# One hidden layer
# Neural Network

## Neural Networks
## Overview

deeplearning.ai

# Neural Network Representation

$a^{[0]} = X$

$a^{[1]}$

$w^{[1]}, b^{[1]}$
$(4,3)$ $(4,1)$

"2 layer NN"

$a_1^{[1]}$

$a_2^{[1]}$

$w^{[2]}, b^{[2]}$
$(1,4)$ $(1,1)$

$x_1$

$x_2$

$x_3$

$a_3^{[1]}$

$a_4^{[1]}$

$a^{[2]}$

$\hat{y} = a^{[2]}$

"$\hat{y} = a$"

$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix}$

→ Input layer

→ Hidden layer

→ Output layer

# Neural Network Representation

$x_1$

$x_2$

$x_3$

$\underbrace{w^T x + b}_{z} \bigg| \underbrace{\sigma(z)}_{a} \longrightarrow a = \hat{y}$

$x_1$

$x_2$

$x_3$

$\hat{y}$

$z = w^T x + b$

$a = \sigma(z)$

# Neural Network Representation

$x_1$
$x_2$
$x_3$

$w^T x + b \mid \sigma(z)$ → $a = \hat{y}$

$z$ $a$

$z = w^T x + b$

$a = \sigma(z)$

$x_1$
$x_2$
$x_3$

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

$a_i^{[2]} \leftarrow$ layer
$a_i \leftarrow$ node in layer.

$x_1$
$x_2$
$x_3$

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

Andrew Ng

---

# Neural Network Representation

$x_1$
$x_2$
$x_3$

$a_1^{[1]}$
$a_2^{[1]}$
$a_3^{[1]}$
$a_4^{[1]}$

→ $\hat{y}$

$(w_1^{[1]})^T x$    $a^{[1]}$

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

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

$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$

$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$

$W^{[1]}$

$z^{[1]} = \begin{bmatrix} - w_1^{[1]T} - \\ - w_2^{[1]T} - \\ - w_3^{[1]T} - \\ - w_4^{[1]T} - \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ w_3^{[1]T} x + b_3^{[1]} \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$

(4,3)    $b^{[1]}$ (4,1)

$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]})$

Andrew Ng

# Neural Network Representation learning



$\omega^T = W^{[2]}$
$b = b^{[2]}$

$x_1$

$a_1^{[1]}$
$a_2^{[1]}$
$a_3^{[1]}$
$a_4^{[1]}$

$x_2$

$x_3$

$\hat{y}$

$\hat{y} = a^{[2]}$

$x = a^{[0]}$

$W^{[2]}, b^{[2]}$
$(1,4) \quad (1,1)$

$z = \omega^T x + b$
$\hat{y} = a = \sigma(z)$

Given input x:

$$z^{[1]} = W^{[1]}\underset{a^{[0]}}{x} + b^{[1]}$$
$$(4,1) \quad (4,3) \quad (3,1) \quad (4,1)$$

$$a^{[1]} = \sigma(z^{[1]})$$
$$(4,1) \qquad\qquad (4,1)$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
$$(1,1) \quad (1,4) \quad (4,1) \quad (1,1)$$

$$a^{[2]} = \sigma(z^{[2]})$$
$$(1,1) \qquad\qquad (1,1)$$

One hidden layer
Neural Network

deeplearning.ai

Vectorizing across
multiple examples

# Vectorizing across multiple examples



$x_1$
$x_2$
$x_3$

$\hat{y}$

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

$X \longrightarrow a^{[2]} = \hat{y}$

$x^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)}$

$x^{(2)} \longrightarrow a^{[2](2)} = \hat{y}^{(2)}$

$x^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)}$

$a^{[2](i)}$ ← example $i$
layer 2

for $i = 1$ to $m$,

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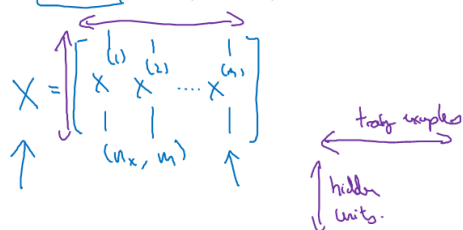$$a^{[2](i)} = \sigma(z^{[2](i)})$$

