

Contents

1	1. shallow neural network	1
1.1	1.1. neural networks overview	1
1.2	1.2. neural network representation	1
1.3	1.3. computing a neural network's output	3
1.4	1.4. vectorizing across multiple examples	3
1.5	1.5. explanation for vectorized implementation	5
1.6	1.6. activation functions	5
1.7	1.7. why do you need non-linear activation functions	6
1.8	1.8. derivatives of activation functions	6
1.9	1.9. gradient descent for neural networks	6
1.10	1.10. backpropagation intuition	6
1.11	1.11. random initialization	11

contents

- 1. shallow neural network
- 1.1. neural networks overview
- 1.2. neural network representation
- 1.3. computing a neural network's output
- 1.4. vectorizing across multiple examples
- 1.5. explanation for vectorized implementation
- 1.6. activation functions
- 1.7. why do you need non-linear activation functions
- 1.8. derivatives of activation functions
- 1.9. gradient descent for neural networks
- 1.10. backpropagation intuition
- 1.11. random initialization

1 1. shallow neural network

1.1 1.1. neural networks overview

其中，每个神经元完成了 $z = w^T x + b$ 以及 $a = \sigma(z)$ 两个操作 (a 表示 activation)，每一层的数据用上标 [i] 表示。

1.2 1.2. neural network representation

图示是一个 2 层 nn (inputlayer 不算在内，有 hidden 和 output 两层)。

如果输入的 x 有 3 维，在 lr 中， $shape(w) = (1, 3)$, $shape(b) = (1, 1)$ 。

而在 nn 中， $shape(w^{[1]}) = (4, 3)$ 因为有 4 个神经元，输入是 3 维。同理 $shape(b^{[1]}) = (4, 1)$ 。

而 $shape(w^{[2]}) = (1, 4)$ ，因为只有 1 个神经元，输入是 3 维。同理 $shape(b^{[2]}) = (1, 1)$ 。

What is a Neural Network?

