

Approximation Theory of Neural Networks

⊙ Introduction

⊙ An example

⊙ Review of two classical approximation thems (Stone-Weierstrass, Kolmogorov Superposition Thm)

⊙ Universal Approximation Thm.

⊙ Error bounds of NN approximations

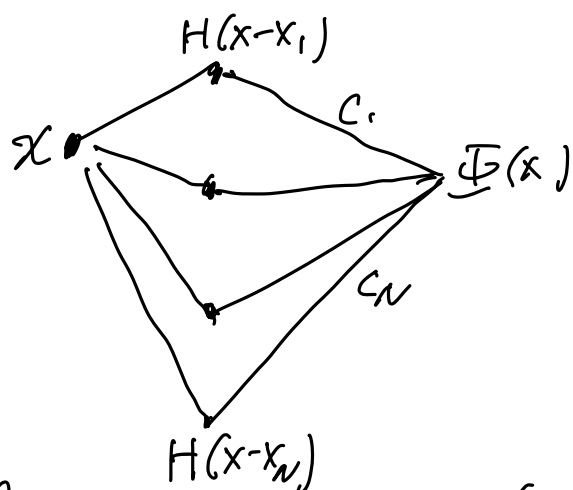
⊙ Other topics

An example Given $g \in C([0,1])$. Construct a simple NN, $\Phi \in \mathcal{NN}(H, L=2, N_0=1, N_L=1)$, to approximate g . Here, $H(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$ — the Heaviside function. ($H(0)=1$)

Thm $\forall \varepsilon > 0, \exists N \in \mathcal{N} (N \gg 1), \exists 0 \leq x_1 < \dots < x_N \leq 1,$
 $\exists c_1, \dots, c_N \in \mathbb{R}, \text{ s.t. } |\Phi_N(x) - g(x)| < \varepsilon \quad \forall x \in [0,1],$

where

$$\Phi_N(x) = \Phi_N(x) = \sum_{j=1}^N c_j \cdot H(x - x_j) \in \mathcal{NN}(H, 2, 1, 1), \quad x \in [0,1].$$



$-x_j = \theta_j$: threshold / bias.

c_j — weights in the output

$$A_1 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_N, \quad b_1 = \begin{bmatrix} -x_1 \\ \vdots \\ -x_N \end{bmatrix}.$$

$$A_2 = [c_1 \dots c_N]_N, \quad b_2 = 0 \in \mathbb{R}.$$

Proof Since $g \in C([0,1])$ is unif. cont. on $[0,1]$, $\exists \delta > 0$ s.t. $|x-y| < \delta \implies |g(x) - g(y)| < \varepsilon$. Choose $N \in \mathcal{N}$ s.t. $N > \frac{1}{\delta}$.

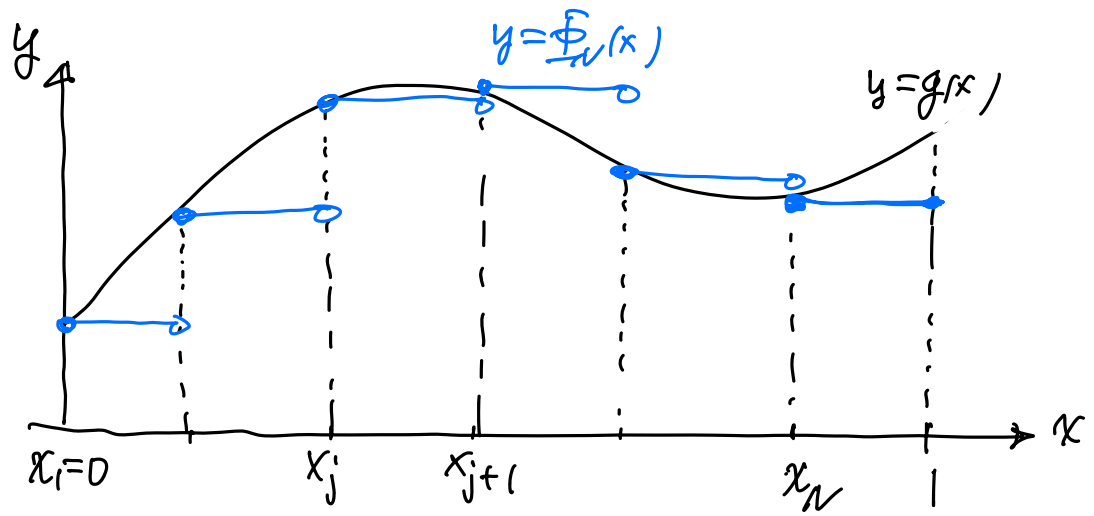
Define $x_j = (j-1)/N, j=1, \dots, N$. On each $[x_j, x_{j+1}]$,

