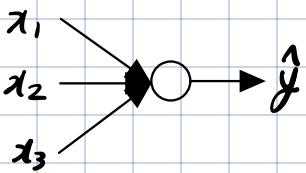


WEEK 4. Deep Neural Network

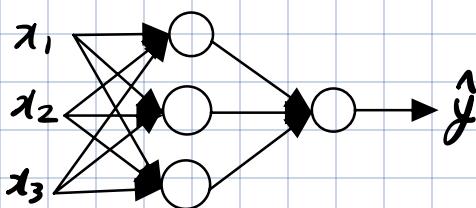
- Deep L-Neural Network

- What is a Deep Neural Network?



< Logistic Regression >

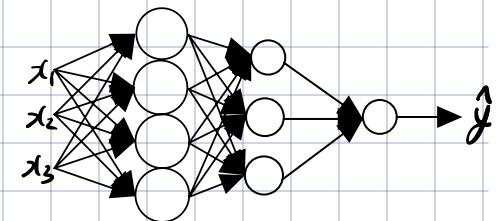
\Rightarrow Very shallow model
(1 layer NN)



< 1 hidden layer >

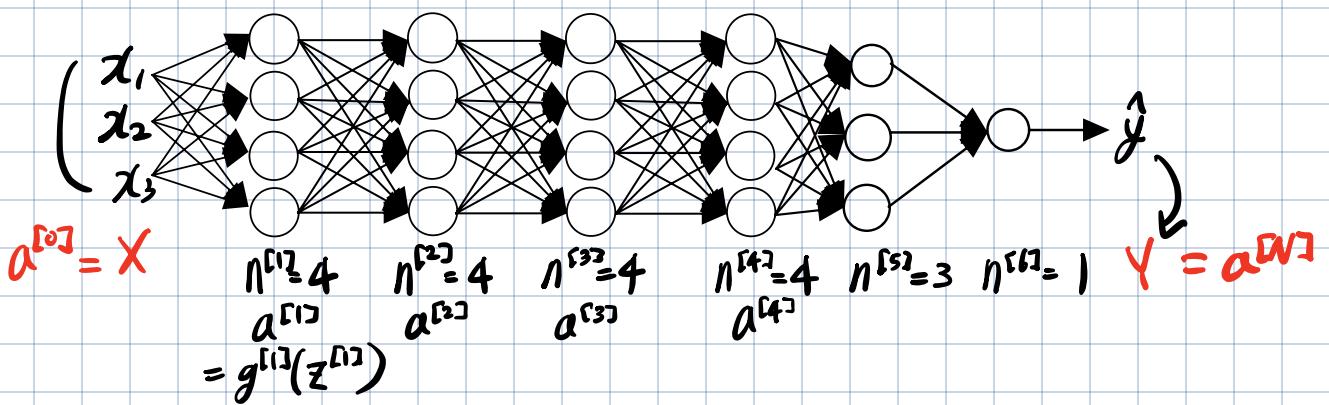
\Rightarrow 2 layer NN

(don't count input layer)



< 2 hidden layers >

\Rightarrow 3 layer NN



< 5 hidden layers >

\Rightarrow deep Model

can learn that shallow models can't learn

\Rightarrow Increase the # of hidden layer to find best model (during validation)

* Notation

↳ hyperparameter

L: # of layer

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

$a^{[l]} = g^{[l]}(z^{[l]})$: activations in layer l

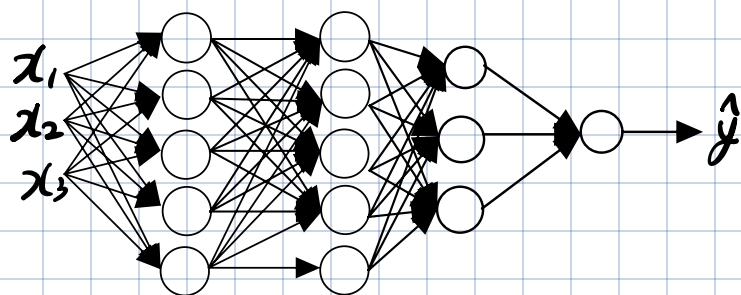
$w^{[l]}$: weights for $z^{[l]}$, $b^{[l]}$

