

1 课程简介

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. 回顾已学内容

2. 问题

3. 学习目标

4. 交叉熵

回顾已学内容

1. 数据集类型：

- 训练集 (Training data) 样本量 n : $\{(\mathbf{x}_i, y_i) : i = 1, \dots, n\}$
- 测试集 (Test data)

2. 回归类型：

- 线性回归: $y_i = b_0 + \mathbf{x}_i^T \mathbf{w}_0 + \epsilon_i \quad (i = 1, \dots, n)$
- 逻辑回归: $\text{logit}P(y_i = 1 \mid \mathbf{x}_i) = b_0 + \mathbf{x}_i^T \mathbf{w}_0 \quad (i = 1, \dots, n)$
 - ▷ $\text{logit}p = \log(p) - \log(1 - p), \quad \forall p \in (0, 1)$
- 多元逻辑回归 (Softmax regression)

符号

1. n : 训练集样本量

2. d : 特征向量维度

3. $\mathbf{X} \in \mathbb{R}^{n \times d}$: 设计矩阵

- $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$

- $\mathbf{x}_i \in \mathbb{R}^{d \times 1}$: 第*i*个训练样本的特征向量

4. $\mathbf{Y} \in \mathbb{R}^{n \times 1}$: 标签向量

- $\mathbf{Y} = (y_1, \dots, y_n)^T$

- $y_i \in \mathbb{R}$: 第*i*个样本的标签

线性回归--模型设定

1. 训练集

- 特征向量, 标签: $\mathbf{x}_i \in \mathbb{R}^{d \times 1}$, $y_i \in \mathbb{R}$ ($i = 1, \dots, n$)

2. 目标

- 学习 \mathbf{x} 和 y 之间的关系

3. 真实模型 (用于产生数据)

- $y_i = \mathbf{b}_0 + \mathbf{x}_i^T \mathbf{w}_0 + \epsilon_i$ ($i = 1, \dots, n$)

▷ $\boldsymbol{\theta}_0 = (\mathbf{b}_0, \mathbf{w}_0^T)^T$: 真实模型参数

▷ ϵ_i : 满足 $E(\epsilon_i | \mathbf{x}_i) = 0$, 以及 $\text{var}(\epsilon_i | \mathbf{x}_i) = \sigma^2 > 0$ 的残差项

线性回归--模型设定

1. 所提出模型（用于拟合数据）

- $E(y_i \mid \mathbf{x}_i) = \mathbf{b} + \mathbf{x}_i^T \mathbf{w} \quad (i = 1, \dots, n)$
 - ▷ $\boldsymbol{\theta} = (\mathbf{b}, \mathbf{w}^T)^T$: (提出) 模型的参数
 - ▷ 与真实模型形式相同，但我们需要确定模型参数
 - ▷ 训练模型是指，估计所提出模型的参数

线性回归--参数估计

1. 损失函数: 对于第 i 个训练样本,

- l_2 -损失: $\mathcal{L}(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2$
- $\hat{y}_i = b + \mathbf{x}_i^T \mathbf{w}$
- \hat{y}_i 是模型参数 $\boldsymbol{\theta} = (b, \mathbf{w}^T)^T$ 的函数, 但为了简单起见, 我们忽略模型参数
- $\mathcal{L}(y_i, \hat{y}_i)$ 实际上是模型参数 $\boldsymbol{\theta}$ 的函数

2. 代价函数

$$\mathcal{J}(\boldsymbol{\theta}) = n^{-1} \sum_{i=1}^n \mathcal{L}(y_i, \hat{y}_i) = n^{-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = n^{-1} \sum_{i=1}^n (y_i - b - \mathbf{x}_i^T \mathbf{w})^2$$

3. 我们通常不区分损失函数和代价函数

线性回归--参数估计

1. 求解

$$\frac{\mathcal{J}}{\partial b} = 0, \quad \frac{\mathcal{J}}{\partial \mathbf{w}} = 0$$

2. 问题：求和运算导致算法计算效率低下

向量化

1. 考虑

$$\sum_{i=1}^m a_i b_i$$

- $a_i \in \mathbb{R}, b_i \in \mathbb{R}$
2. 对于 Python 或者其他软件，用 for 循环计算加和是非常低效的
 3. 然而，矩阵运算是高效的
 4. 因此，我们可对加进行向量化处理

向量化

1. 向量化结果

$$\sum_{i=1}^m \textcolor{red}{a}_i b_i = \textcolor{red}{a}^T \cdot \textcolor{blue}{b}$$

- $\textcolor{red}{a} = (\textcolor{red}{a}_1, \dots, \textcolor{red}{a}_m)^T \in \mathbb{R}^{m \times 1}$
- $\textcolor{blue}{b} = (\textcolor{blue}{b}_1, \dots, \textcolor{blue}{b}_m)^T \in \mathbb{R}^{m \times 1}$

向量化

1. 例 1: 考虑

$$\sum_{i=1}^m \textcolor{red}{a}_i$$

- $\textcolor{red}{a}_i \in \mathbb{R}$

2. 向量化结果

$$\sum_{i=1}^n \textcolor{red}{a}_i = \sum_{i=1}^n \textcolor{red}{a}_i \times \mathbf{1} = \mathbf{a}^T \cdot \mathbf{1} (= \mathbf{1}^T \cdot \mathbf{a})$$

- $\mathbf{a} = (\textcolor{red}{a}_1, \dots, \textcolor{red}{a}_m)^T \in \mathbb{R}^{m \times 1}$
- $\mathbf{1} = (\mathbf{1}, \dots, \mathbf{1})^T \in \mathbb{R}^{m \times 1}$: 由 1 组成的长度为 m 的列向量

