



# 第六章 神经网络

李旻先

智能科学与技术系

计算机科学与工程学院

南京理工大学

[minxianli@njust.edu.cn](mailto:minxianli@njust.edu.cn)

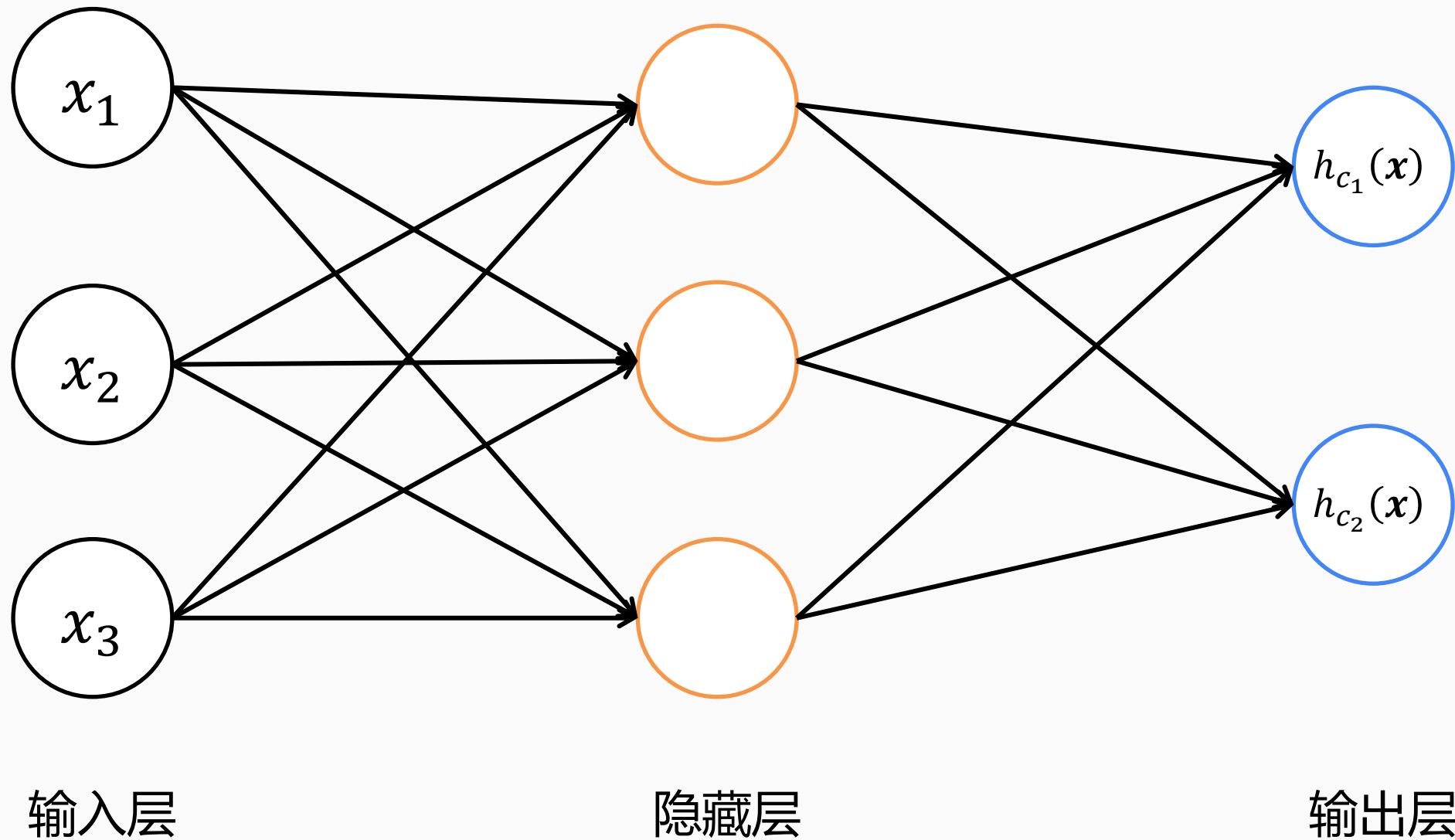


# 神经网络

---

## 网络结构

# 多类神经网络结构



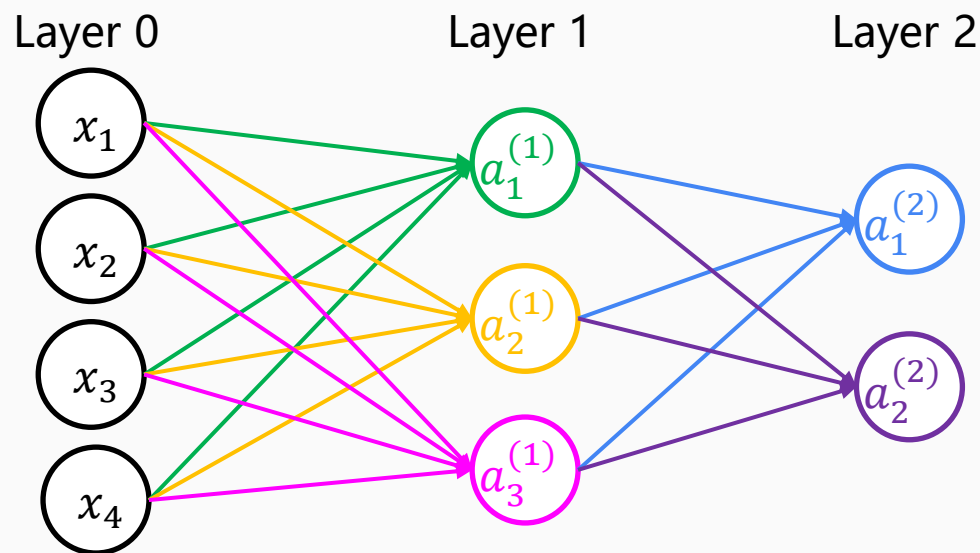


# 神经网络

---

## 假设函数

# 神经网络的模型表示



$L$ : 神经网络的层数 (一般只考虑隐藏层和输出层)

$s_l$ : 第 $l$ 层的神经元数量

$W^{(l)} \in \mathbb{R}^{s_{l-1} \times s_l}$ : 第 $l-1$ 层向第 $l$ 层传播的**权重矩阵**

$b^{(l)} \in \mathbb{R}^{s_l}$ : 第 $l-1$ 层向第 $l$ 层传播的**偏置**

$z^{(l)} \in \mathbb{R}^{s_l}$ : 第 $l$ 层神经元的**净输入值**

$a^{(l)} \in \mathbb{R}^{s_l}$ : 第 $l$ 层神经元的输出值 (**激活值**)

