

Deep Neural Networks

1. Deep L-layer neural networks

Deep neural network notation

$$L = 4 \text{ (# layers)}$$

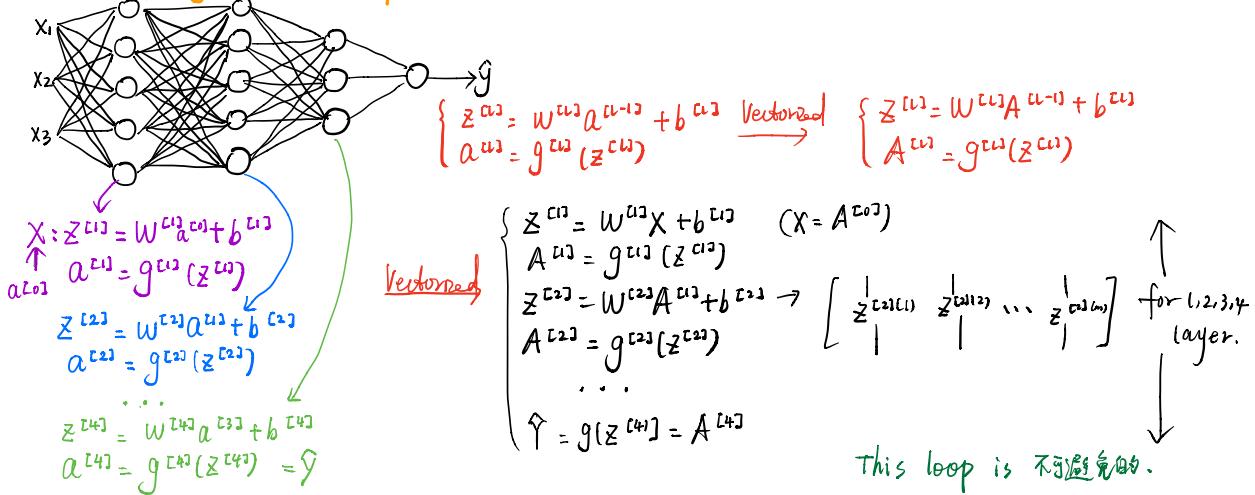
$$n^{[L]} = \# \text{ units of layer } L, \text{ eg: } n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3, n^{[4]} = n^{[L]} = 1, n^{[0]} = n_x = 3$$

$$a^{[l]} = \text{activations of layer } l$$

$$a^{[l]} = g^{[l]}(z^{[l]}) \quad w^{[l]} = \text{weights for } z^{[l]}$$

$$\downarrow X = a^{[0]}, \quad a^{[L]} = \hat{y}$$

2. Forward Propagation in a Deep Network



3. Getting your matrix dimensions right.

When $L=5, n^{[1]}=3, n^{[2]}=5, n^{[3]}=4, n^{[4]}=2, n^{[5]}=1$

$$\begin{array}{lll} z^{[1]} = W^{[1]}X + b^{[1]} & z^{[2]} = W^{[2]} \cdot a^{[1]} + b^{[2]} & z^{[3]} = W^{[3]} \cdot a^{[2]} + b^{[3]} \\ (3,1) \quad (3,2) \quad (2,1) & (5,1) \quad (5,3) \quad (3,1) & (4,1) \quad (4,2) \quad (2,1) \\ (n^{[0]}, 1) \leftarrow (n^{[0]}, n^{[1]}) \times (n^{[1]}, 1) & (n^{[1]}, n^{[2]}) \times (n^{[2]}, 1) & (n^{[2]}, n^{[3]}) \times (n^{[3]}, 1) \end{array}$$

$$\text{For dimensions: } \begin{cases} W^{[l]} : (n^{[l]}, n^{[l-1]}) \\ b^{[l]} : (n^{[l]}, 1) \end{cases} \quad z^{[l]}, a^{[l]} : (n^{[l]}, 1)$$

$$\begin{cases} dw^{[l]} : (n^{[l]}, n^{[l-1]}) \\ db^{[l]} : (n^{[l]}, 1) \end{cases}$$

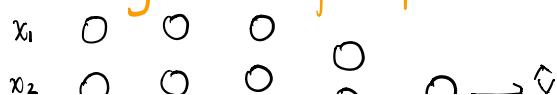
Vectorized implementation

$$\begin{array}{ccc} z^{[1]} = W^{[1]}X + b^{[1]} & \rightarrow [n^{[1]}, 1] & [\text{through broadcasting, it will transform} \\ & (n^{[0]}, m) & \text{to } (n^{[1]}, m)] \\ (n^{[1]}, m) & (n^{[1]}, n^{[1]}) & \end{array}$$

$$\begin{cases} z^{[l]}, A^{[l]} : (n^{[l]}, m) \\ \text{when } l=0, A^{[0]} = X = (n^{[0]}, m) \\ (dz^{[l]}, dA^{[l]} : (n^{[l]}, m)) \end{cases}$$