向量化

1. 例 2: 考虑

$$\sum_{i=1}^m \textcolor{red}{a}_i \textcolor{blue}{b}_i$$

- $\textcolor{red}{a}_i \in \mathbb{R}$: 一个可能取值为 1 的标量
- $\textcolor{blue}{b}_i \in \mathbb{R}^{k \times 1}$

2. 向量化结果

$$\sum_{i=1}^n \textcolor{red}{a}_i \textcolor{blue}{b}_i = \textcolor{blue}{B} \cdot \textcolor{red}{a}$$

- $\textcolor{red}{a} = (\textcolor{red}{a}_1, \dots, \textcolor{red}{a}_m)^T \in \mathbb{R}^{m \times 1}$
- $\textcolor{blue}{B} = (\textcolor{blue}{b}_1, \dots, \textcolor{blue}{b}_m) \in \mathbb{R}^{k \times m}$

向量化

1. 例 3: 考虑

$$\sum_{i=1}^m \color{red}{a_i} \color{blue}{b_i}^\top$$

- $\color{red}{a_i} \in \mathbb{R}$: 一个可能取值为 1 的标量
- $\color{blue}{b_i} \in \mathbb{R}^{k \times 1}$

2. 向量化结果

$$\sum_{i=1}^n \color{red}{a_i} \color{blue}{b_i} = \color{red}{a}^\top \cdot \color{blue}{B}$$

- $\color{red}{a} = (\color{red}{a_1}, \dots, \color{red}{a_m})^\top \in \mathbb{R}^{m \times 1}$
- $\color{blue}{B} = (\color{blue}{b_1}, \dots, \color{blue}{b_m})^\top \in \mathbb{R}^{m \times k}$

向量化

1. 例 4: 考虑

$$\sum_{i=1}^m \color{red}{\mathbf{a}_i} \color{blue}{\mathbf{b}_i}^T$$

- $\color{red}{\mathbf{a}_i} \in \mathbb{R}^{k \times 1}$
- $\color{blue}{\mathbf{b}_i} \in \mathbb{R}^{l \times 1}$

2. 向量化结果

$$\sum_{i=1}^n \color{red}{\mathbf{a}_i} \color{blue}{\mathbf{b}_i}^T = \color{red}{\mathbf{A}} \cdot \color{blue}{\mathbf{B}}$$

- $\color{red}{\mathbf{A}} = (\color{red}{\mathbf{a}_1}, \dots, \color{red}{\mathbf{a}_m}) \in \mathbb{R}^{k \times m}$

- $\color{blue}{\mathbf{B}} = (\color{blue}{\mathbf{b}_1}, \dots, \color{blue}{\mathbf{b}_m})^T \in \mathbb{R}^{m \times k}$

向量化

1. 记 $\hat{\mathbf{Y}} = b\mathbf{1}_n + \mathbf{X}\mathbf{w}$

- $\mathbf{1}_n = (1, \dots, 1)^T \in \mathbb{R}^{n \times 1}$: 由 1 组成的长度为 n 的向量

2. 代价函数

$$\begin{aligned}\mathcal{J}(\boldsymbol{\theta}) &= n^{-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = n^{-1} (\mathbf{Y} - \hat{\mathbf{Y}})^T (\mathbf{Y} - \hat{\mathbf{Y}}) \\ &= n^{-1} (\mathbf{Y} - b\mathbf{1}_n - \mathbf{X}\mathbf{w})^T (\mathbf{Y} - b\mathbf{1}_n - \mathbf{X}\mathbf{w}) \\ &= n^{-1} (\mathbf{Y} - \tilde{\mathbf{X}}\boldsymbol{\theta})^T (\mathbf{Y} - \tilde{\mathbf{X}}\boldsymbol{\theta})\end{aligned}$$

- $\tilde{\mathbf{X}} = (\tilde{x}_1, \dots, \tilde{x}_n)^T$ (将 b 和 \mathbf{w} 放在一起考虑)
- $\tilde{\mathbf{x}}_i = (1, \mathbf{x}_i^T)^T$

