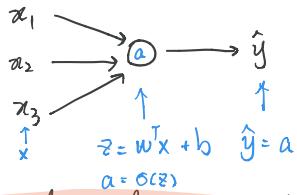


We have seen how to compute one layer of NN in lesson 1.2. Here we will show how to compute a shallow NN with more layers, but before that let's review the structure again.

### One layer of neural network structure (logistic Regression)

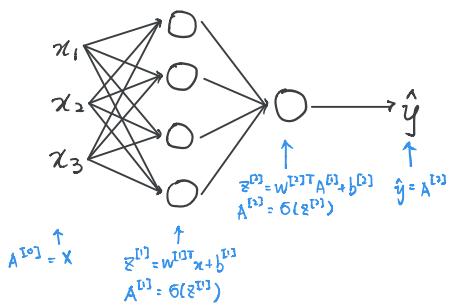


$$w, b: z = w^T x + b \rightarrow a = \sigma(z) \rightarrow J(a, y)$$

$$\text{update } w := w - \alpha \partial_w \quad b := b - \alpha \partial_b$$

matrix size:  
 $X: 3 \times m$   
 $w: 3 \times 1$   
 $b: \text{scalar}$   
 $z, a, y: 1 \times m$

### Two layers of neural network structure (Shallow NN)



$$w^{[1]}, w^{[2]}, b^{[1]}, b^{[2]}: z^{[1]} = w^{[1]T} x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = w^{[2]T} a^{[1]} + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow J(a^{[2]}, y)$$

$$\text{update } w^{[1]}: w^{[1]} - \alpha \partial_w^{[1]} \quad b^{[1]}: b^{[1]} - \alpha \partial_b^{[1]}$$

$$w^{[2]}: w^{[2]} - \alpha \partial_w^{[2]} \quad b^{[2]}: b^{[2]} - \alpha \partial_b^{[2]}$$

\* can be  $4 \times 3$  or  $1 \times 4$  depending on how  $z^{[1]}$  and  $z^{[2]}$  are defined.  
 Here  $z^{[1]} = w^{[1]T} x + b^{[1]}$  and  $z^{[2]} = w^{[2]T} a^{[1]} + b^{[2]}$ .

matrix size:  
 $X: 3 \times m$   
 $w^{[1]}: 3 \times 4$   
 $b^{[1]}: 4 \times 1$   
 $w^{[2]}: 4 \times 1$   
 $b^{[2]}: \text{scalar}$   
 $\underbrace{z^{[1]}, a^{[1]}: 4 \times m}_{\text{first layer}} \quad \underbrace{z^{[2]}, a^{[2]}, \hat{y}: 1 \times m}_{\text{output layer}}$

So to vectorize these, we have:

$$w^{[1]} = \begin{bmatrix} | & | & | & | \\ w_1^{[1]} & w_2^{[1]} & w_3^{[1]} & w_4^{[1]} \\ | & | & | & | \end{bmatrix}_{3 \times 4}$$

$$b^{[1]} = \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}_{4 \times 1}$$

$$\text{So } z^{[1]} = w^{[1]T} x + b^{[1]} = \begin{bmatrix} | & | & | \\ w_1^{[1]T} & \dots & w_4^{[1]T} \\ | & | & | \end{bmatrix}_{4 \times 3} \begin{bmatrix} | & | & | \\ x^{(1)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}_{3 \times m} + \begin{bmatrix} b_1^{[1]} \\ \vdots \\ b_4^{[1]} \end{bmatrix}_{4 \times 1} = \begin{bmatrix} w_1^{[1]} x^{(1)} + b_1^{[1]} & \dots & w_1^{[1]} x^{(m)} + b_1^{[1]} \\ \vdots & \ddots & \vdots \\ w_k^{[1]} x^{(1)} + b_k^{[1]} & \dots & w_k^{[1]} x^{(m)} + b_k^{[1]} \end{bmatrix}_{4 \times m}$$

$$= \begin{bmatrix} | & | & | \\ z^{(1)} & \dots & z^{(m)} \\ | & | & | \end{bmatrix}_{4 \times m} = \begin{bmatrix} | & | & | \\ \hat{y}^{(1)} & \dots & \hat{y}^{(m)} \\ | & | & | \end{bmatrix}_{4 \times m}$$

↑ different examples      ↑ different nodes

$$A^{[1]} = \begin{bmatrix} | & | & | & | \\ a_1^{[1](1)} & \dots & a_1^{[1](m)} & | \\ | & | & | & | \end{bmatrix}_{4 \times m} = \begin{bmatrix} -a_1^{[1]} - \\ -a_4^{[1]} - \end{bmatrix}_{4 \times m} = g^{[1]}(z^{[1]})_{4 \times m}$$

By convention, we use  $g^{[i]}$  to denote activation function for the  $i$ th layer.

