

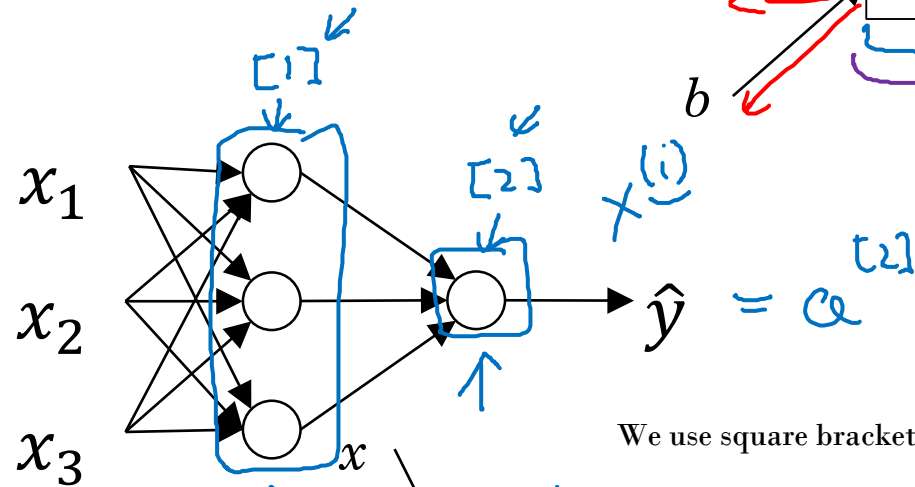
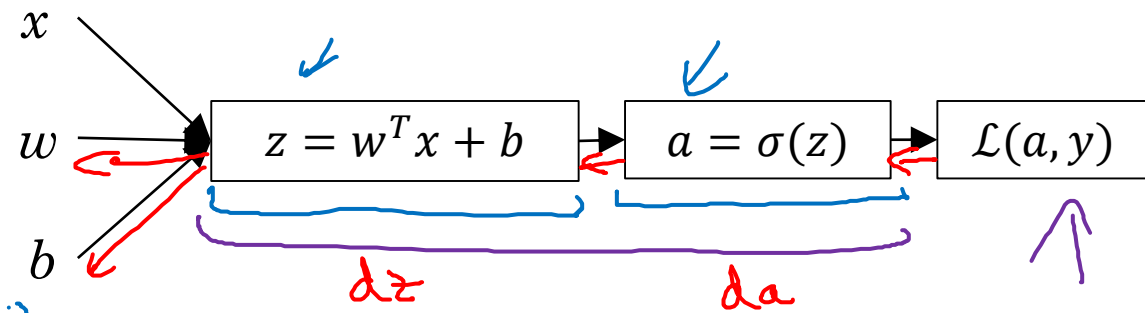
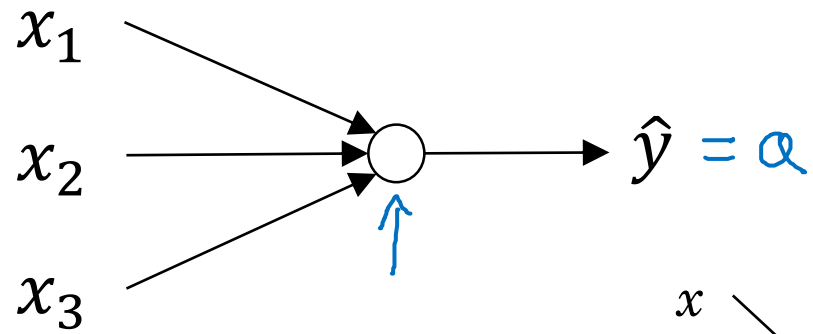


deeplearning.ai

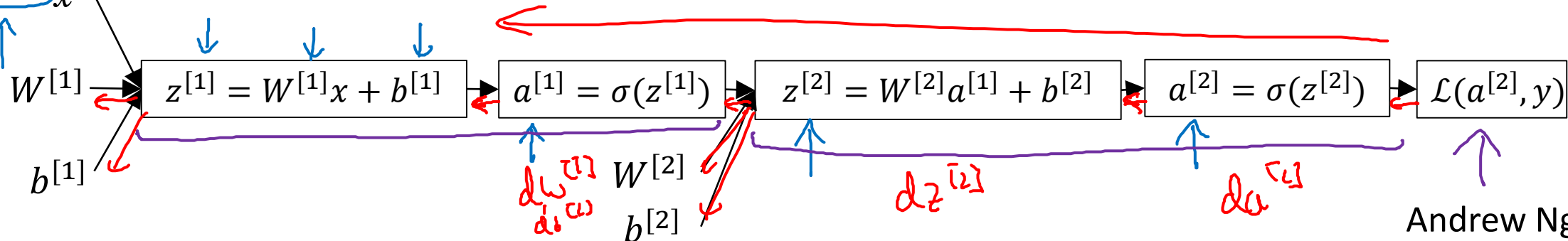
One hidden layer
Neural Network

Neural Networks Overview

What is a Neural Network?