X: input features = activations at layer 0 = $a^{[0]}$

Y: Output = activations at the last layer L = $a^{[L]}$

• Forward Propagation in a Deep Network

- Single training sample x



General Rule of forward propagation

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(Z^{[l]})$$

- ① Compute the activation of the first layer

$$Z^{[1]} = W^{[1]} x + b^{[1]} \rightarrow a^{[0]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

↳ activation function g depends on layer

- ② Compute the activation of the second layer

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

$$A^{[2]} = g^{[2]}(Z^{[2]})$$

:

$$\hat{y} = g^{[4]}(Z^{[4]}) = a^{[4]}$$

- Vectorized forward Propagation for entire training samples

- ① Compute the activation of the first layer

$$Z^{[1]} = W^{[1]} X + b^{[1]} \rightarrow a^{[0]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

- ② Compute the activation of the second layer

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]})$$

:

$$\hat{y} = g(Z^{[4]}) = A^{[4]}$$

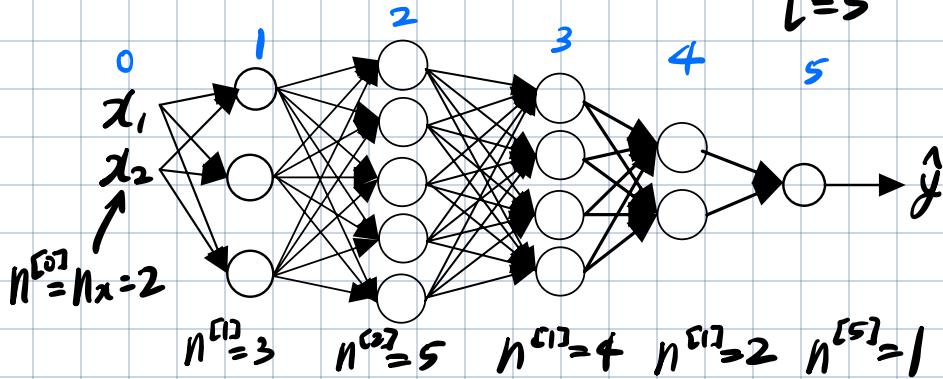
⇒ need for-loop

$$l = 1 \dots 4$$

be careful of matrix dimensions!

• Getting your matrix dimensions right

- parameters $W^{[l]}$ and $b^{[l]}$



General Rule of dimensions

$$Z = (n^{[L]}, 1) = a^{[L]}$$

$$X = (n^{[0]}, 1) \rightarrow g^{[0]}(z^{[0]})$$

$$W = (n^{[L]}, n^{[L-1]}) = dW$$

$$b = (n^{[L]}, 1) = db$$

① Compute Forward propagation at 1st layer

$$Z^{[1]} = W^{[1]} X + b^{[1]} \xrightarrow{\text{③ Something makes}} \text{④ to get output}$$

① activations of first hidden layer
 $\Rightarrow 3 \times 1$ vector
 $(n^{[1]}, 1)$

② 2 input features
 $\Rightarrow (2, 1)$ vector
 $(n^{[0]}, 1)$

$$W(n^{[0]}, 1) = (n^{[1]}, 1)$$

$$\begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \end{bmatrix}$$

$$\therefore (n^{[1]}, n^{[0]})$$

dimension of $(3, 1)$, it should be $(3, 1)$, too.