$g_l(\cdot)$ : 第 $l$ 层的神经元的**激活函数**

**传播规则:**

$$\begin{aligned} z^{(l)} &= (W^{(l)})^T \mathbf{a}^{(l-1)} + b^{(l)} \\ \mathbf{a}^{(l)} &= g_l(z^{(l)}) \end{aligned}$$
$$\mathbf{a}^{(l-1)} \in \mathbb{R}^{s_{l-1} \times 1} \longrightarrow \mathbf{a}^{(l)} \in \mathbb{R}^{s_l \times 1}$$

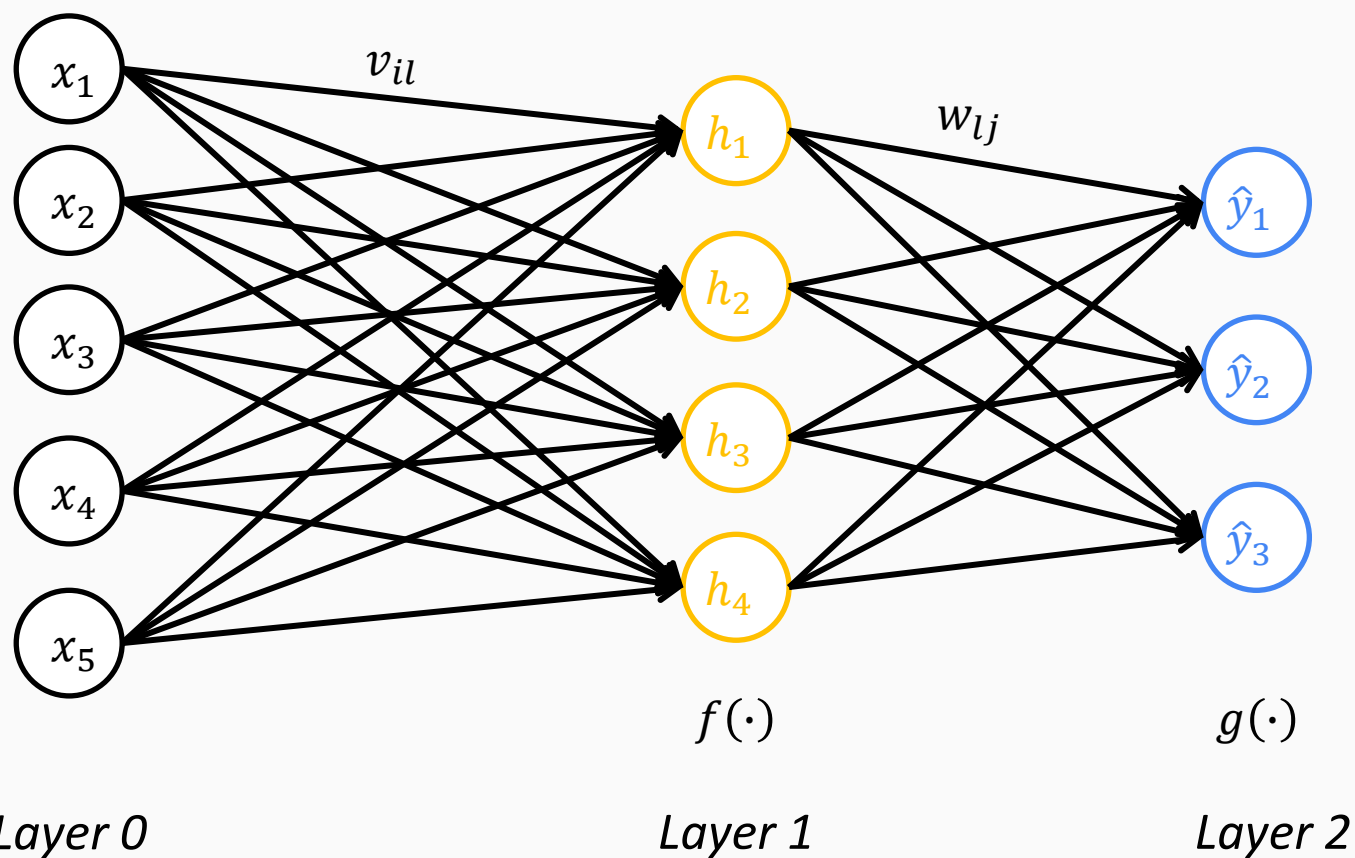
**前向传播:**

$$x = a^{(0)} \rightarrow (z^{(1)} \rightarrow a^{(1)}) \rightarrow \dots \rightarrow (z^{(L)} \rightarrow a^{(L)}) = h_{W,b}(x)$$

$W, b$ 表示所有层的权重矩阵和偏置

# 例题

三层前馈神经网络进行三分类问题建模，其中 $\mathbf{x} \in \mathbb{R}^5$ ,  $\mathbf{h} \in \mathbb{R}^4$ ,  $\hat{\mathbf{y}} \in \mathbb{R}^3$ 分别是输入层、隐藏层和输出层的向量表示。假设输入层到隐藏层为全连接， $h_l = f(p_l)$ ,  $p_l = \sum_{i=1}^5 x_i v_{il}$ ,  $f$ 为sigmoid函数；隐藏层到输出层为全连接， $\hat{y}_j = g(o_j)$ ,  $o_j = \sum_{l=1}^4 h_l w_{lj}$ ,  $g$ 为softmax函数。求该网络的模型假设（即输入 $\mathbf{x}$ 和输出 $\hat{\mathbf{y}}$ 之间的关系）。



$$p_l = \sum_{i=1}^5 x_i v_{il}$$

$$h_l = f(p_l) = \frac{1}{1 + e^{-p_l}}$$

$$o_j = \sum_{l=1}^4 h_l w_{lj}$$

$$\hat{y}_j = g(o_j) = \frac{e^{o_j}}{\sum_{c=1}^C e^{o_c}}$$

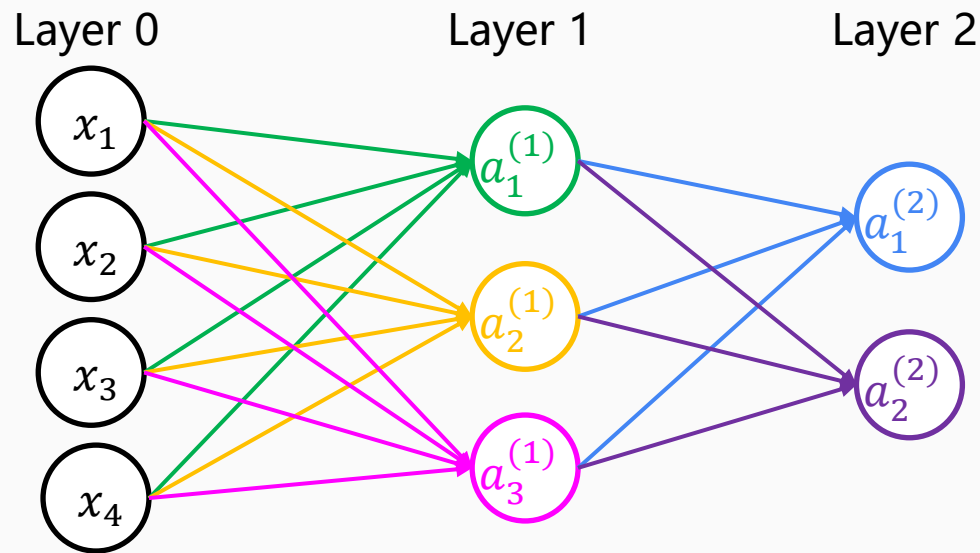


# 神经网络

---

## 损失函数

# 多分类



$L$ : 神经网络的层数 (一般只考虑隐藏层和输出层)

