

分段线性神经网络的逼近理论*

吴心宇¹ 陈天平² 卢文联³

摘要 随着分段线性函数的广泛应用,本文尝试研究浅层和深层的分段线性神经网络的逼近理论。作者将应用于三层感知机模型的万能逼近定理拓展到分段线性神经网络中,并给出与隐藏神经元个数相关的逼近误差估计。利用分段线性函数构造锯齿函数的显式方法,证明解析函数可以通过分段线性神经网络的深度堆叠以指数速率逼近,并辅以相应的数值实验。

关键词 分段线性, 神经网络, 逼近理论

MR (2000) 主题分类 41A30

中图法分类 O29

文献标志码 A

文章编号 1000-8314(2024)01-0053-18

§1 引 言

神经网络的万能逼近定理可以追溯到 20 世纪 80 年代^[1–2]。单隐层的前馈神经网络,也被称为三层感知机简写为 MLPs, 可以表示为

$$\left\{ f(\mathbf{x}) = \sum_{i=1}^N \beta_i \sigma(\mathbf{w}_i \cdot \mathbf{x} + b_i), \quad \mathbf{x} \in \mathbb{R}^d, \mathbf{w}_i \in \mathbb{R}^d, b_i, \beta_i \in \mathbb{R}, i = 1, \dots, N \right\},$$

其中 $\sigma(\cdot)$ 为 S 型 (sigmoidal) 函数, 当 $t \rightarrow +\infty$ 时, $\sigma(t) \rightarrow 1$, 当 $t \rightarrow -\infty$ 时, $\sigma(t) \rightarrow 0$ ^[1], N 是隐藏神经元个数。众所周知, 三层感知机可以任意精度逼近连续函数^[1–8]。其中, 文 [3] 强调了文 [1] 中对 sigmoidal 函数施加的连续性假设是非必要的。此外, 还有一些与隐藏神经元个数 N 有关的逼近速率的结果。例如, 在一维和上确界范数的情形下, 对于有界变差的连续函数, 误差上界为 $O(\frac{1}{N})$ ^[9]。在积分均方误差的情形下, 对于傅里叶变换振幅一阶矩有界的函数, 误差上界为 $O(\frac{1}{N})$ ^[10]。

随着计算能力的增长, 深度神经网络的应用越来越广泛^[11]。激活函数 $\text{ReLU}(x) = \max(0, x)$ 因其较低的前向和后向计算复杂度而被广泛使用。 ReLU 函数并不满足 sigmoidal 函数的条件, 然而, 由于其在子区域上的线性性质, 可以在深度神经网络上建立逼近理论。锯齿函数在文 [12] 中被提出, 以表明 x^2 可以通过深度和复杂度为 $O(\ln(\frac{1}{\varepsilon}))$ 的 ReLU 神经网络进行误差为 ε 的逼近。该构造方法使用深层网络而不是浅层网络构建了 x^2 的线性插

本文 2022 年 11 月 28 日收到, 2023 年 10 月 10 日收到修改稿。

¹复旦大学数学科学学院, 上海 200433. E-mail: xywu19@fudan.edu.cn

²复旦大学数学科学学院, 上海 200433; 上海数学中心, 上海 200438; 上海市现代应用数学重点实验室, 上海 200433. E-mail: tchen@fudan.edu.cn

³通信作者。复旦大学数学科学学院, 上海 200433; 上海数学中心, 上海 200438; 上海市现代应用数学重点实验室, 上海 200433. E-mail: wenlian@fudan.edu.cn

*本文受到国家重点研发计划 (No. 2018AAA0100303), 国家自然科学基金 (No. 62072111), 上海市科技重大专项 (No. 2018SHZDZX01) 和张江实验室项目的资助。

值. 此外, 文 [13] 指出深度 ReLU 神经网络对解析函数的逼近呈指数收敛, 并给出了具有跳跃连接的深度网络的简明表示.

然而, “ReLU 神经元死亡”是训练深度神经网络的障碍之一^[14], 这意味着当 ReLU 神经元变得不活跃时, 没有梯度通过神经元回传. 在这种情况下, 一组 ReLU 变体激活函数被提出, 例如 Leaky ReLU 函数^[15], APL 函数^[16], SReLU 函数^[17] 和 ReLU6 函数^[18]. 它们都属于分段线性(简写为 PWL) 函数, 并作为分段线性神经网络(简写为 PWLNN) 的激活函数. PWLNN 结合了线性和非线性的特点, 并且比其他非线性方法更适合建模、学习和分析^[19]. 尽管如此, 目前仍然缺乏对 PWLNN 的理论分析. 为了解决这个问题, 我们探索了浅层分段线性神经网络和深层分段线性神经网络(简写为 PWL-DNN) 的逼近能力.

本文的主要内容如下:

- (1) 揭示以分段线性函数为激活函数的三层感知机模型的万能逼近能力, 并通过隐藏神经元的数量控制误差.
- (2) 提供从分段线性函数构造锯齿函数的显式方法, 并表明解析函数可以通过分段线性神经网络深度的堆叠而非宽度的增加以实现指数速率的逼近.
- (3) 通过数值实验验证本文的结论.

§2 预备知识

§2.1 激活函数

我们首先给出文 [3] 中激活函数的定义.

定义 2.1 $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ 被称为 sigmoidal 函数, 如果极限 $\lim_{x \rightarrow +\infty} \sigma(x)$ 和 $\lim_{x \rightarrow -\infty} \sigma(x)$ 存在, 且

$$\lim_{x \rightarrow +\infty} \sigma(x) = 1, \quad (2.1)$$

$$\lim_{x \rightarrow -\infty} \sigma(x) = 0. \quad (2.2)$$

应该强调的是 sigmoidal 函数不一定是连续的或单调递增的. 在本文中, 我们将此定义扩展到以下可叠加 S 型 (superpositioned-sigmoidal) 函数.

定义 2.2 $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ 被称为 superpositioned-sigmoidal 函数, 如果存在正整数 P 以及常数 $a \neq \frac{1}{2}$, $a_j \neq 0$, $b_j > 0$ 与 c_j , $j = 1, \dots, P$, 使得

$$a + \sum_{j=1}^P a_j \sigma(b_j x + c_j) \quad (2.3)$$

是有界的 sigmoidal 函数.

例 2.1 (ReLU 函数) $\text{ReLU}(x) = \max(0, x)$ 是 superpositioned-sigmoidal 函数, 由于

$$\text{ReLU}(x) - \text{ReLU}(x-1) = \sigma_1(x), \quad (2.4)$$

其中

$$\sigma_1(x) = \begin{cases} 0, & x < 0, \\ x, & 0 \leq x \leq 1, \\ 1, & x > 1 \end{cases} \quad (2.5)$$

是有界的 sigmoidal 函数.

例 2.2 为了方便起见, 我们定义函数 $\sigma_2(\cdot) \in C([0, 1])$ 为

$$\sigma_2(x) = \begin{cases} 2x, & x < \frac{1}{2}, \\ 2(1-x), & x \geq \frac{1}{2}. \end{cases} \quad (2.6)$$

直接验证可得

$$\sigma_2(x) = 2\sigma_1(x) - 4\sigma_1\left(x - \frac{1}{2}\right) = 2\text{ReLU}(x) - 4\text{ReLU}\left(x - \frac{1}{2}\right).$$

单位分解. 显然, $\sigma_1(\cdot)$ 和 $\sigma_2(\cdot)$ 存在以下单位分解:

$$\sigma_1(x) + \sigma_1(1-x) = 1, \quad \sum_{i \in \mathbb{Z}} \tilde{\sigma}_2\left(x - \frac{i}{2}\right) = 1, \quad (2.7)$$

其中 $\tilde{\sigma}_2 : \mathbb{R} \rightarrow \mathbb{R}$ 是 σ_2 的零延拓, 即当 $x \in [0, 1]$ 时, $\tilde{\sigma}_2(x) = \sigma_2(x)$, 其余情况下 $\tilde{\sigma}_2(x) = 0$.

