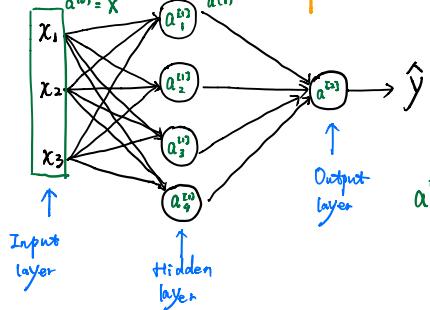


## 1. Neural Network Representation



2 layer Network (don't count input layer)

$$a^{L2} = \begin{bmatrix} a_1^{L2} \\ a_2^{L2} \\ \vdots \\ a_4^{L2} \end{bmatrix}, \quad a^{L2} = X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad a^{L2} = \hat{y}$$

## 2. Computing a NN's output.

$$\text{e.g. } \begin{cases} z_1^{L1} = w_1^{L1T} x + b_1^{L1} \\ a_1^{L1} = \theta(z_1^{L1}) \end{cases} \quad \begin{cases} z_2^{L1} = w_2^{L1T} x + b_2^{L1} \\ a_2^{L1} = \theta(z_2^{L1}) \end{cases} \quad \dots$$

$\downarrow$   
 $w_i^{L1T} \leftarrow \text{layer } i \quad a_i^{L1} \leftarrow \text{node in layer } i$

$$= \theta[w_2^{L1T} x + b_2^{L1}]$$

If we vectorize the results:

$$\begin{bmatrix} w_1^{L1T} \\ w_2^{L1T} \\ w_3^{L1T} \\ w_4^{L1T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{L1} \\ b_2^{L1} \\ b_3^{L1} \\ b_4^{L1} \end{bmatrix}_{(4 \times 1)} = \begin{bmatrix} w_1^{L1T} x + b_1^{L1} \\ w_2^{L1T} x + b_2^{L1} \\ w_3^{L1T} x + b_3^{L1} \\ w_4^{L1T} x + b_4^{L1} \end{bmatrix} = \begin{bmatrix} z_1^{L1} \\ z_2^{L1} \\ z_3^{L1} \\ z_4^{L1} \end{bmatrix} = Z^{L1}$$

$\triangle$

$$a^{L1} = \begin{bmatrix} a_1^{L1} \\ \vdots \\ a_4^{L1} \end{bmatrix} = \theta(Z^{L1})$$

$\triangle$

Given input  $X$ :

$$\begin{cases} z^{L1} = W^{L1} X + b^{L1} \\ a^{L1} = \theta(z^{L1}) \\ z^{L2} = W^{L2} X + b^{L2} \\ a^{L2} = \theta(z^{L2}) \end{cases}$$

## 3. Vectorizing across multiple examples

$$\begin{cases} X \longrightarrow a^{L2} = \hat{y} \\ X^{(1)} \longrightarrow a^{L2(1)} = \hat{y}^{(1)} \\ X^{(2)} \longrightarrow a^{L2(2)} = \hat{y}^{(2)} \\ \vdots \\ X^{(m)} \longrightarrow a^{L2(m)} = \hat{y}^{(m)} \end{cases}$$

$a^{L2(i)} \leftarrow \text{example } i$

$\text{layer 2}$

for  $i = 1$  to  $m$ ,  
 $z^{L2(i)} = W^{L2} X^{(i)} + b^{L2}$   
 $a^{L2(i)} = \theta(z^{L2(i)})$   
 $z^{L2(i)} = W^{L2} a^{L1(i)} + b^{L2}$   
 $a^{L2(i)} = \theta(z^{L2(i)})$

$$X = \begin{bmatrix} X^{(1)} & X^{(2)} & \dots & X^{(m)} \end{bmatrix}_{n \times m}$$

↑ training examples  
hidden units

$$\left\{ \begin{array}{l} Z^{(1)} = W^{(1)} X + b^{(1)} \\ A^{(1)} = \theta(Z^{(1)}) \\ Z^{(2)} = W^{(2)} A^{(1)} + b^{(2)} \\ A^{(2)} = \theta(Z^{(2)}) \end{array} \right.$$

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

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

4. Explanation for vectorized implementation

$$Z^{(1)(1)} = W^{(1)} X^{(1)} + b^{(1)}, Z^{(1)(2)} = W^{(1)} X^{(2)} + b^{(1)}, Z^{(1)(3)} = W^{(1)} X^{(3)} + b^{(1)} \quad (3 \text{ training examples})$$

$$W^{(1)} = \begin{bmatrix} \dots \\ \dots \end{bmatrix}, W^{(1)} X^{(1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}, W^{(1)} X^{(2)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}, W^{(1)} X^{(3)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$W^{(1)} \begin{bmatrix} X^{(1)} & X^{(2)} & X^{(3)} & \dots \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} Z^{(1)(1)} & Z^{(1)(2)} & Z^{(1)(3)} \end{bmatrix} = Z^{(1)}$$

More training examples, more columns.