$s_l$ : 第 $l$ 层的神经元数量

$W^{(l)} \in \mathbb{R}^{s_{l-1} \times s_l}$ : 第 $l-1$ 层向第 $l$ 层传播的**权重矩阵**

$b^{(l)} \in \mathbb{R}^{s_l}$ : 第 $l-1$ 层向第 $l$ 层传播的**偏置**

$z^{(l)} \in \mathbb{R}^{s_l}$ : 第 $l$ 层神经元的**净输入值**

$a^{(l)} \in \mathbb{R}^{s_l}$ : 第 $l$ 层神经元的输出值 (**激活值**)

$g_l(\cdot)$ : 第 $l$ 层的神经元的**激活函数**

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

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

**传播规则:**

$$a^{(l-1)} \in \mathbb{R}^{s_{l-1} \times 1} \longrightarrow a^{(l)} \in \mathbb{R}^{s_l \times 1}$$

**前向传播:**

$$x = a^{(0)} \rightarrow (z^{(1)} \rightarrow a^{(1)}) \rightarrow \dots \rightarrow (z^{(L)} \rightarrow a^{(L)} = \text{softmax}(z^{(L)}) = h_{W,b}(x))$$

$W, b$ 表示所有层的权重矩阵和偏置



# 多分类的损失函数——交叉熵

$$a_c^{(L)} = \left( h_{W,b}(\mathbf{x}^{(i)}) \right)_c = \frac{e^{z_c^{(L)}}}{\sum_{c'=1}^C e^{z_{c'}^{(L)}}}$$

假设神经网络有 $C$ 个输出，则 $h_{W,b}(\mathbf{x}) \in \mathbb{R}^C$ ， $(h_{W,b}(\mathbf{x}))_c$ 表示第 $c$ 个输出神经元的预测概率值。

$$J_{W,b}(h_{W,b}(\mathbf{x}), \mathbf{y}) = -\frac{1}{m} \sum_{i=1}^m \sum_{c=1}^C [\mathbf{y}^{(i)}]_c \log \left( (h_{W,b}(\mathbf{x}^{(i)}))_c \right)$$

**根据任务和学习准则的不同，可以设计不同的损失函数**



# 神经网络

---

## 优化方法

# 梯度下降法

损失函数

$$J_{\mathbf{W}, \mathbf{b}}(h_{\mathbf{W}, \mathbf{b}}(\mathbf{x}), \mathbf{y})$$

梯度下降法

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} - \alpha \frac{\partial J_{\mathbf{W}, \mathbf{b}}}{\partial \mathbf{W}^{(l)}}$$

$$\mathbf{b}^{(l)} := \mathbf{b}^{(l)} - \alpha \frac{\partial J_{\mathbf{W}, \mathbf{b}}}{\partial \mathbf{b}^{(l)}}$$

# 梯度

因为 $\mathbf{z}^{(l-1)}$ 与 $\mathbf{z}^{(l)}$ 有这样的关系:  $\mathbf{a}^{(l-1)} = g(\mathbf{z}^{(l-1)})$

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

$$\frac{\partial J_{W,b}}{\partial \mathbf{W}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{W}^{(l)}} \frac{\partial J_{W,b}}{\partial \mathbf{z}^{(l)}} = \mathbf{a}^{(l-1)} \left( \frac{\partial J_{W,b}}{\partial \mathbf{z}^{(l)}} \right)^T$$

直接求导，繁琐低效

$$\frac{\partial J_{W,b}}{\partial \mathbf{b}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{b}^{(l)}} \frac{\partial J_{W,b}}{\partial \mathbf{z}^{(l)}} = \mathbf{I} \frac{\partial J_{W,b}}{\partial \mathbf{z}^{(l)}} = \frac{\partial J_{W,b}}{\partial \mathbf{z}^{(l)}}$$

怎么求？

# 误差

将 $\frac{\partial J_{W,b}}{\partial \mathbf{z}^{(l)}}$ 记为 $\delta^{(l)}$ ，又因为 $\mathbf{z}^{(l)}$ 与 $\mathbf{z}^{(l+1)}$ 有这样的关系：

$$\mathbf{a}^{(l)} = g(\mathbf{z}^{(l)})$$

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

当前层的误差

$$\boxed{\delta^{(l)}} = \frac{\partial J_{W,b}}{\partial \mathbf{z}^{(l)}} = \frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{z}^{(l)}} \frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{a}^{(l)}} \frac{\partial J_{W,b}}{\partial \mathbf{z}^{(l+1)}} = \boxed{\text{diag}(g'(\mathbf{z}^{(l)})) \mathbf{W}^{(l+1)}} \boxed{\delta^{(l+1)}}$$
$$= g'(\mathbf{z}^{(l)}) \odot (\mathbf{W}^{(l+1)} \delta^{(l+1)})$$

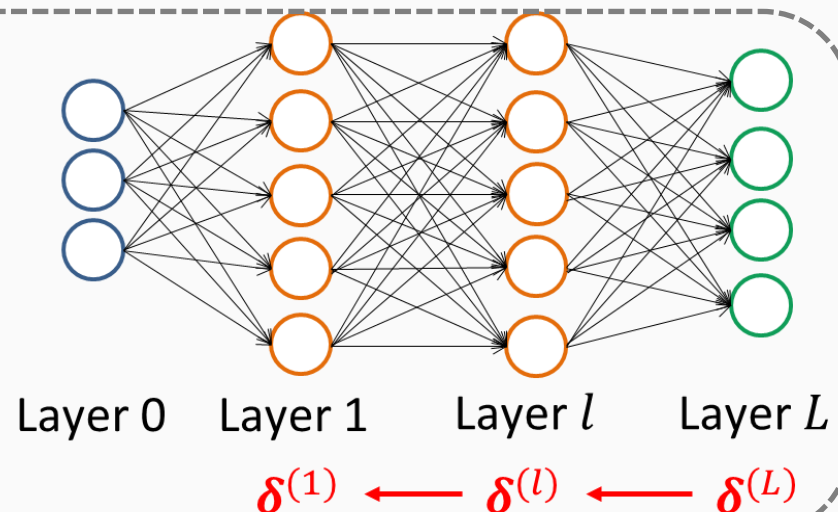
$\odot$ 为Hadamard乘积

当前层的误差 $\delta^{(l)}$ 可由后一层误差 $\delta^{(l+1)}$ 计算得到，则可求当前层的梯度

# 反向（误差）传播算法

- ① 利用损失函数求得模型的最终误差  $\delta^{(L)} = \frac{\partial J_{W,b}}{\partial \mathbf{z}^{(L)}}$ ，接着再将误差自后向前层层传递，获取每层神经元的误差

$$\delta^{(l)} = g'(\mathbf{z}^{(l)}) \odot (\mathbf{W}^{(l+1)} \delta^{(l+1)}).$$



