

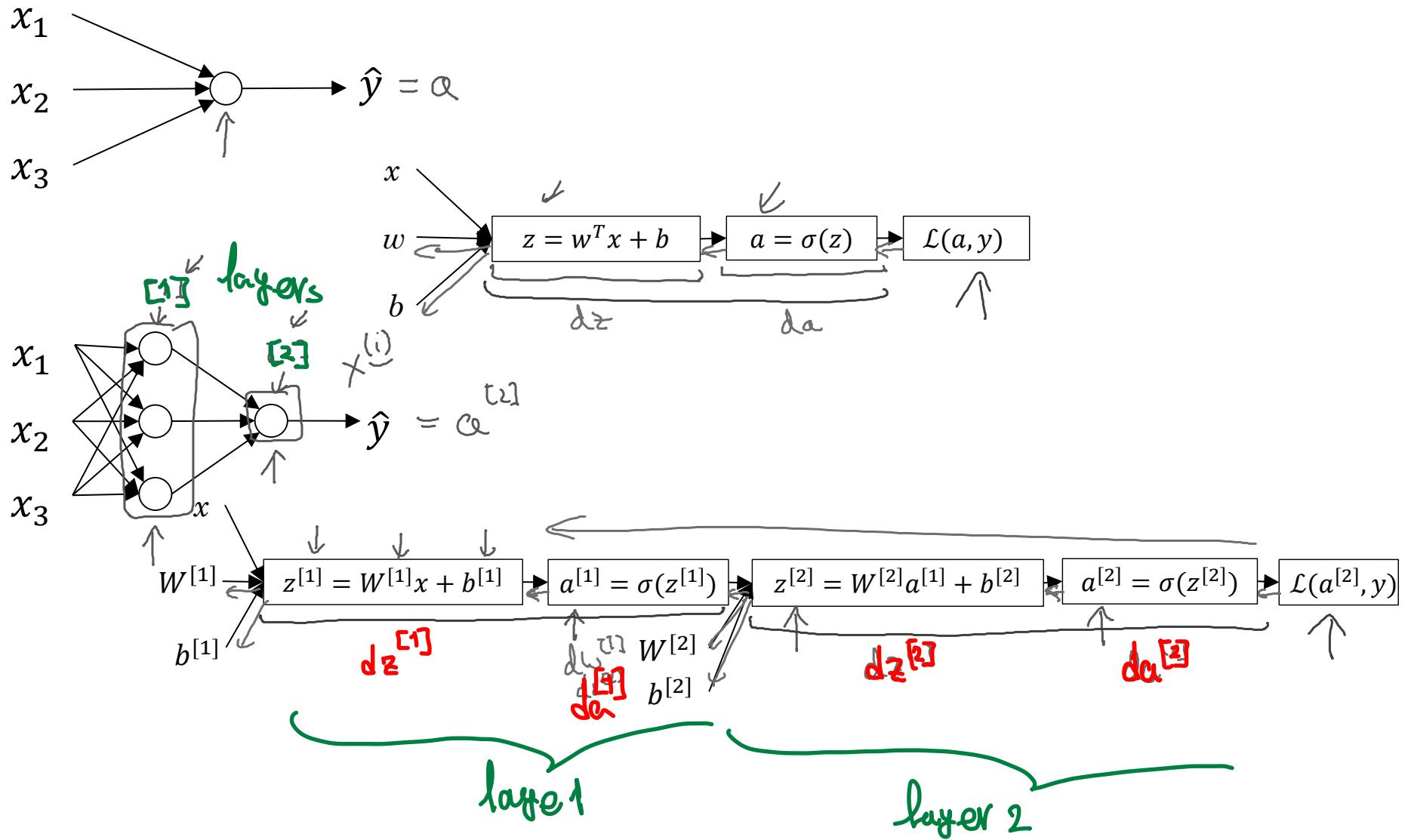


deeplearning.ai

One hidden layer
Neural Network

Neural Networks
Overview

What is a Neural Network?



