

2-2 只有一个隐藏层的神经网络 II

(多个训练样本)

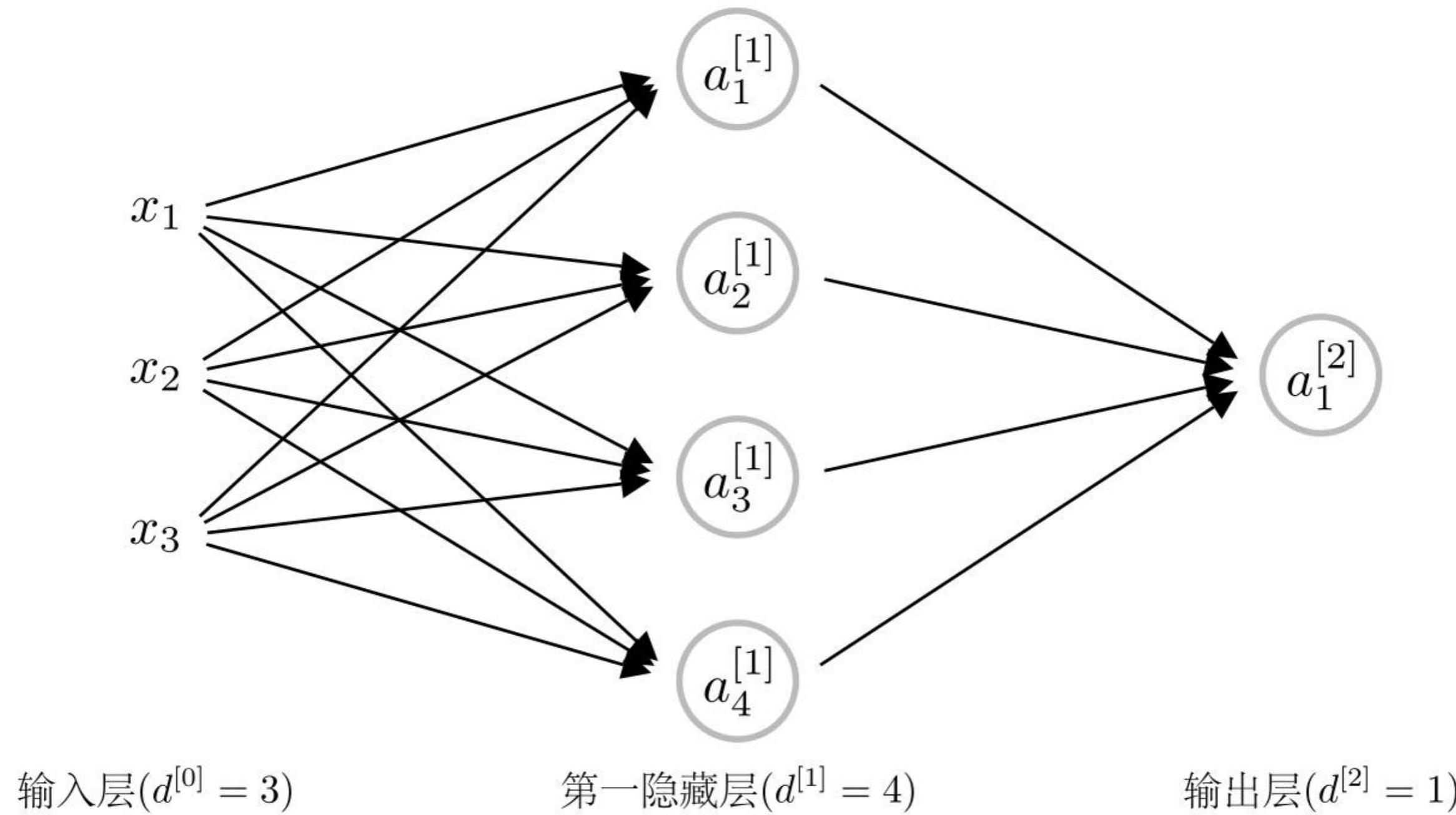
王中雷

厦门大学王亚南经济研究院和经济学院, 2025

内容摘要

1. 前向传播
2. 后向传播
3. 批量梯度下降法

回顾



回顾

1. 回顾

- L : 神经网络模型的层数
- $d^{[l]}$: 第 l 层神经元的个数 ($l = 0, \dots, L$)
- $\mathbf{a}^{[l]} = (a_1^{[l]}, \dots, a_{d^{[l]}}^{[l]})^T \in \mathbb{R}^{d^{[l]} \times 1}$
- $\mathbf{W}^{[l]} = (\mathbf{w}_1^{[l]}, \dots, \mathbf{w}_{d^{[l]}}^{[l]})^T \in \mathbb{R}^{d^{[l]} \times d^{[l-1]}}$
- $\mathbf{b}^{[l]} = (b_1^{[l]}, \dots, b_{d^{[l]}}^{[l]})^T \in \mathbb{R}^{d^{[l]} \times 1}$

回顾

1. 如果我们只有一个训练样本 (\mathbf{x}, y) , 前向传播的计算如下

$$\mathbf{z}^{[1]} = \mathbf{b}^{[1]} + \mathbf{W}^{[1]}\mathbf{x}$$

$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$

$$z^{[2]} = b^{[2]} + \mathbf{W}^{[2]}\mathbf{a}^{[1]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

2. 代价函数为

$$\mathcal{J} = \mathcal{L} = - \left\{ y \log a^{[2]} + (1 - y) \log (1 - a^{[2]}) \right\}$$

前向传播

1. 假设我们有 n 个训练样本 $\{(\mathbf{x}_i, y_i) : i = 1, \dots, n\}$

2. 记 $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times d}$

3. 基于向量化的前向传播

