



deeplearning.ai

One hidden layer
Neural Network

Neural Networks
Overview

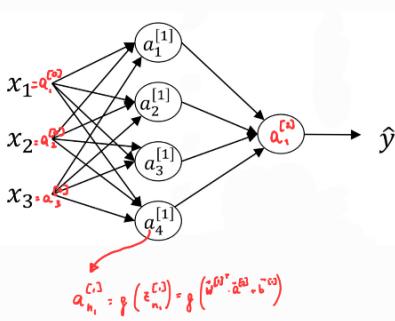
This lesson is about :

- GRADIENT DESCENT in NN via forward + backpropagation
- how to vectorize the algorithm (represent the network in terms of matrices and formulate forward and backpropagation as linear algebra operations) to optimize the computation
- what activation functions to choose (ReLU for hidden, sigmoid for output)
- how to initialize the weights (randomly).

BRIEFLY, WHAT IS BACKPROPAGATION?

- GRADIENT DESCENT optimizes the NN's parameters based on the gradient of the cost fn wrt each parameter. The following applies similarly to other numerical optimization problems that use the gradient of the cost function.
- BACKPROPAGATION is an efficient algo to compute the gradient of the cost fn wrt each parameter. Backpropagation is efficient because it first computes the loss gradient wrt "later" parameters (closer to the output) and uses them to compute the gradient wrt "earlier" parameters.

INFORMALLY SKETCH BACKPROPAGATION FOR A 2-layer NN



- Each neuron implements a linear concatenation of the inputs followed by a nonlinear function:

$$a_k^{[i]} = g(z_k^{[i]}) = g(w^{[i]^T} \cdot a^{[i-1]} + b^{[i]})$$

("activation of the k-th neuron on the i-th layer").

- We operate a FORWARD PROPAGATION pass to get the prediction of the network $\hat{y} = \vec{a}^{[2]} = a_0^{[2]}$ for the given input sample(s), and for the loss $L(\hat{y}, y)$. If we have many samples, we should instead compute the cost function $J = \frac{1}{m} \sum_{i=1}^m L(a^{[2]}, y^{(i)})$. However that is redundant because in the end we will need $\frac{\partial L}{\partial w}$ which is $\frac{1}{m} \sum_{i=1}^m \frac{\partial L(a^{[2]}, y^{(i)})}{\partial w}$ (orange of the gradient of the loss wrt variable). So, for each sample we compute the loss wrt variable of interest, we accumulate it, and we average its value after having cycled over all the samples.

- We operate a BACKPROPAGATION pass to compute the gradient of the loss (averaged over the input samples because we actually need the gradient of the cost function) wrt the parameters of the network using the chain rule for efficient computation. For example in this case we want $\frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{12}}, \frac{\partial L}{\partial w_{13}}, \frac{\partial L}{\partial w_{14}}$ and to get those we proceed backwards from the last layer to the first. We can do this because the loss function depends on the output of the network $a^{[2]}$ which corresponds to a concatenation of functions

$$\begin{aligned} L(a^{[2]}, y) &= L(g(w^{[2]} \cdot a^{[1]} + b^{[2]}), y) \\ &= L(g(w^{[2]} \cdot (w^{[1]^T} \cdot a^{[0]} + b^{[1]}) + b^{[2]}), y) \end{aligned}$$

- After iterating over all the samples and accumulating the gradient of the loss wrt the neurons of interest we compute its average (gradient of the cost function) and use that to update the parameters according to the chosen optimization algorithm, eg, gradient descent.

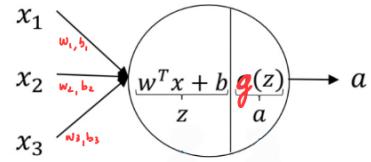
averaging the gradients of the loss over the mini-batch is equal to computing the gradient of the cost fn.

$w^{[2]} \rightarrow m$	$db^{[2]} \rightarrow m$	$dL \rightarrow m$	$dw^{[2]} \rightarrow m$
$b^{[2]} \rightarrow a \cdot db^{[2]}$	$a \cdot db^{[2]}$		
$w^{[1]} \rightarrow a \cdot dw^{[1]}$	$a \cdot dw^{[1]}$		
$b^{[1]} \rightarrow a \cdot db^{[1]}$	$a \cdot db^{[1]}$		

VECTORIZED REPRESENTATION OF THE NETWORK, & EFFICIENT IMPLEMENTATION OF GRADIENT DESCENT FOR A 2-LAYER NETWORK WITH SIGMOID ACTIVATION NEURONS AND LOGISTIC REGRESSION LOSS

NOTATION AND VECTORIZED REPRESENTATION OF THE NETWORK

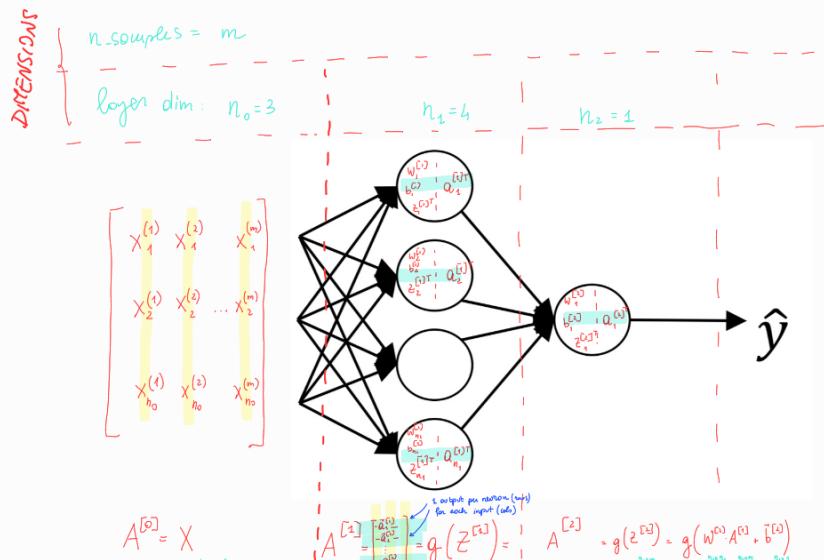
Single neuron



- Assume that each neuron operates a linear combination z of the inputs \vec{x} , and then applies a nonlinear function $g(z)$.

- The linear combination uses the parameters \vec{w} and \vec{b} .
- the parameters are referred to as weights \vec{w} and offsets \vec{b} (one for each input / edge), the nonlinear function as "activation function", while the output of the neuron as "activation value" of the neuron.
- Functionally the neuron is $a = g(z) = g(w^T \cdot x + b)$.

2-layer network 2 many input samples



$$\begin{aligned}
 A^{[0]} &= X \\
 &\quad n_0 \times 1 \\
 A^{[1]} &= g\left(Z^{[1]}\right) = g\left(W^{[1]} \cdot A^{[0]} + b^{[1]}\right) \\
 &= g\left(W^{[1]} \cdot A^{[0]} + b^{[1]}\right) \\
 &= g\left(\begin{bmatrix} -W_{11}^{[1]} \\ -W_{21}^{[1]} \\ -W_{31}^{[1]} \end{bmatrix} \cdot \begin{bmatrix} \vec{a}^{(1)} \\ \vec{a}^{(2)} \\ \vdots \\ \vec{a}^{(n)} \end{bmatrix} + \begin{bmatrix} b_{11}^{[1]} \\ b_{21}^{[1]} \\ b_{31}^{[1]} \end{bmatrix}\right)
 \end{aligned}$$

Annotations: "n samples = m", "layer dim: $n_0 = 3$ ", "n₁ = 4", "n₂ = 1". Labels: "n₀ rows", "n₁ rows", "n₂ rows", "n₀ columns", "n₁ columns", "n₂ columns".

- assume m input samples with dimension n_0 .
- assume the dimension (= no. neurons) of the layer to be n_k .
- with $\text{var}_k^{[j][i]}$ we refer to a variable that relates to the k^{th} neuron of the i^{th} layer for the j^{th} sample.
- intuition: the matrices that represent layers activations have 1 row for each neuron, and 1 column for each input sample.
- intuition: the parameter matrices (relate to edges) have 1 row for each neuron, and 1 column for each neuron of the previous layer.

- logistic regression loss $L(a, y) = -[y \log(a) + (1-y) \log(1-a)]$
- Sigmoid activation function $g(z) = \frac{1}{1+e^{-z}}$
- Derivatives in the non-vectorized case:
 - $\frac{\partial L}{\partial z} = g(z) \cdot (1-g(z))$; $\frac{\partial \log(a)}{\partial z} = \frac{1}{a}$
 - $\frac{\partial L}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a}$
 - $\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = \left[-\frac{y}{a} + \frac{1-y}{1-a}\right] \cdot [a(1-a)] = a - y$
 - $\frac{\partial L}{\partial W} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial W} = (a^T \cdot y) \cdot a^{T-1}$
 - $\frac{\partial L}{\partial b} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial b} = a^T \cdot y$

- vectorized procedure:
 - REPEAT UNTIL CONVERGENCE:
 - Given $X = A^{[0]}$ ($n_0 \times m$) and y ($1 \times m$), init $W^{[1]} = \vec{0}$ ($n_1 \times n_0$), $b^{[1]} = \vec{0}$ ($n_1 \times 1$), $W^{[2]} = \vec{0}$ ($n_2 \times n_1$), $b^{[2]} = \vec{0}$ ($n_2 \times 1$)
 - F.P. $A^{[1]} = g\left(W^{[1]} \cdot A^{[0]} + b^{[1]}\right)$
 - $A^{[2]} = g\left(W^{[2]} \cdot A^{[1]} + b^{[2]}\right)$
 - we may skip the explicit computation of the loss ($L(A^{[0]}, y) = y \log(A^{[0]}) + (1-y) \log(1-A^{[0]})$) because we already know its gradient w.r.t z !
 $\frac{\partial L}{\partial z} = A^{[0]} y$ ($n_0 \times m$) each col is the grad generated by one sample
 $\frac{\partial L}{\partial z} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial W^{[1]}} = dL/dz$ ($n_2 \times m$) but we want the mean grad over all samples so = np.mean((dL/dz) , axis=1, keepdims=True) $n_2 \times 1$
 - B.P. $dW^{[1]} = \frac{\partial L}{\partial W^{[1]}} = \frac{\partial L}{\partial z} \cdot A^{[0]T} \cdot A^{[0]T}$ ($n_2 \times n_0$) again, we want the average gradient, but in this case the cols of dL/dz have already been summed up by the dot product, so to obtain the average it is enough to divide by m
 $dL/dz \cdot A^{[0]T} \cdot A^{[0]T}$ ($n_2 \times n_0$) $= W^{[0]T} \cdot dL/dz \otimes A^{[0]T}$ ($n_2 \times n_0$)
 - $db^{[1]} = \frac{\partial L}{\partial b^{[1]}} = \frac{\partial L}{\partial z} \cdot A^{[0]T}$ ($n_2 \times 1$) we want to average again = $\frac{1}{m} dL/dz \cdot A^{[0]T}$ ($n_2 \times 1$)
 - $db^{[2]} = np.mean(dL/dz, axis=1, keepdims=True)$
 - G.D. $W^{[1]} = -\alpha dW^{[1]}$; $b^{[1]} = -\alpha db^{[1]}$; $W^{[2]} = -\alpha dW^{[2]}$; $b^{[2]} = -\alpha db^{[2]}$

GENERAL FORM OF FP & BP (for the general layer l) , provided in the last slide of week 4.

