



第4讲：前馈神经网络

师 佳

化学化工学院

化学工程与生物工程系

纲要

➤神经网络概述

- 概述
- 神经元
- 网络结构

➤前馈神经网络

- 模型
- 通用近似定理
- 反向传播算法
- 自动微分计算
- 编程实现方法
- 优化问题

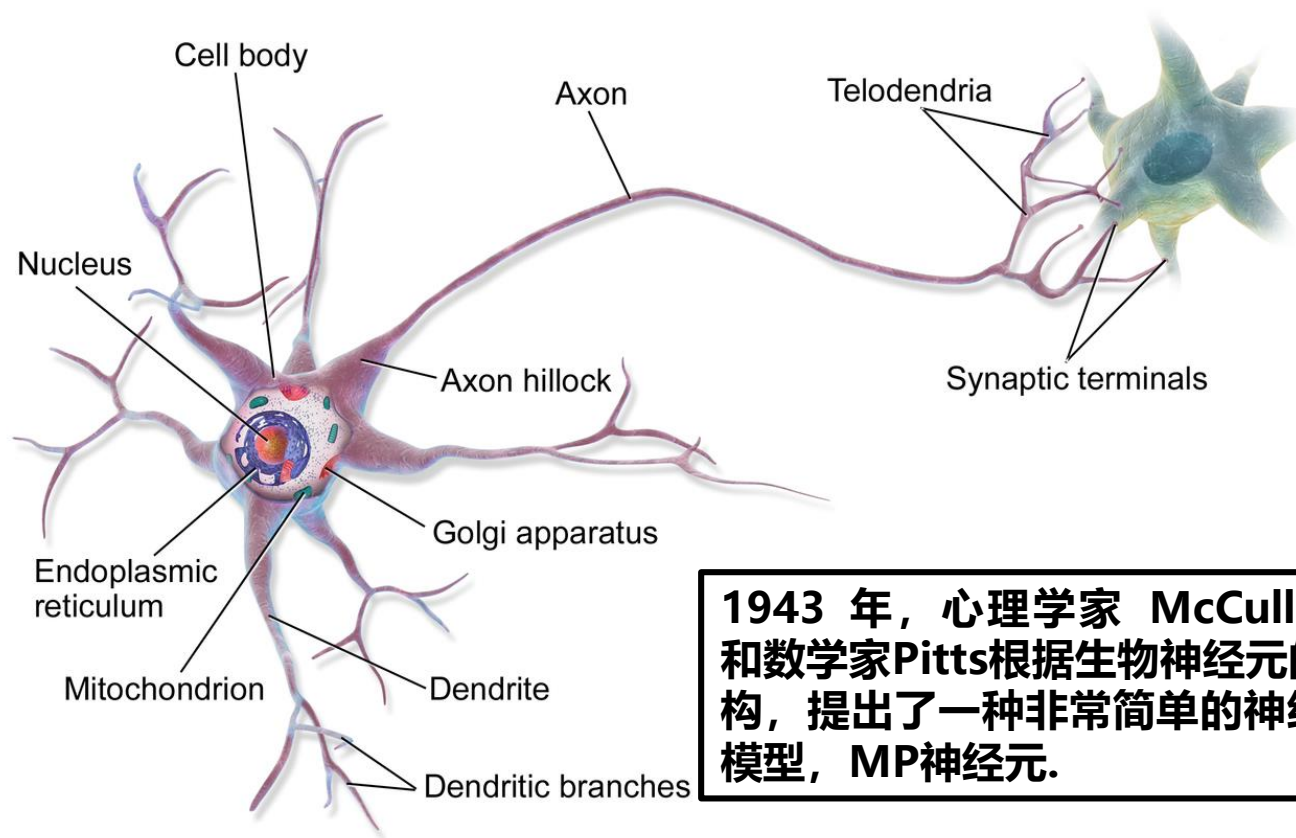
神经网络概述

➤概述

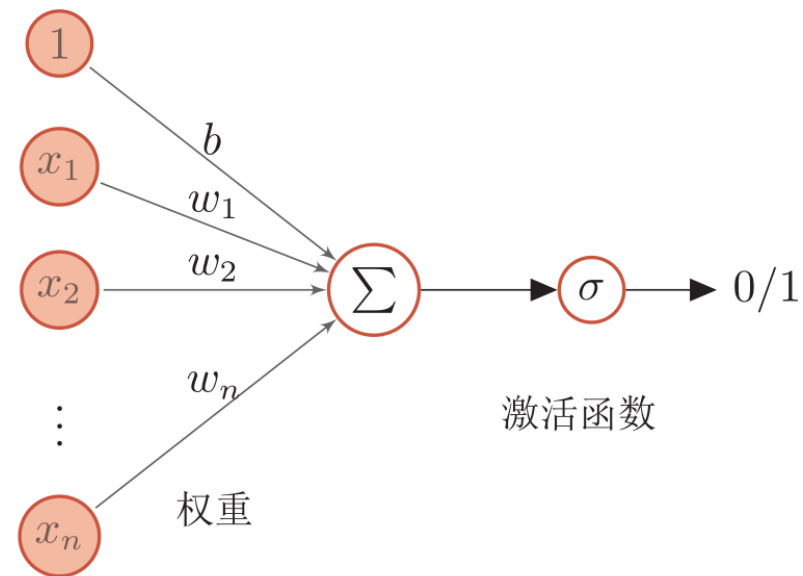
- 神经网络最早是作为一种主要的**连接主义模型**。诞生于20世纪80年代中后期，是最流行的一种**分布式并行处理(Parallel Distributed Processing, PDP)模型**，具有3个主要特征：
 - 信息表示是分布式的(非局部的)
 - 记忆和知识是存储在单元之间的连接上
 - 通过逐渐改变单元之间的连接强度来学习新的知识
- 从机器学习的角度来看，神经网络一般可以看作一个**非线性模型**，其基本组成单元为具有**非线性激活函数**的神经元，通过**大量神经元之间的连接**，使得神经网络成为一种**高度非线性的模型**。
- 在引入**误差反向传播**来改进其学习能力之后，神经网络也越来越多地应用在各种机器学习任务上。