4. Why deep representations?

5. Building blocks of deep neural networks



$$\begin{matrix} x_3 & \textcircled{1} & \textcircled{0} & \textcircled{0} & \textcircled{0} \\ x_4 & \textcircled{0} & \textcircled{0} & \textcircled{0} & \textcircled{0} \end{matrix}$$

layer $l: W^{[l]}, b^{[l]}$

Forward: Input $a^{[l-1]}$, output $a^{[l]}$

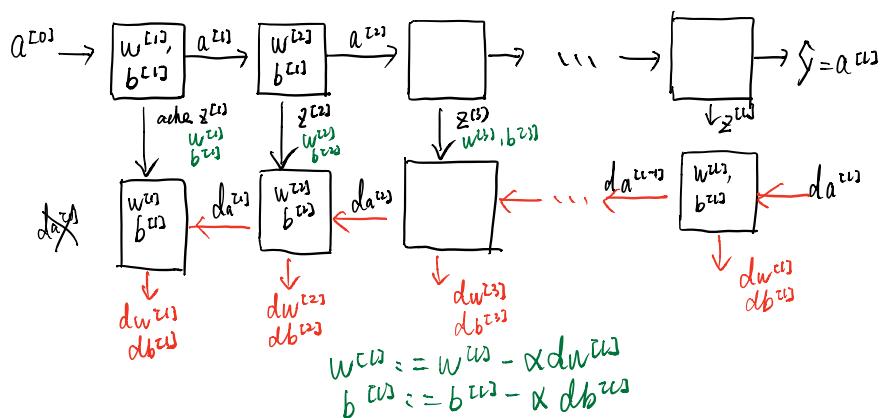
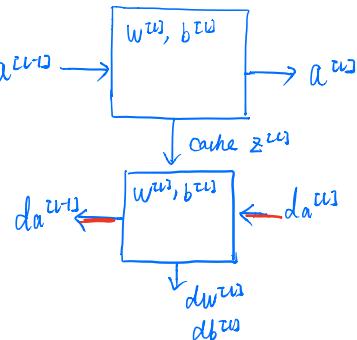
$$\begin{cases} z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]} \\ a^{[l]} = g^{[l]}(z^{[l]}) \end{cases}$$

cache $z^{[l]}$ (把 $z^{[l]}$ 存储起来)

Backward: Input $da^{[l]}$, output
cache($z^{[l]}$)

$$\begin{matrix} da^{[l-1]} \\ dw^{[l]} \\ db^{[l]} \end{matrix}$$

layer l



6. Forward and Backward Propagation

Forward propagation for layer l :

Input $a^{[l-1]}$

Output $a^{[l]}$, cache($z^{[l]}$, $W^{[l]}$, $b^{[l]}$)

$$z^{[l]} = W^{[l]} \cdot A^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

$$X = A^{[0]} \rightarrow \square \rightarrow \square \rightarrow \dots$$

Backward propagation for layer l :

Input $da^{[l]}$

Output $da^{[l-1]}, dw^{[l]}, db^{[l]}$

$$dZ^{[l]} = dA^{[l]} \cdot g^{[l]'}(z^{[l]}) \rightarrow dZ^{[l]} = W^{[l+1]T} \cdot dA^{[l+1]} \cdot g^{[l+1]'}(z^{[l]})$$

$$dw^{[l]} = dZ^{[l]} \cdot A^{[l-1]T}$$

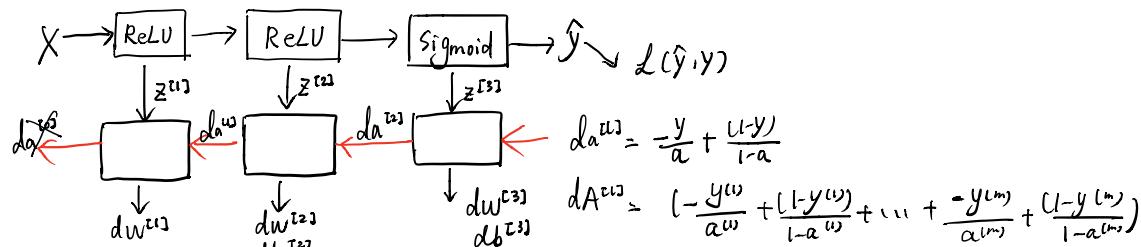
$$db^{[l]} = dZ^{[l]}$$

$$dA^{[l-1]} = W^{[l]T} \cdot dZ^{[l]}$$

↙ Vectorized

$$\star \left\{ \begin{array}{l} dZ^{[l]} = dA^{[l]} \cdot g^{[l]'}(z^{[l]}) \\ dW^{[l]} = \frac{1}{m} dZ^{[l]} \cdot A^{[l-1]T} \\ db^{[l]} = \frac{1}{m} \text{np.sum}(dZ^{[l]}, \text{axis}=1, \text{keepdims=True}) \\ dA^{[l-1]} = W^{[l]T} \cdot dZ^{[l]} \end{array} \right.$$

Summary: (3-layer network)



$db^{(l)}$ $dw^{(l)}$

7. Parameters vs. Hyperparameters

parameters: $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, W^{(3)}, b^{(3)}, \dots$

hyperparameters: { learning rate α
iterations
hidden layers L
hidden units $n^{(1)}, n^{(2)}, \dots$
choice of activation function

这些系数是控制 parameters 和最终结果的

other hyperparameters: Momentum, minibatch size, regularizations ...