# Neural Network Representation

$a^{[0]} = X$

$a^{[1]}$

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

"2 layer NN"

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

$x_1$

$a_1^{[1]}$

$x_2$

$a_2^{[1]}$

$a_3^{[1]}$

$a^{[2]}$

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

$x_3$

$a_4^{[1]}$

"$\hat{y} = a$"

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

→ Input layer

→ Hidden layer

→ Output

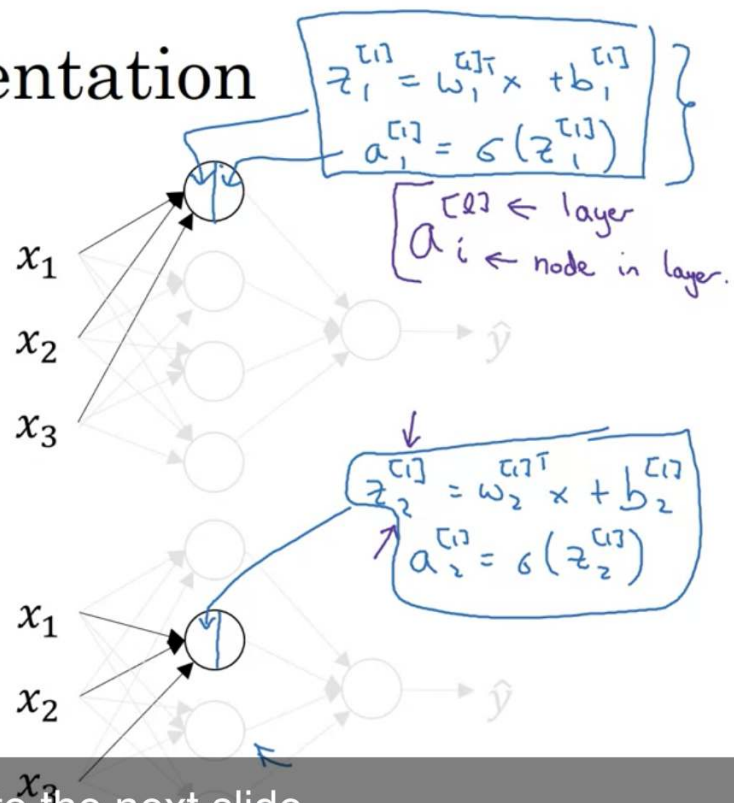Let's go more deeply into exactly what this neural network computes.

Andrew Ng

# Neural Network Representation

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

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

$x_1$

$x_2$ —→ $\boxed{w^T x + b \mid \sigma(z)}$ —→ $a = \hat{y}$

$x_3$

$z = w^T x + b$

$a = \sigma(z)$

$x_1$

$x_2$

$x_3$

$\hat{y}$

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

$x_1$

$x_2$

$x_3$

$\hat{y}$

and let's copy them to the next slide.

Andrew Ng

# Neural Network Representation

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



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

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

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

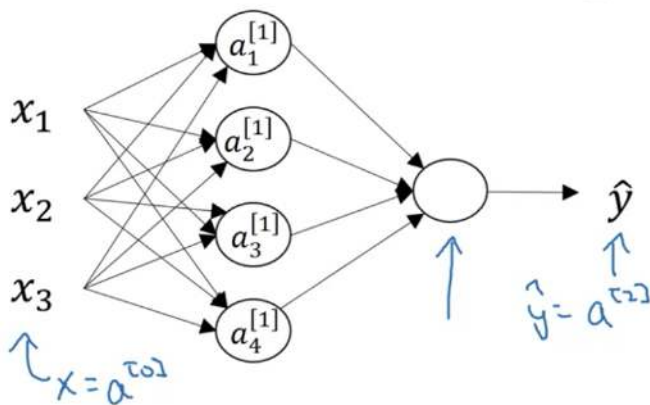$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

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

$(4,3)$

which is taken by stacking up these individuals of z's into a column vector.