Forward and backward propagation

FP

$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ Z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(Z^{[2]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) = \hat{Y} \end{aligned}$$

BP

$$\begin{aligned} dZ^{[L]} &= A^{[L]} - Y \\ dW^{[L]} &= \frac{1}{m} dZ^{[L]} A^{[L]T} \\ db^{[L]} &= \frac{1}{m} np.\text{sum}(dZ^{[L]}, \text{axis} = 1, \text{keepdims} = \text{True}) \\ * \quad dZ^{[L-1]} &= dW^{[L]T} dZ^{[L]} g'^{[L]}(Z^{[L-1]}) \\ dZ^{[1]} &= dW^{[L]T} dZ^{[2]} g'^{[1]}(Z^{[1]}) \\ dW^{[1]} &= \frac{1}{m} dZ^{[1]} A^{[1]T} \\ db^{[1]} &= \frac{1}{m} np.\text{sum}(dZ^{[1]}, \text{axis} = 1, \text{keepdims} = \text{True}) \end{aligned}$$

* Some formulas in the BP are not explained , they basically say to trust them or try verify it by yourself.

WHICH ACTIVATIONS TO USE:

- In general use ReLU (or maybe try leaky ReLU) because they it is less affected by the vanishing gradient problem (at least for large inputs). The derivative of the ReLU although analytically incorrect , works well in practice .

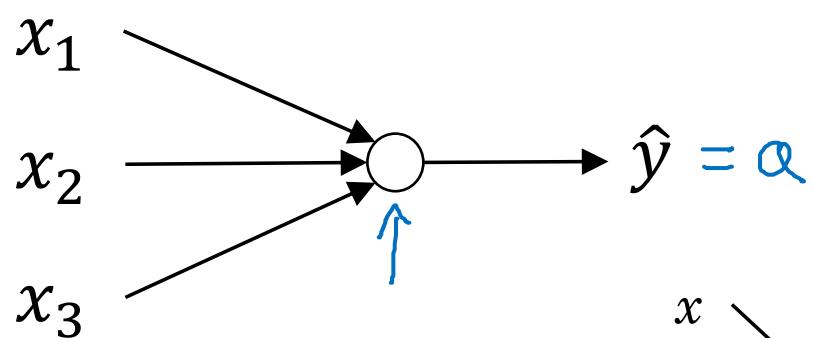
$$\frac{d \max(0, x)}{dx} := \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

- Never use linear activation : the resulting NN would be linear and boil down to simple linear regression/classification.
- A sigmoid activation can be used in the last layer if binary classification is desired .

HOW TO INITIATE THE WEIGHTS :

- To small random number , especially if a sigmoid activation is used anywhere in the network

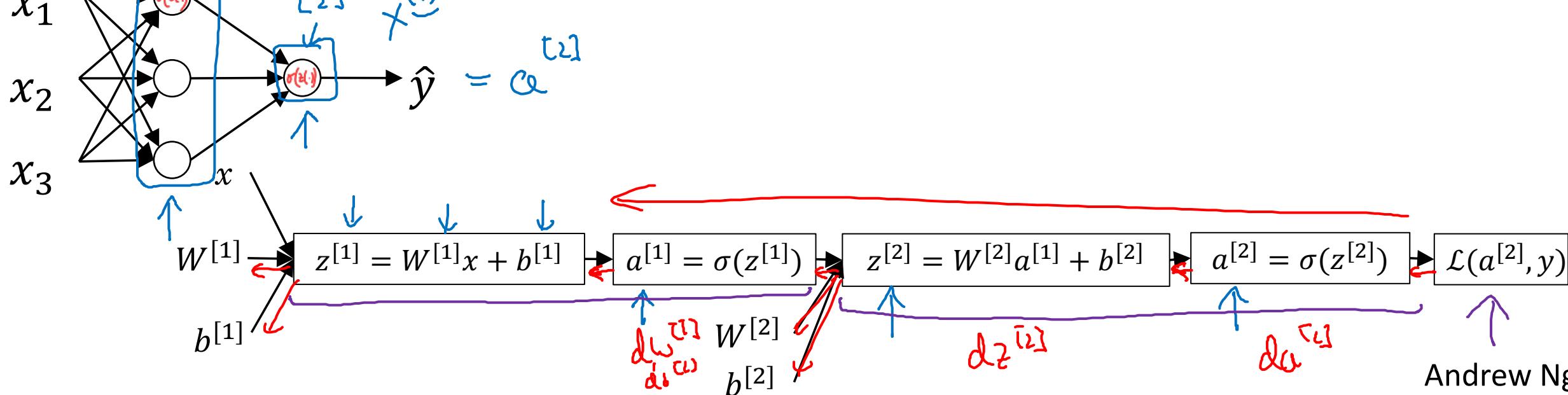
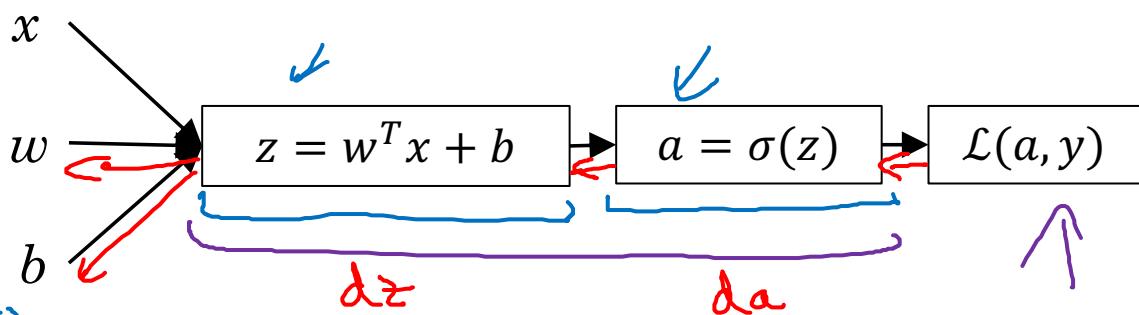
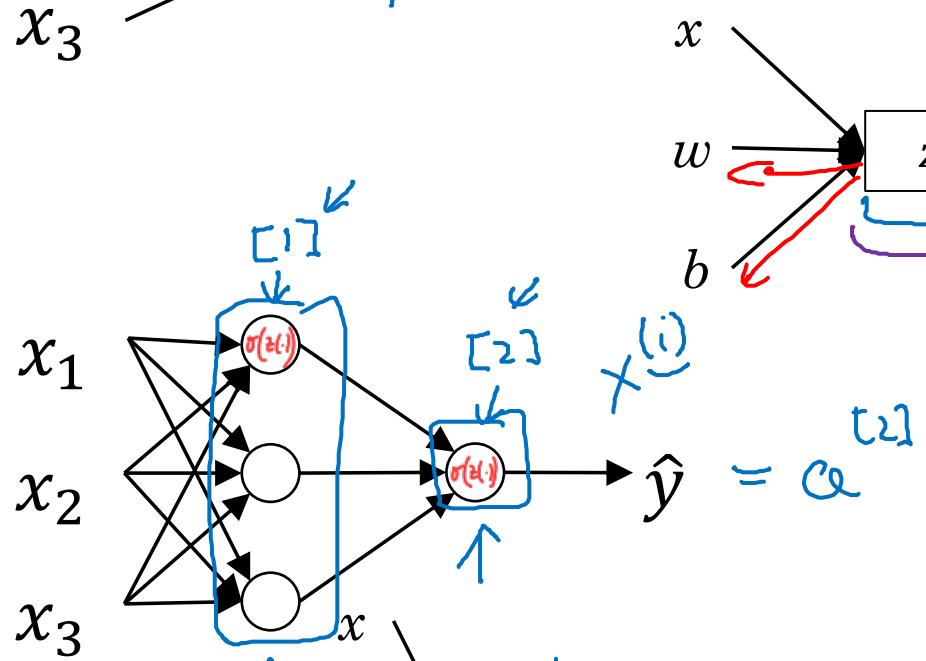
What is a Neural Network?



- Notation :

- superscript $[i]$ refers to the i :th layer
- superscript (i) refers to the i :th training sample