$$|g(x) - g(x_j)| < \epsilon \quad \forall x \in [x_j, x_{j+1}].$$

Define $c_1 = g(x_1), c_j = g(x_j) - g(x_{j-1}), j=2, \dots, N$. Then,

$$\begin{aligned} \Phi_N(x) &= \sum_{j=1}^N c_j H(x-x_j) \\ &= \begin{cases} g(x_1) & \text{if } x \in [x_1, x_2), \\ g(x_2) & \text{if } x \in [x_2, x_3), \\ \dots & \dots \\ g(x_{j+1}) & \text{if } x \in [x_j, x_{j+1}), \\ \dots & \dots \\ g(x_N) & \text{if } x \in [x_N, 1]. \end{cases} \end{aligned}$$

Then $|\Phi_N(x) - g(x)| < \epsilon, \forall x \in [0, 1]$. Q.E.D



The First Weierstrass Approximation Theorem. Let $a, b \in \mathbb{R}$ with $a < b$. Let $f \in C([a, b])$ and $\epsilon > 0$. Then there exists a polynomial p such that

$$\max_{a \leq x \leq b} |f(x) - p(x)| < \epsilon. \quad (*)$$

Remarks

○ Define for $g \in C([a, b])$,

$$\|g\| = \|g\|_{cb} = \|g\|_{d([a, b])} = \|g\|_c = \max_{a \leq x \leq b} |g(x)|.$$

Then $(C([a, b]), \|\cdot\|)$ is a Banach space. The inequality (*) is the same as $\|f - p\| < \epsilon$.

3

⊙ Equivalently, $\lim_{n \rightarrow \infty} E_n(f) = 0$, where
 $E_n(f) = \min_{g \in P_n} \|f - g\|$, (min is attained) uniquely
 $P_n = \{ \text{all polynomials of deg} \leq n \}$.

⊙ The theorem provides no info about ρ .

The Second Weierstrass Approximation Theorem Let

$f \in C_{2\pi}$ and $\varepsilon > 0$ then there exists a trigonometric polynomial T such that
 $\max_{-\pi \leq x \leq \pi} |f(x) - T(x)| < \varepsilon$.

The Stone-Weierstrass Theorem Let X be a compact Hausdorff space and $C(X, \mathbb{R})$ or $C(X)$ the space of all real continuous functions equipped with the uniform norm. Let $\mathcal{A} \subseteq C(X, \mathbb{R})$. Assume

⊙ \mathcal{A} is a subalgebra;

⊙ \mathcal{A} separates points of X ; and

⊙ \mathcal{A} contains the constant functions.

Then $\overline{\mathcal{A}} = C(X, \mathbb{R})$ (i.e., \mathcal{A} is dense in $C(X, \mathbb{R})$).

Remarks

⊙ \mathcal{A} is a subalgebra means: \mathcal{A} is a vector subspace of $C(X, \mathbb{R})$ and \mathcal{A} is closed in multiplication: $f, g \in \mathcal{A} \Rightarrow fg \in \mathcal{A}$; $f \in \mathcal{A}, \alpha \in \mathbb{R} \Rightarrow \alpha f \in \mathcal{A}$; $f, g \in \mathcal{A} \Rightarrow f+g \in \mathcal{A}$.