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

② Compute Forward propagation at 2nd layer

$$Z^{[2]} = W^{[2]} \cdot a^{[1]} + b^{[2]}$$

$$\begin{matrix} \uparrow & \uparrow & \uparrow & \uparrow \\ (5, 1) & (5, 3) & (3, 1) & (5, 1) \\ (n^{[2]}, 1) & (n^{[2]}, n^{[1]}) & (n^{[1]}, 1) \end{matrix}$$

- The dimension of derivatives should be same during back propagation.

$$W \leftrightarrow dW, b \leftrightarrow db$$

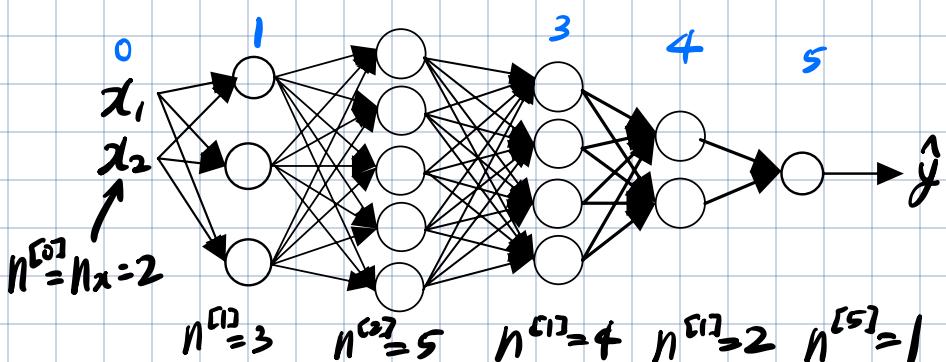
- The other important dimension checks

Z

x

$$a^{[4]} = g^{[4]}(Z^{[4]}) \Rightarrow Z \& a \text{ have same dimensions}$$

- In vectorized implementation on entire training set
 $\Rightarrow z, a, x$ will be changed



General rule of dimensions in vectorized implementation

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

$$l=0: A^{[0]} = X = (n^{[0]}, m)$$

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

① Compute Forward propagation at 1st layer

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

$$(n^{[1]}, 1) \quad (n^{[0]}, n^{[0]}) \quad (n^{[0]}, 1) \quad (n^{[0]}, 1)$$

$$z^{[1]} = W^{[1]}X + b^{[1]} \rightarrow \textcircled{3} (n^{[0]}, m)$$

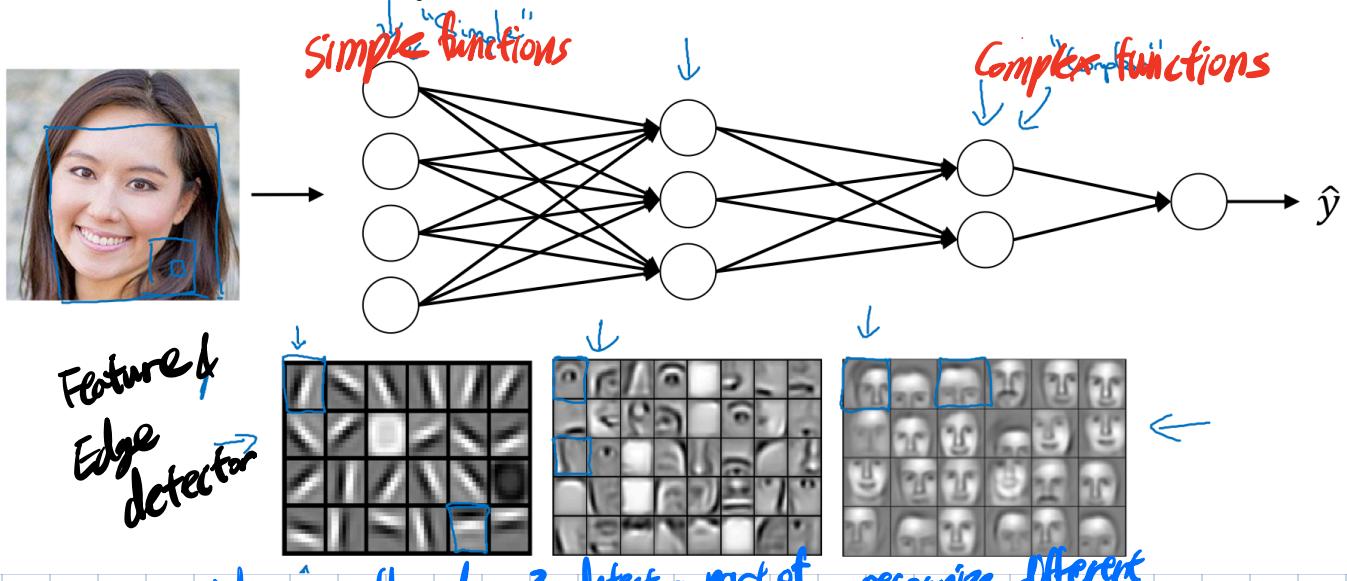
$$\textcircled{1} [z^{1}, z^{[1](2)}, \dots, z^{[1](m)}] \rightarrow \textcircled{2} \text{ stays the same } (n^{[1]}, n^{[0]}) \rightarrow \textcircled{4} \text{ still } (n^{[1]}, 1)$$