$$\mathbf{Z}^{[1]} = (\mathbf{b}^{[1]})^T + \mathbf{X}(\mathbf{W}^{[1]})^T \in \mathbb{R}^{n \times d^{[1]}}$$

$$\mathbf{A}^{[1]} = \sigma(\mathbf{Z}^{[1]}) \in \mathbb{R}^{n \times d^{[1]}}$$

$$\mathbf{Z}^{[2]} = b^{[2]} + \mathbf{A}^{[1]}(\mathbf{W}^{[2]})^T \in \mathbb{R}^{n \times d^{[2]}}$$

$$\mathbf{A}^{[2]} = \sigma(\mathbf{Z}^{[2]}) \in \mathbb{R}^{n \times d^{[2]}}$$

- Python 中的广播机制被用于计算 $\mathbf{A}^{[1]}$ 以及 $\mathbf{A}^{[2]}$

前向传播

1. 对于以上例子，我们有 $d^{[0]} = 3$, $\textcolor{blue}{d}^{[1]} = 4$, $\textcolor{red}{d}^{[2]} = 1$

2. 因此，我们有

$$\mathbf{Z}^{[1]} = (\mathbf{b}^{[1]})^T + \mathbf{X}(\mathbf{W}^{[1]})^T \in \mathbb{R}^{n \times \textcolor{blue}{4}}$$

$$\mathbf{A}^{[1]} = \sigma(\mathbf{Z}^{[1]}) \in \mathbb{R}^{n \times \textcolor{blue}{4}}$$

$$\mathbf{Z}^{[2]} = b^{[2]} + \mathbf{A}^{[1]}(\mathbf{W}^{[2]})^T \in \mathbb{R}^{n \times \textcolor{red}{1}}$$

$$\mathbf{A}^{[2]} = \sigma(\mathbf{Z}^{[2]}) \in \mathbb{R}^{n \times \textcolor{red}{1}}$$

前向传播

1. 代价函数为

$$\mathcal{J} = -\frac{1}{n} \sum_{i=1}^n \left\{ y_i \log a_i^{[2]} + (1 - y_i) \log (1 - a_i^{[2]}) \right\}$$

- $a_i^{[2]}$: 对应于特征向量 \mathbf{x}_i 的“激活值”