向量化

1. 考虑

$$\frac{\partial \mathcal{J}}{\partial \boldsymbol{\theta}} = 0$$

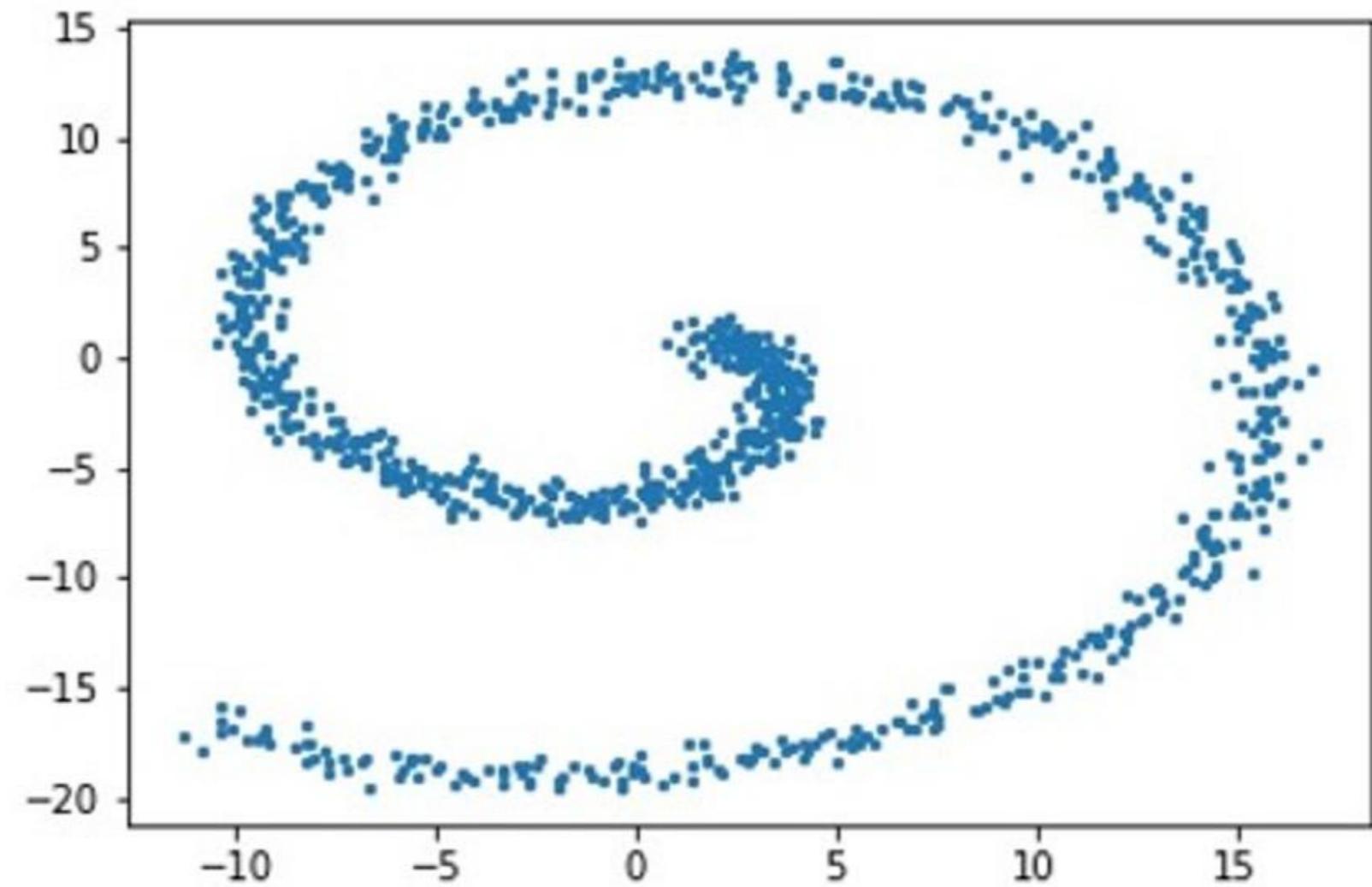
2. 得到正则方程 (normal equation)

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \boldsymbol{\theta} = \tilde{\mathbf{X}}^T \mathbf{Y}$$

3. 解 (显式解)

$$\hat{\boldsymbol{\theta}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

例子



逻辑回归--模型设定

1. 训练集

- 特征向量, 标签: $\mathbf{x}_i \in \mathbb{R}^{d \times 1}$, $y_i \in \{0, 1\}$ ($i = 1, \dots, n$)

2. 目标

- 学习 \mathbf{x} 和 y 之间的关系

3. 真实模型 (用于产生数据)

$$E(y_i | \mathbf{x}_i) = P(y_i = 1 | \mathbf{x}_i) = \sigma(\mathbf{b}_0 + \mathbf{x}_i^T \mathbf{w}_0) \quad (i = 1, \dots, n)$$

- $\sigma(z) = \{1 + \exp(-z)\}^{-1}$ (我们为什么需要这步变换?)

逻辑回归--模型设定

1. 所提出模型（用于拟合数据）

- $P(y_i = 1 \mid \mathbf{x}_i) = \sigma(\mathbf{b} + \mathbf{x}_i^T \mathbf{w}) \quad (i = 1, \dots, n)$
 - ▷ $\boldsymbol{\theta} = (\mathbf{b}, \mathbf{w}^T)^T$: (提出) 模型的参数
 - ▷ 与真实模型形式相同，但我们需要确定模型参数
 - ▷ 训练模型是指，估计所提出模型的参数

逻辑回归--参数估计

1. 损失函数：对于第 i 个训练样本，

- 交叉熵（负对数似然函数）

$$\mathcal{L}(y_i, a_i) = -\{y_i \log a_i + (1 - y_i) \log(1 - a_i)\}$$

- $a_i = \sigma(z_i)$ (估计概率)
 - ▷ $a_i = \sigma(z_i)$ (估计概率)
 - ▷ $z_i = b + \mathbf{x}_i^T \mathbf{w}$
- a_i 是模型参数 $\boldsymbol{\theta} = (b, \mathbf{w}^T)^T$ 的函数，但简单起见，我们忽略该模型参数
- 交叉熵实际上是模型参数 $\boldsymbol{\theta}$ 的函数

逻辑回归--参数估计

1. 代价函数

$$\begin{aligned}\mathcal{J}(\boldsymbol{\theta}) &= n^{-1} \sum_{i=1}^n \mathcal{L}\{y_i, a_i\} \\ &= -n^{-1} \sum_{i=1}^n \{y_i \log a_i + (1 - y_i) \log\{1 - a_i\}\}\end{aligned}$$

- $\boldsymbol{\theta} = (b, \mathbf{w}^T)^T$
- $a_i = \sigma(b + \mathbf{x}_i^T \mathbf{w})$

