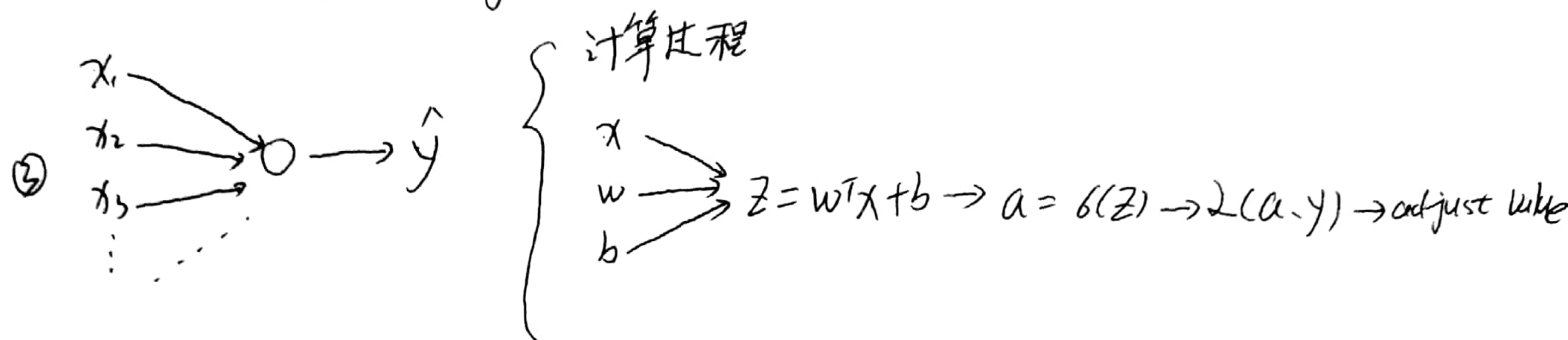


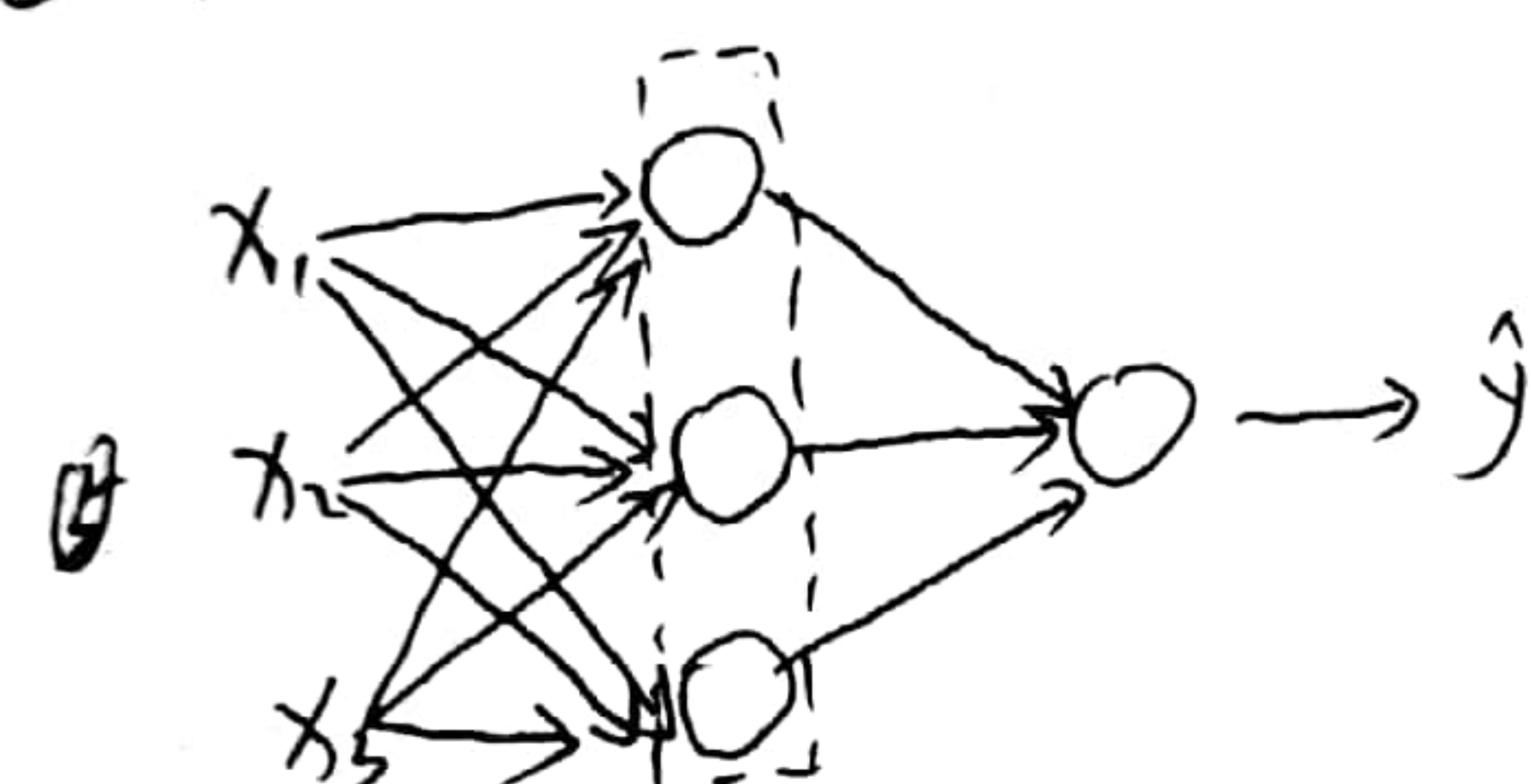
3.1 神经网络概览

①在这一节中(第三单元)我们将学习如何搭建一个神经网络模型. 在此之前, 我们先来看一下神经网络的大概. 有些看不懂没关系, 后面会有解释

②在上一周我们学习了 logistic 回归, 它的模型如下



③之前我们介绍过, 神经网络模型像下面这样



④很明显, 神经网络模型就是一个 logistic 回归堆叠起来的. 它的每个节点都包含两步运算

即 ① $z = w^T x + b$ ② $a = \sigma(z)$

⑤之后我们会使用 [] 标注的方式来标识神经网络中的层级.

比如, 在 ⑤ 图中虚线部分为神经网络第一层, 有关它的计算可以表示为

$z^{[1]} = w^{[1]}x + b^{[1]}$ 用方括号来区分训练集中的样本表示

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

⑥在计算完成 $a^{[1]}$ 之后, 我们就可以计算下一层, 即第二层

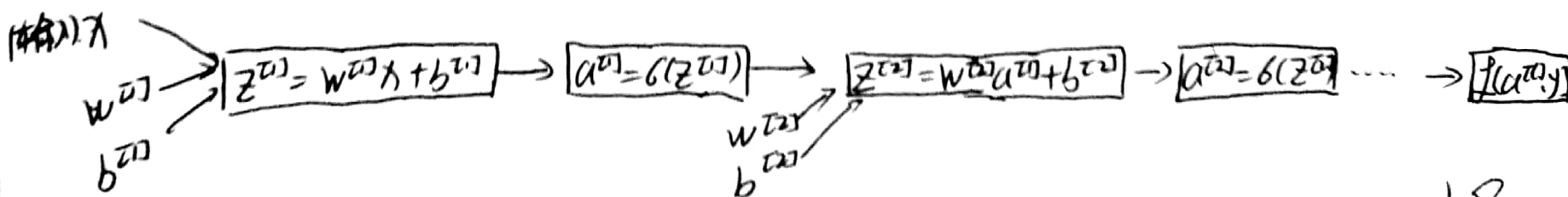
$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$

