

2-1 只有一个隐藏层的神经网络 I

(一个训练样本)

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. 回顾逻辑回归模型

2. 前向传播

3. 后向传播

4. 批量梯度下降法

直观

1. 只用一个训练样本 (\mathbf{x}, y) 展示基本概念 ($y \in \{0, 1\}$)
2. 神经网络模型是特征向量 \mathbf{x} 的函数，对应的模型参数为 θ
3. 利用（批量）梯度下降法得到模型参数 θ 的估计
4. 核心步骤是

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

- $\theta^{(t)}$: 当前模型参数
- 代价函数及其对参数的导数是什么呢?

直观

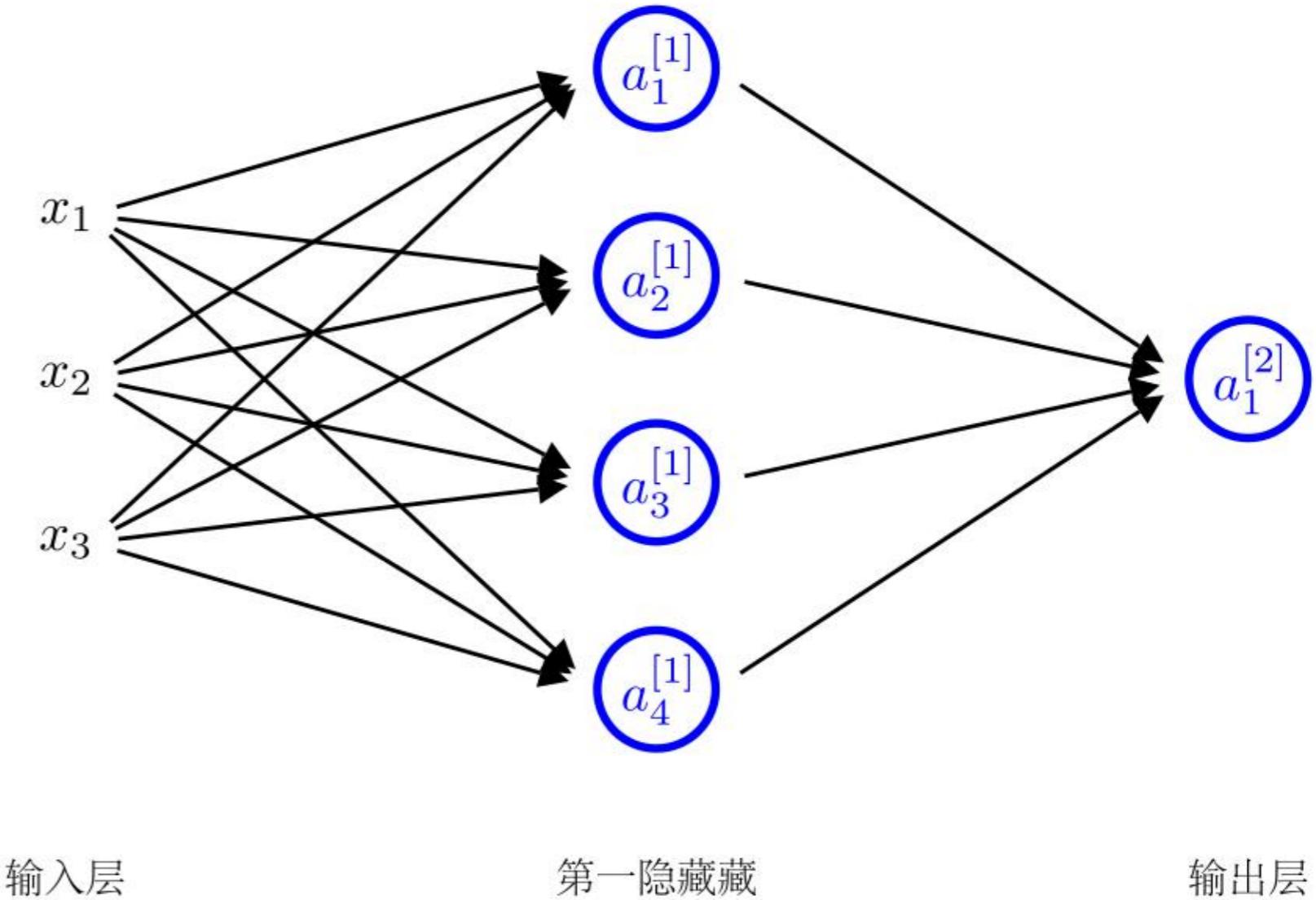
1. 为了得到偏导数 $\partial \mathcal{J}(\theta^{(t)}) / \partial \theta$, 我们需要引入
 - 前向传播: 基于当前模型参数, 计算“激活值”以及代价函数值
 - 后向传播: 基于前向传播的计算结果, 得到偏导数
2. 我们用一个简单的例子介绍前向和后向传播的计算过程

