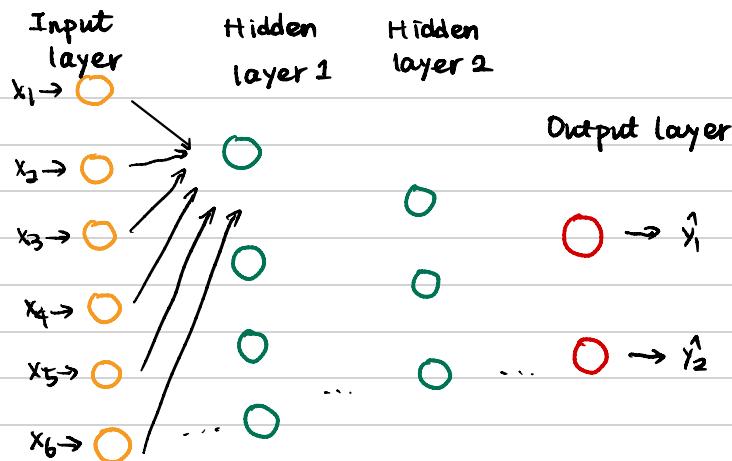


Chapter 4 Feedforward Neural Networks





$$Y = F_{w,b}(X) + \varepsilon$$

$$\hat{Y}(X) = F_{w,b}(X) = f_n^{(4)}(w, b^{(4)}) \circ \dots \circ f_1^{(1)}(w, b^{(1)})(X)$$

	notation	definition
超参数	L	神经网络的层数
参数	M_L	第 L 层神经元的个数
	$f_L(\cdot)$	第 L 层神经元的激活函数
参数	$W^{L-1} \in \mathbb{R}^{M_L \times M_{L-1}}$	第 $L-1$ 层到第 L 层的权重矩阵
	$b^{L-1} \in \mathbb{R}^{M_L}$	第 $L-1$ 层到第 L 层的偏置
活性值	$z^{(L)} \in \mathbb{R}^{M_L}$	第 L 层神经元的净输入(净活性值)
	$a^{(L)} \in \mathbb{R}^{M_L}$	第 L 层神经元的输出(活性值)

激活函数	函数	导数
Logistic 函数	$f(x) = \frac{1}{1 + \exp(-x)}$	$f'(x) = f(x)(1 - f(x))$
Tanh 函数	$f(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$	$f'(x) = 1 - f(x)^2$

ReLU 函数	$f(x) = \max(0, x)$	$f'(x) = I(x > 0)$
---------	---------------------	--------------------

ELU 函数	$f(x) = \max(0, x) + \min(0, \gamma(\exp(x) - 1))$	$f'(x) = I(x > 0) + I(x \leq 0) \cdot \gamma \exp(x)$
--------	--	---

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$$

$$a^{(l)} = f_l(z^{(l)})$$

$$a^{(l)} = f_l(W^{(l)} a^{(l-1)} + b^{(l)})$$

$$\begin{aligned} x = a^{(0)} &\rightarrow z^{(1)} \rightarrow a^{(1)} \rightarrow z^{(2)} \rightarrow a^{(2)} \\ &\cdots \rightarrow a^{(L-1)} \rightarrow z^{(L)} \rightarrow a^{(L)} = \phi(x; W, b) \end{aligned}$$

Ex. $\sigma: \mathbb{R} \rightarrow B \subset \mathbb{R}$

$$f_{W^{(l)}, b^{(l)}}: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$f(v) = W^{(l)} \sigma^{(l-1)}(v) + b^{(l)}, \quad W^{(l)} \in \mathbb{R}^{m \times n}, \quad b^{(l)} \in \mathbb{R}^m$$

$$\text{eg. } f(v) = w \tanh(v) + b$$

Ex. One hidden layer, $\sigma^{(1)}$ is the identity function

$$\begin{aligned} \hat{y}(x) &= W^{(2)}(W^{(1)}x + b^{(1)}) + b^{(2)} = W^{(2)}W^{(1)}x + W^{(2)}b^{(1)} + b^{(2)} \\ &= \tilde{w}x + \tilde{b} \end{aligned}$$

Universal Approximation Theorem

对于具有线性输出层和至少一个使用“挤压”性质的激活函数的隐藏层组成的前馈神经网络，只要其隐藏层神经元的数量足够，它可以以任意的精度来近似任何一个定义在实数空间的有界闭集函数。