Andrew Ng

# Neural Network Representation learning



Given input x:

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

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

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

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

In the diagram: $x_1$, $x_2$, $x_3$ with $x = a^{[0]}$; hidden layer $a^{[1]}_1$, $a^{[1]}_2$, $a^{[1]}_3$, $a^{[1]}_4$; output $\hat{y}$ with $\hat{y} = a^{[2]}$.

also be written similarly where what the output layer does is,

# Vectorizing across multiple examples



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

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

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

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

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

$a^{[2](i)} \nwarrow$ example $i$
layer 2

for $i = 1$ to $m$,

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

Andrew Ng

# Vectorizing across multiple examples

```
for i = 1 to m:
```

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

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

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

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

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

$$(n_x, m)$$

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

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

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

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

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

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

6:40 / 9:05

Andrew Ng

# Vectorizing across multiple examples

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

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

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

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

# Justification for vectorized implementation

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

$$w^{[1]} = \begin{bmatrix} \rule[.5ex]{2em}{0.4pt} \\ \rule[.5ex]{2em}{0.4pt} \\ \rule[.5ex]{2em}{0.4pt} \\ \rule[.5ex]{2em}{0.4pt} \end{bmatrix} \qquad w^{[1]} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \vdots \\ \cdot \end{bmatrix} \qquad w^{[1]} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \vdots \\ \cdot \end{bmatrix} \qquad w^{[1]} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

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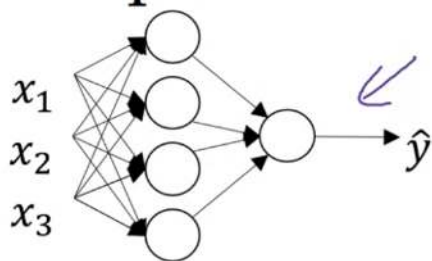\cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} | & | & | \\ z^{[1](1)} & z^{[1](2)} & z^{[1](3)} \cdots \\ | & | & | \end{bmatrix} = z^{[1]}$$

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

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

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

Andrew Ng

# Recap of vectorizing across multiple examples



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

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

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

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

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

# Activation functions

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

$x_1$

$x_2$

$x_3$

$\tanh$

$\tanh$

$\tanh$

$\sigma$

$\hat{y}$

Sigmoid

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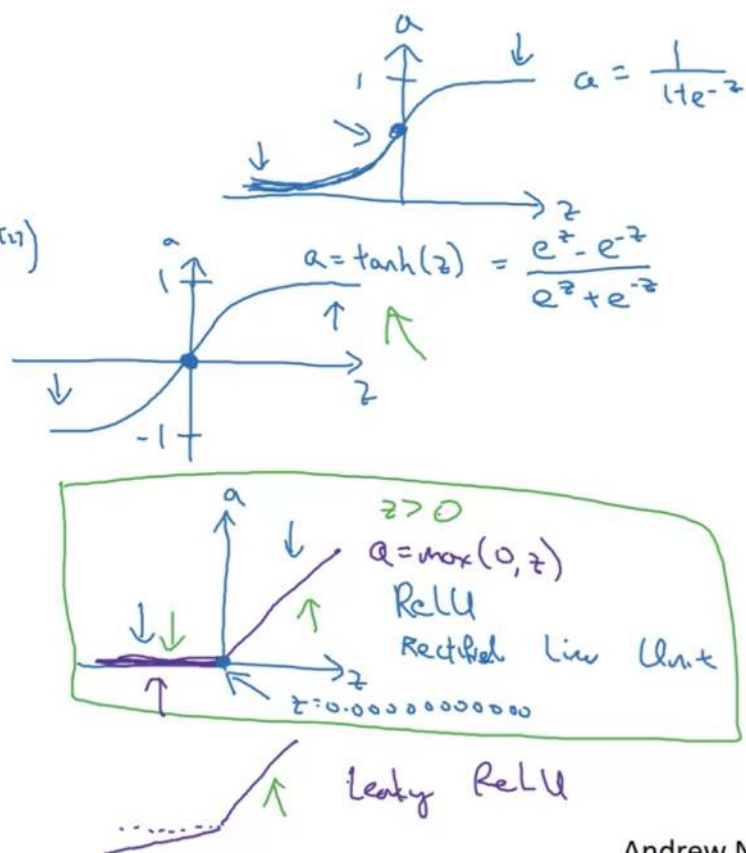
[2]

[1]

