceed in the same way!

For
$$f(x) = \sum_{m \in \mathbb{Z}^d} a_m e^{2\pi i m \cdot x}, \quad x \in [0,1]^d$$

we put $\|f\|_{W^s}^2 = |a_0|^2 + \sum_{m \in \mathbb{Z}^d \setminus \{0\}} \|m\|_2^{2s} \cdot |a_m|^2 < \infty$.

And for
$$f(x) = \alpha_0 + \sum_{k \,\dot{>}\, 0} \Big[ \alpha_k \mathcal{C}_k(x) + \beta_k \mathcal{S}_k(x) \Big], \quad x \in [0,1]^d.$$

we define $\|f\|_{\mathcal{F}^s}^2 := \alpha_0^2 + \sum_{k \,\dot{>}\, 0} \|k\|_2^{2s} (\alpha_k^2 + \beta_k^2) < \infty$

## Multivariate version

For $d \geq 2$ and $s \geq 0$, we proceed in the same way!

For

$$f(x) = \sum_{m \in \mathbb{Z}^d} a_m e^{2\pi i m \cdot x}, \quad x \in [0,1]^d$$

we put $\|f\|_{W^s}^2 = |a_0|^2 + \sum_{m \in \mathbb{Z}^d \setminus \{0\}} \|m\|_2^{2s} \cdot |a_m|^2 < \infty.$

And for

$$f(x) = \alpha_0 + \sum_{k \vartriangleright 0} \Big[ \alpha_k \mathcal{C}_k(x) + \beta_k \mathcal{S}_k(x) \Big], \quad x \in [0,1]^d.$$

we define $\|f\|_{\mathcal{F}^s}^2 := \alpha_0^2 + \sum_{k \vartriangleright 0} \|k\|_2^{2s} (\alpha_k^2 + \beta_k^2) < \infty$

**Theorem:** Let $0 \leq s < 1$. Then, $W^s([0,1]^d) = \mathcal{F}^s([0,1]^d)$ in the sense of equivalent norms. The constants of equivalence of these norms depend on $s$ but are independent of $d$.

# Recovery by ANN's

- Let $f \in W^s([0,1]^d)$. Then $f \in \mathcal{F}^s([0,1]^d)$
  (with equivalent norms)

- It can be decomposed into the Riesz basis of $\mathcal{C}_k$ and $\mathcal{S}_k$

- Fix $R > 0$ real and split

$$f = f_R + f^R = \Big(\alpha_0 + \sum_{k \rhd 0 : \|k\|_2 \leq R} \cdots \Big) + \Big( \sum_{k \rhd 0 : \|k\|_2 > R} \cdots \Big).$$

- Recover $f_R$ exactly(!) by an ANN; the error is given by $f^R$

- We need to:
    - Estimate the size of the sum in $f_R$
    - Estiamte the norm of $f^R$

## Gauss circle problem

Let $d \geq 1$ and $t \geq 0$. Then
$$N(t, d) := \{k \in \mathbb{Z}^d : \|k\|_2 \leq t\}$$
is the number of integer lattice points in a ball with radius $t \geq 0$.

**Lemma:** There exist two absolute constants $c_1, c_2 > 0$ such that for every $d \geq 1$ and $t \geq 0$ the following holds.
(i) If $t \geq \sqrt{d}/2$, then
$$N(t, d) \leq \left(\frac{c_1 t}{\sqrt{d}}\right)^d.$$

(ii) If $0 < t \leq \sqrt{d}/2$, then
$$N(t, d) \leq \left(\frac{c_2 d}{t^2}\right)^{t^2}.$$