逻辑回归--参数估计

1. 求解

$$\frac{\mathcal{J}}{\partial \boldsymbol{\theta}} = 0$$

2. 没有显式解

牛顿-拉裴森 (Newton-Raphson) 算法

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

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

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

$$H(\mathcal{J})(\boldsymbol{\theta}^{(t)}) = \frac{\partial^2 \mathcal{J}}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T}(\boldsymbol{\theta}^{(t)})$$

步骤3. 更新模型参数

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

步骤4. 回到步骤 2 直至收敛

讨论

1. 收敛速度快（得益于 Hessian 矩阵的使用）
2. 计算效率低下（涉及到对 Hessian 矩阵求逆）
3. 适用于当模型参数规模较小的情况
4. 不适用于深度学习模型参数的训练
5. 什么算法适用于深度学习模型的训练呢？

批量梯度下降 (Batch gradient descent) 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

批量梯度下降 (Batch gradient descent) 法

1. α : 学习率 (learning rate)

- 控制算法收敛速度
- 影响训练模型的好坏
- 将在后续学习中展开讨论

算法比较

1. 牛顿-拉裴森算法

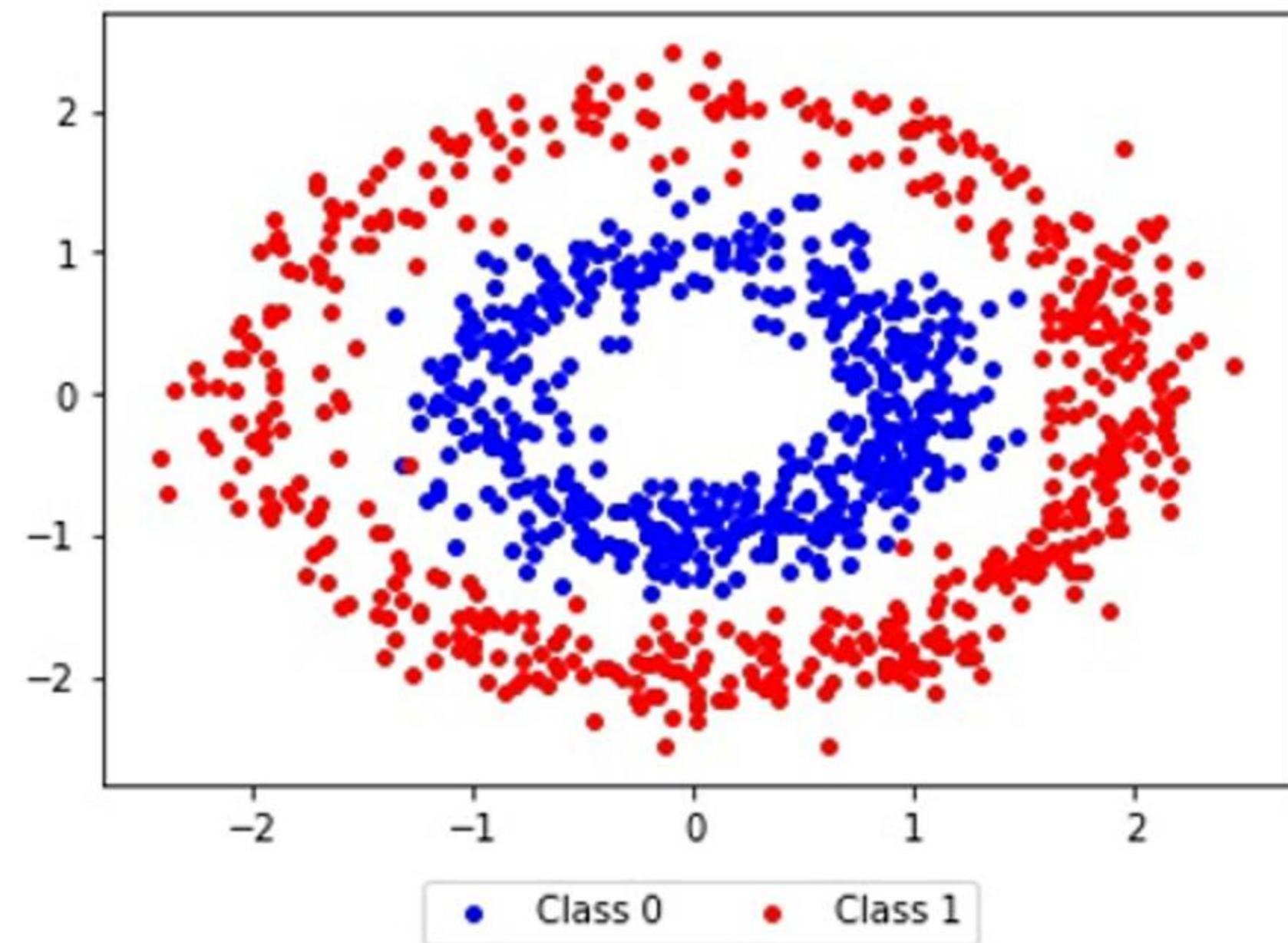
- 优势：
 - ▷ 收敛速度快
 - ▷ 算法精度高
- 缺点：
 - ▷ 计算效率低 (Hessian 矩阵求逆)
 - ▷ 稳定性差 (当 Hessian 矩阵出现病态时)

算法比较

1. 批量梯度下降法

- 优势：
 - ▷ 计算效率高
 - ▷ 便于操作
- 缺点：
 - ▷ 收敛速度较慢
 - ▷ 对学习率较为敏感

例子



总结

1. 我们通常利用如下步骤进行数据建模

问题	模型	损失	算法
$y \in \mathbb{R}$	$E(y \mathbf{x}) = b + \mathbf{x}^T \mathbf{w}$	$n^{-1} \sum_{i=1}^n (y_i - b - \mathbf{x}_i^T \mathbf{w})^2$	正则方程
$y \in \{0, 1\}$	$E(y \mathbf{x}) = \sigma(b + \mathbf{x}^T \mathbf{w})$	$-n^{-1} \sum_{i=1}^n \{y_i \log a_i + (1 - y_i) \log(1 - a_i)\}$	梯度下降