⊙ \mathcal{A} separates points of X : $\forall x, y \in X, x \neq y$,
 $\exists f \in \mathcal{A}$ s.t. $f(x) \neq f(y)$.

⊙ $\overline{\mathcal{A}}$ is the closure of \mathcal{A} w.r.t. the uniform

norm. $\overline{\mathcal{A}} = C(X, \mathbb{R})$ means that $\forall f \in C(X, \mathbb{R})$ 4
 there exist $g_k \in \mathcal{A}$ ($k=1, 2, \dots$) such that $g_k \rightarrow f$
 as $k \rightarrow \infty$. i.e.,

$$\|f - g_k\| = \max_{x \in X} |f(x) - g_k(x)| \rightarrow 0 \text{ as } k \rightarrow \infty.$$

① You may consider X to be a compact subset of \mathbb{R}^n . Compact here = closed + bounded.

Example $X = \prod_{j=1}^n [a_j, b_j]$. ($a_j, b_j \in \mathbb{R}$, $a_j < b_j$, $j=1, \dots, n$)

or $X = \overline{\text{co}}(X_1, \dots, X_m)$ the convex hull of m

points in \mathbb{R}^n ($X_1, \dots, X_m \in \mathbb{R}^n$) $\mathcal{A} = \{ \text{all real multivariable polynomials } p = p(x_1, \dots, x_n) \}$.

Then $\overline{\mathcal{A}} = C(X, \mathbb{R})$. i.e., $\forall f \in C(X)$ $\forall \varepsilon > 0$, there exists a polynomial $p: \mathbb{R}^n \rightarrow \mathbb{R}$ s.t.

$$|f(x) - p(x)| < \varepsilon \quad \forall x \in X.$$

The classical Weierstrass Thm.

Corollary. Let $X = \prod_{j=1}^n [a_j, b_j]$ with all $a_j, b_j \in \mathbb{R}$

and $a_j < b_j$ ($j=1, \dots, n$). Let $f \in C(X)$ and $\varepsilon > 0$.

Then, $\exists N \in \mathbb{N}$, and $f_{ij} \in C([a_j, b_j])$ ($i=1, \dots, N$; $j=1, \dots, n$) such that

$$\max_{x=(x_1, \dots, x_n) \in X} \left| f(x) - \sum_{i=1}^N \prod_{j=1}^n f_{ij}(x_j) \right| < \varepsilon.$$

Proof Let $\mathcal{A} = \left\{ \sum_{i=1}^N \prod_{j=1}^n g_{ij}(x_j) : g_{ij} \in C([a_j, b_j]), \right.$

$\left. j=1, \dots, n, i=1, \dots, N, N=1, 2, \dots \right\}$.

Apply the Stone-Weierstrass Thm. Q.E.D

(Fill in details!)

Remark. We can choose $f_{ij} \in \mathcal{P}$ (the set of all 5
one-variable real polynomials).

The Kolmogorov Theorem (or Kolmogorov-Arnold Theorem)

Kolmogorov, 1956, Arnold 1957, Kolmogorov 1957.

Theorem (Kolmogorov 1957) For any integer $n \geq 2$ and
and $f \in C([0,1]^n)$ there exist $\psi^{p,q} \in C([0,1])$ ($p=1, \dots, n$,
 $q=1, \dots, 2n+1$) and $\chi_q \in C(\mathbb{R})$ ($q=1, \dots, 2n+1$) such that
$$f(x_1, \dots, x_n) = \sum_{q=1}^{2n+1} \chi_q \left(\sum_{p=1}^n \psi^{p,q}(x_p) \right) \quad \forall x = (x_1, \dots, x_n) \in [0,1]^n.$$

Remarks

⊙ Originated from Hilbert's 13th prob. (1900).