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

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

$-1 \qquad 1$
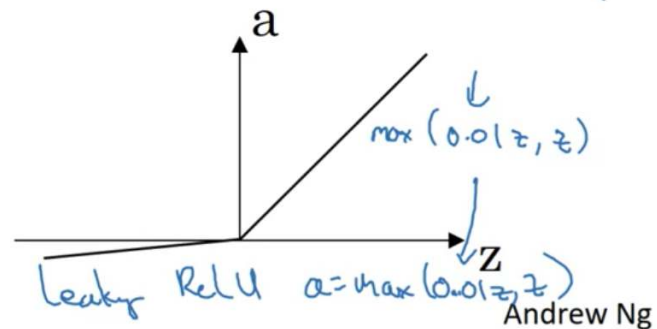
Given x:
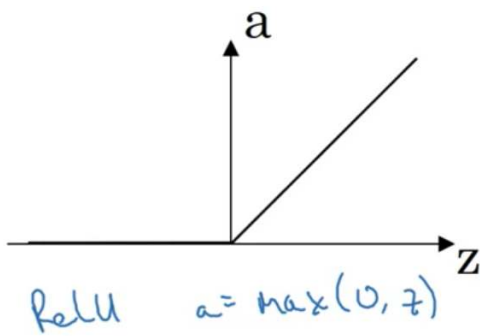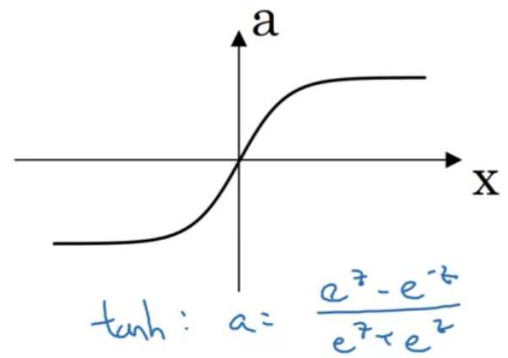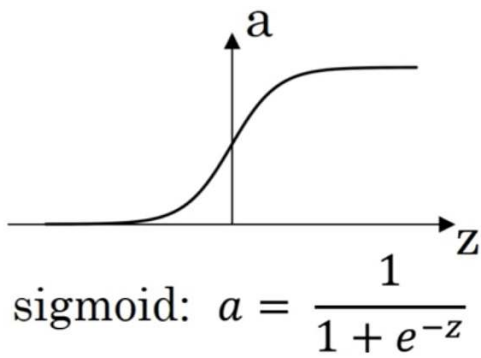
$$z^{[1]} = W^{[1]}x + b^{[1]}$$

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

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

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

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

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

$z > 0$

$a = \max(0, z)$

ReLU

Rectified Linu Unit

$z = 0.0000000000$

Leaty ReLU

Andrew Ng

# 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 $\quad a = \max(0.01z, z)$

$\max(0.01z, z)$

Andrew Ng

# Sigmoid activation function



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

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

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

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

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

$$= \boxed{a\,(1-a)}$$

$$g'(z) = a(1-a)$$
$$\uparrow$$

$z = 10. \quad g(z) \approx 1$
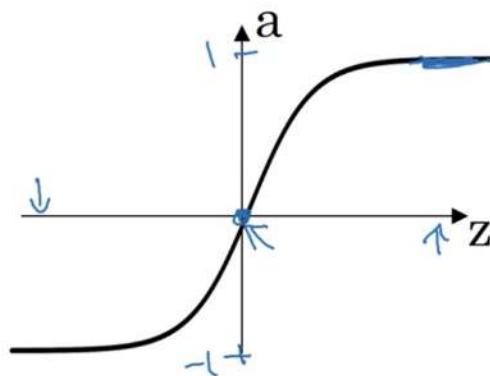
$\frac{d}{dz} g(z) \approx 1\,(1-1) \approx 0$

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

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

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

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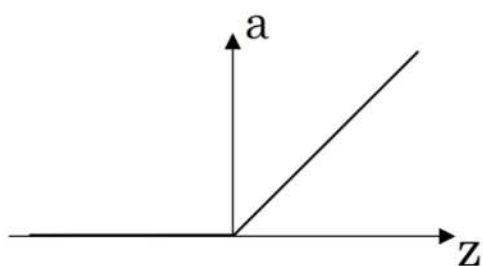
Andrew Ng

# Tanh activation function



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

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

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

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

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

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

# ReLU and Leaky ReLU
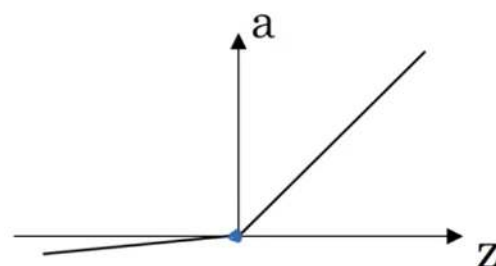


ReLU

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

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

~~undef if z=0~~

$$z = 0.0000\cdots 0$$



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: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$

$(n^{[1]}, n^{[0]})$ $(n^{[1]}, 1)$ $(n^{[2]}, n^{[1]})$ $(n^{[2]}, 1)$

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

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

$\uparrow$ $\uparrow$ $\curvearrowleft a^{[2]}$

Gradient descent:

$\rightarrow$ Repeat {

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

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

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

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

# 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:**

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

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

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

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

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

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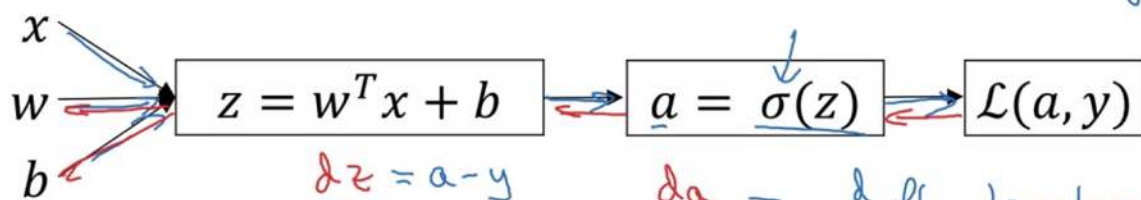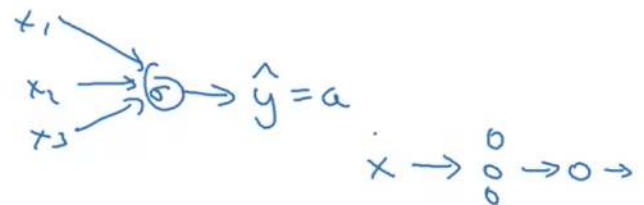
$$(n^{[1]}, 1) \qquad (n^{[1]}, )$$

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

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

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

# Computing gradients

## Logistic regression



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

$dw = dz \cdot x$

$db = dz$

$dz = a - y$

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

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

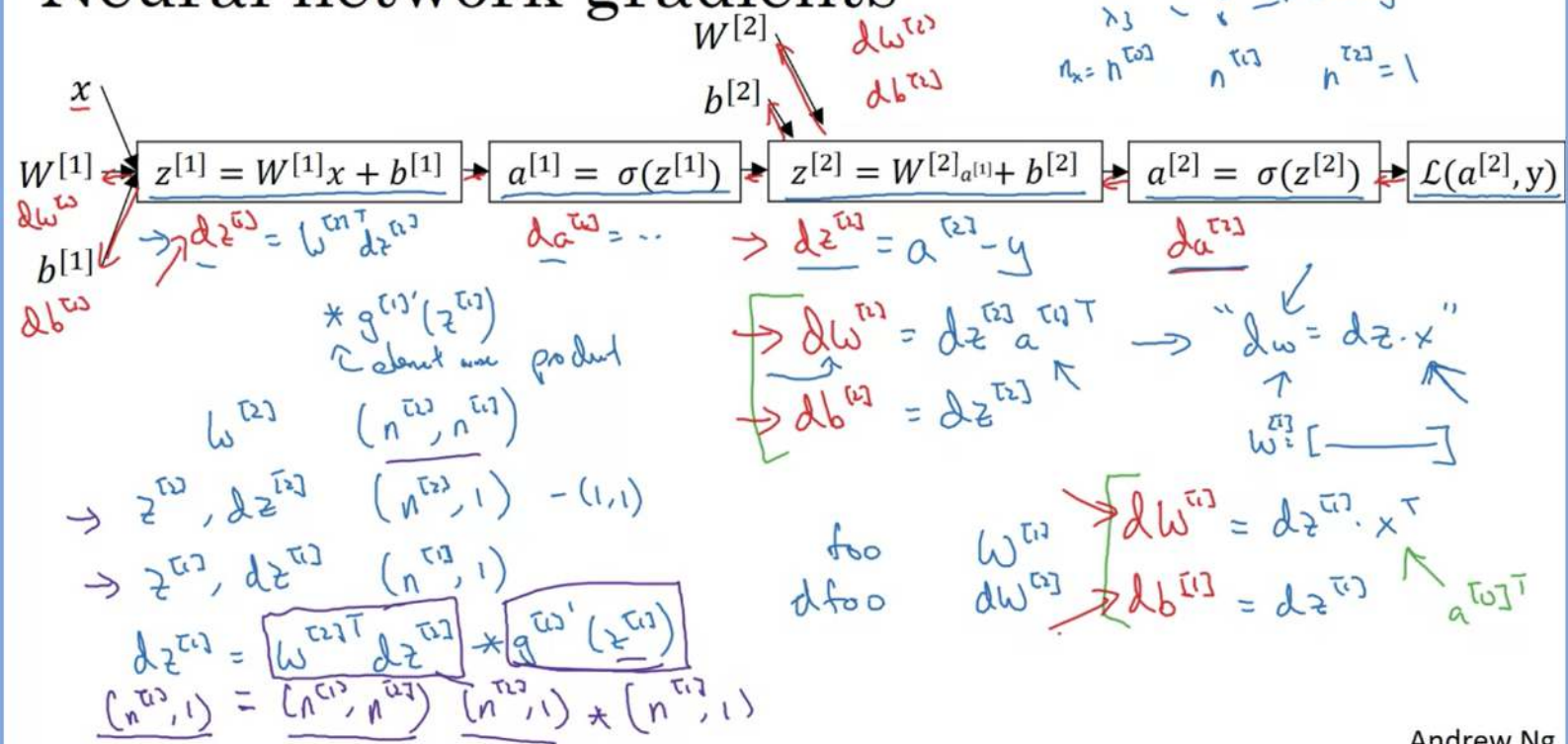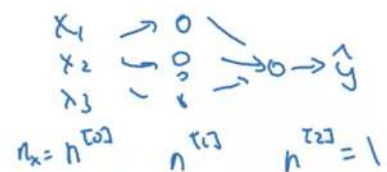$$\frac{\partial \ell}{\partial z} = \frac{\partial \ell}{\partial a} \cdot \frac{da}{dz}$$