2. 梯度下降法广泛应用于人工智能模型中

3. 梯度下降法的关键步骤是计算

$$\frac{\partial \mathcal{J}}{\partial \boldsymbol{\theta}}$$

问题

1. 目前为止，我们假设真实模型和所提出模型具有相同的形式
2. 然而，这种假设在现实生活中通常不成立
3. 模型错误设定
 - 被提出的（统计）模型过于简单
 - 然而，真实模型往往很复杂

学习目标

1. 全连接神经网络（**多层感知机**）
2. 卷积神经网络（**CNN**, ...）
3. 序列模型（**RNN**, **LSTM**, **transformers**, ...）
4. （时间允许情况下的）其他话题（**GNN**, **GCN**, ...）

学习目标

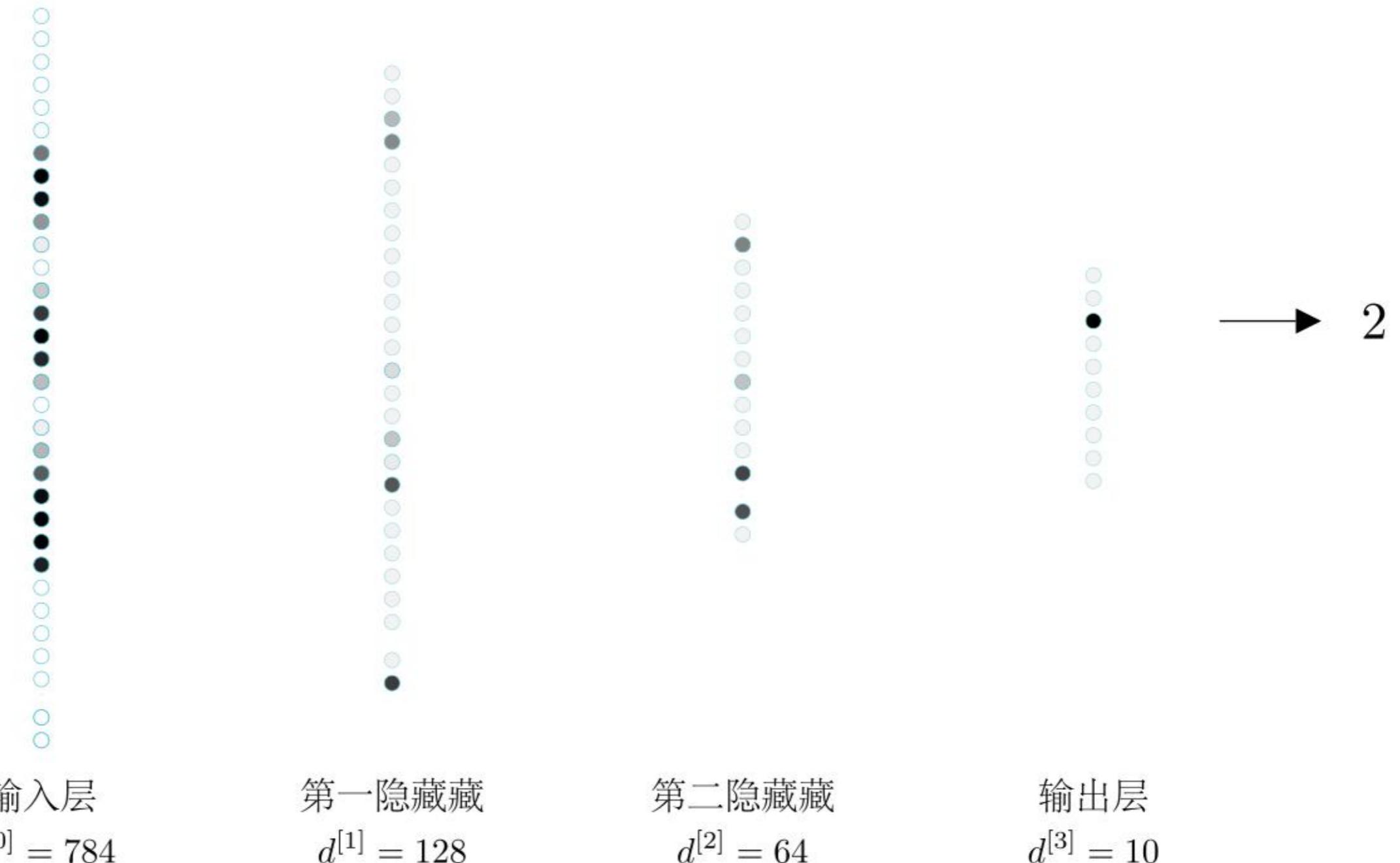
1. 全连接神经网络（**多层感知机**）
2. 卷积神经网络（**CNN**, ...）
3. 序列模型（**RNN**, **LSTM**, **transformers**, ...）
4. （时间允许情况下的）其他话题（**GNN**, **GCN**, ...）

全连接神经网络例子

测试图片



模型 (具有两个隐藏层的全连接神经网络)

