

# 神经网络基础

## 全连接神经网络

$$f(x; \theta) = f_L \circ \phi_\theta(x) \\ = W^L \circ \phi_\theta(x) + b^L$$

随机特征回归:

$$\sum c_i \phi_i(x) \Rightarrow W^L \phi(x) + b^L$$

GELU: Gaussian Error Linear Unit

$$\text{GELU}(x) = x \Phi(x) = x \mathbb{P}(X \leq x)$$

$$\begin{aligned} \text{erf}(z) &= \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt \\ &= x \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \\ &= x \left( \int_0^x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx + \frac{1}{2} \right) \quad \boxed{\frac{x}{\sqrt{2}} = z} \\ &= x \left( \int_0^{\frac{x}{\sqrt{2}}} \frac{1}{\sqrt{2\pi}} e^{-z^2} d\sqrt{2}z + \frac{1}{2} \right) \\ &= \frac{x}{2} \left( \int_0^{\frac{x}{\sqrt{2}}} \frac{2}{\sqrt{\pi}} e^{-z^2} dz + 1 \right) \end{aligned}$$

silu (Sigmoid Linear Unit)

$$\text{silu}(x) = x \cdot \text{sigmoid}(x)$$

## Wierstrass 逼近定理

$$(x+1-x)^n = \sum_k C_n^k x^k (1-x)^{n-k}$$

$$\text{考虑 } \sum C_n^k x^k (1-x)^{n-k} f\left(\frac{k}{n}\right) \rightarrow f(x) \text{ in } [0,1]$$

## 通用逼近定理 (universal approximation)

引理:  $g$  是  $\mathbb{R} \rightarrow \mathbb{R}$  的  $p$ -Lipschitz 函数, 2 层  $\lceil \frac{p}{\varepsilon} \rceil$  个神经元  
使用  $\mathbb{1}_{x \geq 0}$  的神经网络, 能对  $g$  在  $[0,1]$  上达到  $\varepsilon$  近似。

证明

$$\hat{g} = \sum_{i=1}^m a_i \mathbb{1}(x - b_i), \text{ 其中}$$

$$b_i = \frac{(i-1)\varepsilon}{p}, \quad m = \lceil \frac{p}{\varepsilon} \rceil, \quad \text{取}$$

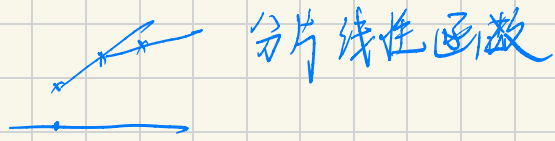
$$a_1 = g(b_1), \quad a_2 = g(b_2) - g(b_1) \dots, \quad a_m = g(b_m) - g(b_{m-1})$$

$$\text{对 } b_j \leq x < b_{j+1}$$

$$\begin{aligned} \hat{g}(x) &= \sum_{i=1}^j a_i \mathbb{1}(x - b_i) \\ &= \sum_{i=1}^j a_i = g(b_j) \end{aligned}$$

$$|\hat{g}(x) - g(x)| = |g(b_j) - g(x)| \leq (b_j - x)p = \frac{\varepsilon}{p} \cdot p = \varepsilon$$

⇒ sigmoid, ReLU - - - - -



引理:  $g$  是  $\mathbb{R}^d \rightarrow \mathbb{R}$  的  $L$ -Lipschitz 函数, 2层神经元使用  $\mathbb{1}_{x \in R_i}$  作为激活函数的神经网络, 能对  $g$  在  $[0, 1]^d$  上达到  $\varepsilon$  近似

其中  $R_1, R_2, \dots, R_m$  是  $[0, 1]^d$  上的边长不超过  $\delta$  的小方块

选取  $x_i \in R_i$

$$f(x) = \sum_{i=1}^m \alpha_i \mathbb{1}_{x \in R_i} \quad \alpha_i = g(x_i)$$

那么

$$\begin{aligned} \sup |g(x) - f(x)| &= \sup_i \sup_{x \in R_i} |g(x) - f(x)| \\ &\leq \sup_i \sup_{x \in R_i} \underbrace{|g(x) - g(x_i)|}_{\leq \delta} + \underbrace{|g(x_i) - f(x)|}_0 \end{aligned}$$

激活函数去近似  $\mathbb{1}_{x \in R_i}$  ( $\delta \sim \text{error}$ )

$$\Rightarrow D = \left(\frac{1}{\delta}\right)^d \quad \text{维度灾难}$$

目标: 指数收敛  $\text{error} = e^{-\left(\frac{1}{\delta}\right)^d}$  计算量  $\left(\frac{1}{\delta}\right)^d = -\log \text{error}$

直观理解: Fourier transform

$$\hat{f}(x) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} f(x) e^{-i\omega^T x} dx$$

$$f(x) = \int \hat{f}(\omega) e^{i\omega^T x} d\omega \quad (\hat{f}(\omega) = |\hat{f}(\omega)| e^{i b(\omega)})$$

$$= \int |\hat{f}(\omega)| e^{i b(\omega) + i\omega^T x} d\omega$$

$$= \int |\hat{f}(\omega)| \cos(b(\omega) + \omega^T x) d\omega$$

定义:  $C_f = \int |\hat{f}(\omega)| d\omega$ ,  $p(\omega) = \frac{|\hat{f}(\omega)|}{C_f}$   
概率密度函数

$$f(x) = C_f E_{\omega \sim p} [\cos b(\omega) + \omega^T x]$$

定义  $f_m = \frac{1}{m} \sum_{j=1}^m C_f \cos(\omega_j^T x + b(\omega_j)) \quad \omega_j \sim p$