$"dz" = "da"$

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

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

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

$\hat{y} = a$

# Neural network gradients

$x_1 \searrow 0 \searrow$
$x_2 \searrow 0 \rightarrow 0 \rightarrow \hat{y}$
$x_3 \nearrow 0 \nearrow$

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

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

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

$x$

$W^{[1]}$
$dw^{[1]}$

$b^{[1]}$
$db^{[1]}$

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

$\rightarrow dz^{[1]} = W^{[1]T} dz^{[1]}$

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

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

$da^{[2]}$

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

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

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

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

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

$\rightarrow dW^{[2]} = dz^{[2]} a^{[1]T}$  $\rightarrow$ "$dw = dz \cdot x$"
$\rightarrow db^{[2]} = dz^{[2]}$

$W^{[1]}: [\quad\quad]$

foo         $W^{[1]}$  $\rightarrow dW^{[1]} = dz^{[1]} \cdot x^T$
dfoo        $dW^{[1]}$  $\rightarrow db^{[1]} = dz^{[1]}$   $a^{[0]T}$

Andrew Ng

# Summary of gradient descent

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

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

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

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

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

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

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

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

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

elementwise product

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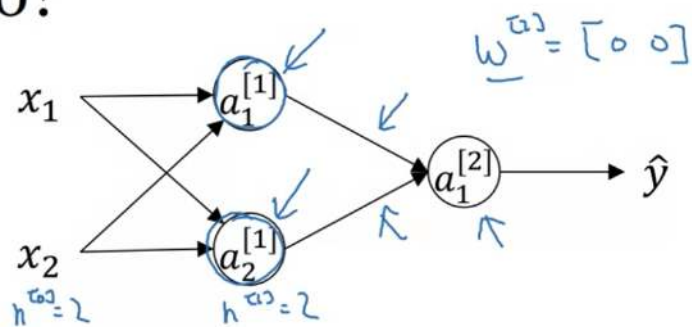
$$dW^{[1]} = \frac{1}{m}dZ^{[1]}X^T$$

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

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

Andrew Ng

# What happens if you initialize weights to zero?

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



$x_1$

$a_1^{[1]}$

$x_2$

$a_2^{[1]}$

$a_1^{[2]}$

$\hat{y}$

$n^{[0]} = 2$

$n^{[1]} = 2$

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

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

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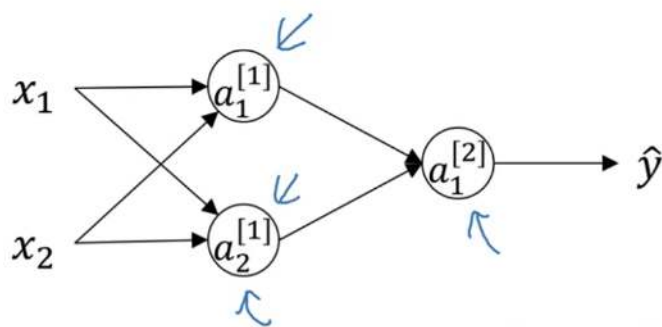
Symmetric

$x_1$

$x_2$

$x_3$

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

# Random initialization



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

$$b^{[1]} = np.zeros((2,1))$$

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

$$b^{[2]} = 0$$

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

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