We use square brackets to denote layer 1



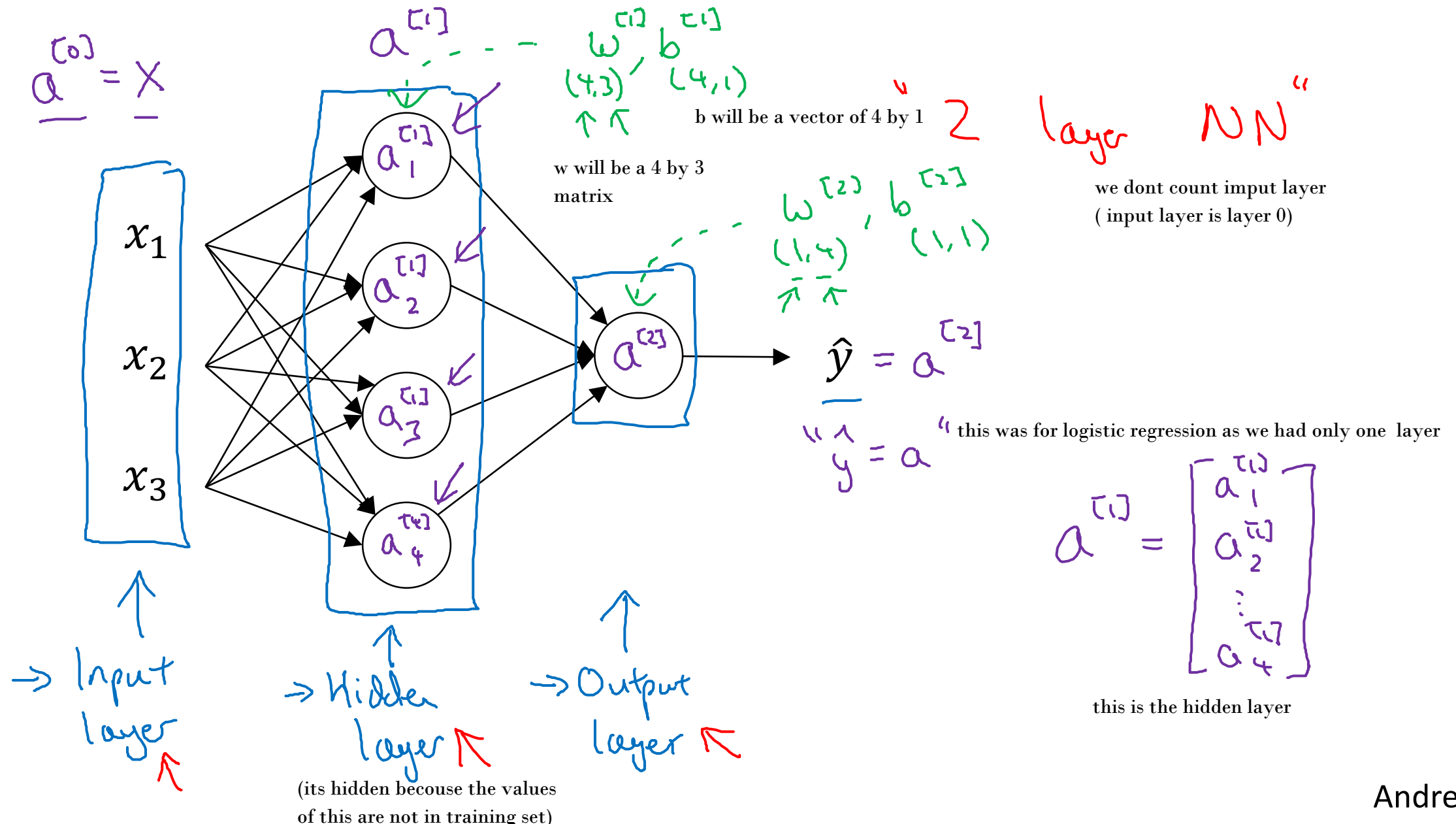


deeplearning.ai

One hidden layer
Neural Network

Neural Network
Representation

Neural Network Representation



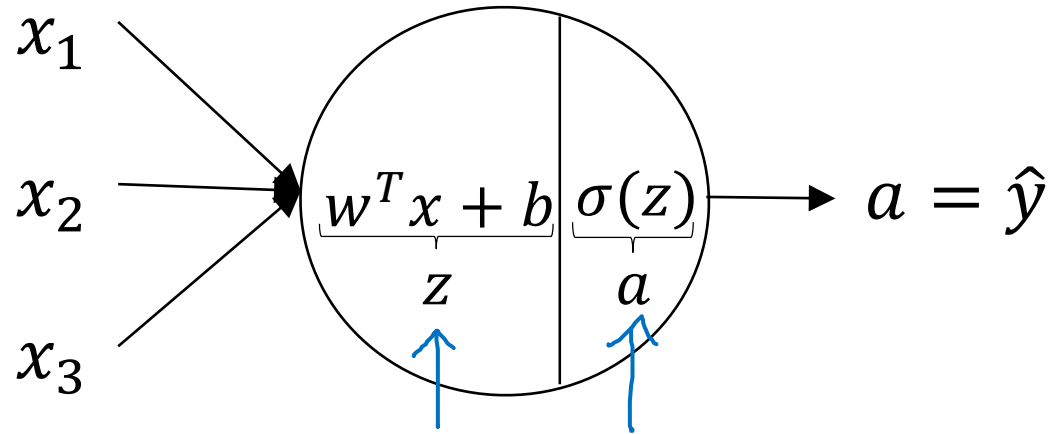


deeplearning.ai

One hidden layer Neural Network

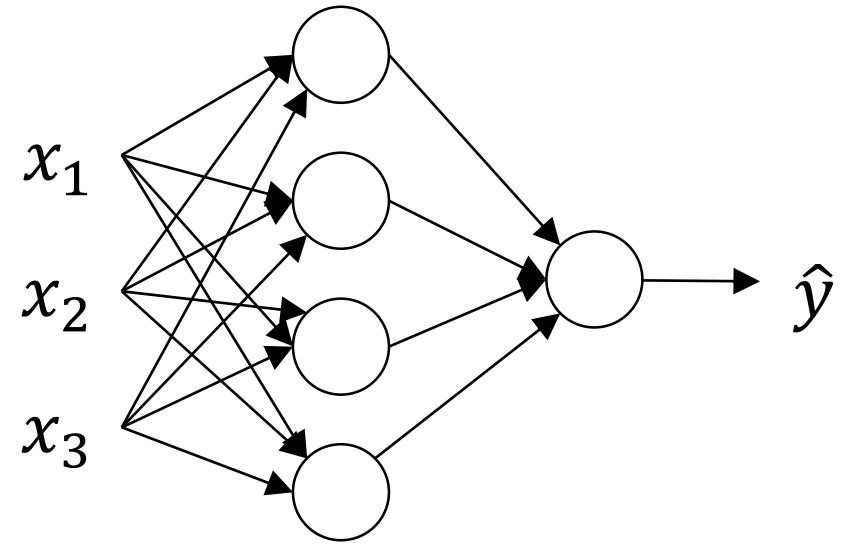
Computing a
Neural Network's
Output

Neural Network Representation

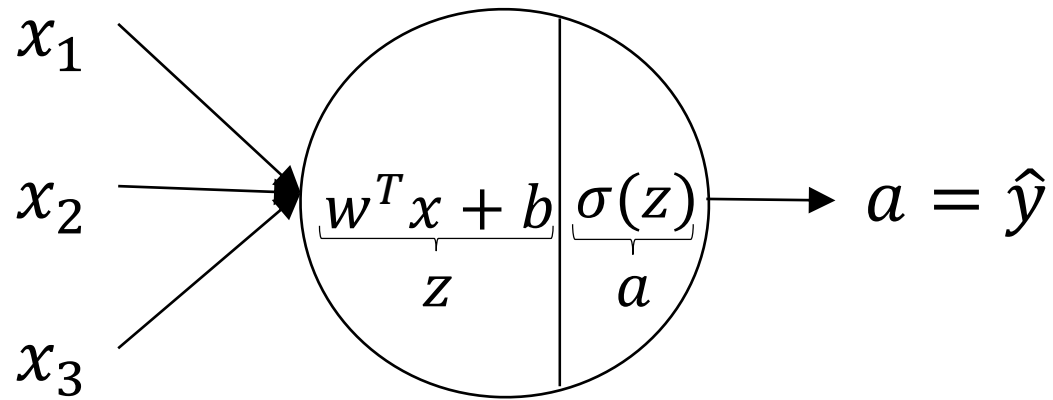


$$z = w^T x + b$$

$$a = \sigma(z)$$

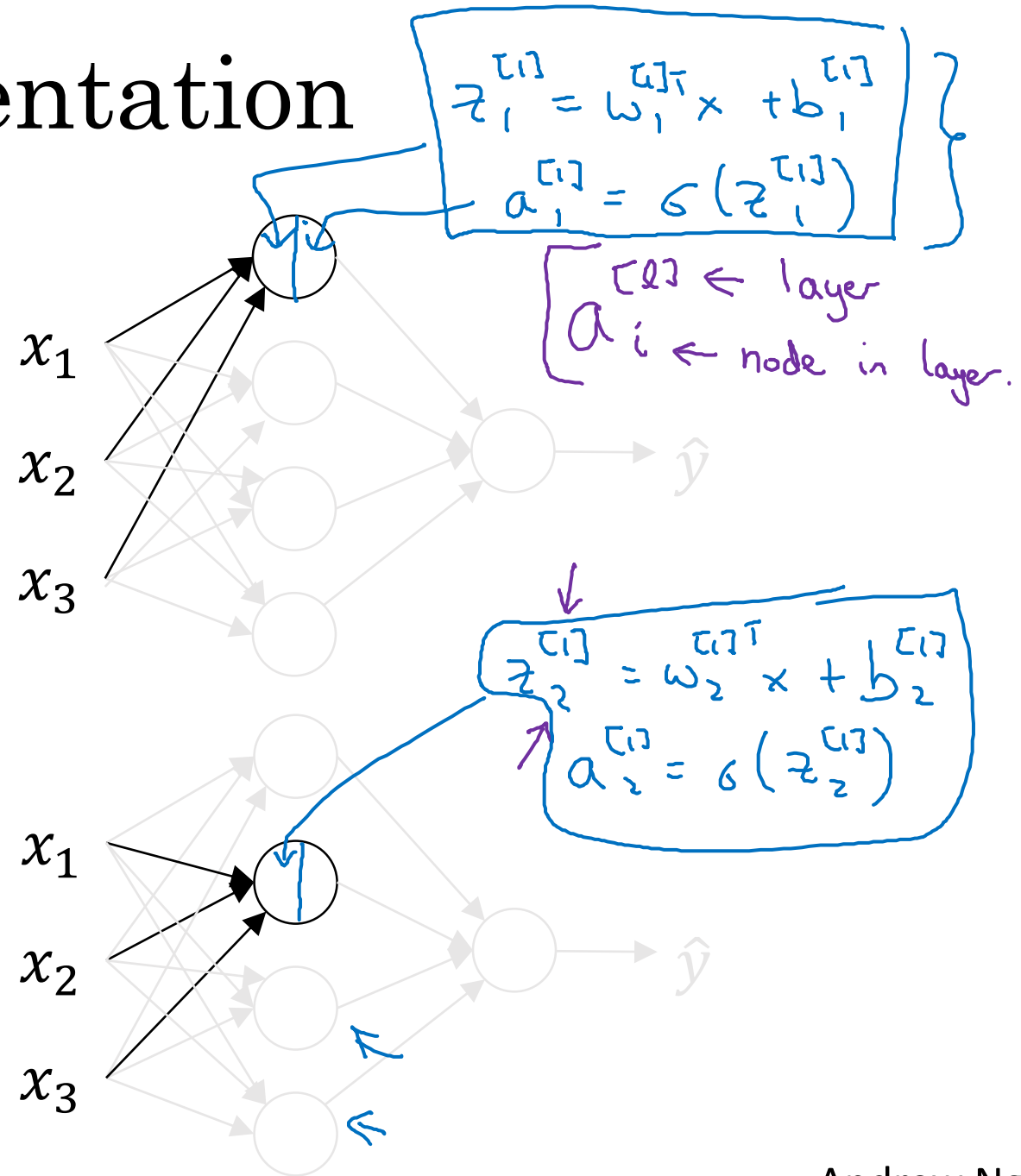


Neural Network Representation

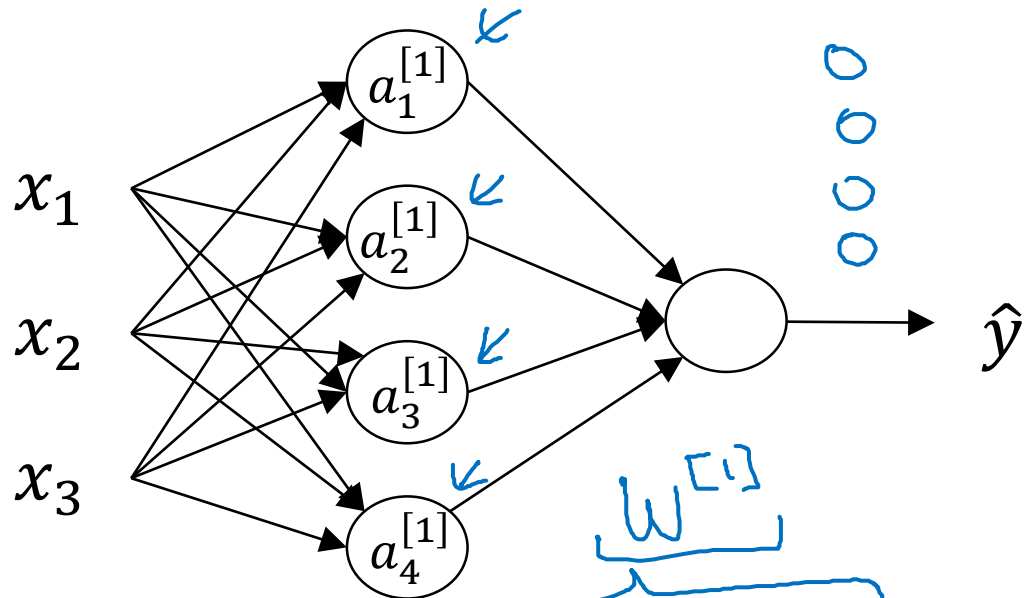


$$z = w^T x + b$$

$$a = \sigma(z)$$



Neural Network Representation



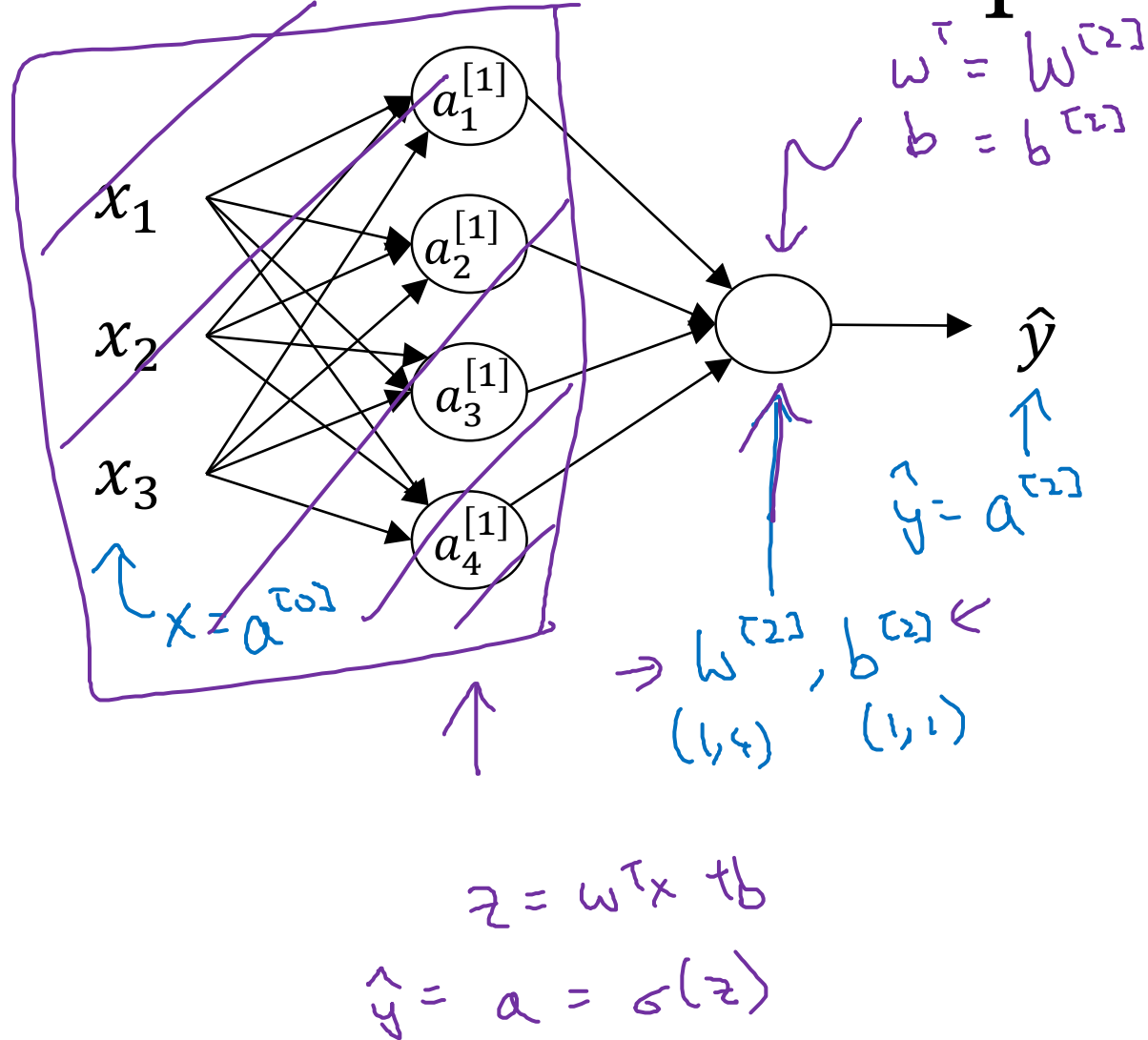
$$\begin{aligned}
 z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]} & a_1^{[1]} &= \sigma(z_1^{[1]}) \\
 z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]} & a_2^{[1]} &= \sigma(z_2^{[1]}) \\
 z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]} & a_3^{[1]} &= \sigma(z_3^{[1]}) \\
 z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]} & a_4^{[1]} &= \sigma(z_4^{[1]})
 \end{aligned}$$

Handwritten notes: $(w_1^{[1]})^T x$ and $Q^{[1]}$ are written above the first equation. A red box highlights the activation function part $a_i^{[1]} = \sigma(z_i^{[1]})$.