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

⑦在所有层的运算结束之后, 我们会计算损失函数

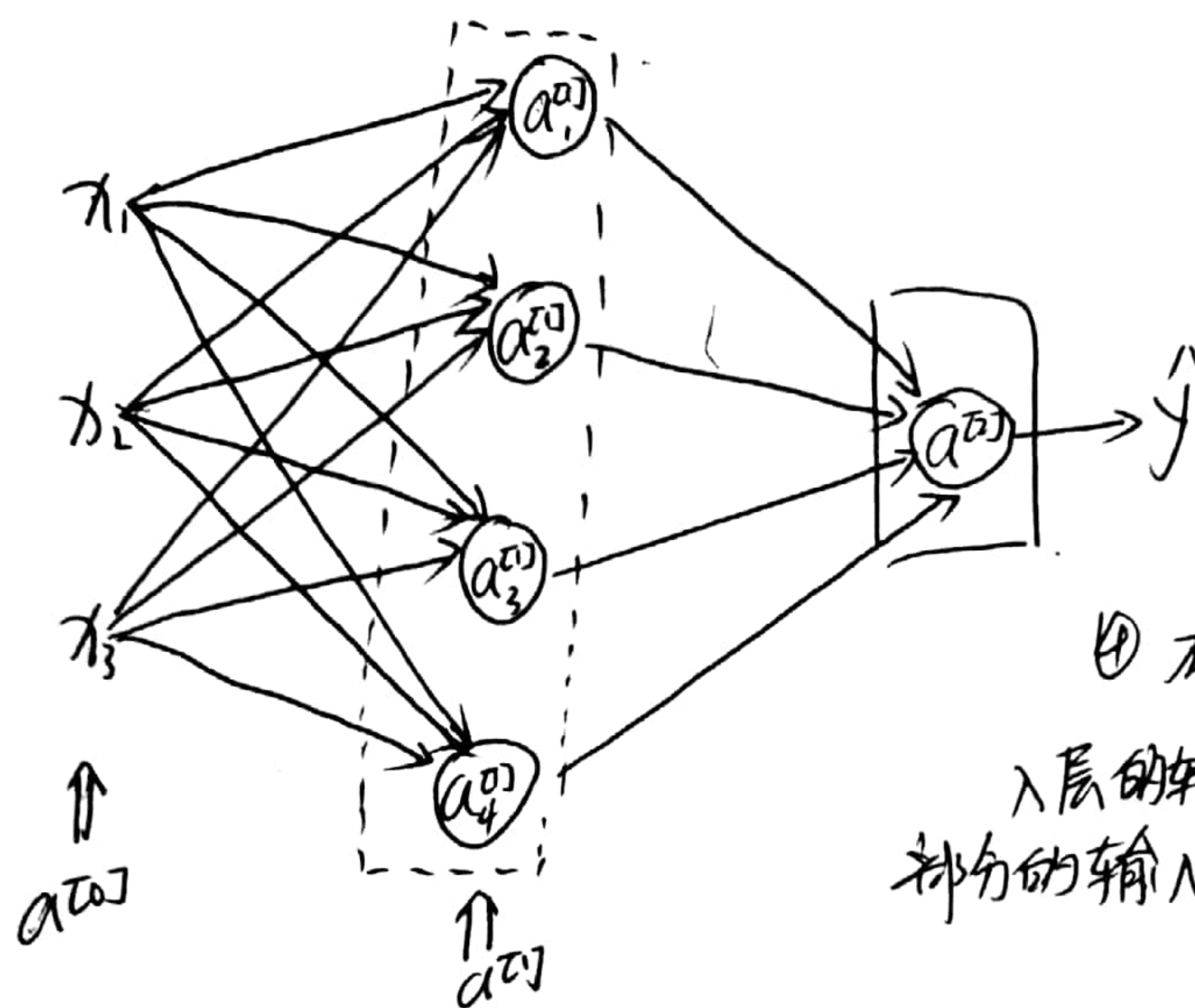
$L(a^{[last layer]}, y)$

⑧上面的运算可以用下面的流程图来表示



3.2 神经网络表示

①我们先来看一个只有一个隐藏层的神经网络。



① x_1, x_2, x_3 (输入层)

② $\boxed{}$ 隐藏层

③ $\boxed{}$ 输出层

④ 在我们运行模型时，我们能看到输入层的输入，能看到输出层的输出，但是中间虚线部分的输入输出我们是看不到的，因此称为——

⑤ 我们应该清楚，上面的每一层都是带有参数的，
即第一层有 $w^{(1)}$ 和 $b^{(1)}$ 。

第二层有 $w^{(2)}$ 和 $b^{(2)}$

⑥ 在上面的例子中， $w^{(1)}$ 是一个 4×3 矩阵， $b^{(1)}$ 是 4×1 矩阵

$w^{(1)}$ 是 4×3 矩阵，其 4 是因为第一层有四个节点

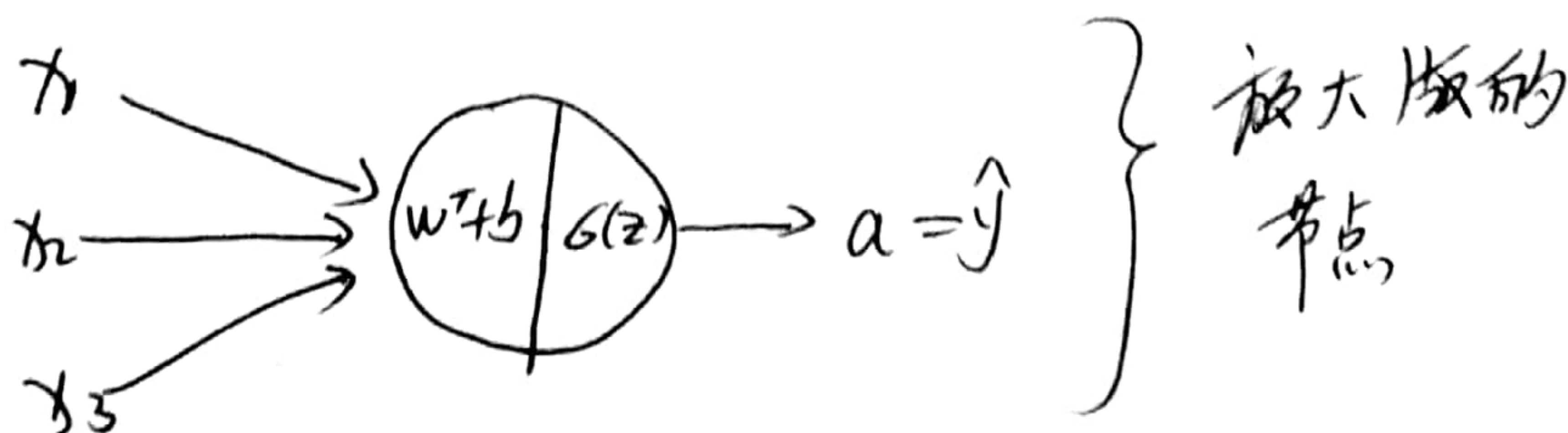
3 是因为有 3 个输入特征

同理： $w^{(2)}$ 是一个 1×3 矩阵， $b^{(2)}$ 是 1×1 矩阵

3.3 计算神经网络的输出

承接 3.2 中的图。

在 logistic 回归中，节点的运算如下。



对隐藏层第一个节点、计算过程如下

$$z_1^{(1)} = w_1^{(0)T} x + b_1^{(1)} \quad a_1^{(1)} = \sigma(z_1^{(1)})$$

} 和 logistic 回归很像。

同理、对于第二个节点、计算过程如下：

$$z_2^{(1)} = w_2^{(0)T} x + b_2^{(1)} \quad a_2^{(1)} = \sigma(z_2^{(1)})$$

下面是第一个隐藏层的所有节点的计算公式

$$\left\{ \begin{array}{l} z_1^{(1)} = w_1^{(0)T} x + b_1^{(1)} \quad a_1^{(1)} = \sigma(z_1^{(1)}) \\ z_2^{(1)} = w_2^{(0)T} x + b_2^{(1)} \quad a_2^{(1)} = \sigma(z_2^{(1)}) \\ z_3^{(1)} = w_3^{(0)T} x + b_3^{(1)} \quad a_3^{(1)} = \sigma(z_3^{(1)}) \\ z_4^{(1)} = w_4^{(0)T} x + b_4^{(1)} \quad a_4^{(1)} = \sigma(z_4^{(1)}) \end{array} \right.$$

- 一个想法自然是使用 for 循环计算各节点值。
但使用向量化会使计算更加快速，让我们看看怎么向量化它们

要向量化，首先应该想，第一层的作用是什么？它最终产生的数据格式应该是什么样的
① 第一层作用就是对输入层数据做一次操作，产生输出作为下一层的输入。
既然是作为下一层的输入，那么格式应该和输入层数据格式一样
因此，最终，第一层的输出应该像下面这样

$$\begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{bmatrix} = \sigma \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix} = \begin{bmatrix} w_1^{(0)T} x + b_1^{(1)} \\ w_2^{(0)T} x + b_2^{(1)} \\ w_3^{(0)T} x + b_3^{(1)} \\ w_4^{(0)T} x + b_4^{(1)} \end{bmatrix} = \begin{bmatrix} w_1^{(0)T} \\ w_2^{(0)T} \\ w_3^{(0)T} \\ w_4^{(0)T} \end{bmatrix} \cdot x + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix}$$

因此我们知道，在单个节点中

$w_i^{(0)} = \begin{bmatrix} w_{i1} \\ w_{i2} \\ w_{i3} \end{bmatrix}$ 是 3×1 矩阵，而 $w_i^{(0)T}$ 就是 1×3 矩阵。

这样向量化的 $w^{(0)}$ 即为 4×3 矩阵

向量化的 b 为 4×1 矩阵

$$W^{(0)} = \begin{bmatrix} \dots & w_1^{(0)T} & \dots \\ \dots & w_2^{(0)T} & \dots \\ \dots & w_3^{(0)T} & \dots \\ \dots & w_4^{(0)T} & \dots \end{bmatrix}$$

$$b^{(1)} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix}$$

也就是，在 $W^{(0)}$ 中，
行数是节点数，
列数是输入特征维度。

因此对于4.2节中的两层神经网络来说, 计算步骤如下

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

$[4 \times 1]$ $[4 \times 3]$ $[3 \times 1]$ $[4 \times 1]$

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

$[4 \times 1]$ $[4 \times 1]$

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

$[1 \times 1]$ $[1 \times 4]$ $[4 \times 1]$ $[1 \times 1]$

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

$[4 \times 1]$ $[1 \times 1]$

没错.

