

$L = 4$  (# layers)

$n^{[l]}$  = # units in layer  $l$

$a^{[l]}$  = activations in layer  $l$

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

$w^{[l]}$  = weights for  $z^{[l]}$  :  $(n^{[l]}, n^{[l-1]})$

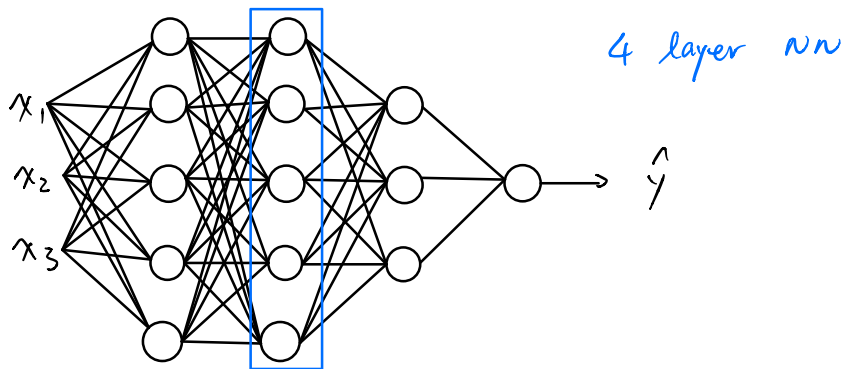
$m$  for all samples  
 $(n^{[l]}, \text{ })$

---


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

$\Rightarrow$  Vectorized:

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



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

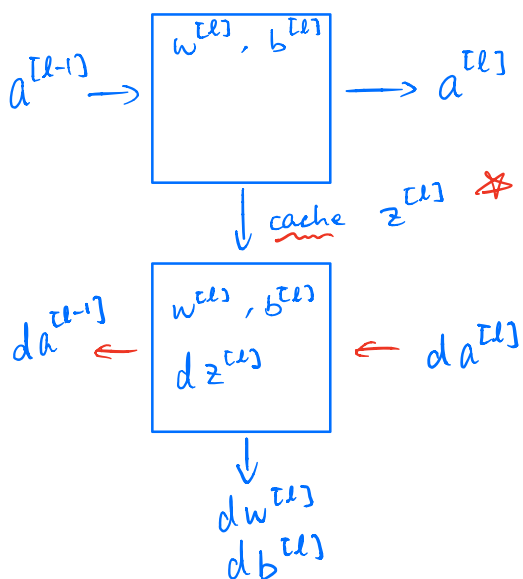
Forward: Input  $a^{[l-1]}$ , output  $a^{[l]}$

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]} \Leftarrow \text{cache } z^{[l]}$$