### 5. Activation functions

$$a^{(1)} = \theta(z^{(1)}) \rightarrow a^{(1)} = g(z^{(1)})$$

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

For the input layer to hidden layer, tangent function always works much better.

However, if you want to do binary classification, since  $y$  should  $\in [0, 1]$ , this time we use sigmoid function.

$$\begin{cases} a^{(1)} = g(z^{(1)}) = \tanh(z^{(1)}) \\ a^{(2)} = \theta(z^{(2)}) \end{cases} \rightarrow \begin{cases} a^{(1)} = g^{(1)}(z^{(1)}) = \tanh(z^{(1)}) \\ a^{(2)} = g^{(2)}(z^{(2)}) = \theta(z^{(2)}) \end{cases}$$

$\downarrow$  sigmoid

Activation functions differ with the different layers.

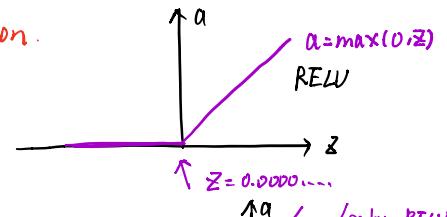
One downside of sigmoid & tangent function is that:

if  $z$  is very large / very small, the gradient descent will be very slow

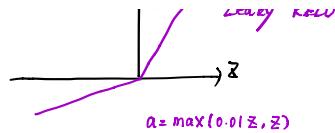
Here is another activation function: RELU function.

Some rules about activation functions:

- 1) Output is binary classification  $\in \{0, 1\}$ , output layer use sigmoid function



2) other <sup>U</sup> layers use RELU function, more efficient.



### 6. Why do we need non-linear activation functions?

e.g. for  $a = g(z)$ , we just use the linear activation,  $g(z) = z$ .

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

If we use linear activation function, the output just the linear combination of inputs.

∴ No matter how many hidden layers the NN have, it will always compute the linear combination of inputs. Then the hidden layers may be useless.

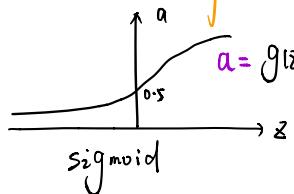
*Hint: Here is an exception when we need linear functions.*

*If you are doing machine learning on the regression problem.*

*Then  $y$  is a real number. e.g. if  $y$  is a house price.*

Also, since house price should always  $\geq 0$ , we can also use RELU function.

### 7. Derivatives of activation functions



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

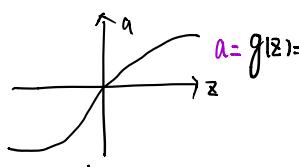
$$\frac{d}{dz}g(z) = \text{slope of } g(x) \text{ at } z$$

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

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

$$= a(1-a)$$

e.g.  $z=10, g(z) \approx 1$   
 $\text{slope} \approx 1 \times (1-1) \approx 0$   
 $z=-10, g(z) \approx 0$   
 $\text{slope} \approx 0 \times (1-0) \approx 0$   
 $z=0, g(z) = \frac{1}{2}$   
 $\text{slope} \approx \frac{1}{2} \times (1-\frac{1}{2}) = \frac{1}{4}$

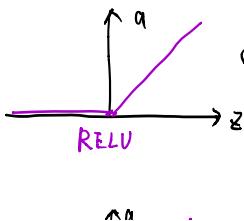


$$a = g(z) = \tanh(z)$$

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

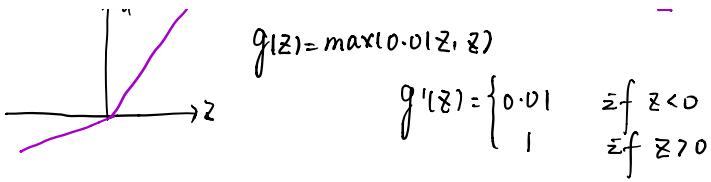
$$= 1 - (\tanh(z))^2$$

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



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

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \\ \text{undefined} & \text{if } z=0 \\ z=0.00\ldots & \end{cases}$$



## 8. Gradient descent for Neural Networks

One hidden layer.

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

cost function:  $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i)$

Gradient descent:

Repeat { Compute predicts ( $\hat{y}^{(i)}$ ,  $i=1, 2, \dots, m$ )  
 $d w^{[2]} = \frac{\partial J}{\partial w^{[2]}}$ ,  $d b^{[2]} = \frac{\partial J}{\partial b^{[2]}}$ , ...  
 $w^{[2]} = w^{[2]} - \alpha d w^{[2]}$   
 $b^{[2]} = b^{[2]} - \alpha d b^{[2]}$   
 $\vdots$   
Until parameters reach convergence}

Formulas for computing derivatives.