回顾逻辑回归模型

1. 在处理实际问题时，这个模型过于简单
2. 为什么不加入更多的“圆圈”呢？

神经网络的例子

1. 假设 $x = (x_1, x_2, x_3)^T$



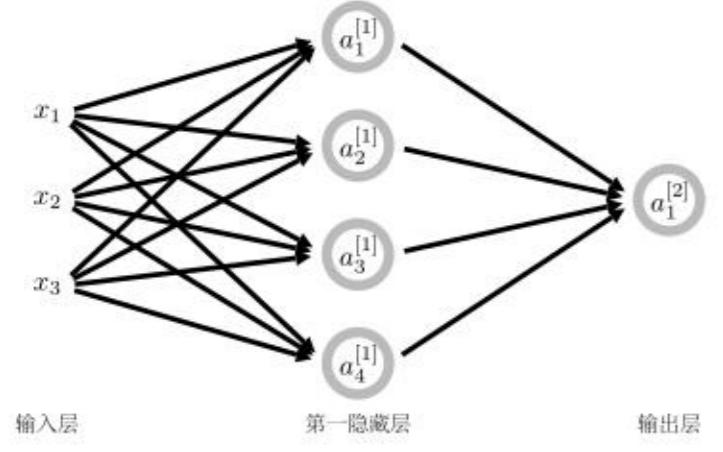
2. 每个圆圈包含如下两个运算

- 线性变换 并对应两个模型参数（偏置以及权重）
- 激活变换（非线性变换），其不对应模型参数

神经网络的例子

1. 注意：

- 不同的“圆圈”对应不同模型参数，用来提取数据中的不同信息
- 我们在隐藏层中“构建”了一组“新的”特征向量
- 相比于逻辑回归模型，神经网络多了一个隐藏层用来提取更多信息
- 对于输出层而言，隐藏层可被认为是它的特征向量



计算过程

模型参数

$$z_1^{[1]} = b_1^{[1]} + \mathbf{x}^T \mathbf{w}_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$b_1^{[1]} \quad \mathbf{w}_1^{[1]}$$

$$z_2^{[1]} = b_2^{[1]} + \mathbf{x}^T \mathbf{w}_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$b_2^{[1]} \quad \mathbf{w}_2^{[1]}$$

$$z_3^{[1]} = b_3^{[1]} + \mathbf{x}^T \mathbf{w}_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$b_3^{[1]} \quad \mathbf{w}_3^{[1]}$$

$$z_4^{[1]} = b_4^{[1]} + \mathbf{x}^T \mathbf{w}_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

$$b_4^{[1]} \quad \mathbf{w}_4^{[1]}$$

$$z_1^{[2]} = b_1^{[2]} + (\mathbf{a}^{[1]})^T \mathbf{w}_1^{[2]}, \quad a_1^{[2]} = \sigma(z_1^{[2]})$$