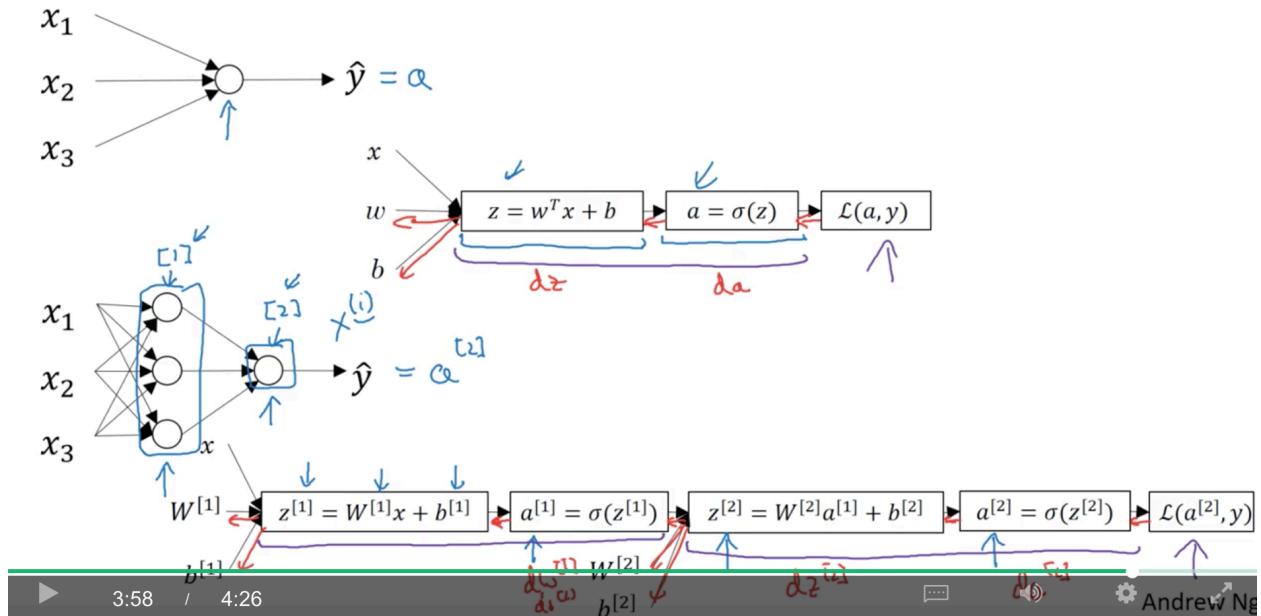


Figure 1: neural-networks-overview

Neural Network Representation

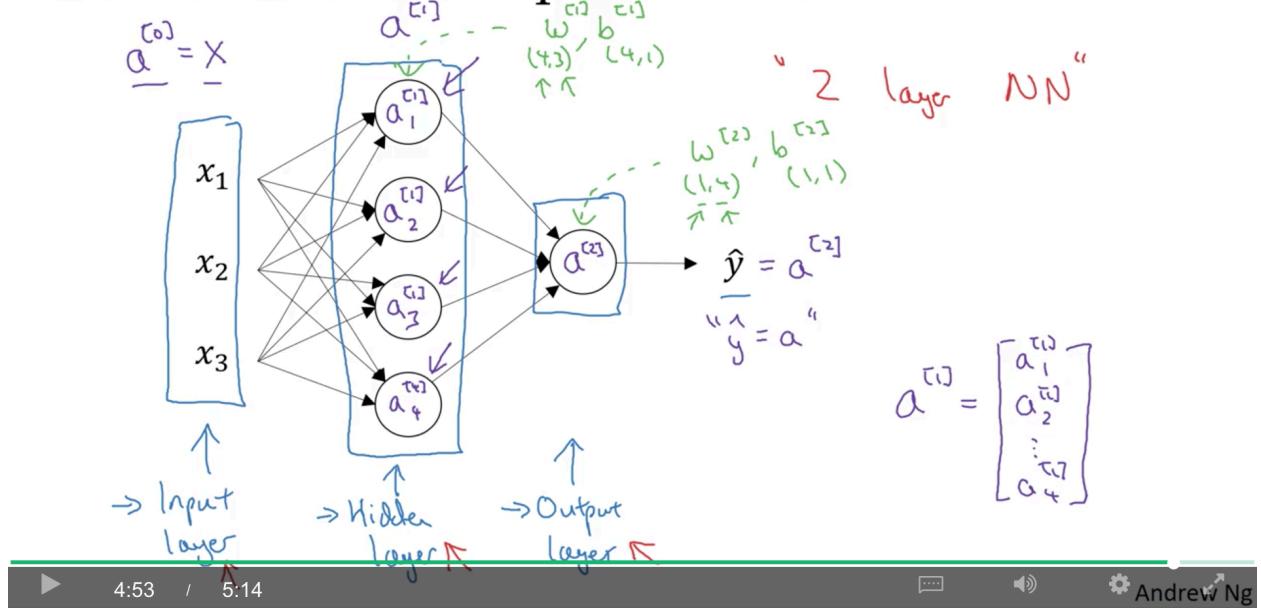
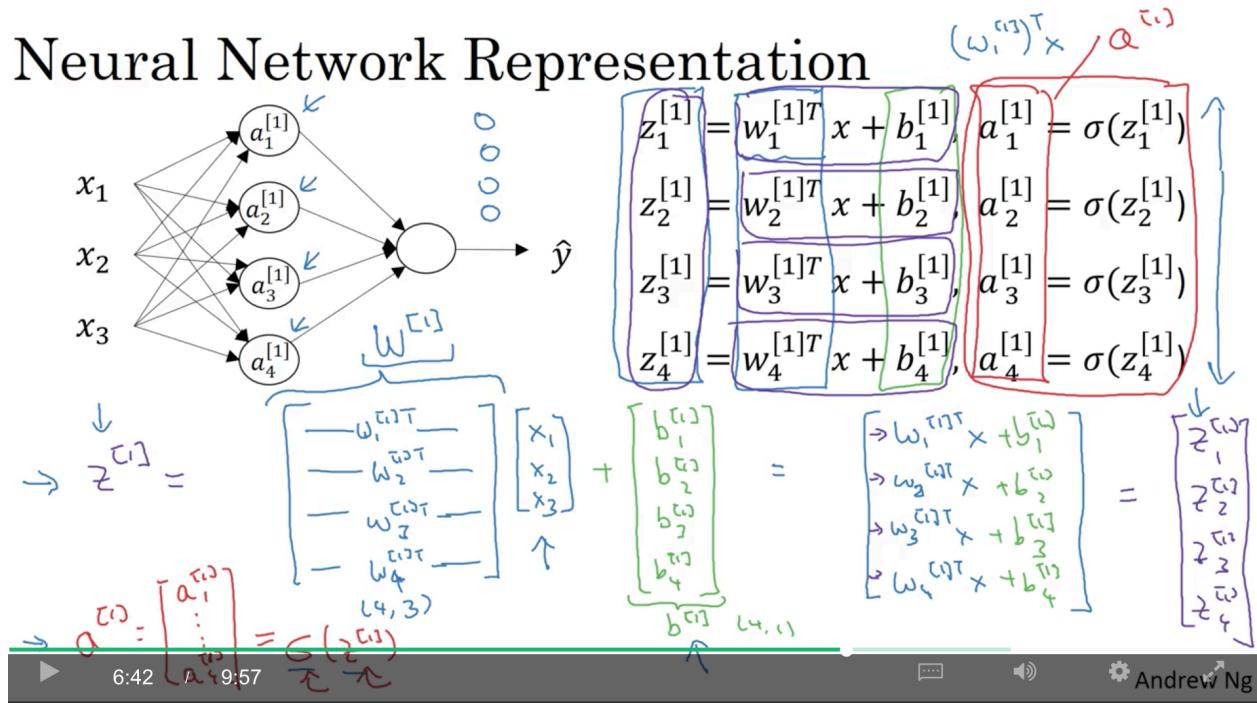


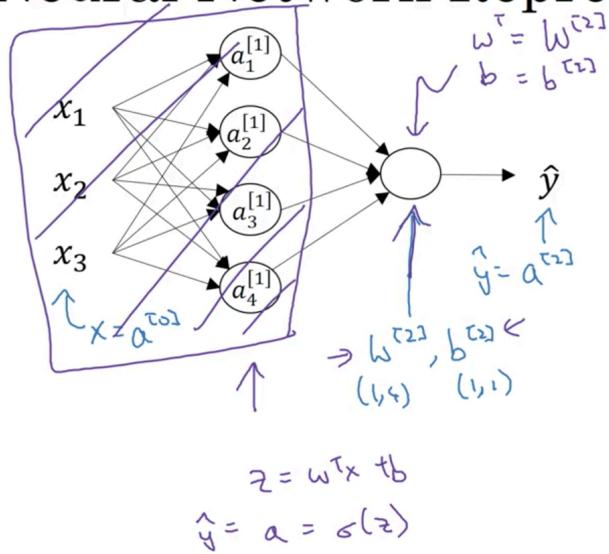
Figure 2: neural-network-representation

1.3 1.3. computing a neural network's output

Neural Network Representation



Neural Network Representation learning



Given input x :

$$\begin{aligned} &\rightarrow z^{[1]} = W^{[1]} x^{[0]} + b^{[1]} \\ &\rightarrow a^{[1]} = \sigma(z^{[1]}) \\ &\rightarrow z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \\ &\rightarrow a^{[2]} = \sigma(z^{[2]}) \end{aligned}$$