- They use sigmoid neurons : each neuron implements $\sigma(w^T x + b)$



Andrew Ng

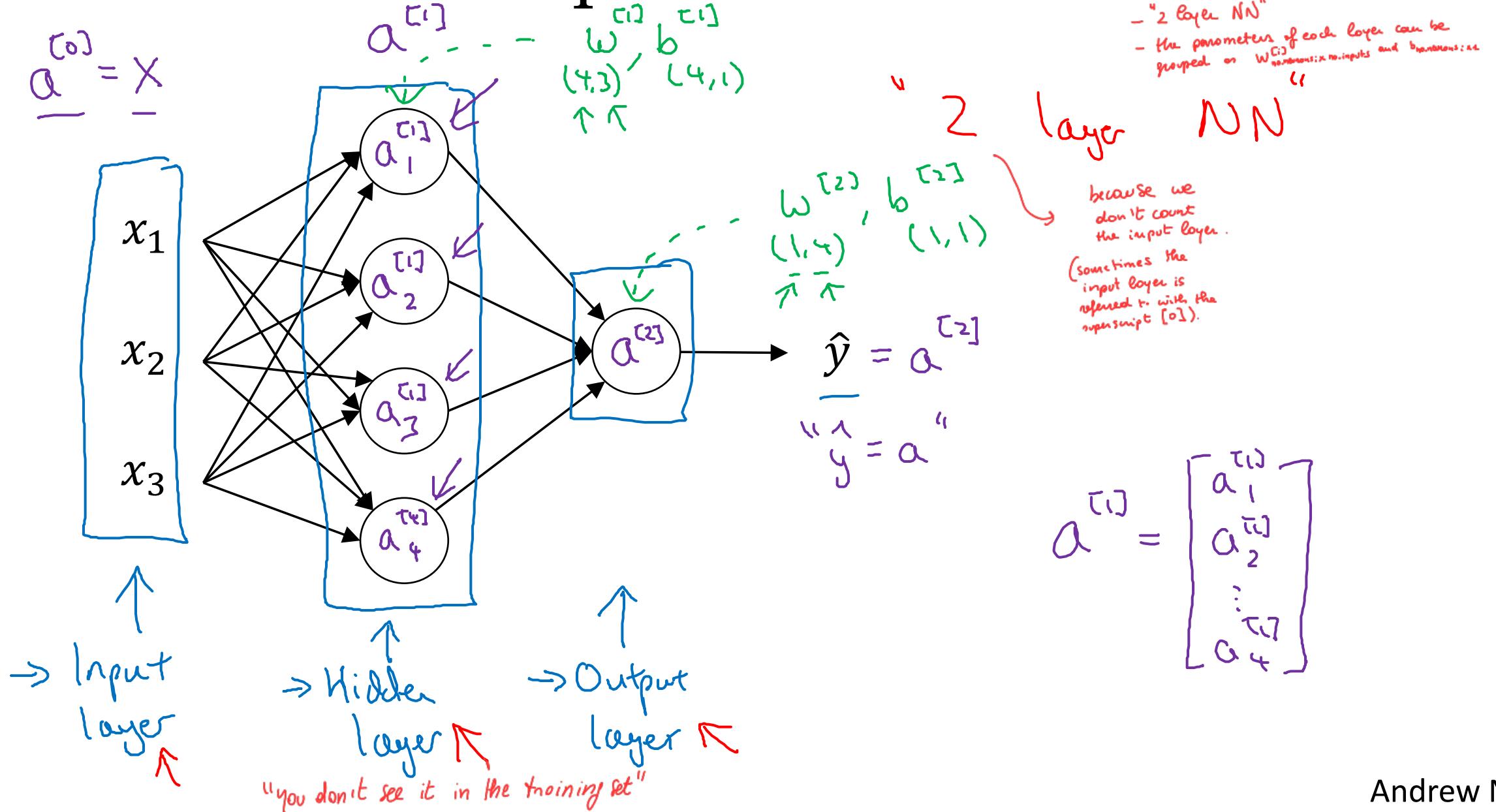


deeplearning.ai

One hidden layer
Neural Network

Neural Network
Representation

Neural Network Representation



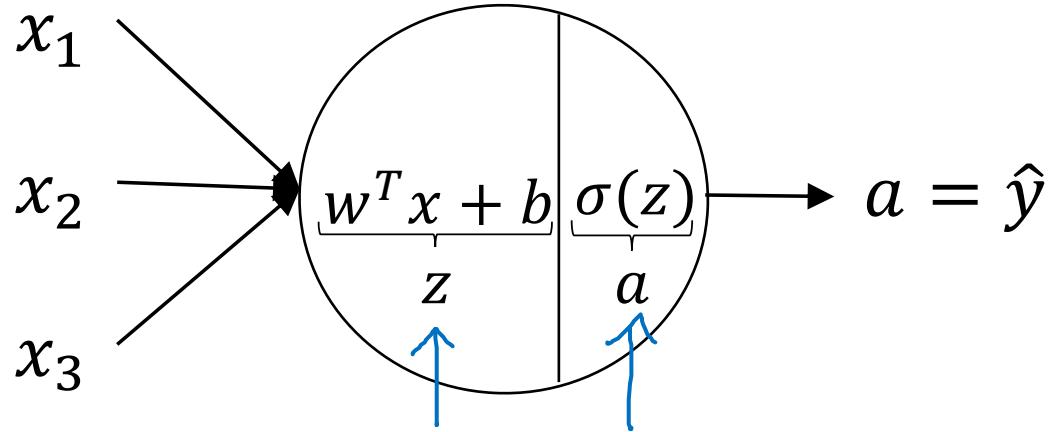


deeplearning.ai

One hidden layer Neural Network

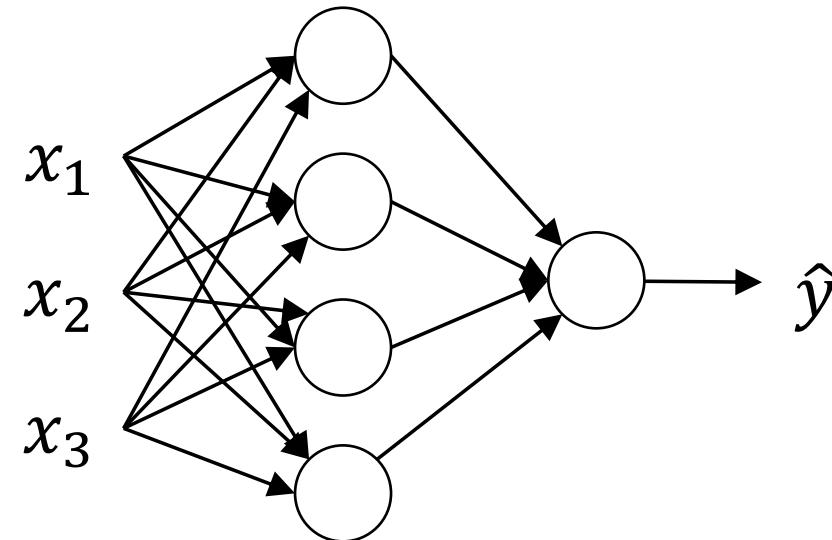
Computing a Neural Network's Output

Neural Network Representation

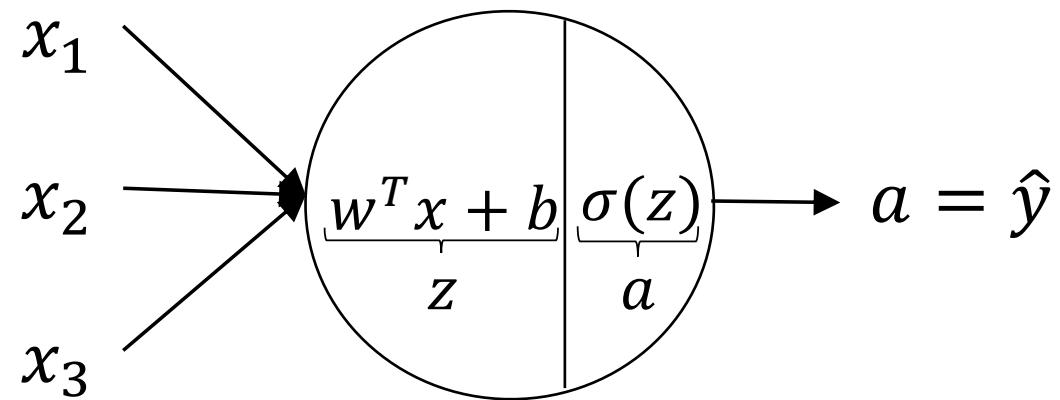


$$z = w^T x + b$$

$$a = \sigma(z)$$

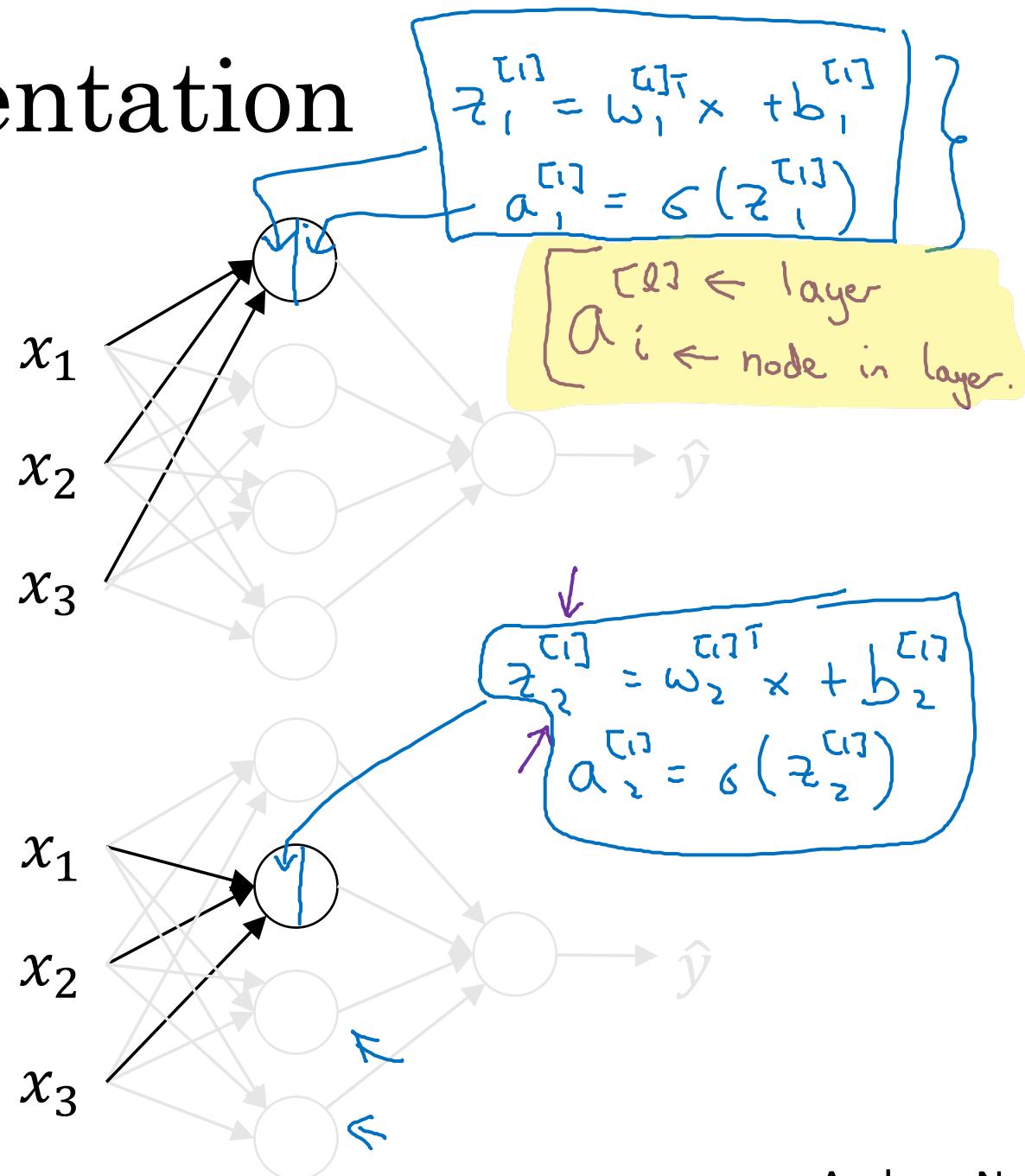


Neural Network Representation



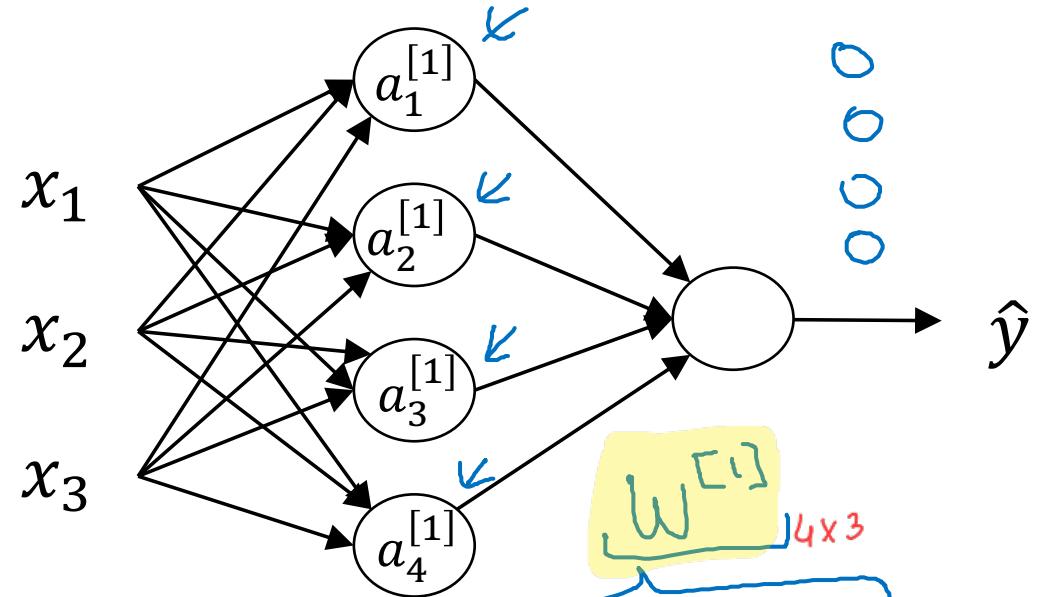
$$z = w^T x + b$$

$$a = \sigma(z)$$



- how to vectorize the operations of the neurons of one layer.

Neural Network Representation

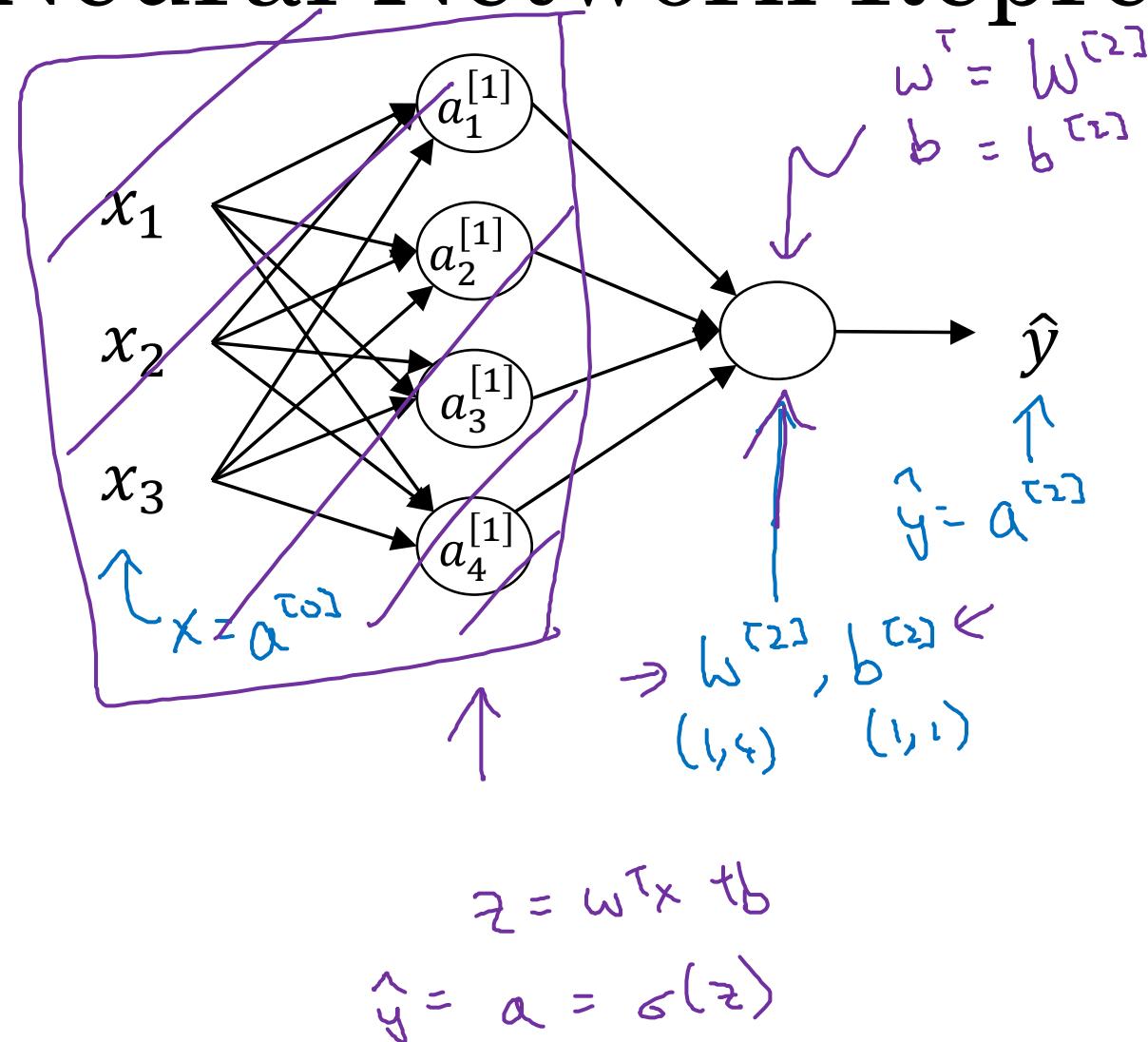


$$\rightarrow z^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$$
$$\rightarrow a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\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