- ② 根据每层神经元的误差  $\delta^{(l)}$  对  $\mathbf{W}^{(l)}$  和  $\mathbf{b}^{(l)}$  求偏导，

求得梯度  $\frac{\partial J_{W,b}}{\partial \mathbf{W}^{(l)}}$  和  $\frac{\partial J_{W,b}}{\partial \mathbf{b}^{(l)}}$ 。

$$\frac{\partial J_{W,b}}{\partial \mathbf{W}^{(l)}} = \mathbf{a}^{(l-1)} (\delta^{(l)})^T$$

$$\frac{\partial J_{W,b}}{\partial \mathbf{b}^{(l)}} = \delta^{(l)}$$

- ③ 根据梯度，更新参数  $\mathbf{W}^{(l)}$  和  $\mathbf{b}^{(l)}$ 。

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} - \alpha \frac{\partial J_{W,b}}{\partial \mathbf{W}^{(l)}}$$

$$\mathbf{b}^{(l)} := \mathbf{b}^{(l)} - \alpha \frac{\partial J_{W,b}}{\partial \mathbf{b}^{(l)}}$$

# 损失函数——基于Softmax函数的交叉熵

$$J_{W,b}(h_{W,b}(x), \mathbf{y}) = \sum_{c=1}^C -[\mathbf{y}]_c \log(h_{W,b}(x))_c = \sum_{c=1}^C -[\mathbf{y}]_c \log a_c^{(L)} = -\log(a_y^{(L)})$$

$$\mathbf{a}^{(L)} = g(\mathbf{z}^{(L)}) = \text{softmax}(\mathbf{z}^{(L)})$$

$$\delta_c^{(L)} = \frac{\partial J_{W,b}}{\partial z_c^{(L)}} = \frac{\partial a_y^{(L)}}{\partial z_c^{(L)}} \frac{\partial J_{W,b}}{\partial a_y^{(L)}} = - \frac{\partial a_y^{(L)}}{\partial z_c^{(L)}} \frac{1}{a_y^{(L)}}$$

Softmax函数的导数

$$\frac{\partial g(\mathbf{z})_i}{\partial z_j} = \begin{cases} g(\mathbf{z})_i(1 - g(\mathbf{z})_i) & i = j \\ -g(\mathbf{z})_i g(\mathbf{z})_j & i \neq j \end{cases}$$

# 损失函数——交叉熵

Softmax函数的导数  $\frac{\partial g(\mathbf{z})_i}{\partial z_j} = \begin{cases} g(\mathbf{z})_i(1 - g(\mathbf{z})_i) & i = j \\ -g(\mathbf{z})_i g(\mathbf{z})_j & i \neq j \end{cases}$

$$\delta_c^{(L)} = \frac{\partial J_{W,b}}{\partial z_c^{(L)}} = - \frac{\partial a_y^{(L)}}{\partial z_c^{(L)}} \frac{1}{a_y^{(L)}} = \begin{cases} a_y^{(L)} (a_y^{(L)} - 1) \frac{1}{a_y^{(L)}} = a_c^{(L)} - 1 & y = c \\ a_y^{(L)} a_c^{(L)} \frac{1}{a_y^{(L)}} = a_c^{(L)} & y \neq c \end{cases}$$

$$= a_c^{(L)} - [\mathbf{y}]_c$$

$$\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - \mathbf{y}$$



# 损失函数——平方误差

$$J_{W,b}(h_{W,b}(\mathbf{x}), \mathbf{y}) = \frac{1}{2} \|h_{W,b}(\mathbf{x}) - \mathbf{y}\|_2^2 = \frac{1}{2} \sum_{c=1}^C (a_c^{(L)} - [\mathbf{y}]_c)^2$$

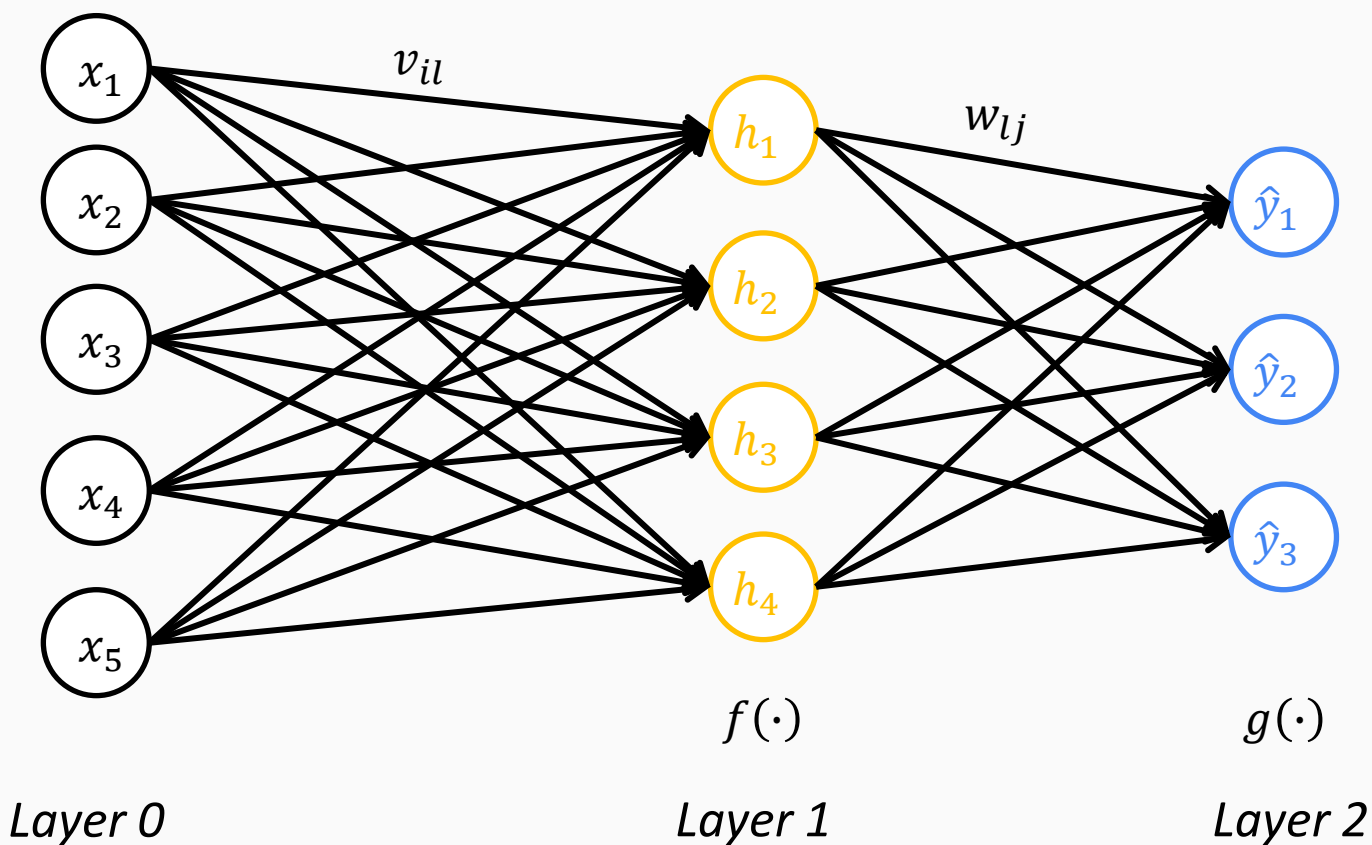
$$\mathbf{a}^{(L)} = g(\mathbf{z}^{(L)}) = \text{sigmoid}(\mathbf{z}^{(L)})$$