$$b_1^{[2]} \quad \mathbf{w}_1^{[2]}$$

向量化

1. 记

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

2. 模型参数

- $\{(\mathbf{b}^{[l]}, \mathbf{W}^{[l]}): l = 1, \dots, L\}$

向量化

1. 对于之前考虑的例子，我们有

- $d^{[0]} = 3, d^{[1]} = 4, d^{[2]} = 1$
- $\mathbf{W}^{[1]} = (\mathbf{w}_1^{[1]}, \mathbf{w}_2^{[1]}, \mathbf{w}_3^{[1]}, \mathbf{w}_4^{[1]})^T \in \mathbb{R}^{4 \times 3}, \mathbf{W}^{[2]} = (\mathbf{w}_1^{[2]})^T \in \mathbb{R}^{1 \times 4}$
- $\mathbf{b}^{[1]} = (b_1^{[1]}, b_2^{[1]}, b_3^{[1]}, b_4^{[1]})^T \in \mathbb{R}^{4 \times 1}, b^{[2]} \in \mathbb{R}^{1 \times 1}$

前向传播

1. 前向传播利用当前参数模型计算“激活值”以及代价函数值
2. 假设当前模型参数为 $\mathbf{b}^{[1]}, \mathbf{W}^{[1]}, b^{[2]}, \mathbf{W}^{[2]}$
3. (简化的) 计算过程为

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

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

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

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

$$\mathbf{a}^{[1]} = \begin{pmatrix} \sigma(z_1^{[1]}) \\ \sigma(z_2^{[1]}) \\ \sigma(z_3^{[1]}) \\ \sigma(z_4^{[1]}) \end{pmatrix}$$

前向传播

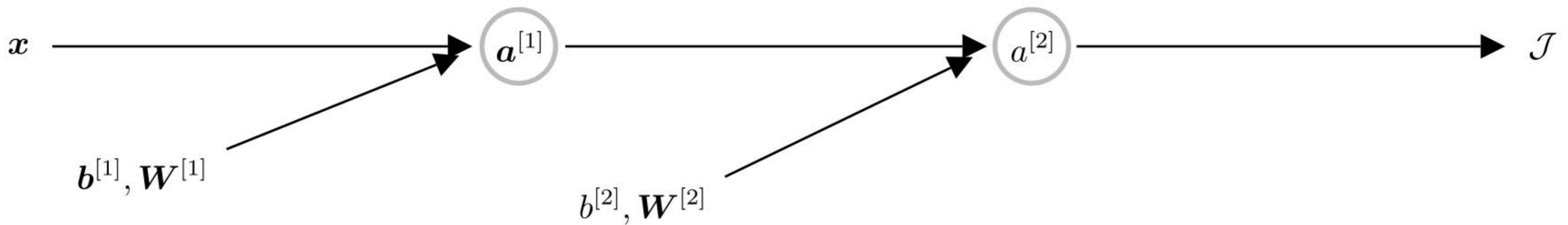
1. 回顾：我们考虑二元回归问题： $y \in \{0, 1\}$
2. $a^{[2]}$ 是基于当前模型参数，对条件概率 $P(y = 1 | \mathbf{x})$ 的估计
3. 考虑损失函数

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

4. 如果 $y \in \mathbb{R}$ ，那么 $a^{[2]}$ 和损失函数 \mathcal{L} 的形式是什么？
5. 如果标签 y 的取值只能非负呢？

前向传播

1. 将计算过程可视化



前向传播

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

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

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

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

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

后向传播

1. 基于当前模型参数及前向传播计算结果，后向传播用于计算

$$\frac{\partial \mathcal{J}}{\partial \mathbf{b}^{[l]}}, \quad \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{[l]}} \quad (l = 1, \dots, L)$$

2. 核心技术: **链式法则**

- 记 $f(\mathbf{x}) = g \circ h(\mathbf{x})$
- 我们有

$$\frac{\partial f}{\partial \mathbf{x}} = \frac{\partial g}{\partial h} \cdot \frac{\partial h}{\partial \mathbf{x}}$$