⊙ Constructive proof is given by Braun and Griebel (Constr. Approx. 2009), after Sprecher (1996, 1997) and Köppen (2002).

Universal Approximations

Question: Given a function/map: $f: X \rightarrow \mathbb{R}^m$,
are there NNWs Φ_k ($k=1, 2, \dots$) such that
 $\Phi_k \rightarrow f$? Here $X \subseteq \mathbb{R}^d$ and the convergence
is w.r.t. some norm or metric or some
topology.

More questions

⊙ what is the class of functions that
can be approximated well by NNWs?

⊙ The norm/metric? Uniform norm?

L^1, L^2, L^p ($1 \leq p < \infty$), L^∞ ? $W^{k,p}$? 6

① Sequence of NNs w.r.t. depth $L \gg 1$, with $N \gg 1$?

② NNW approximations vs. classical approximations (by polynomials, piecewise polynomials, trigonometric polynomials, splines, wavelets, finite elements, etc. and data fitting, e.g., interpolation, least-squares, etc.)?

Some observations/considerations

① Only need to consider the output dimension = 1? Yes, for fixed L . But generally?

Suppose $\Phi^{(j)} \in \text{NN}(0, L, N_0, N_j=1)$ and $\Phi^{(j)} \approx f^{(j)}: X (\subseteq \mathbb{R}^{N_0}) \rightarrow \mathbb{R}^1, j=1, \dots, m$. Then by parallelization with shared inputs, we can construct $\Phi \in \text{NN}(0, L, N_0, N_L=m)$ such that $(\Phi(x))_j = \Phi^{(j)}(x) \forall x \in X$ for $j=1, \dots, m$. Thus, $\Phi \approx f$ with $f = \begin{bmatrix} f^{(1)} \\ \vdots \\ f^{(m)} \end{bmatrix}$.

In deep NN approximations, we may have $\Phi^{(j)} \in \text{NN}(0, L_j, N_0, N_j=1)$ ($j=1, \dots, m$). How to

construct Φ from these $\Phi^{(j)}$ so that Φ can approximate $f: X (\subseteq \mathbb{R}^{N_0}) \rightarrow \mathbb{R}^{\max(N_{L_1}, \dots, N_{L_m})}$? This

leads to the question: how to add one layer to an existing NN to get a new NN which is the same function? For some σ , it may be

fine. But in general, universal or may not 7
be possible.

The next is: for each component of $f: X \rightarrow \mathbb{R}^m$,
we approximate it by a sequence of NNs.

① Any polynomial will not work (as an
activation function). Since, if σ is a polynomial
of degree n . then any $\Phi \in \mathcal{NN}(0, N_0=1, N_L=1)$ is
also a polynomial of degree $\leq n$. Therefore
 $\mathcal{NN}(0, N_0=1, N_L=1)$ is a finite-dimensional space
and cannot approximate $C([0,1])$ which is infinitely
dimensional.

② Proof of universal approximations — abstract
vs. constructive.

Early works

- ① G. Cybenko, Math. Control Signals Systems, 1989.
- ② K.-I. Funahashi, Neural Networks, 1989.
- ③ K. Hornik, M. Stinchcombe, and H. White, Neural
Networks, 1989.
- ④ K. Hornik, Neural Networks, 1991.
- ⑤ Leshno et al. Neural Networks, 1993.

Some notations / definitions.

① The space $C(K)$.

② $d \in \mathbb{N}$: input dimension.

③ $K \subseteq \mathbb{R}^d$: compact (i.e., bounded + closed,
e.g., $K = [0,1]^d$, $K =$ a closed ball in \mathbb{R}^d , etc.)