神经网络概述

➤ 神经元



1943 年，心理学家 McCulloch 和数学家 Pitts 根据生物神经元的结构，提出了一种非常简单的神经元模型，MP 神经元。



➤ **净输入 (Net Input) :**

$$z = \sum_{d=1}^D w_d x_d + b = \mathbf{w}^T \mathbf{x} + b$$

➤ **激活函数 (Activation Function) :**

$$a = f(z)$$

神经网络概述

➤ 神经元

➤ 激活函数必须具备的性质：

- **连续并可导**（允许少数点上不可导）的**非线性函数**。可导的激活函数可以直接利用数值优化的方法来学习网络参数。
- 激活函数及其导函数要**尽可能的简单**，有利于提高网络计算效率。
- 激活函数的导函数的值域要在一个**合适的区间内**，不能太大也不能太小，否则会影响训练的效率和稳定性。

➤ 常用的激活函数

➤ Sigmoid型函数

1) Logistic函数

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

2) Tanh函数

$$\tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$$

➤ 性质：

- 1) 饱和函数
- 2) 连续可导
- 3) Tanh是零中心化
- 4) Logistic具有偏置偏移

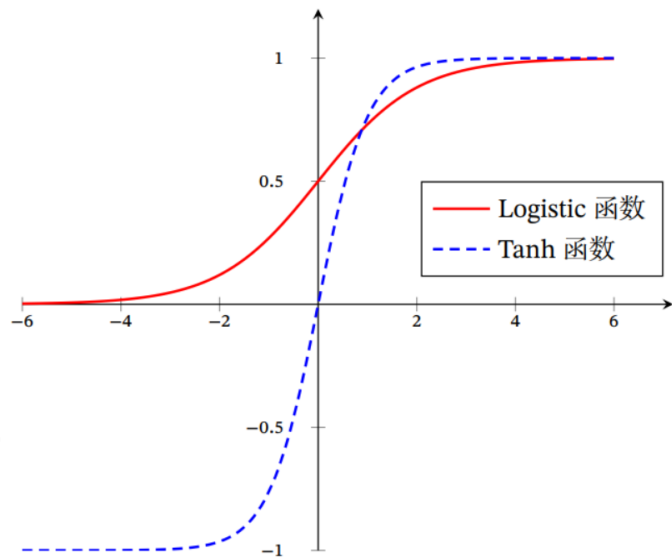


图 4.2 Logistic 函数和 Tanh 函数

神经网络概述

➤ 神经元

➤ 常用的激活函数

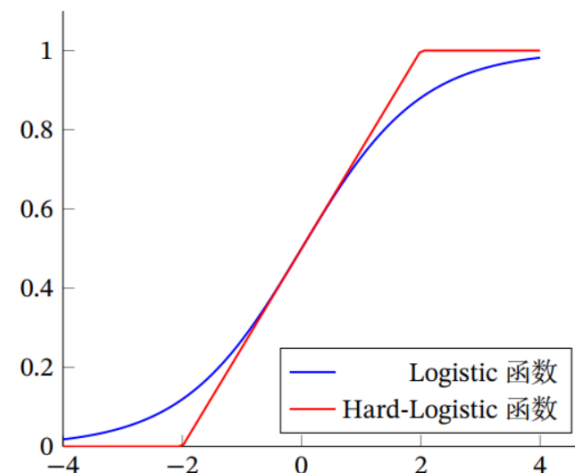
➤ Sigmoid型函数

1) Hard-Logistic函数

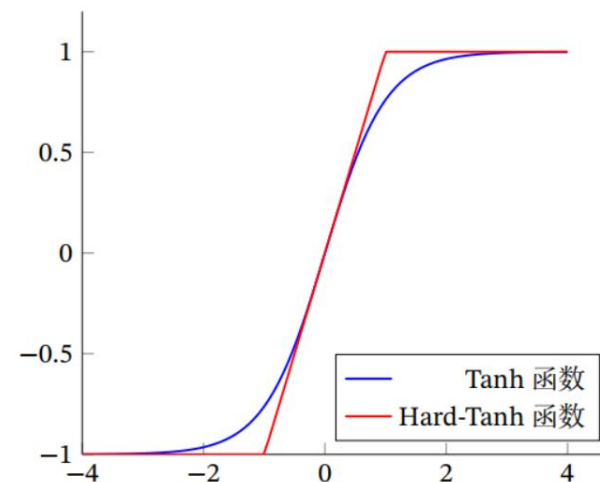
$$\text{hard-logistic}(x) = \max(\min(0.25x + 0.5, 1), 0)$$

2) Hard-Tanh函数

$$\text{hard-tanh}(x) = \max(\min(x, 1), -1)$$



(a) Hard Logistic 函数



(b) Hard Tanh 函数

神经网络概述

➤ 神经元

➤ 常用的激活函数

- **ReLU**(Rectified Linear Unit, 修正线性单元) 函数

$$\text{ReLU}(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases} = \max(0, x)$$

- **带泄露的ReLU(Leak ReLU)**函数

$$\text{LeakyReLU}(x) = \begin{cases} x & \text{if } x > 0 \\ \gamma x & \text{if } x \leq 0 \end{cases}$$
$$= \max(0, x) + \gamma \min(0, x)$$

- **带参数的ReLU(Parametric ReLU)**函数

$$\text{PReLU}_i(x) = \begin{cases} x & \text{if } x > 0 \\ \gamma_i x & \text{if } x \leq 0 \end{cases}$$

- **Softplus**函数

$$\text{Softplus}(x) = \log(1 + \exp(x))$$

- **ELU**(Exponential Linear Unit, 指数线性单元)函数

$$\text{ELU}(x) = \begin{cases} x & \text{if } x > 0 \\ \gamma(\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$
$$= \max(0, x) + \min(0, \gamma(\exp(x) - 1))$$