就是这么简单

其中 $z^{[1]} (4 \times 1)$ $W^{[1]} [4 \times 3]$ $x (3 \times 1)$ $b^{[1]} (4 \times 1)$

$z^{[2]} (1 \times 1)$ $W^{[2]} (1 \times 4)$ $a^{[1]} (4 \times 1)$ $a^{[1]} (1 \times 1)$

为了公式更加统一, 也可以把输入层看作是 $z^{[0]}$, 因此 x 即为 $a^{[0]}$.

3.4. 多个例子中的向量化.

在上-节中, 我们看一下输入 x , 了解到它是一个 3×1 的矩阵, 很明显这是一个单一样本的输入. 在前章节中, 对于单个节点 logistic 回归, 我们已做到了将多个 (m) 个样本封装到一个向量里面 $(n \times m)$. 这样做避免了显式的 for 循环, 提升了运算效率. 那么让我们看一下在多层神经网络中如何将输入 x 封装起来, 达到向量化的目的.

怎么向量化? 先回忆一下 logistic 回归中的情况.

对于输入 x^1, x^2, \dots, x^m , 我们期望的输出为 y^1, y^2, \dots, y^m .

同时由于 $x^{(i)}$ 是 $n \times 1$ 矩阵, 很明显

我们不可能将 $x^{(1-m)}$ 封装成下面的样子

$$\begin{bmatrix} x^1 \\ x^1 \\ \vdots \\ x^m \end{bmatrix}$$

如果这样封装的话
则令 $x^1 - x^m$ 就成为一个很麻烦的事情.

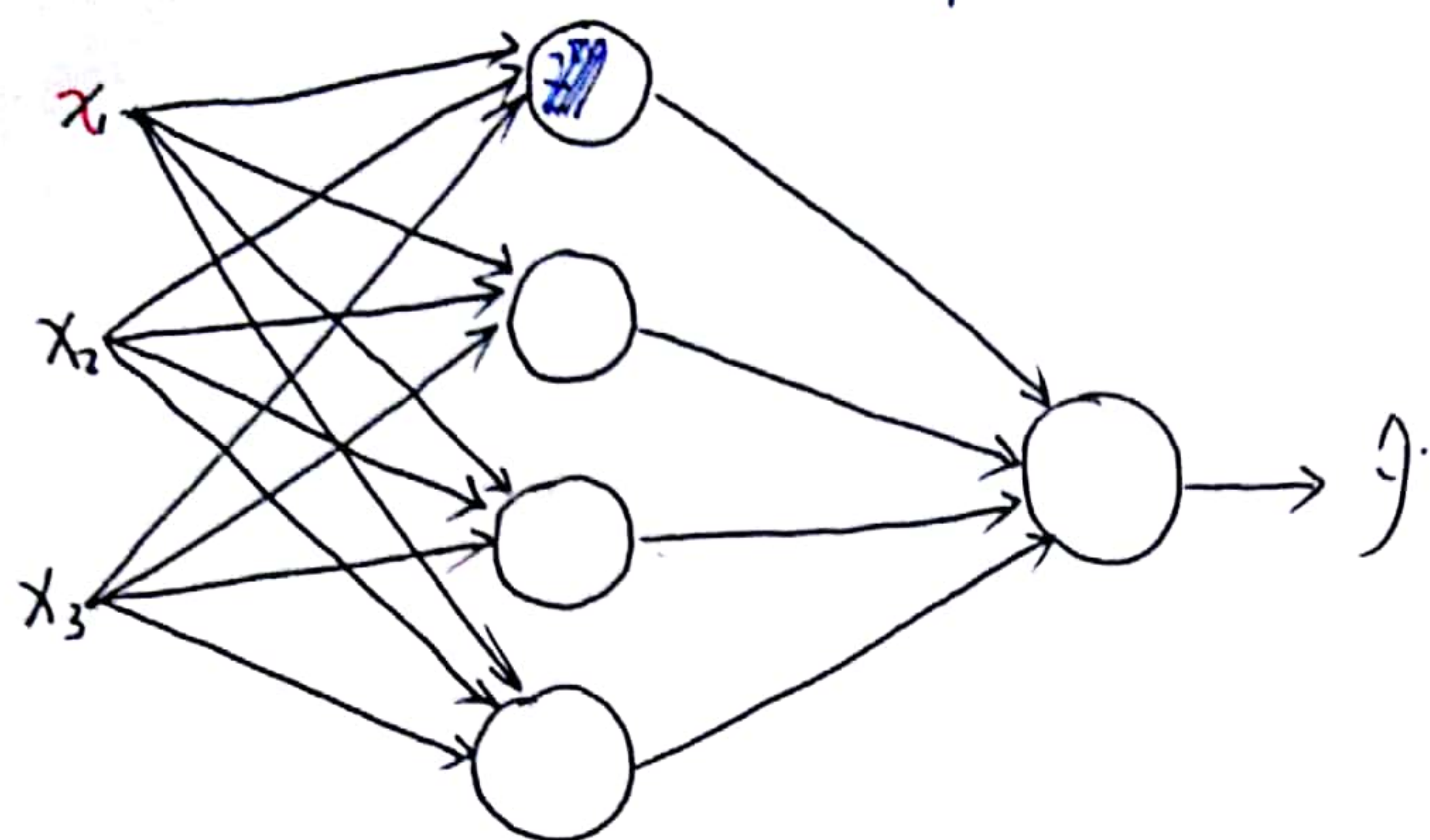
因此, 理想的封装方式应该和 logistic 回归中的封装方式一样.

$$X = \begin{bmatrix} \vdots & \vdots & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

让我们看一下, 这样封装后.

我们上面的两层神经网络的计算结果是什么?

在以前,把之前的神经网络计算总结如下



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

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

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

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

之后让我们用输入 x 来进行计算, 先计算第一层.

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

$$\begin{bmatrix} W_1^{[1]T} \\ W_2^{[1]T} \\ W_3^{[1]T} \\ W_4^{[1]T} \end{bmatrix} \cdot \begin{bmatrix} x \end{bmatrix} + b^{[1]} \Rightarrow \begin{bmatrix} W_1^{[1]T}x \\ W_2^{[1]T}x \\ W_3^{[1]T}x \\ W_4^{[1]T}x \end{bmatrix} + b^{[1]}$$

我们来研究 $W_i^{[1]T}x$ 的结果是什么?

$$W_i^{[1]T}x \Rightarrow W_i^{[1]T} \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix} \Rightarrow \begin{bmatrix} W_i^{[1]T}x^{(1)} & W_i^{[1]T}x^{(2)} & \dots & W_i^{[1]T}x^{(m)} \end{bmatrix}$$

其中 $W_i^{[1]T}x^{(n)}$ 就是第一层网络的第 i 个节点对输入第 n 个样样的输入 $z_i^{[1](n)}$ 是此. $W_i^{[1]T}x \Rightarrow \begin{bmatrix} z_i^{1} & z_i^{[1](2)} & \dots & z_i^{[1](m)} \end{bmatrix}$

因此可知

$$W^{[1]}x \Rightarrow \begin{bmatrix} z_1^{1} & z_1^{[1](2)} & \dots & z_1^{[1](m)} \\ z_2^{1} & z_2^{[1](2)} & \dots & z_2^{[1](m)} \\ z_3^{1} & z_3^{[1](2)} & \dots & z_3^{[1](m)} \\ z_4^{1} & z_4^{[1](2)} & \dots & z_4^{[1](m)} \end{bmatrix}$$

其中, 红线部分 $\begin{bmatrix} z_1^{1} \\ z_1^{[1](2)} \\ z_1^{[1](3)} \\ z_1^{[1](m)} \end{bmatrix}$ 即为对第一个样样的输入即 z_1^{1}

因此，最终的输出为

$$W^{(1)} X = [z^{(1)(1)} \ z^{(1)(2)} \ \dots \ z^{(1)(m)}] = Z^{(1)} \quad (\text{大写})$$

因此最终的第一层输出为

$$\sigma(W^{(1)} X) = [a^{(1)(1)} \ a^{(1)(2)} \ \dots \ a^{(1)(m)}] = A^{(1)} \quad (\text{大写})$$