⑥ Jones 1990
constructive
proof.
⑦ Carruth-
Dickson
(1990)
Also constructions

① $C(K) = \{ \text{all real-valued continuous functions } K \rightarrow \mathbb{R} \}$ 8

② $\|f\| = \|f\|_C = \|f\|_{C(K)} = \|f\|_\infty = \max_{x \in K} |f(x)|$.
Unif. norm, or max norm.

Proposition $(C(K), \|\cdot\|)$ is a Banach space.

Let $\mathcal{A} \subseteq C(K)$. \mathcal{A} is dense in $C(K)$ if: $\forall f \in C(K)$ there exist $\Phi_k \in \mathcal{A}$ ($k=1, 2, \dots$) such that $\|\Phi_k - f\| \rightarrow 0$ as $k \rightarrow \infty$. Equivalently, $\forall f \in C(K) \forall \varepsilon > 0 \exists \Phi \in \mathcal{A}$ such that $\|\Phi - f\| < \varepsilon$. Equivalently, $f \in \overline{\mathcal{A}}$, the closure of \mathcal{A} in $C(K)$.

Remark A NN function is defined on the entire space of the input space \mathbb{R}^d . If we want to approximate $f \in C(K)$, $K \subseteq \mathbb{R}^d$: compact, we can approximate f by some $\tilde{f} \in C(Q)$, $Q = [a, b]^d \supseteq K$, $-\infty < a < b < \infty$, by using mollifiers and approximation to get $\tilde{f} \in C_c(\mathbb{R}^d)$, even that $\tilde{f} \in C_c^\infty(\mathbb{R}^d)$, $\text{supp } \tilde{f} \subseteq Q$. So, we can consider Q instead of K .

Question: If NNs approx. $C(K)$, can they approx $C(K')$ for any different K' ? (K, K' : compact).

① $NN(\sigma, L, d, 1) =$ all NNs with the activation function σ , depth L (#hidden layers = $L-1$), input dimension d , and output dimension 1.

A Universal Approximation Theorem

Theorem (Cybenko 1989, Funahashi 1989 and Hornik, Stinchcombe, and White 1989) Let $\sigma \in C(\mathbb{R})$ be such that $\sigma(-\infty) = 0$ and $\sigma(+\infty) = 1$. Then $NN(\sigma, 2, d, 1)$ is dense in $C(K)$ for any compact $K \subseteq \mathbb{R}^d$ and any $d \in \mathbb{N}$.

Remarks

⊙ Three papers published in the same year! Proved the same result using different methods — so there are three different methods. Note: Funahashi (1989) and Hornik, Stinchcombe, & White (1989) require σ to be increasing. Also, Carroll and Dickinson (1990) gave a constructive proof of the result using Radon transforms.

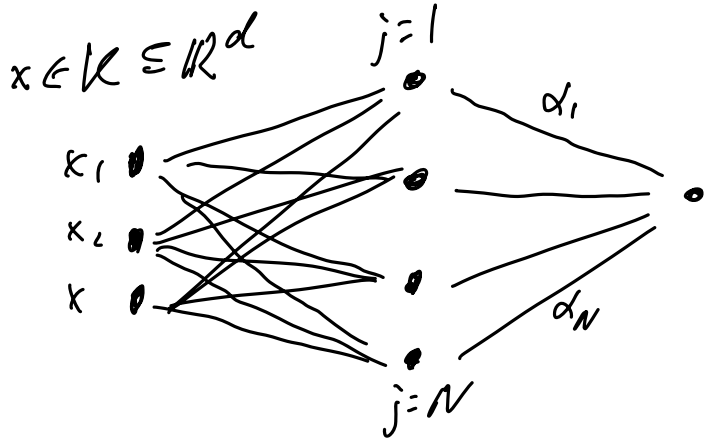
Here, we detail Cybenko's proof. We also sketch Funahashi's proof which is based on an integral formula of Irie-Miyake (1988). We also briefly describe the ideas of proof by Hornik, Stinchcombe, and White.