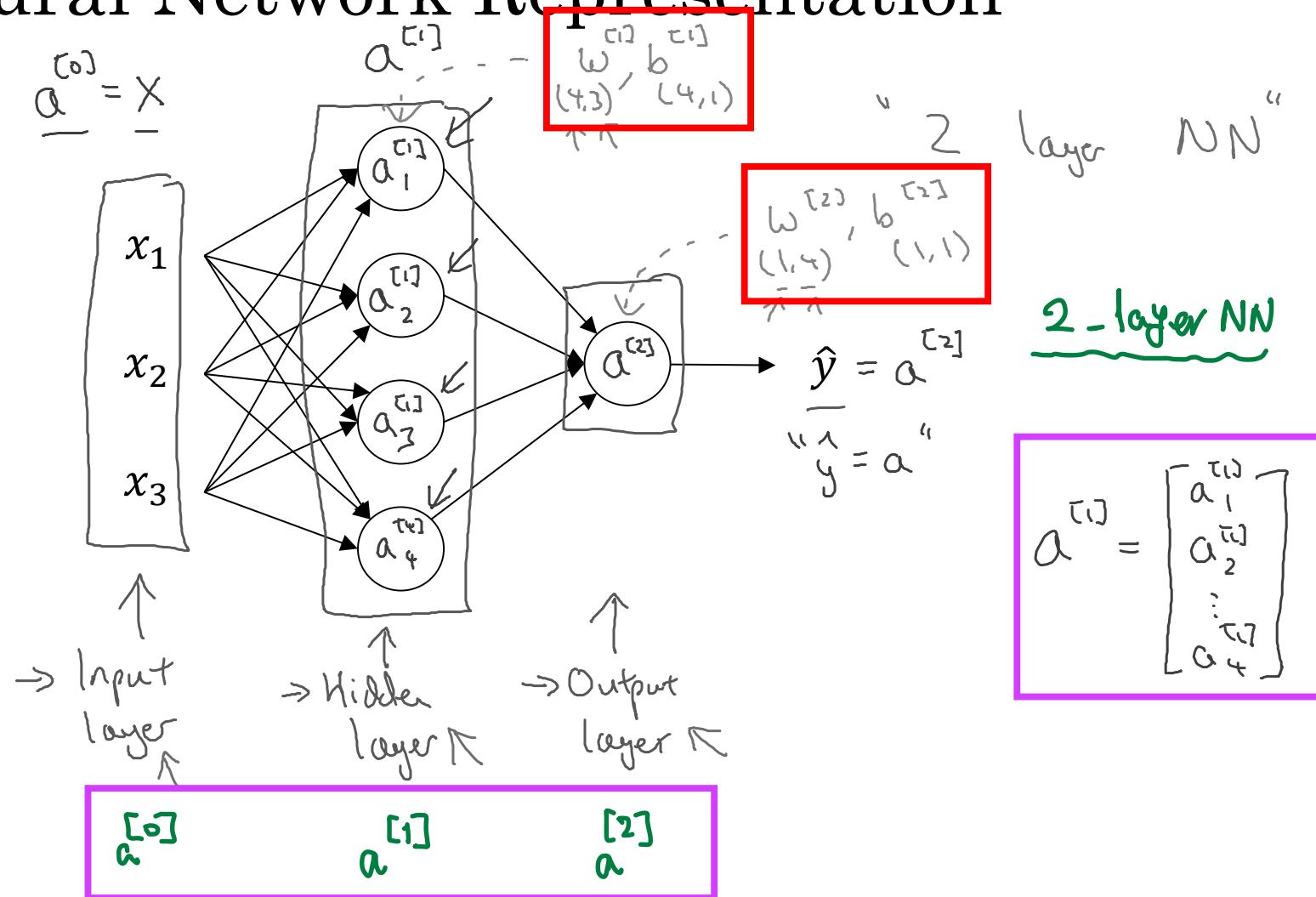


deeplearning.ai

One hidden layer
Neural Network

Neural Network
Representation

Neural Network Representation



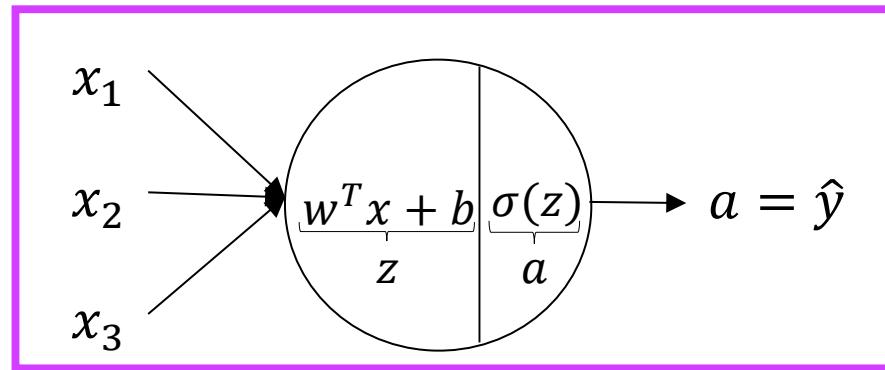


deeplearning.ai

One hidden layer
Neural Network

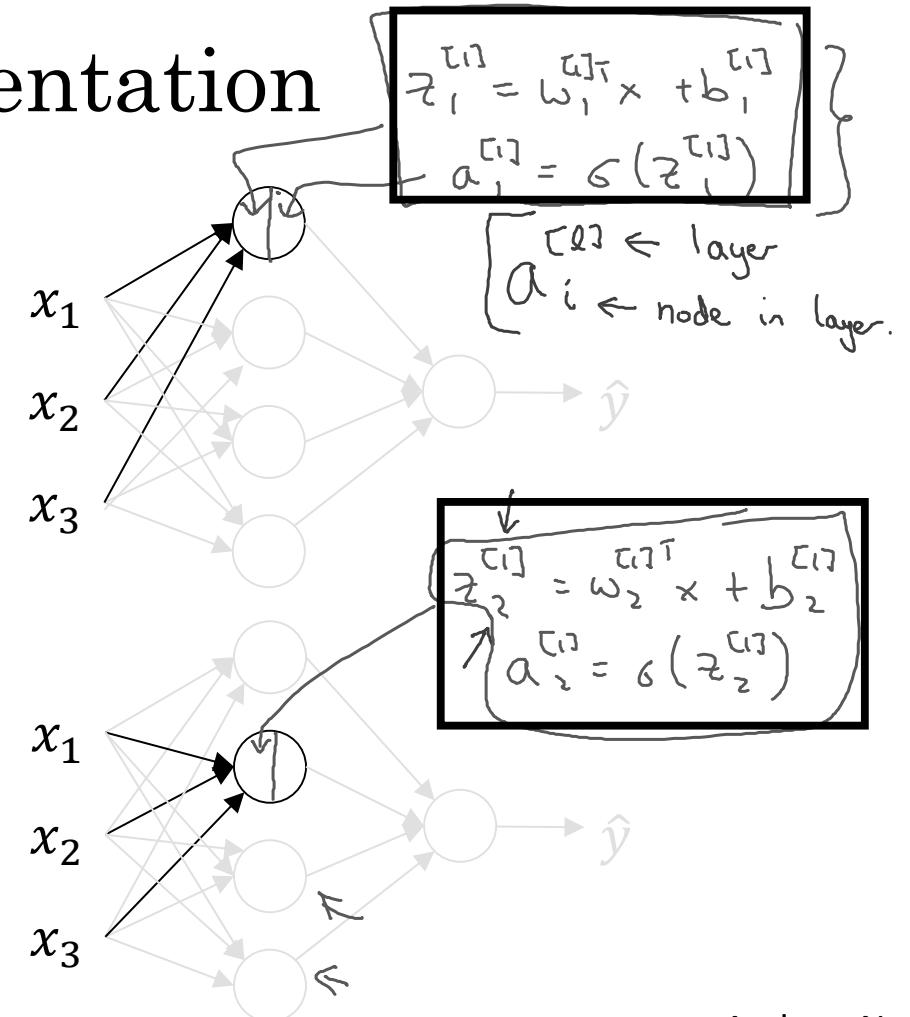
Computing a
Neural Network's
Output

Neural Network Representation



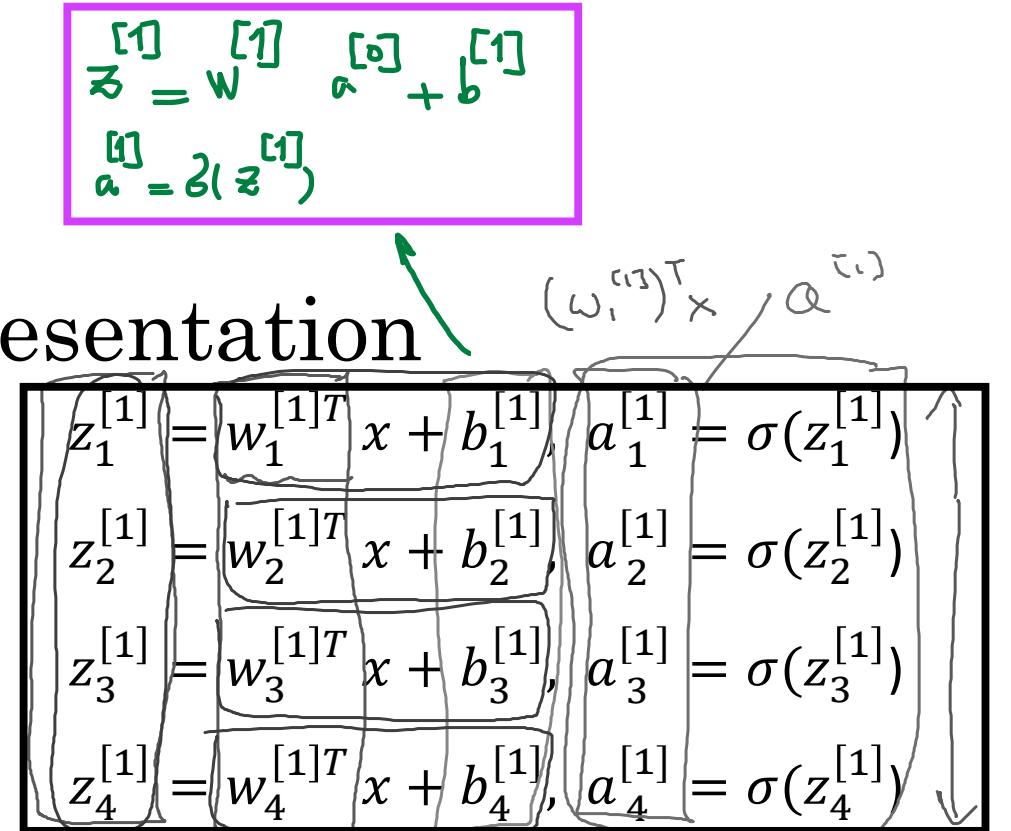