\Rightarrow broadcasting to $(n^{[1]}, m)$

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

• Why deep representations?

- Intuition about deep representation



Speech recognition system: low level audio wave \rightarrow words \rightarrow sentence/phrase
(phonemes)

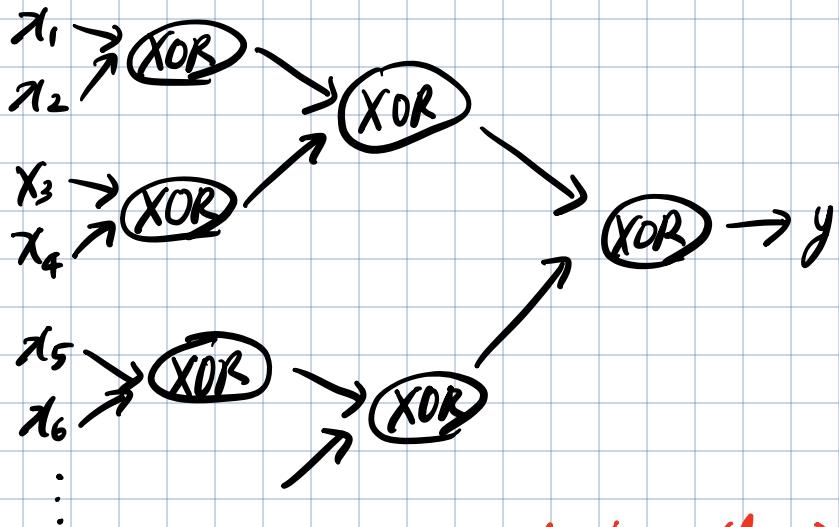
- Circuit theory & deep learning

Informally: There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

Low # of hidden / layer

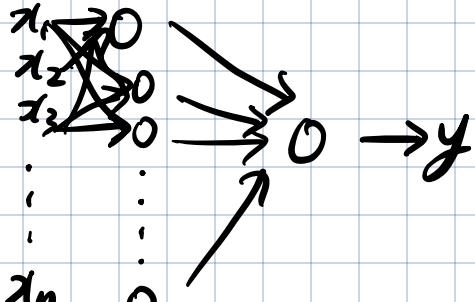
Let's say you compute 'XOR' of all the input features

$$x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } \dots \text{ XOR } x_n$$



You don't need
Many gates (node)

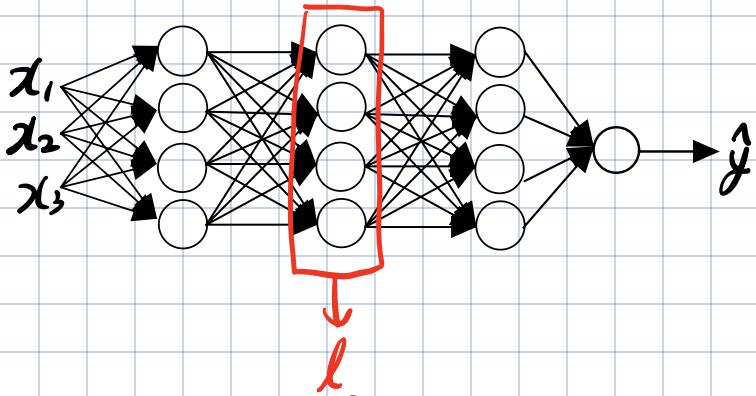
✓ If you use 1 hidden layer...