再来看第二层网络的输出(计算)

第二层计算公式如下，

$$Z^{(2)} = W^{(2)} \cdot A^{(1)} + b^{(2)} \quad \text{将 } a^{(1)} \text{ 换作 } A^{(1)}$$

$$\text{有 } W^{(2)} \cdot A^{(1)} \Rightarrow W^{(2)} \cdot [a^{(1)(1)} \ a^{(1)(2)} \ \dots \ a^{(1)(m)}]$$

$$\text{有 } \Rightarrow [W^{(2)} a^{(1)(1)} \ W^{(2)} a^{(1)(2)} \ \dots \ W^{(2)} a^{(1)(m)}]$$

最终再加 $b^{(2)}$ 根据上面的式子有

$$W^{(2)} A^{(1)} + b^{(2)} \Rightarrow [z^{(2)(1)} \ z^{(2)(2)} \ \dots \ z^{(2)(m)}] = Z^{(2)} \quad (\text{大写})$$

最终有

$$\begin{bmatrix} \vdots & \vdots & \dots & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \vdots & \vdots & & \vdots \end{bmatrix} \Rightarrow [a^{(1)(1)} \ a^{(1)(2)} \ \dots \ a^{(1)(m)}] \Rightarrow [\hat{y}^{(1)} \ \hat{y}^{(2)} \ \dots \ \hat{y}^{(m)}]$$

可以看到，和之前 logistic、回归中的输入输出形式是一样的，这也就印证我们的想法是正确的。

$$\text{即 } X = \begin{bmatrix} \vdots & \vdots & \dots & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \vdots & \vdots & & \vdots \end{bmatrix} \quad \text{为 } n \times m \text{ 的矩阵。}$$

因此,从上面的总结中,我们可以看到一个规律,对于各层的输出来说,

它的横向各项是对应每个节点的输出

而纵向各项是对应每个训练样本对各节点的输出

就拿第一层的A来说

$$A^{[0]} = \begin{pmatrix} a_1^{[0](1)} & a_1^{[0](2)} & \dots & a_1^{[0](m)} \\ a_2^{[0](1)} & a_2^{[0](2)} & \dots & a_2^{[0](m)} \\ a_3^{[0](1)} & a_3^{[0](2)} & \dots & a_3^{[0](m)} \\ a_4^{[0](1)} & a_4^{[0](2)} & \dots & a_4^{[0](m)} \end{pmatrix}$$

第一个节点对m个样本的输出

第一个样本4个节点的输出

最终我们的神经网络计算公式如下

$$Z^{[0]} = W^{[0]} X + b^{[0]}$$

$(4 \times m)$ (4×3) $(3 \times m)$ (4×1)

$$A^{[0]} = \sigma(Z^{[0]})$$

$(4 \times m)$ $(4 \times m)$

$$Z^{[1]} = W^{[1]} A^{[0]} + b^{[1]}$$

$(1 \times m)$ (1×4) $(4 \times m)$ (1×1)