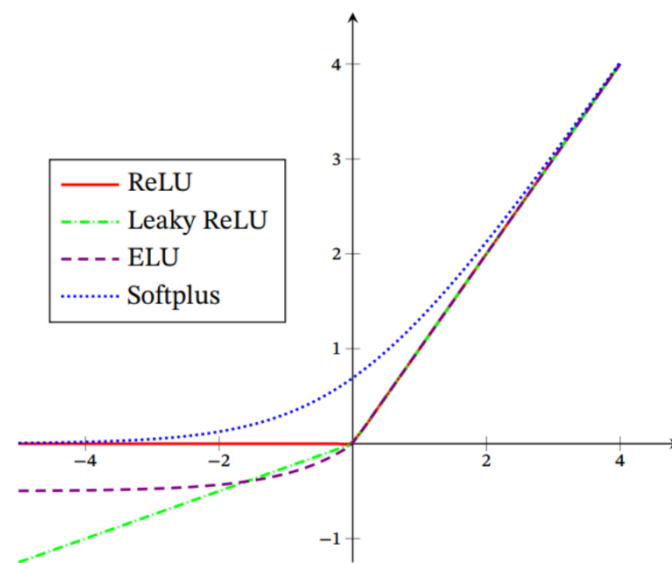


图 4.4 ReLU、Leaky ReLU、ELU 以及 Softplus 函数

神经网络概述

➤ 神经元

➤ 常用的激活函数

➤ Swish 函数

$$\text{swish}(x) = x\sigma(\beta x)$$

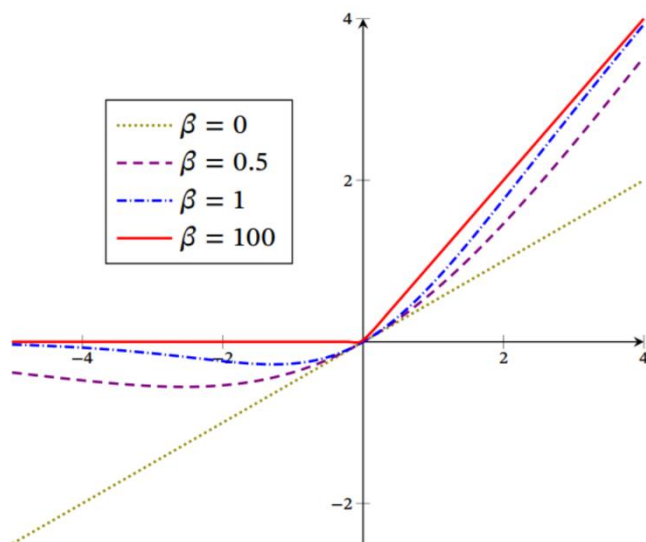


图 4.5 Swish 函数

线性函数和ReLU函数之间的非线性插值函数，其非线性程度由参数 β 控制

➤ GELU(Gaussian Error Linear Unit, 高斯误差线性单元)

$$\text{GELU}(x) = xP(X \leq x)$$

其中 $P(X \leq x)$ 是高斯分布 $\mathcal{N}(\mu, \sigma^2)$ 的累积分布函数

➤ Maxout单元

$$\text{maxout}(\mathbf{x}) = \max_{k \in [1, K]} (z_k)$$

$$z_k = \mathbf{w}_k^\top \mathbf{x} + b_k$$

Maxout激活函数可以看作任意凸函数的分段线性近似，并且在有限的点上是不可微的。

神经网络概述

➤ 神经元

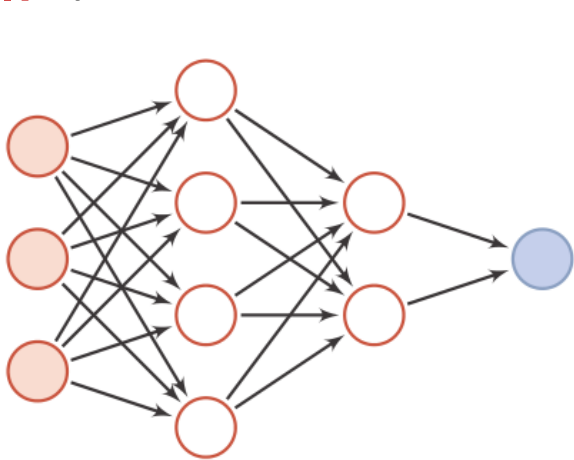
➤ 常用激活函数的导数

激活函数	函数	导数	课后练习：证明
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	$f(x) = \max(0, x) + \min(0, \gamma(\exp(x) - 1))$	$f'(x) = I(x > 0) + I(x \leq 0) \cdot \gamma \exp(x)$	
SoftPlus 函数	$f(x) = \log(1 + \exp(x))$	$f'(x) = \frac{1}{1+\exp(-x)}$	

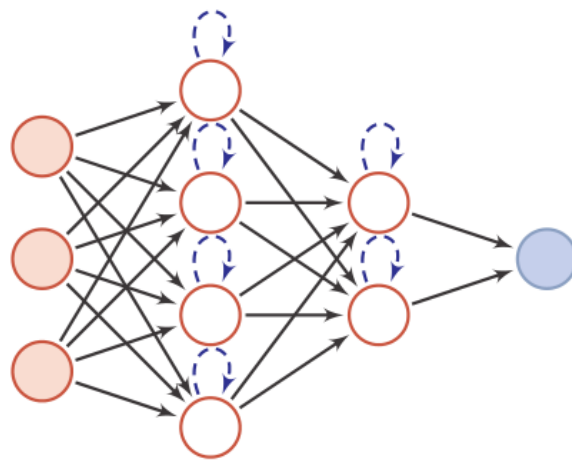
神经网络概述

➤ 网络结构

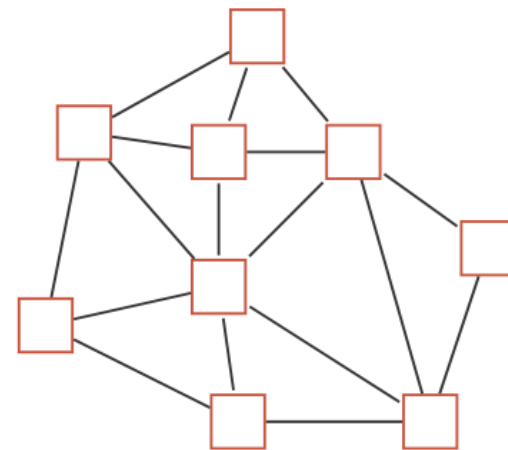
- 人工神经网络由神经元模型构成，这种由许多神经元组成的信息处理网络具有并行分布结构。