# Vectorizing across multiple examples
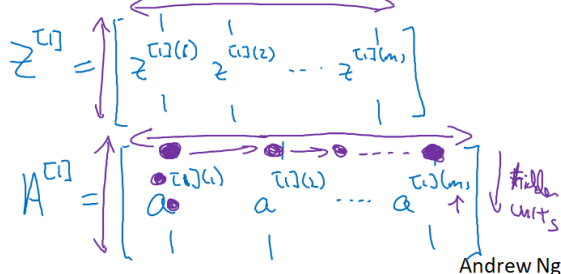
for i = 1 to m:

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

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$
$$\rightarrow A^{[1]} = \sigma(Z^{[1]})$$
$$\rightarrow Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$
$$\rightarrow A^{[2]} = \sigma(Z^{[2]})$$

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix} \quad (n_x, m)$$

training examples

hidden units.

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

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

hidden units

One hidden layer
Neural Network

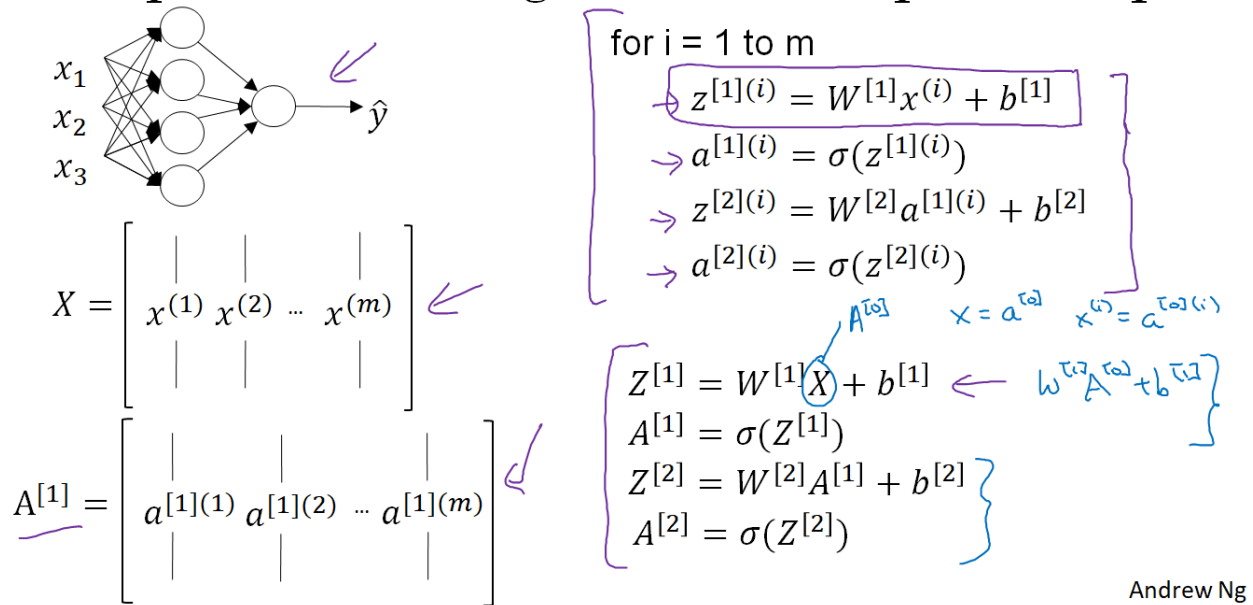Explanation
for vectorized
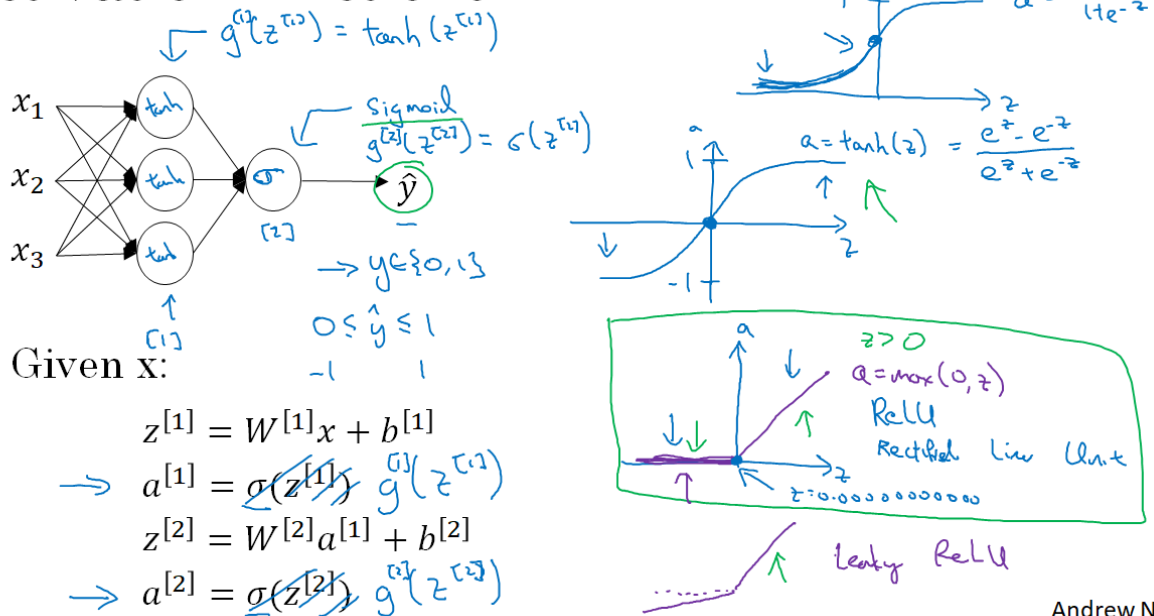implementation

