

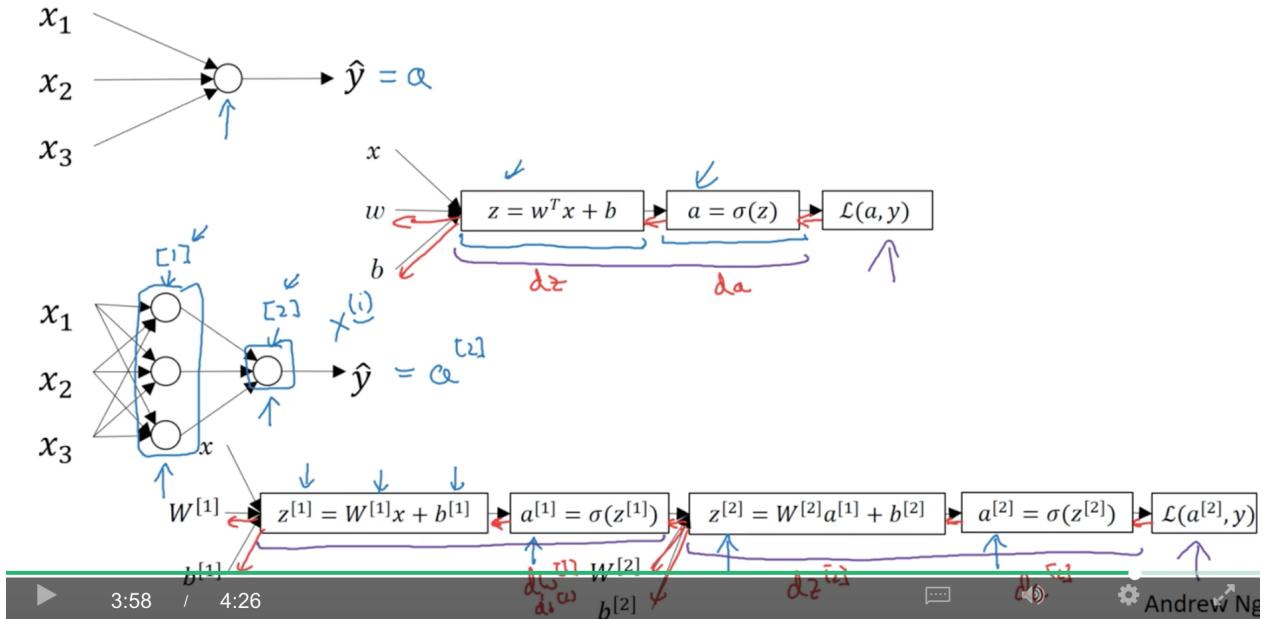
118 lines (70 sloc) 6.7 KB

contents

- neural networks overview
- neural network representation
- computing a neural network's output
- vectorizing across multiple examples
- explanation for vectorized implementation
- activation functions
- why do you need non-linear activation functions
- derivatives of activation functions
- gradient descent for neural networks
- backpropagation intuition
- random initialization

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。

What is a Neural Network?**neural network representation**

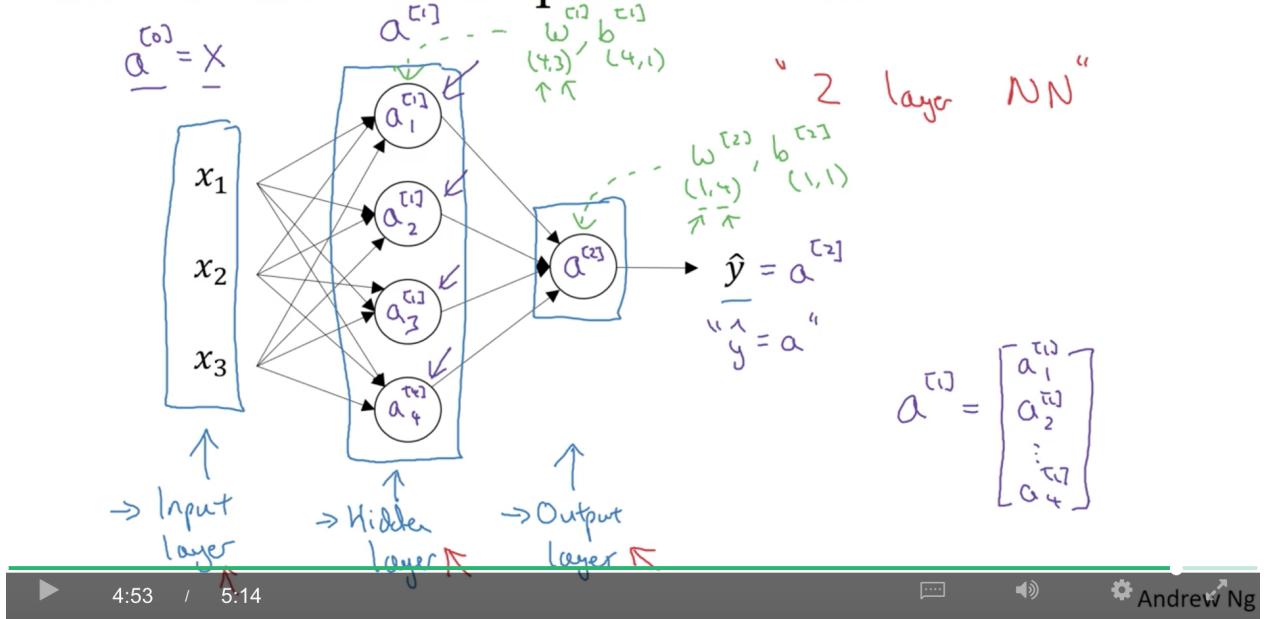
图示是一个2层nn (inputlayer不算在内，有hidden和output两层)。

