

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>

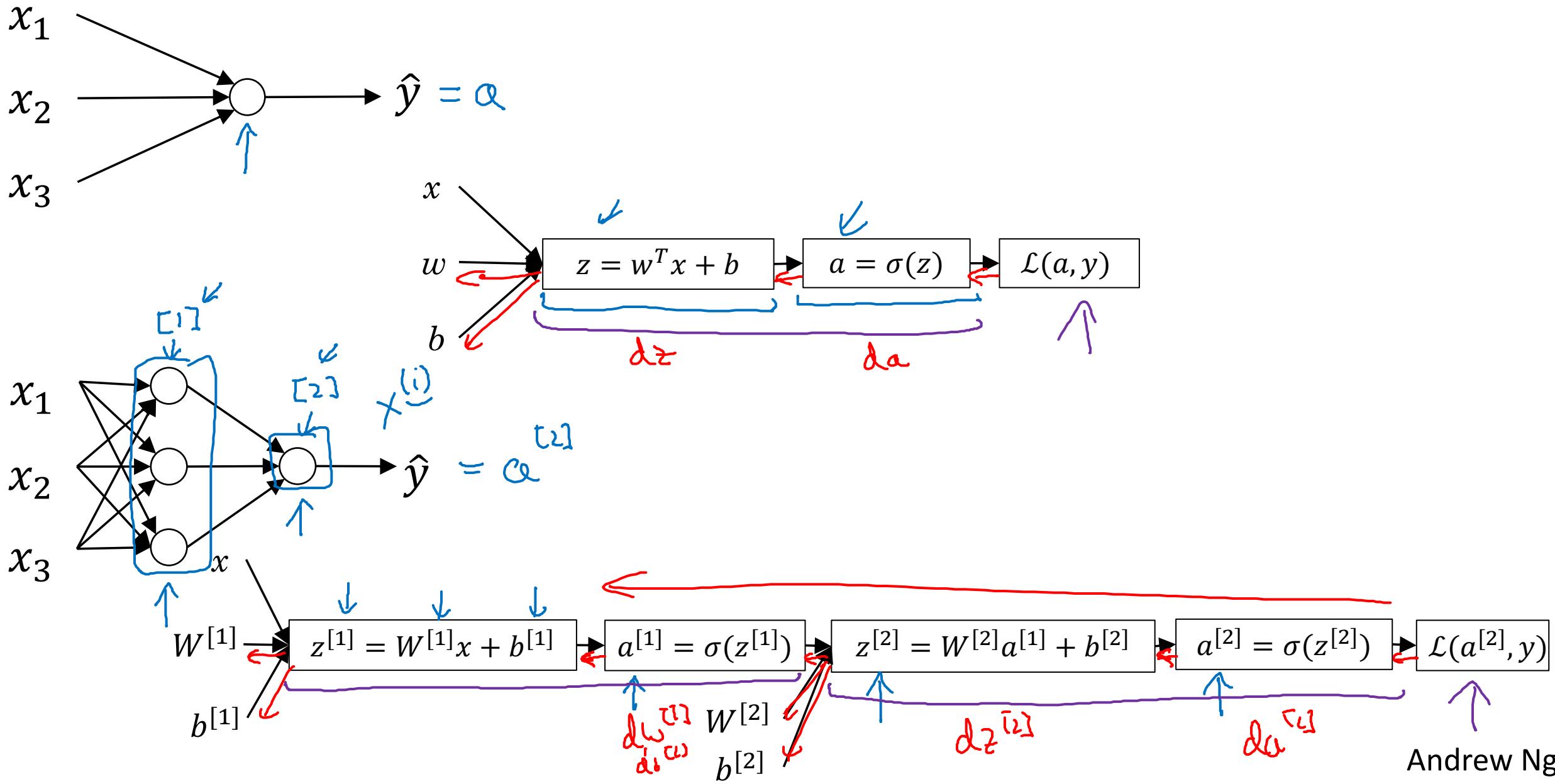


deeplearning.ai

One hidden layer
Neural Network

Neural Networks
Overview

What is a Neural Network?