定义 2.3 (PWL 函数) $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ 被称为 PWL 函数, 如果 σ 连续, 且存在 M ($1 \leq M < \infty$) 个间断点 t_1, \dots, t_M , 在每一区间 $(-\infty, t_1), (t_1, t_2), \dots, (t_M, \infty)$ 中, σ 都是线性函数, 并且任意两个相邻区间内的斜率不同.

在以下条件下, PWL 函数也是 superpositioned-sigmoidal 函数.

命题 2.1 假设具有 M 个间断点的 PWL 函数 σ 表示为 $\sigma(x) = k_0x + b_0 + k_1(x - t_1)\mathbf{1}_{\{x>t_1\}} + \dots + k_M(x - t_M)\mathbf{1}_{\{x>t_M\}}$, 则 σ 是 superpositioned-sigmoidal 函数, 当且仅当存在 c , 使得 $(k_0 + k_\infty)c + b_0 + b_\infty = 0$ 不成立且 $(k_0 - k_\infty)c + b_0 - b_\infty = 0$ 不成立, 其中 $k_\infty = k_0 + \sum_{i=1}^M k_i$, $b_\infty = b_0 - \sum_{i=1}^M k_i t_i$.

证 由定义 2.2, 我们需要找到正整数 P , 以及常数 $a \neq \frac{1}{2}$, $a_j \neq 0$, $b_j > 0$, $c_j, j = 1, \dots, P$, 使得 $a + \sum_{j=1}^P a_j \sigma(b_j x + c_j)$ 是有界的 sigmoidal 函数, 即

$$\begin{aligned} a + \sum_{j=1}^P a_j b_j k_\infty x + \sum_{j=1}^P a_j c_j k_\infty + \sum_{j=1}^P a_j b_\infty &= 1, \quad x \rightarrow +\infty, \\ a + \sum_{j=1}^P a_j b_j k_0 x + \sum_{j=1}^P a_j c_j k_0 + \sum_{j=1}^P a_j b_0 &= 0, \quad x \rightarrow -\infty. \end{aligned}$$

这等价于解以下的方程组:

$$\sum_{j=1}^P a_j b_j = 0, \quad (2.8)$$

$$\sum_{j=1}^P a_j (c_j k_\infty + b_\infty) = 1 - a, \quad (2.9)$$

$$\sum_{j=1}^P a_j (c_j k_0 + b_0) = -a. \quad (2.10)$$

方程组 (2.9) 和 (2.10) 有解的充要条件是

$$\text{rank} \begin{pmatrix} c_1 k_\infty + b_\infty, & \cdots, & c_P k_\infty + b_\infty \\ c_1 k_0 + b_0, & \cdots, & c_P k_0 + b_0 \end{pmatrix} = \text{rank} \begin{pmatrix} c_1 k_\infty + b_\infty, & \cdots, & c_P k_\infty + b_\infty, & 1 - a \\ c_1 k_0 + b_0, & \cdots, & c_P k_0 + b_0, & -a \end{pmatrix}.$$

充分性. 如果存在 c , 使得 $(k_0 + k_\infty)c + b_0 + b_\infty = 0$ 不成立且 $(k_0 - k_\infty)c + b_0 - b_\infty = 0$ 不成立, 我们可以令 $P > 2$, $b_j = j$, $c_j = c$, $j = 1, \dots, P$, 从而得到相应的 $\{a_j\}_{j=1}^P$ 和 $a \neq \frac{1}{2}$.

必要性. 我们往证如下论断: 如果对于任何 c , $(k_0 + k_\infty)c + b_0 + b_\infty = 0$ 或 $(k_0 - k_\infty)c + b_0 - b_\infty = 0$ 成立, 则 PWL 函数 σ 不是 superpositioned-sigmoidal 函数.

(1) 当 $k_0 + k_\infty = 0$ 且 $b_0 + b_\infty = 0$ 时, 将 (2.9) 与 (2.10) 相加可以得到 $a = \frac{1}{2}$.

(2) 当 $k_0 - k_\infty = 0$ 且 $b_0 - b_\infty = 0$ 时,

$$\begin{aligned} & \text{rank} \begin{pmatrix} c_1 k_\infty + b_\infty, & \cdots, & c_P k_\infty + b_\infty, & 1 - a \\ c_1 k_0 + b_0, & \cdots, & c_P k_0 + b_0, & -a \end{pmatrix} \\ &= \text{rank} \begin{pmatrix} c_1 k_\infty + b_\infty, & \cdots, & c_P k_\infty + b_\infty \\ c_1 k_0 + b_0, & \cdots, & c_P k_0 + b_0 \end{pmatrix} + 1. \end{aligned}$$

综上所述, 线性方程组无解.

注 2.1 在有限区间内, 这个条件可以被放宽, 如引理 4.3 的证明所示.

§2.2 神经网络架构

与文 [13] 类似, 我们采用如下符号:

(1) 令 $\{x_{m:n}\} = \{x_i : i = m, m+1, \dots, n\}$ 和 $\{x_{m_1:n_1, m_2:n_2}\} = \{x_{i,j} : i = m_1, m_1+1, \dots, n_1, j = m_2, m_2+1, \dots, n_2\}$.

(2) 称 $y \in \mathcal{L}(x_1, \dots, x_n)$, 如果存在 $\beta_i \in \mathbb{R}$, $i = 0, 1, \dots, n$, 使得 $y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$.

(3) 称 $y_\sigma \in \mathcal{L}_\sigma(x_1, \dots, x_n)$, 如果存在 $y \in \mathcal{L}(x_1, \dots, x_n)$ 并且 $y_\sigma = \sigma(y)$.

定义 2.4 函数 $f(x_1, \dots, x_d)$ 被称为属于激活函数为 σ 的全连接神经网络函数族 $\mathcal{F}_{L,M}^\sigma(\mathbb{R}^d)$, 如果存在变量 $\{y_{1:L,1:M}\}$, 使得

$$y_{1,m} \in \mathcal{L}_\sigma(x_{1:d}), \quad y_{l+1,m} \in \mathcal{L}_\sigma(y_{l,1:M}), \quad f \in \mathcal{L}(y_{L,1:M}), \quad (2.11)$$

其中 $m = 1, \dots, M$, $l = 1, \dots, L-1$, 并且 $\{y_{1:L,1:M}\}$ 被称为 f 的隐变量.

更具体地, $f \in \mathcal{F}_{L,M}^\sigma(\mathbb{R}^d)$ 可以由下式给出:

$$f(\mathbf{x}) = \begin{cases} W(\sigma \circ (W_1 \mathbf{x} + \mathbf{b}_1)) + b, & L = 1, \\ W(\sigma \circ (W_L(\sigma \circ \dots \circ (W_1 \mathbf{x} + \mathbf{b}_1)) + \mathbf{b}_L)) + b, & L \geq 2, \end{cases}$$

其中 $\mathbf{x} \in \mathbb{R}^d$, $W_1 \in \mathbb{R}^{M \times d}$, $W_l \in \mathbb{R}^{M \times M}$, $l = 2, \dots, L$, $W \in \mathbb{R}^{1 \times M}$, $\mathbf{b}_l \in \mathbb{R}^M$, $l = 1, \dots, L$, $b \in \mathbb{R}$. 符号 \circ 表示函数的复合, 激活函数 σ 作用于向量中的每个元素.

注 2.2 传统的三层 MLP 函数族 $\mathcal{M}_N^\sigma(\mathbb{R}^d)$ 可表示为

$$\left\{ f(\mathbf{x}) = \sum_{i=1}^N \beta_i \sigma(\mathbf{w}_i \cdot \mathbf{x} + b_i), \quad \mathbf{x} \in \mathbb{R}^d, \mathbf{w}_i \in \mathbb{R}^d, b_i, \beta_i \in \mathbb{R}, i = 1, \dots, N \right\}.$$

由于 f 是隐变量的纯线性组合 (不含 β_0), 因此 $\mathcal{M}_N^\sigma(\mathbb{R}^d) \subsetneq \mathcal{F}_{1,N}^\sigma(\mathbb{R}^d)$. 我们称 $\mathcal{M}_N^\sigma(\mathbb{R}^d)$ 为浅层神经网络函数族.

