

Approximations by Deep Neural Networks with Rectified Power Units

Li Lin

Peking University

September 10, 2019

- x^2 can be approximated within any error $\varepsilon > 0$ by a ReLU network having the depth, the number of weights and computation units all of order $\mathcal{O}(\log \frac{1}{\varepsilon})$.
- Rectified power units (RePUs) are defined as ($s \in \mathbb{N}$)

$$\sigma_s(x) = \begin{cases} x^s, & x \geq 0, \\ 0, & x < 0. \end{cases} \quad (1)$$

- σ_2 :ReQU.
- σ_3 :ReCU.

Some notations

Denote a neural network Φ with input of dimension d , number of layer L , by a matrix-vector sequence

$$\Phi = ((A_1, b_1), \dots, (A_L, b_L)), \quad (2)$$

where $N_0 = d, N_1, \dots, N_L \in \mathbb{N}$, A_k are $N_k \times N_{k-1}$ matrices, and $b_k \in \mathbb{R}^{N_k}$. If Φ is a neural network, and $\rho : \mathbb{R} \rightarrow \mathbb{R}$ is an arbitrary activation function, then define

$$R_\rho(\Phi) : \mathbb{R}^d \rightarrow \mathbb{R}^{N_L}, \quad R_\rho(\Phi)(\mathbf{x}) = \mathbf{x}_L, \quad (3)$$

where $R_\rho(\Phi)(\mathbf{x})$ is defined as

$$\begin{cases} \mathbf{x}_0 := \mathbf{x}, \\ \mathbf{x}_k := \rho(A_k \mathbf{x}_{k-1} + b_k), \quad k = 1, 2, \dots, L-1, \\ \mathbf{x}_L := A_L \mathbf{x}_{L-1} + b_L. \end{cases} \quad (4)$$

- number of hidden layers: $L - 1$
- number of nodes: $\sum_{k=1}^{L-1} N_k$
- number of nonzero wights: $\sum_{k=1}^L (|A_k|_0 + |b_k|_0)$

- The function x, x^2 and xy can be exactly represented with no approximation error using networks having just a few nodes and nonzero weights.

Lemma

For $\forall x, y \in \mathbb{R}$ the following identities hold:

$$x^2 = \beta_2^T \sigma_2(\omega_2 x), \quad (5)$$

$$x = \beta_1^T \sigma_2(\omega_1 x + \gamma_1), \quad (6)$$

$$xy = \beta_1^T \sigma_2(\omega_1 x + \gamma_1 y), \quad (7)$$

where

$$\beta_1 = \frac{1}{4}[1, 1, -1, -1]^T, \beta_2 = [1, 1]^T, \quad (8)$$

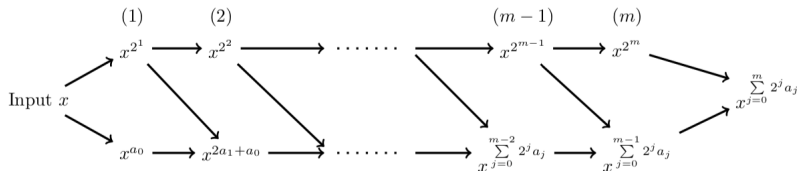
$$\omega_1 = [1, -1, 1, -1]^T, \omega_2 = [1, -1]^T, \gamma_1 = [1, -1, -1, 1]^T. \quad (9)$$

- The Realizations are not unique!

Optimal realizations of polynomials by deep ReQU networks

Theorem

- The monomial $x^n, n \in \mathbb{N}$ defined on \mathbb{R} can be represented exactly by a σ_2 network. The number of network layers, number of nodes and number of weights required to realize x^n are at most $\lfloor \log_2 n \rfloor + 2, 5\lfloor \log_2 n \rfloor + 5$ and $25\lfloor \log_2 n \rfloor + 14$, respectively. Here $\lfloor x \rfloor$ represents the largest integer not exceeding x for $x \in \mathbb{R}$.
- For any $n > 2$, x^n can not be represented exactly by any ReQU network with only one hidden layer.



Univariate polynomials

Theorem

If $f(x)$ is a polynomial of degree n on \mathbb{R} , then it can be represented exactly by a σ_2 neural network with $\lfloor \log_2 n \rfloor + 1$ hidden layers, and number of nodes and nonzero weights are both of order $\mathcal{O}(n)$. To be more precise, the number of nodes is bounded by $9n$, and number of nonzero weights is bounded by $61n$.