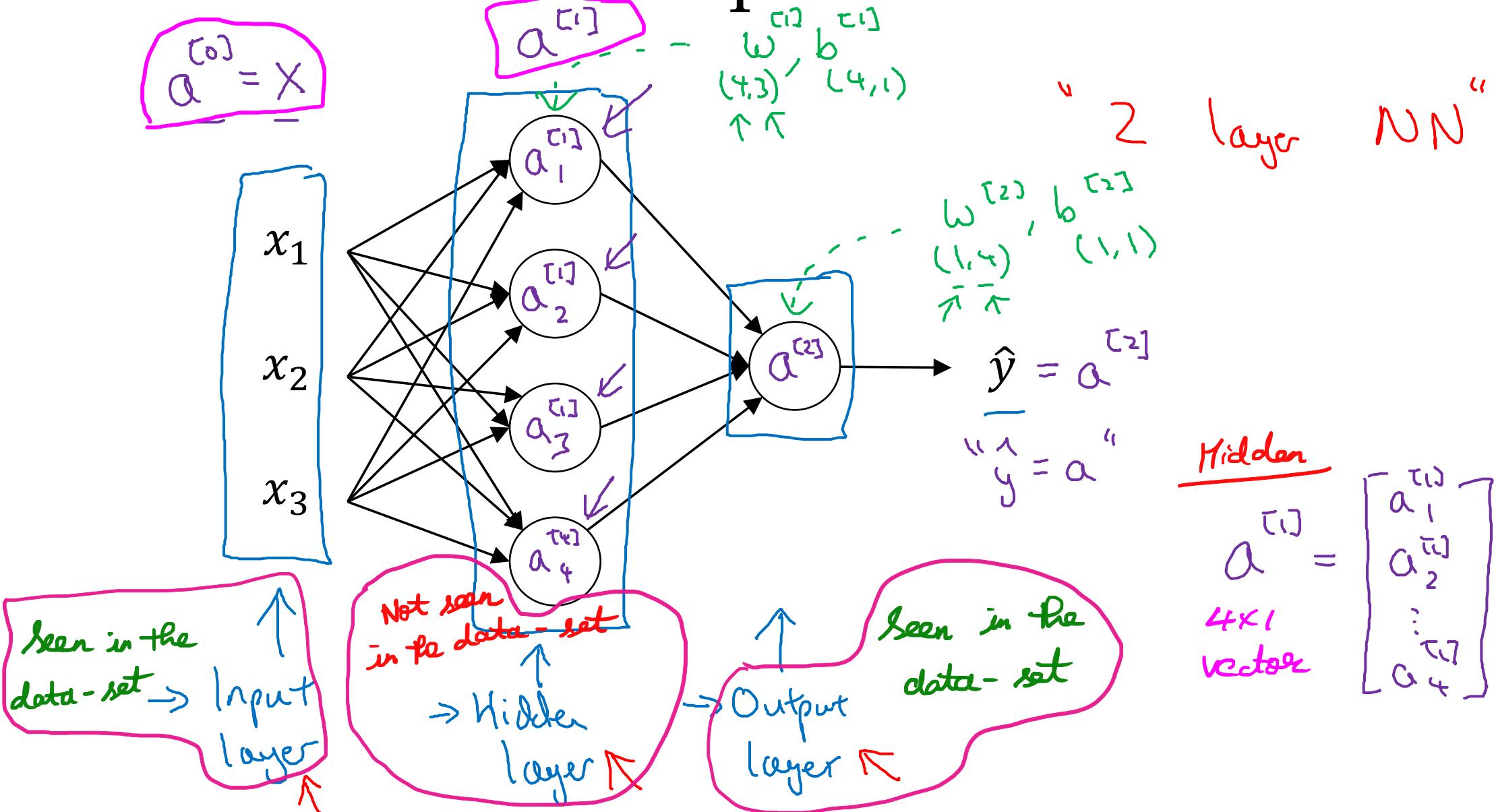


deeplearning.ai

One hidden layer
Neural Network

Neural Network
Representation

Neural Network Representation



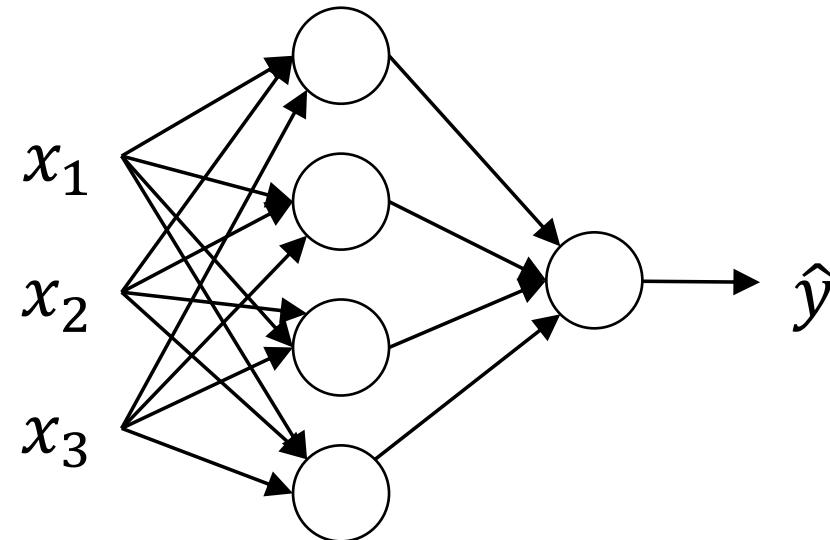
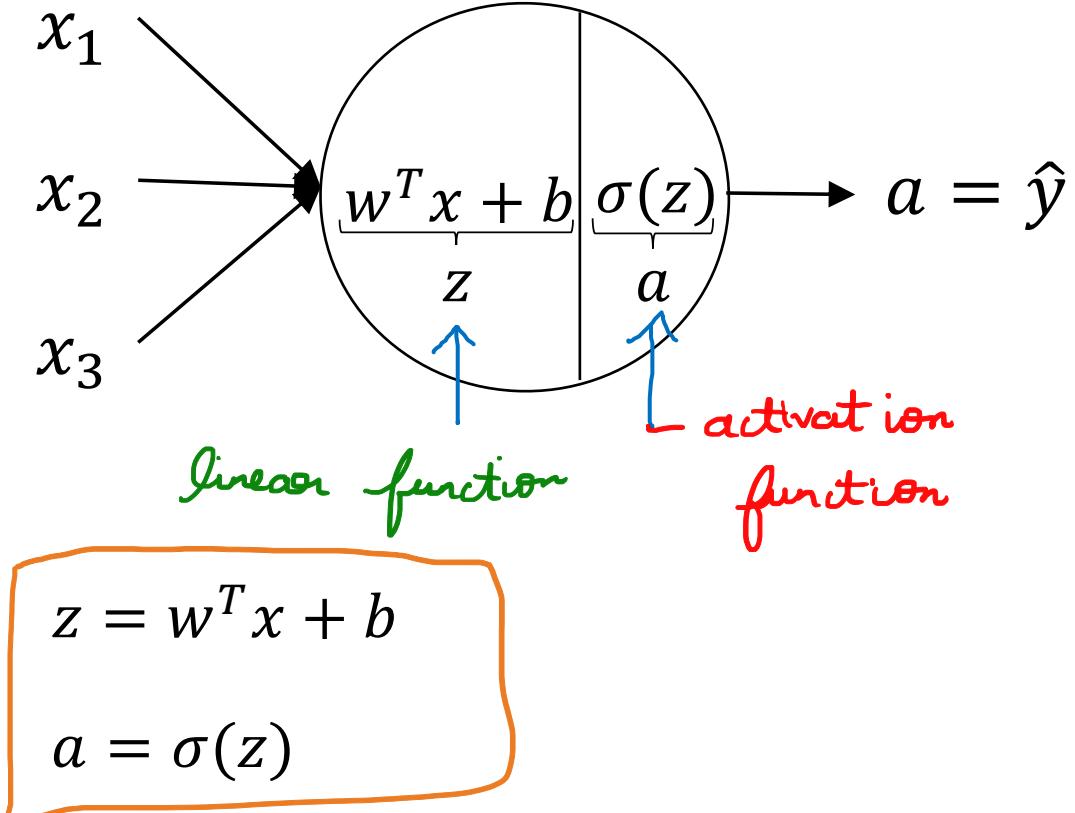


deeplearning.ai

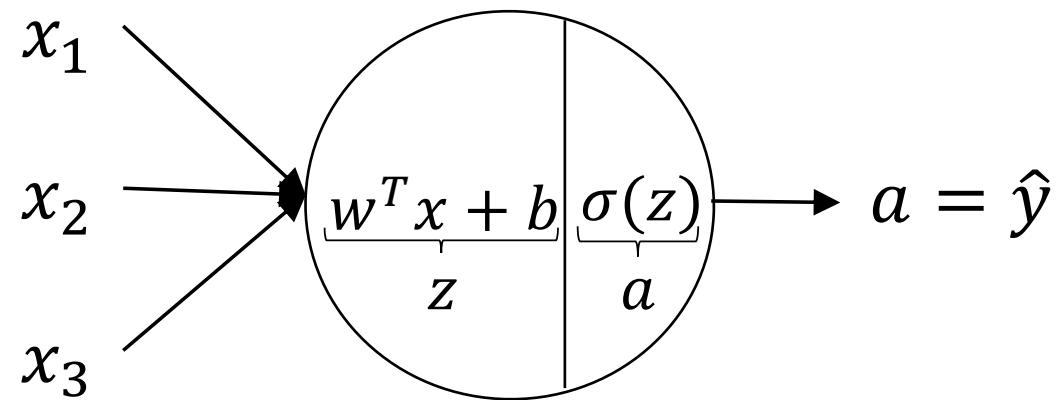
One hidden layer Neural Network

Computing a Neural Network's Output

Neural Network Representation

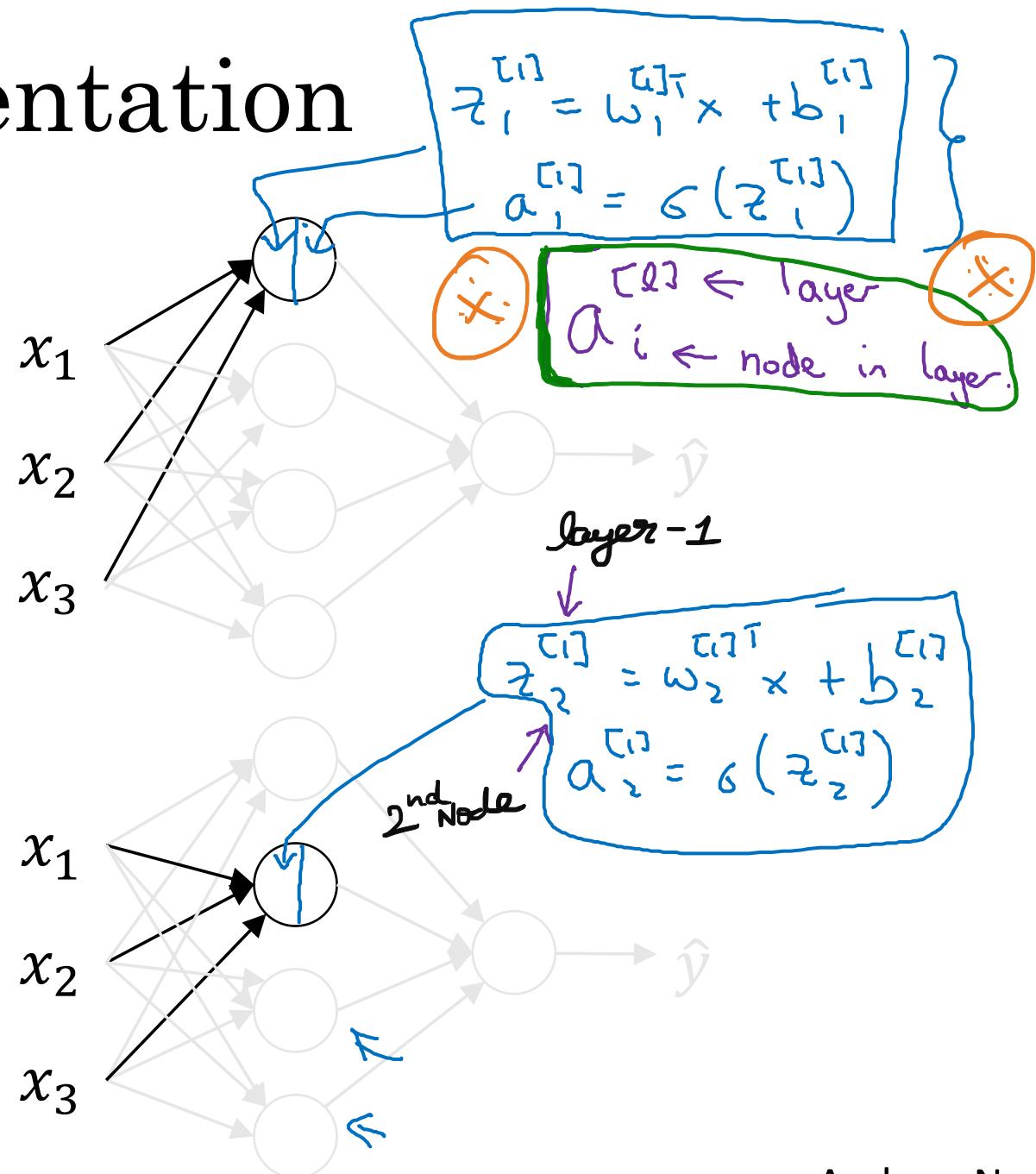


Neural Network Representation

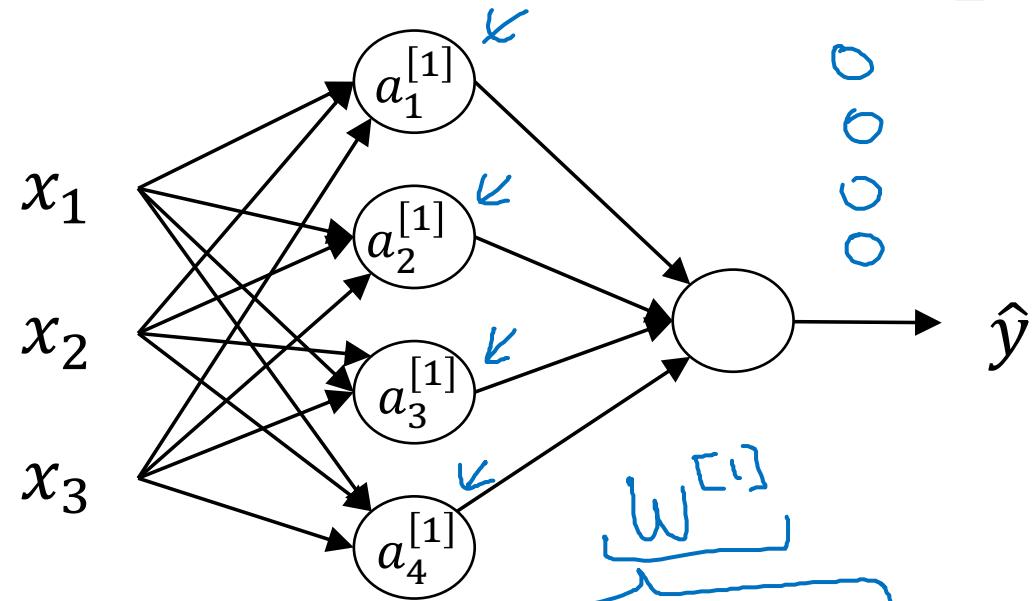


$$z = w^T x + b$$

$$a = \sigma(z)$$



Neural Network Representation



$$\rightarrow z^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$$

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

Diagram illustrating the mathematical representation of the neural network layers:

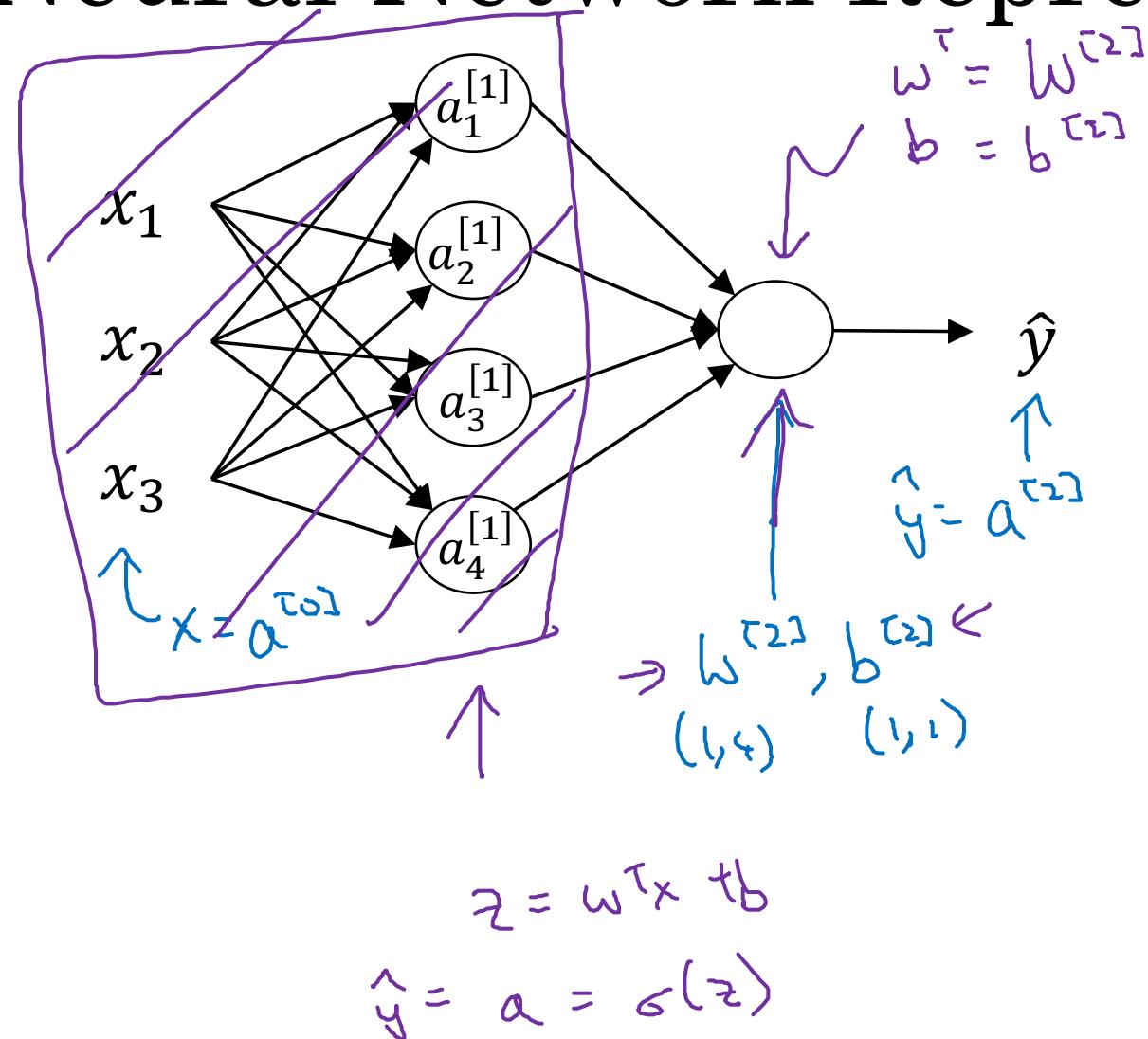
- Layer 1: $z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$, $a_1^{[1]} = \sigma(z_1^{[1]})$
- Layer 2: $z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$, $a_2^{[1]} = \sigma(z_2^{[1]})$
- Layer 3: $z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}$, $a_3^{[1]} = \sigma(z_3^{[1]})$
- Layer 4: $z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}$, $a_4^{[1]} = \sigma(z_4^{[1]})$

Annotations:

- $(w_i^{[1]T} x)$ is highlighted in blue.
- $a^{[1]}$ is highlighted in red.
- $\sigma(z^{[1]})$ is highlighted in red.
- $z^{[1]}$ is highlighted in purple.
- $b^{[1]}$ is highlighted in green.
- $w^{[1]}$ is highlighted in blue.

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

Neural Network Representation learning



Given input x :

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

$(4, 1) \quad (4, 3) \quad (3, 1) \quad (4, 1)$

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

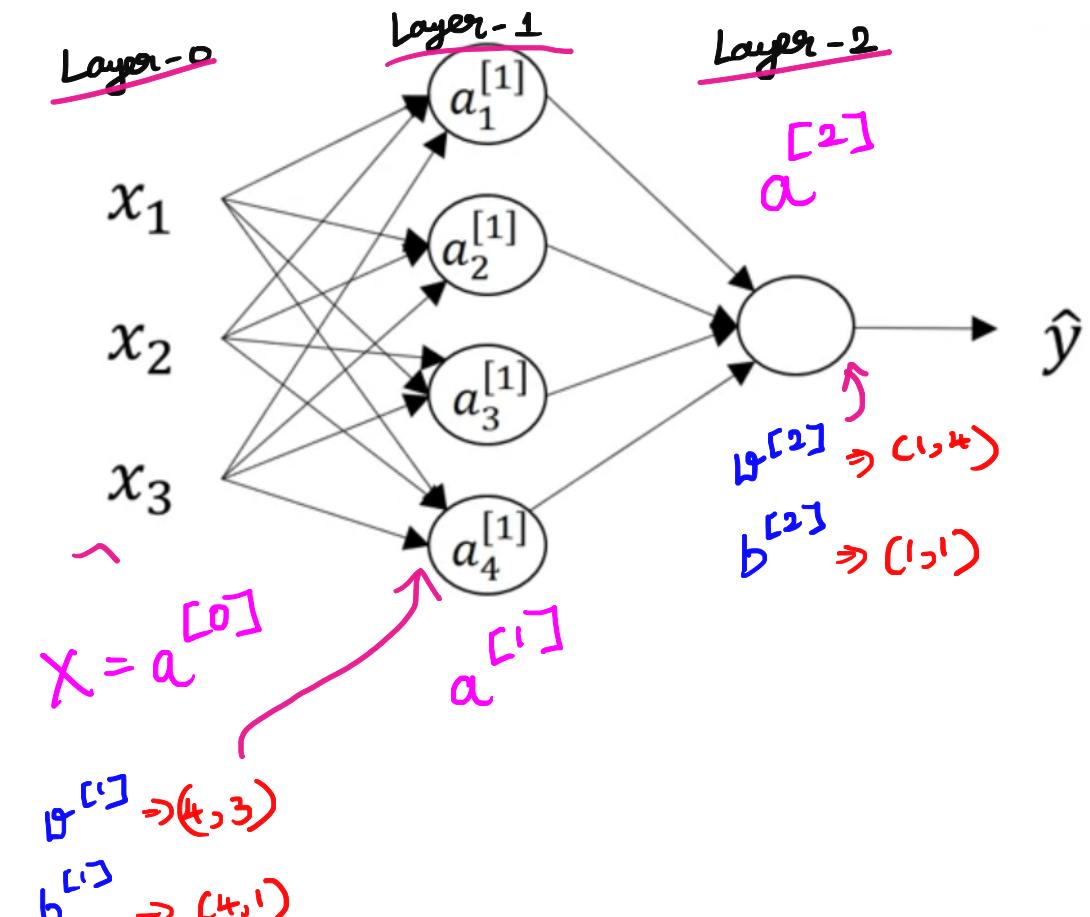
$(4, 1) \quad (4, 1)$

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

$(1, 1) \quad (1, 4) \quad (4, 1) \quad (1, 1)$

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

$(1, 1) \quad (1, 1)$



- 1) Layer-0 (Input features) along with $w^{[1]}$ and $b^{[1]}$ gives layer-1.
- 2) Layer-1 along $w^{[2]}$ and $b^{[2]}$ gives Layer-2.

$$1) z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

$(4 \times 1) \quad \underbrace{(4 \times 3)(3 \times 1)}_{(4 \times 1)} \quad (4 \times 1)$

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

$(4 \times 1) \quad (4 \times 1)$

$$3) z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$(1 \times 1) \quad \underbrace{(1 \times 4)(4 \times 1)}_{(1 \times 1)} \quad (1 \times 1)$

$$4) a^{[2]} = \sigma(z^{[2]})$$

$(1 \times 1) \quad (1 \times 1)$

$$z \Rightarrow w^T x + b$$

$$\hat{y} \Rightarrow a \Rightarrow \sigma(z)$$

Layer-1
(Hidden layer)

Layer-2
(Output layer)