In order to compute XOR,
the hidden layer will need to
be exponentially large

$\hookrightarrow 2^{n-1}$ possible configurations $\Rightarrow O(2^n)$

• Building blocks of deep neural networks



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

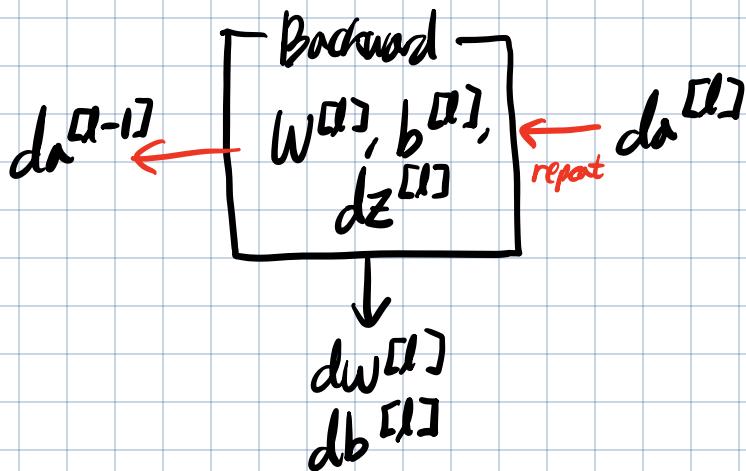
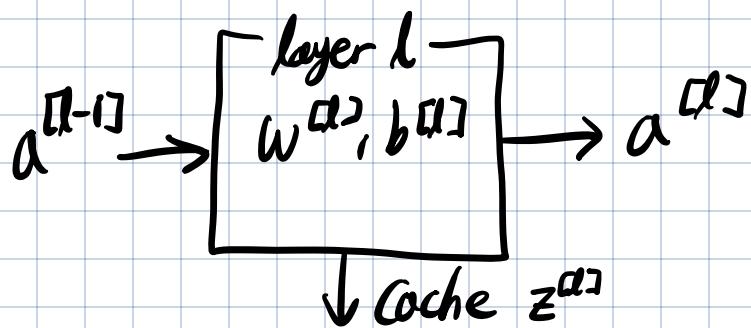
Forward: Input $a^{[l-1]}$, output $a^{[l]}$

$Z^{[l]}: W^{[l]}a^{[l-1]} + b^{[l]} \rightarrow$ cache the value $Z^{[l]}$ for later back propagation