## Approximation by ANN's from the Riesz basis

**Theorem:**

Let $0 < s < 1$ and $0 < \varepsilon < 1$. Let $R := (C_s/\varepsilon)^{1/s}$. Then, for every $f \in W^s([0,1]^d)$ there is an ANN $\mathcal{N} \in \Upsilon_d^{W,L}$ with

$$W = 4 \cdot N(R,d) \quad \text{and} \quad L \leq 4 + \log_2\left(R \cdot \sqrt{\min(R,d)}\right)$$

such that

$$\|f - \mathcal{N}\|_2 \leq \varepsilon \cdot \|f\|_{W^s}.$$

**Remarks:**

- The architecture does not depend on $f$

- $C_s$ is independent on $d$

- If $\varepsilon > 0$ is fixed and $d \to \infty$, then $L$ is bounded and $W$ grows polynomially in $d \implies$ we avoid the *curse of dimensionality*

## Barron spaces

Let $s \geq 0$ and define Fourier-analytic Barron spaces

$$\mathbb{B}^s([0,1]^d) : \|f\|_{\mathbb{B}^s} := |a_0| + \sum_{m \in \mathbb{Z}^d \setminus \{0\}} \|m\|_2^s \cdot |a_m|$$

and Barron spaces with respect to $\mathcal{C}_k$ and $\mathcal{S}_k$

$$\mathcal{B}^s([0,1]^d) : \|f\|_{\mathcal{B}^s} := |\alpha_0| + \sum_{k \,\hat{\gt}\, 0} \|k\|_2^s \cdot (|\alpha_k| + |\beta_k|).$$

**Theorem:** Let $0 \leq s < 1$. Then $\mathbb{B}^s([0,1]^d) = \mathcal{B}^s([0,1]^d)$ in the sense of equivalent norms. (The constants of equivalence of these norms depend on $s$ but are independent of $d$.)

## Barron spaces

We start again with

$$f(x) = \alpha_0 + \sum_{k \mathrel{\dot{>}} 0} \Big[ \alpha_k \mathcal{C}_k(x) + \beta_k \mathcal{S}_k(x) \Big], \quad x \in [0,1]^d$$

with

$$\|\alpha\|_{b_s^d} = \sum_{k \mathrel{\dot{>}} 0} \|k\|_2^s \cdot |\alpha_k| < \infty$$

(and same for $\beta$)

**Best $n$-term approximation**: We keep the largest terms, so that the remainder is small in $\ell_2$

**Remarks:**

- Therefore, the architecture depends on $f$
- Again, we can avoid the curse of dimensionality.

## Barron spaces

**Theorem:** Let $0 < s < 1$ and $0 < \varepsilon < 1$. Let $R := (C/\varepsilon)^{1/s}$ and $n$ such that $\sigma_n(b_s^d, \ell_2) < c\varepsilon$. Then, for every $f \in \mathbb{B}^s([0,1]^d)$ there is an ANN $\mathcal{N} \in \Upsilon_d^{W,L}$ with

$$W = 4n \quad \text{and} \quad L \le 4 + \log_2\left(R \cdot \sqrt{\min(R,d)}\right)$$

such that

$$\|f - \mathcal{N}\|_2 \le \varepsilon \cdot \|f\|_{\mathbb{B}^s}.$$

## Literature

- M. Telgarsky, *Benefits of depth in neural networks*, Conference on learning theory, PMLR, 2016

- D. Yarotsky, *Error bounds for approximations with deep ReLU networks*, Neural networks 94 (2017), 103–114

- D. Yarotsky, *Optimal approximation of continuous functions by very deep ReLU networks*, Conference on learning theory, PMLR, 2018

- P. Bartlett, V. Maiorov, and R. Meir, *Almost linear VC dimension bounds for piecewise polynomial networks*, Advances in neural information processing systems 11 (1998)

- A. R. Barron, *Universal approximation bounds for superpositions of a sigmoidal function*, IEEE Transactions on Information Theory 39.3 (1993), 930–945

## ChatGPT vs. DeepSeek

Forty percent of families in a city have one child, thirty percent of families have two children, twenty percent of families have three children, and ten percent of families have four children.
What is the probability that a randomly selected child has a sister?

ChatGPT: 38.75%
Are you sure?: 38.75%
You are wrong!: 38.75%
Think again: 110%, 50%, 45%, 50,83%, 80%, 65%, 41,9%

DeepSeek: 60%
Are you sure?: 60%
You are wrong!: 60%

# Thank you for your attention!