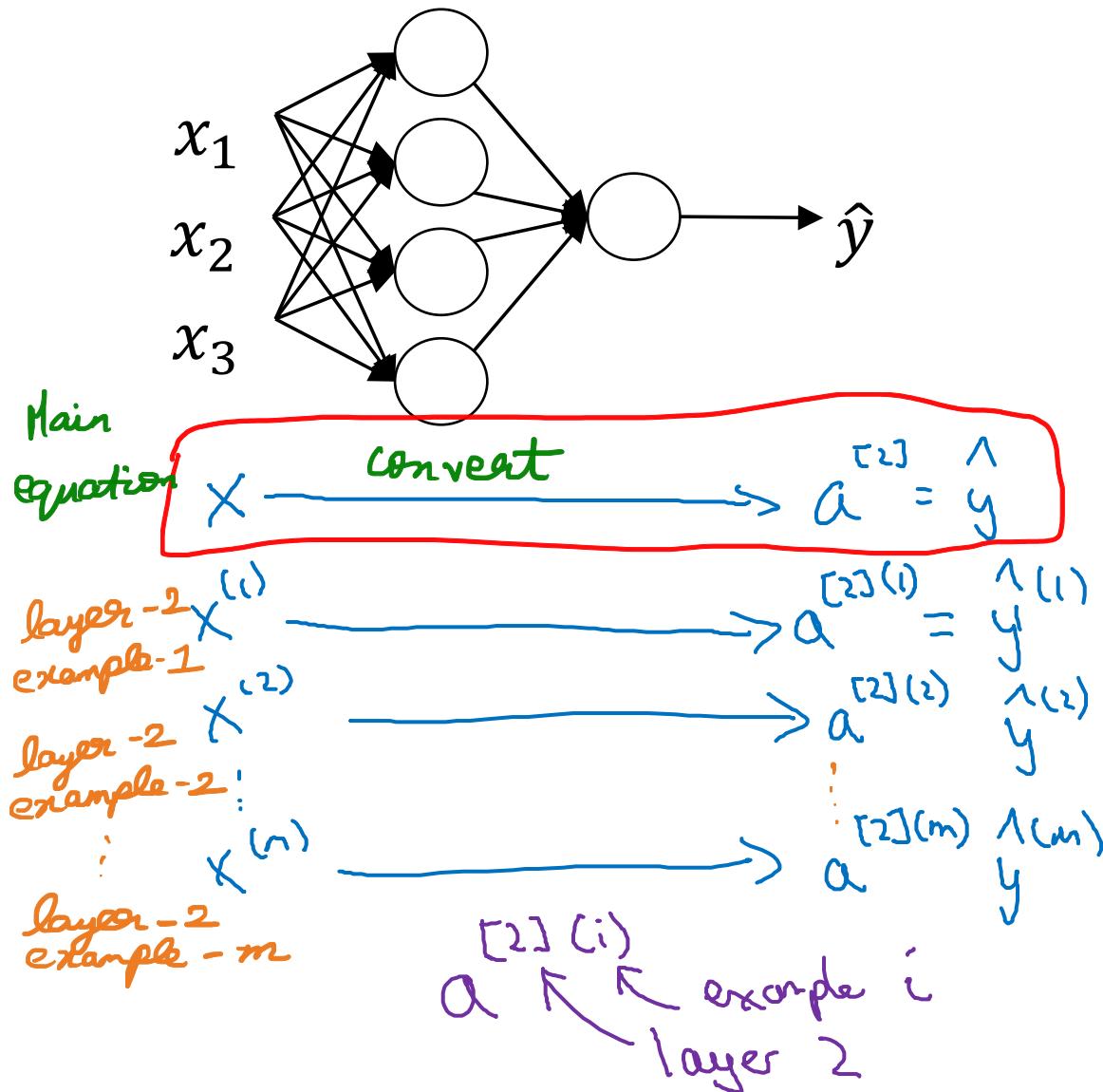


deeplearning.ai

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 m {

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

layer-2 {

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

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

Linear function } 1st hidden unit of
function } 1st train example

Input train example

Independent variables

$$X \Rightarrow \begin{bmatrix} \vdots & \vdots & \vdots \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \quad (n_x, m)$$

$$Z \Rightarrow \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ Z^{(1)} & Z^{(2)} & \dots & Z^{(m)} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

13-10-2021

hidden units

$$A \Rightarrow \begin{bmatrix} \vdots & \vdots & \vdots \\ a^{(1)} & a^{(2)} & \dots & a^{(m)} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

Horizontal Index \Rightarrow training examples.

Vertical Index \Rightarrow Represents hidden units
in the training examples.

$$\begin{bmatrix} \vdots & \vdots & \vdots \\ a^{(1)} & a^{(2)} & \dots & a^{(m)} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

$$\begin{bmatrix} \vdots & \vdots & \vdots \\ a^{(1)} & a^{(2)} & \dots & a^{(m)} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

Activation function } 2nd hidden of
function } 1st train example

Vectorizing across multiple examples

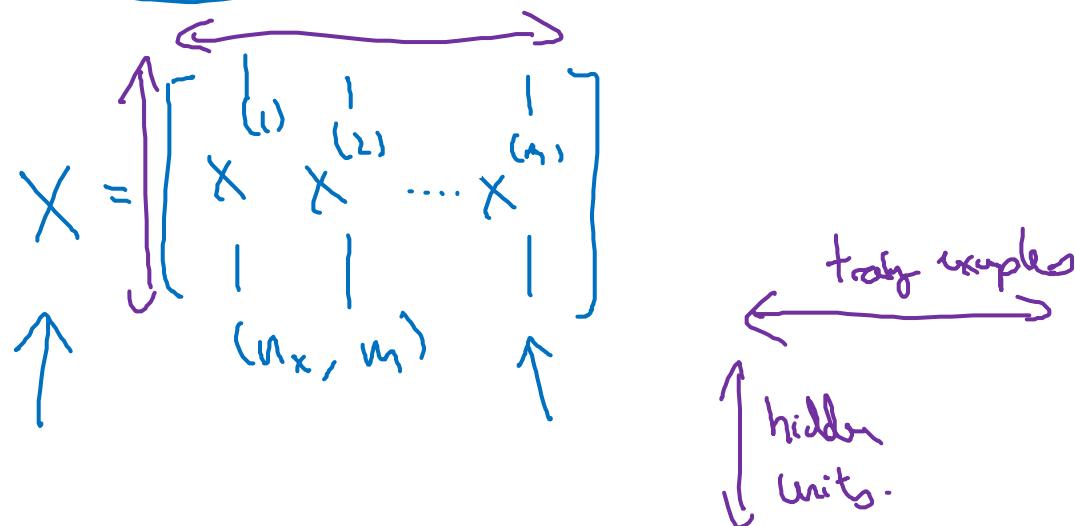
for $i = 1$ to m :

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

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

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

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

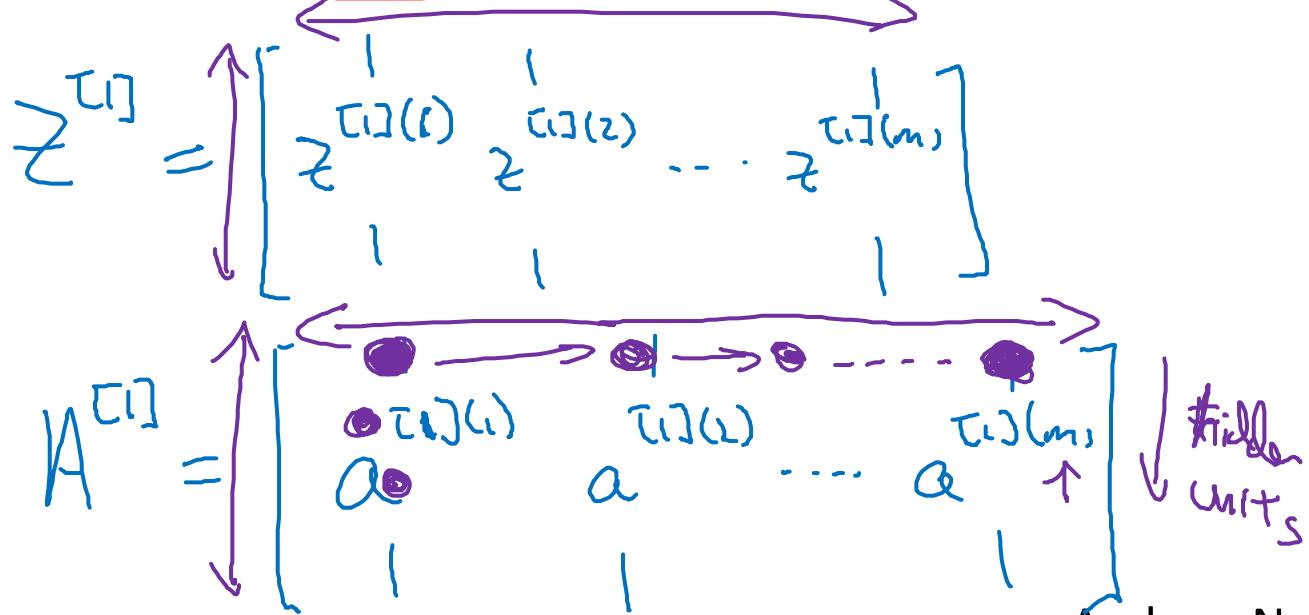


$$z^{[1]} = W^{[1]}X + b^{[1]}$$

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

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

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





deeplearning.ai

One hidden layer Neural Network

Explanation for vectorized implementation

Justification for vectorized implementation

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

$$z^{(i)(2)} = w^{(2)} x^{(2)} + b^{(2)}$$

$$z^{(i)(3)} = w^{(3)} x^{(3)} + b^{(3)}$$

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

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

*1st feature
all x's*

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

*2nd feature
all x's*

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

*3rd feature
all x's*

Take $w^{(1)}$ as common

$$w^{(1)}$$

$$\begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & x^{(3)} \dots \end{bmatrix}$$

$$\Rightarrow$$

$$\begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$\begin{bmatrix} z^{(1)(1)} \\ z^{(1)(2)} \\ z^{(1)(3)} \dots \end{bmatrix}$$

$$\begin{bmatrix} 1 & 1 & 1 \\ z^{(1)(1)} & z^{(1)(2)} & z^{(1)(3)} \dots \end{bmatrix}$$

$$z^{(1)}$$

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

~~Σ~~

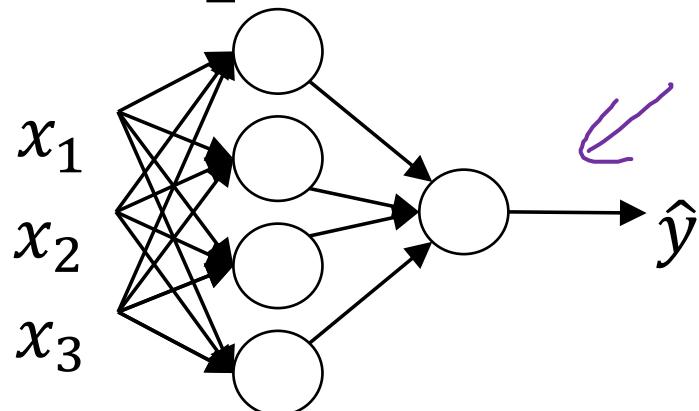
$$w^{(1)} x^{(1)} = z^{(1)(1)}$$

$$+ b^{(1)}$$

$$+ b^{(2)}$$

$$+ b^{(3)}$$

Recap of vectorizing across multiple examples



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

A vertical stack of input vectors $x^{(1)}, x^{(2)}, \dots, x^{(m)}$. A purple arrow points from the text "vectorizing across multiple examples" down to this equation.