# Justification for vectorized implementation

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

$$W^{[1]} = \begin{bmatrix} \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \end{bmatrix}$$

$$W^{[1]}x^{(1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad W^{[1]}x^{(2)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad W^{[1]}x^{(3)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

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

$$W^{[1]} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \cdots \\ | & | & | \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix} = \begin{bmatrix} | & | & | \\ z^{[1](1)} & z^{[1](2)} & z^{[1](3)} \cdots \\ | & | & | \end{bmatrix} = Z^{[1]}$$

$$X$$

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

$$+b^{[1]} \quad +b^{[1]} \quad +b^{[1]}$$

# Recap of vectorizing across multiple examples

$x_1$
$x_2$
$x_3$
$\hat{y}$

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

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

for i = 1 to m

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

$A^{[0]}$      $x = a^{[0]}$      $x^{(i)} = a^{[0](i)}$

$Z^{[1]} = W^{[1]}X + b^{[1]}$   ←   $W^{[1]}A^{[0]} + b^{[1]}$
$A^{[1]} = \sigma(Z^{[1]})$
$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$
$A^{[2]} = \sigma(Z^{[2]})$

---

# Activation functions

$g^{[1]}(z^{[1]}) = \tanh(z^{[1]})$

$x_1$  tanh
$x_2$  tanh   $\sigma$   $\hat{y}$
$x_3$  tanh

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

[2]

[1]

$\Rightarrow y \in \{0, 1\}$

$0 \le \hat{y} \le 1$

$-1 \qquad 1$

$a = \dfrac{1}{1+e^{-z}}$

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

$-1$

$z > 0$
$a = \max(0, z)$
ReLU
Rectified Linear Unit
$z = 0.0000000000$

Leaky ReLU

Given x:

$z^{[1]} = W^{[1]}x + b^{[1]}$
$a^{[1]} = \sigma(z^{[1]}) \quad g^{[1]}(z^{[1]})$
$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$
$a^{[2]} = \sigma(z^{[2]}) \quad g^{[1]}(z^{[2]})$

# Pros and cons of activation functions



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

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

ReLU $\quad a = \max(0, z)$

Leaky ReLU $\quad a = \max(0.01z, z)$

$\max(0.01z, z)$

---

One hidden layer
Neural Network

deeplearning.ai

Why do you
need non-linear
activation functions?