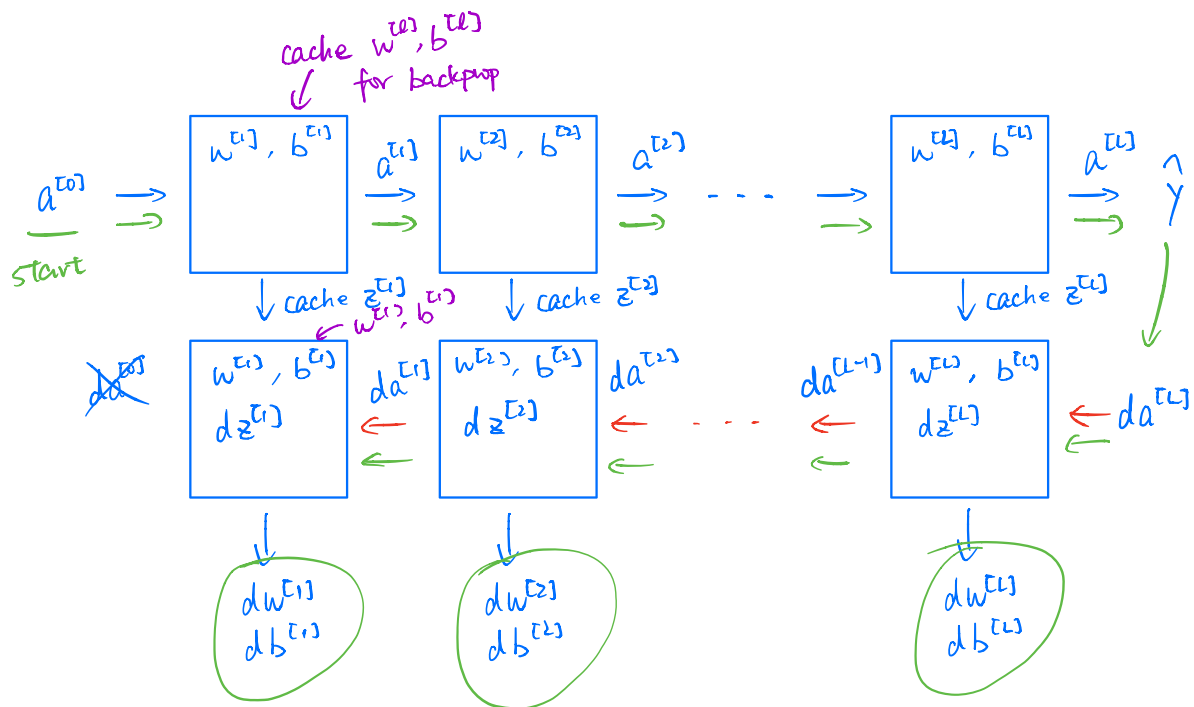
$$a^{[l]} = \sigma^{[l]}(z^{[l]})$$

Backward: Input  $da^{[l]}$ , output  $\frac{da^{[l-1]}}{dz^{[l]}}$   
 $\text{cache } (z^{[l]})$   $\frac{dw^{[l]}}{db^{[l]}}$

layer  $l$



## Forward and Backward Functions



$$w^{[L]} = w^{[L]} - \alpha dw^{[L]}$$

$$b^{[L]} = b^{[L]} - \alpha db^{[L]}$$

### Forward

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

Output  $a^{[L]}$ , cache  $z^{[L]}$

$$\Rightarrow \begin{cases} z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \\ a^{[L]} = \sigma^{[L]}(z^{[L]}) \end{cases}$$

### Backward

Input  $da^{[L]}$

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

$$\Rightarrow \begin{cases} dz^{[L]} = da^{[L]} * \sigma'^{[L]}(z^{[L]}) \\ dw^{[L]} = dz^{[L]} \cdot a^{[L-1]T} \\ db^{[L]} = dz^{[L]} \\ da^{[L-1]} = \underline{w^{[L]T} dz^{[L]}} \end{cases}$$

~~Practice !!~~

Vectorized

Forward

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

Output  $A^{[l]}$ , cache  $z^{[l]}$

also  $w^{[l]}, b^{[l]}$

$$\Rightarrow \begin{cases} z^{[l]} = w^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} = \sigma^{[l]}(z^{[l]}) \end{cases}$$

Backward

Input  $dA^{[l]}$

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

$$\Rightarrow \begin{cases} dz^{[l]} = dA^{[l]} * \sigma'^{[l]}(z^{[l]}) \\ dw^{[l]} = \frac{1}{n} dz^{[l]} \cdot A^{[l-1]T} \\ db^{[l]} = \frac{1}{n} \text{np.sum}(dz^{[l]}, \text{axis}=1, \text{keepDim}=\text{True}) \\ dA^{[l-1]} = \underline{w^{[l]T} dz^{[l]}} \end{cases}$$

Practice

$$dz^{[l]} = dA^{[l]} * \frac{\partial A^{[l]}}{\partial z^{[l]}} \rightarrow \sigma^{[l]}(z^{[l]})$$

$(n^{[l]}, m)$     $(n^{[l]}, m)$     $(n^{[l]}, m)$

$$dw^{[l]} = \frac{1}{n} dz^{[l]} \cdot A^{[l-1]T}$$

$(n^{[l]}, n^{[l-1]})$     $(n^{[l]}, m)$     $(m, n^{[l-1]})$

$$db^{[l]} = \frac{1}{n} dz^{[l]} \xrightarrow{\text{Sum axis}=1} (n^{[l]}, 1)$$

$$\star dA^{[l-1]} = \frac{\partial z^{[l]}}{\partial A^{[l-1]}} \cdot dw^{[l]T}$$

$(n^{[l-1]}, 1)$     $(n^{[l]}, n^{[l-1]})$     $(n^{[l]}, 1)$

## Hyper parameters

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

Hyperparameters: learning rate  $\alpha$   
# iterations  
# hidden layers  $L$   
# hidden units  $n^{(1)}, n^{(2)}, \dots$   
Choice of activation function

Later: momentum  
minibatch size  
regularizations  
...

Try different values  $\rightarrow$  empirical process