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

Handwritten notes: $(4, 3)$ is written below the weight matrix. $b^{[1]} (4, 1)$ is written below the bias vector. $\sigma(z^{[1]})$ is written below the activation vector. A red box highlights the activation function part $a_i^{[1]} = \sigma(z_i^{[1]})$. The text "This is a vector" is written below the final vector equation.

Neural Network Representation learning



Given input x :

you can replace x with $a[0]$ as its the layer 0 (its same)

$$\begin{aligned} \rightarrow z^{[1]} &= W^{[1]} a^{[0]} + b^{[1]} \\ &\quad (4,1) \quad (4,3) \quad (3,1) \quad (4,1) \\ \rightarrow a^{[1]} &= \sigma(z^{[1]}) \\ &\quad (4,1) \quad (4,1) \\ \rightarrow z^{[2]} &= W^{[2]} a^{[1]} + b^{[2]} \\ &\quad (1,1) \quad (1,4) \quad (4,1) \quad (1,1) \\ \rightarrow a^{[2]} &= \sigma(z^{[2]}) \\ &\quad (1,1) \quad (1,1) \end{aligned}$$



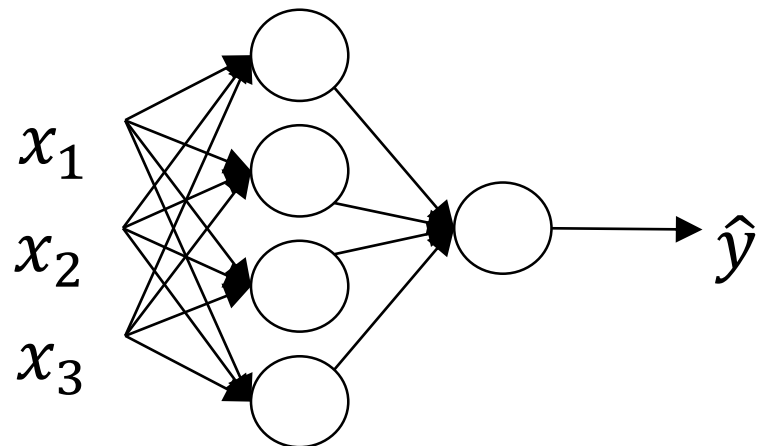
deeplearning.ai

One hidden layer Neural Network

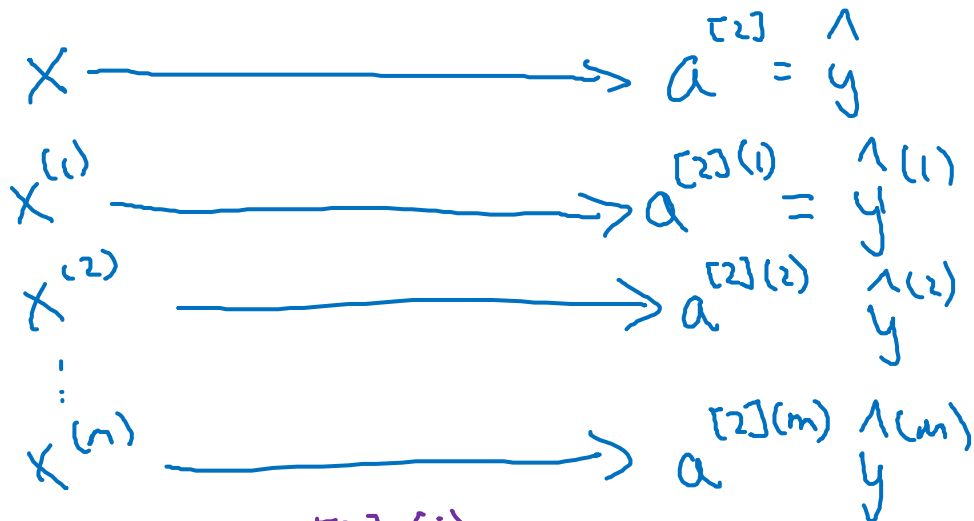
Vectorizing across multiple examples

Vectorizing across multiple examples

We saw previously how to calculate the output with one training example, now we will see for m training examples.