---

# Activation function

ReLU, tanh

$x_1$

$x_2 \quad \rightarrow \hat{y} \in \mathbb{R}$

$x_3$

ReLU

$y \in \mathbb{R}$

$\hat{y} \geq 0.$

$\$0 \quad \cdots \quad \$1,000,000s$

Given x:

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

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

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

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

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

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

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

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

$\underbrace{\qquad}_{a^{[1]}}$

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

$= W'x + b'$

$g(z) = z$

deeplearning.ai

# Sigmoid activation function



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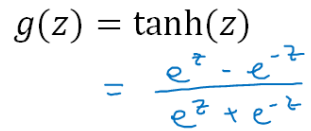
$$a = g(z) = \frac{1}{1 + e^{-z}}$$

$$\boxed{g'(z) = }\boxed{\frac{d}{dz} g(z)} = \text{slope of } g(x) \text{ at } z$$

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

$$= g(z)\left(1 - g(z)\right) \leftarrow$$

$$= \boxed{a(1-a)} \qquad \Big| g'(z) = a(1-a)$$

$z = 10. \quad g(z) \approx 1$

$\frac{d}{dz} g(z) \approx 1(1-1) \approx 0$

$z = -10 \quad g(z) \approx 0$

$\frac{d}{dz} g(z) \approx 0 \cdot (1-0) \approx 0$

$z = 0 \qquad g(z) = \frac{1}{2}$

$\frac{d}{dz} g(z) = \frac{1}{2}(1 - \frac{1}{2}) = \frac{1}{4}$

Andrew Ng

# Tanh activation function



$$g(z) = \tanh(z)$$

$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z$$

$$= 1 - (\tanh(z))^2 \leftarrow$$

$$a = g(z), \qquad g'(z) = 1 - a^2$$

$z = 10 \quad \tanh(z) \approx 1$
$\qquad\qquad g'(z) \approx 0$
$z = -10 \quad \tanh(z) \approx -1$
$\qquad\qquad g'(z) \approx 0$
$z = 0 \quad \tanh(z) = 0$
$\qquad\qquad g'(z) = 1$

# ReLU and Leaky ReLU



## ReLU

$$g(z) = \max(0, z)$$

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

$z = 0.0000\cdots 0$

## Leaky ReLU

$$g(z) = \max(0.01z, z)$$

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

# Gradient descent for neural networks

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
$(n^{[1]}, n^{[0]}) \quad (n^{[1]}, 1) \quad (n^{[2]}, n^{[1]}) \quad (n^{[2]}, 1)$

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

Cost function: $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^{n} \mathcal{L}(\hat{y}, y) \quad \leftarrow a^{[2]}$

Gradient descent:

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

$\quad dW^{[1]} = \frac{\partial J}{\partial W^{[1]}}, \quad db^{[1]} = \frac{\partial J}{\partial b^{[1]}}, \ \dots$

$\quad W^{[1]} := W^{[1]} - \alpha \, dW^{[1]}$

$\quad b^{[1]} := b^{[1]} - \alpha \, db^{[1]}$

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

# Formulas for computing derivatives

Forward propagation:

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

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

$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$

$A^{[2]} = g^{[2]}(Z^{[2]}) = \sigma(Z^{[2]})$

Back propagation:

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

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

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

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

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

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

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

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

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

One hidden layer
Neural Network

Backpropagation
intuition (Optional)

# Computing gradients

### Logistic regression

$$z = w^T x + b \quad\longrightarrow\quad a = \sigma(z) \quad\longrightarrow\quad \mathcal{L}(a,y)$$

$x$
$w$
$b$

$dz = a - y$

$dw = dz \cdot x$
$db = dz$

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

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

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

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

$\hat{y} = a$

$da = \frac{d}{da}\mathcal{L}(a,y) = -y\log a - (1-y)\log(1-a)$

$= -\frac{y}{a} + \frac{1-y}{1-a}$

Andrew Ng

# Neural network gradients

$$z^{[1]} = W^{[1]}x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = W^{[2]}x + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow \mathcal{L}(a^{[2]}, y)$$