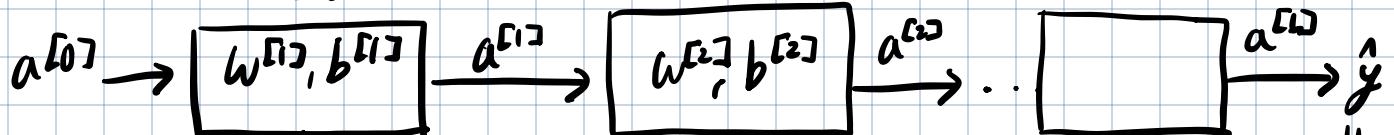
$a^{[l]}: g^{[l]}(Z^{[l]})$

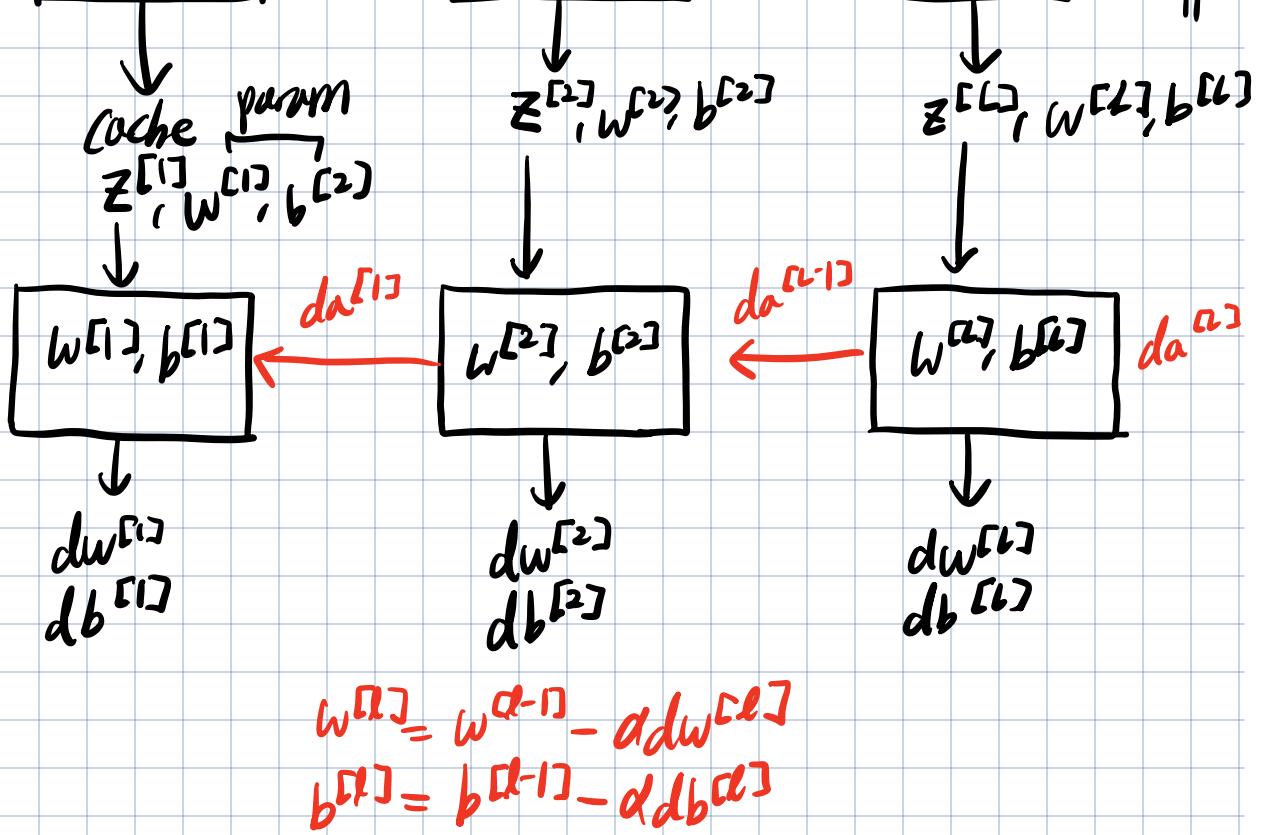
Backward: Input $da^{[l]}$, output $da^{[l-1]}$

use cache($Z^{[l]}$) $\frac{dw^{[l]}}{db^{[l]}}$



- Forward & Backward functions

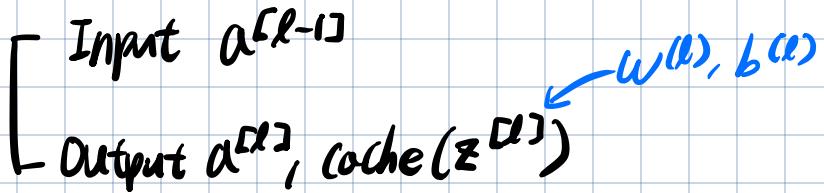




* Cache: storing value of Z for backpropagation

• Forward and backward propagation

- Forward propagation for layer l



$$Z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

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

<Vectorized>

$$Z^{[l]} = (W^{[l]} A^{[l-1]} + b^{[l]})$$

$$A^{[l]} = g^{[l]}(Z^{[l]})$$

- Backward propagation

[Input: $da^{[l]}$]