$$A^{[1]} = \sigma(Z^{[1]})$$

$(1 \times m)$ $(1 \times m)$

$$Z^{[0]} = W^{[0]} A^{[0]} + b^{[0]}$$

$$A^{[0]} = \sigma(Z^{[0]})$$

$$Z^{[1]} = W^{[1]} A^{[0]} + b^{[1]}$$

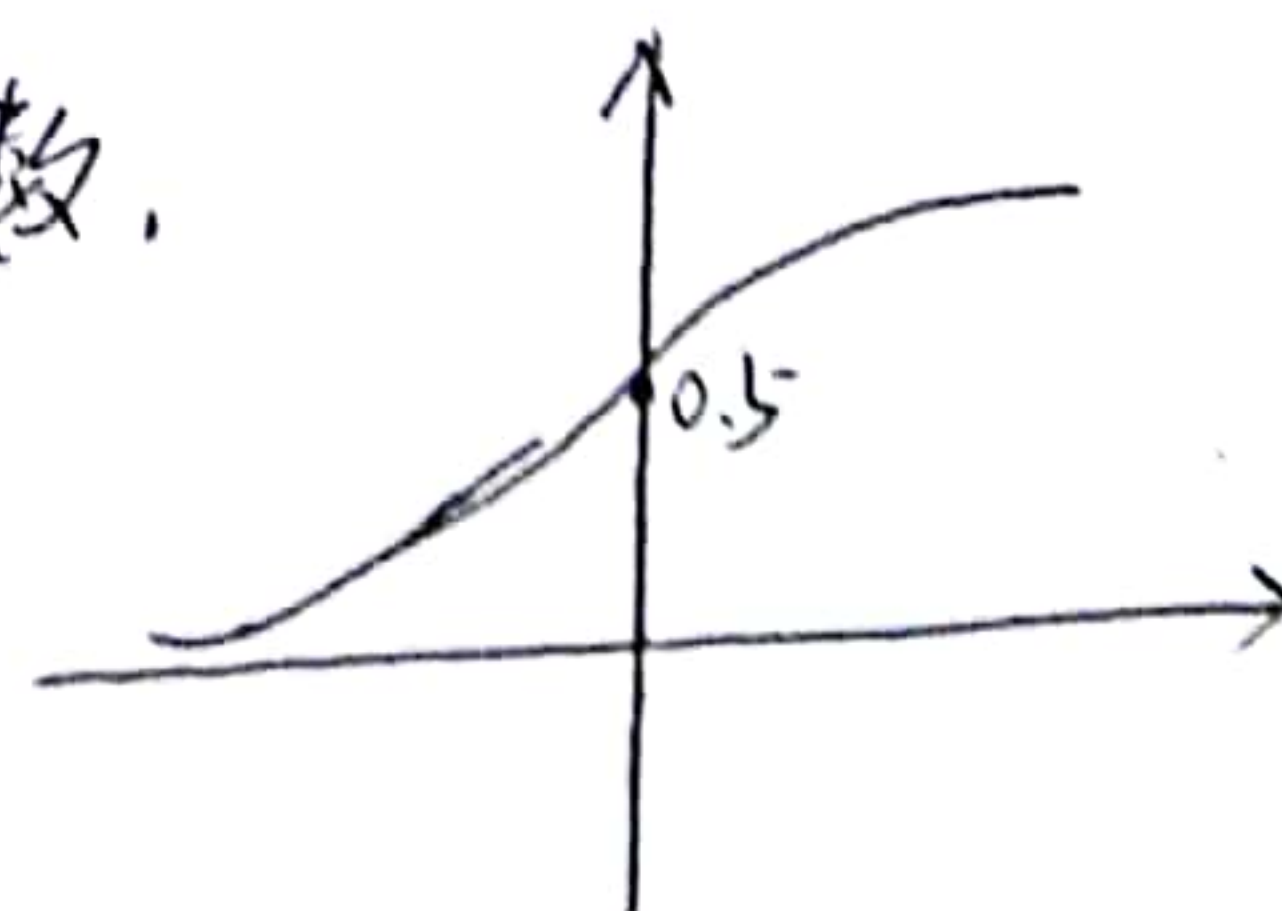
$$A^{[1]} = \sigma(Z^{[1]})$$

$$X = A^{[0]}$$

3.6 激活函数

在你搭建神经网络的时候,你可能会想使用什么激活函数,因为激活函数有很多种,到目前为止,我们使用的都是sigmoid函数。

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



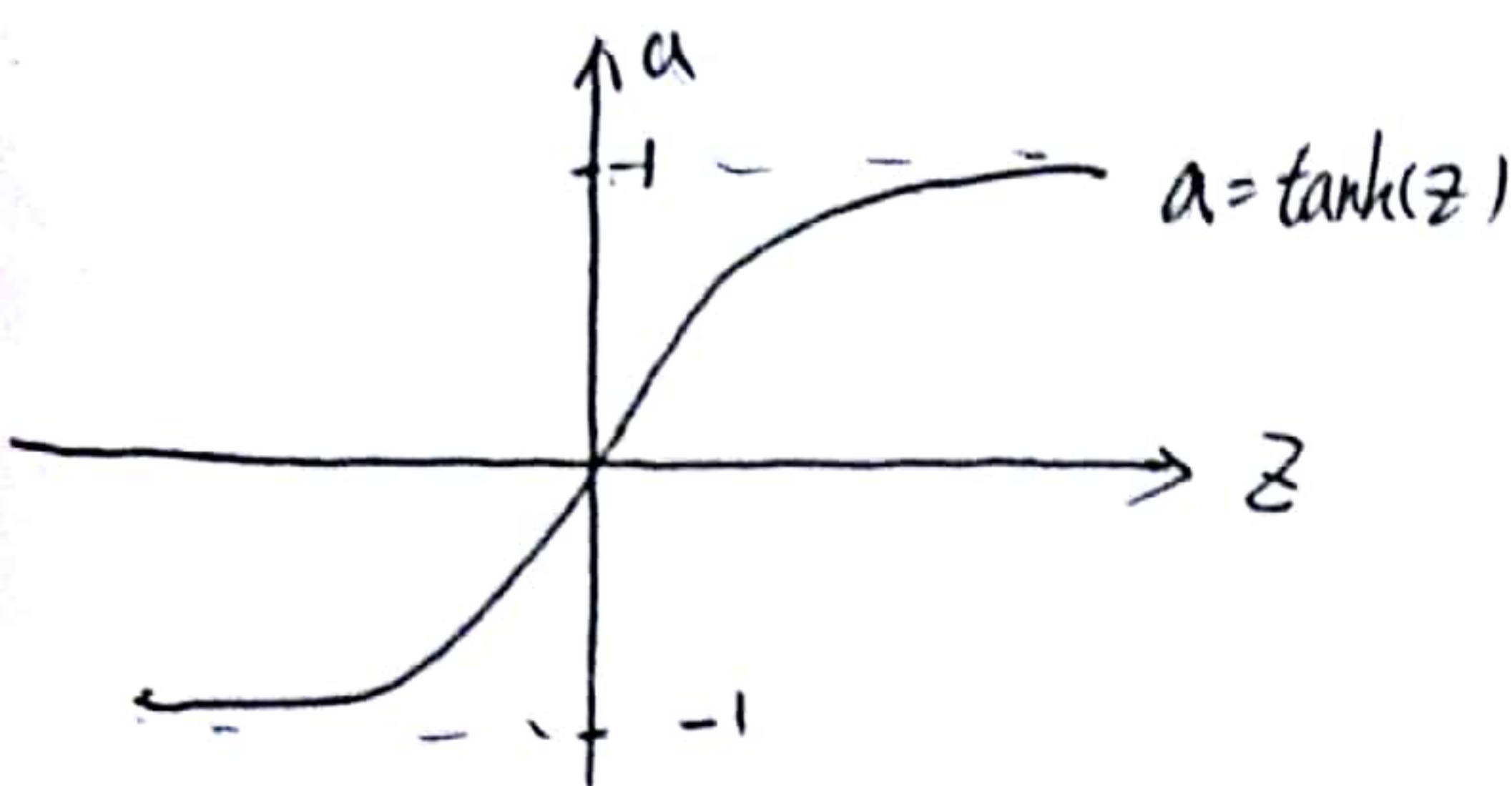
在上面提到的神经网络中,我们有两处用到了激活函数,即

$$A^{[0]} = \sigma(Z^{[0]})$$

$$A^{[1]} = \sigma(Z^{[1]})$$

在一般情况下我们用 $g(z)$ 来表示激活函数。

$g(z)$ 不一定是 Sigmoid 函数。更多情况下它可以是 tanh 函数，即双曲正切函数



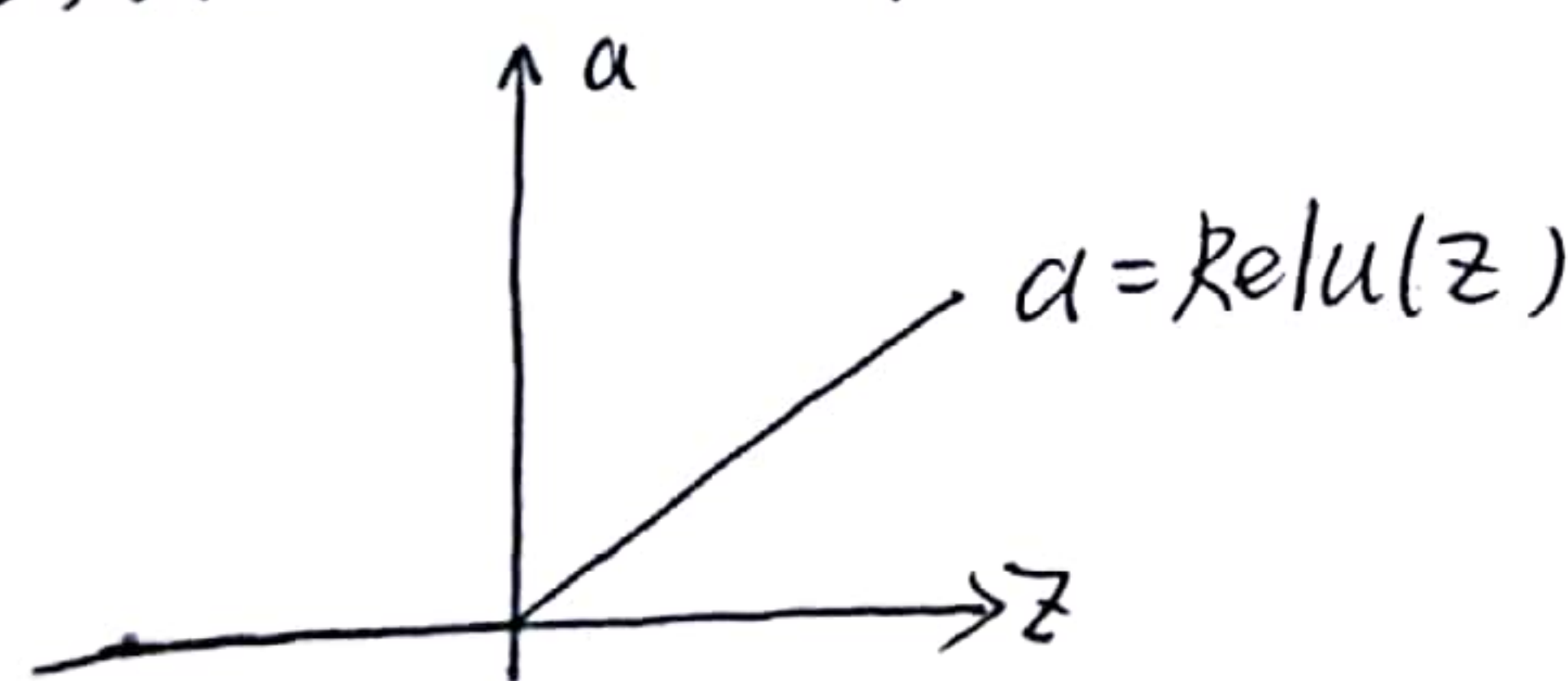
$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

事实证明，如果是隐藏单元，让 $g(z) = \tanh(z)$ 效果会更好。

因为此时激活函数在 $(-1, 1)$ 之间，它的平均值为零。使用 tanh 函数而不是 σ 函数有类似数据中心化的效果，使数据平均值更接近 0，而不是 0.5。实际上让下一层的学习更方便一点。

tanh 在很多场合都比 σ 函数更为优秀，但有一个例外，就是在输出层。如果输出层考虑二分类的话，很明显使用 σ 函数要更为适合（它可以表示为概率）。对于 tanh 函数来说，它和 Sigmoid 函数一样有一个共同的缺点，就是当我们使用梯度下降法来进行计算的时候，当 z 非常大时，效率会很低（因为斜率趋于零）。