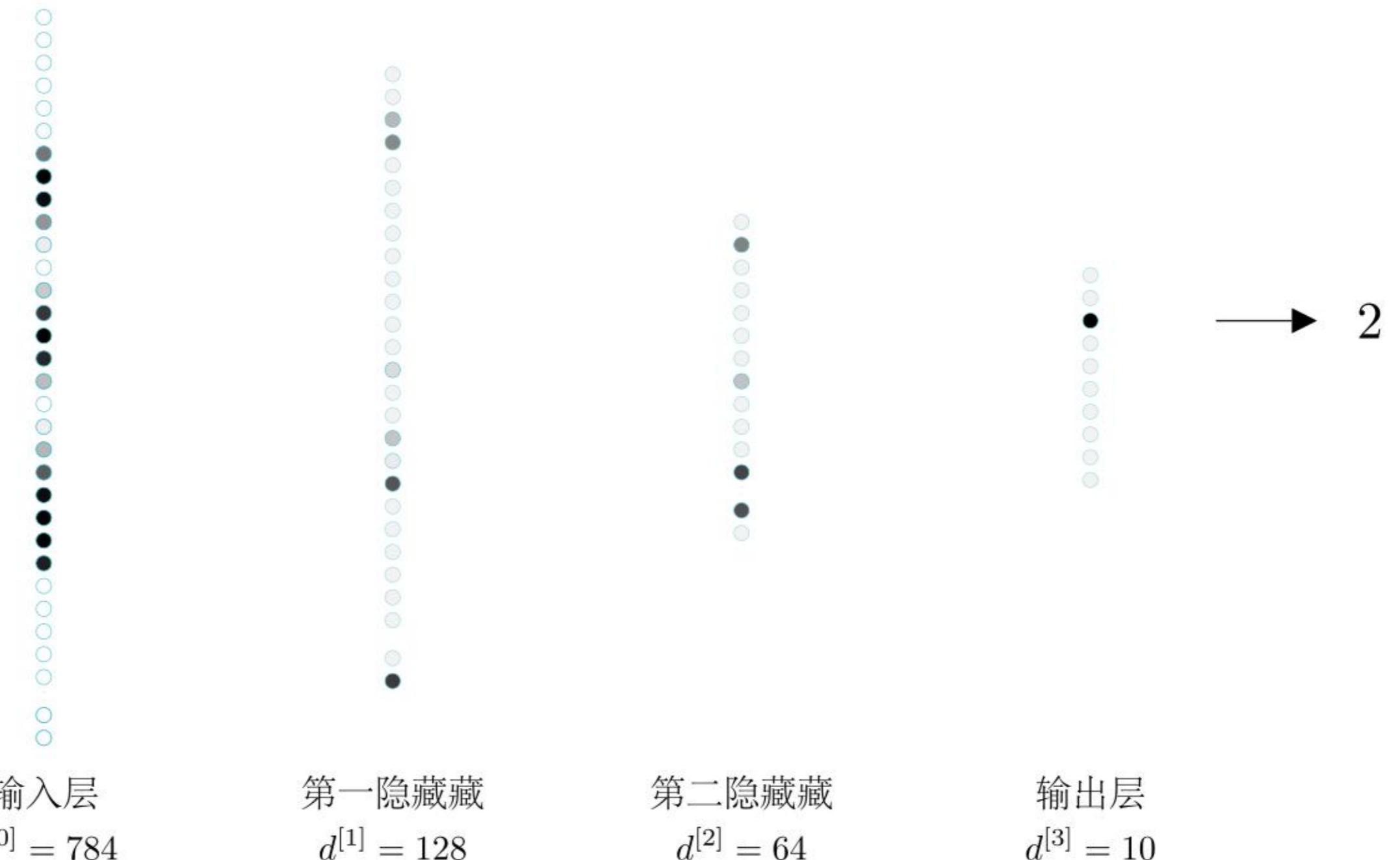


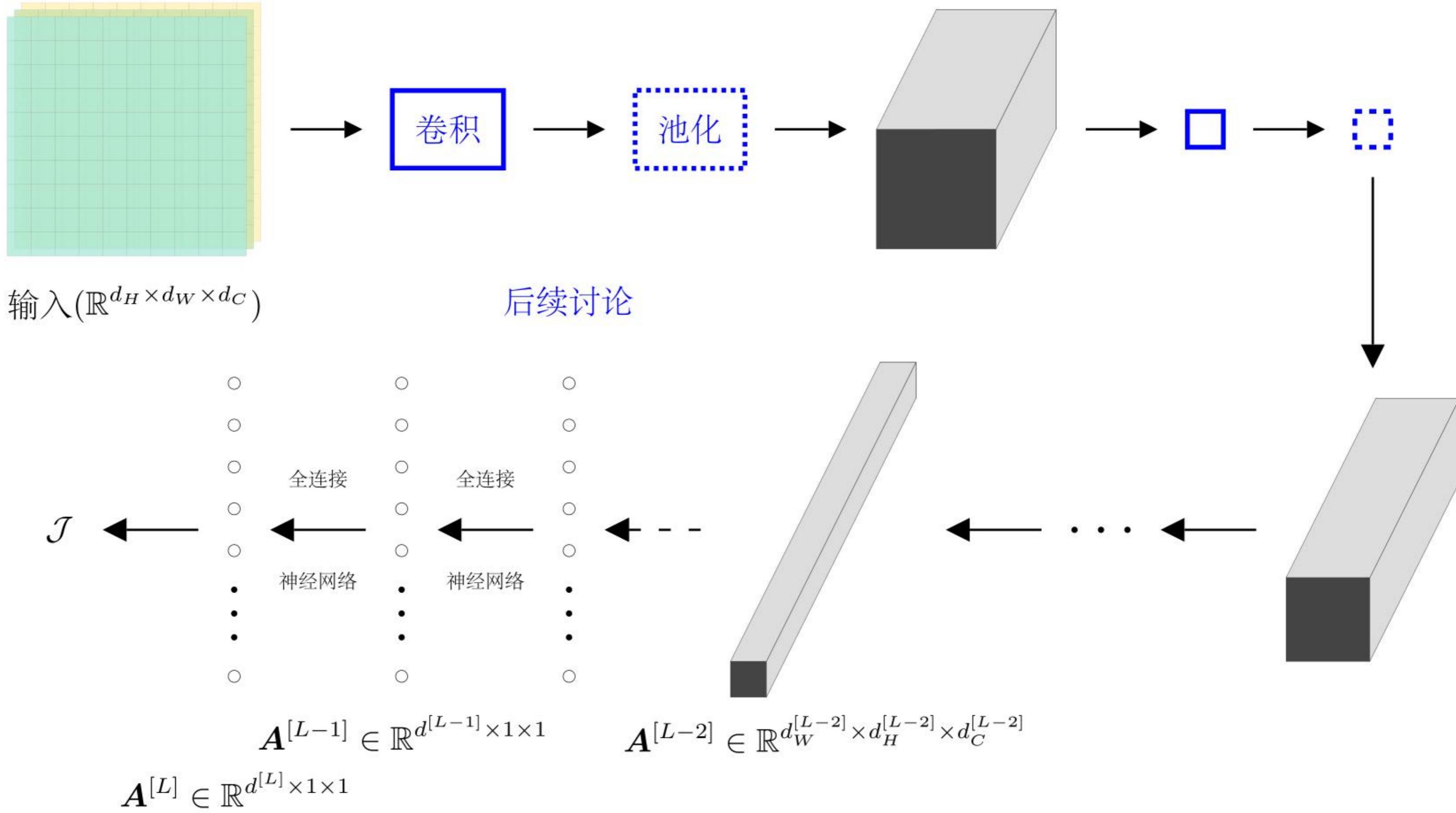
全连接神经网络例子

测试图片



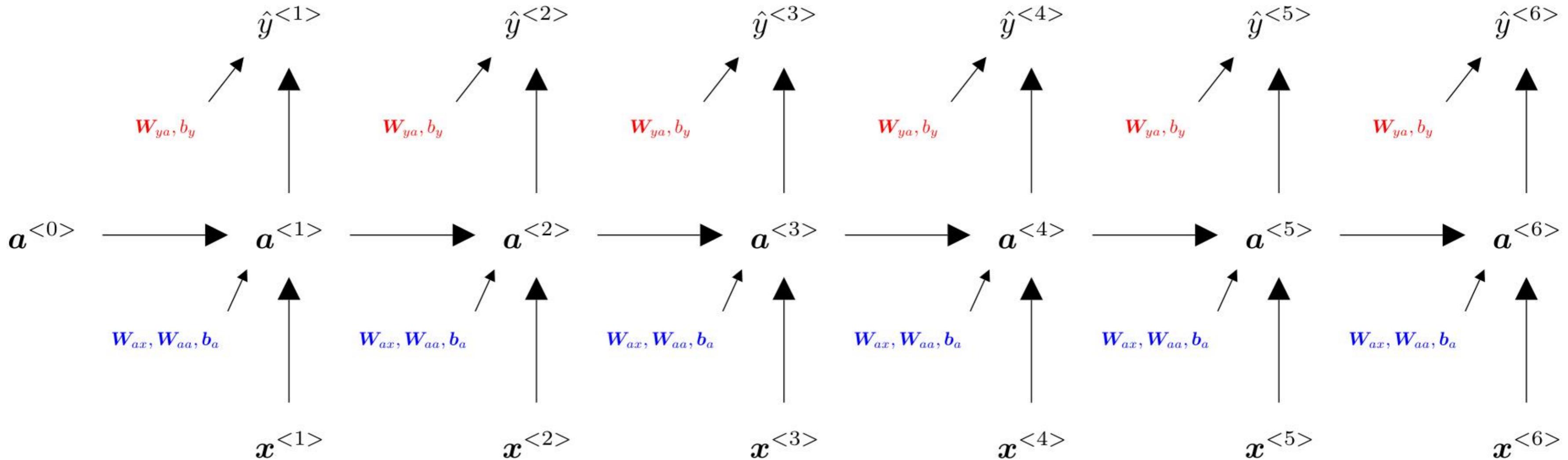
模型 (具有两个隐藏层的全连接神经网络)