(a) 前馈网络



(b) 记忆网络



(c) 图网络

前馈网络可以看作一个函数，通过简单非线性函数的多次复合，实现输入空间到输出空间的复杂映射。

具有更强的计算和记忆能力.包括循环神经网络、Hopfield网络、玻尔兹曼机、受限玻尔兹曼机等。

是前馈网络和记忆网络的泛化，包含很多不同的实现方式，比如图卷积网络、图注意力网络、消息传递神经网络等。

前馈神经网络

模型

➤ **前馈神经网络**中，各神经元分别属于不同的**层**。每一层的神经元可以接收前一层神经元的信号，并产生信号输出到下一层。

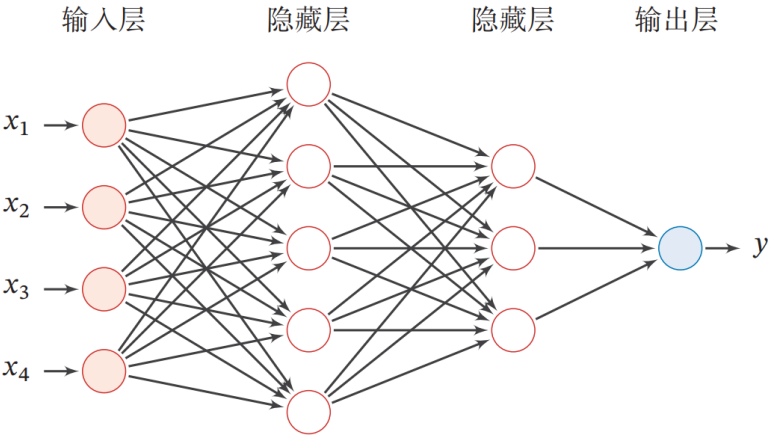


图 4.7 多层前馈神经网络

表 4.1 前馈神经网络的记号

记号	含义
L	神经网络的层数
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
$\mathbf{W}^{(l)} \in \mathbb{R}^{M_l \times M_{l-1}}$	第 $l-1$ 层到第 l 层的权重矩阵
$\mathbf{b}^{(l)} \in \mathbb{R}^{M_l}$	第 $l-1$ 层到第 l 层的偏置
$\mathbf{z}^{(l)} \in \mathbb{R}^{M_l}$	第 l 层神经元的净输入（净活性值）
$\mathbf{a}^{(l)} \in \mathbb{R}^{M_l}$	第 l 层神经元的输出（活性值）

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

$$\mathbf{a}^{(l)} = f_l(\mathbf{z}^{(l)}). \quad \mathbf{a}^{(0)} = \mathbf{x}$$

$$\mathbf{x} = \mathbf{a}^{(0)} \rightarrow \mathbf{z}^{(1)} \rightarrow \mathbf{a}^{(1)} \rightarrow \mathbf{z}^{(2)} \rightarrow \dots \rightarrow \mathbf{a}^{(L-1)} \rightarrow \mathbf{z}^{(L)} \rightarrow \mathbf{a}^{(L)} = \phi(\mathbf{x}; \mathbf{W}, \mathbf{b})$$

前馈神经网络

➤ 通用近似定理

定理 4.1 – 通用近似定理 (Universal Approximation Theorem) [Cybenko, 1989; Hornik et al., 1989]: 令 $\phi(\cdot)$ 是一个非常数、有界、单调递增的连续函数, \mathcal{I}_D 是一个 D 维的单位超立方体 $[0, 1]^D$, $C(\mathcal{I}_D)$ 是定义在 \mathcal{I}_D 上的连续函数集合. 对于任意给定的一个函数 $f \in C(\mathcal{I}_D)$, 存在一个整数 M , 和一组实数 $v_m, b_m \in \mathbb{R}$ 以及实数向量 $\mathbf{w}_m \in \mathbb{R}^D, m = 1, \dots, M$, 以至于我们可以定义函数

$$F(\mathbf{x}) = \sum_{m=1}^M v_m \phi(\mathbf{w}_m^\top \mathbf{x} + b_m), \quad (4.37)$$

作为函数 f 的近似实现, 即

$$|F(\mathbf{x}) - f(\mathbf{x})| < \epsilon, \forall \mathbf{x} \in \mathcal{I}_D, \quad (4.38)$$

其中 $\epsilon > 0$ 是一个很小的正数.

➤ 只要隐藏层神经元的数量足够, 可以以任意的精度来近似任何一个定义在实数空间 \mathbb{R}^D 中的有界闭集函数.

➤ 神经网络可以作为一个“万能”函数来使用, 可以用来进行复杂的特征转换, 或逼近一个复杂的条件分布.