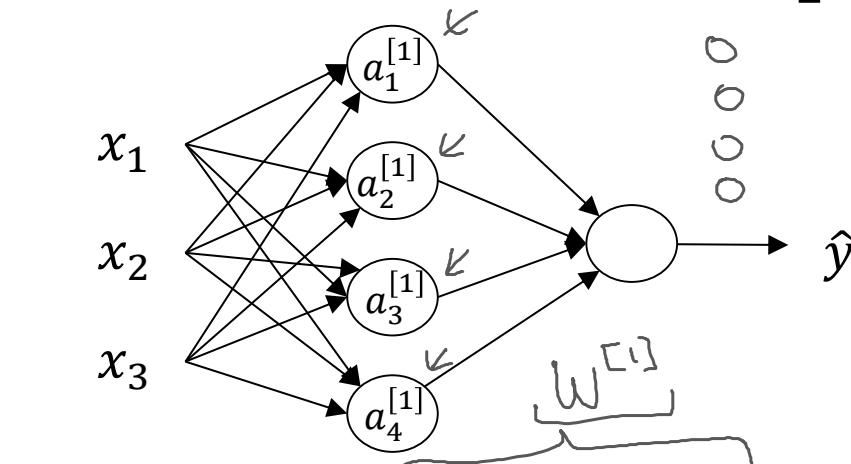
$$z = w^T x + b$$

$$a = \sigma(z)$$



Andrew Ng

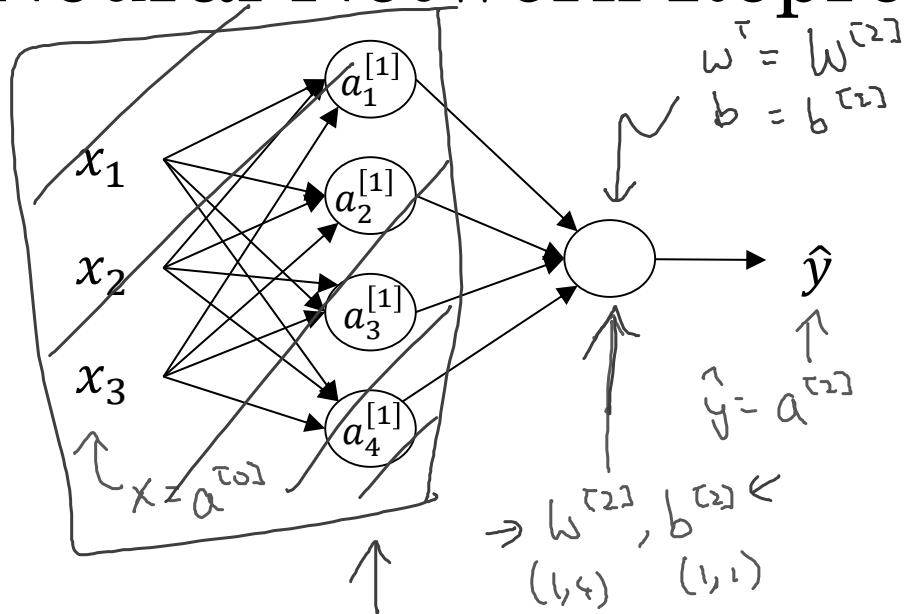
Neural Network Representation



$$\begin{aligned}
 \rightarrow z^{[1]} &= \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} \\
 \rightarrow a^{[1]} &= \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]}) = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}
 \end{aligned}$$

