

## Week 4 Deep Neural Network

### Deep k-layer neural network

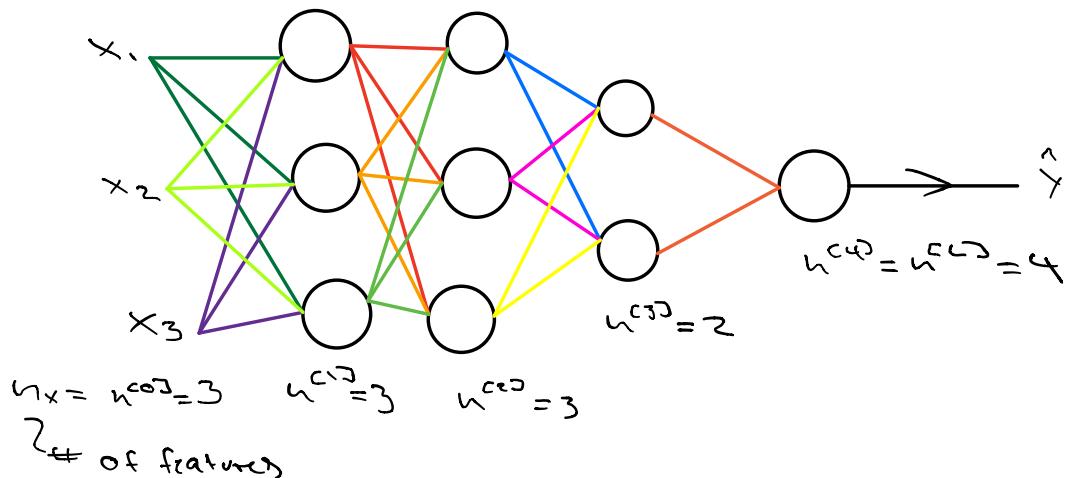
- What is a deep neural network?

↳ the # of hidden layers in a neural network is another hyperparameter that can be tweaked.

↳ " $L=4$ " denotes the # of layers in a neural network with 3 hidden layers + 1 output layer.

↳ " $n^{(j)}$ " describes the units in a particular layer

↳ anything larger than 2-layer is categorized as deep neural net.



### Forward propagation in Deep neural networks

Computing forward propagation

$$\begin{aligned} \text{Given: } z^{(l)} &= w^{(l)}x + b^{(l)} & x &= a^{(0)} \\ a^{(l)} &= g^{(l)}(z^{(l)}) \end{aligned}$$

$$h^{(2)}: \bar{z}^{(2)} = w^{(2)}a^{(1)} + b^{(2)}$$

$$a^{(2)} = g^{(2)}(\bar{z}^{(2)})$$

:

$$h^{(4)}: \bar{z}^{(4)} = w^{(4)}a^{(3)} + b^{(4)}$$

$$a^{(4)} = g^{(4)}(\bar{z}^{(4)}) = \hat{y} \approx \text{estimated output}$$

General rule:  $\bar{z}^{(l)} = w^{(l)}a^{(l-1)} + b^{(l)}$

$$a^{(l)} = g^{(l)}(\bar{z}^{(l)})$$

• Vectorized:

$$\hookrightarrow \bar{z}^{(l)} = w^{(l)}X + b^{(l)} \quad X = A^{(0)}$$

$$\hookrightarrow A^{(l)} = g^{(l)}(\bar{z}^{(l)})$$

:

:

$$\hookrightarrow \hat{y} = g(\bar{z}^{(4)}) = A^{(4)}$$

\* once again, each example represents a column vector and they're stacked together to form a matrix of  $m$  examples.

\* there's a for loop required to create the layers required in the neural net

$\hookrightarrow$  for  $l=1 \dots 4$  (in the example above)  
in forward propagation.

## Getting the dimensions Right

- often used as a debugging tool when writing code.

↳  $W^{[l]} : (n^{[l]}, n^{[l-1]})$   
?  $\underbrace{\quad \quad \quad}_{\# \text{ of neurons}}$  index from previous layer

↳  $b^{[l]} : (n^{[l]}, 1)$

↳  $dW^{[l]} : (n^{[l]}, n^{[l-1]})$

↳  $db^{[l]} : (n^{[l]}, 1)$

\* for single example implementation



## Why Deep Representations?

- it's not only that they have to be big, they must be able with many hidden layers.