定义 2.5 函数 $f(x_1, \dots, x_d)$ 被称为属于激活函数为 σ 的跳跃连接神经网络函数族 $\mathcal{S}_{L,M}^\sigma(\mathbb{R}^d)$, 如果存在变量 $\{y_{1:L,1:M}\}$, 使得

$$y_{1,m} \in \mathcal{L}_\sigma(x_{1:d}), \quad y_{l+1,m} \in \mathcal{L}_\sigma(x_{1:d}, y_{l,1:M}), \quad f \in \mathcal{L}(x_{1:d}, y_{1:L,1:M}), \quad (2.12)$$

其中 $m = 1, \dots, M$, $l = 1, \dots, L-1$, 并且 $\{y_{1:L,1:M}\}$ 被称为 f 的隐变量.

更具体地, $f \in \mathcal{S}_{L,M}^\sigma(\mathbb{R}^d)$ 可以由下式给出:

$$f(\mathbf{x}) = \sum_{l=1}^L W_l \mathbf{y}_l + W \mathbf{x} + b,$$

$$\mathbf{y}_l = \begin{cases} \sigma \circ (W_0^1 \mathbf{x} + \mathbf{b}_1), & l = 1, \\ \sigma \circ (W_1^l \mathbf{y}_{l-1} + W_0^l \mathbf{x} + \mathbf{b}_l), & l \geq 2, \end{cases}$$

其中 $\mathbf{x} \in \mathbb{R}^d$, $W_l \in \mathbb{R}^{1 \times M}$, $W_0^l \in \mathbb{R}^{M \times d}$, $W_1^l \in \mathbb{R}^{M \times M}$, $\mathbf{b}_l \in \mathbb{R}^M$, $l = 1, \dots, L$, $W \in \mathbb{R}^{1 \times d}$, $b \in \mathbb{R}$. 符号 \circ 表示函数的复合, 激活函数 σ 作用于向量中的每个元素.

注 2.3 本质上 $\mathcal{S}_{L,M}^\sigma(\mathbb{R}^d) \subseteq \mathcal{F}_{L,M+d+1}^\sigma(\mathbb{R}^d)$. 然而, 当 $L \gg 2$ 时, 符号 $\mathcal{S}_{L,M}^\sigma(\mathbb{R}^d)$ 为深度神经网络的分析提供了便利.

§3 主要结论

在本节中, 我们分别研究了浅层 PWLNN 和 PWL-DNN 的逼近能力. 详细的证明在第 4 节中.

基于命题 2.1, 我们首先给出了激活函数为 superpositioned-sigmoidal 函数的三层 MLP 在 \mathbb{R} 上的万能逼近定理.

定理 3.1 假设 σ 是一个 superpositioned-sigmoidal 函数, $f(x)$ 在 \mathbb{R} 上连续, 满足

$$\lim_{x \rightarrow -\infty} f(x) = A, \quad (3.1)$$

$$\lim_{x \rightarrow +\infty} f(x) = B, \quad (3.2)$$

其中 A, B 是常数. 则对于任意 $\varepsilon > 0$, 存在 $N \in \mathbb{Z}$ 和函数 $\hat{f} \in \mathcal{M}_N^\sigma(\mathbb{R})$, 使得 $|f(x) - \hat{f}(x)| < \varepsilon$ 对任意 $x \in \mathbb{R}$ 成立.

与文 [3] 的证明框架类似, 我们基于定理 3.1 和 Bochner-Riesz Means^[20] 建立了 \mathbb{R}^d 上的近似. 此外, 我们用隐藏神经元个数 N 和空间维度 d 给出了误差估计.

定理 3.2 假设 σ 是一个 PWL 函数, $f(\mathbf{x})$ 在 $[0, 1]^d$ 上连续. 则对于任意 $\varepsilon > 0$, 存在 $N \in \mathbb{Z}$ 和函数 $\hat{f} \in \mathcal{M}_N^\sigma(\mathbb{R}^d)$, 使得 $|f(\mathbf{x}) - \hat{f}(\mathbf{x})| < \frac{\varepsilon}{2} + O\left(\frac{R^{3d+2}}{N^2}\right)$, 其中 $R = R(\varepsilon) > 0$.

与文 [3] 中的 sigmoidal 激活函数设置不同, 这里的上界 $O\left(\frac{R^{3d+2}}{N^2}\right)$ 是通过 PWL 函数插值得到的, 主要根据引理 4.3. 对于 sigmoidal 函数, 上界为 $O\left(\frac{R^{2d+1}}{N}\right)$.

除了研究浅层 PWLNN 之外, 我们还给出 PWL-DNN 逼近解析函数的误差估计. 受文 [13] 的启发, 我们证明解析函数可以通过 PWL-DNN 堆叠深度以指数收敛速率逼近.

定理 3.3 假设 σ 是一个 PWL 函数, 解析函数 $f(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{N}^d} a_{\mathbf{k}} \mathbf{x}^{\mathbf{k}}$ 在 $[-1, 1]^d$ 中绝对收敛. 则对任意 $\delta \in (0, 1)$ 与 $\varepsilon \in (0, \frac{1}{e})$, 存在函数 $\hat{f} \in \mathcal{S}_{L,3}^\sigma(\mathbb{R}^d)$, 其中 $L = \lceil [e(\frac{1}{d\delta} \log \frac{1}{\varepsilon} + 1)]^{\max(2d, 3)} \rceil$, 使得 $|f(\mathbf{x}) - \hat{f}(\mathbf{x})| < 2\varepsilon \sum_{\mathbf{k}} |a_{\mathbf{k}}|$ 对于任意 $\mathbf{x} \in [-1 + \delta, 1 - \delta]^d$ 成立, 其中 $\lceil \cdot \rceil$ 表示向上取整函数.

注 3.1 从第 4 节中的证明可以发现, 浅层 PWLNN 和 PWL-DNN 都能根据 PWL 函数的线性性从而构造出线性插值函数. 然而, 对于 x^2 , 宽度为 2、深度为 L 的 PWL-DNN 可以等价于具有 2^L 隐藏神经元的浅层 PWLNN. 因此, PWL-DNN 比浅层 PWLNN 具有更强的逼近能力.

与文 [13] 和其他工作不同, 我们的结论适用于所有 PWL 函数而不仅仅是 ReLU 函数. 关键在于我们在引理 4.3 中给出了网络的显式表达.

在这里, 我们举一些例子.

例 3.1 (Leaky ReLU 函数^[15]) Leaky ReLU 定义为

$$\sigma(x) = \begin{cases} \alpha x, & x < 0, \\ x, & x \geq 0, \end{cases}$$

其中 $0 < \alpha < 1$.

显然 σ_1 和 σ_2 可以被写作

$$\begin{aligned} \sigma_1(x) &= -\frac{\alpha}{1-\alpha} + \frac{\sigma(x)}{1-\alpha} - \frac{\sigma(x-1)}{1-\alpha}, \quad x \in \mathbb{R}, \\ \sigma_2(x) &= -\frac{2\alpha}{1-\alpha} + \frac{2+2\alpha}{1-\alpha} \sigma(x) - \frac{4\sigma(x-\frac{1}{2})}{1-\alpha}, \quad x \in [0, 1]. \end{aligned}$$

例 3.2 (ReLU6 函数^[18]) ReLU6 定义为

$$\sigma(x) = \min(6, \max(0, x)),$$

显然 σ_1 和 σ_2 可以被写作

$$\begin{aligned}\sigma_1(x) &= \frac{1}{6}\sigma(6x), \quad x \in \mathbb{R}, \\ \sigma_2(x) &= 2\sigma(x) - 4\sigma\left(x - \frac{1}{2}\right), \quad x \in [0, 1].\end{aligned}$$

例 3.3 (APL 函数^[16]) 自适应分段线性 (APL) 激活函数是合页型函数的和, 定义为

$$\sigma(x) = \max(0, x) + \sum_{s=1}^S a^s \max(0, -x + b^s).$$

合页型函数的数量 S 是预先设置的超参数, 而变量 $a^s, b^s, s = 1, \dots, S$ 是可学习的参数.