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

Neural Network Representation learning



Given input x :

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

(4,1) (4,3) (3,1) (4,1)

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

(4,1)

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

(1,1) (1,4) (4,1) (1,1)

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

(1,1)

4 lines of code
(one sample)

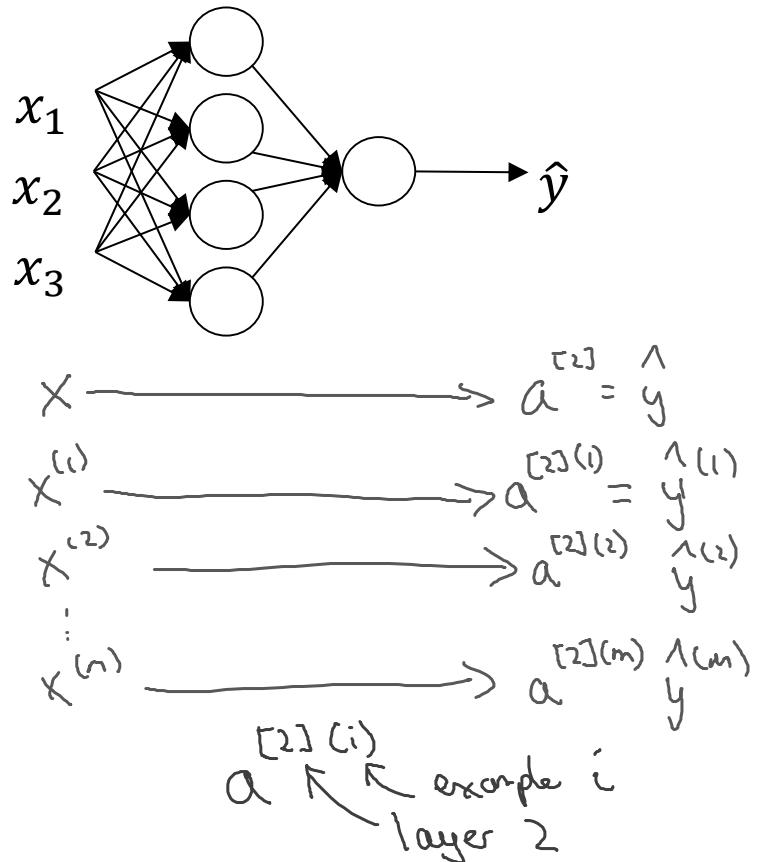


deeplearning.ai

One hidden layer
Neural Network

Vectorizing across
multiple examples

Vectorizing across multiple examples



$$\left\{ \begin{array}{l} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{array} \right.$$

```
for i = 1 to n,
     $z^{[1](i)} = w^{[1]}x^{(i)} + b^{[1]}$ 
     $a^{[1](i)} = \sigma(z^{[1](i)})$ 
     $z^{[2](i)} = w^{[2]}a^{[1](i)} + b^{[2]}$ 
     $a^{[2](i)} = \sigma(z^{[2](i)})$ 
```

m Samples, 2 layers

درباره می خواهیم از سر loop ها شوییم ... نیاز است

Vectorizing across multiple examples

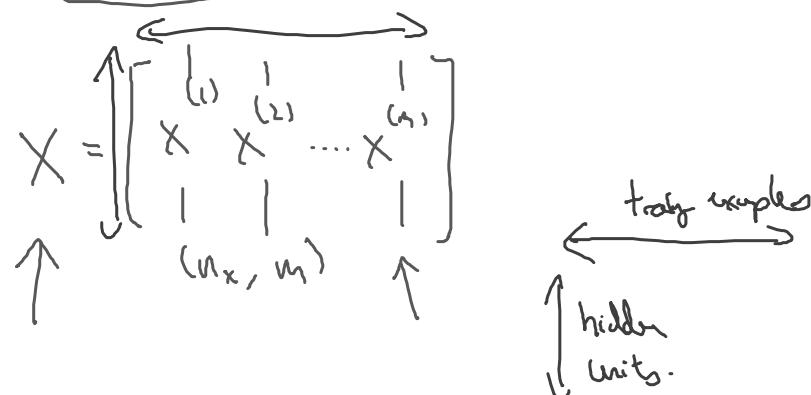
for $i = 1$ to m :

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

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

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

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



$\boxed{\begin{array}{l} z^{[1]} = W^{[1]}X + b^{[1]} \\ \rightarrow A^{[1]} = \sigma(z^{[1]}) \\ \rightarrow z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} \\ \rightarrow A^{[2]} = \sigma(z^{[2]}) \end{array}}$

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

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

hidden units.

Justification for vectorized implementation

$$z^{(1)(1)} = w^{(1)} x^{(1)} + b^{(1)}, \quad z^{(1)(2)} = w^{(1)} x^{(2)} + b^{(1)}, \quad z^{(1)(3)} = w^{(1)} x^{(3)} + b^{(1)}$$

$$w^{(1)} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

$$w^{(1)} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

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

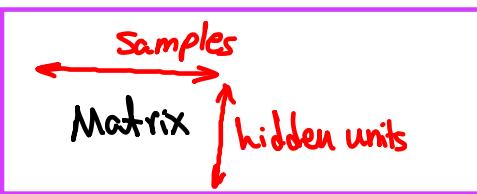
$$w^{(1)} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

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

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

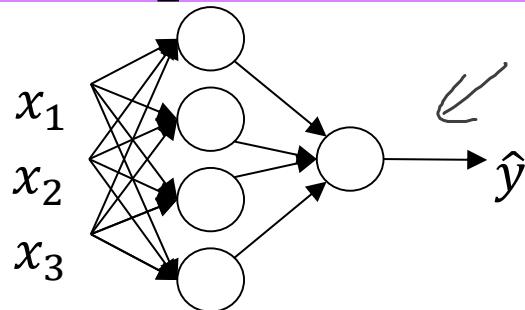
\times

$$w^{(1)} x^{(1)} = z^{(1)(1)}$$



$a_k^{[i](j)}$ = node k , layer i , sample j

Recap of vectorizing across multiple examples



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

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

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

for $i = 1$ to m

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

$\leftarrow A^{[1]}$ $x = a^{[1]}$ $x^{(i)} = a^{[1](i)}$

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

Fully vectorized version for m samples in 2 layer NN.