Forward propagation:

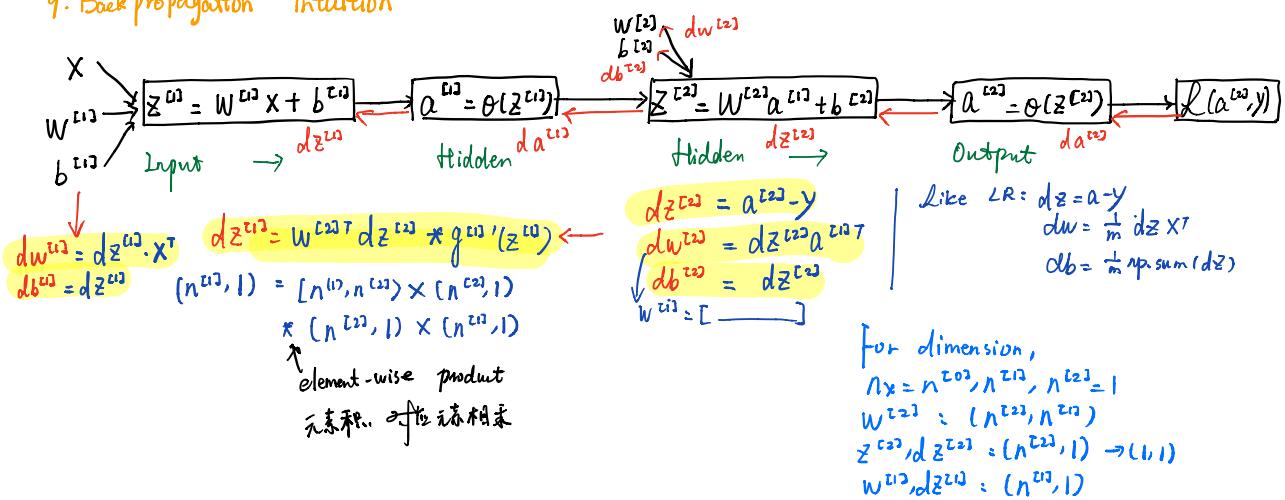
$$\begin{aligned} z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= g^{[1]}(z^{[1]}) \\ z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(z^{[2]}) = \theta(z^{[2]}) \end{aligned}$$

Back propagation

$$\begin{aligned} dz^{[2]} &= A^{[2]} - Y \\ dw^{[2]} &= \frac{1}{m} dz^{[2]} A^{[1]T} \\ db^{[2]} &= \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims=True}) \\ dz^{[1]} &= W^{[2]T} dz^{[2]} * g^{[2]'}(z^{[1]}) \\ dw^{[1]} &= \frac{1}{m} dz^{[1]} X^T \\ db^{[1]} &= \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims=True}) \rightarrow (n^{[1]}, 1) \end{aligned}$$

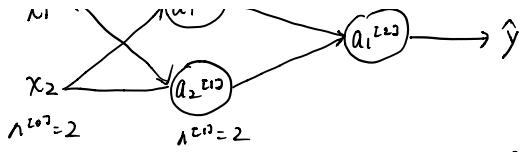
make sure array( $n^{[2]}, 1$ )

## 9. Back propagation intuition



## 10. Random Initialization

→ What happens if you initialize weights to 0?



If we initialize  $w^{L+1} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$ ,  $b^{L+1} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

Then  $a_1^{L+1} = a_2^{L+1}$ ,  $d\mathbf{z}_1^{L+1} = d\mathbf{z}_2^{L+1} \rightarrow$  Then completely symmetric

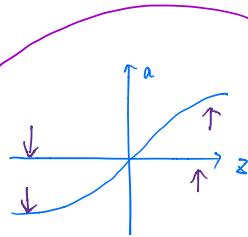
It turns out:

after every single iteration of training your two hidden units are still computing exactly the same function.

Then  $d\mathbf{w} = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$ ,  $\mathbf{w}^{L+1} = \mathbf{w}^{L+1} - \alpha d\mathbf{w}$ ,  $\mathbf{w}$  的每-行都是相同的, 所以每个  $x_2$  对应的 weights 是相同的. 不要有多个 layers.

Random Initialization

e.g.  $\left\{ \begin{array}{l} w^{L+1} = np.random.randn(2, 2) * \frac{0.01}{100} \\ b^{L+1} = np.zeros((2, 1)) \\ w^{L+2} = np.random.randn((1, 2)) * 0.01 \\ b^{L+2} = 0 \end{array} \right.$



If  $w$  is large,  $z$  is large, a 的值会 large / small  
gd 会慢 slow.