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

$$dZ^{[l]} = da^{[l]} * g'^{[l]}(Z^{[l]})$$

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

$$db^{[l]} = \frac{1}{m} \sum_{i=1}^m dZ^{[l]}$$

→ cache

$$dZ^{[l]} = dA^{[l]} * g'^{[l]}(Z^{[l]})$$

$$dw^{[l]} = \frac{1}{m} dZ^{[l]} * A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} \sum_{i=1}^m dZ^{[l]}$$

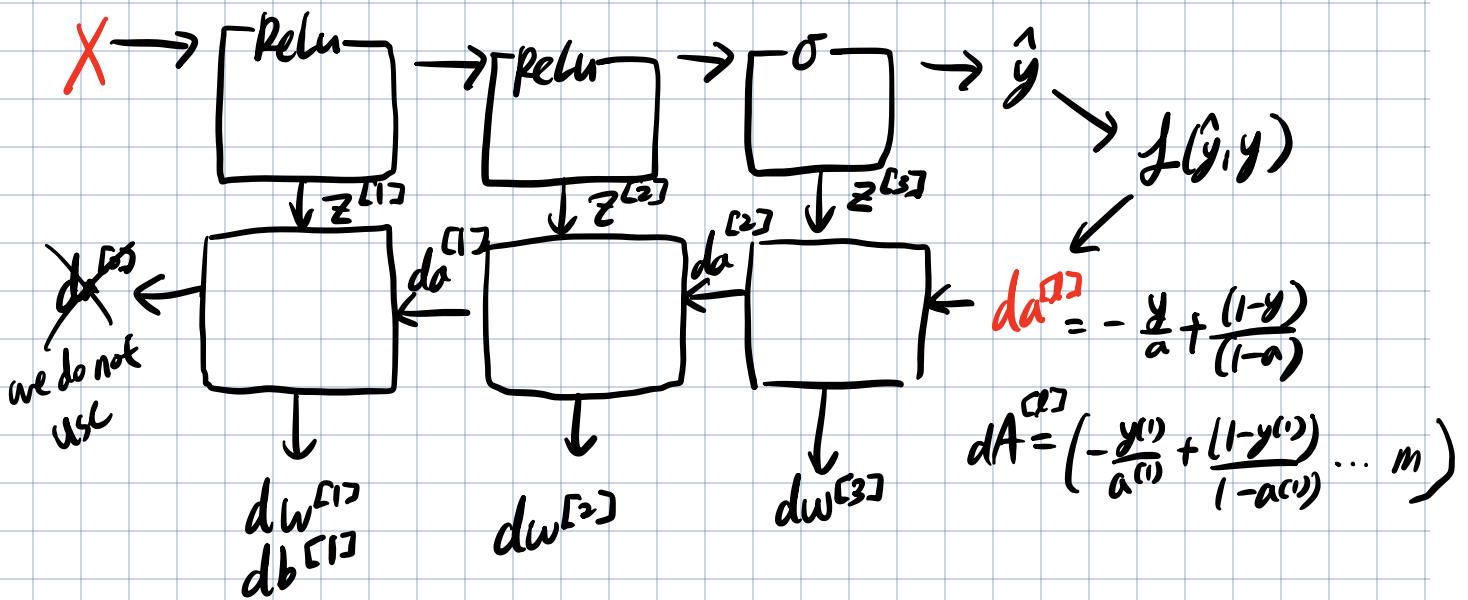
$$da^{[l-1]} = W^{[l]} T \cdot dZ^{[l]}$$

$$\hookrightarrow dZ^{[l]} = W^{[l+1] T} \cdot dZ^{[l+1]} * g^{[l]}(z^{[l]})$$

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

keepdims = True)

- Summary



• Parameters vs Hyperparameters

- Hyperparameters?

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, \dots$

Hyperparameters: learning rate α , # iterations, # hidden layers, control W & b

⇒ determine final values of **parameters** choice of activation functions, momentum, minibatch size, ...

- Applied deep learning is a very empirical process

⇒ If you are not sure what's the best value?



try a single value of α and check if J decreases and try larger value for α

⇒ find α for fast learning & minimal J

- Try a range of values and see if it works
- There is no always -best hyperparameters (\because CPU, GPU change)