For the second layer:

$$w^{[2]} = \begin{bmatrix} w_1^{[2]} \\ \vdots \\ w_4^{[2]} \end{bmatrix}_{4 \times 1}$$

$$z^{[2]} = w^{[2]T} A^{[1]} + b^{[2]} = \begin{bmatrix} w_1^{[2]} & \dots & w_4^{[2]} \end{bmatrix}_{1 \times 4} \begin{bmatrix} | & | & | & | \\ a_1^{[1](1)} & \dots & a_1^{[1](m)} & | \\ | & | & | & | \end{bmatrix}_{4 \times m} + b^{[2]} = \begin{bmatrix} -w_1^{[2]} a^{[1]} + b^{[2]} - \\ \vdots \\ -w_4^{[2]} a^{[1]} + b^{[2]} - \end{bmatrix}_{1 \times m}$$

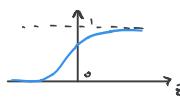
$$A^{[2]} = g^{[2]}(z^{[2]})_{1 \times m} = \hat{Y} = \begin{bmatrix} \hat{y}^{(1)} & \dots & \hat{y}^{(m)} \end{bmatrix}_{1 \times m}.$$

## Choosing Activation functions

So far we've only used sigmoid function  $\sigma$ , but almost always there are other better choices for activation functions. The following will compare different activation functions.

Sigmoid

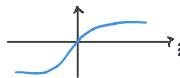
$$\frac{1}{1 + e^{-z}}$$



- almost never use, except for output layer for binary classification.
- slow learning rate as the slope is close to 0.

Tanh

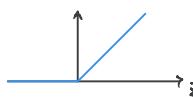
$$\frac{e^z - e^{-z}}{e^z + e^{-z}}$$



- almost always superior to sigmoid
- has the effect of centering data to 0 as the mean is 0
- still slow learning rate for the same reason as sigmoid's

Relu

$$\max(0, z)$$

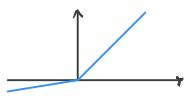


- default activation function for most NNs
- fast learning rate as the slope = 1 (for  $z > 0$ )
- slope = 0 for  $z \leq 0$  but in practise  $z$  rarely becomes negative.

leaky Relu

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

↑  
can be  
other constant



- not as oftenly used as ReLU
- slope = 0.01 for  $z \leq 0$

The reason we use non-linear functions as activation functions is because if we use linear function (or identity activation function), a NN is the same as linear regression.

Prove: suppose we have a 2-layers NN, let  $A^{[1]} = Z^{[1]}$ ,  $A^{[2]} = Z^{[2]}$ .

$$\begin{aligned} \hat{y} &= A^{[2]} : Z^{[2]} = W^{[2]T} A^{[1]} + b^{[2]} = W^{[2]T} (W^{[1]T} X + b^{[1]}) + b^{[2]} = \underbrace{W^{[2]T} W^{[1]T} X}_{= w} + \underbrace{W^{[2]T} b^{[1]} + b^{[2]}}_{= b} \\ &= wX + b' \end{aligned}$$

## Differentiation of Activation Functions

To update  $w$  and  $b$  by gradient descent, we need to first find out  $dW^{[1]}$ ,  $dW^{[2]}$ ,  $dB^{[1]}$  and  $dB^{[2]}$  given our choice of activation functions and thus we need to know  $g'(z)$  first.

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

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

ReLU

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

not exactly right mathematically,  
but it works just fine in practise

Tanh

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

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

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

## Shallow NN

To generalise our model, suppose we have  $n_n = n^{[0]}$  input features,  $n^{[1]}$  hidden layer features,  $n^{[2]} = 1$  output layer feature and  $m$  training examples. Then:

$$X : n_n \times m \quad W^{[1]} : n_n \times n^{[1]} \quad b^{[1]} : n^{[1]} \times 1 \quad W^{[2]} : n^{[1]} \times n^{[2]} \quad b^{[2]} : n^{[2]} \times 1 \quad Z^{[1]}, A^{[1]} : n^{[1]} \times m \quad Z^{[2]}, A^{[2]} : n^{[2]} \times m$$

First thing we need to do is to initialise  $w$  and  $b$ .

We need to initialise  $w^{[1]}$  and  $w^{[2]}$  randomly and initialise  $b$  to zeros, because initialising  $w$  to zero will let all hidden units compute the same activation function, but initialising  $b$  to zero will not have this problem, often known as symmetry breaking.