repeat for
n training
examples



$a^{[2](i)}$
← example i
layer 2

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

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

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

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

← This where the four
equations on how u
compute it for just
one training example

To compute the four above eq for m training examples

→ for $i = 1$ to m , we want to get rid of that for loop

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

Vectorizing across multiple examples

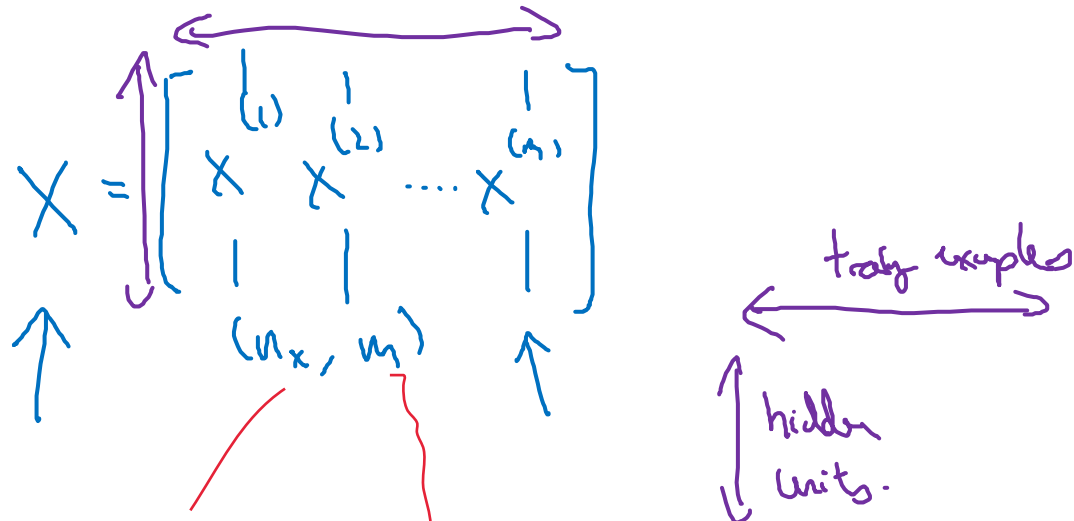
for $i = 1$ to m :

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

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

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

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

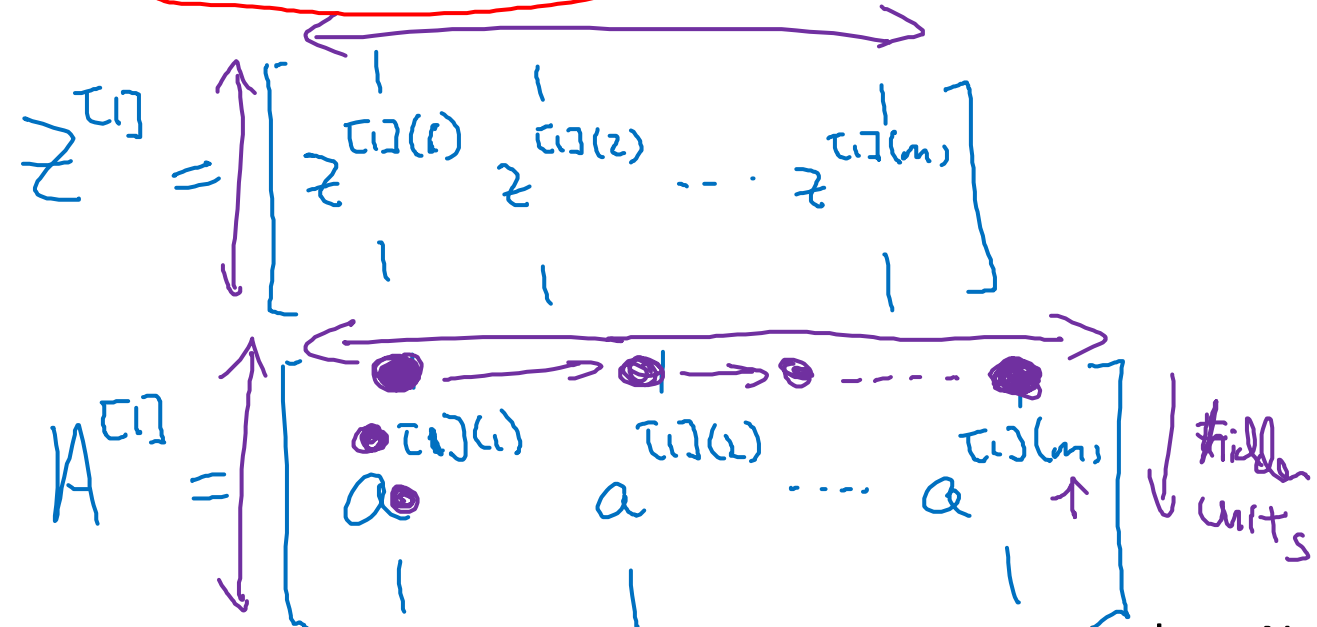


n_x number
of features

m is number
of training examples

This is how we vectorize it!!!

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





deeplearning.ai

One hidden layer Neural Network

Explanation
for vectorized
implementation

Justification for vectorized implementation

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

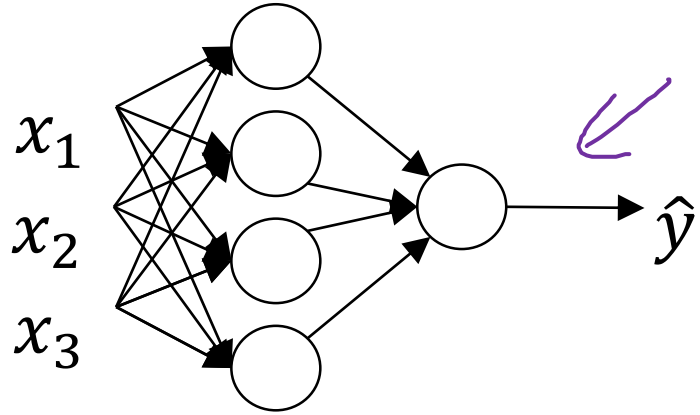
to simplify we assume $b = 0$

$$w^{[1]} = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{bmatrix} \quad w^{[1]} x^{(1)} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix} \quad w^{[1]} x^{(2)} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix} \quad w^{[1]} x^{(3)} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix}$$

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

here we have just justified that the above is a correct vectorization.

Recap of vectorizing across multiple examples



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

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

for $i = 1$ to m

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

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

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

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

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

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

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

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

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

$$W^{[1]}A^{[0]} + b^{[1]}$$



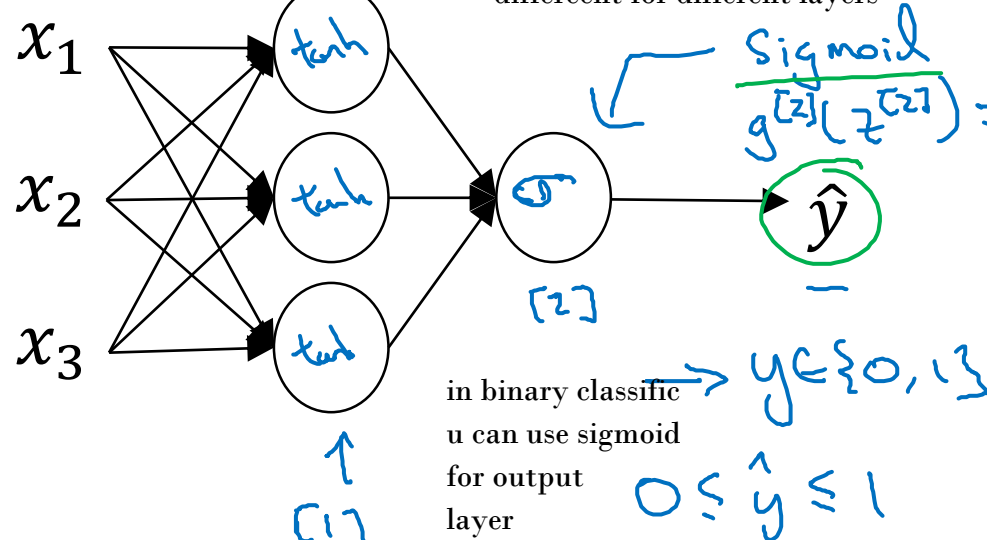
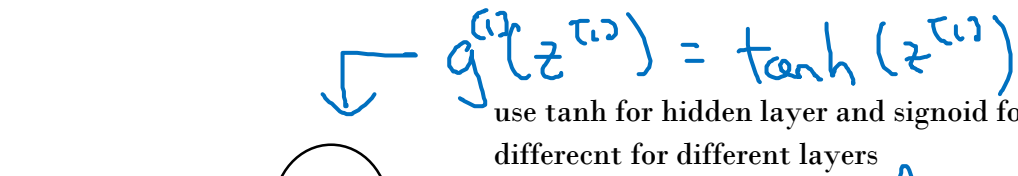
deeplearning.ai

One hidden layer Neural Network

Activation functions

Activation functions

tanh is centred at 0, when u train u NN u center data with 0 mean, using tanh instead of sigmoid has the effect of centering your data so that mean is closer to 0 rather than 0.5 and this makes learning a bit easier.