$$\delta_c^{(L)} = \frac{\partial J_{W,b}}{\partial z_c^{(L)}} = \frac{\partial a_c^{(L)}}{\partial z_c^{(L)}} \frac{\partial J_{W,b}}{\partial a_c^{(L)}} = a_c^{(L)} (1 - a_c^{(L)}) [a_c^{(L)} - [\mathbf{y}]_c]$$

$$\boldsymbol{\delta}^{(L)} = \text{diag}\left(a_c^{(L)} (1 - a_c^{(L)})\right) [\mathbf{a}^{(L)} - \mathbf{y}]$$

# 例题

三层前馈神经网络进行三分类问题建模，其中 $\mathbf{x} \in \mathbb{R}^5$ ,  $\mathbf{h} \in \mathbb{R}^4$ ,  $\hat{\mathbf{y}} \in \mathbb{R}^3$ 分别是输入层、隐藏层和输出层的向量表示。假设输入层到隐藏层为全连接， $h_l = f(p_l)$ ,  $p_l = \sum_{i=1}^5 x_i v_{il}$ ,  $f$ 为sigmoid函数；隐藏层到输出层为全连接， $\hat{y}_j = g(o_j)$ ,  $o_j = \sum_{l=1}^4 h_l w_{lj}$ ,  $g$ 为softmax函数。若损失函数为交叉熵损失，求 $w_{32}$ 和 $v_{43}$ 的更新规则。



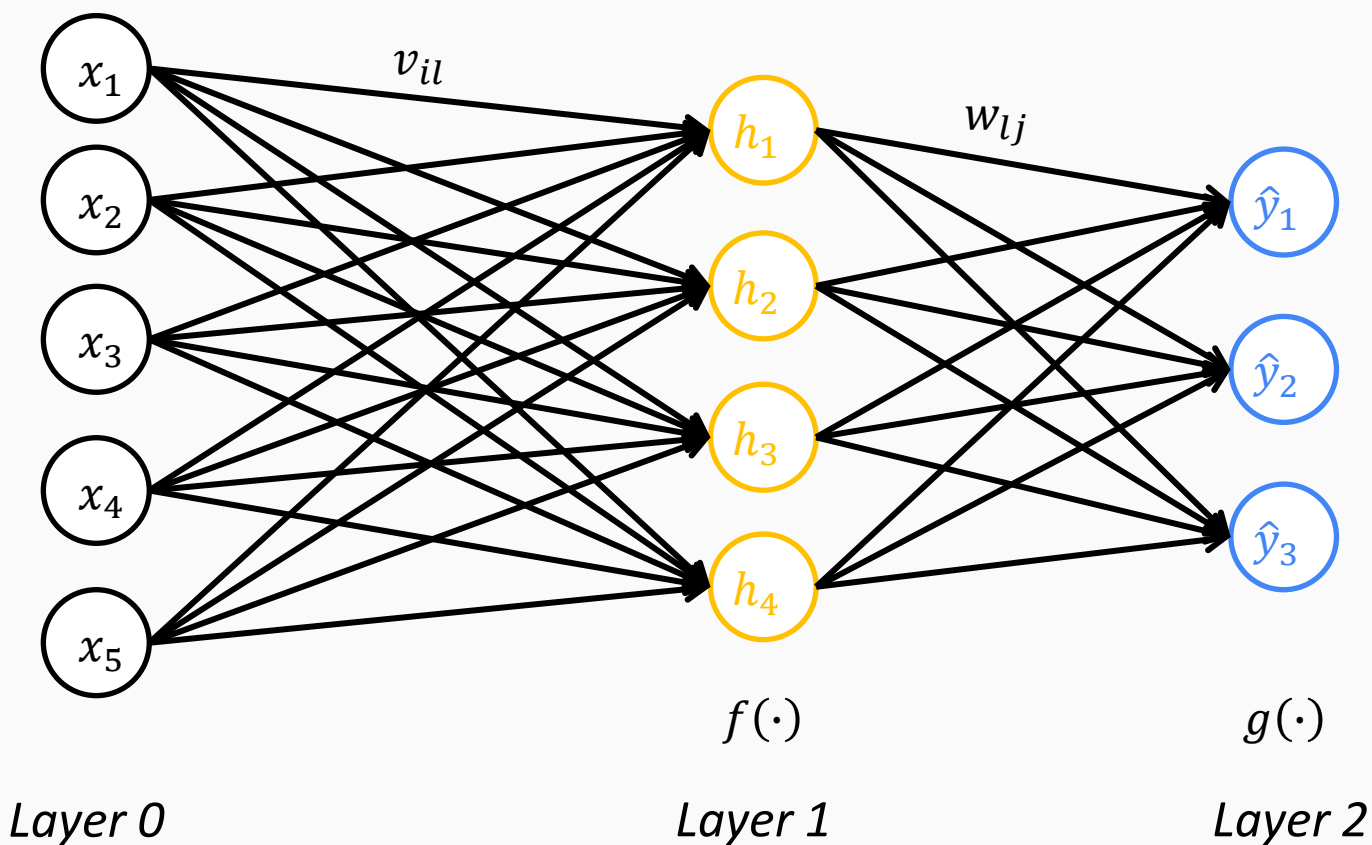
$\mathbf{W}^{(l)} \in \mathbb{R}^{s_{l-1} \times s_l}$ : 第 $l-1$ 层向第 $l$ 层传播的权重矩阵