在机器学习领域，最受欢迎的一个函数是 ReLU 函数

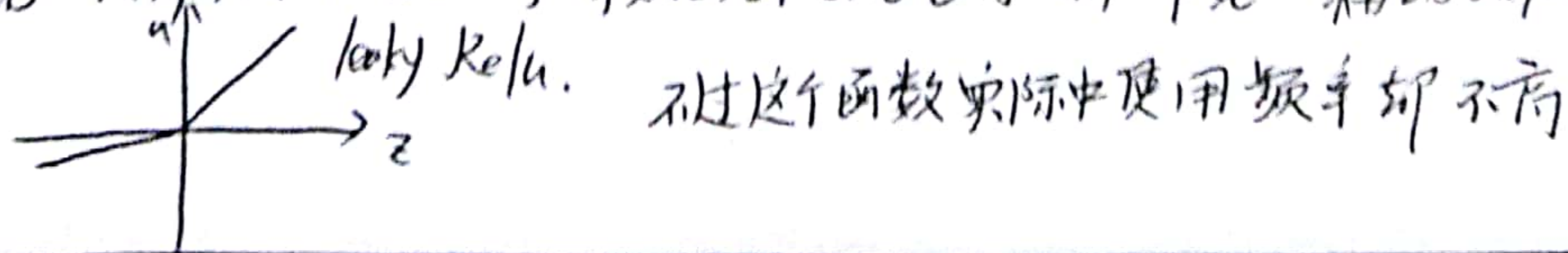


$$\text{Relu}(z) = \max(0, z)$$

在选择激活函数时有些经验法则：

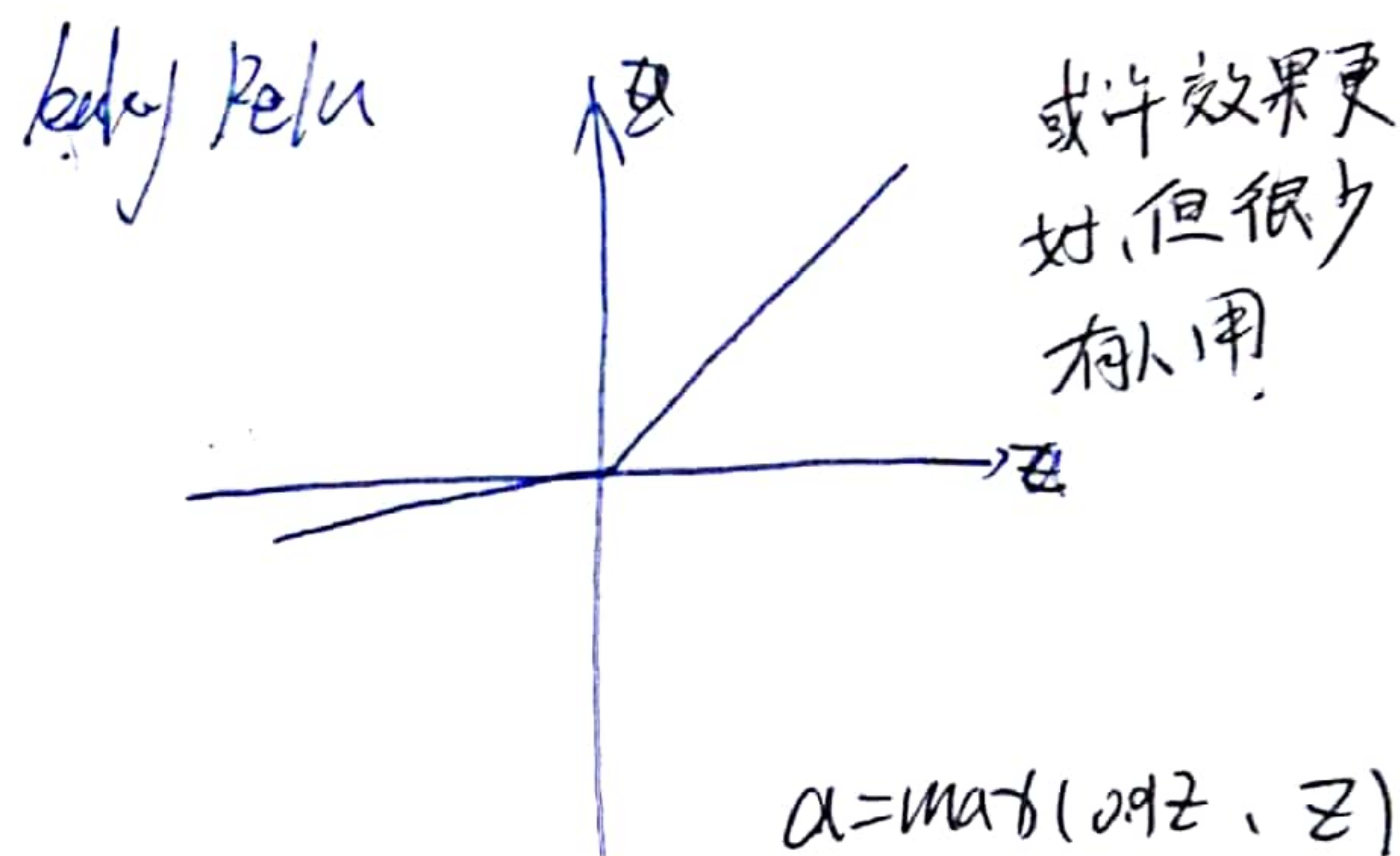
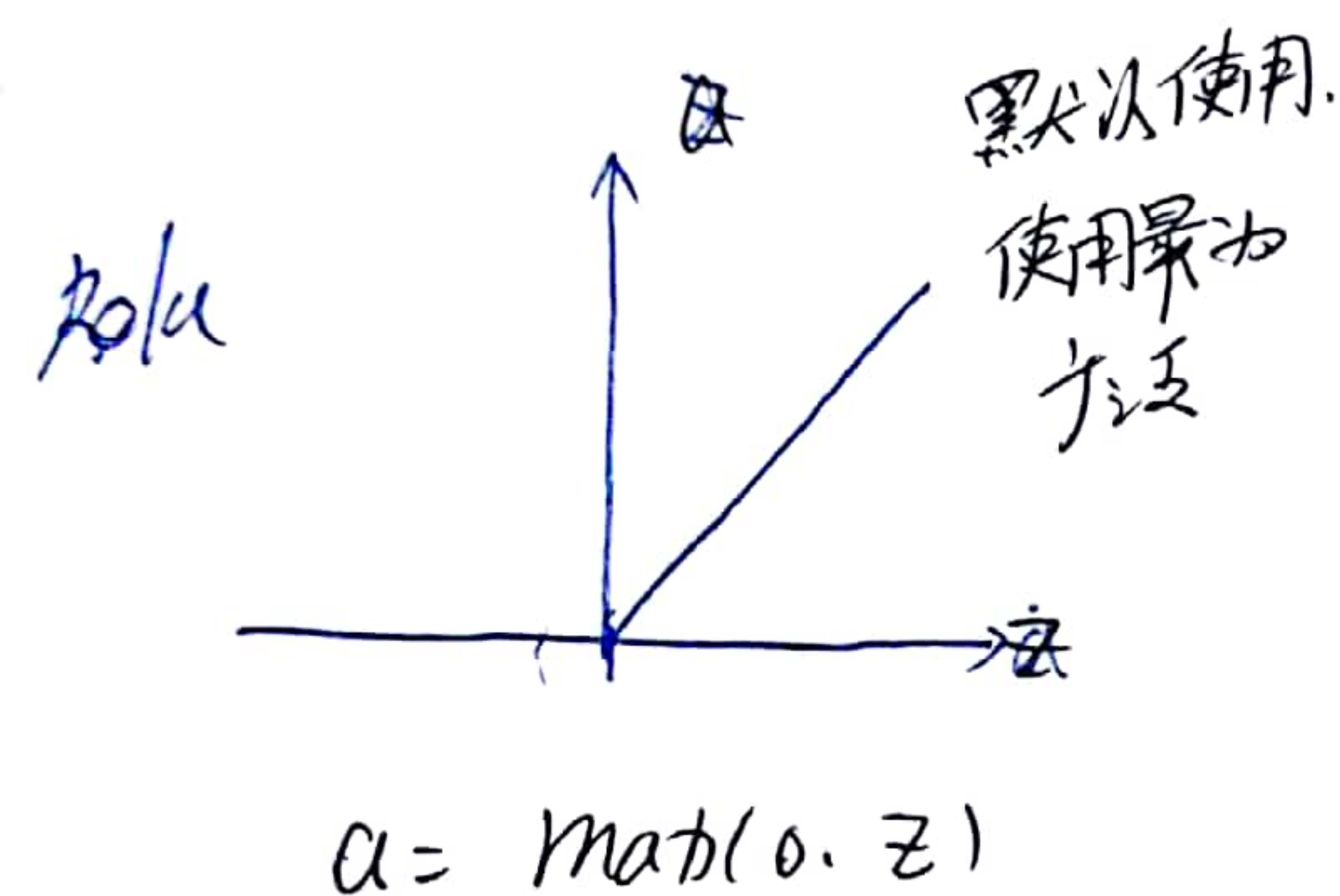
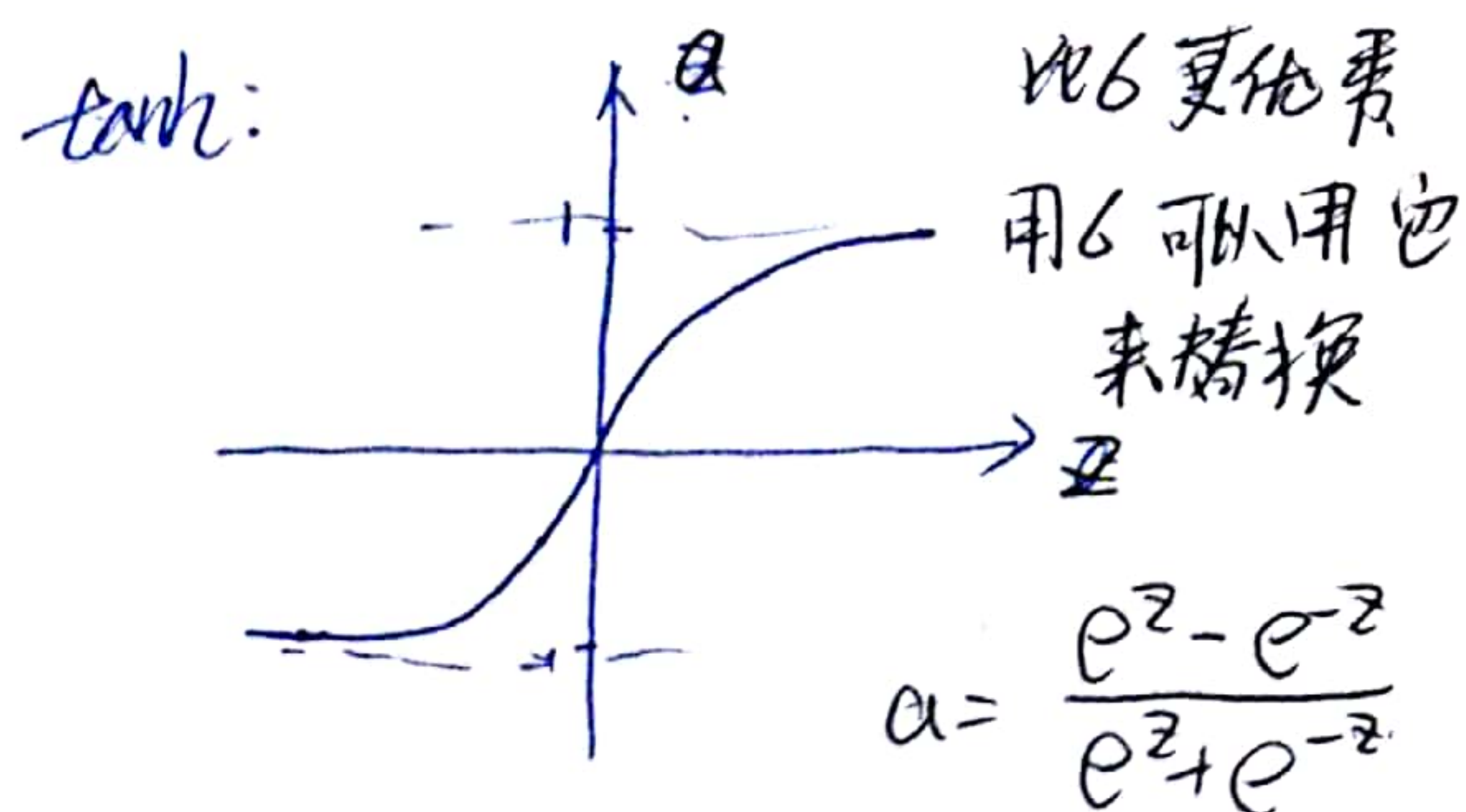
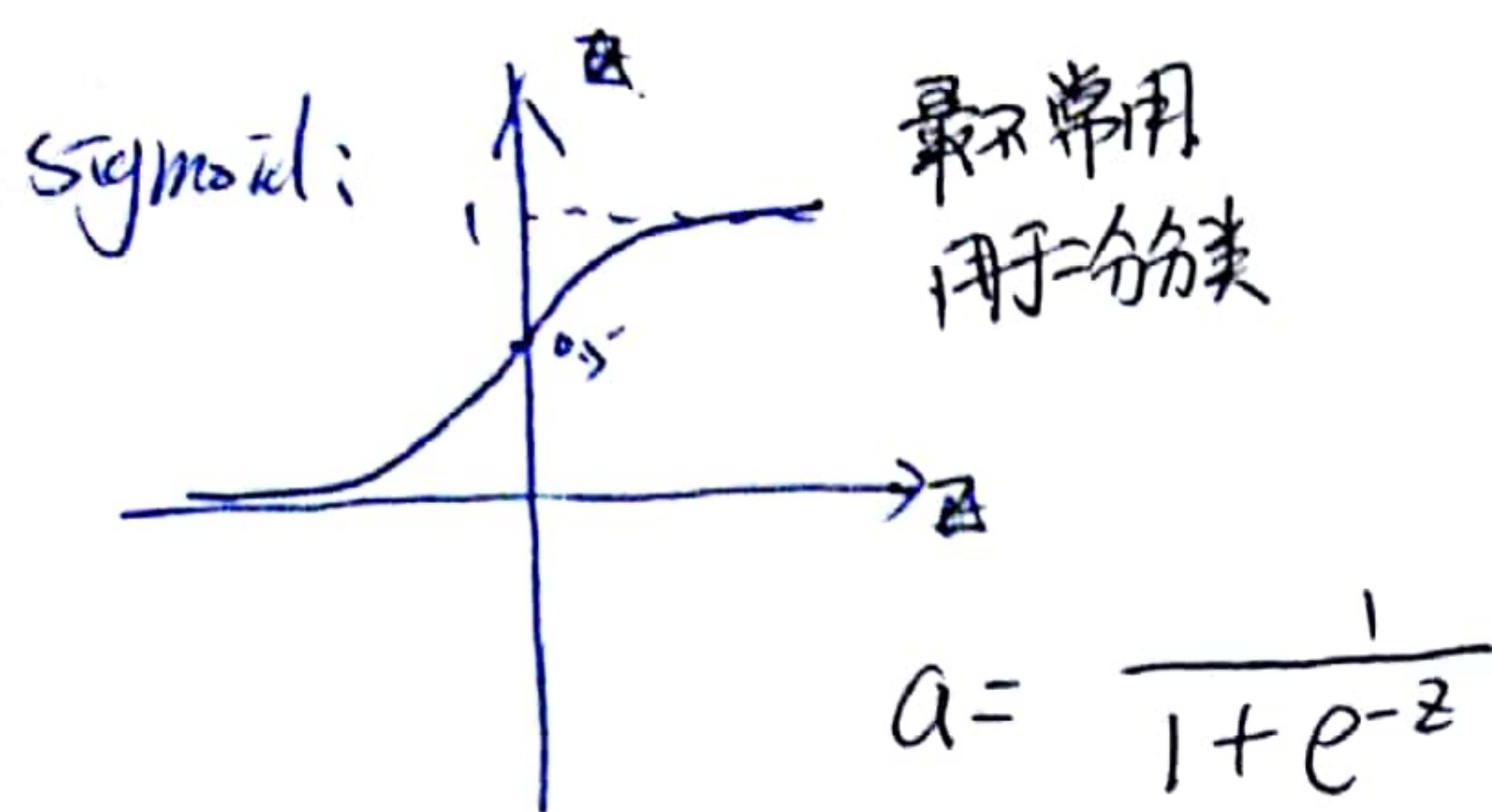
- ① 如果你的输出值是 0 和 1，那么 σ 函数很适合做输出层的激活函数
- ② 之后其他单元用 ReLU
- ③ 如果不确定隐藏层用什么激活函数，就用 ReLU 函数