- What is a deep network computing?

↳ On an image (for facial recognition)

↳ The first layer might be computing the features or detecting the edges.

↳ The next stage will likely compute a more complex feature (such as an eye) by putting together lots of edges.

↳ Finally, by putting together the more complex features (from the previous layer), the network can try to recognize different types of faces.

↳ Similar functions are composed onto more

Complex features as we move down the neural network. (Similar to context hierarchical representation).

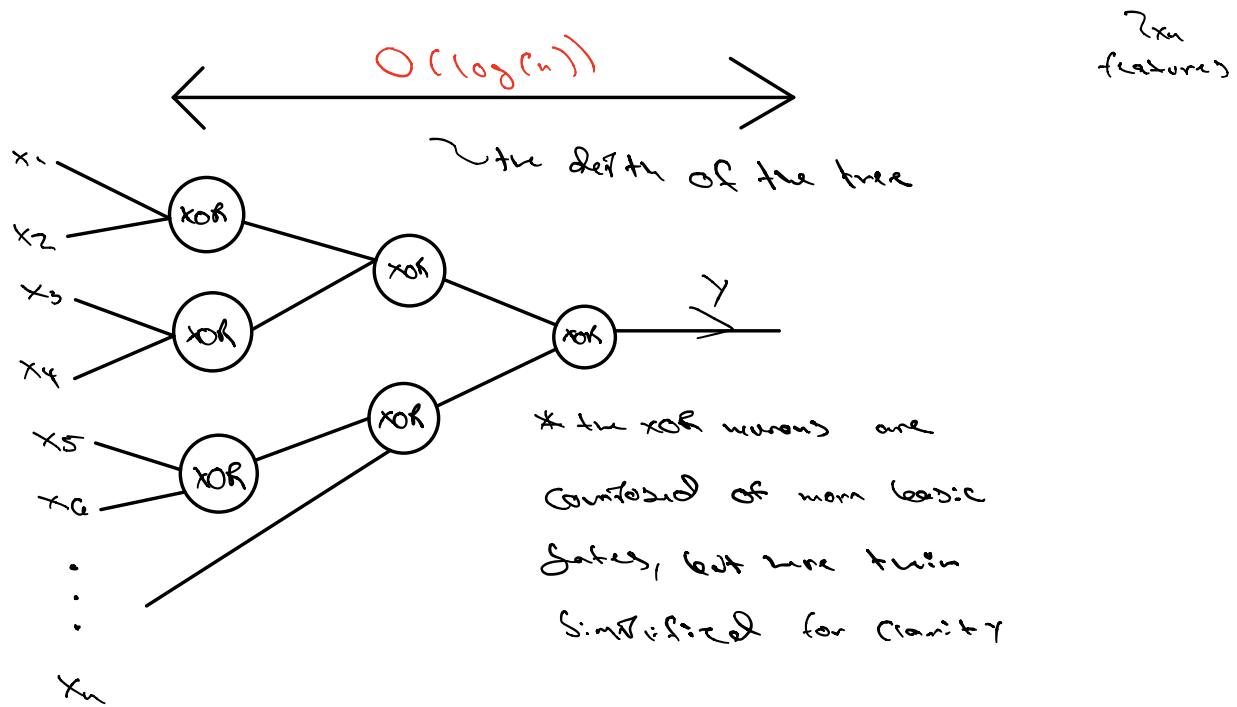
- ↳ in a speech recognition system,
  - ↳ the first layer might learn to detect low level audio waveforms (features) by composing basic low level waveforms
  - we might then learn to detect basic units of sound (phonemes)
  - ↳ combining the phonemes we could learn to combine words
  - ↳ combining the words together we could learn to recognize phrases or sentences.

\* Once again a deep neural net learns in a hierarchical composition way from low level features to more abstract complex constructs.