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

A vertical stack of hidden layer activations $a^{1}, a^{[1](2)}, \dots, a^{[1](m)}$. A purple arrow points from the text "vectorizing across multiple examples" down to this equation.

```

for i = 1 to m
    z[1](i) = W[1]x(i) + b[1]
    → a[1](i) = σ(z[1](i))
    → z[2](i) = W[2]a[1](i) + b[2]
    → a[2](i) = σ(z[2](i))

```

Handwritten annotations: $x = a^{[1]}$ and $x^{(i)} = a^{[1](i)}$ are written next to the equations for $a^{[1]}$ and $a^{[2]}$ respectively. Brackets group the first two equations and the last two equations.

```

Z[1] = W[1]X + b[1] ← wT,1A[1]+bT,1
A[1] = σ(Z[1])
Z[2] = W[2]A[1] + b[2]
A[2] = σ(Z[2])

```

Handwritten annotations: $A^{[1]}$ is circled in blue. Brackets group the first two equations and the last two equations.

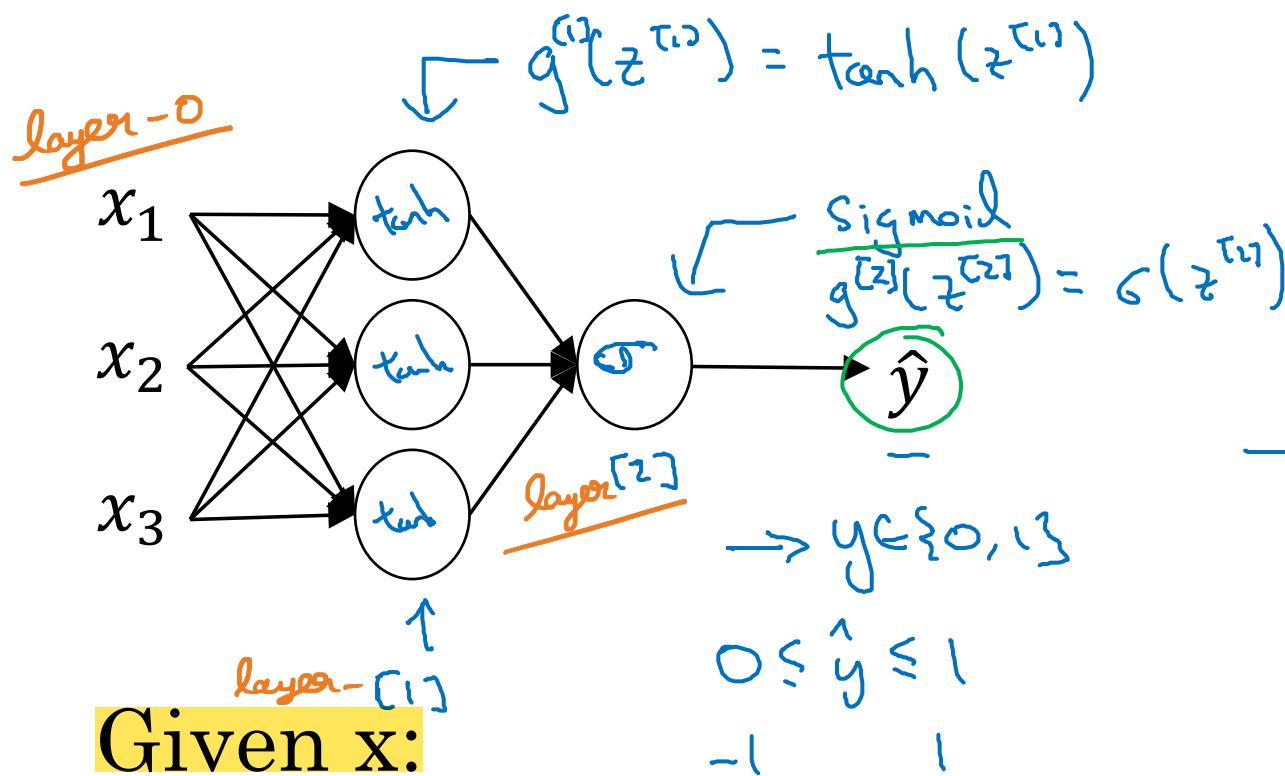


deeplearning.ai

One hidden layer
Neural Network

Activation functions

Activation functions

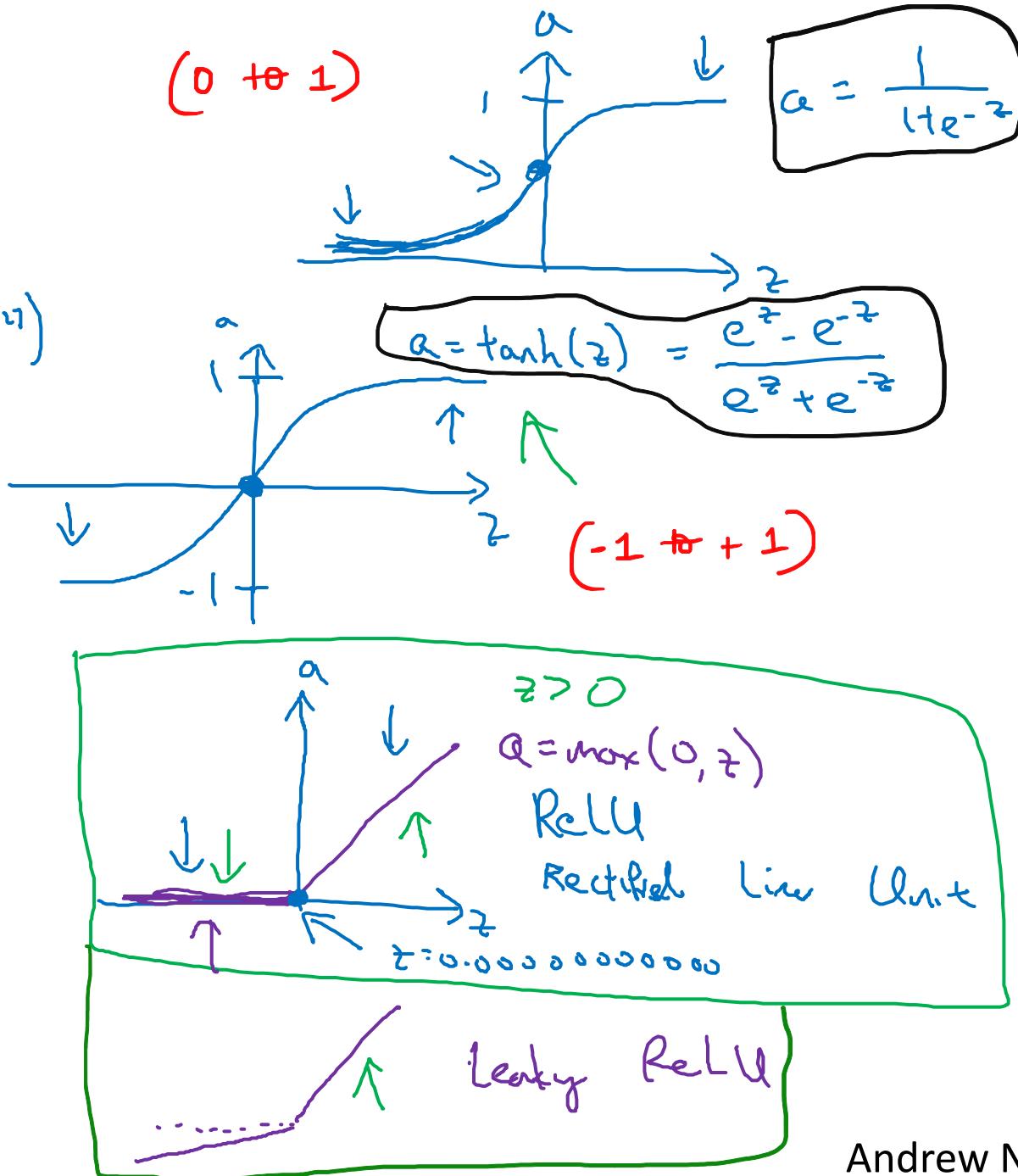


$$z^{[1]} = W^{[1]}x + b^{[1]}$$

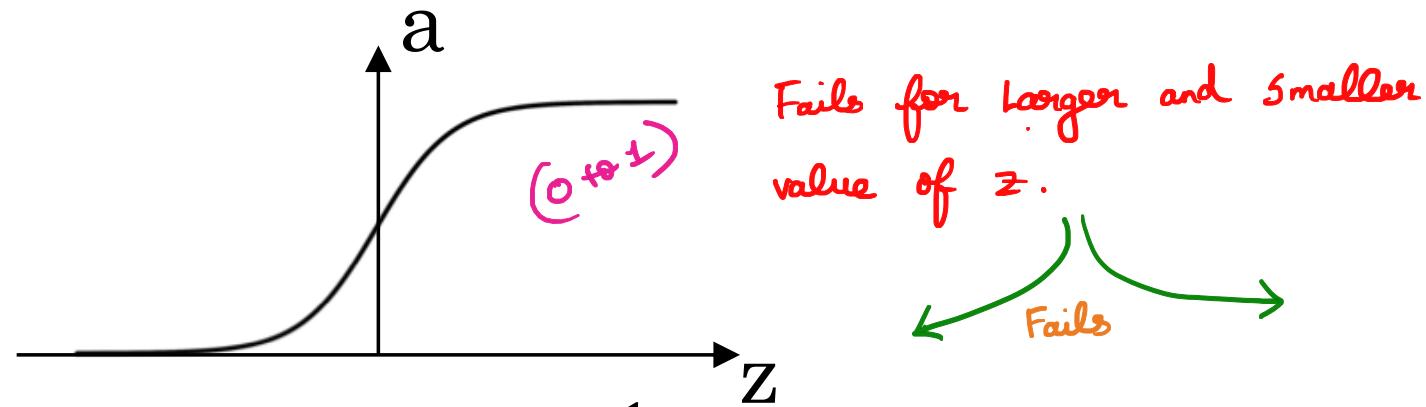
→ ~~$a^{[1]} = \sigma(z^{[1]})$~~ $\overset{[1]}{g}(z^{[1]})$