1.4 1.4. vectorizing across multiple examples

矩阵 X 纵向是 x 的维数（行数），横向是 training examples 的个数（列数）。

矩阵 Z 、 A 纵向是 hidden units 的个数（行数），横向是 training examples 的个数（列数）。

Vectorizing across multiple examples

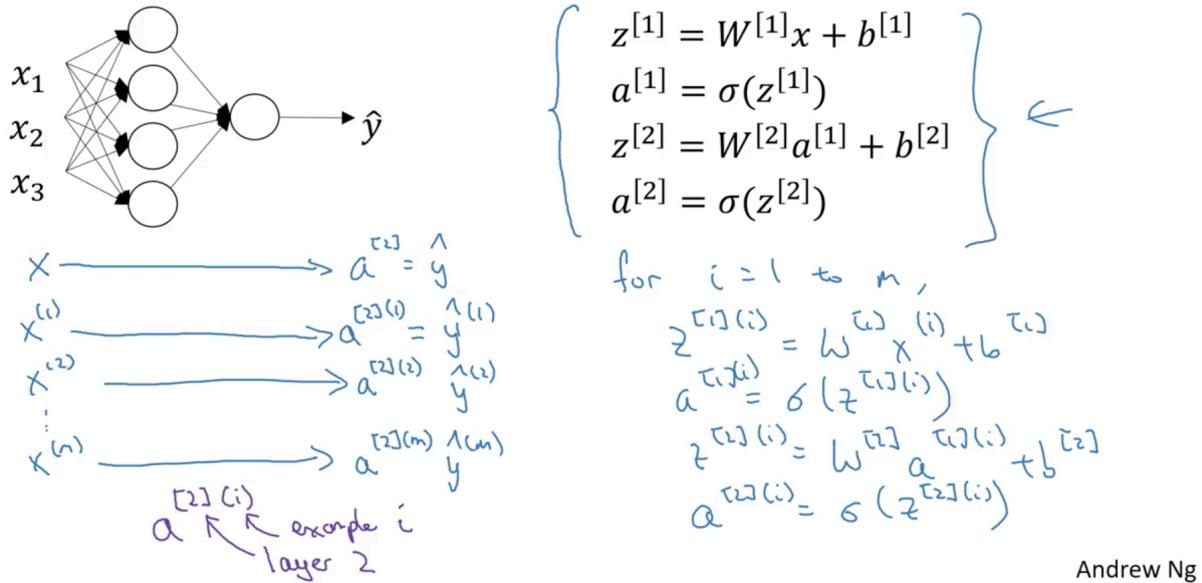


Figure 3: vectorizing-across-multiple-examples-for-loop

Vectorizing across multiple examples

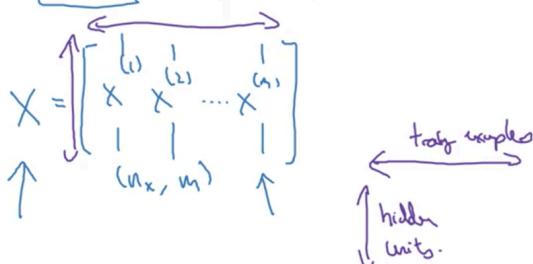
for $i = 1$ to m :

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

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

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

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



$$\begin{aligned} z^{[1]} &= W^{[1]}X + b^{[1]} \\ \Rightarrow A^{[1]} &= \sigma(z^{[1]}) \\ \Rightarrow z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ \Rightarrow A^{[2]} &= \sigma(z^{[2]}) \end{aligned}$$

$$z^{[1]} = \begin{bmatrix} z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \end{bmatrix}$$