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 $shape(b^{[1]}) = (4, 1)$ 。

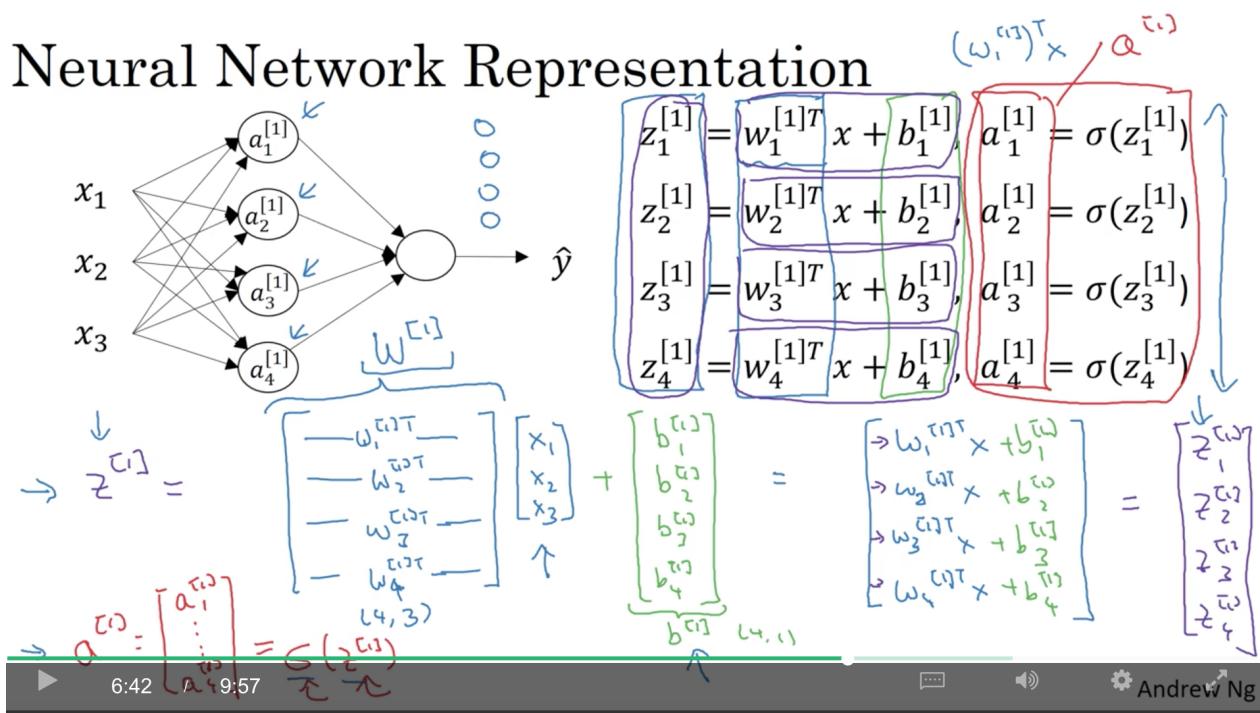
而 $shape(w^{[2]}) = (1, 4)$ ，因为只有1个神经元，输入是3维。同理 $shape(b^{[2]}) = (1, 1)$ 。