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

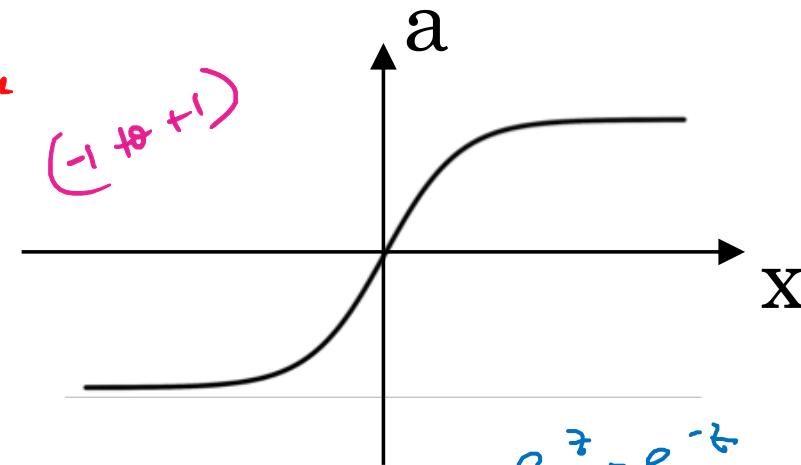
→ ~~$a^{[2]} = \sigma(z^{[2]})$~~ $\overset{[2]}{g}(z^{[2]})$



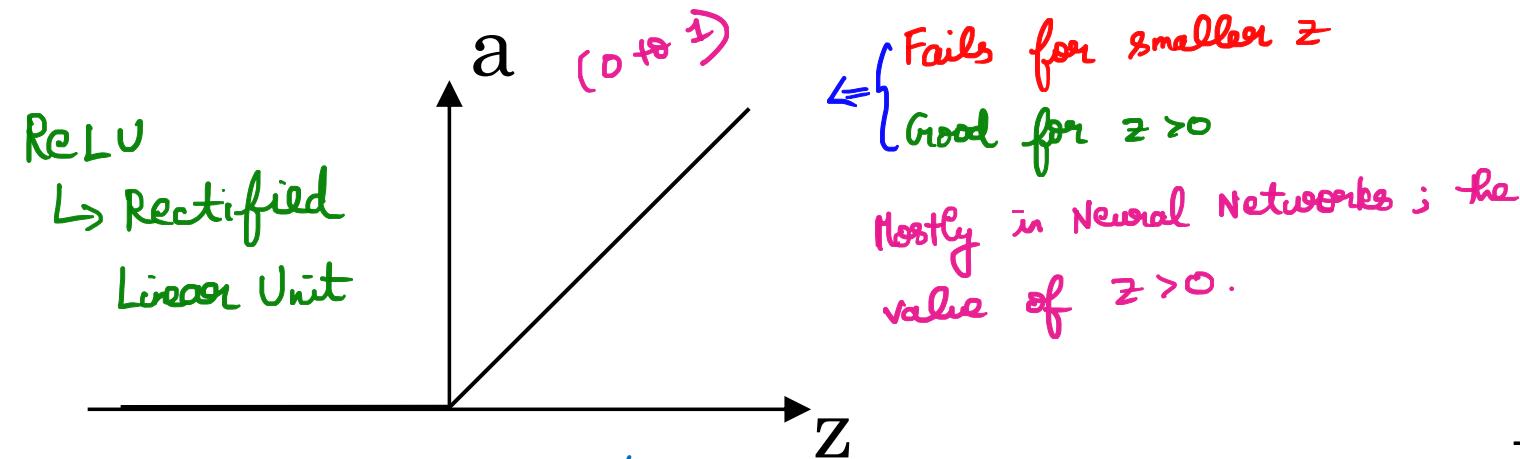
Pros and cons of activation functions



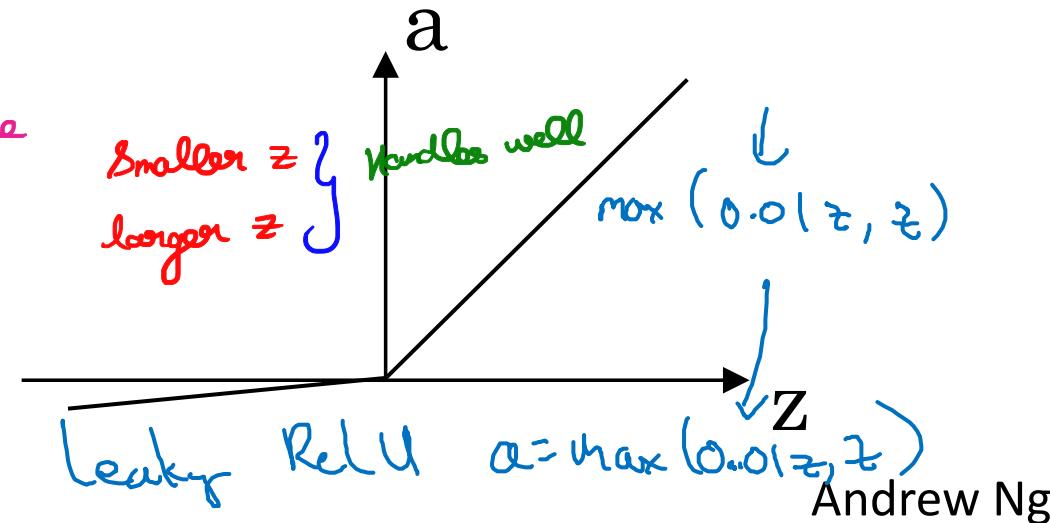
$$\text{sigmoid: } a = \frac{1}{1 + e^{-z}}$$



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



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



Andrew Ng

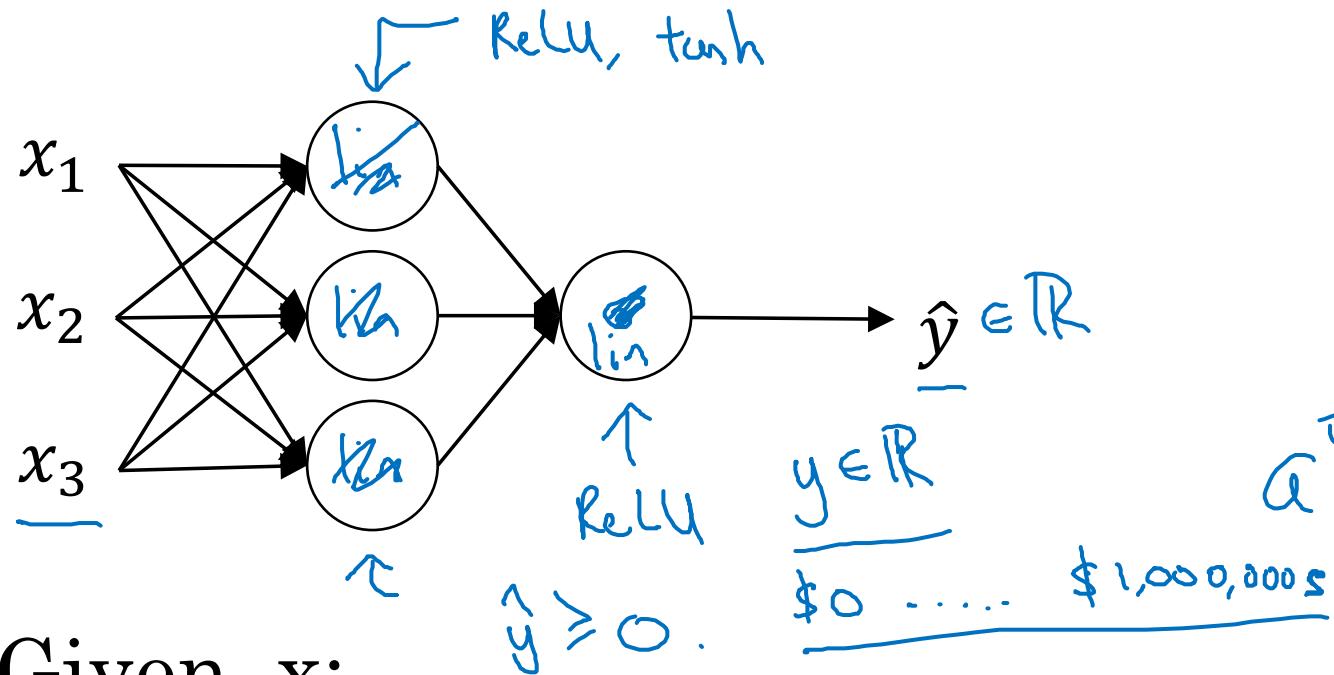


deeplearning.ai

One hidden layer Neural Network

Why do you
need non-linear
activation functions?