例子

1. 例 1: 考虑 $f(\mathbf{x})$

- $\mathbf{x} = (x_1, \dots, x_m)^T$: 列向量
- $f(\mathbf{x})$: 一维

2. 那么, 我们有

$$\frac{\partial f}{\partial \mathbf{x}} = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_m} \right)^T$$

3. $\partial f / \partial \mathbf{x}$ 与 \mathbf{x} 具有相同维度

例子

1. 例 1: 考虑 $f(\mathbf{x})$

- $\mathbf{x} = (x_1, \dots, x_m)^T$: 行向量
- $f(\mathbf{x})$: 一维

2. 那么, 我们有

$$\frac{\partial f}{\partial \mathbf{x}} = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_m} \right)$$

3. $\partial f / \partial \mathbf{x}$ 与 \mathbf{x} 具有相同维度

例子

1. 考虑 $f(\mathbf{x}, b\mathbf{y}) = \mathbf{x}^T \mathbf{y}$

- $\mathbf{x} = (x_1, \dots, x_m)^T$
- $\mathbf{y} = (y_1, \dots, y_m)^T$

2. 那么，我们有

$$\begin{array}{ll} \frac{\partial f}{\partial \mathbf{x}} = \mathbf{y} & \frac{\partial f}{\partial \mathbf{y}} = \mathbf{x} \\ \frac{\partial f}{\partial \mathbf{x}^T} = \mathbf{y}^T & \frac{\partial f}{\partial \mathbf{y}^T} = \mathbf{x}^T \end{array}$$

例子

1. 例 3: 考虑 $f(\mathbf{X})$

- $\mathbf{X} = (x_{ij}) \in \mathbb{R}^{m \times k}$: 一个 $m \times k$ 维矩阵
- $f(\mathbf{x})$: 一维

2. 那么, 我们有

$$\frac{\partial f}{\partial \mathbf{X}} = \left(\frac{\partial f}{\partial x_{ij}} \right)$$

3. $\partial f / \partial \mathbf{X}$ 与 \mathbf{X} 具有相同维度

例子

1. 例 3: 考虑 $f(\mathbf{x})$

- $\mathbf{x} = (x_1, \dots, x_{\textcolor{blue}{m}})^T$: 一个 $\textcolor{blue}{m}$ 维列向量
- $f(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_{\textcolor{red}{k}}(\mathbf{x}))^T$: 一个 $\textcolor{red}{k}$ 维列向量

2. 那么, 我们有

$$\frac{\partial \mathbf{f}^T}{\partial \mathbf{x}} = \left(\frac{\partial f_1}{\partial \mathbf{x}}, \dots, \frac{\partial f_{\textcolor{red}{k}}}{\partial \mathbf{x}} \right) \in \mathbb{R}^{\textcolor{blue}{m} \times \textcolor{red}{k}}$$

3. $\partial \mathbf{f}^T / \partial \mathbf{x}$ 为一个 $\textcolor{blue}{m} \times \textcolor{red}{k}$ 维矩阵

4. 我们将在本课程中使用以上定义, 但不要和高等数学课中的定义相混淆

例子

1. 例 1 (链式法则): 考虑

$$f(\boldsymbol{x}) = g \circ \boldsymbol{h}(\boldsymbol{x})$$

- $\boldsymbol{h}(\boldsymbol{x}) = (h_1(\boldsymbol{x}), \dots, h_m(\boldsymbol{x}))$
- $m \geq 2$

2. 那么, 我们有

$$\frac{\partial f}{\partial \boldsymbol{x}} = \sum_{j=1}^m \frac{\partial g}{\partial h_j} \frac{\partial h_j}{\partial \boldsymbol{x}}$$