例 3.4 (SReLU 函数^[17]) SReLU 是三个线性函数的组合:

$$\sigma(x) = \begin{cases} t^r + a^r(x - t^r), & x \geq t^r, \\ x, & t^r > x > t^l, \\ t^l + a^l(x - t^l), & x \leq t^l, \end{cases}$$

其中 t^r, a^r, t^l, a^l 是可学习的参数.

对于 APL 函数和 SReLU 函数, 我们可以利用 §4.3 节中的 (4.25) 式构造 σ_2 . 这些例子都表明我们的理论适用于所有 PWL 函数, 无论是具有预先指定的还是可学习参数的函数.

§4 证 明

§4.1 定理 3.1 的证明

定理 3.1 的证明 我们应用文 [3] 中给出的构造性证明. 由于 σ 是 superpositioned-sigmoidal 函数, 因此存在正整数 P 和常数 $a \neq \frac{1}{2}, a_j \neq 0, b_j > 0$ 以及 $c_j, j = 1, \dots, P$, 使得

$$\tilde{\sigma}(x) = a + \sum_{j=1}^P a_j \sigma(b_j x + c_j) \tag{4.1}$$

是有界的 sigmoidal 函数.

记

$$F(a, x', x'') = \frac{af(x') + (1-a)f(x'')}{1-2a}. \tag{4.2}$$

由于 f 是有界的, 所以存在常数 C , 使得 $|F(a, x', x'')| \leq C$ 对任意 x', x'' 成立.

由假设 (3.1)–(3.2), 对任意 $\varepsilon > 0$, 我们可以找到整数 $M > \lceil \sqrt{\frac{9C}{\varepsilon}} \rceil$, 使得当 $x < -M$ 时 $|f(x) - A| < \frac{\varepsilon}{6}$; 当 $x > M$ 时 $|f(x) - B| < \frac{\varepsilon}{6}$; 并且当 $|x'| \leq M, |x''| \leq M$ 以及 $|x' - x''| \leq \frac{1}{M}$ 时 $|f(x'') - f(x')| < \frac{\varepsilon}{6}$.

将 $[-M, M]$ 等分为 $2M^2$ 个小区间, 并令

$$-M = x_0 < x_1 < \dots < x_{M^2} = 0 < x_{M^2+1} < \dots < x_{2M^2} = M,$$

且 $t_i = \frac{1}{2}(x_i + x_{i+1})$, $i = 0, \dots, 2M^2 - 1$.

定义 $\hat{f}(x)$ 与 $\tilde{f}(x)$ 如下:

$$\begin{aligned}\hat{f}(x) &= \frac{F(a, M, -M)}{2M^2} \sum_{i=1}^{2M^2} \sum_{j=1}^P a_j \sigma(Kb_j(x - t_{i-1}) + c_j) \\ &\quad + \frac{F(a, M, -M)}{2M^2} \sum_{i=1}^{2M^2} \sum_{j=1}^P a_j \sigma(-Kb_j(x - t_{i-1}) + c_j) \\ &\quad + \sum_{i=1}^{2M^2} [f(x_i) - f(x_{i-1})] \sum_{j=1}^P a_j \sigma(Kb_j(x - t_{i-1}) + c_j),\end{aligned}\tag{4.3}$$

$$\tilde{f}(x) = f(-M) + \sum_{i=1}^{2M^2} [f(x_i) - f(x_{i-1})] \tilde{\sigma}(K(x - t_{i-1})),\tag{4.4}$$

其中 K 是待确定的足够大的常数.

从 (4.1), (4.3) 以及 (4.4) 式中可知,

$$\hat{f}(x) - \tilde{f}(x) = \frac{F(a, M, -M)}{2M^2} \sum_{i=1}^{2M^2} [\tilde{\sigma}(K(x - t_{i-1})) + \tilde{\sigma}(-K(x - t_{i-1})) - 1].\tag{4.5}$$

不失一般性, 我们假设 $|\tilde{\sigma}(x)| \leq 1$. 由于 $\tilde{\sigma}$ 是 sigmoidal 函数, 所以存在 $W > 0$, 使得若 $u > W$, 则有 $|\tilde{\sigma}(u) - 1| < \min(\frac{1}{2M^2}, \frac{\varepsilon}{6C})$, 若 $u < -W$ 则有 $|\tilde{\sigma}(u)| < \min(\frac{1}{2M^2}, \frac{\varepsilon}{6C})$. 取 $K > 0$ 满足 $\frac{K}{2M} > W$.

(a) 当 $x < -M$ 时, 有 $x - t_{i-1} < -\frac{1}{2M}$, $K(x - t_{i-1}) < -W$ 且 $-K(x - t_{i-1}) > W$ 对 $i = 1, \dots, 2M^2$ 成立. 因此

$$\begin{aligned}|\hat{f}(x) - \tilde{f}(x)| &\leq \frac{C}{2M^2} \sum_{i=1}^{2M^2} |\tilde{\sigma}(K(x - t_{i-1}))| + |\tilde{\sigma}(-K(x - t_{i-1})) - 1| \\ &\leq \frac{C}{2M^2} \cdot 2M^2 \cdot \left(\frac{\varepsilon}{6C} + \frac{\varepsilon}{6C} \right) < \frac{\varepsilon}{2},\end{aligned}\tag{4.6}$$

$$\begin{aligned}|\tilde{f}(x) - f(x)| &\leq |f(-M) - f(x)| + \sum_{i=1}^{2M^2} |[f(x_i) - f(x_{i-1})] \tilde{\sigma}(K(x - t_{i-1}))| \\ &\leq \frac{\varepsilon}{3} + 2M^2 \cdot \frac{\varepsilon}{6} \frac{1}{2M^2} = \frac{\varepsilon}{2}.\end{aligned}\tag{4.7}$$

(b) 当 $x > M$ 时, 有 $x - t_{i-1} > \frac{1}{2M}$, $K(x - t_{i-1}) > W$ 且 $-K(x - t_{i-1}) < -W$ 对

$i = 1, \dots, 2M^2$ 成立. 因此

$$\begin{aligned} |\hat{f}(x) - \tilde{f}(x)| &\leq \frac{C}{2M^2} \sum_{i=1}^{2M^2} |\tilde{\sigma}(K(x - t_{i-1})) - 1| + |\tilde{\sigma}(-K(x - t_{i-1}))| \\ &\leq \frac{C}{2M^2} \cdot 2M^2 \cdot \left(\frac{\varepsilon}{6C} + \frac{\varepsilon}{6C} \right) < \frac{\varepsilon}{2}, \end{aligned} \quad (4.8)$$

$$\begin{aligned} |\tilde{f}(x) - f(x)| &\leq |f(M) - f(x)| + \sum_{i=1}^{2M^2} |f(x_i) - f(x_{i-1})| |\tilde{\sigma}(K(x - t_{i-1})) - 1| \\ &\leq \frac{\varepsilon}{3} + 2M^2 \cdot \frac{\varepsilon}{6} \cdot \frac{1}{2M^2} = \frac{\varepsilon}{2}. \end{aligned} \quad (4.9)$$

(c) 当 $x \in [x_{k-1}, x_k]$ 时, 若 $i = k$ 则 $|x - t_{i-1}| \leq \frac{1}{2M}$; 若 $i < k$ 则 $x - t_{i-1} > \frac{1}{2M}$; 若 $i > k$ 则 $x - t_{i-1} < -\frac{1}{2M}$. 因此

$$\begin{aligned} |\hat{f}(x) - \tilde{f}(x)| &\leq \frac{C}{2M^2} \sum_{i=1}^{k-1} |\tilde{\sigma}(K(x - t_{i-1})) - 1| + |\tilde{\sigma}(-K(x - t_{i-1}))| \\ &\quad + \frac{C}{2M^2} \sum_{i=k+1}^{2M^2} |\tilde{\sigma}(K(x - t_{i-1}))| + |\tilde{\sigma}(-K(x - t_{i-1})) - 1| \\ &\quad + \frac{C}{2M^2} |\tilde{\sigma}(K(x - t_{k-1})) + \tilde{\sigma}(-K(x - t_{k-1})) - 1| \\ &< \frac{C}{2M^2} \cdot (2M^2 - 1) \cdot \left(\frac{\varepsilon}{6C} + \frac{\varepsilon}{6C} \right) + \frac{3C}{2M^2} < \frac{\varepsilon}{2}, \end{aligned} \quad (4.10)$$

$$\begin{aligned} |\tilde{f}(x) - f(x)| &\leq |f(x_{k-1}) - f(x)| + \sum_{i=1}^{k-1} |f(x_i) - f(x_{i-1})| |\tilde{\sigma}(K(x - t_{i-1})) - 1| \\ &\quad + |f(x_k) - f(x_{k-1})| |\tilde{\sigma}(K(x - t_{k-1}))| \\ &\quad + \sum_{i=k+1}^{2M^2} |f(x_i) - f(x_{i-1})| |\tilde{\sigma}(K(x - t_{i-1}))| \\ &< \frac{\varepsilon}{6} + k \cdot \frac{1}{2M^2} \frac{\varepsilon}{6} + \frac{\varepsilon}{6} + (2M^2 - k) \cdot \frac{1}{2M^2} \frac{\varepsilon}{6} < \frac{\varepsilon}{2}. \end{aligned} \quad (4.11)$$