Handwritten annotations explain the vectorized representation of the layer. The input x is multiplied by the transpose of the weight matrix $(w_i^{[1]T})x$ and passed through the activation function $\sigma(z)$ to produce the output $a_1^{[1]}$. This process is repeated for all neurons in the layer, resulting in the vector $a^{[1]}$.

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

Neural Network Representation learning



Given input x :

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

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

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

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

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

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

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

$(1,1)$ $(1,1)$

vectorized representation of the nn on the left.



deeplearning.ai

$X = A^{[s]}$ are $n_x \times m$
↑ input dim ↑ no. samples

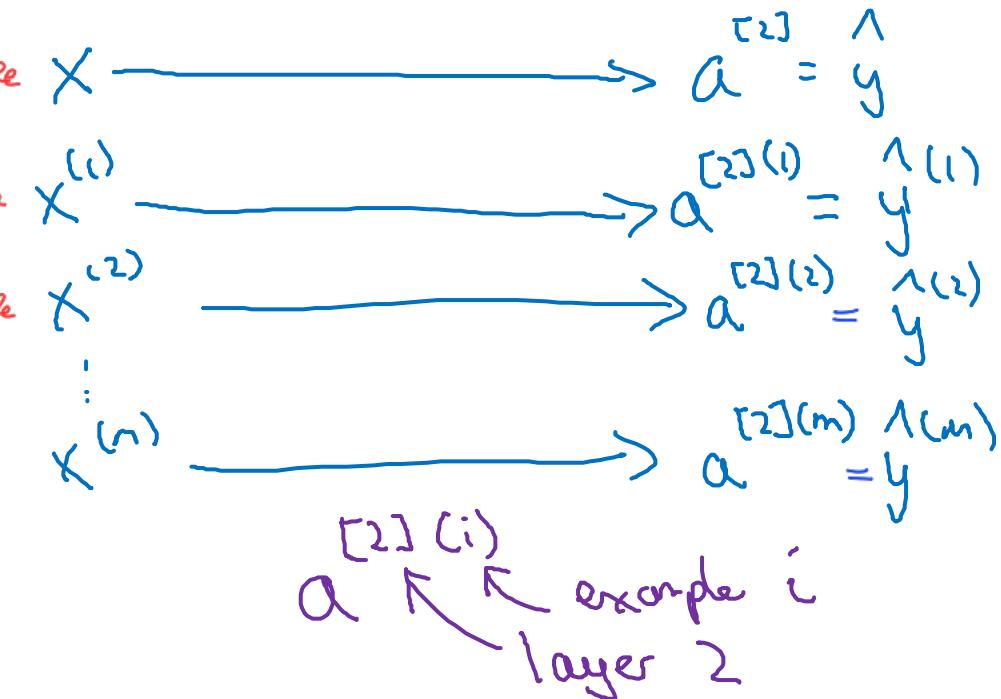
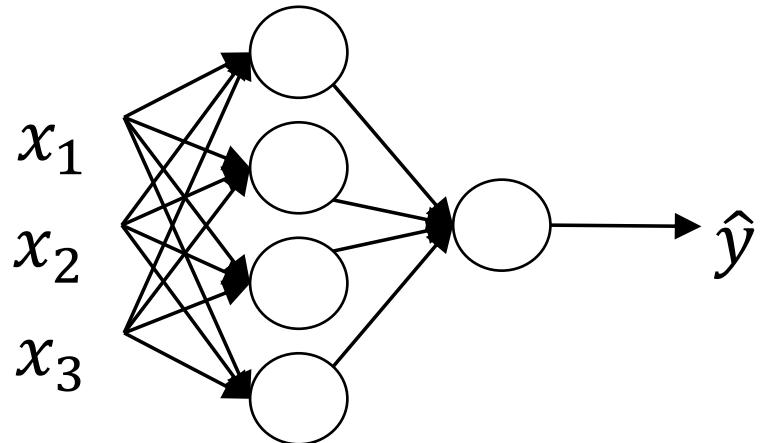
One hidden layer Neural Network