hidden units

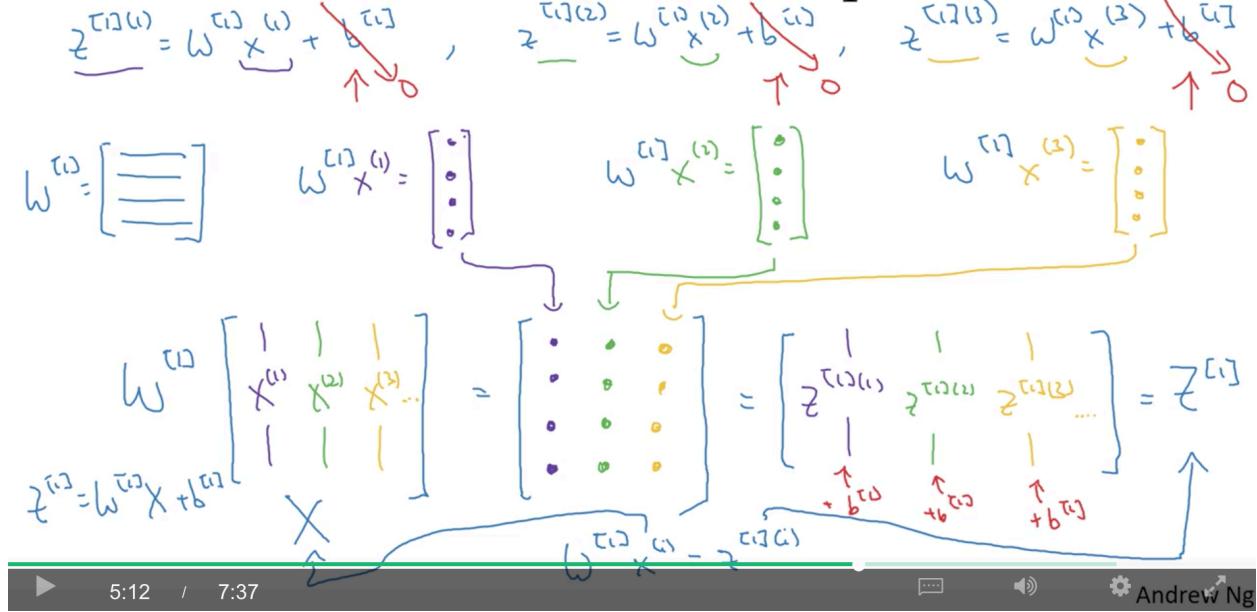
8:52 / 9:05

Andrew Ng

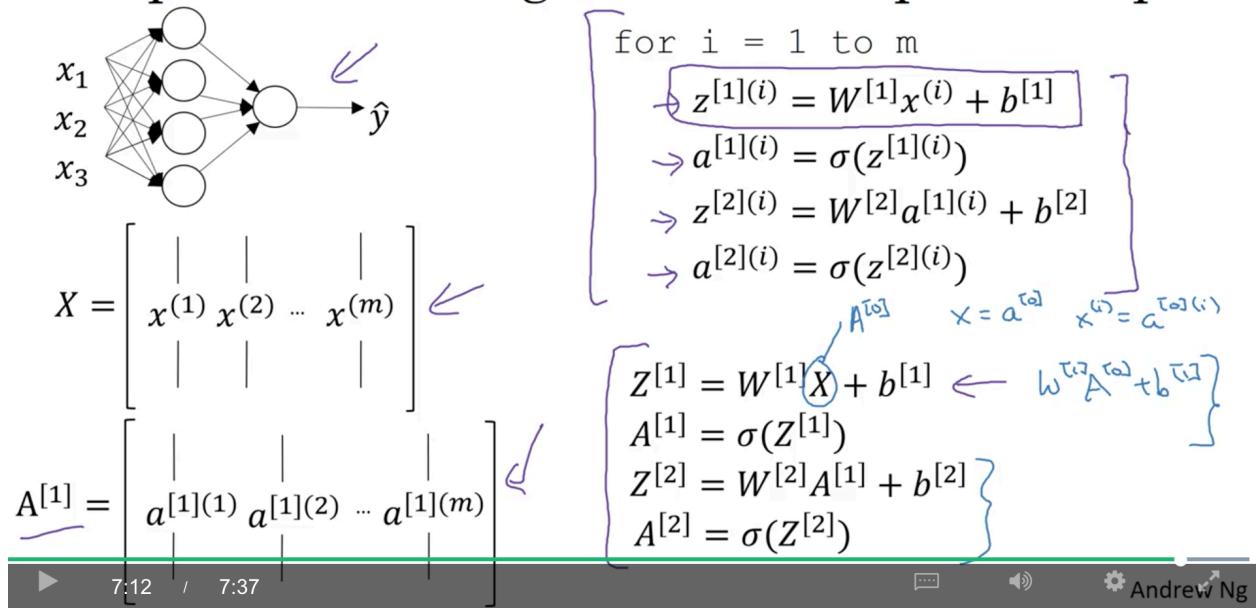
Figure 4: vectorizing-across-multiple-examples-vectorization

1.5 1.5. explanation for vectorized implementation

Justification for vectorized implementation



Recap of vectorizing across multiple examples



1.6 1.6. activation functions

一般来说, $tanh$ 效果比 $sigmoid$ 好, 因为均值是 0。但对于 outputlayer 而言, 一般 $y \in \{0, 1\}$, 所以希望 $0 \leq \hat{y} \leq 1$, 所以会用 $sigmoid$ 。

$ReLU(z) = \max(0, z)$ 比 $tanh$ 好, 因为当 $x \leq 0$ 时, 梯度为 0。而 $leakyReLU$ 在 $x \leq 0$ 时, 梯度是接近 0, 效果会好一点, 但在实践中还是 $ReLU$ 居多。当 $x > 0$ 时, 梯度和 $x \leq 0$ 差很多, 所以训练速度会加快。

Activation functions

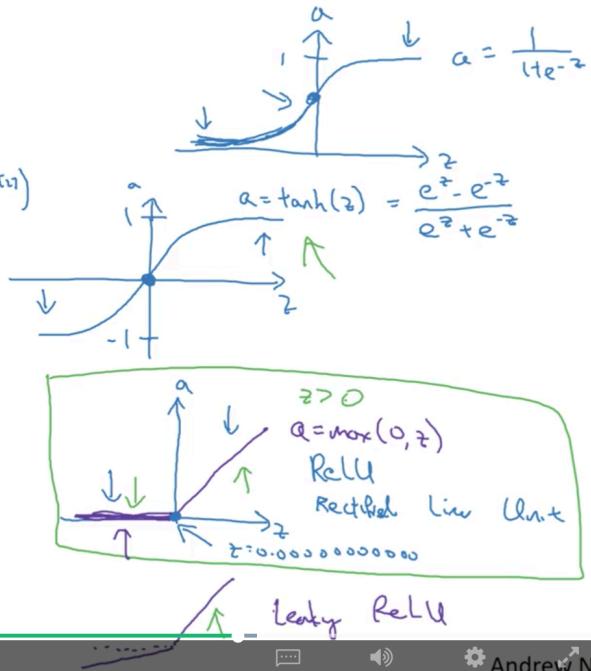
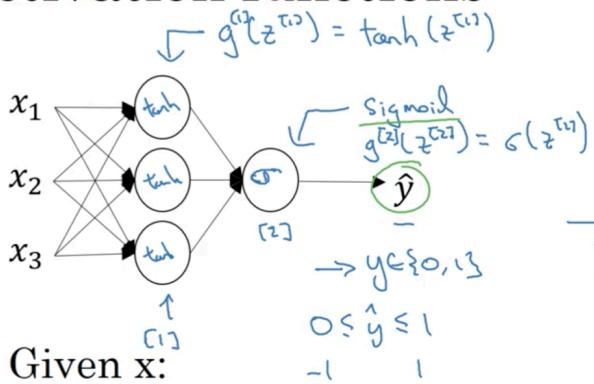


Figure 5: activation-functions

1.7 why do you need non-linear activation functions

linear activation: 因为 $z = wx + b$, 所以激活函数 $g(z) = z = wx + b$ 就叫 linear activation function, 也叫 identity activation function.

不要在 hidden layer 用 linear activation functions, 因为多个 linear 嵌套, 实质上还是 linear.

例外, 当进行回归时, $y \in R$, 可以 hidden layer 用 ReLU, 但 output layer 用 linear activation.

1.8 derivatives of activation functions

sigmoid 的导数

tanh 的导数

relu 和 leaky relu 的导数 ($z=0$ 时不可导, 但在工程上, 直接归入 $z>0$ 部分)

1.9 gradient descent for neural networks

记住每个 W/b 的 shape!