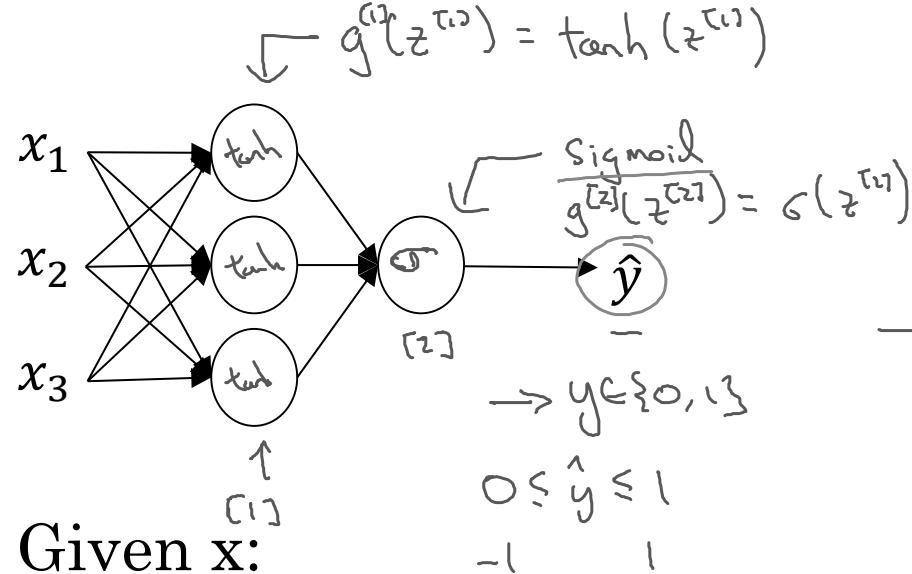


deeplearning.ai

One hidden layer
Neural Network

Activation functions

Activation functions

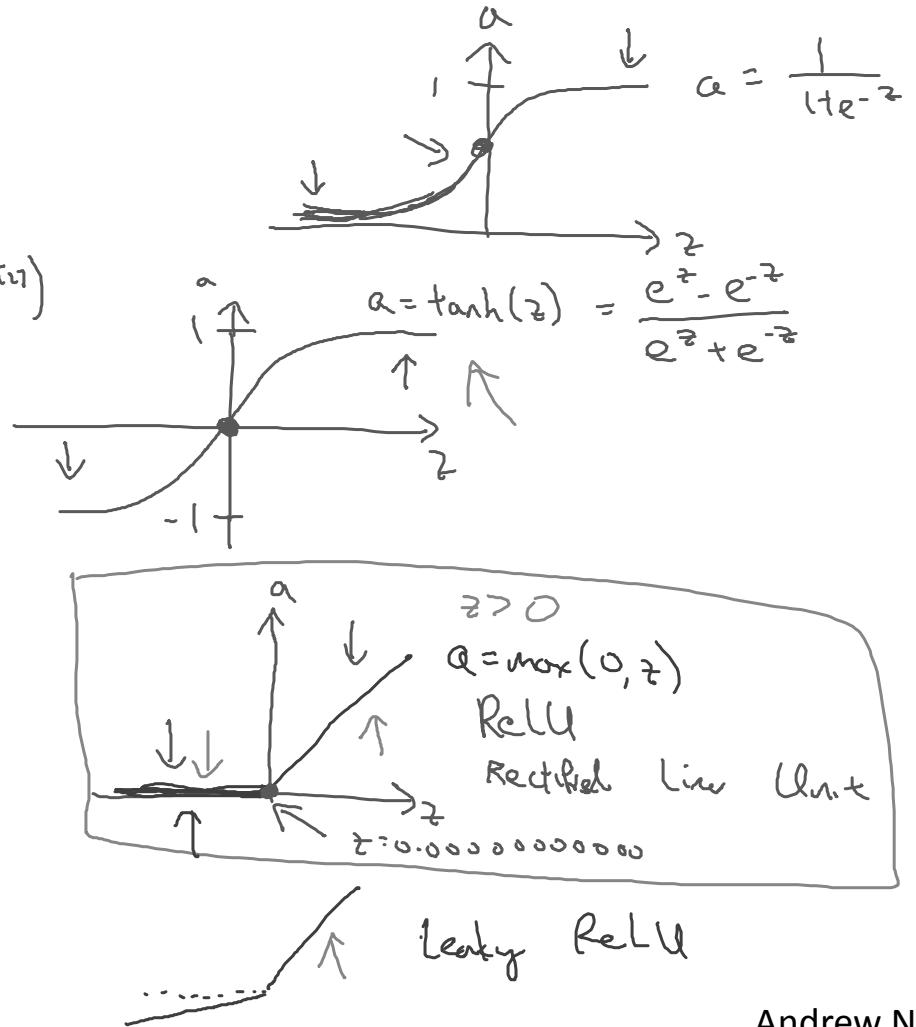


$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$\rightarrow a^{[1]} = \underline{\sigma}(z^{[1]}) \quad g^{(1)}(z^{(1)})$$