Vectorizing across
multiple examples

$Z^{[i]}, A^{[i]}$ are no-neurons^[i] $\times m$

$W^{[i]} = \begin{bmatrix} -w^{[i]T} \end{bmatrix}$ is no-neurons^[i] $\times n_x$, $b^{[i]} = \begin{bmatrix} b^{[i]} \end{bmatrix}$ is no-neurons^[i] $\times 1$

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 ,

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

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

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

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

Vectorizing across multiple examples

for $i = 1$ to m :

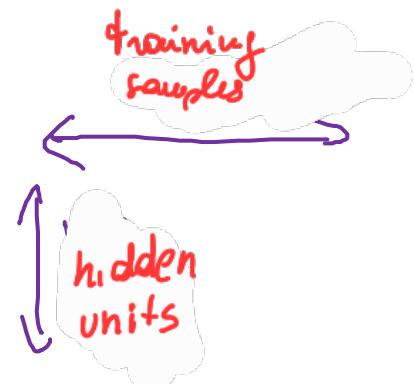
$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$

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

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

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

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



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

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

middle units

A diagram showing the vectorized representations of hidden units and output units. It consists of two large grey cloud-like shapes. The top cloud is labeled "middle units" and contains several smaller circles. The bottom cloud is labeled "output units" and also contains smaller circles. Arrows point from the text labels to their respective cloud shapes. Below each cloud is a corresponding matrix equation where the columns represent the individual hidden or output units for each training sample.



deeplearning.ai

One hidden layer Neural Network

Explanation for vectorized implementation

Justification for vectorized implementation

$$\underline{z}^{1} = \omega^{[1]} \underline{x}^{(1)} + \cancel{v^{[1]}} ,$$

$$\underline{z}^{(1)(2)} = \underline{\omega}^{(1)(2)} \underline{x}^{(2)} + \underline{b}$$

$$\underline{z}^{(1)(3)} = \omega^{(1)} x^{(3)} + b^{(1)}$$

$$\omega^{(i)} = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}$$

$$\omega^{[1]} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$\omega^{(i)} x^{(i)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

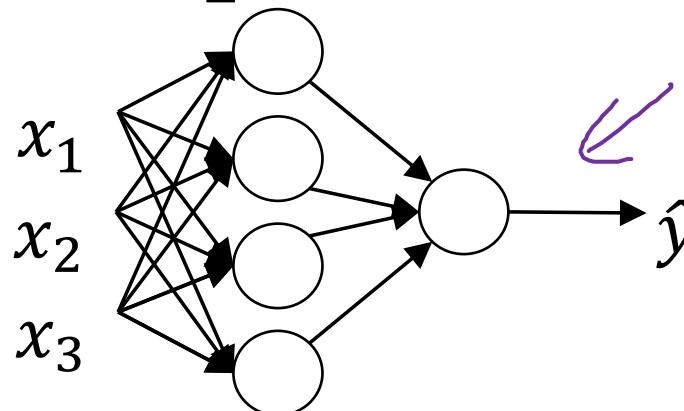
$$\omega^{(1)} x^{(3)} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$z^{[i]} = \bar{w}^{[i]} X + b^{[i]}$$

$\bar{w}^{[i]}$ $\left[\begin{array}{c|c|c|c} 1 & | & | & | \\ X^{(1)} & X^{(2)} & X^{(3)} \dots & | \\ | & | & | & \dots \end{array} \right] = \left[\begin{array}{cccc} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{array} \right] = \left[\begin{array}{c|c|c|c} 1 & | & | & | \\ z^{1} & z^{[1](2)} & z^{[1](3)} \dots & | \\ | & | & | & \dots \end{array} \right] = z^{[i]}$

\cancel{X} \uparrow $w^{[i]} x^{[i]} = z^{[i]}$

Recap of vectorizing across multiple examples



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

A purple arrow points from the text "vectorizing across multiple examples" to this equation.