其中的 keepdims=True 表示, 输出的 shape = $(n^{[2]}, 1)$ 而非 $(n^{[2]},)$ 另外求 dz 时, 两项之前是 element-wise product(np.multiply), 其中第二项就是对激活函数求在 $z^{[1]}$ 的导数

1.10 backpropagation intuition

先回顾一下 lr:

Activation function

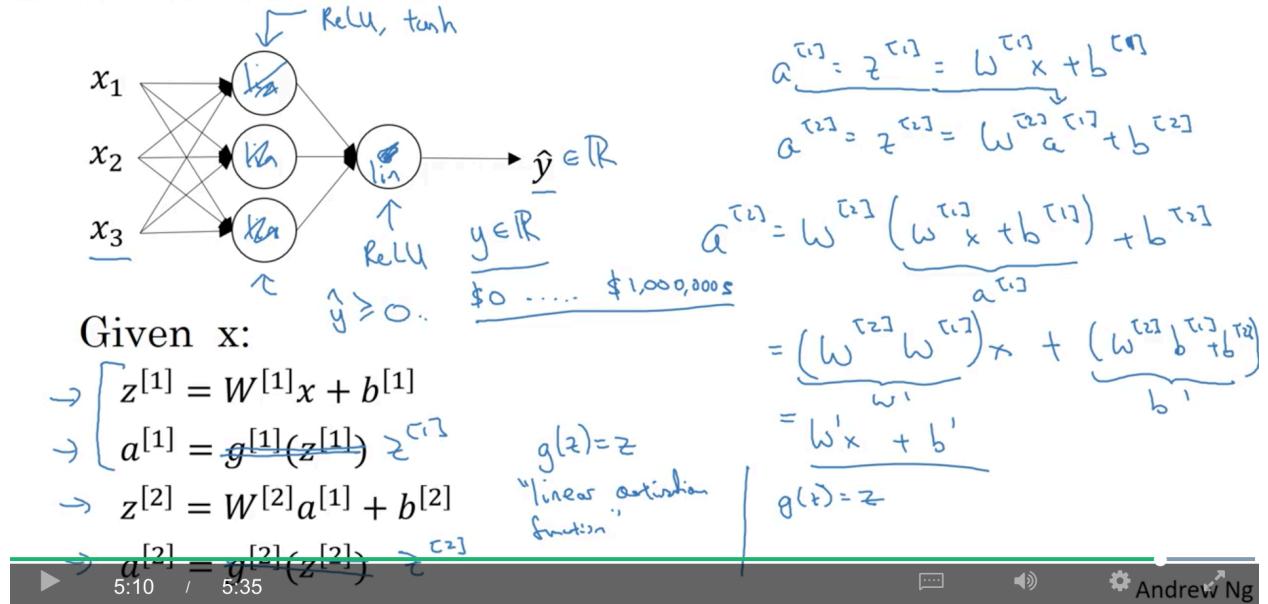


Figure 6: why-use-non-linear-functions.png

Sigmoid activation function

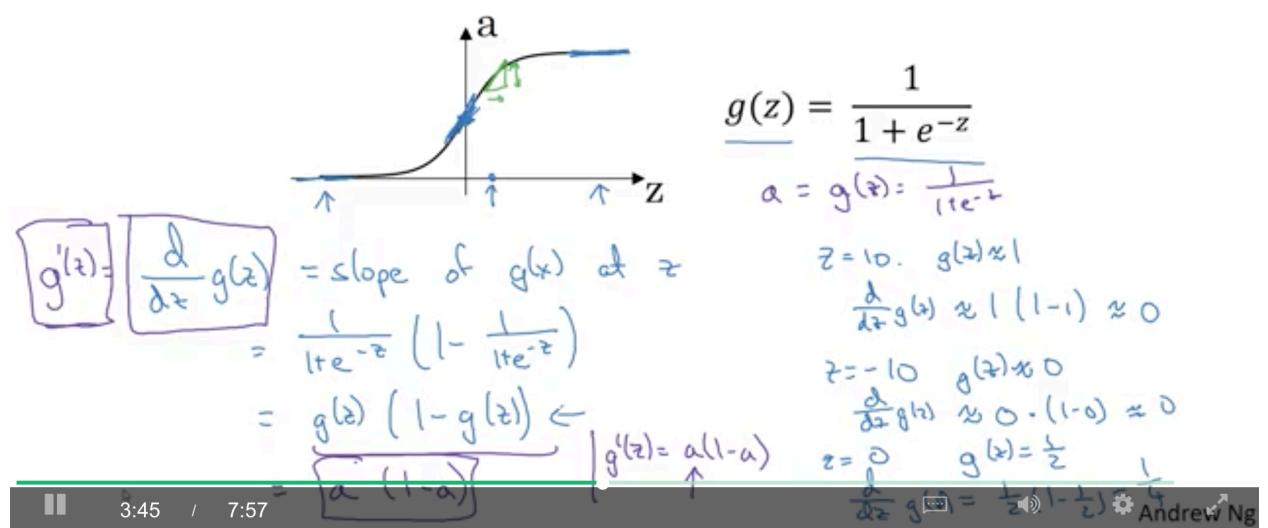


Figure 7: derivative-of-sigmoid.png

Tanh activation function

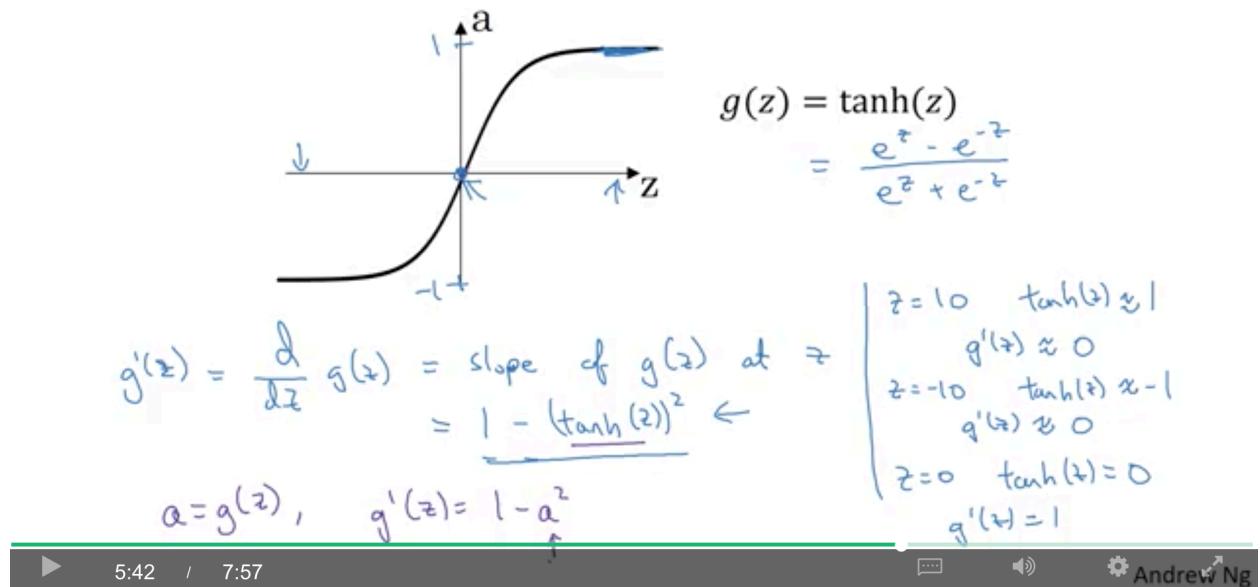


Figure 8: derivative-of-tanh.png

Tanh activation function

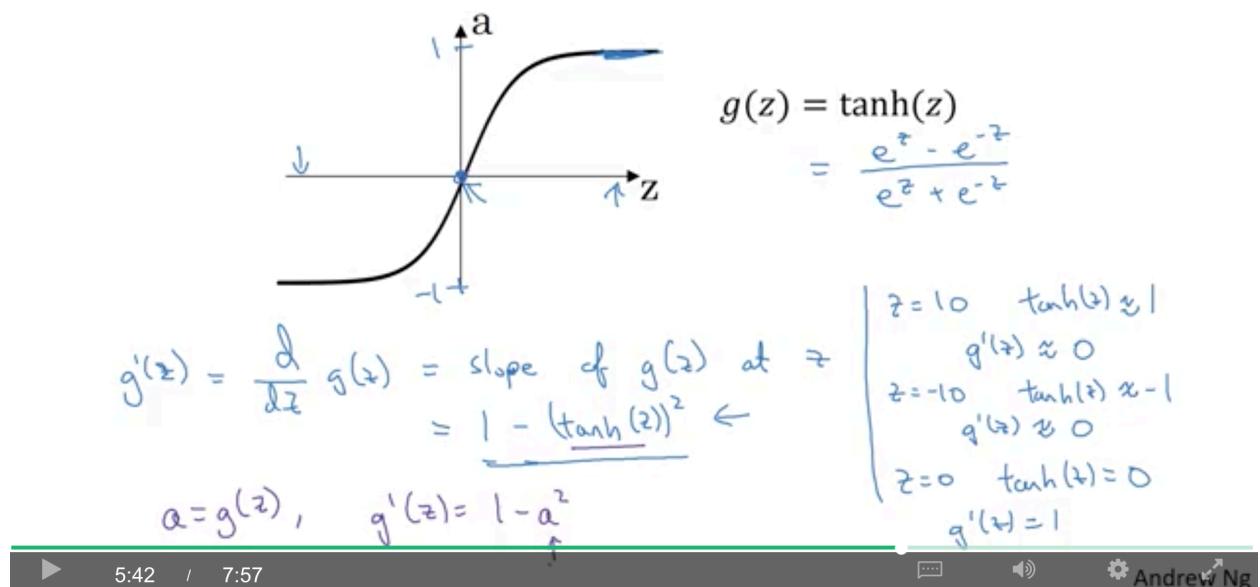


Figure 9: derivative-of-tanh.png

Gradient descent for neural networks

Parameters: $\underbrace{w^{[0]}, b^{[0]}}_{(n^{[0]}, n^{[0]})}, \underbrace{w^{[1]}, b^{[1]}}_{(n^{[1]}, 1)}, \dots, \underbrace{w^{[L]}, b^{[L]}}_{(n^{[L]}, 1)}$ $n_x = n^{[0]}, n^{[1]}, \dots, \underline{n^{[L]} = 1}$

Cost function: $J(w^{[0]}, b^{[0]}, w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$

Gradient descent:

→ Repeat {

- Compute predicts $(\hat{y}^{(i)}, i=1, \dots, m)$