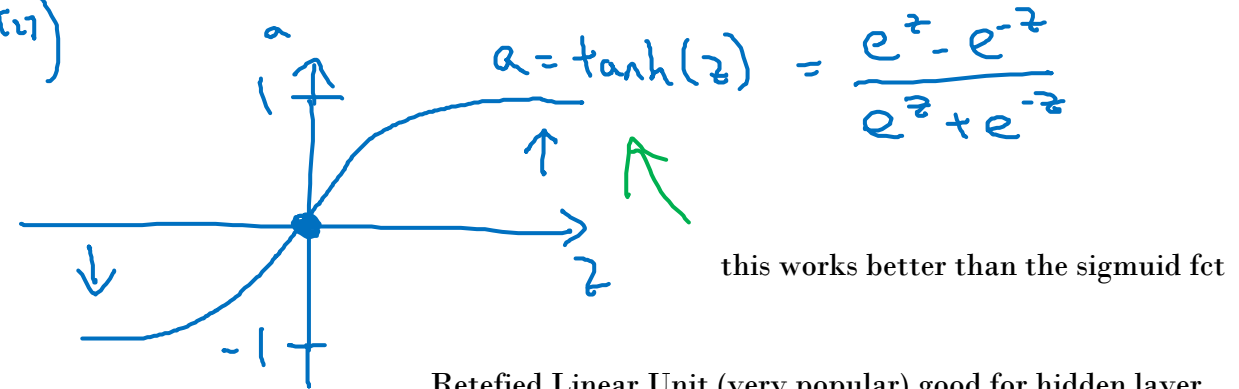
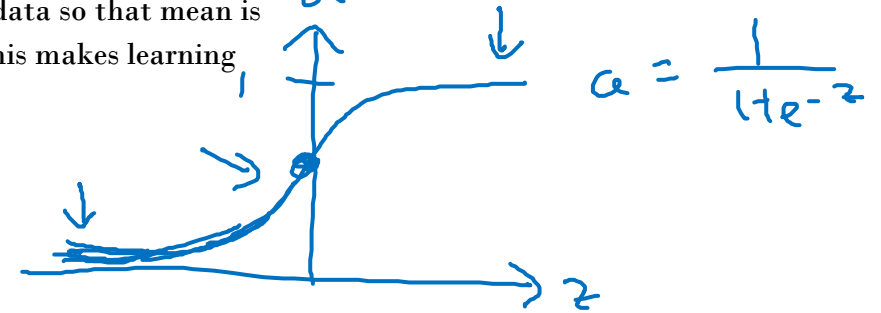


$$z^{[1]} = W^{[1]}x + b^{[1]}$$

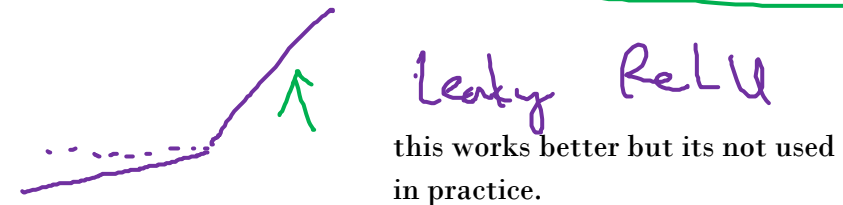
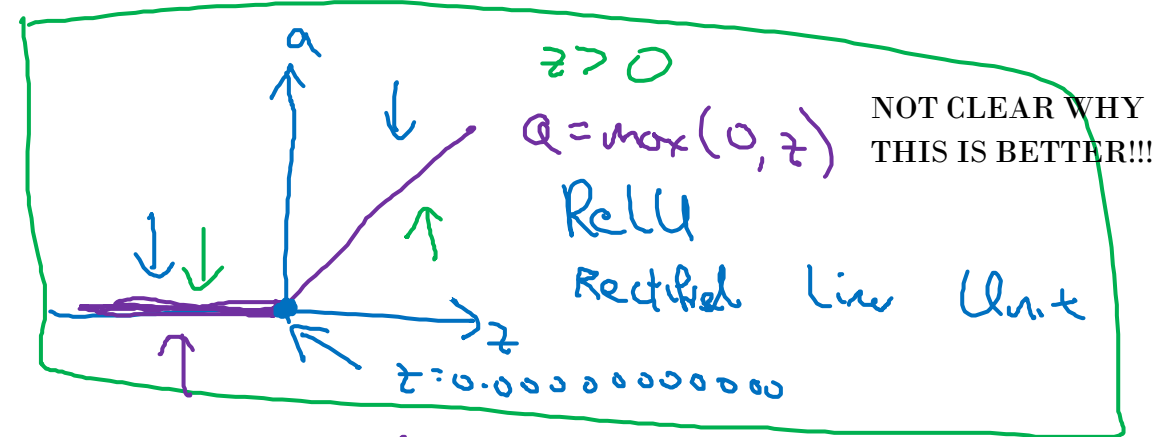
$$\rightarrow a^{[1]} = \sigma(z^{[1]}) \quad g^{(1)}(z^{(1)})$$

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

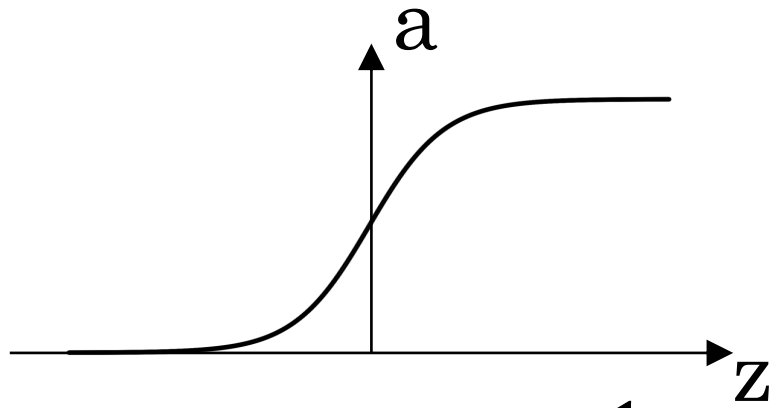
$$\rightarrow a^{[2]} = \sigma(z^{[2]}) \quad g^{(2)}(z^{(2)})$$



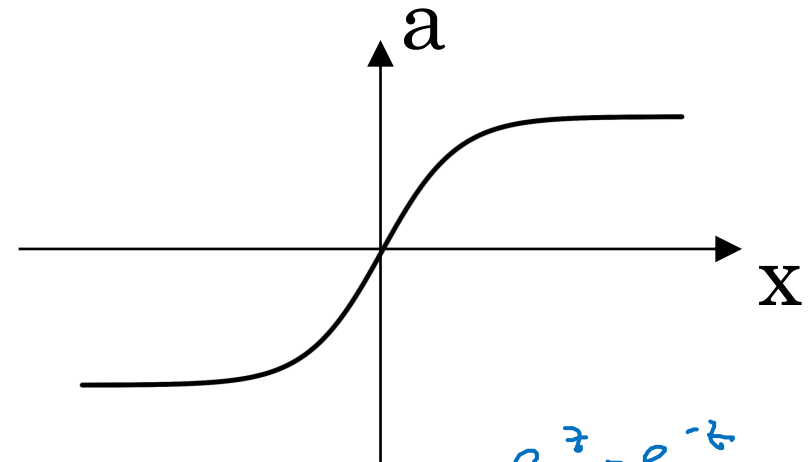
Rectified Linear Unit (very popular) good for hidden layer