综上所述, 存在 $N \in \mathbb{Z}$ 和 $\hat{f} \in \mathcal{M}_N^\sigma(\mathbb{R})$, 使得

$$|f(x) - \hat{f}(x)| < \varepsilon, \quad x \in (-\infty, +\infty). \quad (4.12)$$

§4.2 定理 3.2 的证明

引理 4.1 若函数 $f(x) \in C^2[0, 1]$, 则存在正整数 N 和常数 $c_i, y_i, \theta_i, i = 1, \dots, N$, 使得

$$\left| f(x) - \sum_{i=1}^N c_i \sigma_1(y_i x + \theta_i) \right| = O\left(\frac{1}{N^2}\right) \quad (4.13)$$

对所有 $x \in [0, 1]$ 成立.

引理 4.1 的证明 令 $x_i = \frac{i}{N}$, $i = 0, 1, \dots, N$. 记

$$\tilde{f}(x) = f(x_0) + \sum_{i=1}^N [f(x_i) - f(x_{i-1})] \sigma_1\left(\frac{x - x_{i-1}}{x_i - x_{i-1}}\right). \quad (4.14)$$

当 $x \in [x_{k-1}, x_k]$ 时, $\tilde{f}(x)$ 可写作

$$\begin{aligned} \tilde{f}(x) &= f(x_0) + \sum_{i=1}^N [f(x_i) - f(x_{i-1})] \sigma_1\left(\frac{x - x_{i-1}}{x_i - x_{i-1}}\right) \\ &= f(x_0) + \sum_{i=1}^{k-1} [f(x_i) - f(x_{i-1})] + [f(x_k) - f(x_{k-1})] \frac{x - x_{k-1}}{x_k - x_{k-1}} \\ &= f(x_{k-1}) + \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}(x - x_{k-1}). \end{aligned}$$

由于 $f(x) \in C^2[0, 1]$,

$$f(x) = f(x_{k-1}) + f'(x_{k-1})(x - x_{k-1}) + f''(x_{k-1})(x - x_{k-1})^2 + o((x - x_{k-1})^2). \quad (4.15)$$

因此 $|f(x) - \tilde{f}(x)| = O((x - x_{k-1})^2) = O(\frac{1}{N^2})$.

引理 4.2 若函数 $f(x) \in C^2[0, 1]$, 则存在正整数 N 和常数 $c_i, y_i, \theta_i, i = 1, \dots, N$, 使得

$$|f(x) - \sum_{i=1}^N c_i \text{ReLU}(y_i x + \theta_i)| = O\left(\frac{1}{N^2}\right) \quad (4.16)$$

对所有 $x \in [0, 1]$ 成立.

引理 4.2 的证明 对于 $\sigma_1(x) = \text{ReLU}(x) - \text{ReLU}(x - 1)$, 由 (2.7) 式可得 $\sigma_1(x) + \sigma_1(1 - x) = 1$. 将其代入 (4.14) 式得到

$$\begin{aligned} \tilde{f}(x) &= f(x_0)[\text{ReLU}(x) - \text{ReLU}(x - 1) + \text{ReLU}(1 - x) - \text{ReLU}(-x)] \\ &\quad + \sum_{i=1}^M [f(x_i) - f(x_{i-1})] \left[\text{ReLU}\left(\frac{x - x_{i-1}}{x_i - x_{i-1}}\right) - \text{ReLU}\left(\frac{x - x_i}{x_i - x_{i-1}}\right) \right]. \end{aligned} \quad (4.17)$$

因此可取 $N = 2M + 4$, 以及系数 $c_i, y_i, \theta_i, i = 1, \dots, N$, 使得 $\sum_{i=1}^N c_i \text{ReLU}(y_i x + \theta_i) = \tilde{f}(x)$ 和 (4.16) 式成立.

引理 4.3 若函数 $f(x) \in C^2[0, 1]$, σ 是 PWL 函数. 则存在正整数 N 以及常数 $c_i, y_i, \theta_i, i = 1, \dots, N$, 使得

$$\left| f(x) - \sum_{i=1}^N c_i \sigma(y_i x + \theta_i) \right| = O\left(\frac{1}{N^2}\right) \quad (4.18)$$

对所有 $x \in [0, 1]$ 成立. 此外, 若

$$\left| f(x) - \sum_{i=1}^N c_i \sigma(y_i x + \theta_i) \right| = o\left(\frac{1}{N^2}\right) \quad (4.19)$$

对所有 $x \in [0, 1]$ 成立, 则 $f(x)$ 一定是线性函数, 可写作 $f(x) = ax + b$.

引理 4.3 的证明 由文 [12], 我们可以用任一 PWL 函数表示 ReLU 函数,

$$\text{ReLU}(x) = \frac{\sigma(a + \frac{r_0}{2}x) - \sigma(a + \frac{r_0}{2}x - \frac{r_0}{2}) - \sigma(a) + \sigma(a - \frac{r_0}{2})}{(\sigma'(a+) - \sigma'(a-))\frac{r_0}{2}}, \quad x \in [-1, 1], \quad (4.20)$$

其中 a 是 σ 的间断点, 满足 $\sigma'(a+) \neq \sigma'(a-)$, r_0 是 a 与其他间断点的距离. 由于误差估计由 (4.14) 式中 σ_1 构造的 \tilde{f} 得到, 且 $\sigma_1(x) = \text{ReLU}(x) - \text{ReLU}(x-1)$, 因此引理 4.1 中的结论仍然成立. 将 (4.20) 式代入 (4.17) 式, 则存在正整数 N 和常数 $c_i, y_i, \theta_i, i = 1, \dots, N$, 使得

$$\tilde{f}(x) = \sum_{i=1}^N c_i \sigma(y_i x + \theta_i) \quad (4.21)$$

对 PWL 函数 σ 成立. 此外, 若 $|f(x) - \tilde{f}(x)| = o(\frac{1}{N^2})$, 当 $N \rightarrow \infty$ 时, 有 $f''(x) = 0$, 因此 $f(x)$ 是线性函数.

定理 3.2 的证明 类似于文 [3] 中的证明方法, 将 $f(\mathbf{x})$ 用以下的方式延拓为 $[-1, 1]^d$ 上的函数 $g(\mathbf{x})$: $g(x_1, \dots, -x_k, \dots, x_d) = g(x_1, \dots, x_k, \dots, x_d)$, 且当 $\mathbf{x} \in [0, 1]^d$ 时, $g(\mathbf{x}) = f(\mathbf{x})$. 因此根据 Bochner-Riesz Means [20] 的结果, 对于 $\alpha > \frac{d-1}{2}$,

$$B_R^\alpha(f)(\mathbf{x}) = \sum_{\substack{\mathbf{m} \in \mathbb{Z}^d \\ |\mathbf{m}| \leq R}} \left(1 - \frac{|\mathbf{m}|^2}{R^2}\right)^\alpha \widehat{f}(\mathbf{m}) e^{2\pi i \mathbf{m} \cdot \mathbf{x}} \quad (4.22)$$

在 $R \rightarrow \infty$ 上一致收敛到 f , 其中 $\widehat{f}(\mathbf{m})$ 是 f 的 \mathbf{m} -th 傅里叶系数.