Activation function



Given x :

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

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

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

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

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

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

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

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

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

$$= \underline{w'x + b'}$$

$$g(z) = z$$

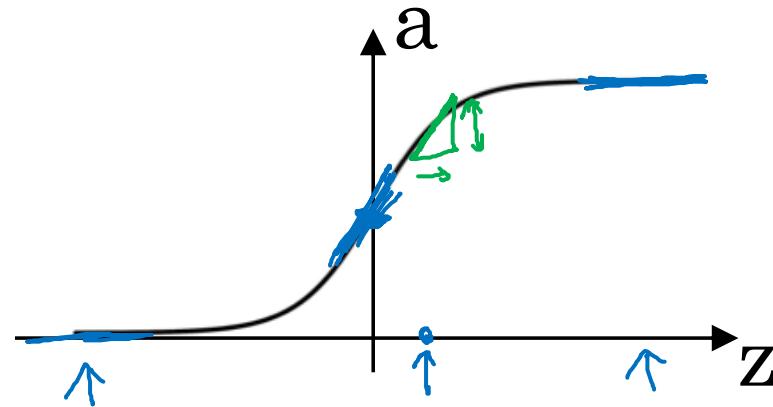


deeplearning.ai

One hidden layer Neural Network

Derivatives of activation functions

Sigmoid activation function



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

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

$$\begin{aligned} g'(z) &= \boxed{\frac{d}{dz} g(z)} \\ &= \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}} \right) \\ &= g(z) \left(1 - g(z) \right) \leftarrow \\ &= \boxed{a(1-a)} \quad \left| \begin{array}{l} g'(z) = a(1-a) \\ \uparrow \end{array} \right. \end{aligned}$$

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

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

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

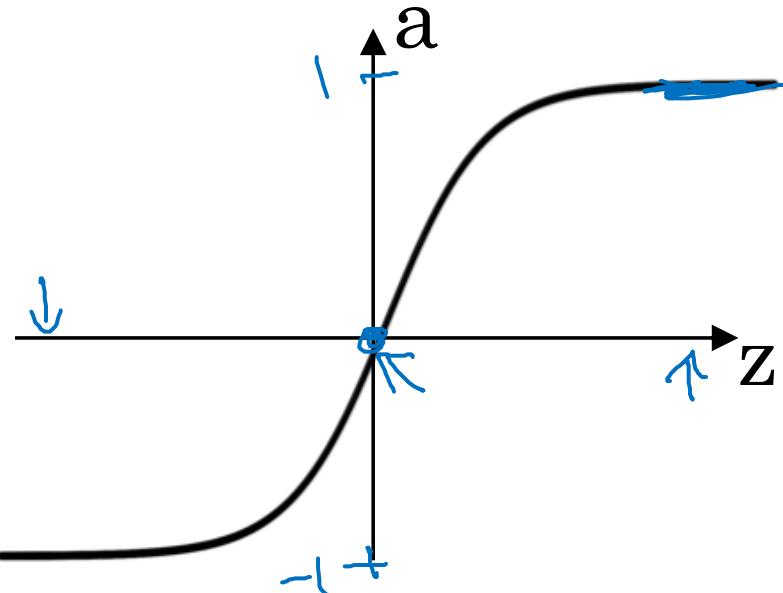
$$\frac{d}{dz} g(z) \approx 0 \cdot (1-0) \approx 0$$

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

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

Andrew Ng

Tanh activation function



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

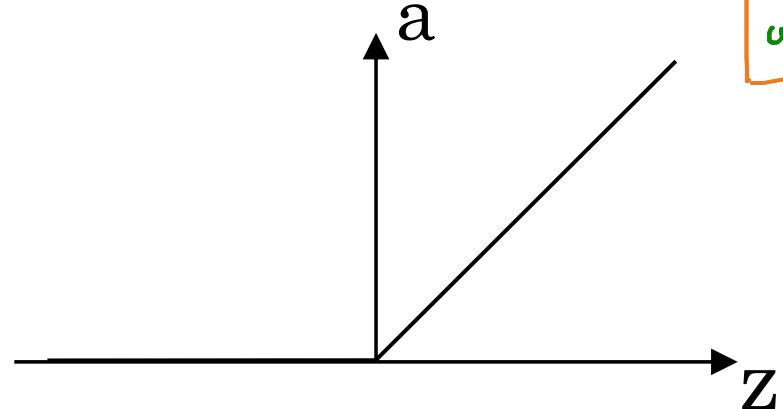
$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\begin{aligned} g'(z) &= \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z \\ &= \underline{\underline{1 - (\tanh(z))^2}} \leftarrow \end{aligned}$$

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

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

ReLU and Leaky ReLU

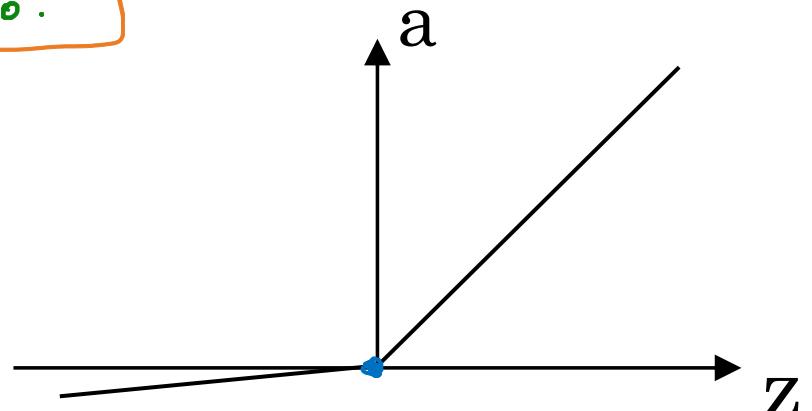


ReLU

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

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

Gradient is technically Not defined when z is exactly equal to 0.



Leaky ReLU

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

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



deeplearning.ai

One hidden layer
Neural Network

Gradient descent for
neural networks

Gradient descent for neural networks

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

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

Gradient Descent:

→ Repeat {

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

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

$$\omega^{[1]} := \omega^{[1]} - \alpha \frac{\partial J}{\partial \omega^{[1]}}$$

$$b^{[1]} := b^{[1]} - \alpha \frac{\partial J}{\partial b^{[1]}}$$

↳

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

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

Formulas for computing derivatives

Forward propagation:

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

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

$$z^{[2]} = w^{[2]}A^{[1]} + b^{[2]}$$

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

Back propagation:

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

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

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

$$dz^{[1]} = \underbrace{(w^{[2]T} dz^{[2]})}_{(n^{[2]}, m)} \times \underbrace{g^{[2]'}(z^{[2]})}_{\text{element-wise product}} \underbrace{(n^{[1]}, m)}$$

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

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

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

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

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



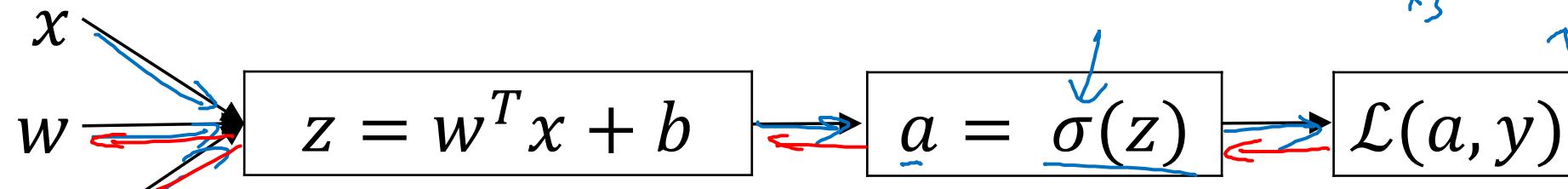
deeplearning.ai

One hidden layer
Neural Network

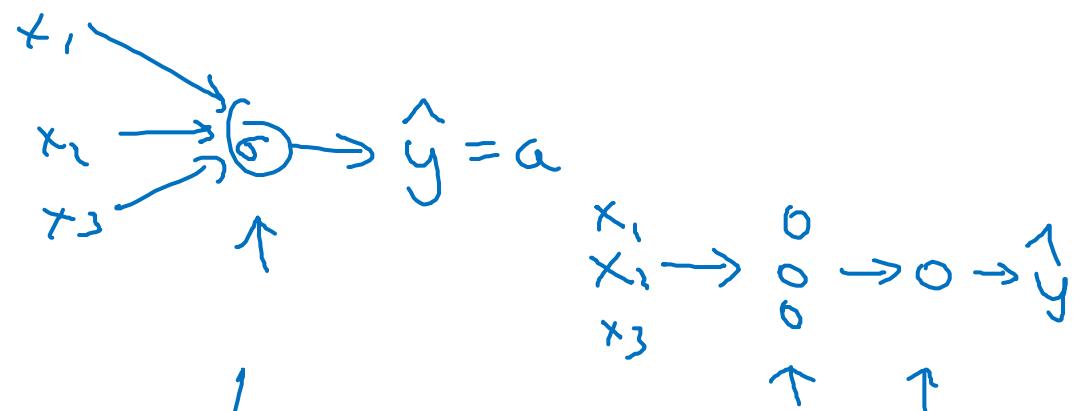
Backpropagation
intuition (Optional)