So :

$$W^{[1]} = np.random.randn(n_0, n^{[1]}) \cdot 0.01 \quad \text{smaller value of } W \text{ allows for smaller values of } z \text{ and therefore larger slope of activation function, i.e. faster learning rate.}$$

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

$$W^{[2]} = np.random.randn(n^{[1]}, 1) \cdot 0.01$$

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

For forward propagation:

$$\bullet X : n_x \times m \quad W^{[1]} : n_x \times n^{[1]} \quad b^{[1]} : n^{[1]} \times 1 \quad w^{[1]} : n^{[1]} \times 1 \quad b^{[1]} : 1 \times 1$$

$$z^{[1]} = \text{np.dot}(w^{[1]}, T, X) + b^{[1]} = w^{[1]T} X + b^{[1]} \quad A^{[1]} = g^{[1]}(z^{[1]}) \Rightarrow z^{[1]}, A^{[1]} : n^{[1]} \times m$$

$$z^{[2]} = \text{np.dot}(w^{[2]}, T, A^{[1]}) + b^{[2]} = w^{[2]T} A^{[1]} + b^{[2]} \quad A^{[2]} = g^{[2]}(z^{[2]}) \Rightarrow z^{[2]}, A^{[2]} : 1 \times m$$

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

$$\bullet \text{the cost function } J = \frac{1}{m} \sum_{i=1}^m L(a^{[2]}_i, y_i) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log a^{[2](i)} + (1-y^{(i)}) \log(1-a^{[2](i)})] g^{[2](i)}(1-g^{[2](i)}(z^{[2](i)}))$$

for binary classification

$$= -\frac{1}{m} \cdot (\text{np.sum}(y \cdot \log a^{[2]}) + \text{np.sum}(1-y \cdot \log(1-a^{[2]}))) \text{ or } \begin{cases} \text{logprob} = \text{np.multiply}(\text{np.log}(A^{[2]}), Y) \\ -\frac{1}{m} \text{np.sum(logprob)} \end{cases}$$

For back propagation:

$$\bullet d\hat{y} = A^{[2]} - Y : 1 \times m$$

$$\bullet dW^{[2]} = \frac{1}{m} A^{[1]} (d\hat{y})^T = \frac{1}{m} \text{np.dot}(A^{[1]T}, d\hat{y}) : n^{[1]} \times 1$$

$$\bullet db^{[2]} = \frac{1}{m} \text{np.sum}(d\hat{y}, \text{axis}=1, \text{keepdims=True}) : 1 \times 1$$

$$\bullet dz^{[2]} = \underbrace{\text{np.dot}(W^{[2]}, d\hat{y})}_{n^{[2]} \times m} \cdot \underbrace{g^{[2]'}(z^{[2]})}_{n^{[2]} \times m} : n^{[2]} \times m$$

$$\bullet dW^{[1]} = \frac{1}{m} X (dz^{[2]})^T = \frac{1}{m} \text{np.dot}(X, dz^{[2]}) : n_x \times n^{[1]}$$

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

The following will show the prove for back propagation formulae we used above:

$$\bullet d\hat{y}^{(i)} = \frac{\partial L^{(i)}}{\partial z^{[2](i)}} = \frac{\partial L^{(i)}}{\partial z^{[2](i)}} \frac{\partial a^{[2](i)}}{\partial z^{[2](i)}} = \frac{\partial}{\partial a^{[2](i)}} - [y^{(i)} \log a^{[2](i)} + (1-y^{(i)}) \log(1-a^{[2](i)})] g^{[2](i)}(z^{[2](i)}) \text{ for binary classification}$$

$$= -\left(\frac{y^{(i)}}{a^{[2](i)}} - \frac{1-y^{(i)}}{1-a^{[2](i)}}\right) (a^{[2](i)}(1-a^{[2](i)})) = -y^{(i)}(1-a^{[2](i)}) + (1-y^{(i)})a^{[2](i)}$$

from sigmoid function differentiation.

$$= a^{[2](i)} - y^{(i)} \quad (1 \times 1)$$

$$dz^{[2]} = \frac{\partial}{\partial z^{[2]}} \sum_{i=1}^m L^{(i)} = \left[ \frac{\partial}{\partial z^{[2](1)}} L^{(1)} \dots \frac{\partial}{\partial z^{[2](m)}} L^{(m)} \right]_{1 \times m} = [a^{[2](1)} - y^{(1)} \dots a^{[2](m)} - y^{(m)}]_{1 \times m}$$

$$= \underline{(A^{[2]} - Y)_{1 \times m}}$$

$$\bullet dW^{[2](i)} = \frac{\partial L^{(i)}}{\partial W^{[2]}} = \frac{\partial L^{(i)}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial W^{[2]}} = dz^{[2]} \frac{\partial}{\partial W^{[2]}} a^{[2](i)} + b^{[2]} = dz^{[2]} \left[ \begin{array}{c} \frac{\partial}{\partial W^{[2]}} a^{[2](1)} \\ \vdots \\ \frac{\partial}{\partial W^{[2]}} a^{[2](m)} \end{array} \right]_{n^{[2]} \times 1} = dz^{[2]} a^{[2](i)}_{n^{[2]} \times 1}$$

