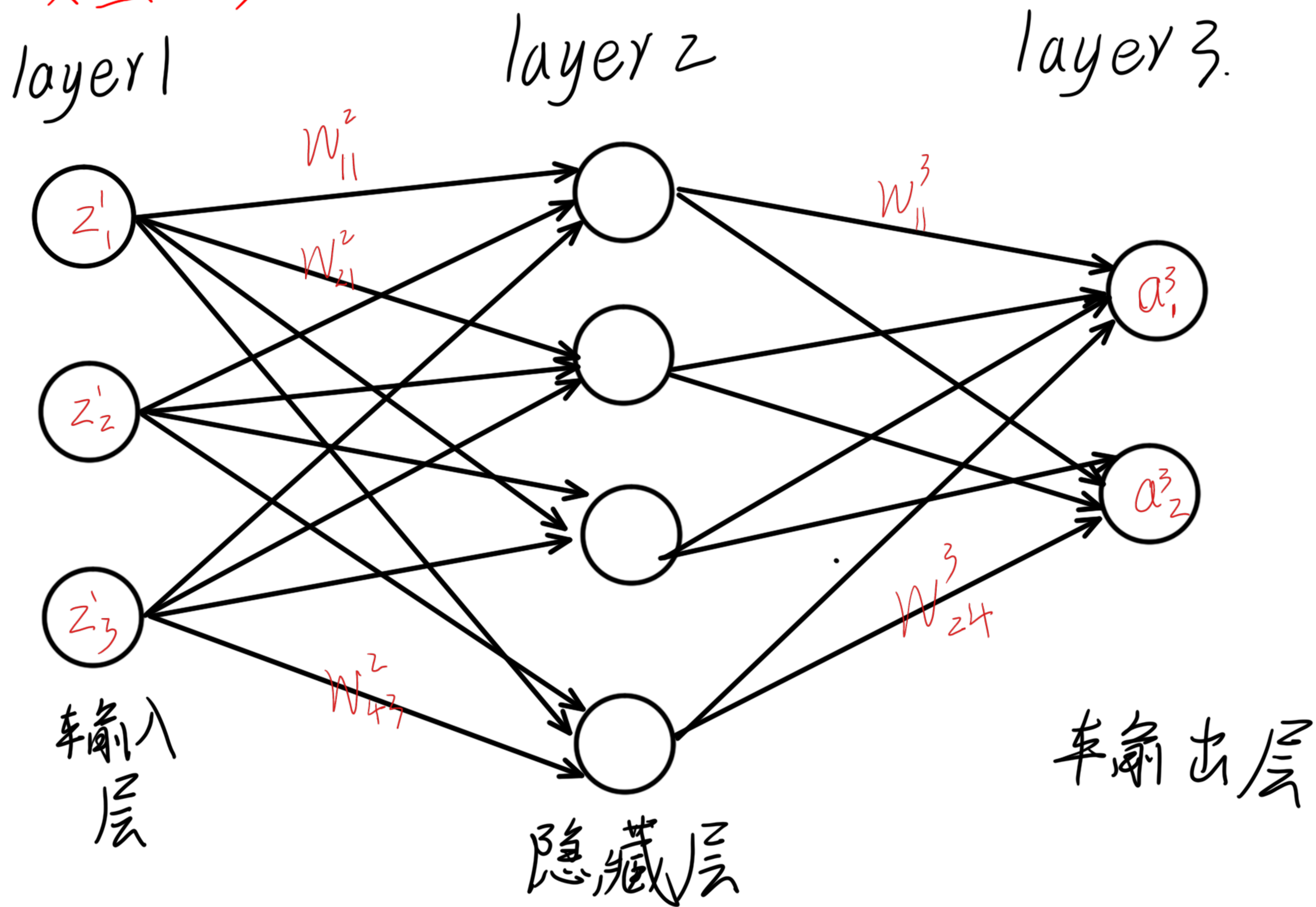上图是一个三层神经网络，对变量进行定义：

$w_{jk}^l$ 表示第 $(l-1)$ 层的第 $K$ 个神经元连接到第 $l$ 层第 $j$ 个神经元的权重。

$b_j^l$ 表示第 $l$ 层第 $j$ 个神经元的偏置。

$z_j^l$ 表示第 $l$ 层第 $j$ 个神经元的输入，即：$z_j^l = \sum_K w_{jk}^l a_k^{l-1} + b_j^l$

$a_j^l$ 表示第 $l$ 层第 $j$ 个神经元的输入，即：$a_j^l = 6(\sum_K w_{jk}^l a_k^{l-1} + b_j^l)$

$6$ 表示激活函数，常见的例如 Sigmoid 函数、Relu 函数等。 $= 6(z_j^l)$

## 2. 公式推导.

首先，将第 $l$ 层第 $j$ 个神经元中产生的错误（即实际值与预测值之间的误差）定义为：

$$\delta_j^l = \frac{\partial c}{\partial z_j^l}$$

以一个输入样本为例进行说明，此时损失函数与代价函数的均方误差表示形式均为：

$$C = \frac{1}{2} \sum_j (y_j - a_j^L)^2$$

$y_j$表示真实值，$a_j^L$表示预测的输出值

<span style="color:cyan">公式11计算最后一层神经网络产生的误差 )</span>

⭐ $\delta^L = \nabla_a C \odot \sigma'(z^L) = \frac{\partial L}{\partial z^L} \odot \frac{\partial a^L}{\partial z^L}$