Neural Network Representation

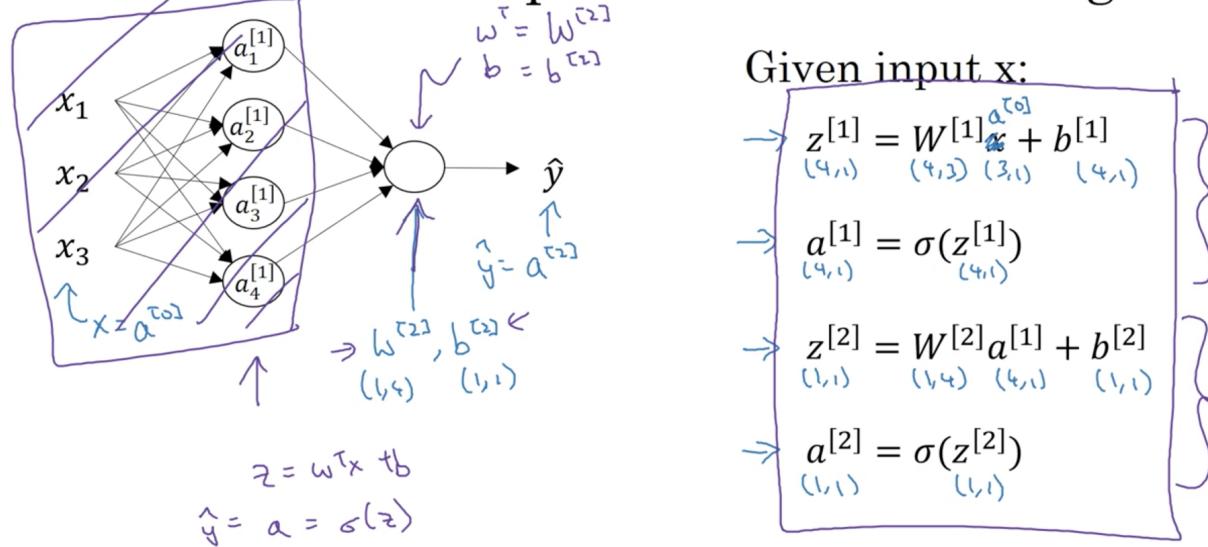


computing a neural network's output

Neural Network Representation

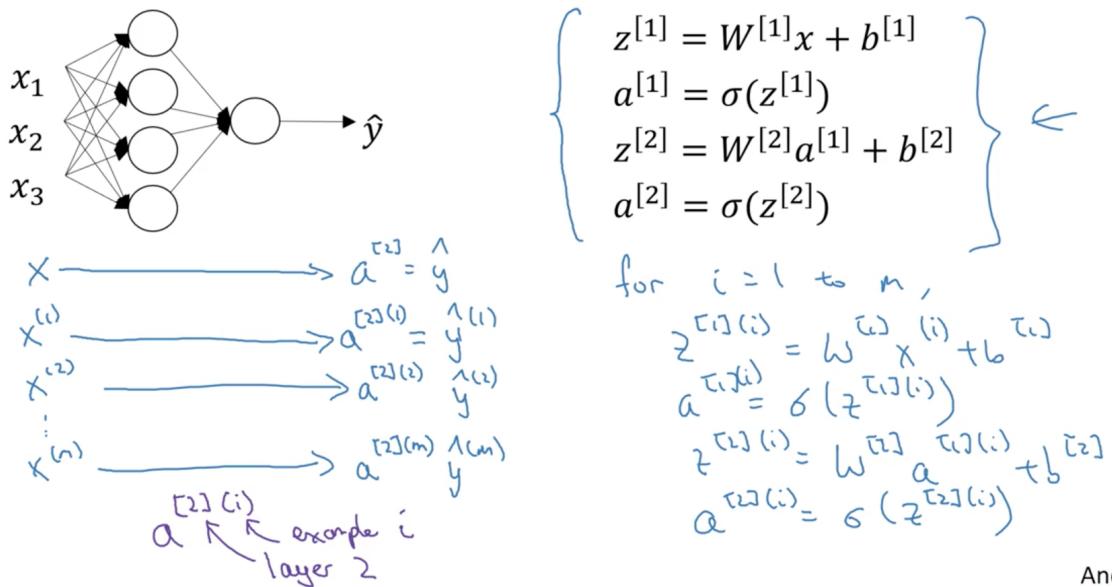


Neural Network Representation learning



vectorizing across multiple examples

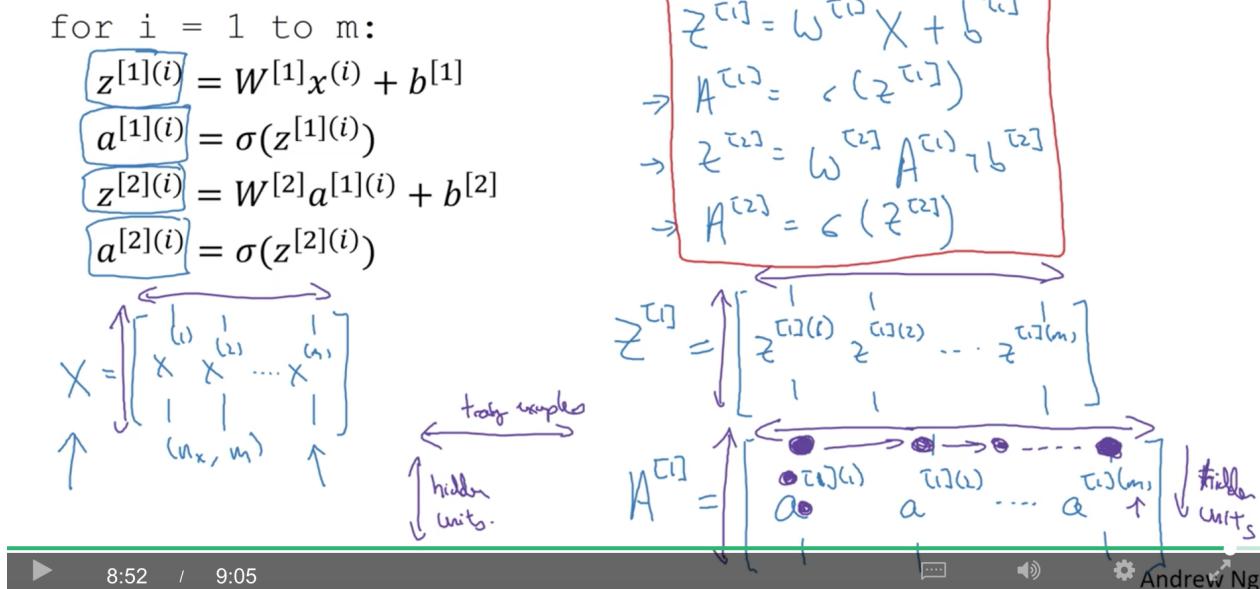
Vectorizing across multiple examples



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



explanation for vectorized implementation

Justification for vectorized implementation

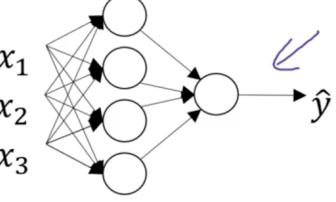
$$\begin{aligned}
 z^{(1)(1)} &= w^{(1)} x^{(1)} + b^{(1)}, \\
 z^{(1)(2)} &= w^{(1)} x^{(2)} + b^{(1)}, \\
 z^{(1)(3)} &= w^{(1)} x^{(3)} + b^{(1)}
 \end{aligned}$$

$$\begin{aligned}
 w^{(1)} &= \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix}, & w^{(1)} x^{(1)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, & w^{(1)} x^{(2)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, & w^{(1)} x^{(3)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}
 \end{aligned}$$

$$\begin{aligned}
 w^{(1)} \begin{bmatrix} 1 & x^{(1)} & x^{(2)} & x^{(3)} \dots \end{bmatrix} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} z^{(1)(1)} \\ z^{(1)(2)} \\ z^{(1)(3)} \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} = z^{(1)}
 \end{aligned}$$

$\boxed{z^{(1)} = w^{(1)} X + b^{(1)}}$

Recap of vectorizing across multiple examples



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

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

`for i = 1 to m`

$$\begin{cases} \rightarrow z^{[1](i)} = W^{[1]} x^{(i)} + b^{[1]} \\ \rightarrow a^{[1](i)} = \sigma(z^{[1](i)}) \end{cases}$$

$$\begin{cases} \rightarrow z^{[2](i)} = W^{[2]} a^{[1](i)} + b^{[2]} \\ \rightarrow a^{[2](i)} = \sigma(z^{[2](i)}) \end{cases}$$