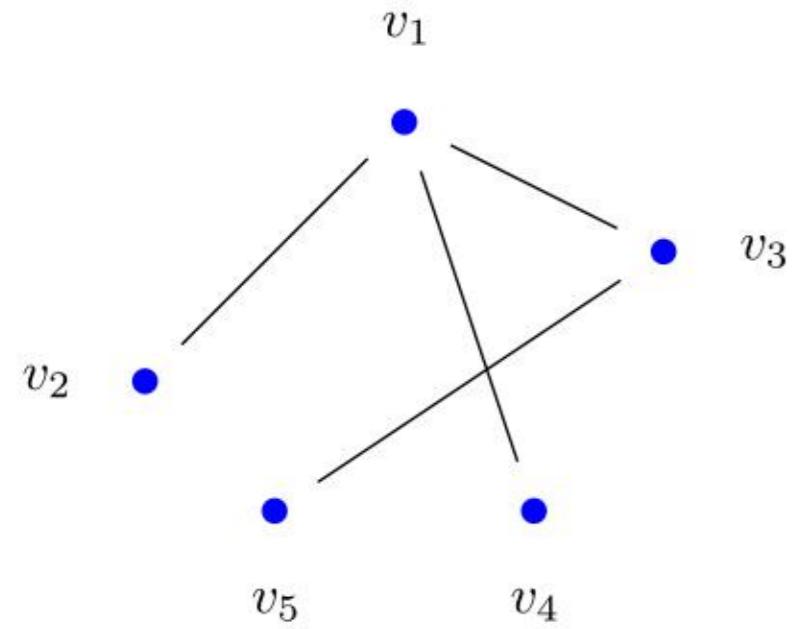
循环神经网络模块

$$\hat{y}^{*} = \sigma_{ya}(\mathbf{W}_{ya}\mathbf{a}^{*} + b_y) \quad (i = 1, \dots, 6)**$$