因此对任意的 $\varepsilon > 0$, 存在 R , 使得对任意的 $x = (x_1, \dots, x_d) \in [-1, 1]^d$, $|B_R^\alpha(f)(\mathbf{x}) - f(\mathbf{x})| < \frac{\varepsilon}{2}$ 成立. 由傅里叶系数的定义以及 $g(\mathbf{x})$ 的对称性, 我们可以将上述不等式重写为

$$\left| \sum_{|\mathbf{m}| \leq R} d_{\mathbf{m}} \cos(\mathbf{m} \cdot \mathbf{x}) - g(\mathbf{x}) \right| < \frac{\varepsilon}{2}, \quad (4.23)$$

其中 $d_{\mathbf{m}}$ 是实数. 固定 \mathbf{m} , 令 $u = \mathbf{m} \cdot \mathbf{x}$, 对任意的 $x \in [0, 1]^d$, $|u| \leq \sum_{i=1}^d |m_i| \leq \sqrt{d}R$. 因此由引理 4.3, 对于分段线性函数 σ , 有 $\tilde{f}_{\mathbf{m}}(u) = \sum_{i=1}^N c_{i,\mathbf{m}} \sigma(y_{i,\mathbf{m}} u + \theta_{i,\mathbf{m}})$, 使得 $|\tilde{f}_{\mathbf{m}}(u) - d_{\mathbf{m}} \cos(u)| = O(\frac{R^2}{N^2})$. 归纳得到 $\tilde{f} \in \mathcal{M}_N^\sigma(\mathbb{R}^d)$ 以及 $|\tilde{f}(\mathbf{x}) - g(\mathbf{x})| < \frac{\varepsilon}{2} + O(\frac{R^{3d+2}}{N^2})$.

§4.3 定理 3.3 的证明

为了记号的简洁, 这里我们将 $\sigma_2(x)$ 重写为 $m(x)$, 并定义 $m_l(x) = \underbrace{m \circ \dots \circ m}_{l \text{ 个}}(x)$, 其中符号 \circ 表示函数的复合. 我们有

$$\max_{x \in [0, 1]} \left| x - \sum_{i=1}^l \frac{1}{4^i} m_i(x) - x^2 \right| = \frac{1}{4^{l+1}}. \quad (4.24)$$

对于任意 PWL 函数 σ , 假设 a 是 σ 的间断点, 满足 $\sigma'(a+) \neq \sigma'(a-)$, r_0 是 a 与其他间断点的距离. 我们可以利用 σ 构造 m .

例如, 当 $\sigma'(a+) \neq 0$ 时,

$$m(x) = k_1 \sigma\left(a + \frac{r_0}{2}x\right) + k_2 \sigma\left(a + \frac{r_0}{2}\left(x - \frac{1}{2}\right)\right) + k_2 \sigma'(a-) \frac{r_0}{4} - (k_1 + k_2) \sigma(a), \quad x \in [0, 1], \quad (4.25)$$

其中

$$k_1 = \frac{4}{r_0 \sigma'(a+)} \cdot \frac{\sigma'(a+) + \sigma'(a-)}{\sigma'(a+) - \sigma'(a-)}, \quad k_2 = \frac{-8}{r_0 (\sigma'(a+) - \sigma'(a-))}.$$

当 $\sigma'(a+) = 0$ 时,

$$m(x) = k_1 \sigma\left(a - \frac{r_0}{2}x\right) + k_2 \sigma\left(a - \frac{r_0}{2}\left(x - \frac{1}{2}\right)\right) - (k_1 + k_2) \sigma(a), \quad x \in [0, 1], \quad (4.26)$$

其中

$$k_1 = \frac{-4}{r_0 \sigma'(a-)}, \quad k_2 = \frac{8}{r_0 \sigma'(a-)}.$$

另外,

$$|x| = k \left[\sigma\left(a + \frac{r_0}{2}x\right) + \sigma\left(a - \frac{r_0}{2}x\right) - 2\sigma(a) \right], \quad x \in [-1, 1], \quad (4.27)$$

其中

$$k = \frac{2}{r_0 (\sigma'(a+) - \sigma'(a-))}.$$

基于上述准备, 我们给出以下引理.

引理 4.4 假设 σ 是 PWL 函数, 则对任意整数 $L \geq 3$, 存在 $\hat{f} \in \mathcal{S}_{L,2}^\sigma([-1, 1])$, 使得 $|x^2 - \hat{f}(x)| \leq 4^{-L}$ 对所有 $x \in [-1, 1]$ 成立.

引理 4.4 的证明 在此, 我们仅利用 (4.25) 式给出的显式构造方法展开证明, 与利用 (4.26) 式的证明过程类似. 定义隐变量 $\{y_{1:L,1:2}\}$ 如下:

$$\begin{aligned} y_{1,1} &= \sigma\left(a + \frac{r_0}{2}x\right), \quad y_{1,2} = \sigma\left(a - \frac{r_0}{2}x\right), \\ y_{2,1} &= \sigma\left(a + \frac{r_0}{2}k(y_{1,1} + y_{1,2} - 2\sigma(a))\right), \\ y_{2,2} &= \sigma\left(a - \frac{r_0}{4} + \frac{r_0}{2}k(y_{1,1} + y_{1,2} - 2\sigma(a))\right), \\ y_{l,1} &= \sigma\left(a + \frac{r_0}{2}(k_1 y_{l-1,1} + k_2 y_{l-1,2} + c)\right), \quad l = 3, \dots, L, \\ y_{l,2} &= \sigma\left(a - \frac{r_0}{4} + \frac{r_0}{2}(k_1 y_{l-1,1} + k_2 y_{l-1,2} + c)\right), \quad l = 3, \dots, L, \end{aligned} \quad (4.28)$$

其中 $c = \frac{k_2 \sigma'(a-) r_0}{4} - (k_1 + k_2) \sigma(a)$. 由归纳可知, $|x| = k(y_{1,1} + y_{1,2} - 2\sigma(a))$, $m_l(|x|) = k_1 y_{l+1,1} + k_2 y_{l+1,2} + c$, $l = 1, \dots, L-1$ 对任意 $x \in [-1, 1]$ 成立. 令

$$\hat{f}(x) = |x| - \sum_{i=1}^{L-1} \frac{1}{4^i} m_i(|x|),$$

则 $\hat{f} \in \mathcal{S}_{L,2}^\sigma([-1, 1])$, 且 $|x^2 - \hat{f}| \leq 4^{-L}$ 对所有 $x \in [-1, 1]$ 成立.

与文 [13] 中的引理 10 稍有不同, 我们给出以下引理.

引理 4.5 假设 σ 是 PWL 函数, 则对任意的 p 阶多项式 $P_p(\mathbf{x}) = \sum_{|\mathbf{k}| \leq p} a_{\mathbf{k}} \mathbf{x}^{\mathbf{k}}$, $\mathbf{x} \in [-1, 1]^d$, $\mathbf{k} = (k_1, \dots, k_d) \in \mathbb{N}^d$, 有

$$\text{dist}(P_p, \mathcal{S}_{2\binom{p+d}{d}(p-1)L,3}^\sigma) < 2(p-1) \cdot 4^{-L} \sum_{|\mathbf{k}| \leq p} |a_{\mathbf{k}}|, \quad (4.29)$$

其中

$$\text{dist}(\phi, \mathcal{S}) = \inf_{f \in \mathcal{S}} \max_{\mathbf{x} \in [-1, 1]^d} |\phi(\mathbf{x}) - f(\mathbf{x})|.$$

引理 4.5 的证明 由于

$$xy = \left(\frac{x+y}{2}\right)^2 - \left(\frac{x-y}{2}\right)^2,$$

因此 $\text{dist}(xy, \mathcal{S}_{2L,2}^\sigma) \leq 2 \cdot 4^{-L}$. 由归纳法,

$$\text{dist}(M_p(\mathbf{x}) = x_{i_1}x_{i_2} \cdots x_{i_p}, \mathcal{S}_{2(p-1)L,3}^\sigma) \leq 2(p-1) \cdot 4^{-L}, \quad \mathbf{x} \in [-1, 1]^d,$$

其中 $i_1, \dots, i_p \in \{1, 2, \dots, d\}$. 由于 $\sum_{|\mathbf{k}| \leq p} = \binom{p+d}{d}$, (4.29) 式成立.