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

A purple arrow points from the text "vectorizing across multiple examples" 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))
```

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

Handwritten annotations: $A^{[1]}(i)$ is circled in blue. To its right, $x = a^{[1]}$ is written above $x^{(i)} = a^{[1](i)}$. Below the equations, there are yellow brackets grouping the operations for each layer.

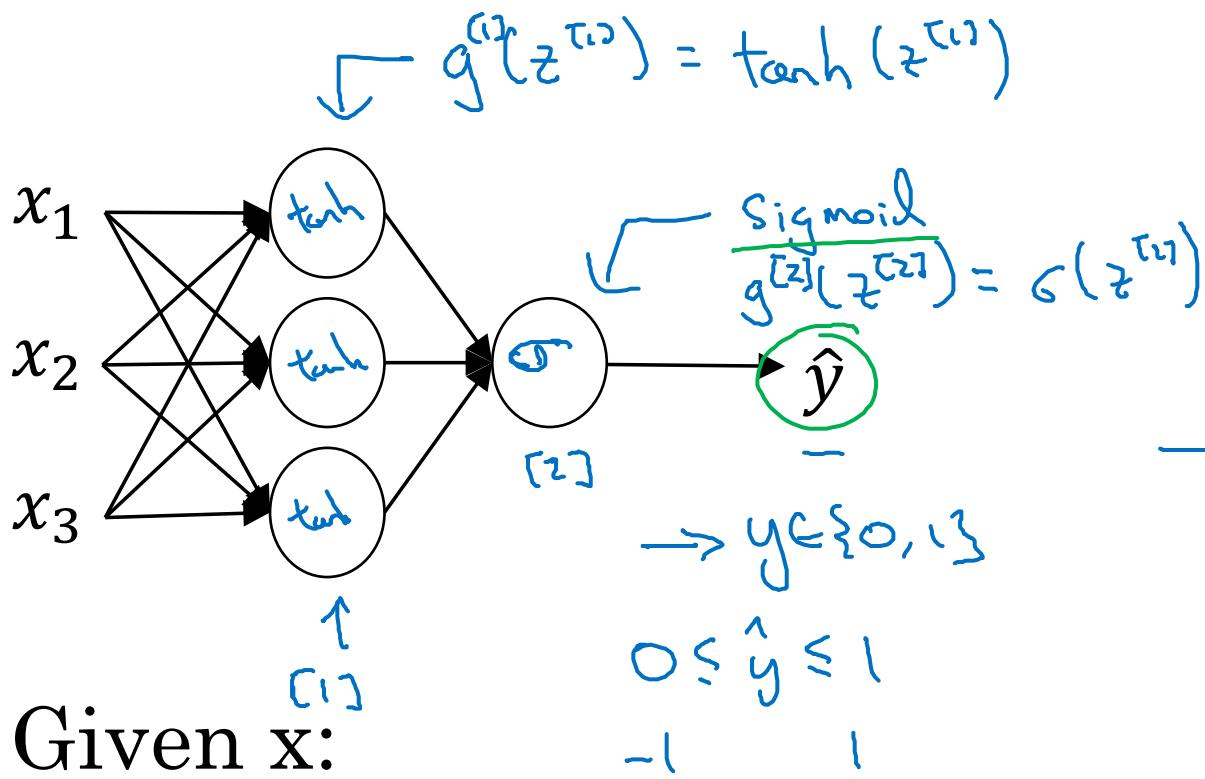


deeplearning.ai

One hidden layer
Neural Network

Activation functions

Activation functions

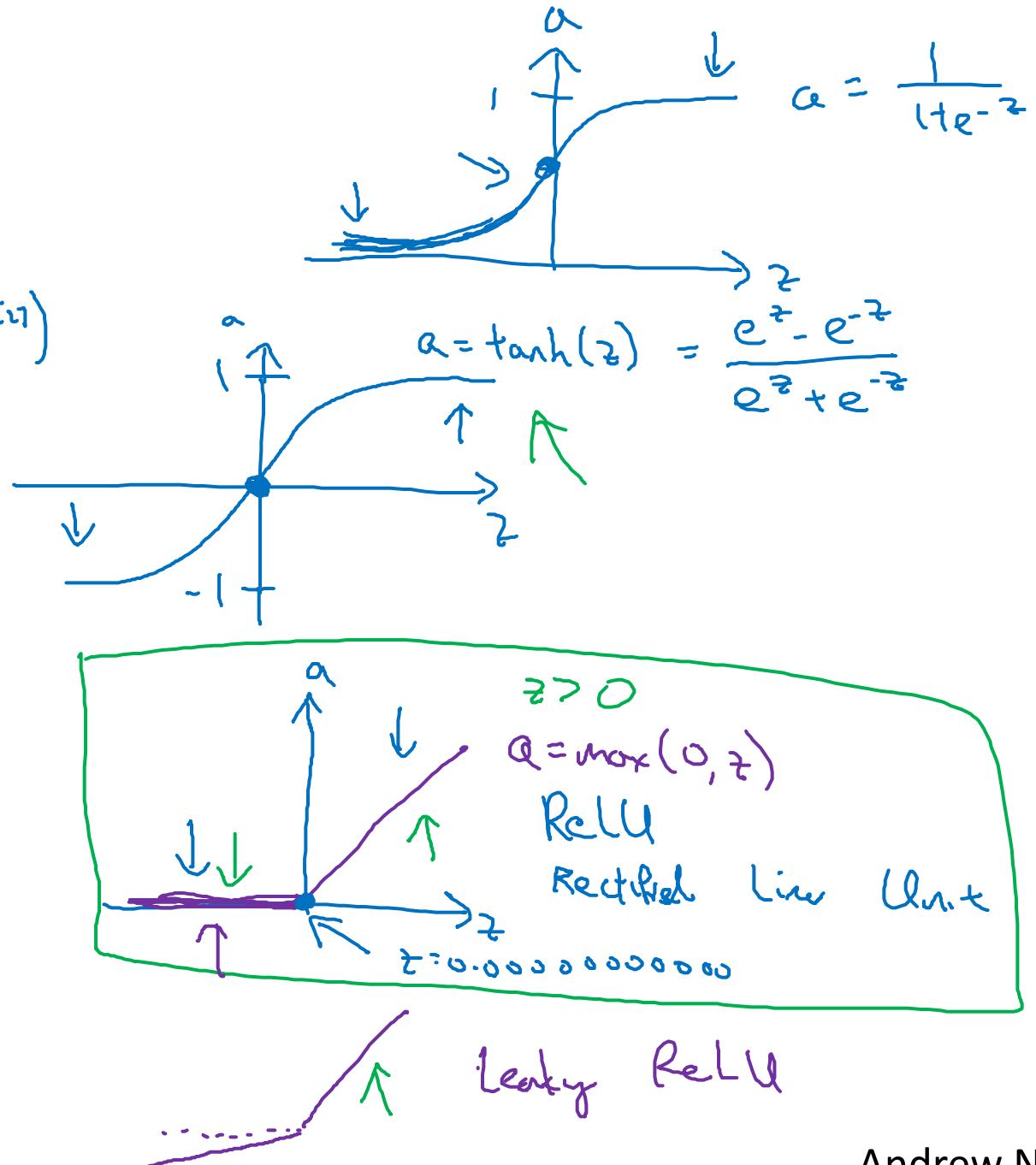


$$z^{[1]} = W^{[1]}x + b^{[1]}$$

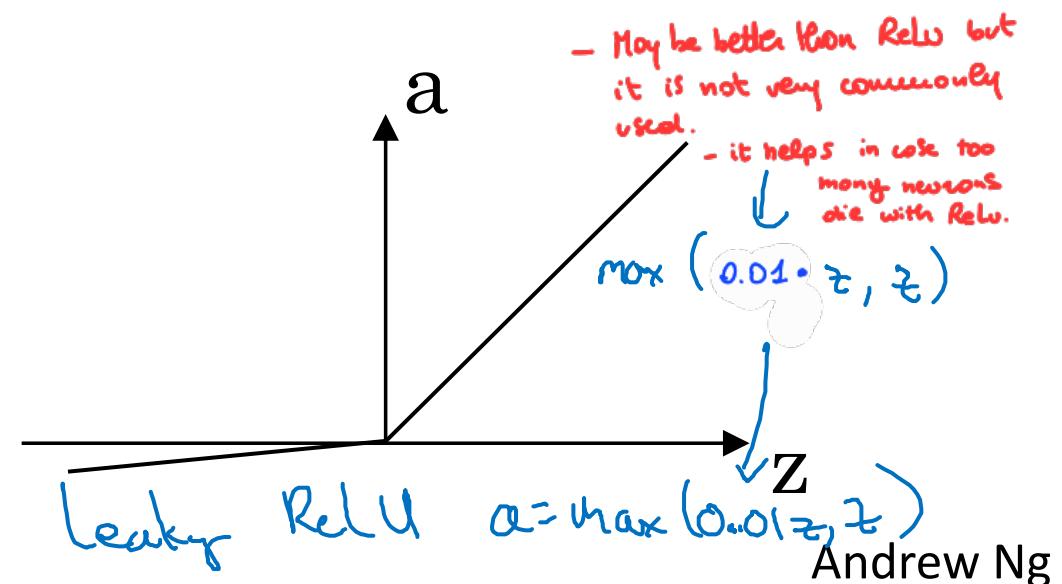
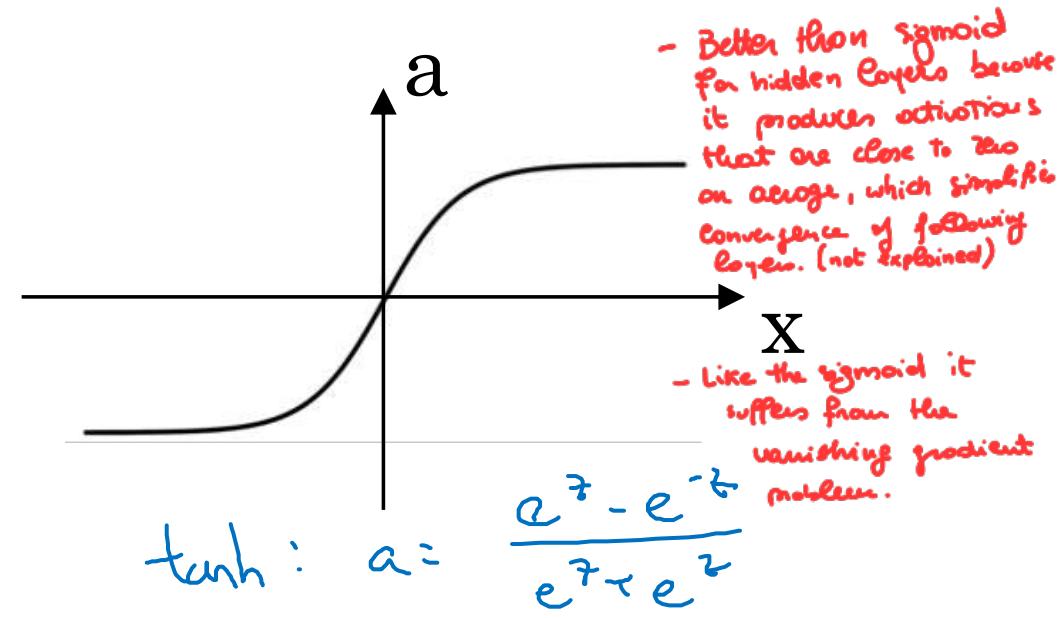
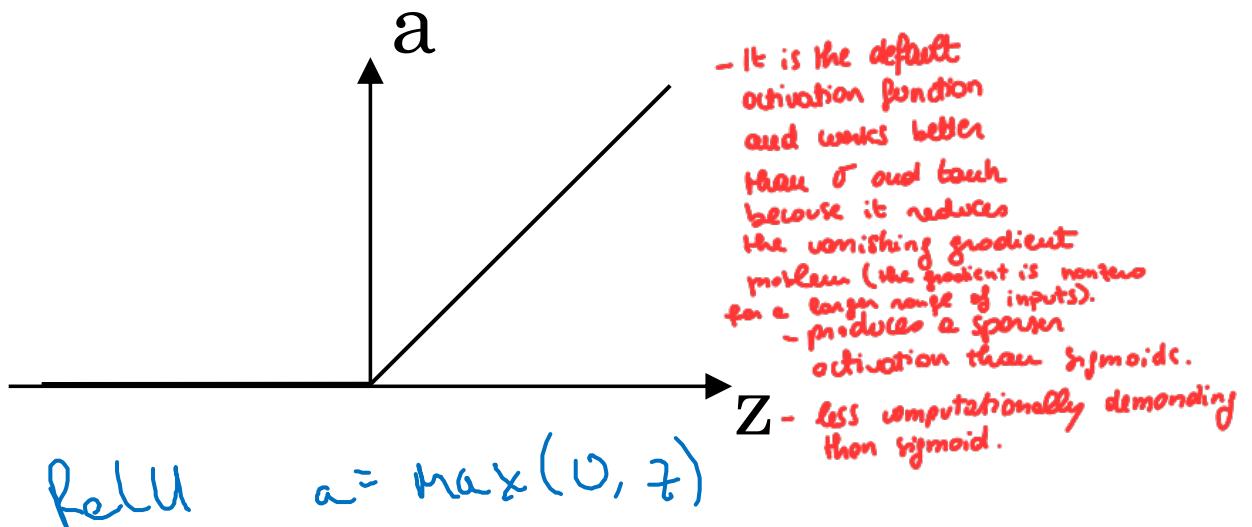
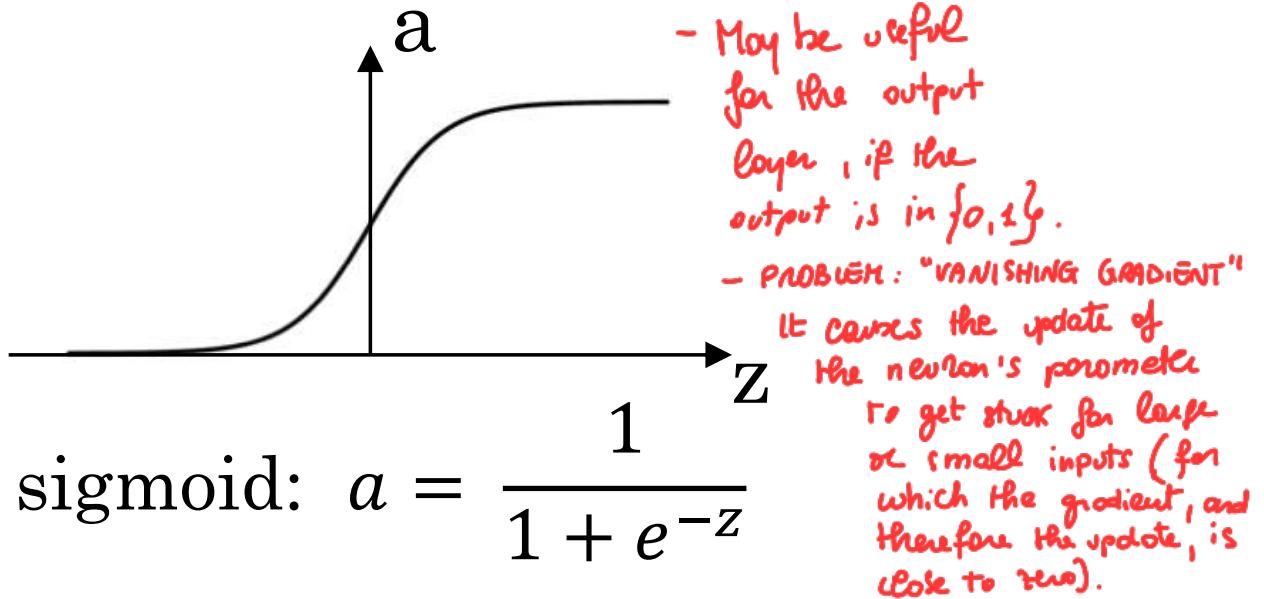
$$\rightarrow a^{[1]} = \sigma(\cancel{z^{[1]}}) \quad \cancel{g^{[1]}(z^{[1]})}$$

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

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



Pros and cons of activation functions



Because otherwise the NN computes a linear function overall...
(the composition of linear functions is a linear function).

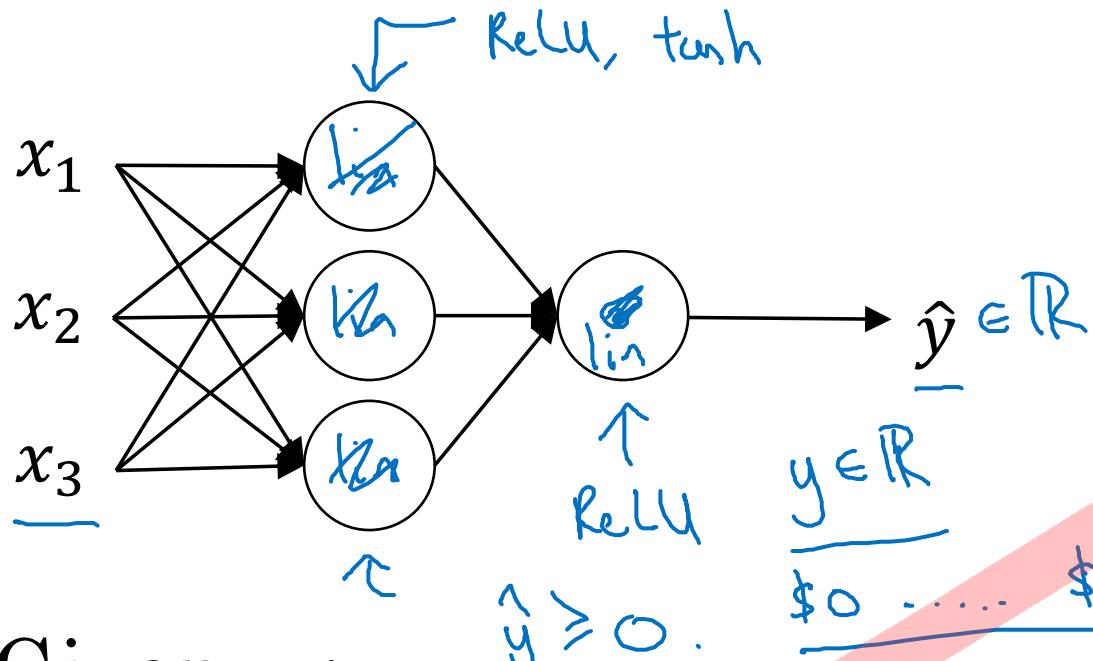


deeplearning.ai

One hidden layer Neural Network

Why do you need non-linear activation functions?

Activation function



Given x :

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

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

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

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

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

$$\begin{aligned} a^{[1]} &= z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[2]} &= z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} &= W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]} \\ &= (\underbrace{W^{[2]}W^{[1]}}_{w'})x + (\underbrace{W^{[2]}b^{[1]} + b^{[2]}}_{b'}) \\ &= \underbrace{w'x}_{\hat{w}} + \underbrace{b'}_{\hat{b}} \\ g(z) &= z \end{aligned}$$

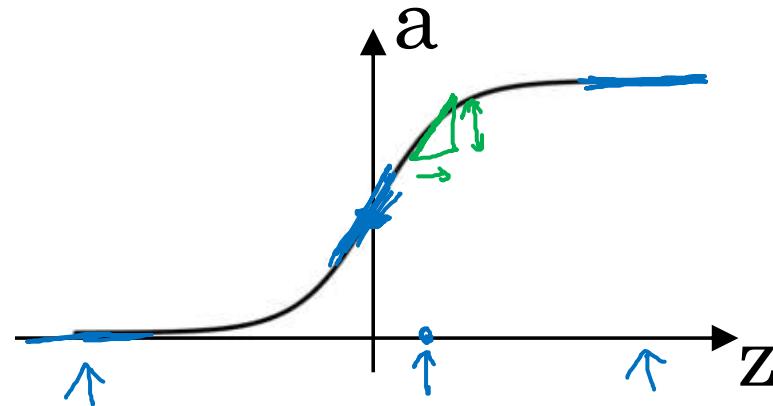


deeplearning.ai

One hidden layer Neural Network

Derivatives of
activation functions

Sigmoid activation function



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

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

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

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

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

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

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

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

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