$$Z^{[1]} = W^{[1]} X + b^{[1]} \leftarrow \boxed{w^{[1]} A^{[1]} + b^{[1]}}$$

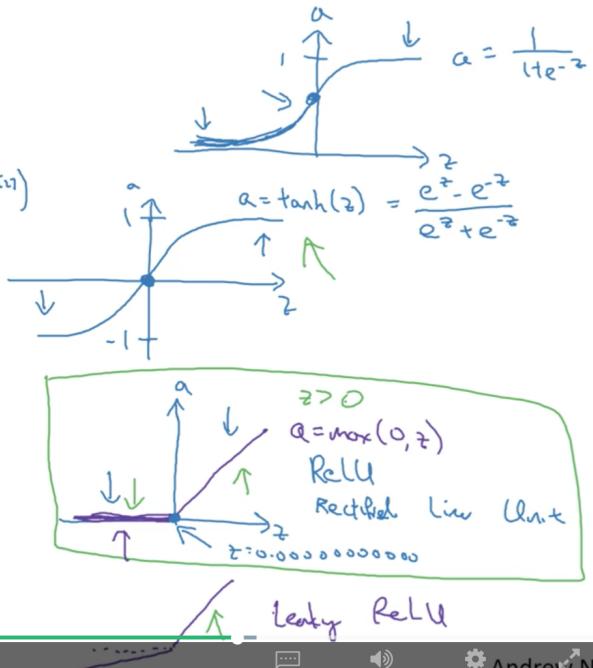
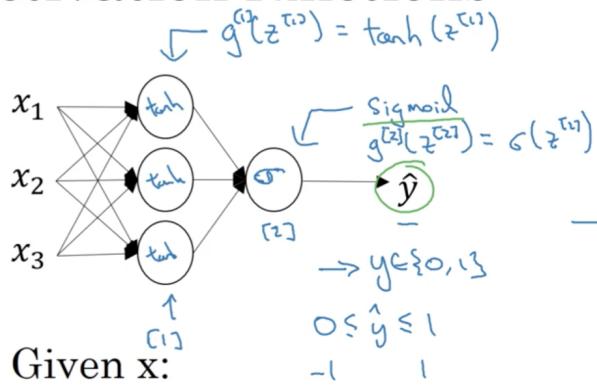
$$\begin{cases} A^{[1]} = \sigma(Z^{[1]}) \\ Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} = \sigma(Z^{[2]}) \end{cases}$$

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



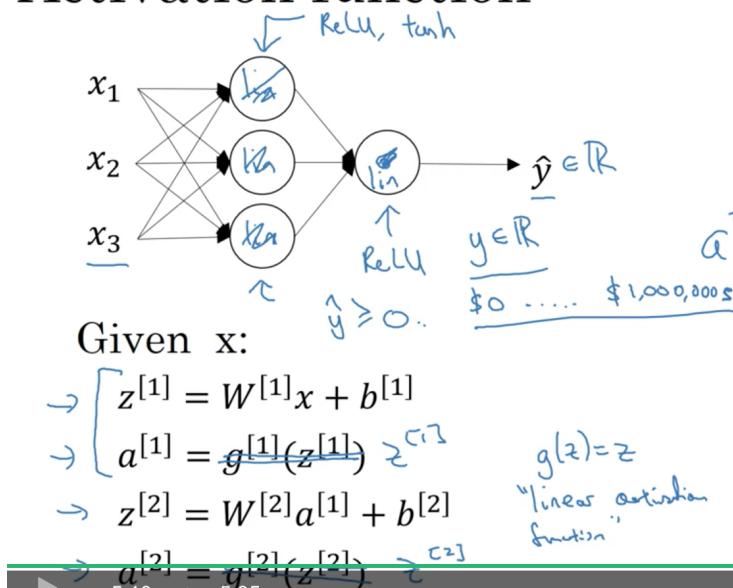
why do you need non-linear activation functions