定理 3.3 的证明 证明过程中, $\lfloor \cdot \rfloor$ 和 $\lceil \cdot \rceil$ 分别表示向下取整和向上取整函数. 对任意的 $\delta \in (0, 1)$ 以及 $\varepsilon \in (0, \frac{1}{e})$, 令 $p = \lceil \frac{1}{\delta} \ln \frac{1}{\varepsilon} \rceil \geq 2$. 不失一般性, 假设 $\sum_{\mathbf{k}} |a_{\mathbf{k}}| = 1$, 下面证存在 $\hat{f} \in \mathcal{S}_{L,3}^\sigma(\mathbb{R}^d)$, 使得 $\|f - \hat{f}\|_\infty < 2\varepsilon$, 其中 $L = \lceil [e(\frac{1}{d\delta} \ln \frac{1}{\varepsilon} + 1)]^{\max(2d, 3)} \rceil$. 记

$$f(\mathbf{x}) = P_p(\mathbf{x}) + R(\mathbf{x}) := \sum_{|\mathbf{k}| \leq p} a_{\mathbf{k}} \mathbf{x}^{\mathbf{k}} + \sum_{|\mathbf{k}| > p} a_{\mathbf{k}} \mathbf{x}^{\mathbf{k}}.$$

由于 $\ln(1 - \delta) < -\delta$, 可以得到

$$|R(\mathbf{x})| < (1 - \delta)^p < \varepsilon, \quad \mathbf{x} \in [-1 + \delta, 1 - \delta]^d.$$

由引理 4.5, 有 $\text{dist}(P_p, \mathcal{S}_{L,3}^\sigma) < 2(p-1) \cdot 4^{-L'}$, 其中 $L' = \lfloor 2^{-1} L \binom{p+d}{p}^{-1} (p-1)^{-1} \rfloor$.

在 $d = 1$ 的情形下, $L' = \lfloor \frac{1}{2} \frac{1}{p^2-1} [e^3 (\frac{1}{\delta} \ln \frac{1}{\varepsilon} + 1)^3] \rfloor$, 因此

$$\begin{aligned} L' &> \frac{1}{2} \frac{1}{p^2-1} e^3 p^3 - 1 > 10p - 1 \\ &\geq 10 \frac{1}{\delta} \log \frac{1}{\varepsilon} - 1 > 1 + \log \left(\frac{1}{\delta} \log \frac{1}{\varepsilon} \right) + \log \frac{1}{\varepsilon}. \end{aligned}$$

从而

$$2(p-1) \cdot 4^{-L'} < 2(p-1) \left(e \cdot \frac{1}{\delta} \log \frac{1}{\varepsilon} \cdot \frac{1}{\varepsilon} \right)^{-1} < \varepsilon.$$

当 $d \geq 2$ 时, 用 Stirling 估计可以得到

$$\begin{aligned} L' &> 2^{-1} L \binom{p+d}{p}^{-1} (p-1)^{-1} - 1 \\ &> \frac{1}{2} \frac{\sqrt{2\pi d} \left(\frac{d}{e} \right)^d}{(p+d)^d} \frac{1}{p} L - 1 \\ &> \frac{\sqrt{2\pi d}}{2} \left[e \left(\frac{p-2}{d} + 1 \right) \right]^d p^{-1} - 1. \end{aligned}$$

因此

$$L' \ln 4 > 1 + 2 \frac{1}{\delta} \ln \frac{1}{\varepsilon} > 1 + \ln \left(\frac{1}{\delta} \ln \frac{1}{\varepsilon} \right) + \frac{1}{\delta} \ln \frac{1}{\varepsilon},$$

并且

$$2(p-1) \cdot 4^{-L'} < \varepsilon.$$

所以存在 $\hat{f} \in \mathcal{S}_{L,3}^{\sigma}(\mathbb{R}^d)$, 使得 $\|f - \hat{f}\|_{\infty} \leq \|f - P_p\|_{\infty} + \|P_p - \hat{f}\|_{\infty} < 2\varepsilon$.

§5 数值实验

在本节中, 我们分别进行了浅层 PWLNN 和 PWL-DNN 对 $[-1, 1]$ 上 x^2 函数逼近的数值实验.

§5.1 浅层 PWLNN

首先, 我们研究了 $\hat{f} \in \mathcal{M}_N^{\sigma}(\mathbb{R})$ 对 $x^2, x \in [-1, 1]$ 的逼近能力. 我们选择了例子 3.1 中的 Leaky ReLU 函数 ($\alpha = 0.5$) 和例子 3.2 中的 ReLU6 函数作为神经网络的激活函数.

我们根据例子 3.1 和例子 3.2 中的公式初始化浅层 PWLNN 的权重和偏差,

$$\begin{aligned} \text{Leaky ReLU : } & \begin{cases} \hat{f}(x) = x_0^2 + \sum_{i=1}^{\frac{N}{2}} [x_i^2 - x_{i-1}^2] \sigma_1 \left(\frac{x - x_{i-1}}{x_i - x_{i-1}} \right), \\ \sigma_1(x) = -\frac{\alpha}{1-\alpha} + \frac{\sigma(x)}{1-\alpha} - \frac{\sigma(x-1)}{1-\alpha}, \\ x_i = -1 + \frac{4i}{N}, \quad i = 0, \dots, \frac{N}{2}, \end{cases} \\ \text{ReLU6 : } & \begin{cases} \hat{f}(x) = x_0^2 + \sum_{i=1}^N [x_i^2 - x_{i-1}^2] \sigma_1 \left(\frac{x - x_{i-1}}{x_i - x_{i-1}} \right), \\ \sigma_1(x) = \frac{1}{6} \sigma(6x), \\ x_i = -1 + \frac{2i}{N}, \quad i = 0, \dots, N. \end{cases} \end{aligned}$$

在训练过程中, 我们以均方误差 (MSELoss) 为目标函数, 使用 Adam 优化器迭代网络参数, 设置学习率为 10^{-5} , 批量大小为 128, 并比较训练 100 轮之后的逼近性能. 训练数据集为 $\{(x_i, x_i^2)\}_{i=0}^{1024}$, 其中 $x_i = -1 + \frac{i}{1024}$.

表 1 神经元个数对浅层 PWLNN 逼近性能的影响

N		4	8	16	32	64
ReLU6	E_0	6.250E-02	1.563E-02	3.906E-03	9.766E-04	2.441E-04
	E_{100}	5.901E-02	1.361E-02	3.008E-03	7.554E-04	1.940E-04
Leaky ReLU	E_0	2.500E-01	6.250E-02	1.563E-02	3.906E-03	9.766E-04
	E_{100}	2.445E-01	5.569E-02	1.081E-02	2.732E-03	9.369E-04

实验结果如表 1 所示, 其中 E_0 表示网络初始时的逼近误差, 即 $\max_x |x^2 - \hat{f}(x)|$, E_{100} 表示网络训练 100 轮后的逼近误差.

根据表 1, 我们可以得到以下结论:

(1) 对于激活函数 ReLU6, 网络在训练前近似于 x^2 , 误差为 N^{-2} , 训练 100 轮后逼近误差减小.

(2) 对于激活函数 Leaky ReLU($\alpha = 0.5$), 网络在训练前近似于 x^2 , 误差为 $4N^{-2}$, 训练 100 轮后逼近误差减小.

这些结论与引理 4.3 中 $O(\frac{1}{N^2})$ 的估计一致. ReLU6 函数和 Leaky ReLU 函数的情况之间的逼近差异在于两者 σ_1 的表达方式不同.

§5.2 PWL-DNNs

同时, 我们探索了 $\hat{f} \in \mathcal{S}_{L,2}^\sigma(\mathbb{R})$ 对 $x^2, x \in [-1, 1]$ 的逼近能力. 参数的初始化参考引理 4.4 中的 (4.28) 式. 除了将批量大小调整为 16, 其余细节与训练浅层 PWLNN 时的情形相同. 实验结果如表 2 所示.