①  $v_{43} \in \mathbf{W}^{(1)}, \mathbf{W}^{(1)} \in \mathbb{R}^{5 \times 4}$

②  $w_{32} \in \mathbf{W}^{(2)}, \mathbf{W}^{(2)} \in \mathbb{R}^{4 \times 3}$

# 例题

三层前馈神经网络进行三分类问题建模，其中 $\mathbf{x} \in \mathbb{R}^5$ ,  $\mathbf{h} \in \mathbb{R}^4$ ,  $\hat{\mathbf{y}} \in \mathbb{R}^3$ 分别是输入层、隐藏层和输出层的向量表示。假设输入层到隐藏层为全连接， $h_l = f(p_l)$ ,  $p_l = \sum_{i=1}^5 x_i v_{il}$ ,  $f$ 为sigmoid函数；隐藏层到输出层为全连接， $\hat{y}_j = g(o_j)$ ,  $o_j = \sum_{l=1}^4 h_l w_{lj}$ ,  $g$ 为softmax函数。若损失函数为交叉熵损失，求 $w_{32}$ 和 $v_{43}$ 的更新规则。



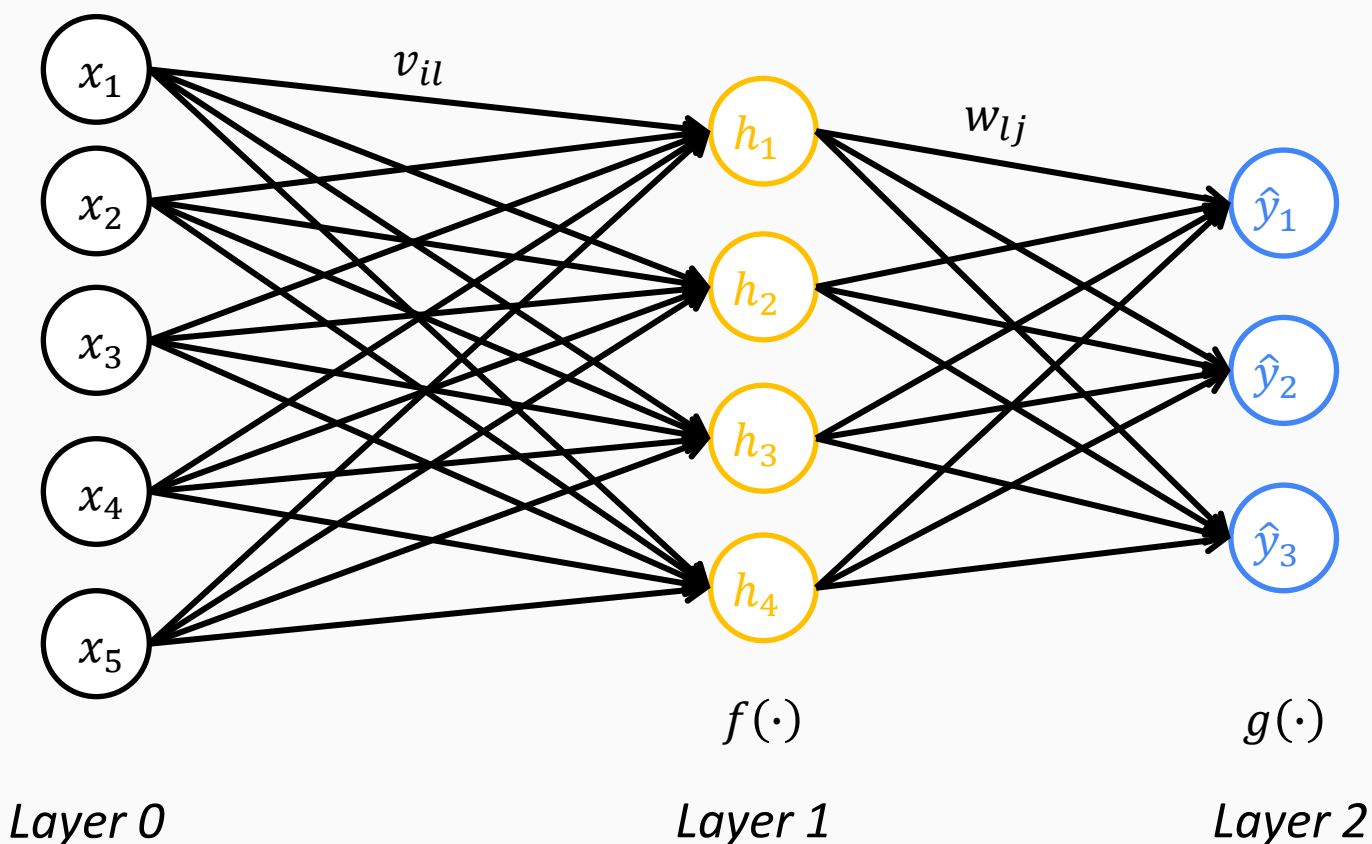
$$\textcircled{1} \mathbf{W}^{(l)} := \mathbf{W}^{(l)} - \alpha \frac{\partial J_{\mathbf{W}, \mathbf{b}}}{\partial \mathbf{W}^{(l)}}$$

$$\textcircled{2} \frac{\partial J_{\mathbf{W}, \mathbf{b}}}{\partial \mathbf{W}^{(l)}} = \mathbf{a}^{(l-1)} \left( \frac{\partial J_{\mathbf{W}, \mathbf{b}}}{\partial \mathbf{z}^{(l)}} \right)^T$$

$$\textcircled{3} \frac{\partial J_{\mathbf{W}, \mathbf{b}}}{\partial \mathbf{z}^{(l)}} = \boldsymbol{\delta}^{(l)} = g'(\mathbf{z}^{(l)}) \odot \mathbf{W}^{(l+1)} \boldsymbol{\delta}^{(l+1)}$$

# 例题

三层前馈神经网络进行三分类问题建模，其中 $\mathbf{x} \in \mathbb{R}^5$ ,  $\mathbf{h} \in \mathbb{R}^4$ ,  $\hat{\mathbf{y}} \in \mathbb{R}^3$ 分别是输入层、隐藏层和输出层的向量表示。假设输入层到隐藏层为全连接， $h_l = f(p_l)$ ,  $p_l = \sum_{i=1}^5 x_i v_{il}$ ,  $f$ 为sigmoid函数；隐藏层到输出层为全连接， $\hat{y}_j = g(o_j)$ ,  $o_j = \sum_{l=1}^4 h_l w_{lj}$ ,  $g$ 为softmax函数。若损失函数为交叉熵损失，求 $w_{32}$ 和 $v_{43}$ 的更新规则。