$$\begin{aligned}
 f(x) &= a_{15}x^{15} + a_{14}x^{14} + \cdots + a_8x^8 + a_7x^7 + a_6x^6 + \cdots + a_1x + a_0 \\
 &= \underbrace{\underbrace{\underbrace{x^8}_{\xi_{3,0}} \left\{ \underbrace{x^4}_{\xi_{2,0}} \left[\underbrace{x^2}_{\xi_{1,0}} \underbrace{(a_{15}x + a_{14})}_{\xi_{1,8}} + \underbrace{(a_{13}x + a_{12})}_{\xi_{1,7}} \right] + \underbrace{x^2}_{\xi_{1,6}} \underbrace{(a_{11}x + a_{10})}_{\xi_{1,5}} + \underbrace{(a_9x + a_8)}_{\xi_{1,5}} \right] }_{\xi_{2,4}} }_{\xi_{3,2}}}_{\xi_{2,3}} \\
 &\quad + \underbrace{\underbrace{\underbrace{x^4}_{\xi_{2,2}} \left[\underbrace{x^2}_{\xi_{1,4}} \underbrace{(a_7x + a_6)}_{\xi_{1,3}} + \underbrace{(a_5x + a_4)}_{\xi_{1,3}} \right] + \underbrace{x^2}_{\xi_{1,2}} \underbrace{(a_3x + a_2)}_{\xi_{1,1}} + \underbrace{(a_1x + a_0)}_{\xi_{1,1}} }_{\xi_{2,1}} }_{\xi_{3,1}} \right\}.
 \end{aligned}$$

Some remarks

- Use a σ_2 network of scale $\mathcal{O}(\log_2 n)$ to represent x^n exactly.
- Any polynomial of degree less than n can be represented exactly by a σ_2 neural network with $\lfloor \log_2 n \rfloor + 1$ hidden layers, and no more than $\mathcal{O}(n)$ nonzero weights.

Some notations

Define Jacobi-weighted Sobolev space $B_{\alpha,\beta}^m(I)$ as

$$B_{\alpha,\beta}^m(I) := \{u : \partial_x^k u \in L_{\omega^{\alpha+k,\beta+k}}^2(I), 0 \leq k \leq m\}, m \in \mathbb{N}, \quad (10)$$

with norm

$$\|f\|_{B_{\alpha,\beta}^m} := \left(\sum_{k=0}^m \|\partial_x^k u\|_{L_{\omega^{\alpha+k,\beta+k}}^2}^p \right)^{\frac{1}{2}}. \quad (11)$$

The weight

$$\omega^{\alpha,\beta} = (1-x)^\alpha (1+x)^\beta, \quad \alpha, \beta > -1. \quad (12)$$

Define the $L_{\omega^{\alpha,\beta}}^2$ -orthogonal projection $\pi_N^{\alpha,\beta} : L_{\omega^{\alpha,\beta}}^2(I) \rightarrow P_N$ as

$$(\pi_N^{\alpha,\beta} u - u, v)_{\omega^{\alpha,\beta}} = 0, \quad \forall v \in P_N. \quad (13)$$

Error bounds of approximating smooth functions

Theorem

Let $\alpha, \beta > -1$. For any $u \in B_{\alpha, \beta}^m(I)$, there exist a ReQU network Φ_N^u with $\lfloor \log_2 N \rfloor + 1$ hidden layers, $\mathcal{O}(N)$ nodes, and $\mathcal{O}(N)$ nonzero weights, satisfying the following estimate

- if $0 \leq l \leq m \leq N + 1$, we have

$$\|\partial_x^l (R_{\sigma_2}(\Phi_N^u) - u)\|_{\omega^{\alpha+l, \beta+l}} \leq c \sqrt{\frac{(N-m+1)!}{(N-l+1)!}} (N+m)^{\frac{(l-m)}{2}} \|\partial_x^m u\|_{\omega^{\alpha+m, \beta+m}} \quad (14)$$

- if $m > N + 1$, we have

$$\|\partial_x^l (R_{\sigma_2}(\Phi_N^u) - u)\|_{\omega^{\alpha+l, \beta+l}} \leq c (2\pi N)^{-\frac{1}{4}} \left(\frac{\sqrt{e}}{N}\right)^{N-l+1} \|\partial_x^{N+1} u\|_{\omega^{\alpha+N+1, \beta+N+1}} \quad (15)$$

where $c \approx 1$ for $N \gg 1$.

Theorem

For any given function $f(x) \in B_{\alpha,\beta}^m(I)$ with norm less than 1, where m is either a fixed positive integer or infinity, there exists a ReQU network Φ_ε^f with number of layers L , number of nonzero weights N satisfying

- if m is a fixed positive integer, then $L = \mathcal{O}(\frac{1}{m} \log_2(\frac{1}{\varepsilon}))$, and $N = \mathcal{O}(\varepsilon^{-\frac{1}{m}})$;
- if $m = \infty$, i.e. f is analytic, then $L = \mathcal{O}(\log_2(\log \frac{1}{\varepsilon}))$, and $N = \mathcal{O}(\frac{1}{\gamma} \log(\frac{1}{\varepsilon}))$, $\gamma \approx \mathcal{O}(\log(\log \frac{1}{\varepsilon}))$,

can approximate f within an error tolerance ε , i.e.

$$\|R_{\sigma_2}(\Phi_\varepsilon)^f - f\|_{\omega^{\alpha,\beta}(I)} \leq \varepsilon. \quad (16)$$