$$\mathbf{a}^{*} = \sigma_{ax}(\mathbf{W}_{ax}\mathbf{x}^{*} + \mathbf{W}_{aa}\mathbf{a}^{} + \mathbf{b}_a) \quad (i = 1, \dots, 6)**$$

无向图



邻接矩阵 ($\{0,1\}$; $[0,1]$)

$$\begin{pmatrix} 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

顶点集 $\mathcal{V} = \{v_i : i = 1, \dots, n\}$

边集: \mathcal{E}

无向图: $\mathcal{G} = \mathcal{V} \cup \mathcal{E}$

度矩阵 (对角阵, 邻接矩阵 A 的每行求和)

$$D = \begin{pmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

交叉熵

1. 一个事实：在一般条件下，

$$E_{X \sim P}\{\log p(X)\} \geq E_{X \sim P}\{\log q(X)\}$$

- $X \sim P$: 随机变量 X 产生于分布 P
- $p(x)$: 分布 P 的概率密度函数
- $q(x)$: 其他任意概率密度函数

2. 对于两个密度函数的 p 和 q 的交叉熵

$$H(p, q) = -E_{X \sim P}\{\log q(X)\} = -E\{\log q(X)\}$$

- 衡量概率密度函数 q 对随机变量 X 分布的近似程度

交叉熵

1. 只观测到随机样本 $\{x_1, \dots, x_n\}$
2. 交叉熵 $H(p, q)$ 可被如下公式近似

$$\hat{H}(p, q) = -\frac{1}{n} \sum_{i=1}^n \log q(x_i)$$

- 以上是经验密度函数 $P_n(x) = n^{-1} \sum_{i=1}^n \delta(x = x_i)$ 和 q 的交叉熵
- 3. 目标: 寻找密度函数 q 用来最小化 $\hat{H}(p, q)$
 - 即为了近似随机变量的密度函数

回顾逻辑回归

1. 假设随机变量 $\{(\mathbf{x}_i, y_i) : i = 1, \dots, n\}$ 是独立同分布的，对应的概率密度函数为

$$p(\mathbf{x}, y) = p(\mathbf{x})p(y | \mathbf{x})$$

- $p(\mathbf{x})$: (不是很重要的) 边际概率密度函数 \mathbf{x}
- $p(y | \mathbf{x}) = \sigma(\mathbf{b}_0 + \mathbf{x}_i^T \mathbf{w}_0)$: 关于 y 的 (我们感兴趣的) 参数化的条件密度函数

2. 目标: 寻找概率密度函数 $q(\mathbf{x}, y)$ 以最小化

$$\hat{H}(p, q) = -\frac{1}{n} \sum_{i=1}^n \log q(\mathbf{x}_i, y_i)$$

回顾逻辑回归

1. 一些数学推到

$$\begin{aligned} -\frac{1}{n} \sum_{i=1}^n \log q(\mathbf{x}_i, y_i) &= -\frac{1}{n} \sum_{i=1}^n \{\log q(\mathbf{x}_i) + \log q(y_i \mid \mathbf{x}_i)\} \\ &= -\frac{1}{n} \sum_{i=1}^n \log q(\mathbf{x}_i) - \frac{1}{n} \sum_{i=1}^n [y_i \log a_i(\boldsymbol{\theta}) + (1 - y_i) \log\{1 - a_i(\boldsymbol{\theta})\}] \end{aligned}$$

- $q(\mathbf{x}, y) = q(\mathbf{x})q(y \mid \mathbf{x})$
- $q(\mathbf{x})$: 不假定具体形式
- $q(y \mid \mathbf{x}) = a(\boldsymbol{\theta})^y \{1 - a(\boldsymbol{\theta})\}^{1-y}$
- $a(\boldsymbol{\theta}) = \sigma(b + \mathbf{x}^T \mathbf{w})$: 模型参数为 $\boldsymbol{\theta} = (b, \mathbf{w}^T)^T$

回顾逻辑回归

1. 因此，我们有

$$-\frac{1}{n} \sum_{i=1}^n \log q(\mathbf{x}_i, y_i) = -\frac{1}{n} \sum_{i=1}^n \log q(\mathbf{x}_i) - \frac{1}{n} \sum_{i=1}^n [y_i \log a_i(\boldsymbol{\theta}) + (1 - y_i) \log \{1 - a_i(\boldsymbol{\theta})\}]$$

- 其中，我们对第一部分不感兴趣
- 我们对第二部分中的模型参数感兴趣

2. 因此，等价于我们需要最小化下面的式子

$$-\frac{1}{n} \sum_{i=1}^n [y_i \log a_i(\boldsymbol{\theta}) + (1 - y_i) \log \{1 - a_i(\boldsymbol{\theta})\}]$$

- 以上便是从交叉熵出发，推导出的逻辑回归模型的代价函数