例子

1. 如果 \mathbf{x} 是列向量，我们还有

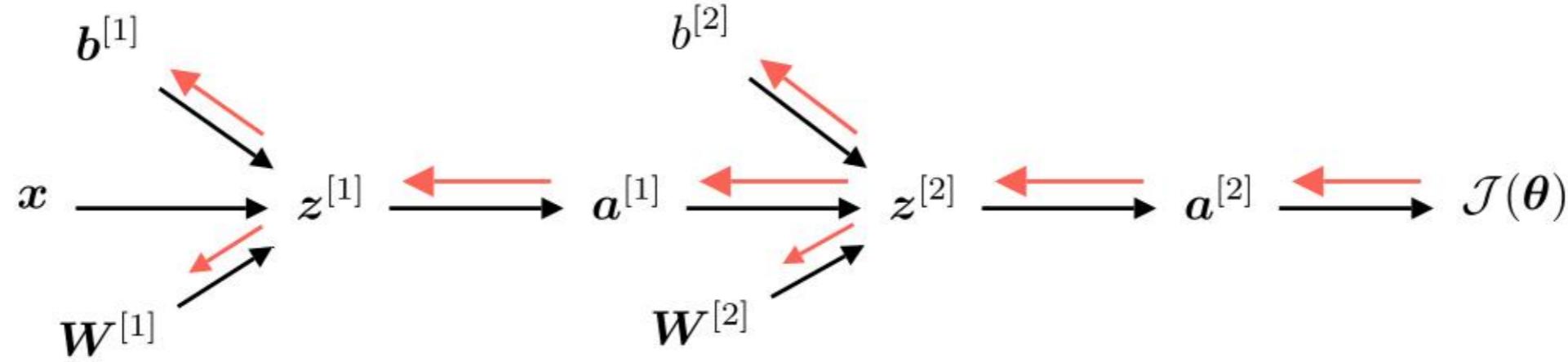
$$\frac{\partial f}{\partial \mathbf{x}} = \sum_{j=1}^m \frac{\partial g}{\partial h_j} \frac{\partial h_j}{\partial \mathbf{x}} = \left(\frac{\partial h_1}{\partial \mathbf{x}}, \dots, \frac{\partial h_m}{\partial \mathbf{x}} \right) \begin{pmatrix} \frac{\partial g}{\partial h_1} \\ \vdots \\ \frac{\partial g}{\partial h_m} \end{pmatrix} = \frac{\partial \mathbf{h}^T}{\partial \mathbf{x}} \frac{\partial g}{\partial \mathbf{h}}$$

- $\partial g / \partial \mathbf{h} = (\partial g / \partial h_1, \dots, \partial g / \partial h_m)^T$
- $\partial \mathbf{h}^T / \partial \mathbf{x} = (\partial h_1 / \partial \mathbf{x}, \dots, \partial h_m / \partial \mathbf{x})$

2. 即链式法则求导可被向量化

3. 然而，核心步骤是得到 $\partial h_j / \partial \mathbf{x}$

4. 当 \mathbf{x} 为矩阵时，向量化失效



$$\frac{\partial \mathcal{J}(\theta)}{\partial b^{[2]}} = a^{[2]} - y$$

$$\frac{\partial \mathcal{J}(\theta)}{\partial \mathbf{W}^{[2]}} = (a^{[2]} - y)(\mathbf{a}^{[1]})^T$$

$$\frac{\partial \mathcal{J}(\theta)}{\partial \mathbf{b}^{[1]}} = (a^{[2]} - y)\mathbf{D}(\mathbf{W}^{[2]})^T$$

$$\frac{\partial \mathcal{J}(\theta)}{\partial \mathbf{W}^{[1]}} = (a^{[2]} - y)\mathbf{D}(\mathbf{W}^{[2]})^T \mathbf{x}^T$$

$$\mathbf{D} = \text{diag}(\{a_j^{[1]}(1 - a_j^{[1]}): j = 1, \dots, d^{[1]}\})$$

我们只需缓存 $\mathbf{a}^{[1]}$ and $a^{[2]}$

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. 更新模型参数

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

步骤4. 回到步骤 2 直至收敛