$$dW^{[2]} = \frac{\partial J}{\partial W^{[2]}} = \frac{\partial}{\partial W^{[2]}} \frac{1}{m} \sum_{i=1}^m L^{(i)} = \frac{1}{m} \sum_{i=1}^m dz^{[2]} a^{[2](i)} = \frac{1}{m} A^{[1]} (dz^{[2]})^T$$

$$\bullet db^{[2](i)} = \frac{\partial L^{(i)}}{\partial b^{[2]}} = \frac{\partial L^{(i)}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial b^{[2]}} = dz^{[2]} \frac{\partial}{\partial b^{[2]}} a^{[2](i)} + b^{[2]} = dz^{[2]} a^{[2](i)}$$

$$db^{[2]} = \frac{\partial J}{\partial b^{[2]}} = \frac{\partial}{\partial b^{[2]}} \frac{1}{m} \sum_{i=1}^m L^{(i)} = \frac{1}{m} \sum_{i=1}^m dz^{[2]} a^{[2](i)} = \frac{1}{m} dz^{[2]} \text{ summed up horizontally.}$$

$$\bullet dz^{[1](i)} = \frac{\partial L^{(i)}}{\partial z^{[1]}} = \frac{\partial L^{(i)}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial z^{[1]}} = dz^{[2]} \frac{\partial}{\partial z^{[1]}} a^{[2](i)} + b^{[2]} = dz^{[2]} g^{[2]'}(z^{[2]}) = dz^{[2]} w^{[2]} g^{[2]'}(z^{[2](i)})$$

$$dz^{[1]} = \frac{\partial}{\partial z^{[1]}} \sum_{i=1}^m L^{(i)} = \left[ \frac{\partial}{\partial z^{[1]}} L^{(1)} \dots \frac{\partial}{\partial z^{[1]}} L^{(m)} \right]_{n^{[1]} \times m} = \left[ dz^{[2]} \frac{\partial}{\partial z^{[1]}} a^{[2](1)} \dots dz^{[2]} \frac{\partial}{\partial z^{[1]}} a^{[2](m)} \right]_{n^{[1]} \times m}$$

element-wise multiplication

$$= w^{[2]} dz^{[2]} * g^{[2]'}(z^{[2]}) \quad n^{[1]} \times m$$

$$\begin{aligned}
 \bullet \quad dW^{[l]} &= \frac{\partial L^{(i)}}{\partial w^{[l]}} = \frac{\partial L^{(i)}}{\partial z^{[l+1]}} \frac{\partial z^{[l+1]}}{\partial w^{[l]}} = dz^{[l+1]} \frac{\partial}{\partial w^{[l]}} w^{[l]T} x^{(i)} + b^{[l]} = dz^{[l+1]} \frac{\partial}{\partial \begin{bmatrix} w_{11}^{[l]} & \dots & w_{m1}^{[l]} \\ \vdots & \ddots & \vdots \\ w_{1n}^{[l]} & \dots & w_{mn}^{[l]} \end{bmatrix}} w^{[l]T} x^{(i)} = dz^{[l+1]} \begin{bmatrix} \frac{\partial}{\partial w_{11}^{[l]}} & \dots & \frac{\partial}{\partial w_{m1}^{[l]}} \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial w_{1n}^{[l]}} & \dots & \frac{\partial}{\partial w_{mn}^{[l]}} \end{bmatrix} \\
 \text{where } W^{[l]T} x^{(i)} &= \begin{bmatrix} w_{11}^{[l]} x_1^{(i)} + \dots + w_{m1}^{[l]} x_n^{(i)} \\ \vdots \\ w_{1n}^{[l]} x_1^{(i)} + \dots + w_{mn}^{[l]} x_n^{(i)} \end{bmatrix} \Rightarrow dW^{[l]} = dz^{[l+1]} \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} = dz^{[l+1]} x^{(i)}
 \end{aligned}$$

$$dW^{[l]} = \frac{\partial}{\partial w^{[l]}} \sum_{i=1}^m L^{(i)} = \frac{1}{m} \sum_{i=1}^m dz^{[l+1]} x^{(i)} = \underline{\frac{1}{m} \times (dz^{[l]})^T}$$

$$\bullet \quad db^{[l]} = \frac{\partial L^{(i)}}{\partial b^{[l]}} = \frac{\partial L^{(i)}}{\partial z^{[l+1]}} \frac{\partial z^{[l+1]}}{\partial b^{[l]}} = dz^{[l+1]}$$

$$db^{[l]} = \frac{\partial}{\partial b^{[l]}} \sum_{i=1}^m L^{(i)} = \frac{1}{m} \sum_{i=1}^m dz^{[l+1]} = \underline{\frac{1}{m} dz^{[l]} \text{ summed up horizontally.}}$$