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

$$\rightarrow a^{[2]} = \underline{\sigma}(z^{[2]}) \quad g^{(2)}(z^{(2)})$$

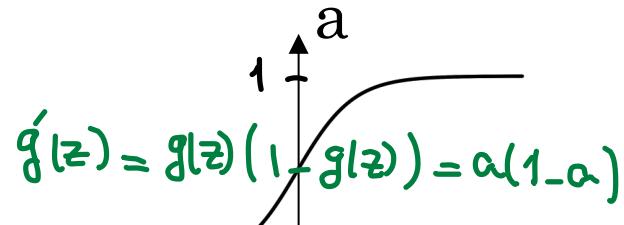


Andrew Ng

Binary Classification: Sigmoid for output layer and ReLU for other units

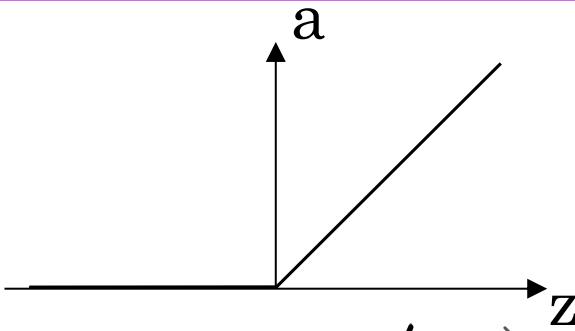
$$a = g(z)$$

Pros and cons of activation functions



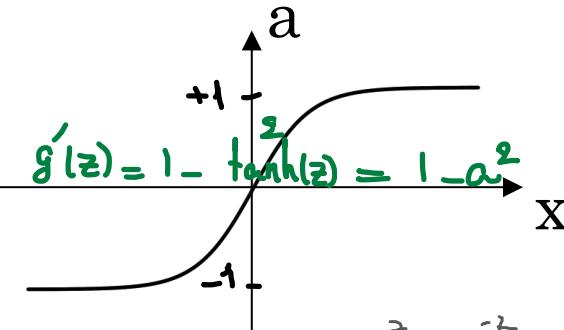
sigmoid: $a = \frac{1}{1 + e^{-z}}$

never use that



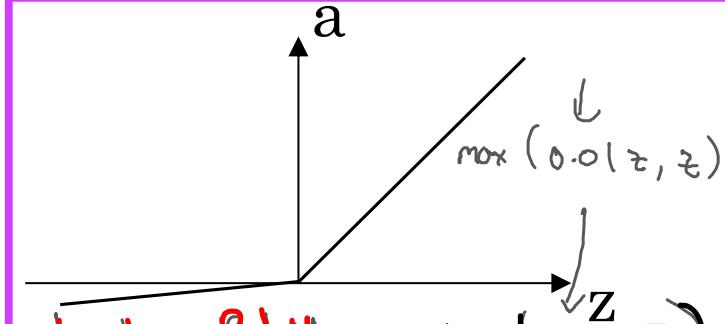
very good

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



tanh: $a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

just for hidden nodes, better than Sigmoid



better than ReLU

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

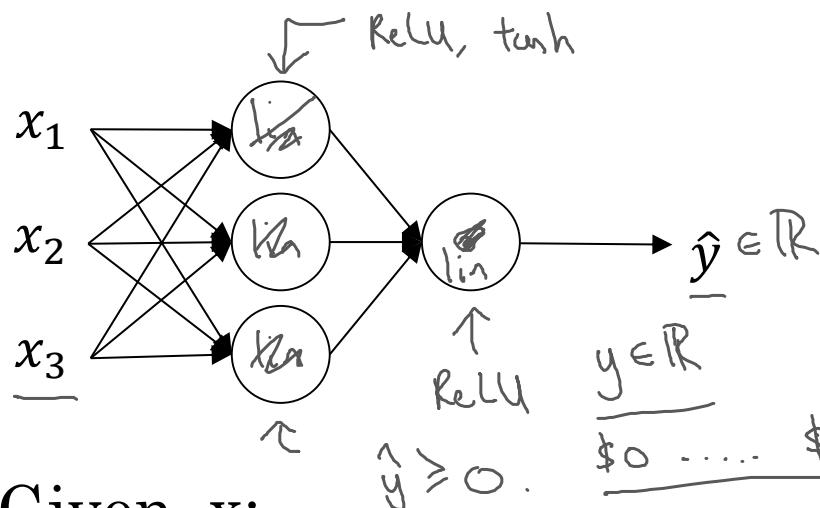


deeplearning.ai

One hidden layer Neural Network

Why do you
need non-linear
activation functions?

Activation function



Given x :

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

$$\rightarrow a^{[1]} = \cancel{g^{[1]}(z^{[1]})} \geq \cancel{z^{[1]}}$$

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

$$\rightarrow a^{[2]} = \cancel{g^{[2]}(z^{[2]})} \geq \cancel{z^{[2]}}$$

$g(z) = z$
"linear activation
function"

$$a^{[1]} = z^{[1]} = \underbrace{W^{[1]}x + b^{[1]}}_{a^{[1]}}$$

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

$$a^{[2]} = W^{[2]} \left(\underbrace{W^{[1]}x + b^{[1]}}_{a^{[1]}} \right) + b^{[2]}$$

$$= (\underbrace{W^{[2]} W^{[1]}}_{w'})x + (\underbrace{W^{[2]} b^{[1]} + b^{[2]}}_{b'})$$

$$= \underbrace{w'x + b'}_{g(z) = z}$$



deeplearning.ai

One hidden layer
Neural Network

Gradient descent for
neural networks

Gradient descent for neural networks

Dimensions Parameters: $\underbrace{(\omega^{[0]}, b^{[0]}, \omega^{[1]}, b^{[1]})}_{(n^{[0]}, n^{[0]})}, \underbrace{(\omega^{[1]}, b^{[1]}, \omega^{[2]}, b^{[2]})}_{(n^{[1]}, n^{[1]})}, \underbrace{(n^{[2]}, b^{[2]})}_{(n^{[2]}, 1)}$

Cost function: $J(\underbrace{\omega^{[0]}, b^{[0]}}, \underbrace{\omega^{[1]}, b^{[1]}}, \underbrace{\omega^{[2]}, b^{[2]}}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i)$

$n_x = n^{[0]}$ $n^{[1]}$ $n^{[2]} = 1$

input units hidden units output units

Gradient descent:

→ Repeat {

→ Compute predicts $(\hat{y}^{(i)}, i=1 \dots m)$

$\frac{\partial J}{\partial \omega^{[l]}} = \frac{\partial J}{\partial \omega^{[l]}} , \frac{\partial J}{\partial b^{[l]}} = \frac{\partial J}{\partial b^{[l]}} , \dots$

$\omega^{[l]} := \omega^{[l]} - \alpha \frac{\partial J}{\partial \omega^{[l]}}$

$b^{[l]} := b^{[l]} - \alpha \frac{\partial J}{\partial b^{[l]}}$

$\omega^{[2]} := \dots$

$b^{[2]} := \dots$

One Iteration

Formulas for computing derivatives

Forward propagation:

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

4 equations

Back propagation:

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

6 equations

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

$$(n^{[1]}) \leftarrow \quad (n^{[2]}, 1) \leftarrow$$

Prevent

(n^{[1]}, m)

(n^{[2]}, 1)

(n^{[1]}, m)

(n^{[2]}, 1)

(n^{[1]}, 1)

(n^{[2]}, 1)

(n^{[1]}, 1)



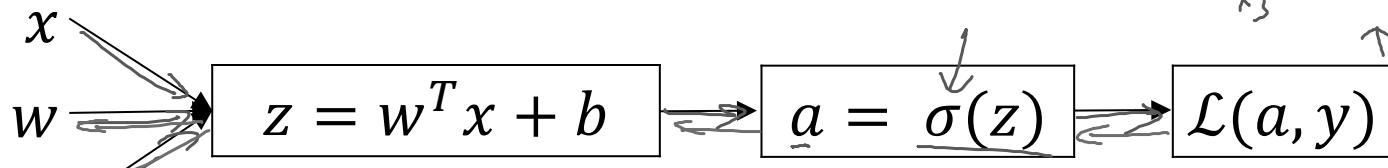
deeplearning.ai

One hidden layer
Neural Network

Backpropagation
intuition (Optional)

Computing gradients

Logistic regression



$$\frac{d\zeta}{dz} = a - y$$

$$dw = dz \cdot x$$

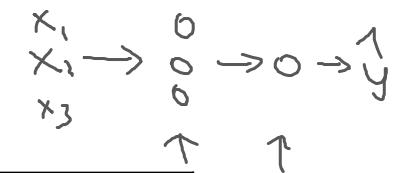
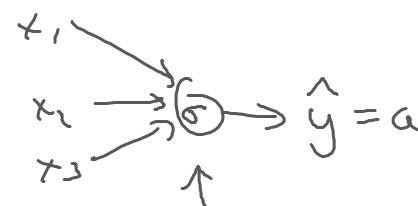
$$db = dz$$

$$dz = da \cdot g'(z)$$

$$g(z) = \sigma(z)$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{da}{dz}$$

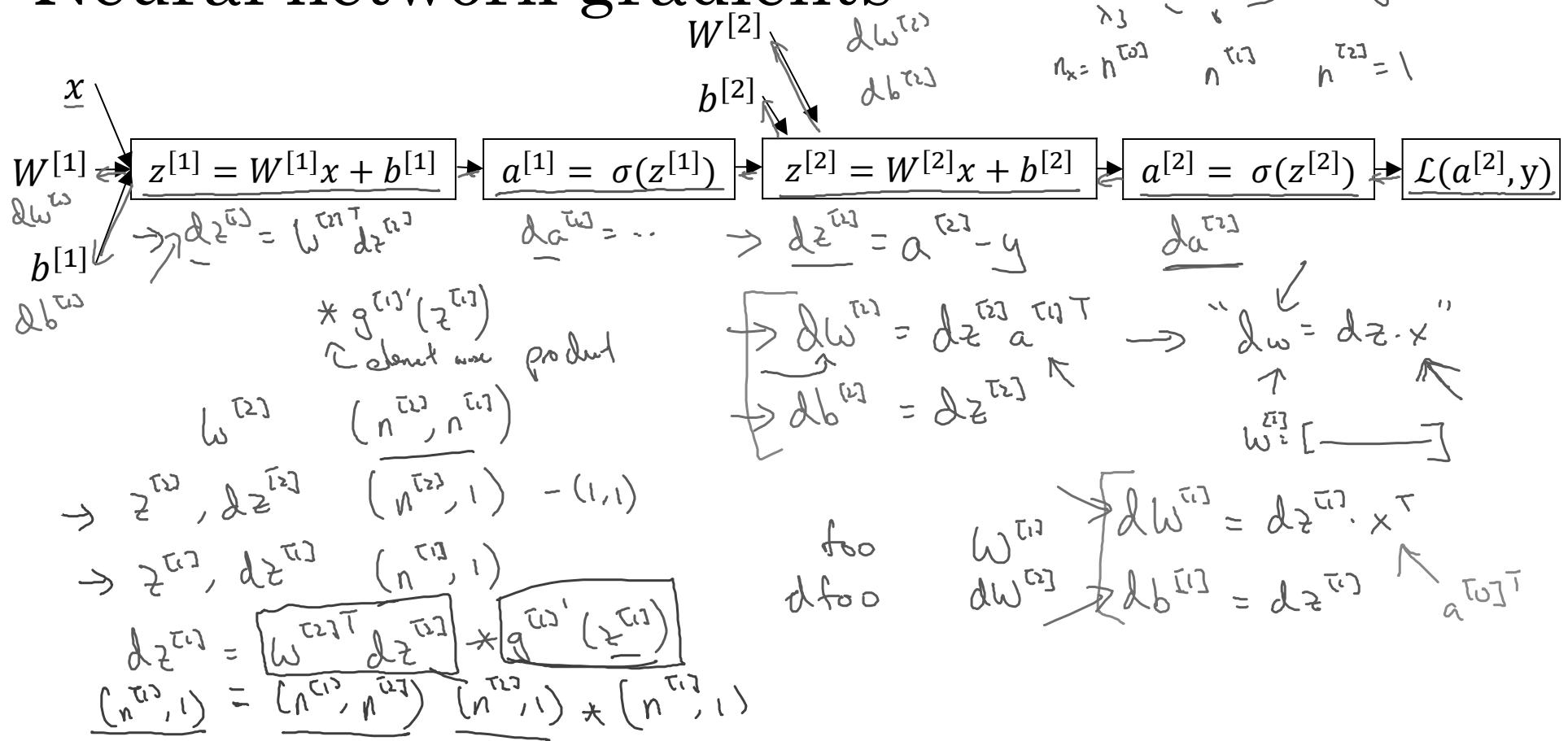
$$"dz" = "da"$$



$$\begin{aligned} \frac{da}{dz} &= \frac{d}{da} L(a, y) = -y \log a - (1-y) \log(1-a) \\ &= -\frac{y}{a} + \frac{1-y}{1-a} \end{aligned}$$