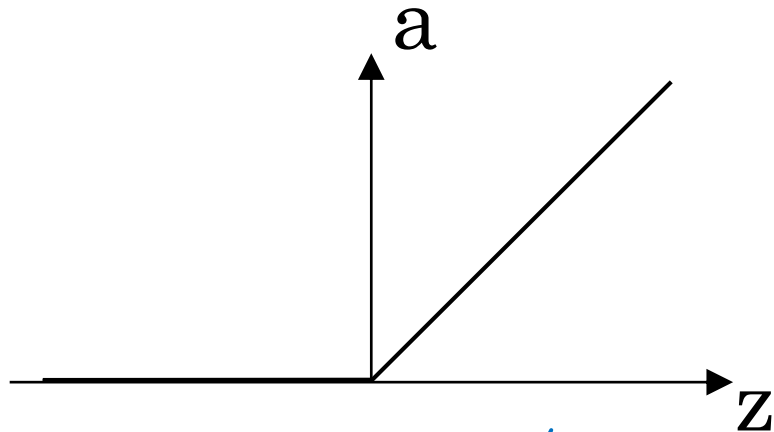
Pros and cons of activation functions



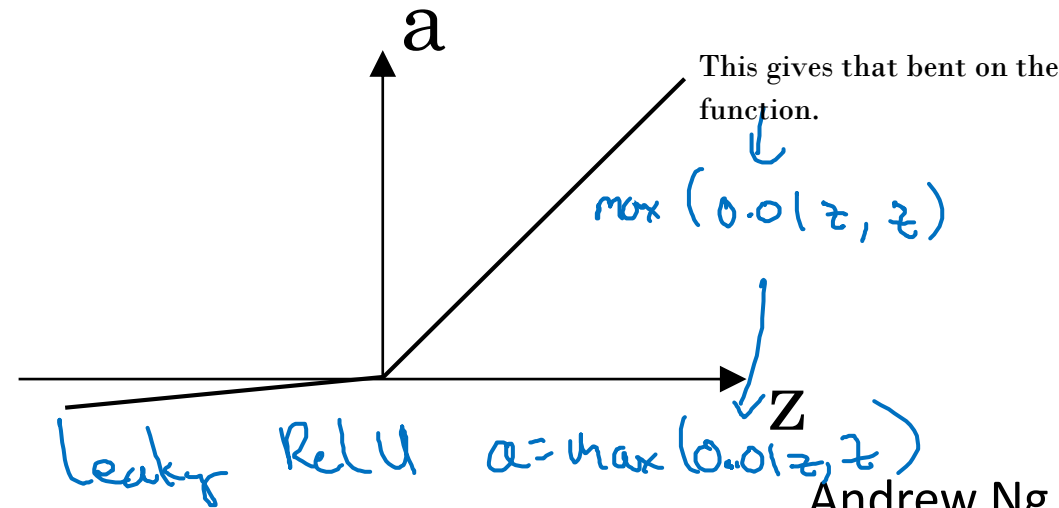
sigmoid: $a = \frac{1}{1 + e^{-z}}$



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



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



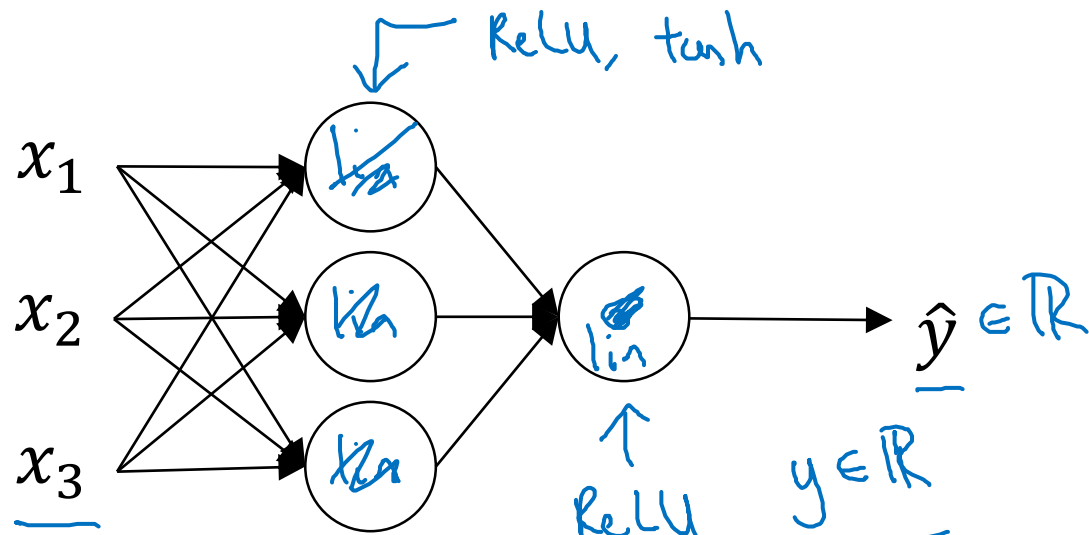


deeplearning.ai

One hidden layer Neural Network

Why do you
need non-linear
activation functions?

Activation function



hidden linear act fct layer is usles.

Given x :

$$\begin{aligned} \rightarrow z^{[1]} &= W^{[1]}x + b^{[1]} \\ \rightarrow a^{[1]} &= \cancel{g^{[1]}(z^{[1]})} z^{[1]} \\ \rightarrow z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\ \rightarrow a^{[2]} &= \cancel{g^{[2]}(z^{[2]})} z^{[2]} \end{aligned}$$

$y \in \mathbb{R}$
\$0 \dots \\$1,000,000\$

for a prob like housing price predict
you can use hidden relu and output layer a
lin activat function

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

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

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

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

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

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

$$g(z) = z$$

if u use a linear activation function, no matter how many layers
all its doing is computing a linear activation fct.

Andrew Ng



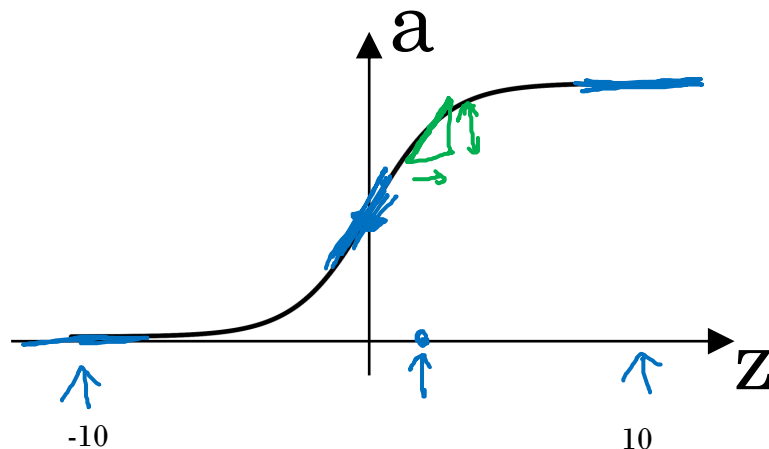
deeplearning.ai

One hidden layer Neural Network

Derivatives of activation functions

When you compute backpropagation u need to compute derivative or slope of the activation function

Sigmoid activation function



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

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

That's check that is the slope of the sigmoid fct:

very large: $z = 10$. $g(z) \approx 1$

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

very large: $z = -10$ $g(z) \approx 0$

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

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

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

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

short hand
denotation

= slope of $g(x)$ at z

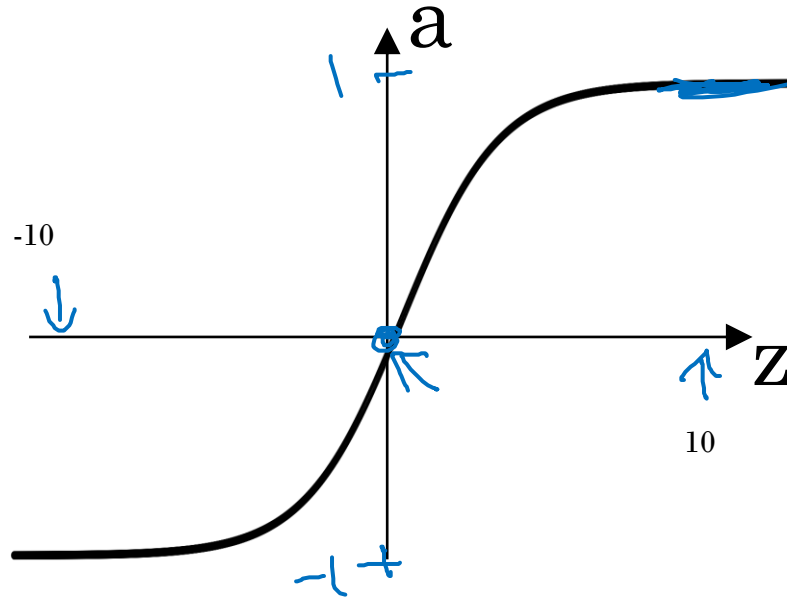
$$= \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}}\right)$$

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

$$= \boxed{a(1-a)} \quad \left| \begin{array}{l} g'(z) = a(1-a) \\ \uparrow \\ a \text{ is simoid fct} \end{array} \right.$$

a is simoid fct

Tanh activation function



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

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

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

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

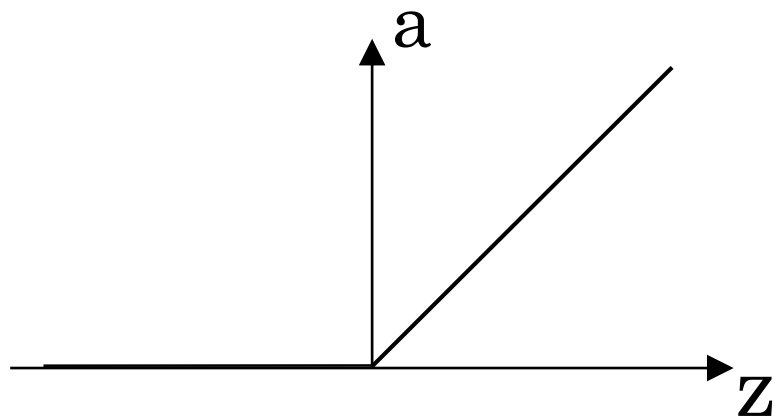
$$a = g(z),$$

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

if you have computed the value of a u can very quickly compute the derivative with that formula

$$\left| \begin{array}{ll} z=10 & \tanh(z) \approx 1 \\ & g'(z) \approx 0 \\ z=-10 & \tanh(z) \approx -1 \\ & g'(z) \approx 0 \\ z=0 & \tanh(z) = 0 \\ & g'(z) = 1 \end{array} \right.$$