Computing gradients

Logistic regression



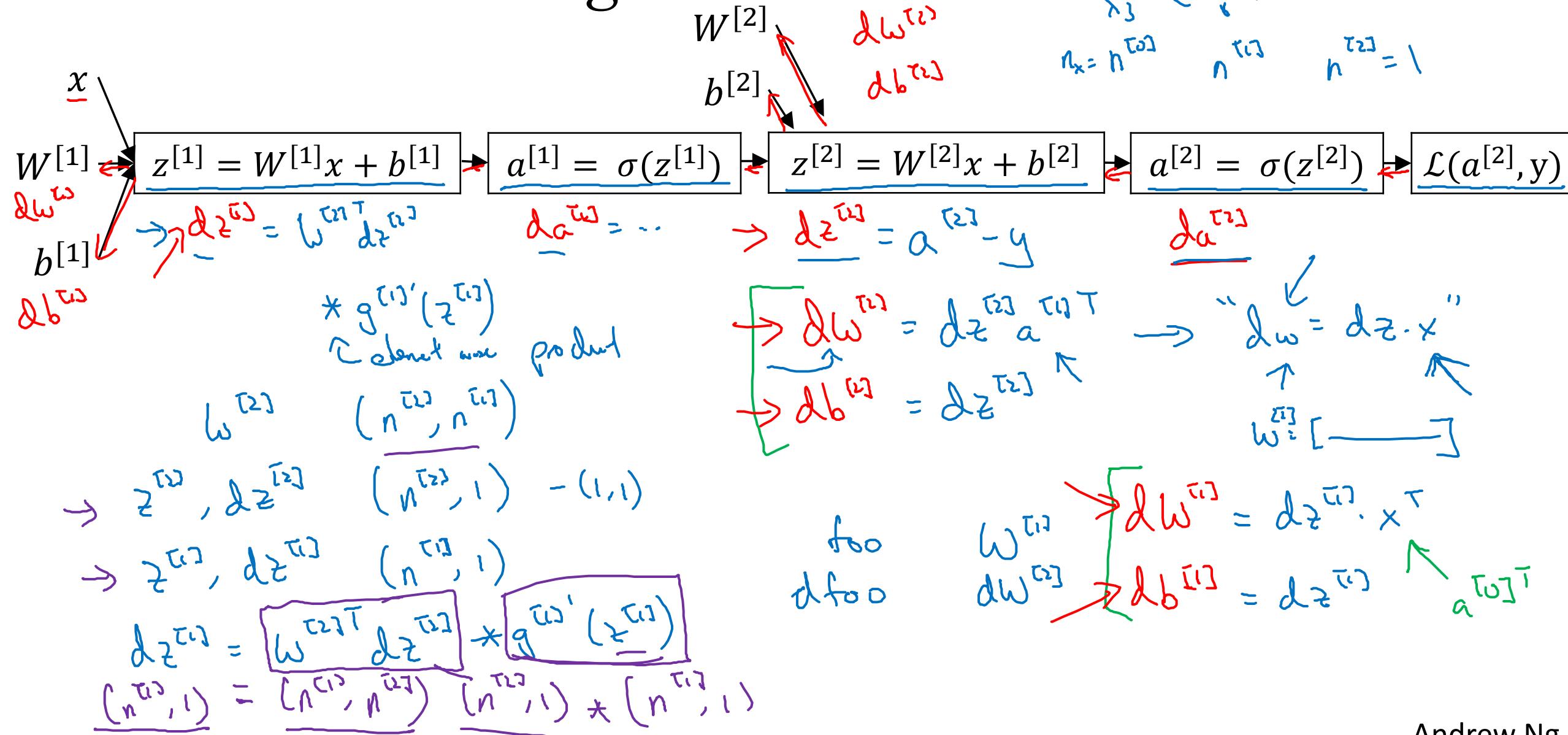
$$\begin{aligned} \frac{\partial z}{\partial w} &= x \\ \frac{\partial z}{\partial b} &= 1 \\ \frac{\partial z}{\partial a} &= g'(z) \\ g(z) &= \sigma(z) \\ \frac{\partial \mathcal{L}}{\partial z} &= \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial a}{\partial z} \\ "dz" &= "da" \end{aligned}$$



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

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

Neural network gradients



Summary of gradient descent

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

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

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

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

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

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

Vectorized implementation:

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

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

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

Summary of gradient descent

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

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

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

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

(n^{T₁}, 1)

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

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

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

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

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

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{(n^{T₂}, m)} * \underbrace{g^{[1]'}(Z^{[1]})}_{(n^{T₁}, m)}$$

elementwise product

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

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

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

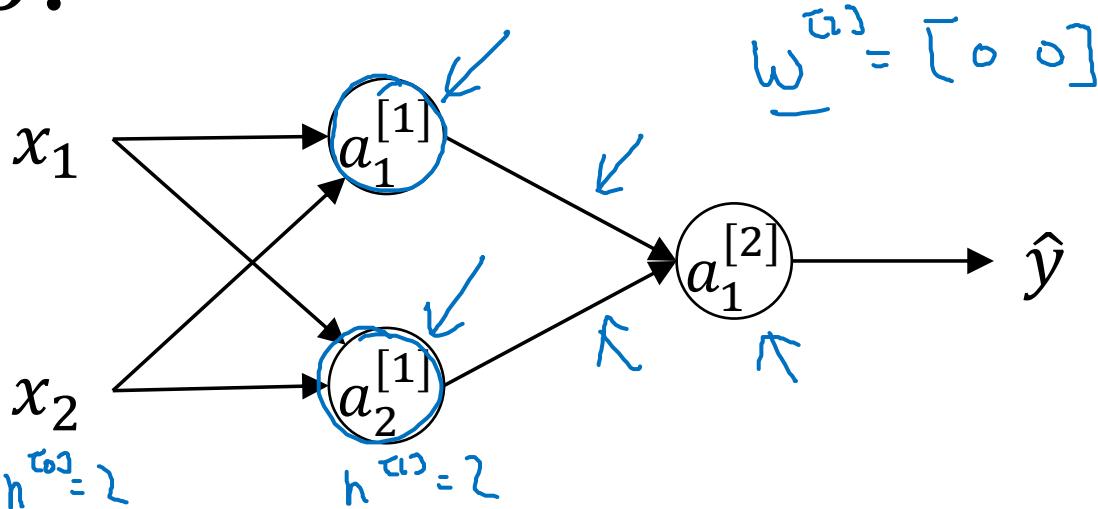


deeplearning.ai

One hidden layer
Neural Network

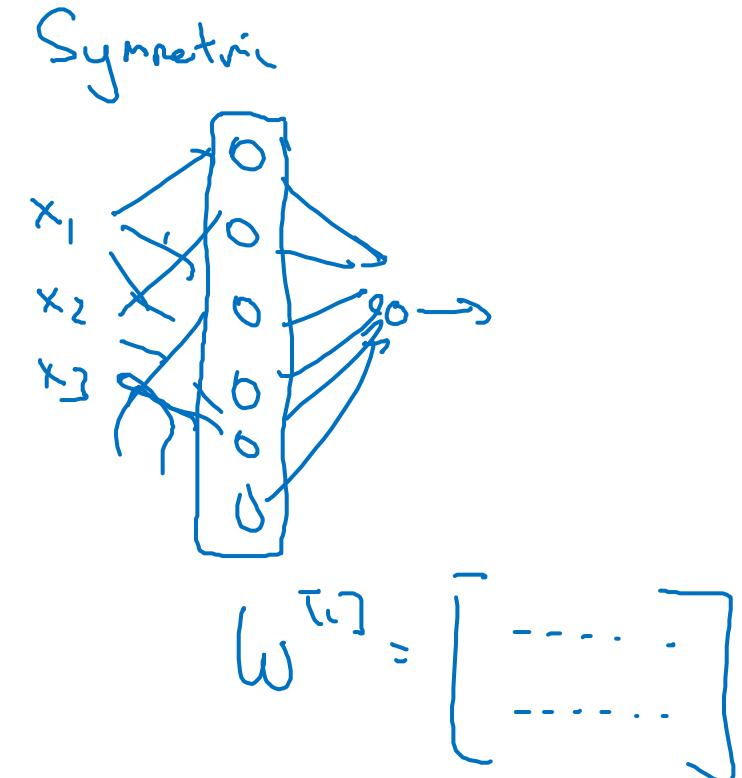
Random Initialization

What happens if you initialize weights to zero?

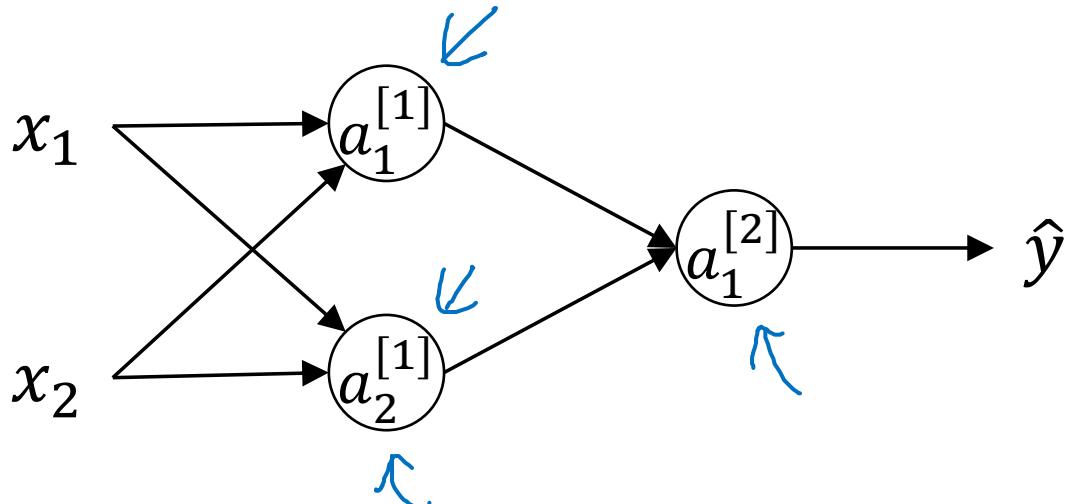


$$\delta w = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

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



Random initialization



$$\begin{aligned} \rightarrow w^{[1]} &= \text{np.random.randn}(2, 2) \times \frac{0.01}{100?} \\ b^{[1]} &= \text{np.zeros}(2, 1) \\ w^{[2]} &= \text{np.random.randn}(1, 2) \times 0.01 \\ b^{[2]} &= 0 \end{aligned}$$

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