= slope of $g(x)$ at z

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

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

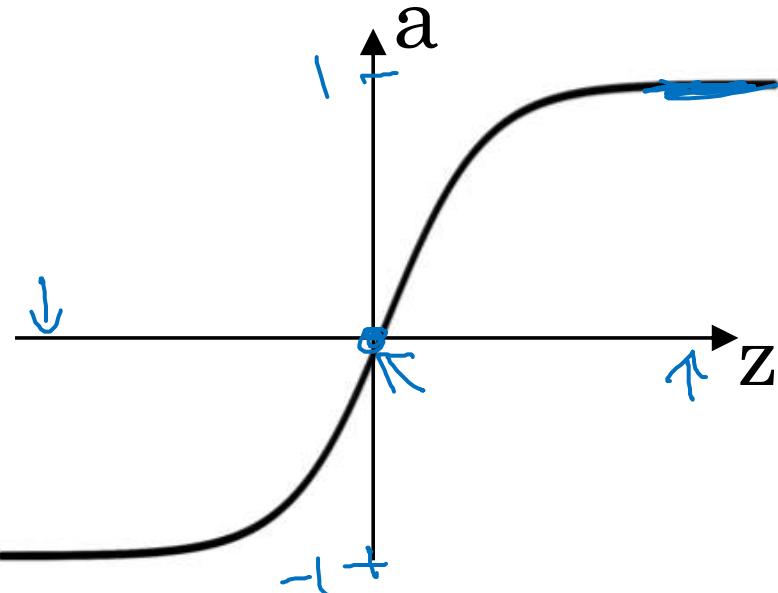
$$g'(z) = a(1-a)$$

being $a = g(z)$

$$= a(1-a)$$

Andrew Ng

Tanh activation function



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

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

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

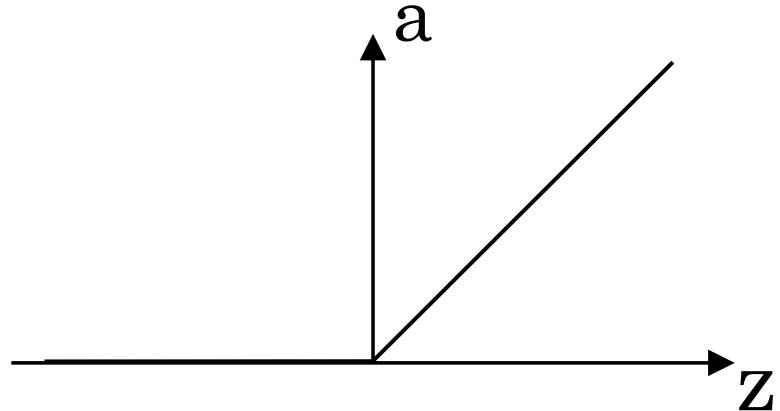
$$= 1 - (\tanh(z))^2$$

$$a = g(z),$$

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

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

ReLU and Leaky ReLU



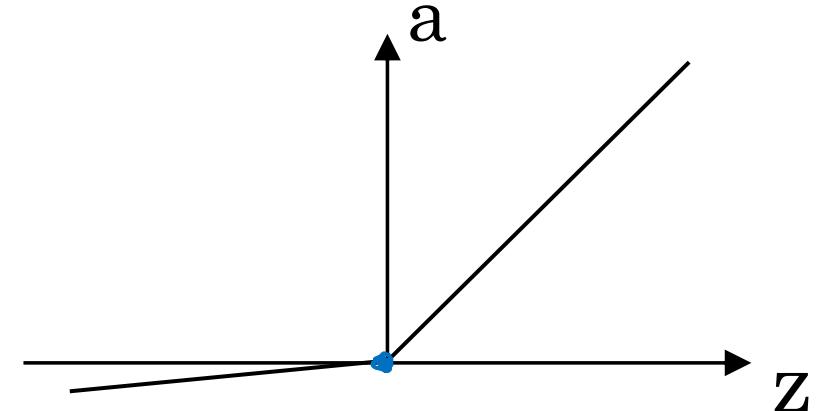
ReLU

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

mathematically, the derivative is not defined in $t=0$. In practice, however, the chance to have β exactly = 0 is so low that we can approximate this with the \geq sign.

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

$z = 0.0000\cdots 0$



Leaky ReLU

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

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



deeplearning.ai

One hidden layer
Neural Network

Gradient descent for
neural networks

Gradient descent for neural networks

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

Activations dimensions:

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

Cost function: $J(\underline{\omega}^{[1]}, \underline{b}^{[1]}, \underline{\omega}^{[2]}, \underline{b}^{[2]}) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}, y)$

$$\leftarrow a^{[2]}$$

Gradient Descent (VECTORIZED FORM)

→ Repeat until convergence {

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

$$\frac{\partial J}{\partial \omega^{[1]}} = \frac{\partial J}{\partial \omega^{[1]}} , \quad \frac{\partial J}{\partial b^{[1]}} = \frac{\partial J}{\partial b^{[1]}} , \quad \frac{\partial J}{\partial \omega^{[2]}} , \quad \frac{\partial J}{\partial b^{[2]}}] \text{ compute derivatives}$$

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

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

↳

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

update the parameters

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:

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

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

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

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

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

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

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

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