$$\frac{d}{dz} g(z) = g'(z)$$

Neural network gradients



Summary of gradient descent

$$\underline{dZ^{[2]}} = \underline{a^{[2]}} - \underline{y}$$

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

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

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

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

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

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

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

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

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

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

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

Backward Propagation Equations



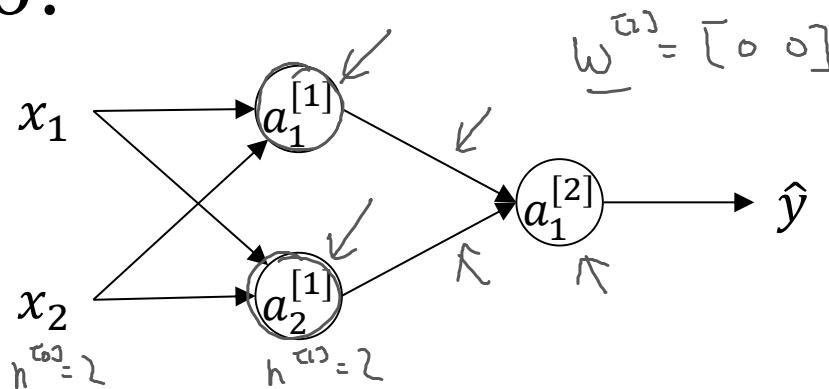
deeplearning.ai

One hidden layer
Neural Network

Random Initialization

Zero Initialization \Rightarrow Symmetry \Rightarrow Canceling Path $\Rightarrow \times$
 \Rightarrow Initial Bias is required or random initialization

What happens if you initialize weights to zero?

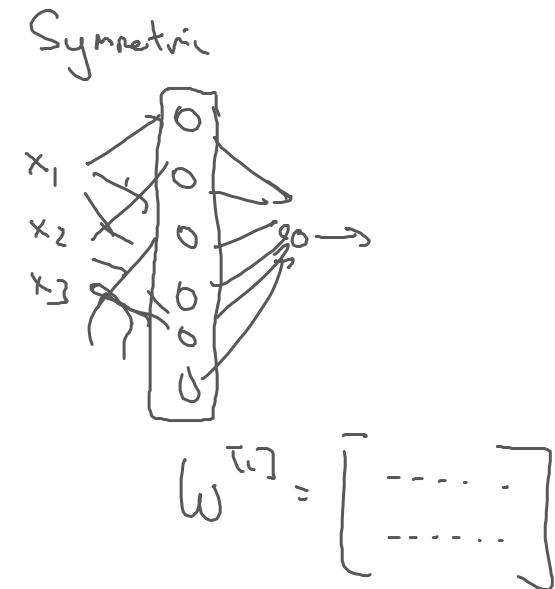


$$w_1^{[1,2]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b_1^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

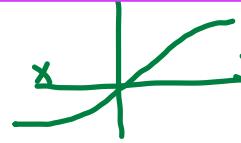
$$a_1^{[1]} = a_2^{[1]} \quad da_1^{[1]} = da_2^{[1]}$$

$$\delta w = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

$$w^{[1]} = w^{[1]} - \lambda \delta w$$

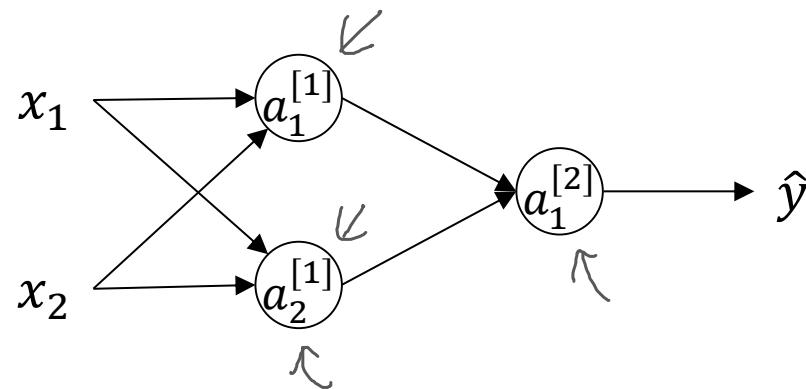


Saturation



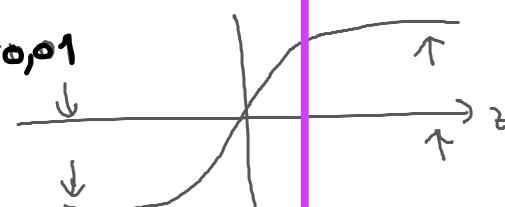
بهره از مقادیر random اولیه خیل بزرگ / کوچک باشد حجون در

Random initialization



$$\rightarrow \omega^{[1]} = \text{np.random.rand}(2, 2) * 0,01$$
$$b^{[1]} = \text{np.zeros}(2, 1)$$
$$\omega^{[2]} = -\text{np.random.rand}(1, 2) * 0,01$$
$$b^{[2]} = 0$$

$$z^{[1]} = \omega^{[1]} x + b^{[1]}$$
$$a^{[1]} = g^{[1]}(z^{[1]})$$



Random Initialization of Parameters