前馈神经网络

通用近似定理

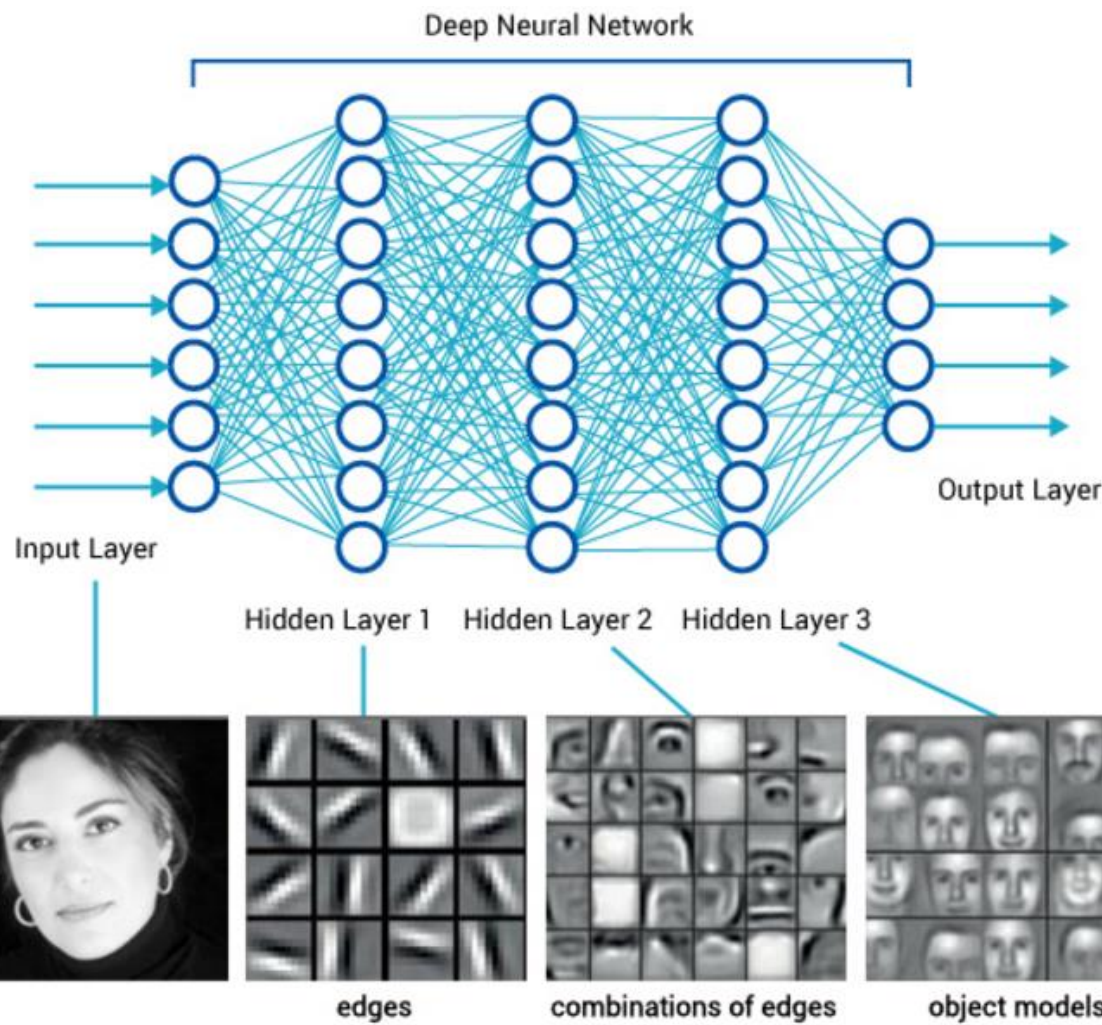
深层前馈神经网络

$$y = f_{W_n}^n \left(\cdots f_{W_3}^3 \left(f_{W_2}^2 \left(f_{W_1}^1 (x) \right) \right) \cdots \right) = f_W(x)$$

多层前馈神经网络可以看作一个非线性复合函数 $\phi: \mathbb{R}^D \rightarrow \mathbb{R}^{D'}$ ，将输入 $x \in \mathbb{R}^D$ 映射到输出 $\phi(x) \in \mathbb{R}^{D'}$ 。因此，多层前馈神经网络也可以看成是一种特征转换方法，其输出 $\phi(x)$ 作为最后一层分类器的输入进行分类。

$$\hat{y} = g(\phi(x); \theta)$$

$$p(y = 1|x) = a^{(L)} \quad \hat{y} = \text{softmax}(z^{(L)})$$



前馈神经网络

➤ 反向传播算法

➤ 学习准则：结构化风险函数

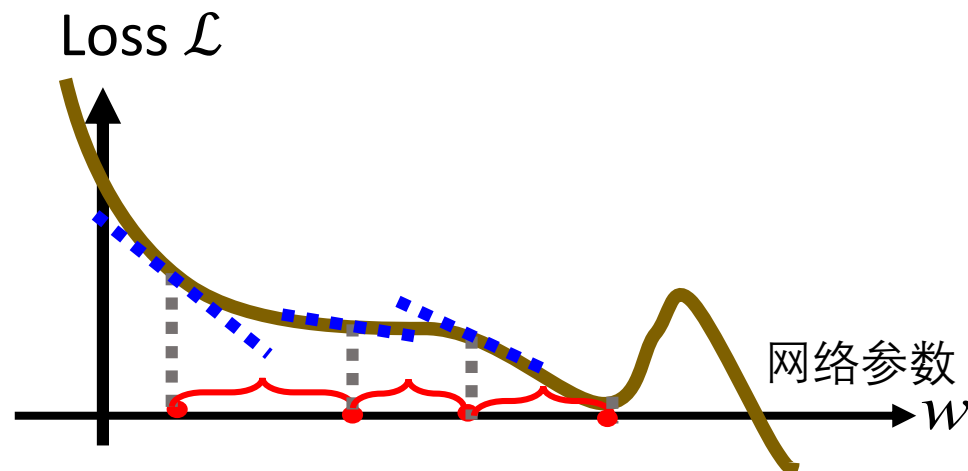
$$\mathcal{R}(\mathbf{W}, \mathbf{b}) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}(\mathbf{y}^{(n)}, \hat{\mathbf{y}}^{(n)}) + \frac{1}{2} \lambda \|\mathbf{W}\|_F^2$$

$$\|\mathbf{W}\|_F^2 = \sum_{l=1}^L \sum_{i=1}^{M_l} \sum_{j=1}^{M_{l-1}} (w_{ij}^{(l)})^2$$