$$\text{损失函数 } J_{W,b} = -\log(a_y^{(2)})$$

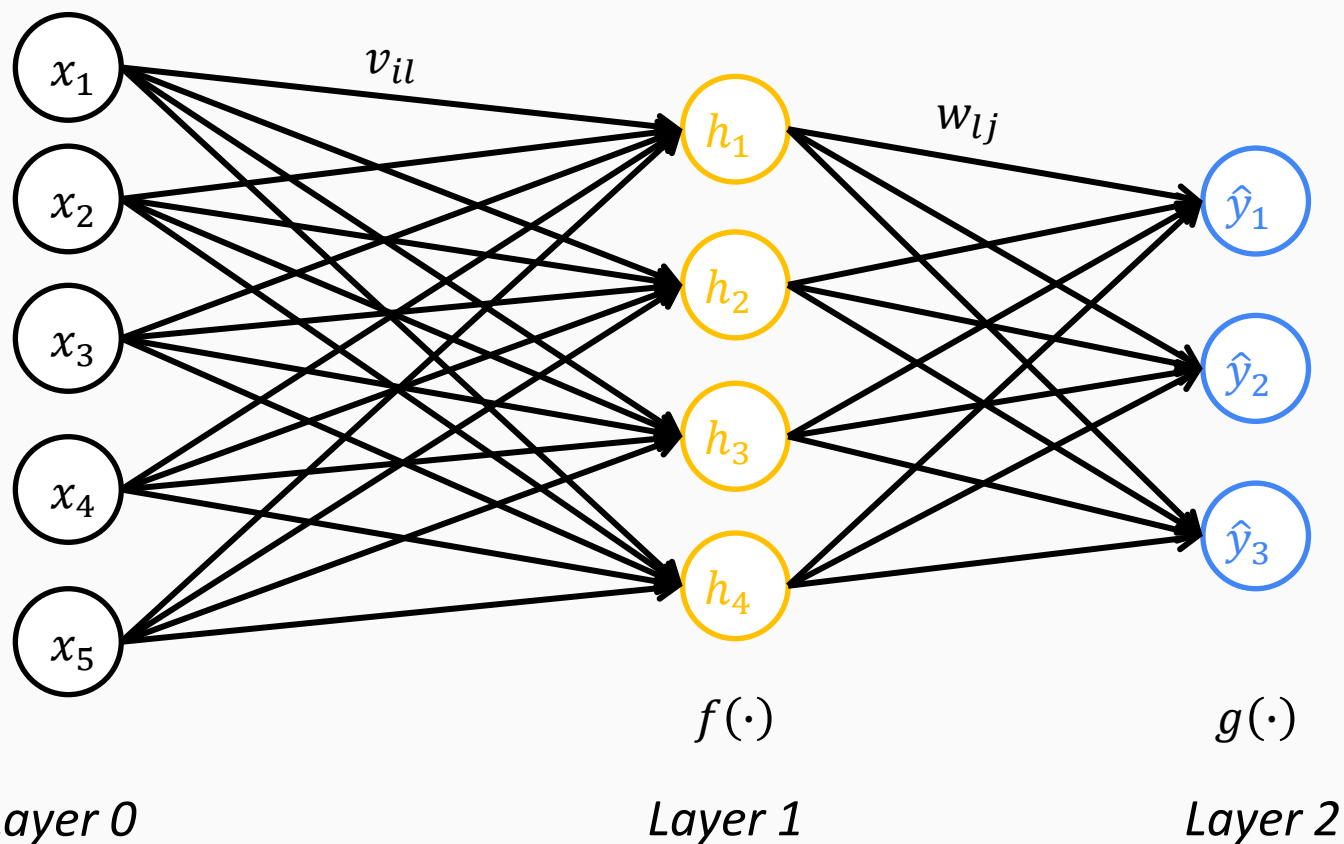
$$\delta^{(2)} = \frac{\partial J_{W,b}}{\partial \mathbf{z}^{(2)}} = \mathbf{a}^{(2)} - \mathbf{y} \\ = \hat{\mathbf{y}} - \mathbf{y}$$

$$\frac{\partial J_{W,b}}{\partial \mathbf{W}^{(2)}} = \mathbf{a}^{(1)} (\delta^{(2)})^T = \mathbf{h} (\hat{\mathbf{y}} - \mathbf{y})^T$$

$$\frac{\partial J_{W,b}}{\partial w_{32}} = h_3 (\hat{y}_2 - y_2)$$

# 例题

三层前馈神经网络进行三分类问题建模，其中 $\mathbf{x} \in \mathbb{R}^5$ ,  $\mathbf{h} \in \mathbb{R}^4$ ,  $\hat{\mathbf{y}} \in \mathbb{R}^3$ 分别是输入层、隐藏层和输出层的向量表示。假设输入层到隐藏层为全连接， $h_l = f(p_l)$ ,  $p_l = \sum_{i=1}^5 x_i v_{il}$ ， $f$ 为sigmoid函数；隐藏层到输出层为全连接， $\hat{y}_j = g(o_j)$ ,  $o_j = \sum_{l=1}^4 h_l w_{lj}$ ， $g$ 为softmax函数。若损失函数为交叉熵损失，求 $w_{32}$ 和 $v_{43}$ 的更新规则。



$$\delta^{(l)} = \text{diag}(g'(\mathbf{z}^{(l)})) \mathbf{W}^{(l+1)} \delta^{(l+1)}$$