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。

Activation function

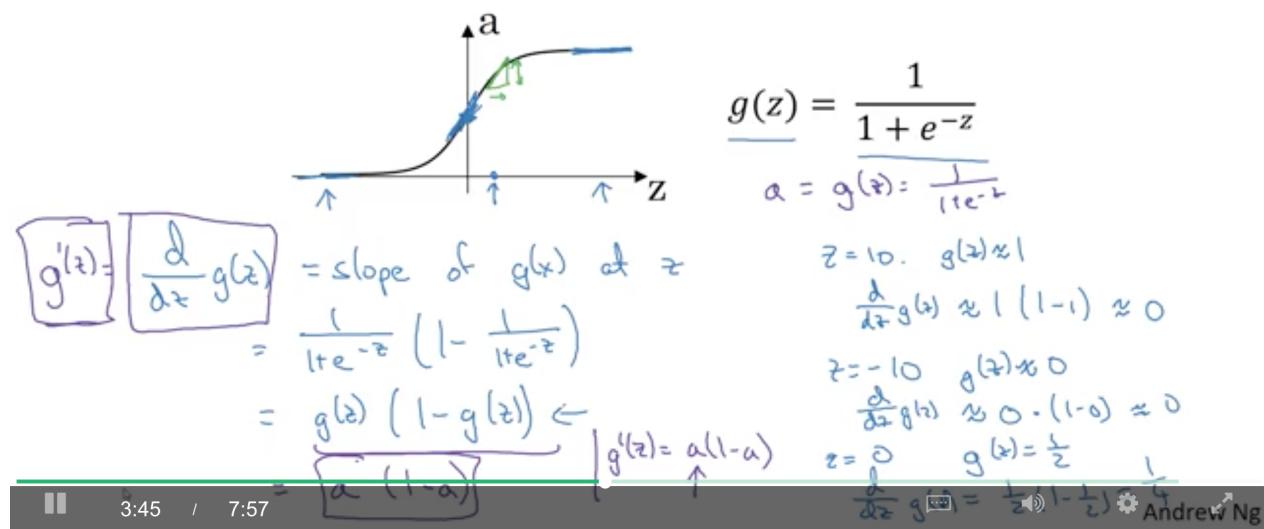


$$\begin{aligned}
 a^{[1]} &= z^{[1]} = W^{[1]}x + b^{[1]} \\
 a^{[2]} &= z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\
 a^{[2]} &= W^{[2]}(\underbrace{W^{[1]}x + b^{[1]}}_{a^{[1]}}) + b^{[2]} \\
 &= (\underbrace{W^{[2]}W^{[1]}}_{w'})x + (\underbrace{W^{[2]}b^{[1]} + b^{[2]}}_{b'}) \\
 &= w'x + b'
 \end{aligned}$$

$g(z) = z$

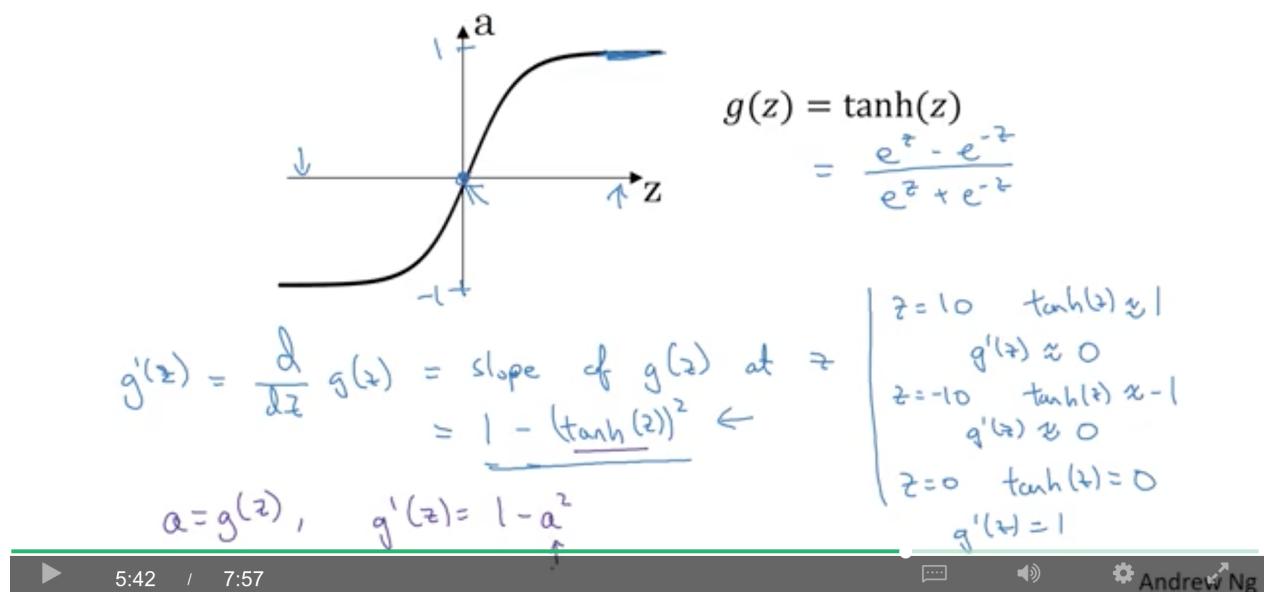
derivatives of activation functions

Sigmoid activation function



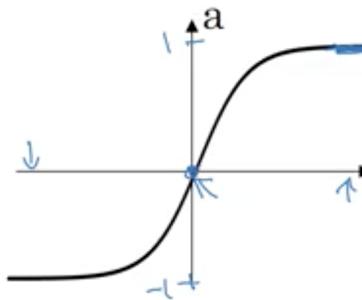
tanh的导数

Tanh activation function



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)

Tanh activation function



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

$$g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z = 1 - (\tanh(z))^2$$
$$a = g(z), \quad g'(z) = 1 - a^2$$

$$\begin{cases} z=10 & \tanh(z) \approx 1 \\ z=-10 & \tanh(z) \approx -1 \\ z=0 & \tanh(z)=0 \\ z=1 & g'(z)=1 \end{cases}$$

▶ 5:42 / 7:57 ⏪ ⏴ ⏵ ⏹ Andrew Ng

gradient descent for neural networks

记住每个W/b的shape!