➤ 优化算法：梯度下降法

$$\begin{aligned} \mathbf{W}^{(l)} &\leftarrow \mathbf{W}^{(l)} - \alpha \frac{\partial \mathcal{R}(\mathbf{W}, \mathbf{b})}{\partial \mathbf{W}^{(l)}} \\ &= \mathbf{W}^{(l)} - \alpha \left(\frac{1}{N} \sum_{n=1}^N \left(\frac{\partial \mathcal{L}(\mathbf{y}^{(n)}, \hat{\mathbf{y}}^{(n)})}{\partial \mathbf{W}^{(l)}} \right) + \lambda \mathbf{W}^{(l)} \right) \end{aligned}$$

$$\begin{aligned} \mathbf{b}^{(l)} &\leftarrow \mathbf{b}^{(l)} - \alpha \frac{\partial \mathcal{R}(\mathbf{W}, \mathbf{b})}{\partial \mathbf{b}^{(l)}} \\ &= \mathbf{b}^{(l)} - \alpha \left(\frac{1}{N} \sum_{n=1}^N \left(\frac{\partial \mathcal{L}(\mathbf{y}^{(n)}, \hat{\mathbf{y}}^{(n)})}{\partial \mathbf{b}^{(l)}} \right) \right) \end{aligned}$$



前馈神经网络

➤ 反向传播算法

➤ 梯度计算

分母布局 分子布局

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial w_{ij}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}} = \boxed{\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}} \boxed{\frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}}}$$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{b}^{(l)}} = \boxed{\frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{b}^{(l)}}} \boxed{\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}}$$

$$= \begin{bmatrix} \frac{\partial z_1^{(l)}}{\partial b_1^{(l)}} & \frac{\partial z_2^{(l)}}{\partial b_1^{(l)}} & \dots & \frac{\partial z_{M_l}^{(l)}}{\partial b_1^{(l)}} \\ \frac{\partial z_1^{(l)}}{\partial b_2^{(l)}} & \frac{\partial z_2^{(l)}}{\partial b_2^{(l)}} & \dots & \frac{\partial z_{M_l}^{(l)}}{\partial b_2^{(l)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_1^{(l)}}{\partial b_{M_l}^{(l)}} & \frac{\partial z_2^{(l)}}{\partial b_{M_l}^{(l)}} & \dots & \frac{\partial z_{M_l}^{(l)}}{\partial b_{M_l}^{(l)}} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial z_1^{(l)}}{\partial w_{ij}^{(l)}} & \frac{\partial z_2^{(l)}}{\partial w_{ij}^{(l)}} & \dots & \frac{\partial z_{M_l}^{(l)}}{\partial w_{ij}^{(l)}} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial L(\mathbf{y}, \hat{\mathbf{y}})}{\partial z_1^{(l)}} & \frac{\partial L(\mathbf{y}, \hat{\mathbf{y}})}{\partial z_2^{(l)}} & \dots & \frac{\partial L(\mathbf{y}, \hat{\mathbf{y}})}{\partial z_{M_l}^{(l)}} \end{bmatrix}^T$$

$$= \begin{bmatrix} \frac{\partial L(\mathbf{y}, \hat{\mathbf{y}})}{\partial z_1^{(l)}} & \frac{\partial L(\mathbf{y}, \hat{\mathbf{y}})}{\partial z_2^{(l)}} & \dots & \frac{\partial L(\mathbf{y}, \hat{\mathbf{y}})}{\partial z_{M_l}^{(l)}} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial z_1^{(l)}}{\partial w_{ij}^{(l)}} & \frac{\partial z_2^{(l)}}{\partial w_{ij}^{(l)}} & \dots & \frac{\partial z_{M_l}^{(l)}}{\partial w_{ij}^{(l)}} \end{bmatrix}^T$$

前馈神经网络

➤ 反向传播算法

➤ 梯度计算

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

$$\begin{aligned} \frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}} &= \left[\frac{\partial z_1^{(l)}}{\partial w_{ij}^{(l)}}, \dots, \frac{\partial z_i^{(l)}}{\partial w_{ij}^{(l)}}, \dots, \frac{\partial z_{M_l}^{(l)}}{\partial w_{ij}^{(l)}} \right] \\ &= \left[0, \dots, \frac{\partial (\mathbf{w}_{i:}^{(l)} \mathbf{a}^{(l-1)} + b_i^{(l)})}{\partial w_{ij}^{(l)}}, \dots, 0 \right] \\ &= \left[0, \dots, a_j^{(l-1)}, \dots, 0 \right] \\ &\triangleq \mathbb{I}_i(a_j^{(l-1)}) \in \mathbb{R}^{1 \times M_l}, \end{aligned}$$

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

$$\frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{a}^{(l)}} = (\mathbf{W}^{(l+1)})^\top \in \mathbb{R}^{M_l \times M_{l+1}}$$

$$\frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{z}^{(l)}} = \frac{\partial f_l(\mathbf{z}^{(l)})}{\partial \mathbf{z}^{(l)}} = \text{diag}(f'_l(\mathbf{z}^{(l)})) \in \mathbb{R}^{M_l \times M_l}$$

$$\delta^{(l)} \triangleq \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}} \quad \text{误差项}$$

$$\begin{aligned} &= \frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{z}^{(l)}} \cdot \frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{a}^{(l)}} \cdot \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l+1)}} \\ &= \text{diag}(f'_l(\mathbf{z}^{(l)})) \cdot (\mathbf{W}^{(l+1)})^\top \cdot \delta^{(l+1)} \\ &= f'_l(\mathbf{z}^{(l)}) \odot ((\mathbf{W}^{(l+1)})^\top \delta^{(l+1)}) \in \mathbb{R}^{M_l} \end{aligned}$$

反向传播 (BP)