⊙ Note that the ReLU is not covered in this theorem. Hornik (1991) extended the result to any $\sigma \in C(\mathbb{R})$ that is bounded and non constant. The most general result is given by Leshno, Lin, Pinkus, and Schocken (1993): $\sigma \in C(\mathbb{R})$ leads to the universal approximation property $\iff \sigma$ is not a polynomial. We will sketch the proof of this general result.

⊙ Representation of $G \in \mathcal{NN}(\sigma, 2, d, 1)$:

$$G(x) = \sum_{j=1}^N \alpha_j \sigma(w_j^T x + \theta_j) + \beta, \quad \forall x \in \mathbb{R}^d, \quad (*)$$

where $\alpha_j \in \mathbb{R}, w_j \in \mathbb{R}^d, \theta_j \in \mathbb{R} \quad (j=1, \dots, N), \beta \in \mathbb{R}$.



$$A_1 = \begin{bmatrix} w_1^T \\ \vdots \\ w_N^T \end{bmatrix}_{N \times d} \quad b_1 = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_N \end{bmatrix}$$

$$A_2 = [\alpha_1 \dots \alpha_N]_{1 \times N}$$

$$b_2 = \beta \in \mathbb{R}^1$$

$$L=2, N_0=d, N_L=1, W=N$$

Note: Since $\sigma \neq \text{const.}$, by choosing all $w_j=0$, we can absorb β . i.e., $\mathcal{NN}(\sigma, 2, d, 1)$ consists of all G with $\beta=0$.

Proof (Cybenko's proof)

Step 1 $N_\sigma := \mathcal{NN}(\sigma, L=2, d, 1)$ is a vector subspace of $C(K)$ by elementary constructions. (Or by direct verification.)

Step 2 If $f \in C(K) \setminus \overline{N_\sigma}$ ($\overline{N_\sigma}$ = the closure of N_σ in $C(K)$). then by the Hahn-Banach Thm, $\exists \varphi \in C(K)^*$ such that $\varphi(f) \neq 0$ and $\varphi = 0$ on $\overline{N_\sigma}$.

By Riesz's Thm, $\exists \mu \in \mathcal{M}(K)$ - the space of Radon measures on K , s.t. $\varphi(g) = \int_K g d\mu \forall g \in C(K)$.

Step 3 For any $g(x) = o(w^T x + \theta)$ ($x \in \mathbb{R}^d$), where $w \in \mathbb{R}^d$, $\theta \in \mathbb{R}$, we have $g \in \overline{N_\sigma}$, hence $\varphi(g) = 0$, i.e.,

$$\varphi(g) = \int_K g d\mu = \int_{\mathbb{R}^d} o(w^T x + \theta) d\mu(x) = 0$$

Extend μ to $\mathcal{M}(\mathbb{R}^d)$ (finite, signed Radon measures on \mathbb{R}^d) trivially, i.e., $\forall A \in \mathcal{B}_{\mathbb{R}^d}$ (all Borel sets of \mathbb{R}^d), $\mu(A) = \mu(A \cap K)$.

Then $\mu \in \mathcal{M}_c(\mathbb{R}^d)$ (i.e., $\mu \in \mathcal{M}(\mathbb{R}^d)$ with compact support).

Moreover
$$\int_{\mathbb{R}^d} o(w \cdot x + \theta) d\mu(x) = 0 \quad \forall w \in \mathbb{R}^d, \forall \theta \in \mathbb{R}.$$

By a lemma (next lecture), the assumptions of σ lead to $\mu = 0$. Hence $\varphi(f) = \int_K f d\mu = 0$, a contradiction. (The continuity of σ is needed for $g(x) = o(w^T x + \theta)$ to be in $C(K)$ so the $\varphi(g) = \int_K g d\mu$.) QED

Remark The continuity of σ is needed in the theorem so that $\mathcal{NN}(\sigma, 2, d, 1) \subseteq C(K)$.