Gradient descent for neural networks

Parameters: $(w^{(0)}, b^{(0)})$, $(w^{(1)}, b^{(1)})$, $(w^{(2)}, b^{(2)})$, \dots , $(w^{(L)}, b^{(L)})$ $n_x = n^{(0)}, n^{(1)}, \dots, n^{(L)} = 1$

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 predictions ($\hat{y}^{(i)}$), $i=1, \dots, m$)
 $\frac{\partial J}{\partial w^{(l)}} = \frac{\partial J}{\partial w^{(l)}}, \quad \frac{\partial J}{\partial b^{(l)}} = \frac{\partial J}{\partial b^{(l)}}, \dots$

$w^{(l+1)} := w^{(l)} - \alpha \frac{\partial J}{\partial w^{(l)}}$

$b^{(l+1)} := b^{(l)} - \alpha \frac{\partial J}{\partial b^{(l)}}$

▶ 3:46 / 9:57 ⏪ ⏴ ⏵ ⏹ ↻

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

Formulas for computing derivatives

Forward propagation:

$$\begin{aligned} z^{[1]} &= w^{[1]} X + b^{[1]} \\ A^{[1]} &= g^{[1]}(z^{[1]}) \leftarrow \\ z^{[2]} &= w^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(z^{[2]}) = \underline{\sigma}(z^{[2]}) \end{aligned}$$

Back propagation:

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

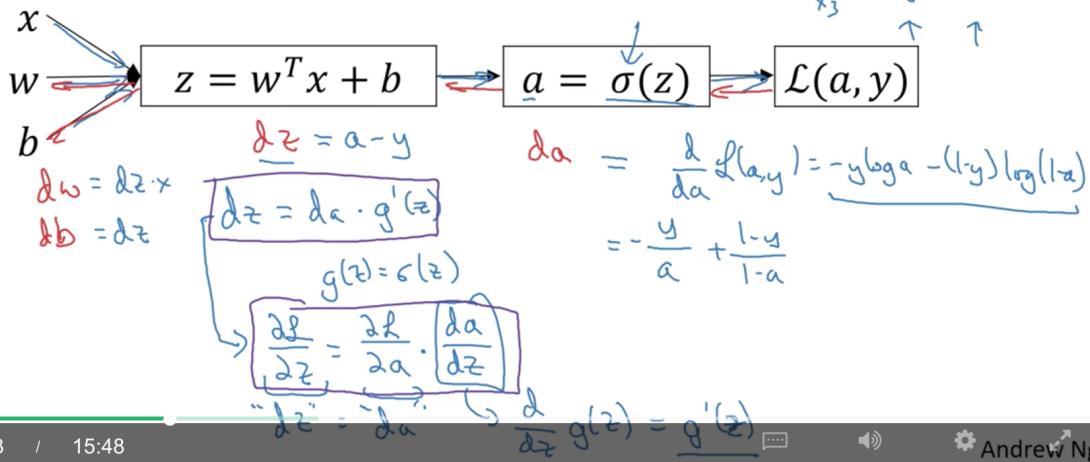


backpropagation intuition

先回顾一下lr:

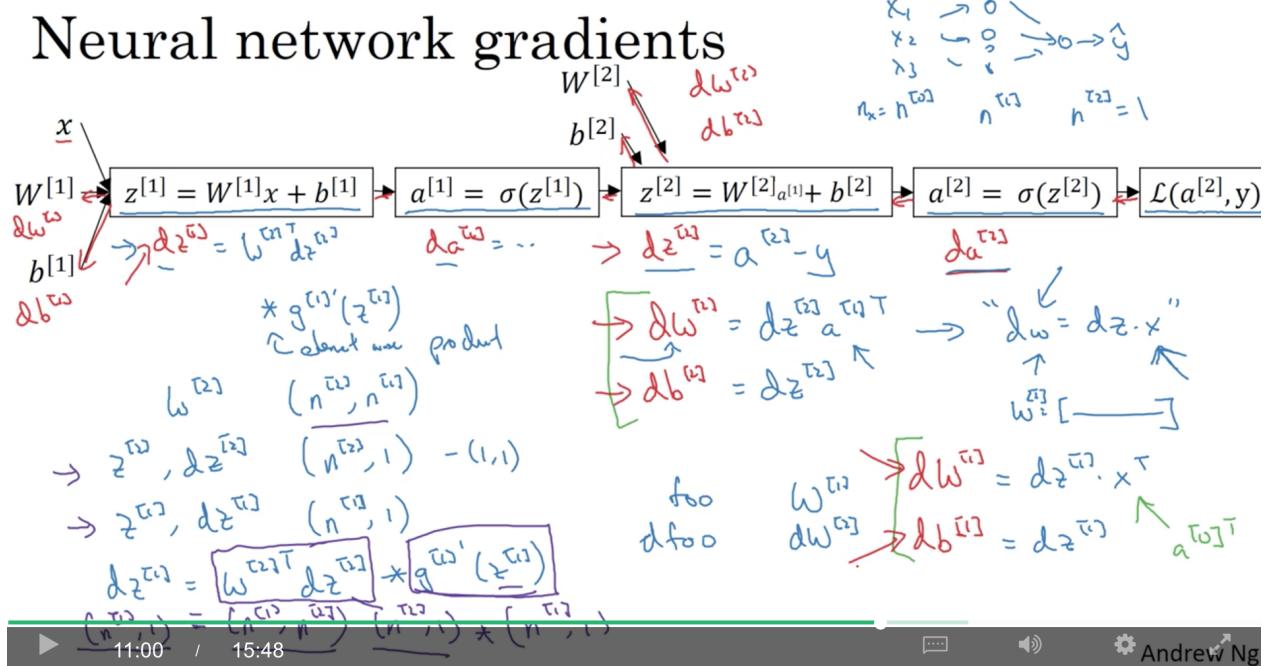
Computing gradients

Logistic regression



Andrew Ng

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:

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

$$Z^{[1]} = \begin{bmatrix} 1 & z_1^{[1]} & z_2^{[1]} & \dots & z_n^{[1]} \end{bmatrix}$$

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

12:27 / 15:48 Andrew Ng

Summary of gradient descent

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

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

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

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

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

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

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

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

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

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

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

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

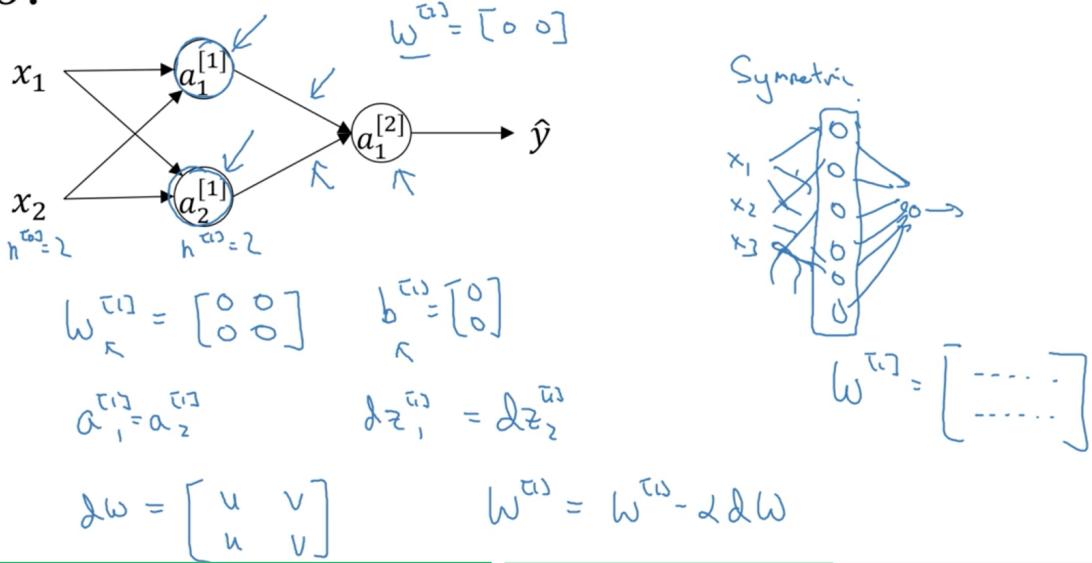
14:22 / 15:48 Andrew Ng

问：为何有的有 $1/m$ ，有的没有。。。。。

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?



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