前馈神经网络

➤ 反向传播算法

➤ 梯度计算

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial w_{ij}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}$$



$$= \mathbb{I}_i(a_j^{(l-1)}) \delta^{(l)}$$

$$= [0, \dots, a_j^{(l-1)}, \dots, 0] [\delta_1^{(l)}, \dots, \delta_i^{(l)}, \dots, \delta_{M_l}^{(l)}]^\top$$

$$= \delta_i^{(l)} a_j^{(l-1)},$$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}^{(l)}} = \delta^{(l)} (\mathbf{a}^{(l-1)})^\top \in \mathbb{R}^{M_l \times M_{l-1}}$$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{b}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{b}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}} = \delta^{(l)} \in \mathbb{R}^{M_l}$$

$$\delta^{(l)} = f_l'(\mathbf{z}^{(l)}) \odot ((\mathbf{W}^{(l+1)})^\top \delta^{(l+1)}) \in \mathbb{R}^{M_l}$$

算法 4.1 使用反向传播算法的随机梯度下降训练过程

输入: 训练集 $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$, 验证集 \mathcal{V} , 学习率 α , 正则化系数 λ , 网络层数 L , 神经元数量 $M_l, 1 \leq l \leq L$.

```
1 随机初始化  $\mathbf{W}, \mathbf{b}$ ;  
2 repeat  
3   对训练集  $\mathcal{D}$  中的样本随机重排序;  
4   for  $n = 1 \dots N$  do  
5     从训练集  $\mathcal{D}$  中选取样本  $(\mathbf{x}^{(n)}, y^{(n)})$ ;  
6     前馈计算每一层的净输入  $\mathbf{z}^{(l)}$  和激活值  $\mathbf{a}^{(l)}$ , 直到最后一层;  
7     反向传播计算每一层的误差  $\delta^{(l)}$ ; // 公式 (4.63)  
      // 计算每一层参数的导数  
8      $\forall l, \quad \frac{\partial \mathcal{L}(\mathbf{y}^{(n)}, \hat{\mathbf{y}}^{(n)})}{\partial \mathbf{W}^{(l)}} = \delta^{(l)} (\mathbf{a}^{(l-1)})^\top$ ; // 公式 (4.68)  
9      $\forall l, \quad \frac{\partial \mathcal{L}(\mathbf{y}^{(n)}, \hat{\mathbf{y}}^{(n)})}{\partial \mathbf{b}^{(l)}} = \delta^{(l)}$ ; // 公式 (4.69)  
      // 更新参数  
10     $\mathbf{W}^{(l)} \leftarrow \mathbf{W}^{(l)} - \alpha (\delta^{(l)} (\mathbf{a}^{(l-1)})^\top + \lambda \mathbf{W}^{(l)})$ ;  
11     $\mathbf{b}^{(l)} \leftarrow \mathbf{b}^{(l)} - \alpha \delta^{(l)}$ ;  
12  end  
13 until 神经网络模型在验证集  $\mathcal{V}$  上的错误率不再下降;  
    输出:  $\mathbf{W}, \mathbf{b}$ 
```

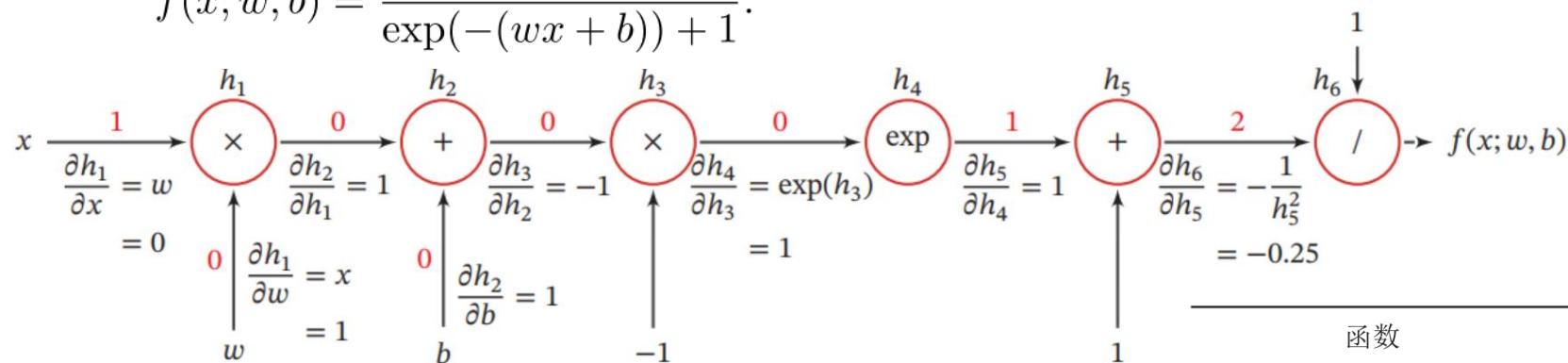
注意: 输出层的误差项要根据损失函数的定义来计算

前馈神经网络

➤ 自动微分计算

➤ **自动微分**的基本原理是所有的数值计算可以**分解为一些基本操作**，然后利用**链式法则**来自动计算一个复合函数的梯度

➤ 例： $f(x; w, b) = \frac{1}{\exp(-(wx + b)) + 1}$.