ReLU 的一个缺点是在 $z < 0$ 时，导数为零，因此也有种带泄漏的 ReLU，即 leaky ReLU。



不过这个函数实际中使用频率却不高

下面来看这四个函数



神经网络一个特点就是可以有很多不同的选择, 比如.

①. 几个隐藏层

②. 层中有几个隐藏节点

③. 使用什么样的激活函数

④. 怎么初始化节点参数 (权重)

...

3.7. 为什么需要非线性修正单元 (激活函数)

让我们来看看神经网络的计算过程

$$z^{[1]} = W^{[0]} \cdot x + b^{[0]}$$

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

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

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

那么, 如果计算中的激活函数是线性的话会怎么样. 让我们以最简单的线性函数为例

$$\text{即 } g(z) = z$$

这个函数也被称为恒等激活函数.

让我们来看看会发生什么?

$$z^{[2]} = W^{[2]} \cdot x + b^{[2]}$$

$$a^{[2]} = z^{[2]}$$

$$z^{[3]} = W^{[3]} a^{[2]} + b^{[3]}$$

$$a^{[3]} = z^{[3]}$$

那么有 $a^{[2]} = W^{[2]} x + b^{[2]}$

因为 $a^{[2]} = z^{[2]}$ 则有

$$a^{[2]} = W^{[2]} \cdot (W^{[1]} x + b^{[1]}) + b^{[2]}$$

$$a^{[2]} = (W^{[2]} W^{[1]}) x + (W^{[2]} b^{[1]} + b^{[2]})$$

我们可以看到上式中 $a^{[2]}$ 就相当于 $a^{[2]} = W' x + b'$

那么这和只有一节点的 logistic 回归有什么区别呢?

很明显毫无区别。因此, 如果我们的激活函数都使用线性激活函数
的话, 那么隐藏层的存在就会变得毫无意义

因此, 只有引入非线性激活函数, 才能让神经网络富于变化
否则, 即使你有再多的隐藏层也是毫无意义。

只有一种情况, 就是在研究某些线性回归问题的时候可以使用。
就比如第一节中研究的房价问题。

但其实在这个问题中我们使用的函数是 ReLU 函数更为适合。

3.8 激活函数的导数

Sigmoid:

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} \text{导数: } \frac{d}{dz} g(z) &= \frac{1}{1 + e^{-z}} \cdot \left(1 - \frac{1}{1 + e^{-z}}\right) \\ &= g(z) \cdot (1 - g(z)) \end{aligned}$$

在神经网络中, $g(z) = a$ 。因此, 就

$$\text{因此 } g'(z) = a \cdot (1 - a)$$

tanh.