ReLU and Leaky ReLU



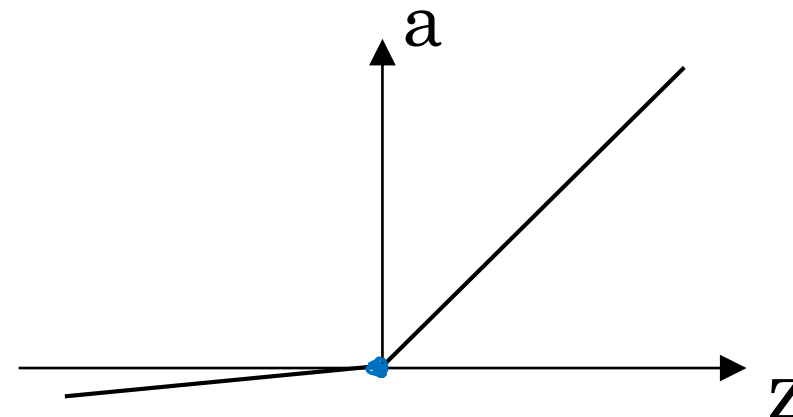
ReLU

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

$\rightarrow g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$
 ~~$undefined \text{ if } z = 0$~~
 $z = 0.0000 \dots 0$

its undefined for $z=0$ but u can set it to 1 or 0
it does not matter

This works better but its not much used.



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

Here we will see how to implement gradient descent for the NN with one hidden layer, next we will see why these equations work.

Gradient descent for neural networks

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

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

Cost function: $J(W^{[1]}, b^{[1]}, \underline{W^{[2]}}, \underline{b^{[2]}}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})$
 $\uparrow \quad \uparrow \quad \uparrow a^{[2]}$

This is the cost function and we are assuming we are doing binary classification.

Gradient descent:

→ Repeat {

we saw previously how
to compute these predictions

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

$\frac{\partial J}{\partial W^{[1]}}$ $= \frac{\partial J}{\partial W^{[1]}}$, $\frac{\partial J}{\partial b^{[1]}}$ $= \frac{\partial J}{\partial b^{[1]}}$, ... the key is to compute these derivative terms.

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

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

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

Formulas for computing derivatives

This is not trivial!!!

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:

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

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

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

summing horizontally in that matrix

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

$$dw^{[1]} = \frac{1}{m} dz^{[1]} x^T$$

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

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

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

this prevents python from outputting $(n,)$ \leftarrow it avoid to output this funny rank array