➤ 当 $x=1, w=0, b=0$ 时

$$\begin{aligned} \frac{\partial f(x; w, b)}{\partial w} \Big|_{x=1, w=0, b=0} &= \frac{\partial f(x; w, b)}{\partial h_6} \frac{\partial h_6}{\partial h_5} \frac{\partial h_5}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w} \\ &= 1 \times -0.25 \times 1 \times 1 \times -1 \times 1 \times 1 \\ &= 0.25. \end{aligned}$$

➤ **反向模式**和反向传播的计算梯度的方式相同

函数	导数	
$h_1 = x \times w$	$\frac{\partial h_1}{\partial w} = x$	$\frac{\partial h_1}{\partial x} = w$
$h_2 = h_1 + b$	$\frac{\partial h_2}{\partial h_1} = 1$	$\frac{\partial h_2}{\partial b} = 1$
$h_3 = h_2 \times -1$	$\frac{\partial h_3}{\partial h_2} = -1$	
$h_4 = \exp(h_3)$	$\frac{\partial h_4}{\partial h_3} = \exp(h_3)$	
$h_5 = h_4 + 1$	$\frac{\partial h_5}{\partial h_4} = 1$	
$h_6 = 1/h_5$	$\frac{\partial h_6}{\partial h_5} = -\frac{1}{h_5^2}$	

前馈神经网络

➤ 自动微分计算

- 输出层的误差项要根据损失函数的定义来计算
- 常见激活函数的导数

表 4.3 常见激活函数及其导数

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	$f(x) = \max(0, x) + \min(0, \gamma(\exp(x) - 1))$	$f'(x) = I(x > 0) + I(x \leq 0) \cdot \gamma \exp(x)$
SoftPlus 函数	$f(x) = \log(1 + \exp(x))$	$f'(x) = \frac{1}{1+\exp(-x)}$

前馈神经网络

➤ 编程实现方法

- 前馈神经网络的训练过程分为以下三步：
 - 前向计算每一层的状态和激活值，直到最后一层
 - 反向计算每一层的参数的偏导数
 - 按照梯度下降算法更新参数
- **静态计算图** (Tensorflow)：在编译时构建计算图，计算图构建好之后在程序运行时不能改变，便于并行计算，灵活性差。
- **动态计算图** (PyTorch)：在程序运行时动态构建，灵活度高，不便于并行计算。
- 实现流程



```
from keras.models import Sequential
from keras.layers import Dense, Activation
from keras.optimizers import SGD
```

```
model = Sequential()
model.add(Dense(output_dim=64, input_dim=100))
model.add(Activation("relu"))
model.add(Dense(output_dim=10))
model.add(Activation("softmax"))
```

```
model.compile(loss='categorical_crossentropy',
              optimizer='sgd', metrics=['accuracy'])
```

```
model.fit(X_train, Y_train, nb_epoch=5, batch_size=32)
```

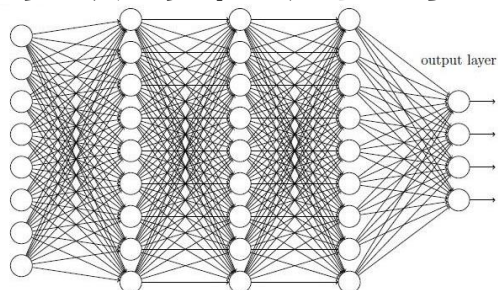
```
loss = model.evaluate(X_test, Y_test, batch_size=32)
```

前馈神经网络

➤ 优化问题

➤ 神经网络的参数学习比线性模型要更加困难：

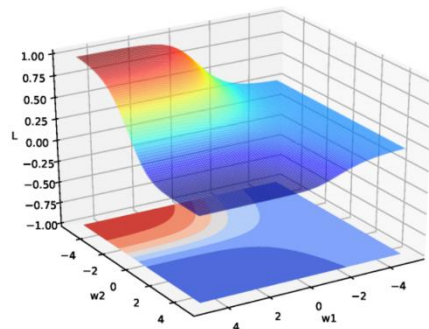
➤ 参数过多，影响训练



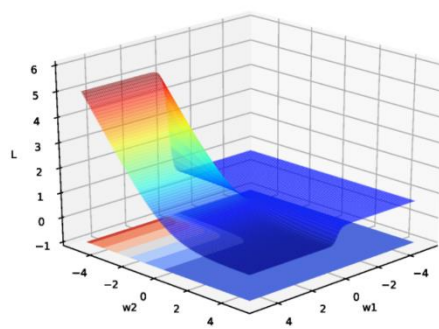
参数数量

$$P = \sum_{l=1}^{N-1} M_l + 1 \quad M_{l+1}$$

➤ 非凸优化问题：即存在局部最优而非全局最优解，影响迭代



(a) 平方误差损失



(b) 交叉熵损失

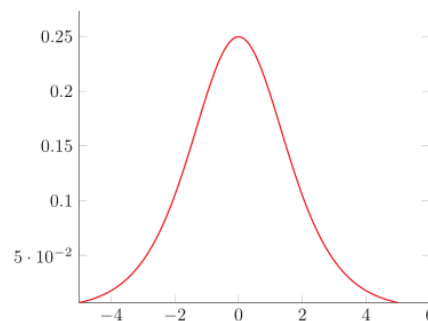
图 4.10 神经网络 $y = \sigma(w_2 \sigma(w_1 x))$ 的损失函数

➤ 梯度消失问题，造成下层参数比较难调

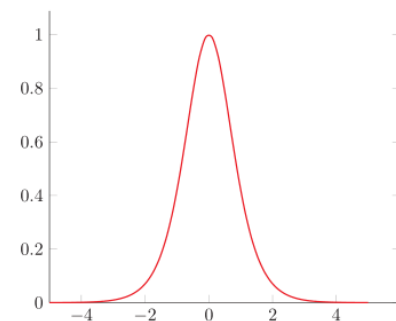
$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}^{(l)}} = \delta^{(l)} (\mathbf{a}^{(l-1)})^\top \in \mathbb{R}^{M_l \times M_{l-1}}$$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{b}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{b}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}} = \delta^{(l)} \in \mathbb{R}^{M_l}$$

$$\delta^{(l)} = f'_l(\mathbf{z}^{(l)}) \odot ((\mathbf{W}^{(l+1)})^\top \delta^{(l+1)}) \in \mathbb{R}^{M_l}$$



(a) logistic 函数的导数



(b) tanh 函数的导数

➤ 可解释性差

课后作业

➤完成

- 习题4-1
- 习题4-2
- 习题4-4

➤思考

- 习题4-7
- 习题4-8
- 习题4-9



谢谢