后向传播

1. 记 $d \cdot = \frac{\partial \mathcal{J}}{\partial \cdot}$
2. 回顾: 对于一个样本点 (\mathbf{x}, y) , 我们有
$$db^{[2]} = a^{[2]} - y, \quad d\mathbf{W}^{[2]} = (a^{[2]} - y)(\mathbf{a}^{[1]})^T$$
$$d\mathbf{b}^{[1]} = (a^{[2]} - y)\mathbf{D}(\mathbf{W}^{[2]})^T, \quad d\mathbf{W}^{[1]} = (a^{[2]} - y)\mathbf{D}(\mathbf{W}^{[2]})^T \mathbf{x}^T$$
 - $\mathbf{D} = \text{diag}(\{a_j^{[1]}(1 - a_j^{[1]}): j = 1, \dots, d^{[1]}\})$
3. 我们有
$$dz^{[2]} = a^{[2]} - y, \quad d\mathbf{z}^{[1]} = (a^{[2]} - y)\mathbf{D}(\mathbf{W}^{[2]})^T = dz^{[2]}\mathbf{D}(\mathbf{W}^{[2]})^T$$

后向传播

1. 此外，我们还有

$$\mathbf{D}(\mathbf{W}^{[2]})^T = (\mathbf{W}^{[2]})^T \circ \sigma'(\mathbf{z}^{[1]})$$

- $\sigma'(\mathbf{z}^{[1]}) = (a_1^{[1]}(1 - a_1^{[1]}), \dots, a_{d^{[1]}}^{[1]}(1 - a_{d^{[1]}}^{[1]}))^T$
- \circ : Hadamard 乘积 (矩阵对应元素相乘)

2. 我们有

$$\begin{aligned} dz^{[2]} &= a^{[2]} - y, & db^{[2]} &= dz^{[2]}, & d\mathbf{W}^{[2]} &= dz^{[2]}(\mathbf{a}^{[1]})^T \\ dz^{[1]} &= dz^{[2]} \mathbf{D}(\mathbf{W}^{[2]})^T, & db^{[1]} &= dz^{[1]}, & d\mathbf{W}^{[1]} &= dz^{[1]} \mathbf{x}^T \end{aligned}$$

后向传播

1. 一般地，如果我们有 n 个训练样本，则我们有

$$d\boldsymbol{\theta} = \frac{1}{n} \sum_{i=1}^n \frac{\partial \mathcal{J}(y_i, a_i)}{\partial \boldsymbol{\theta}}$$

2. 因此，为了得到 $d\boldsymbol{\theta}$ ，我们只需要将对应的偏导数求均值即可

3. 利用向量化

$$\sum_{k=1}^K a_k b_k = \mathbf{ab}$$

- $\mathbf{a} = (a_1, \dots, a_K)$ 以及 $\mathbf{b} = (b_1, \dots, b_K)^T$
- a_k 一个数值或者一个列向量
- b_k 一个数值或者一个行向量

后向传播

1. 基于向量化，我们有

$$d\mathbf{Z}^{[2]} = n^{-1}(\mathbf{A}^{[2]} - \mathbf{Y}) \quad d\mathbf{b}^{[2]} = (d\mathbf{Z}^{[2]})^T \mathbf{1} \quad d\mathbf{W}^{[2]} = (d\mathbf{Z}^{[2]})^T \mathbf{A}^{[1]}$$

$$d\mathbf{Z}^{[1]} = d\mathbf{Z}^{[2]} \mathbf{W}^{[2]} \circ \sigma'(\mathbf{Z}^{[1]}) \quad d\mathbf{b}^{[1]} = (d\mathbf{Z}^{[1]})^T \mathbf{1} \quad d\mathbf{W}^{[1]} = (d\mathbf{Z}^{[1]})^T \mathbf{X}$$

- $\sigma'(z)$ ：对应于 $\sigma(z)$ 的偏导数
- $\sigma'(\mathbf{Z}^{[1]})$ ：利用广播机制，计算 $\mathbf{Z}^{[1]}$ 每个元素对应的导数值

批量梯度下降法

步骤1. 随机初始化 $\boldsymbol{\theta}^{(0)}$

步骤2. 基于 $\boldsymbol{\theta}^{(t)}$ 计算

$$\nabla \mathcal{J}(\boldsymbol{\theta}^{(t)}) = \frac{\partial \mathcal{J}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}^{(t)})$$

步骤3. 更新参数

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \alpha \nabla \mathcal{J}(\boldsymbol{\theta}^{(t)})$$

步骤4. 回到步骤 2 直至收敛