其中，$\odot$表示Hadamard乘积，用于矩阵或向量之间点对点的乘法运算，$\nabla$表示梯度运算符。

公式一推导过程：

$$\because \delta_j^L = \frac{\partial C}{\partial z_j^L} = \frac{\partial L}{\partial a_j^L} \cdot \frac{\partial a_j^L}{\partial z_j^L}$$

<span style="color:violet">⟹ 例. $f = [f_1, f_2 \cdots, f_m]$ $X = [x_1, x_2, \cdots x_n]$</span>

<span style="color:violet">$$\nabla_x f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}, & \frac{\partial f_1}{\partial x_2}, & \cdots\cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1}, & \frac{\partial f_m}{\partial x_2}, & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$</span>

$$\therefore \delta^L = \frac{\partial L}{\partial z^L} \odot \frac{\partial a^L}{\partial z^L} = \nabla_a C \odot \sigma'(z^L)$$

公式 2 ( 由后向前，计算 每 一层神经网络产生的误差 )：

$$☆ \delta^l = ((W^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l) \quad (T表示矩阵转置运算)$$

推导过程：

$$\because \delta_j^l = \frac{\partial C}{\partial z_j^l} = \sum_k \frac{\partial C}{\partial z_k^{l+1}} \cdot \frac{\partial z_k^{l+1}}{\partial a_j^l} \cdot \frac{\partial a_j^l}{\partial z_j^l}$$

$$= \sum_k \delta_k^{l+1} \cdot \frac{\partial (W_{kj}^{l+1} a_j^l + b_k^{l+1})}{\partial a_j^l} \cdot \sigma'(z_j^l)$$

$$= \sum_k \delta_k^{l+1} \cdot W_{kj}^{l+1} \cdot \sigma'(z_j^l)$$

$$\therefore \delta^l = ((W^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

## 公式 3（计算权重梯度）:

✗ $\dfrac{\partial C}{\partial w_{jk}^{l}} = a_k^{l-1} \delta_j^l$

推导过程:

$$\dfrac{\partial C}{\partial w_{jk}^{l}} = \dfrac{\partial C}{\partial z_j^{l}} \cdot \dfrac{\partial z_j^{l}}{\partial w_{jk}^{l}} = \delta_j^l \cdot \dfrac{\partial (w_{jk}^l a_k^{l-1} + b_j^l)}{\partial w_{jk}^{l}} = a_k^{l-1} \delta_j^l$$

## 公式 4（计算偏置的梯度）:

✗ $\dfrac{\partial C}{\partial b_j^{l}} = \delta_j^l$

推导过程:

$$\dfrac{\partial C}{\partial b_j^{l}} = \dfrac{\partial C}{\partial z_j^{l}} \cdot \dfrac{\partial z_j^{l}}{\partial b_j^{l}} = \delta_j^l \cdot \dfrac{\partial (w_{jk}^l a_k^{l-1} + b_j^l)}{\partial b_j^{l}} = \delta_j^l$$

# 更新权重和偏置

设学习率为 $\eta$

$$(W_{jk}^{l})_{新} = W_{j}^{l} - \eta \cdot \frac{\partial C}{\partial W_{jk}^{l}} = W_{jk}^{l} - \eta \, a_{k}^{l-1} \delta_{j}^{l}$$

$$(b_{j}^{l})_{新} = b_{j}^{l} - \eta \cdot \frac{\partial C}{\partial b_{j}^{l}} = b_{j}^{l} - \eta \, \delta_{j}^{l}$$

✳ 总结：四个公式：

① $\delta^{L} = \nabla_{a} C \odot \sigma'(z^{L}) = \frac{\partial L}{\partial z^{L}} \odot \frac{\partial a^{L}}{\partial z^{L}}$

② $\delta^{l} = ((W^{l+1})^{T} \delta^{l+1}) \odot \sigma'(z^{l})$

③ $\frac{\partial C}{\partial W_{jk}^{l}} = a_{k}^{l-1} \delta_{j}^{l}$

④ $\dfrac{\partial c}{\partial b_j^l} = \delta_j^l$