$$z_j = C_f \cos(\omega_j^T x + b(\omega_j))$$

$$E_{\omega} [z_j - f(x)] = 0$$

$$E_{\omega} [(z_j - f(x))^2] = E_{\omega} z_j^2 - f(x)^2$$

$$\leq E_{\omega} z_j^2$$

$$\leq C_f^2$$

对任意概率密度函数  $p$ , 我们有

$$\begin{aligned}
& E_w \|f_m - f\|_{L_2(P)}^2 \\
&= E_w \left\| \frac{1}{m} \sum_j (z_j - f) \right\|_{L_2(P)}^2 \quad z_j \text{ 独立} \\
&= E_w \int \frac{1}{m^2} \left( \sum_j (z_j - f) \right)^2 P(x) dx \\
&= \int \frac{m}{m^2} E_w (z_j - f)^2 P(x) dx
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{m^2} E_{x \sim P} \left[ m E_w (z_j - f)^2 \right] \\
&\leq \frac{C_f^2}{m}
\end{aligned}$$

那么  $\exists w$ , 使

$$\|f_m - f\|_{L_2(P)}^2 \leq \frac{C_f^2}{m}$$

与维度无关

i)  $C_f = \text{poly}(d)$

ii)  $\cos$  一般不用子激活函数

iii) 一般考虑紧支集

Barron space 近似

$$f_e|_{\Omega} = f \quad (\text{假设 } 0 \in \Omega)$$

$$f_e(x) = \int e^{i\omega^T x} \hat{f}_e(\omega) d\omega$$

$$f(x) - f(0) = \int (e^{i\omega^T x} - 1) \hat{f}_e(\omega) d\omega$$

$$= \int \frac{e^{i\omega^T x} - 1}{\|\omega\|} \|\omega\| \hat{f}_e(\omega) d\omega$$

$$= \int \frac{\cos(\omega^T x + b(\omega)) - \cos(b(\omega))}{\|\omega\|} \|\omega\| \hat{f}_e(\omega) d\omega$$

g(x, \omega)

其中  $\|\omega\| = \sup_{x \in \Omega} |\omega^T x|$

定义  $h(t; \omega) = \frac{\cos(\|\omega\|^2 t + b(\omega)) - \cos(b(\omega))}{\|\omega\|}$

$$h\left(\frac{\omega^T x}{\|\omega\|}; \omega\right) = g(x; \omega) \quad (\hat{\omega} = \frac{\omega}{\|\omega\|})$$

由  $\|\omega\|$  的定义,  $\frac{\omega^T x}{\|\omega\|} \in [-1, 1]$

$$\frac{1}{C_f} \int \|\omega\| |\hat{f}_e(\omega)| d\omega = C_f' < +\infty$$

$$\rho \sim \|\omega\| |\hat{f}_e(\omega)| / C_f'$$

$$f(x) - f(0) = C_f' E_{w \sim \rho} [g(x, w)]$$

$$\text{我们有 } g(x, w) = h(\hat{w}^T x; w)$$

$$\text{其中 } \hat{w}^T x \in [-1, 1]$$

注意  $h(-1; w), h'(t; w) \in [-1, 1]$ , 因为

$$h(t; w) = \frac{\cos(\|\omega\|t + b(w)) - \cos(b(w))}{\|\omega\|}$$

$$= \frac{-2 \sin\left(\frac{\|\omega\|t}{2}\right) \sin\left(\frac{\|\omega\|t}{2} + b(w)\right)}{\|\omega\|}$$

$$h'(t; w) = -\sin(\|\omega\|t + b(w))$$

我们有

$$\begin{aligned}h(t; \omega) &= h(-1; \omega) + \int_{-1}^t h'(s; \omega) ds \\ &= h(-1; \omega) + \int_{-1}^1 h'(s; \omega) H(t-s; \omega) ds\end{aligned}$$

Heaviside:

$$H(t) = 1 \quad (t \geq 0)$$



$$\begin{aligned}f(x) &= \underbrace{f(0) + C_f' E_{w \sim p}[h(-1; \omega)]}_{\text{与 } x \text{ 无关的常数}} \\ &\quad + \underbrace{C_f' E_{w \sim p} E_{s \sim U[-1,1]}[h'(s; \omega) H(\hat{w}^T x - s)]}\end{aligned}$$

$$f_D(x) = a_0 + C_f' \frac{1}{D} \sum_{j=1}^D \underbrace{h'(s_j; w_j) H(\hat{w}_j^T x - s_j)}_{z_j}$$

$$\begin{aligned}E_w \|f_D - f\|_{L_2(\mathcal{P})}^2 \\ = E_w \|C_f' \left( \frac{1}{D} \sum_{j=1}^D z_j - E \right)\|_{L_2(\mathcal{P})}^2 \quad E = E z_j\end{aligned}$$

$$= \frac{C_f'^2}{D} E_{x \sim \mathcal{P}} E_w (z_j - E)^2$$

$$\leq \frac{C_f'^2}{D} E z_j^2 \leq \frac{C_f'}{D}$$



## 非凸优化

$$y_k = x_k + \beta(x_k - x_{k-1})$$

$$x_{k+1} = y_k - \alpha \nabla f(x_k)$$

(Polyak's heavy-ball method)

$$y_k = x_k + \beta(\beta_k - x_{k-1})$$

$$x_{k+1} = y_k - \alpha \nabla f(y_k)$$

(Nesterov's accelerated gradient descent)

In Adam, bias correction

$$m_0 = v_0 = 0$$

$$\text{当 } g_t = g$$

$$m_1 = (1 - \beta_1)g \quad m_2 = (1 - \beta_1)g \cdot \beta_1 + (1 - \beta_1)g \\ = (1 - \beta_1^2)g$$

$$m_k = (1 - \beta_1^k)g$$