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



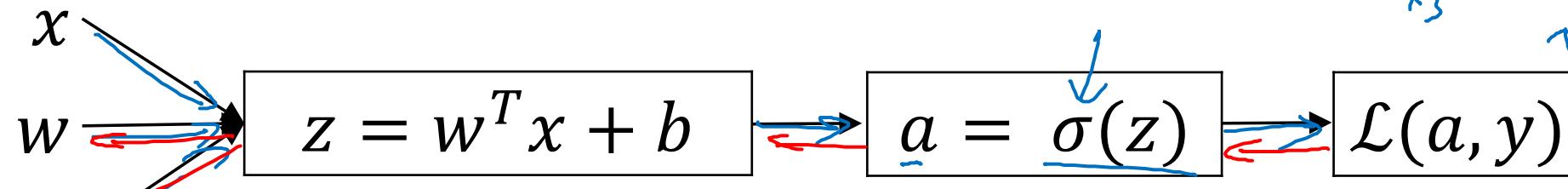
deeplearning.ai

One hidden layer
Neural Network

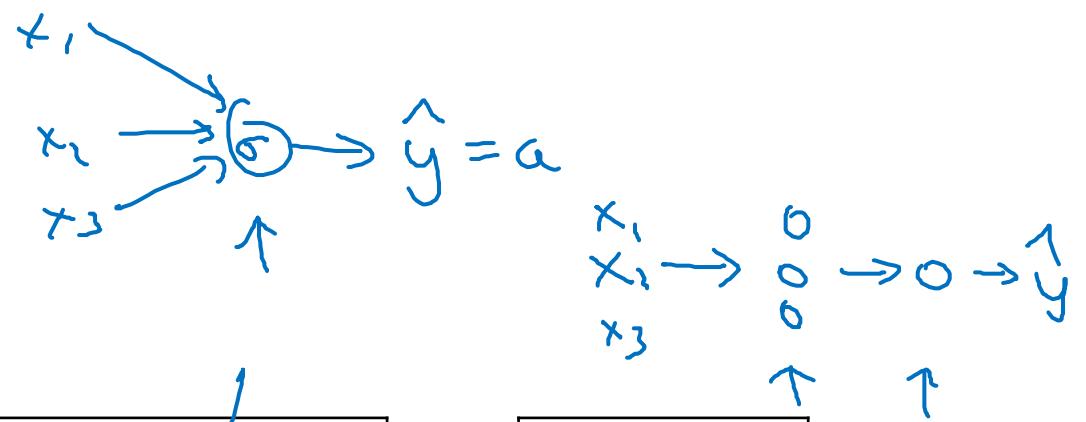
Backpropagation
intuition (Optional)

Computing gradients

Logistic regression



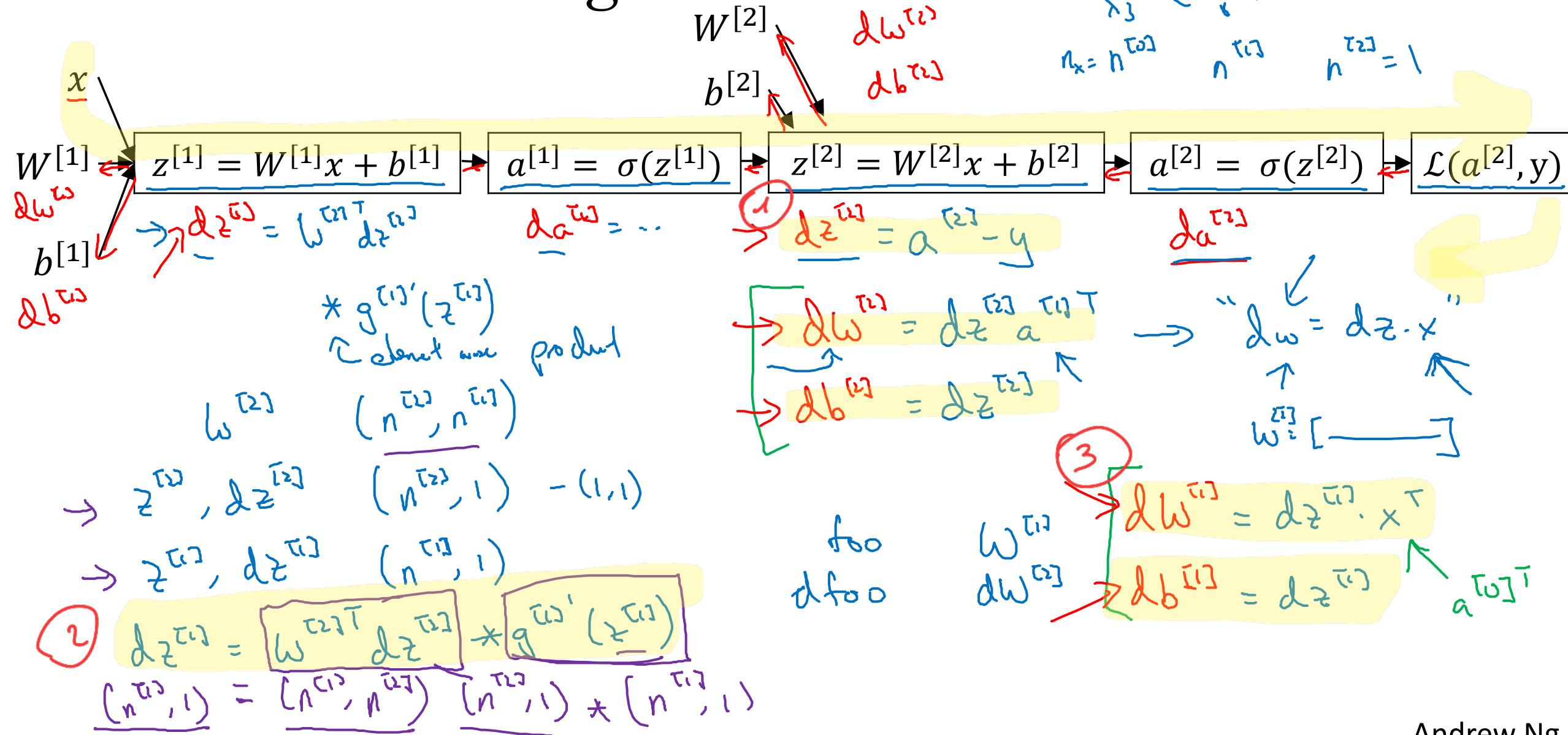
$$\begin{aligned} \frac{\partial z}{\partial w} &= x \\ \frac{\partial z}{\partial b} &= 1 \\ \frac{\partial z}{\partial a} &= g'(z) \\ g(z) &= \sigma(z) \\ \frac{\partial \mathcal{L}}{\partial z} &= \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial a}{\partial z} \\ "dz" &= "da" \end{aligned}$$



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

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

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

Vectorized implementation

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

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

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

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

(n^{T₁}, 1)

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

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

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

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

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

$$dZ^{[1]} = \underbrace{W^{[2]T} dz^{[2]} * g^{[1]'}(Z^{[1]})}_{\substack{\text{elementwise product} \\ (n^{T_1}, m) \quad (n^{T_2}, m) \quad (n^{T_1}, m)}}$$

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

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

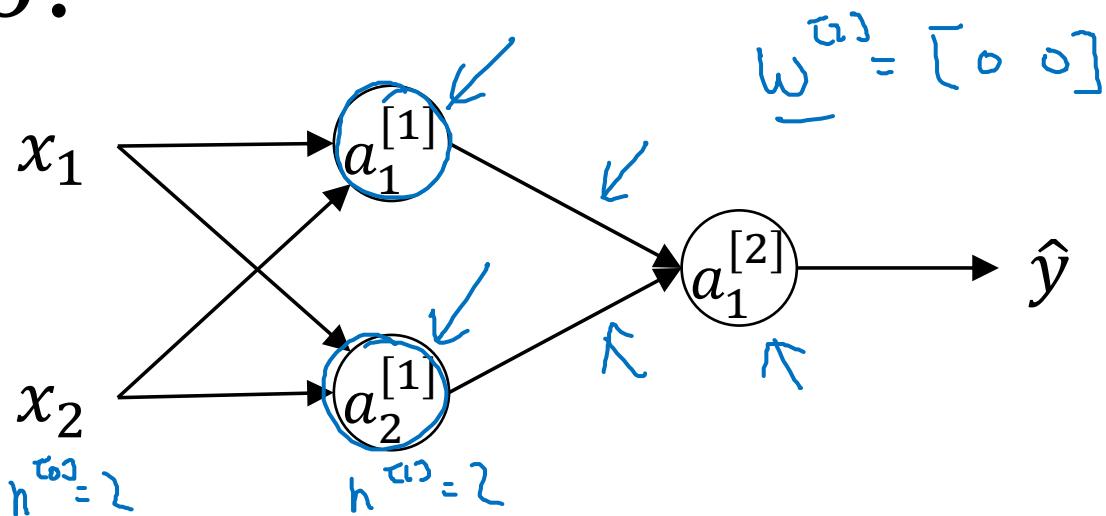


deeplearning.ai

One hidden layer
Neural Network

Random Initialization

What happens if you initialize weights to zero?



$$W^T1 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$b^T1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]}$$

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

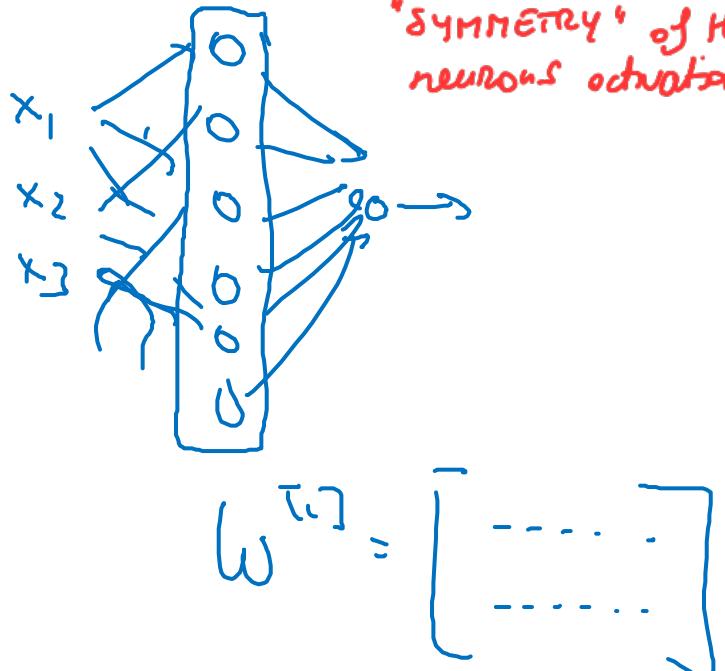
$$\delta w = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

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

- Every neuron in the same layer would have the same activation
=> they would therefore be updated in the same way
=> they would therefore always be identical.

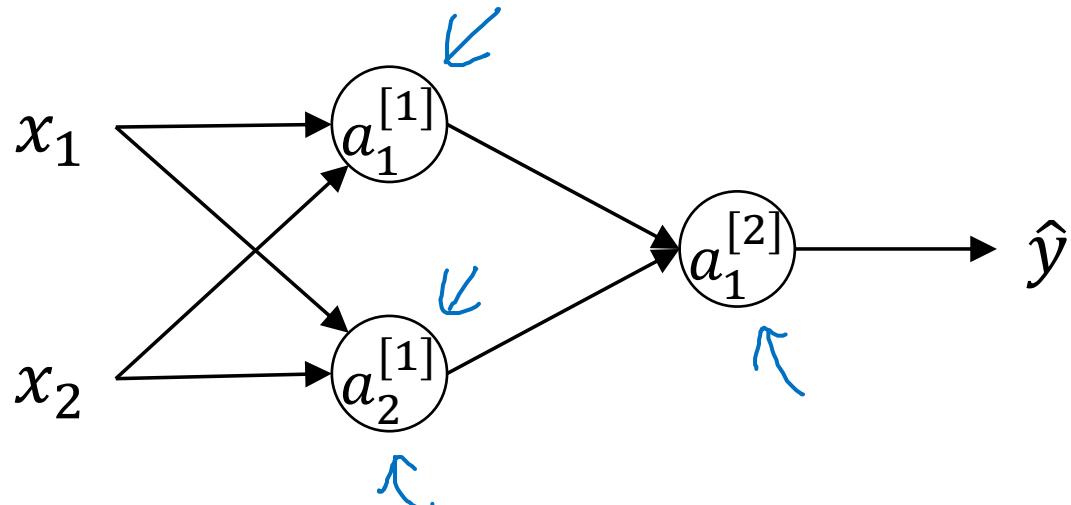
Symmetric

- one talks about "SYMMETRY" of the neurons activations



Random initialization

- Random initialization breaks the symmetry of the neurons activations.



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

- It's often good to scale the random initializations to more than small enough to avoid vanishing gradients (if sigmoid activations are used). This is especially important in shallow networks.