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



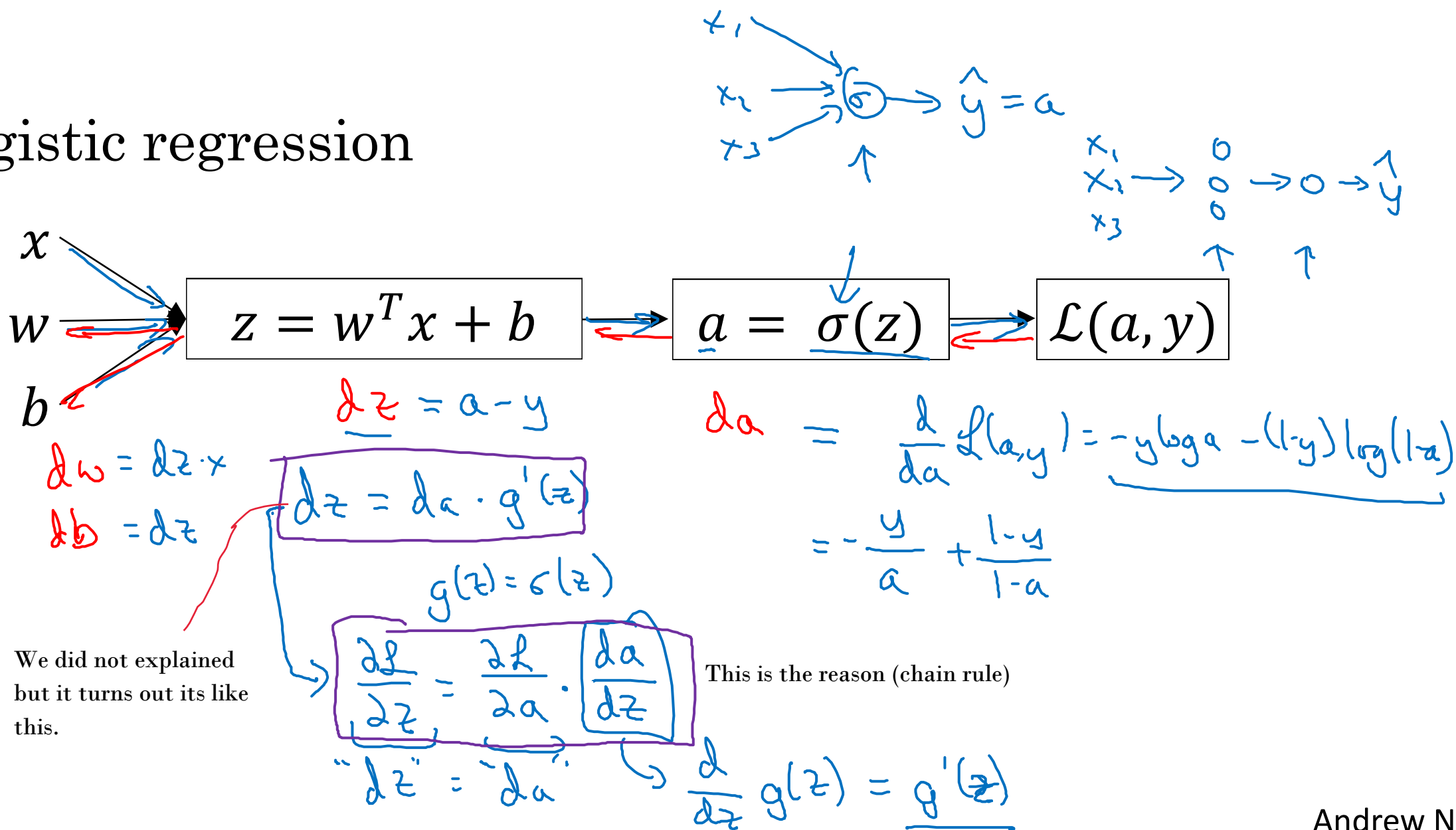
deeplearning.ai

One hidden layer
Neural Network

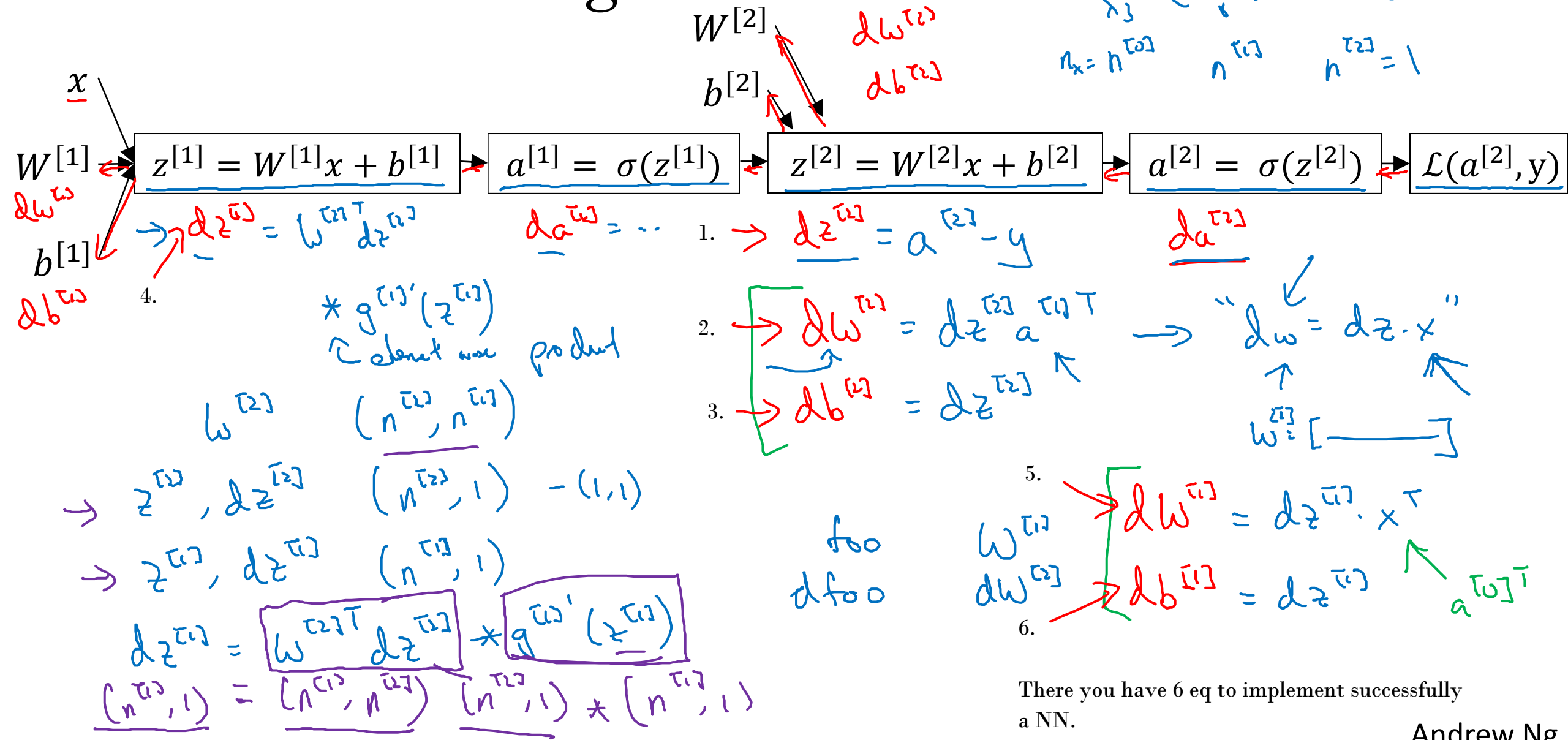
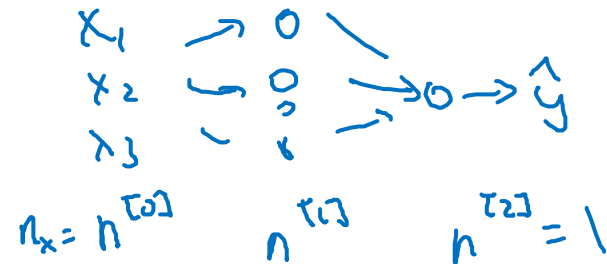
Backpropagation
intuition (Optional)

Computing gradients

Logistic regression



Neural network gradients



There you have 6 eq to implement successfully a NN.

Summary of gradient descent

So here are the 6 equations. They are just for one training example.

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

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

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

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

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

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

Here we vectorize the implementation for m training examples.

Vectorized Implementation:

$$z^{(1)} = W^{(1)} x + b^{(1)}$$
$$a^{(1)} = g^{(1)}(z^{(1)})$$
$$Z^{(1)} = \begin{bmatrix} z^{(1)(1)} \\ z^{(1)(2)} \\ \vdots \\ z^{(1)(m)} \end{bmatrix}$$
$$Z^{(1)} = W^{(1)} X + b^{(1)}$$
$$A^{(1)} = g^{(1)}(Z^{(1)})$$

Summary of gradient descent

Here are the vectorized implementation for m training examples:

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

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

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

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

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

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

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

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

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

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

elementwise product

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

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

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

There is this extra 1/m because the cost fct is the above



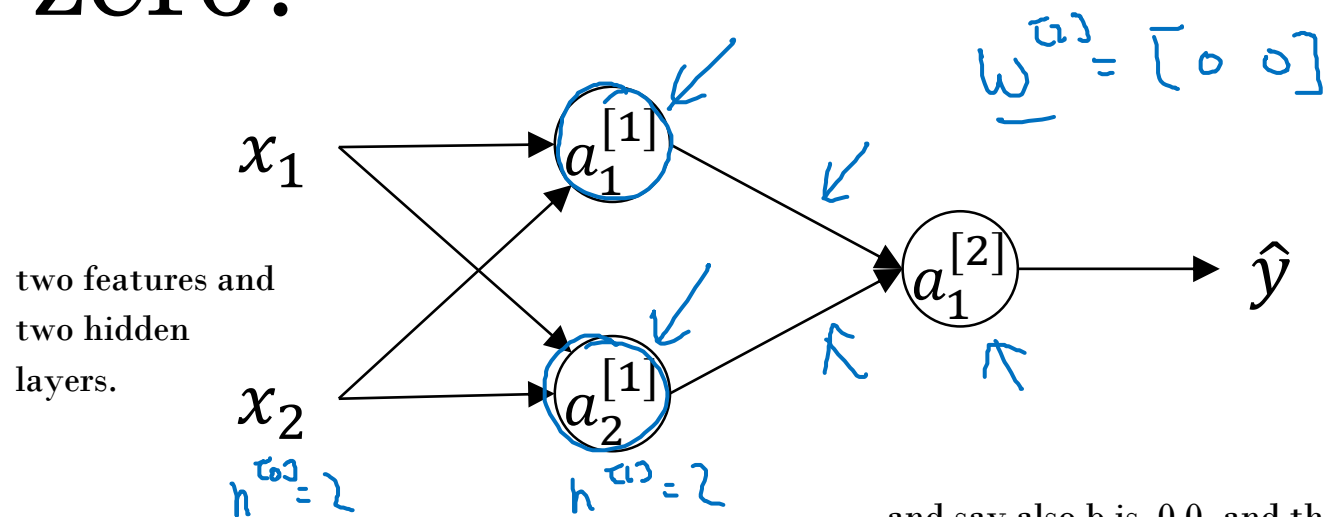
deeplearning.ai

One hidden layer
Neural Network

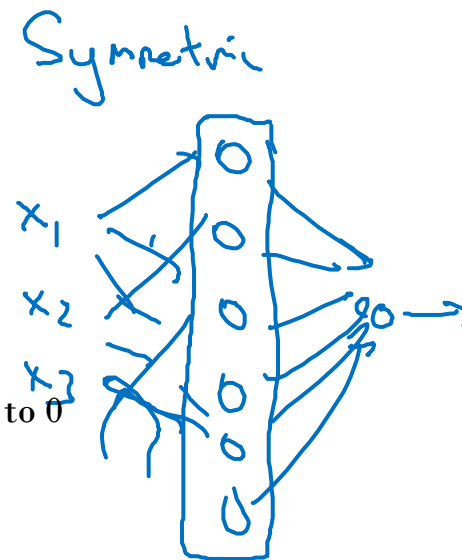
Random Initialization

What happens if you initialize weights to zero?

For a logistic reg is fine to initialize weights to 0 but for a neural network it wont work.



This is also called the Symetry breaking problem.



say we initialize by 0

$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

and say also b is 0 0, and this is ok to initialize to 0 but not for W.

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

For any example that you give it you have:

$a_1^{[1]} = a_2^{[1]}$ Then when you compute backprop

$$\delta z_1^{[1]} = \delta z_2^{[1]}$$

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

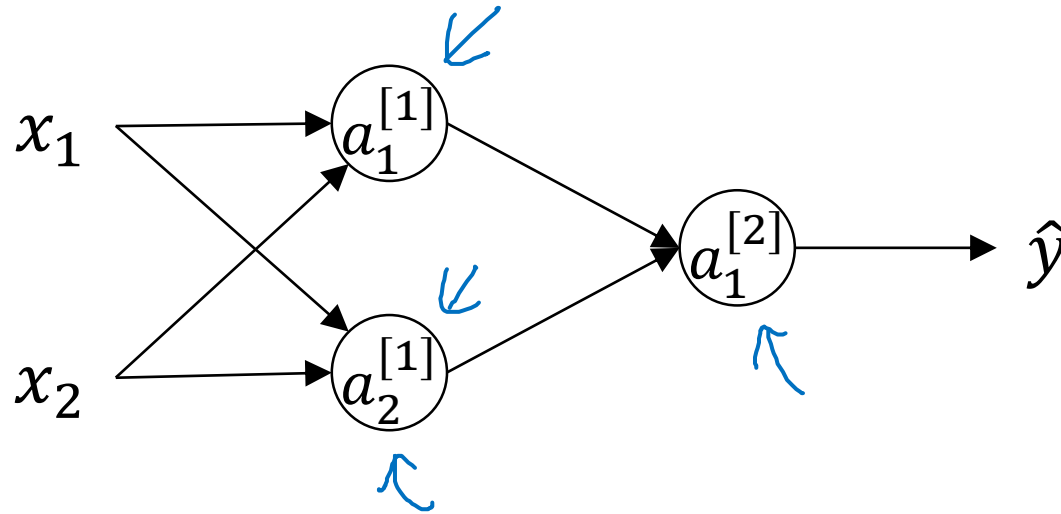
It turns out that after every iteration the hidden units will be the same.

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

$$w^{[1]} = \begin{bmatrix} \dots & \cdot \\ \dots & \cdot \end{bmatrix}$$

it will end up with first row equal to second row.

Random initialization



we prefer to init weight to very small random numbers. In certain cases u might want to choose a different constant, we will talk about this next week.

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

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

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

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

this is fine to init with 0

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