表 2 网络层数对 PWL-DNN 逼近性能的影响

L		2	3	4	5	6
ReLU6	E_0	6.250E-02	1.563E-02	3.906E-03	9.766E-04	2.441E-04
	E_{100}	4.291E-02	1.075E-02	2.479E-03	7.172E-04	1.882E-04
Leaky ReLU	E_0	6.250E-02	1.563E-02	3.906E-03	9.766E-04	2.441E-04
	E_{100}	4.242E-02	1.053E-02	2.641E-03	6.438E-04	1.727E-04

根据表 2, 我们可以得到以下结论:

(1) 对于激活函数 ReLU6, 网络在训练前近似 x^2 , 误差为 4^{-L} , 训练 100 轮后逼近误差减小.

(2) 对于激活函数 Leaky ReLU($\alpha = 0.5$), 网络在训练前近似 x^2 , 误差为 4^{-L} , 训练 100 轮后逼近误差减小.

这些结论与引理 4.4 中 4^{-L} 的估计一致. 由于网络表示的一致性, ReLU6 函数显示出与 Leaky ReLU 函数相似的逼近能力.

§6 讨 论

纵观全文, 我们可以看到, 文 [3] 提供了基于 Bochner-Riesz Means 的分析函数逼近的方法. 此外, 文 [13] 引入了跳跃连接神经网络结构, 为逼近解析函数提供了极大的便利. 在此基础上, 我们将上述参考文献中的激活函数扩展到 PWL 函数 (例如 Leaky ReLU 函数, ReLU6 函数, APL 函数和 SReLU 函数), 并建立浅层和深层分段线性神经网络的逼近定理. 数值实验的结果也与我们的理论结论一致.

分段线性函数将线性性与非线性性结合起来, 以实现出色的模型灵活性, 线性性亦有利于神经网络的逼近. 可以发现, 无论是浅层还是深层分段线性神经网络, 都能通过分段

线性函数的组合构造有效插值. 与具有 sigmoidal 激活函数的浅层神经网络不同, 具有 N 个隐藏神经元的浅层分段线性神经网络的逼近误差上界为 $\frac{\varepsilon}{2} + O(\frac{R^{3d+2}}{N^2})$. 通过深层分段线性神经网络的显式构造, 我们发现对于解析函数, 网络可以通过深度的堆叠而非宽度的增加以实现指数速率的逼近.

本文主要讨论浅层分段线性神经网络和深层分段线性神经网络对函数的逼近能力. 我们将在接下来的研究中讨论分段线性神经网络对泛函和算子的逼近能力. 对于深层分段线性神经网络, 我们利用分段线性函数的一个间断点构造出相应的锯齿函数. 然而, 是否可以利用分段线性函数的子区域(本质上是分段线性函数整体的非线性性)加速逼近仍需要进一步研究.

参 考 文 献

- [1] Cybenko G. Approximation by superpositions of a sigmoidal function [J]. *Mathematics of control, signals and systems*, 1989, 2(4):303–314.
- [2] Hornik K, Stinchcombe M, White H. Multilayer feedforward networks are universal approximators [J]. *Neural networks*, 1989, 2(5):359–366.
- [3] Chen T, Chen H, Liu R. A constructive proof and an extension of Cybenko’s approximation theorem [C]//Computing Science and Statistics, 1992:163–168.
- [4] Chen T, Chen H. Approximations of continuous functionals by neural networks with application to dynamic systems [J]. *IEEE Transactions on Neural networks*, 1993, 4(6):910–918.
- [5] 陈天平. 神经网络及其在系统识别应用中的逼近问题 [J]. 中国科学: A 辑, 1994, 24(1):1–7.
- [6] Chen T, Chen H, Liu R. Approximation capability in $C(\overline{\mathbf{R}}^n)$ by multilayer feedforward networks and related problems [J]. *IEEE Transactions on Neural Networks*, 1995, 6(1):25–30.
- [7] Chen T, Chen H. Universal approximation to nonlinear operators by neural networks with arbitrary activation functions and its application to dynamical systems [J]. *IEEE transactions on neural networks*, 1995, 6(4):911–917.
- [8] Back A D, Chen T. Universal approximation of multiple nonlinear operators by neural networks [J]. *Neural Computation*, 2002, 14(11):2561–2566.
- [9] Debao C. Degree of approximation by superpositions of a sigmoidal function [J]. *Approximation Theory and its Applications*, 1993, 9(3):17–28.
- [10] Barron A R. Universal approximation bounds for superpositions of a sigmoidal function [J]. *IEEE Transactions on Information theory*, 1993, 39(3):930–945.
- [11] He K, Zhang X, Ren S, et al. Deep residual learning for image recognition

- [C]//Proceedings of the IEEE conference on computer vision and pattern recognition. 2016:770–778.
- [12] Yarotsky D. Error bounds for approximations with deep ReLU networks [J]. *Neural Networks*, 2017, 94:103–114.
- [13] Wang Q. Exponential convergence of the deep neural network approximation for analytic functions [J/OL]. ArXiv:1807.00297.
- [14] Agarap A F. Deep learning using rectified linear units (relu) [J/OL]. ArXiv:1803.08375.
- [15] Maas A L, Hannun A Y, Ng A Y. Rectifier nonlinearities improve neural network acoustic models [C]//Proceedings of the International Conference on Machine Learning, 2013, 30(1):3.
- [16] Agostinelli F, Hoffman M, Sadowski P, et al. Learning activation functions to improve deep neural networks [J/OL]. arXiv:1412.6830.
- [17] Jin X, Xu C, Feng J, et al. Deep learning with s-shaped rectified linear activation units [C]//Proceedings of the AAAI Conference on Artificial Intelligence, 2016, 30(1).
- [18] Sandler M, Howard A, Zhu M, et al. Mobilenetv 2: inverted residuals and linear bottlenecks [C]//Proceedings of the IEEE conference on computer vision and pattern recognition, 2018:4510–4520.
- [19] Tao Q, Li L, Huang X, et al. Piecewise linear neural networks and deep learning [J]. *Nature Reviews Methods Primers*, 2022, 2(1):42.
- [20] Grafakos L. Classical fourier analysis [M]. New York: Springer-Verlag, 2008.

Approximation Theory on Piece Wise Linear Neural Networks

WU Xinyu¹ CHEN Tianping² LU Wenlian³

¹School of Mathematical Sciences, Fudan University, Shanghai 200433, China.
E-mail: xywu19@fudan.edu.cn

²School of Mathematical Sciences, Fudan University, Shanghai 200433, China; Shanghai Center for Mathematical Sciences, Fudan University, Shanghai 200438, China; Shanghai Key Laboratory for Contemporary Applied Mathematics, Shanghai 200433, China. E-mail:tchen@fudan.edu.cn

³Corresponding author. School of Mathematical Sciences, Fudan University, Shanghai 200433, China ; Shanghai Center for Mathematical Sciences, Fudan University, Shanghai 200438, China ; Shanghai Key Laboratory for Contemporary Applied Mathematics, Shanghai 200433, China.
E-mail:wenlian@fudan.edu.cn

Abstract With the wide application of Piece Wise Linear (PWL for short) functions, this paper attempts to address the approximation theory on Piece Wise Linear Neural Networks (PWLNNS for short) for both shallow networks and deep neural networks (DNNs for short). The authors extend the universal approximation theorem of three-layer Multi-Layer Perceptrons (MLPs for short) with PWL functions and bound the error by the number of hidden neurons. The authors give an explicit way of constructing sawtooth functions from PWL functions, and thus prove analytic functions can be approximated at an exponentially convergent rate by stacking depth rather than increasing width. Numerical experiments are also provided to verify the conclusions.

Keywords Piece wise linear, Neural network, Approximation theorem

2000 MR Subject Classification 41A30

The English translation of this paper will be published in

Chinese Journal of Contemporary Mathematics, Vol. 45 No. 1, 2024

by ALLERTON PRESS, INC., USA