$x$
$W^{[1]}$
$dw^{[1]}$
$b^{[1]}$
$db^{[1]}$

$W^{[2]}$   $dw^{[2]}$
$b^{[2]}$   $db^{[2]}$

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

$dz^{[1]} = W^{[2]T} dz^{[2]}$

$da^{[1]} = \cdots$

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

$da^{[2]}$

$* g^{[1]'}(z^{[1]})$
element wise product

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

$z^{[2]}, dz^{[2]} \quad (n^{[2]}, 1) - (1,1)$

$z^{[1]}, dz^{[1]} \quad (n^{[1]}, 1)$

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

$(n^{[1]}, 1) = (n^{[1]}, n^{[2]}) (n^{[2]}, 1) * (n^{[1]}, 1)$

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

"$dw = dz \cdot x$"

$w : [\quad\quad\quad]$

$dw^{[1]} = dz^{[1]} \cdot x^T$
$db^{[1]} = dz^{[1]}$

foo
dfoo

$W^{[1]} \quad dW^{[1]}$
$dW^{[2]}$

$a^{[1]T}$

Andrew Ng

# Summary of gradient descent

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

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

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

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

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

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

Vectorized Implementation:

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

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

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

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

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

---

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

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

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

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

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

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

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

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

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

element-wise product

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

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

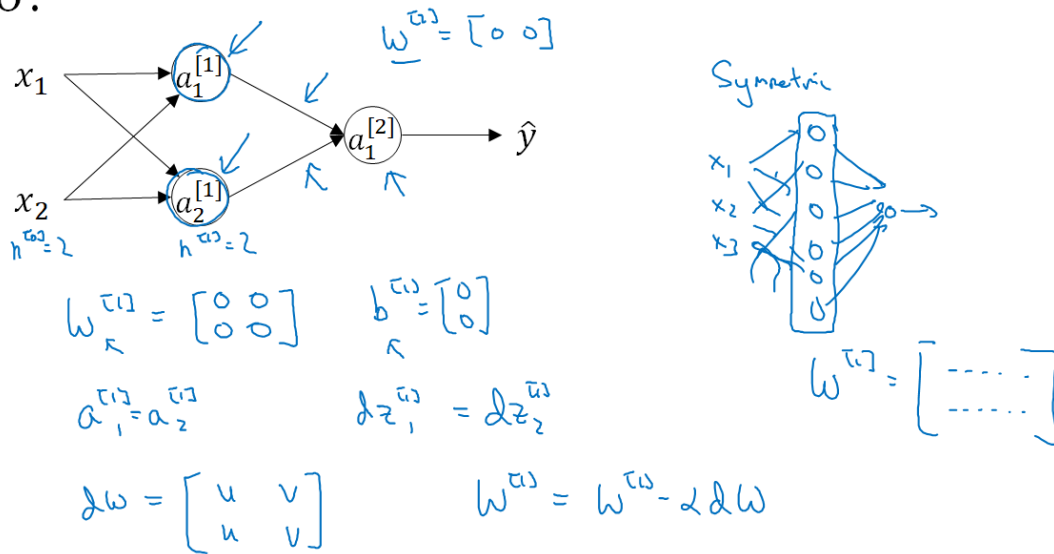$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$$

$$J(\cdots) = \frac{1}{m} \sum_{i=1}^{n} \mathcal{L}(\hat{y}, y)$$

Andrew Ng

# What happens if you initialize weights to zero?



$$W^{[2]} = [0 \ 0]$$

Symmetric

$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \qquad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]} \qquad dz_1^{[1]} = dz_2^{[1]}$$

$$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \qquad W^{[1]} = W^{[1]} - \alpha \, dW$$

$$W^{[1]} = \begin{bmatrix} - - - - - \\ - - - - - \end{bmatrix}$$

# Random initialization



$$\rightarrow W^{[1]} = np.random.randn((2,2)) * 0.01$$
$$100 ?$$

$$b^{[1]} = np.zeros((2,1))$$
$$W^{[2]} = np.random.randn((1,2)) * 0.01$$
$$b^{[2]} = 0$$

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