- $\frac{\partial J}{\partial w^{[l]}} = \frac{\partial J}{\partial w^{[l]}} , \frac{\partial J}{\partial b^{[l]}} = \frac{\partial J}{\partial b^{[l]}} , \dots$
- $w^{[l]} := w^{[l]} - \alpha \frac{\partial J}{\partial w^{[l]}}$
- $b^{[l]} := b^{[l]} - \alpha \frac{\partial J}{\partial b^{[l]}}$

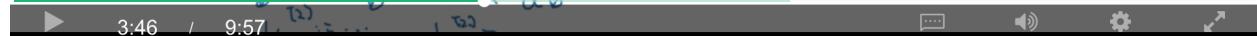


Figure 10: gradient-descent-for-neural-networks

Formulas for computing derivatives

Forward propagation:

$$\begin{aligned} z^{[0]} &= w^{[0]}X + b^{[0]} \\ A^{[0]} &= g^{[0]}(z^{[0]}) \\ z^{[1]} &= w^{[1]}A^{[0]} + b^{[1]} \\ A^{[1]} &= g^{[1]}(z^{[1]}) = \sigma(z^{[1]}) \end{aligned}$$

Back propagation:

$$\begin{aligned} d\hat{z}^{[2]} &= A^{[2]} - Y \\ d\hat{w}^{[2]} &= \frac{1}{m} d\hat{z}^{[2]} A^{[1]T} \\ d\hat{b}^{[2]} &= \frac{1}{m} \text{np.sum}(d\hat{z}^{[2]}, \text{axis}=1, \text{keepdims=True}) \\ d\hat{z}^{[1]} &= \underbrace{(w^{[1]T} d\hat{z}^{[2]})}_{(n^{[0]}, m)} \times \underbrace{g^{[1]'}(z^{[1]})}_{\text{element-wise product}} \\ dw^{[1]} &= \frac{1}{m} d\hat{z}^{[1]} X^T \\ db^{[1]} &= \frac{1}{m} \text{np.sum}(d\hat{z}^{[1]}, \text{axis}=1, \text{keepdims=True}) \end{aligned}$$