多分类问题

0 1

0 2

0 3 - $p(y=c|x)$

$$P(y|x) \Rightarrow \hat{y} = \text{Softmax}(z^{(L)})$$

:

0 C

$$L(y, \hat{y}) = -y^\top \log \hat{y}$$

$$R(w, b) = \frac{1}{N} \sum_{n=1}^N L(y^{(n)}, \hat{y}^{(n)}) + \frac{\lambda}{2} \|w\|_F^2$$

$$\frac{\partial R(w, b)}{\partial w^{(l)}}$$

$$\frac{\partial R(w, b)}{\partial b^{(l)}} \quad (\text{梯度下降})$$

$$\text{Chain Rule } y = f^5(f^4(f^3(f^2(f^1(x)))) \quad \frac{\partial y}{\partial x} = \frac{\partial f^1}{\partial x} \cdot \frac{\partial f^2}{\partial f^1} \cdots \frac{\partial f^5}{\partial f^4}$$

反向传播算法，自动微分 (AD)

An arrangement of $n \geq p$ hyperplanes in \mathbb{R}^p has at most

$$\sum_{j=0}^p \binom{n}{j} \text{ convex regions}$$

$$R(W, b) = \frac{1}{N} \sum_{n=1}^N L(y^{(n)}, \hat{y}^{(n)}) + \frac{1}{2} \lambda \|W\|_F^2$$

$$\frac{\partial R(W, b)}{\partial W^{(l)}}, \quad \frac{\partial R(W, b)}{\partial b^{(l)}}$$

$$\downarrow \frac{\partial L(y, \hat{y})}{\partial W_{ij}^{(l)}}$$

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)} \quad R(a^{(l)}) = f_l(z^{(l)})$$

$$\frac{\partial L(y, \hat{y})}{\partial W_{ij}^{(l)}} = \frac{\partial z^{(l)}}{\partial W_{ij}^{(l)}} \quad \frac{\partial L(y, \hat{y})}{\partial z^{(l)}}$$

$$\frac{\partial z^{(l)}}{\partial W_{ij}^{(l)}} = \left[\frac{\partial z_i^{(l)}}{\partial W_{1j}^{(l)}} \cdots \frac{\partial z_i^{(l)}}{\partial W_{mj}^{(l)}} \cdots \frac{\partial z_m^{(l)}}{\partial W_{ij}^{(l)}} \right]$$

$$= [0, \dots, a_j^{(l-1)}, \dots, 0] \\ \cong \text{Di}(a_j^{(l-1)}) \quad R \in \mathbb{R}^{1 \times M_l}$$

$$\frac{\partial L(y, \hat{y})}{\partial b^{(l)}} = \frac{\partial z^{(l)}}{\partial b^{(l)}} \cdot \frac{\partial L(y, \hat{y})}{\partial z^{(l)}} \quad \downarrow g^{(l)}$$

$$\frac{\partial z}{\partial b^{(l)}} = I_{M_l} \in \mathbb{R}^{M_l \times M_l}$$

$$g^{(l)} \cong \frac{\partial L(y, \hat{y})}{\partial z^{(l)}}$$

$$= \frac{\partial a^{(l)}}{\partial z^{(l)}} \cdot \frac{\partial z^{(l+1)}}{\partial a^{(l)}} \cdot \frac{\partial L(y, \hat{y})}{\partial z^{(l+1)}}$$

$$= \text{diag}(f_L'(z^{(l)})) \cdot (W^{(l+1)})^T \cdot g^{(l+1)}$$

$$= f_L'(z^{(l)}) \odot ((W^{(l+1)})^T \cdot g^{(l+1)}) \in \mathbb{R}^{M_L}$$

$$\Rightarrow \frac{\partial L(y, \hat{y})}{\partial w_{ij}^{(l)}} = [0, \dots, a_j^{(l-1)}, \dots, 0] [s_1^{(l)}, \dots, s_i^{(l)}, \dots, s_{M_l}^{(l)}]^T$$
$$= s_i^{(l)} \cdot a_j^{(l-1)}$$

$$\frac{\partial L(y, \hat{y})}{\partial w^{(l)}} = s^{(l)} (a^{(l-1)})^T \in \mathbb{R}^{M_L \times M_{l-1}}$$

$$\frac{\partial L(y, \hat{y})}{\partial b^{(l)}} = s^{(l)} \in \mathbb{R}^{M_L}$$