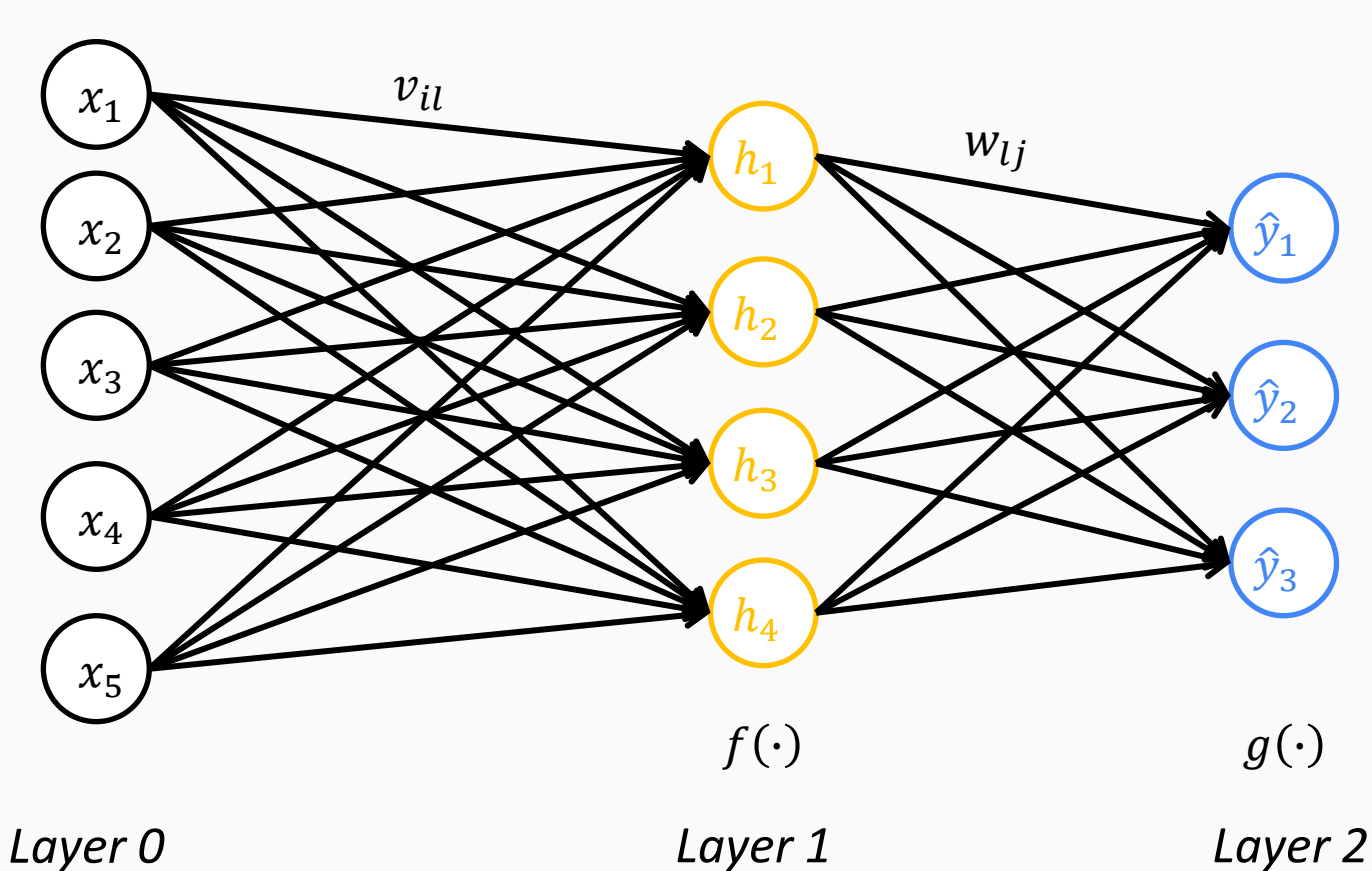
$$\delta^{(1)} = \text{diag}(f'(\mathbf{p})) \mathbf{W}^{(2)} \delta^{(2)}$$

$$= \begin{bmatrix} f'(p_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & f'(p_4) \end{bmatrix} \begin{bmatrix} w_{11} & \cdots & w_{13} \\ \vdots & \ddots & \vdots \\ w_{41} & \cdots & w_{43} \end{bmatrix} (\hat{\mathbf{y}} - \mathbf{y})$$

$$= \begin{bmatrix} f'(p_1) \sum_{j=1}^3 w_{1j} (\hat{y}_j - y_j) \\ f'(p_2) \sum_{j=1}^3 w_{2j} (\hat{y}_j - y_j) \\ f'(p_3) \sum_{j=1}^3 w_{3j} (\hat{y}_j - y_j) \\ f'(p_4) \sum_{j=1}^3 w_{4j} (\hat{y}_j - y_j) \end{bmatrix}$$

# 例题

三层前馈神经网络进行三分类问题建模，其中 $\mathbf{x} \in \mathbb{R}^5$ ,  $\mathbf{h} \in \mathbb{R}^4$ ,  $\hat{\mathbf{y}} \in \mathbb{R}^3$ 分别是输入层、隐藏层和输出层的向量表示。假设输入层到隐藏层为全连接， $h_l = f(p_l)$ ,  $p_l = \sum_{i=1}^5 x_i v_{il}$ ,  $f$ 为sigmoid函数；隐藏层到输出层为全连接， $\hat{y}_j = g(o_j)$ ,  $o_j = \sum_{l=1}^4 h_l w_{lj}$ ,  $g$ 为softmax函数。若损失函数为交叉熵损失，求 $w_{32}$ 和 $v_{43}$ 的更新规则。

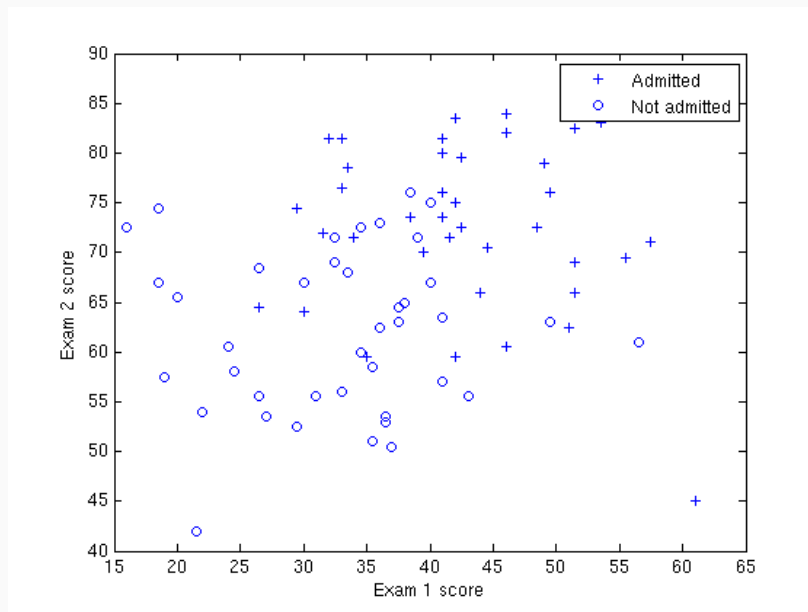


$$\frac{\partial J_{\mathbf{W}, \mathbf{b}}}{\partial \mathbf{W}^{(1)}} = \mathbf{a}^{(0)} (\boldsymbol{\delta}^{(1)})^T = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} \begin{bmatrix} f'(p_1) \sum_{j=1}^3 w_{1j} (\hat{y}_j - y_j) \\ f'(p_2) \sum_{j=1}^3 w_{2j} (\hat{y}_j - y_j) \\ f'(p_3) \sum_{j=1}^3 w_{3j} (\hat{y}_j - y_j) \\ f'(p_4) \sum_{j=1}^3 w_{4j} (\hat{y}_j - y_j) \end{bmatrix}^T$$

$$\frac{\partial J_{\mathbf{W}, \mathbf{b}}}{\partial v_{43}} = x_4 h_3 (1 - h_3) \sum_{j=1}^3 w_{3j} (\hat{y}_j - y_j)$$

# 作业5: 基于BP算法的三层前向神经网络

- 给出训练数据:



<http://openclassroom.stanford.edu/MainFolder/DocumentPage.php?course=DeepLearning&doc=exercises/ex4/ex4.html>

- 使用BP算法实现三层前向神经网络 (自己编码, 不要使用Tensorflow/Pytorch等框架), 并对结果进行5倍交叉验证;
- 将其与Logistic回归和Softmax回归进行比较。