Figure 11: forward-propagation-and-back-propagation

Computing gradients

Logistic regression

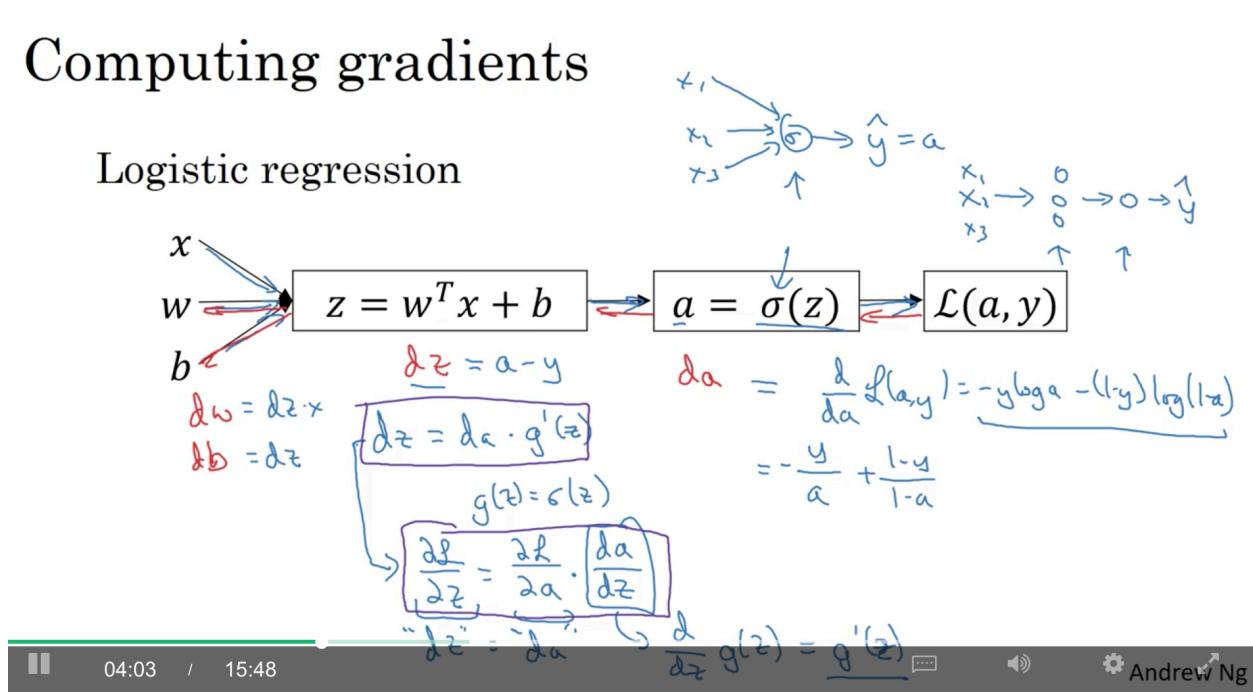


Figure 12: back-propagation-lr

Neural network gradients

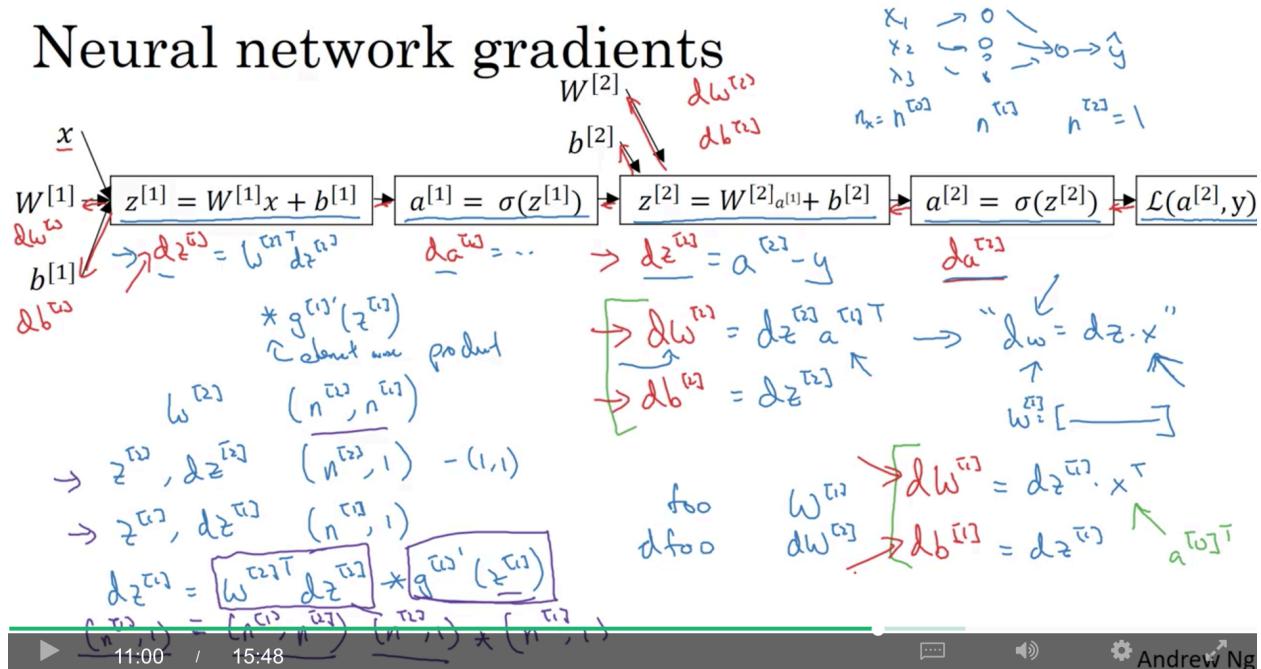


Figure 13: back-propagation-nn

然后看 nn:

扩展到 m 个 examples，并进行 vectorize:

Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

Vectorized implementation:

$$\begin{aligned} z^{[1]} &= w^{[2]} x + b^{[2]} \\ a^{[1]} &= g^{[1]}(z^{[1]}) \end{aligned}$$

$$Z^{[1]} = \begin{bmatrix} 1 & z^{[1]_{11}} & z^{[1]_{12}} & \dots & z^{[1]_{1m}} \\ 1 & | & | & \dots & | \\ 1 & z^{[1]_{m1}} & z^{[1]_{m2}} & \dots & z^{[1]_{mm}} \end{bmatrix}$$

$$\begin{aligned} Z^{[2]} &= w^{[1]} X + b^{[1]} \\ A^{[2]} &= g^{[2]}(Z^{[2]}). \end{aligned}$$

12:27 / 15:48

Andrew Ng

Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

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

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

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

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]'}(Z^{[1]})}_{\text{elementwise product}}$$

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

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

14:22 / 15:48

Andrew Ng

问：为何有的有 $1/m$ ，有的没有。。。

1.11 1.11. random initialization

lr 的训练，参数一般都初始化为 0。但 nn，如果初始化为 0，会发现算出来的 $a_1^{[1]} = a_2^{[1]}$, $dz_1^{[1]} = dz_2^{[1]}$ ，所以相当于每个神经元的 influence 是一样的 (symmetric)，所以 $dW^{[1]}$ 这个矩阵的每一行都相等。“Each neuron in the first hidden layer will perform the same

computation. So even after multiple iterations of gradient descent each neuron in the layer will be computing the same thing as other neurons.”

What happens if you initialize weights to zero?

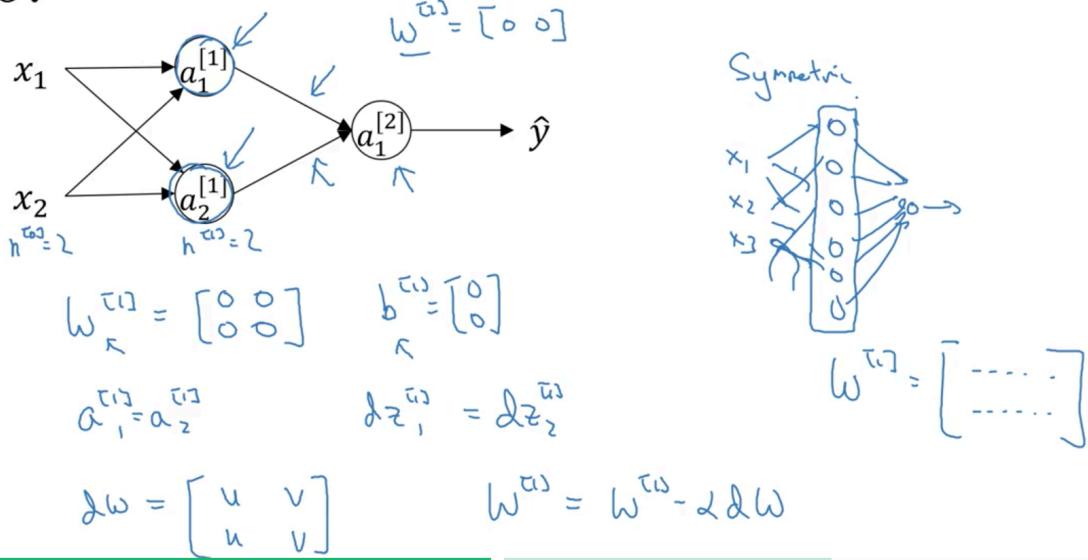


Figure 14: if-initialize-zero

解决：随机初始化只要 W 随机初始化了， b 其实影响不大，0 就可以了。如果用的是 sigmoid/tanh 的话，随机初始化时，尽量小，因为如果大的话，激活后会接近两端（无论是 + 无穷还是 - 无穷，梯度都接近 0），导致学习过程变得很慢。对于浅层的网络，0.01 就可以了，但如果是更深的网络，可能需要其他系数，这个选择过程，在后面的课程会讲。。。

Random initialization

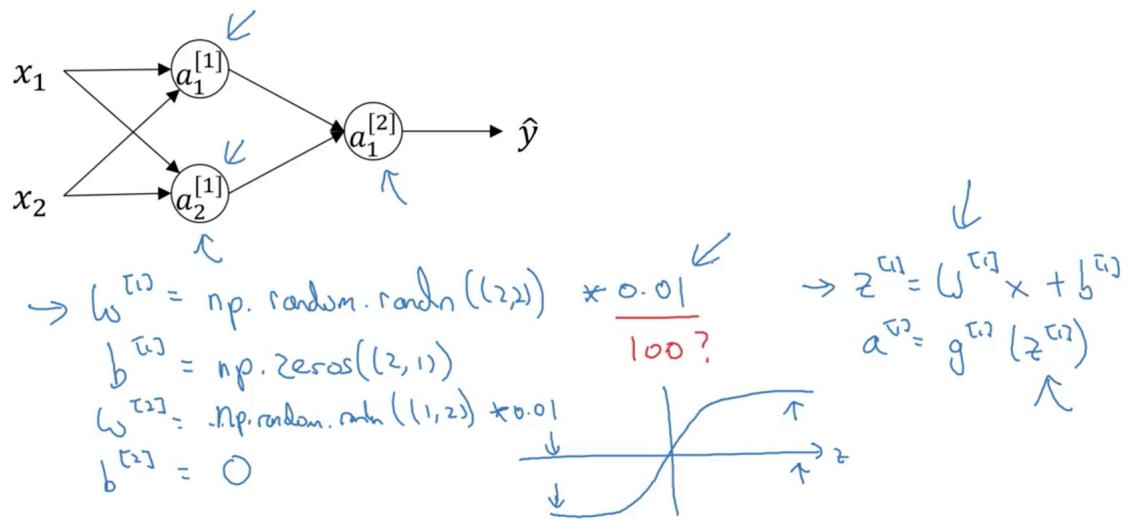


Figure 15: random-initialization