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\frac{dg(z)}{dz} = \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2} = 1 - (\tanh(z))^2$$

因此在神经网络中 $a = g(z)$ 因此对tanh来说 $g'(z) = 1 - a^2$

对于Relu函数不用过多解释.

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

leaky Relu. $g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$ $\text{leaky Relu}(z) = \max(0.01z, z)$

3.9. 神经网络的梯度下降法

在单隐层神经网络中, 我们有如下参数.

$$[W^{[0]}, b^{[0]}, W^{[1]}, b^{[1]}]$$

因此, 要调整这些参数, 使用梯度下降法, 就需要知道它们的导数

我们也用 n 来表示输入参数的维度.

例如 $n^{[0]}$ 代表输入 x 的维度. $n^{[1]}$ 代表第一层的输出维度

$n^{[2]}$ 代表第二层的输出维度

因此, 对于第一层来说.

$$W^{[0]} \text{ 的维度为 } (n^{[0]}, n^{[1]}) \quad b \text{ 为 } (n^{[1]}, 1)$$

对于第二层来说

$$W^{[1]} \text{ 的维度为 } (n^{[1]}, n^{[2]}) \quad b \text{ 为 } (n^{[2]}, 1)$$

在只有一个隐藏层的单元来说.

$$n^{[0]} = 1$$

要使用梯度下降法, 就需要一个成本函数. 这里我们假定在做二分类

因此 cost function $J(W^{[0]}, W^{[1]}, b^{[0]}, b^{[1]}) = \frac{1}{m} \sum_{i=1}^m \ell(y_i, \hat{y}_i)$ 很明显 $\hat{y} = a^{[2]}$

27

因此, 根据 = 命名类来看下面的计算

① 计算 $\hat{y} \Rightarrow (\hat{y}^{(1)} \hat{y}^{(2)} \dots \hat{y}^{(m)})$

② 计算第一层参数导数, 更新参数

$$dw^{[1]} = \frac{dJ}{dw^{[1]}} \quad db^{[1]} = \frac{dJ}{db^{[1]}}$$

$$dW^{[1]} = \alpha \cdot dw^{[1]}$$

$$b^{[1]} -= \alpha \cdot db^{[1]}$$

③ 计算第二层参数导数更新参数

$$dw^{[2]} = \frac{dJ}{dw^{[2]}} \quad db^{[2]} = \frac{dJ}{db^{[2]}}$$

$$W^{[2]} -= \alpha dw^{[2]}$$

$$b^{[2]} -= \alpha db^{[2]}$$

$$\left\{ \begin{array}{l} z^{[1]} = W^{[0]} X + b^{[0]} \\ A^{[1]} = g(z^{[1]}) \\ z^{[2]} = W^{[1]} A^{[1]} + b^{[1]} \\ A^{[2]} = \sigma(z^{[2]}) \\ \uparrow \\ (\hat{y}^{(1)}, \hat{y}^{(2)} \dots \hat{y}^{(m)}) \end{array} \right.$$

④ 重复上面 ①-③ 的计算

这里面最重要的就是计算

$$dw^{[1]}, dw^{[2]}, db^{[1]}, db^{[2]}$$

下面直接给出结论.

$$dz^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims}=\text{True})$$

$$dz^{[1]} = \frac{1}{m} W^{[0]T} dz^{[2]} \odot g'(z^{[1]})$$

(n^[1], m) (n^[1], m)

(之前少写了个括号, 推导后发现结果不一致, 粗心坑人呀!)

$$dW^{[1]} = \frac{1}{m} dz^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims}=\text{True})$$

和 logistic 回归很像.
因为早就看输出层的话, 就是一个 logistic 模型

✍

$$[1]^{[2]}, n^{[2]}$$

$$[n^{\pm 1}, n^{\pm 2}] = [n^{\pm 2}, 1]$$

$$= -70$$

$$[12]z|c = [12]p.p$$

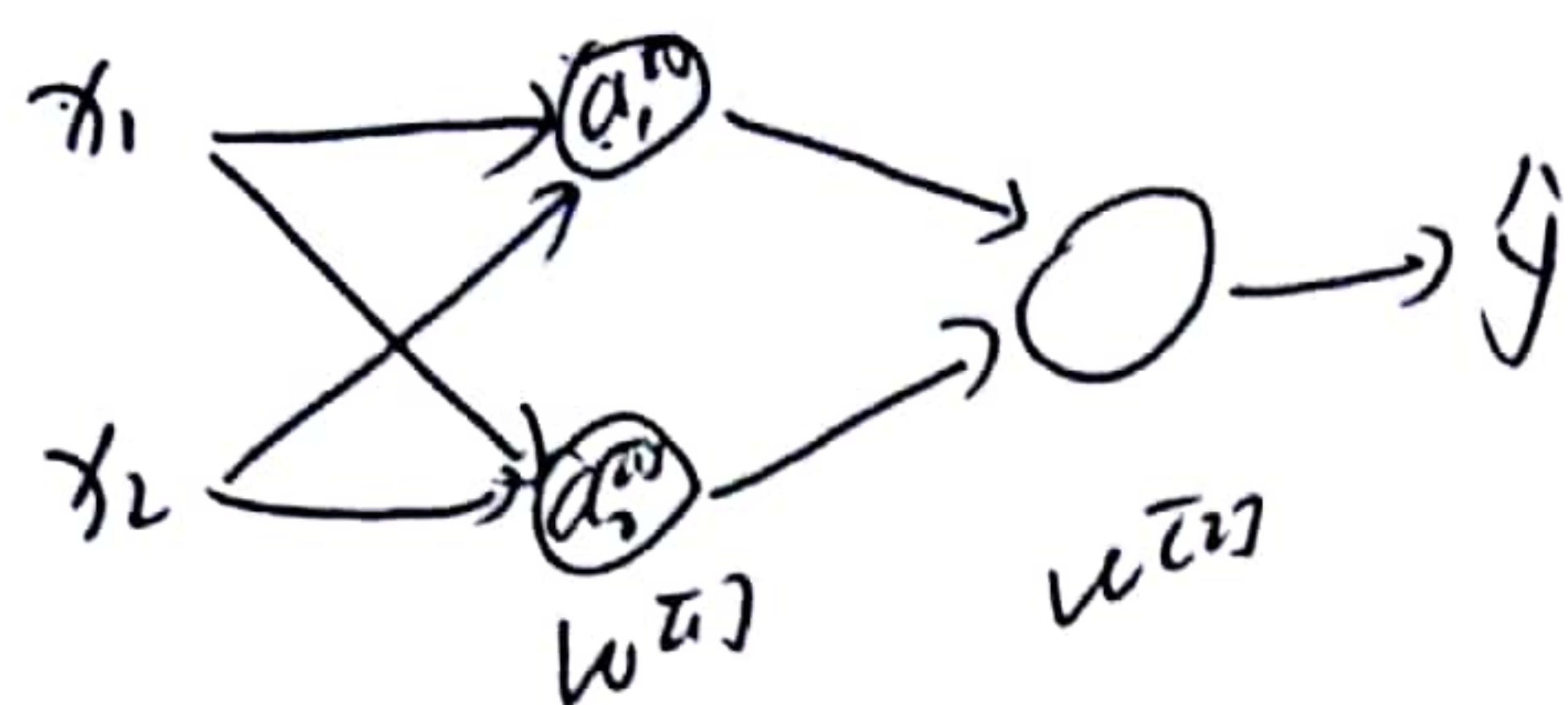
00-11-00

11. 57-1

30

3.11. 随机初始化.

为什么要随机初始化参数呢? 看下面的例子



如果 $w^{(1)}$ 中 $w_1^{(1)}, w_2^{(1)}$ 的参数相同 会出现什么情况呢?

结果是 $a_1^{(1)}$ 和 $a_2^{(1)}$ 计算结果相同. 最后导致它们计算出来的参数也相同. 当我们用梯度下降法来更新时它们的值依然会相同, 也就是说它们会永远相同, 这样就起不到作用了

因此需要随机初始化

如 $w^{(1)} = \text{np.random.randn}(2, 2) * 0.01$.

$w^{(2)} = \text{np.random.randn}(1, 2) * 0.01$

$b^{(1)} = \text{np.zeros}(2, 1)$.

$b^{(2)} = 0$

上面乘以 0.01 是因为我们应该让初始化参数尽量小.

这样计算出的 $z^{(1)}$ 也就相对较小, 如果使用 sigmoid 函数的

话, 会保证较高的学习效率