- For a fixed m , or $N \gg m$,

$$\|R_{\sigma_2}(\Psi_N^u) - u\|_{\omega^{\alpha,\beta}(I)} \leq cN^{-m} \|\partial_x^m u\|_{\omega_{\alpha+m,\beta+m}}. \quad (17)$$

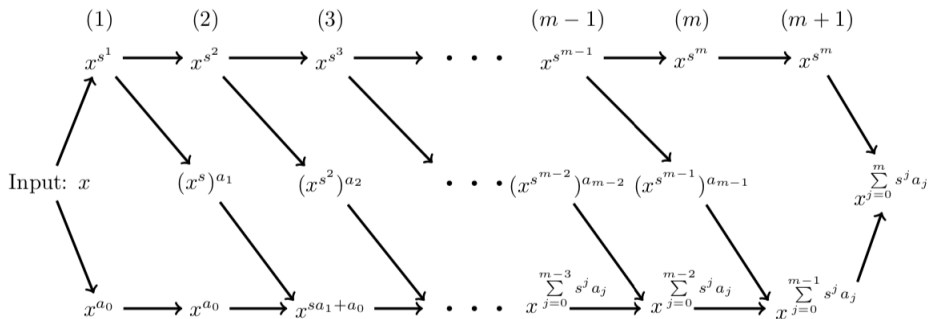
- For analytic function, by taking $m = \infty$

$$\begin{aligned} \|R_{\sigma_2}(\Phi_N^u) - u\|_{\omega^{\alpha,\beta}(I)} &\leq c(2\pi N)^{-\frac{1}{4}} \left(\frac{\sqrt{e/2}}{N}\right)^{N+1} \|u\|_{B_{\alpha,\beta}^\infty} \\ &\leq c'e^{-\gamma N} \|u\|_{B_{\alpha,\beta}^\infty} \end{aligned} \quad (18)$$

Theorem

Regarding the problem of using $\sigma_s(x)$ ($2 \leq s \in \mathbb{N}$) neural networks to exactly represent monomial x^n , $n \in \mathbb{N}$, we have the following results:

- If $s = n$, the monomial x^n can be realized exactly using a σ_s networks having only 1 hidden layer with two nodes.*
- If $1 \leq n < s$, the monomial x^n can be realized exactly using a σ_s networks having only 1 hidden layer with no more than $2s$ nodes.*
- If $n > s \geq 2$, the monomial x^n can be realized exactly using a σ_s networks having only $\lfloor \log_s n \rfloor + 2$ hidden layer with no more than $(6s + 2)(\lfloor \log_s n \rfloor + 2)$ nodes, no more than $\mathcal{O}(25s^2 \lfloor \log_s n \rfloor)$ nonzero weights.*



Approximation of multivariate smooth functions

Theorem

If $f(x)$ is a multivariate polynomial with total degree n on \mathbb{R}^d , then there exists a σ_2 neural network having $d\lfloor \log_2 n \rfloor + d$ hidden layers with no more than $\mathcal{O}(C_d^{n+d})$ activation functions and nonzero weights, can represent f with no error.

$$f(x, y) = \sum_{i=0}^n \left(\sum_{j=0}^{n-i} a_{ij} y^j \right) x^i =: \sum_{i=0}^n a_i^y x^i, \text{ where } a_i^y = \sum_{j=0}^{n-i} a_{ij} y^j \quad (19)$$

Theorem

For a polynomials f_N in a tensor product space $Q_N^d(I_1 \times \cdots \times I_d) := P_N(I_1) \otimes \cdots \otimes P_N(I_d)$, there exists a σ_2 network having $d\lfloor \log_2 N \rfloor + d$ hidden layers with no more than $\mathcal{O}(N^d)$ activation functions and nonzero weights, can represent f_N with no error.

For $u \in B_{\alpha,\beta}^m(I^d)$, we have the following error estimate

$$\|\pi_N^{\alpha,\beta} u - u\|_{L_{\omega,\alpha,\beta}^2(I^d)} \leq cN^{-m}|u|_{B_{\alpha,\beta}^m}, 1 \leq m \leq N. \quad (20)$$

Theorem

For any $u \in B_{\alpha,\beta}^m(I^d)$, with $|u|_{B_{\alpha,\beta}^m(I^d)} \leq 1$, there exists a σ_2 neural network Φ_ε^u having $\mathcal{O}(\frac{d}{m} \log_2 \frac{1}{\varepsilon} + d)$ layers with no more than $\mathcal{O}(\varepsilon^{-\frac{d}{m}})$ nodes and nonzero weights, approximate u with $L_{\omega,\alpha,\beta}^2(I^d)$ -error less than ε , i.e.

$$\|R_{\sigma_2}(\Phi_\varepsilon^u) - u\|_{L_{\omega,\alpha,\beta}^2(I^d)} \leq \varepsilon. \quad (21)$$

Thanks!