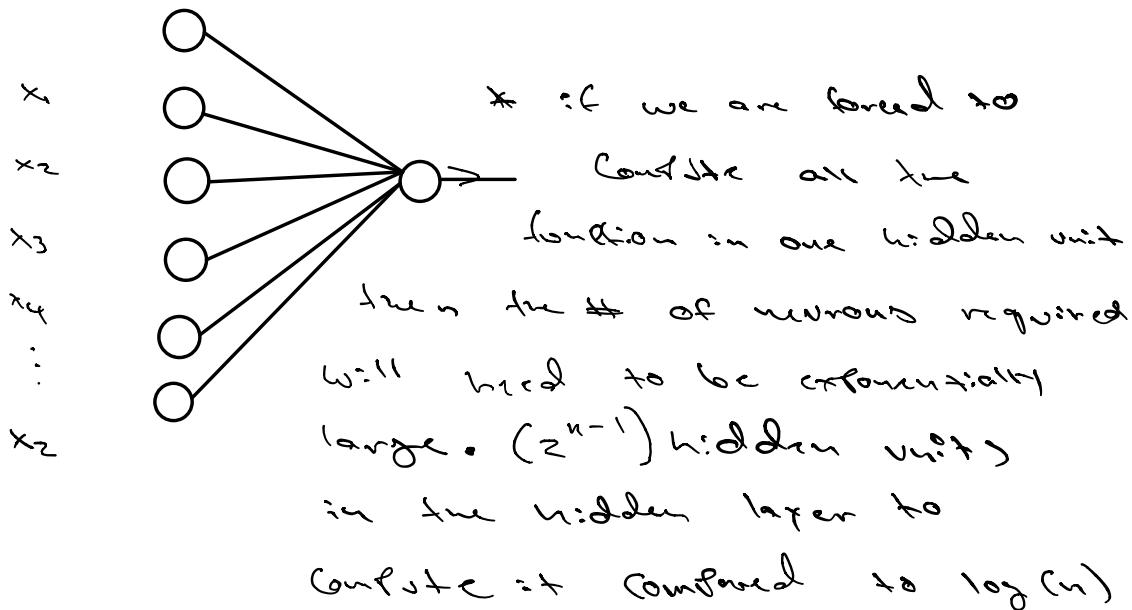
- Circuit theory and deep learning

- ↳ there are functions that can be computed with a "small" L-layer network that shallower networks require exponentially more units to compute (AND/OR gates and such functions)

Example:  $y = x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } \dots \text{ XOR } x_n$

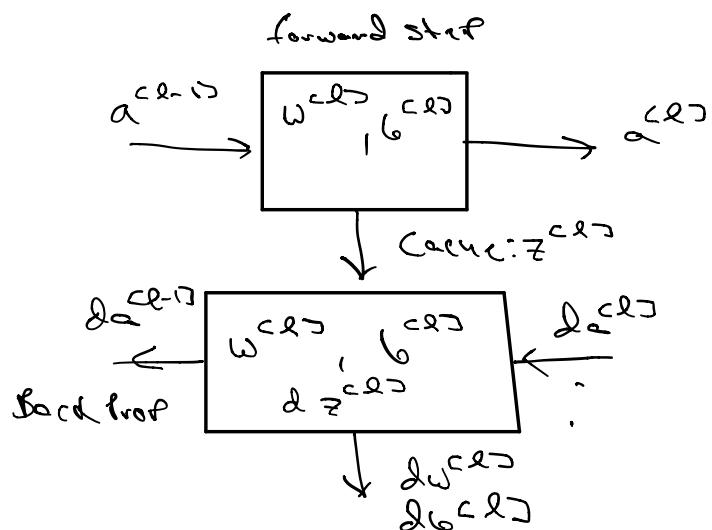


\* the # of gates in this tree setup is not that large as the neurons half every time we advance to the next layer.



## Building Blocks of Deep neural networks

- forward and back word function
- Layer  $l$ :  $W^{cl}, b^{cl}$  activation from previous layer
- forward propagation: Input  $a^{l-1}$ , output  $a^l$
- $$\begin{aligned} L \quad z^{cl} &= W^{cl} a^{l-1} + b^{cl} \\ L \quad a^{cl} &= g^{cl}(z^{cl}) \end{aligned}$$
- \* it is useful to cache the value  $z^{cl}$  for later use. (for the back prop step)
- backward propagation: Input  $d_a^{cl}$  and store cache  $z^{cl}$ . Output  $d_a^{l-1}, dW^{cl}, db^{cl}$



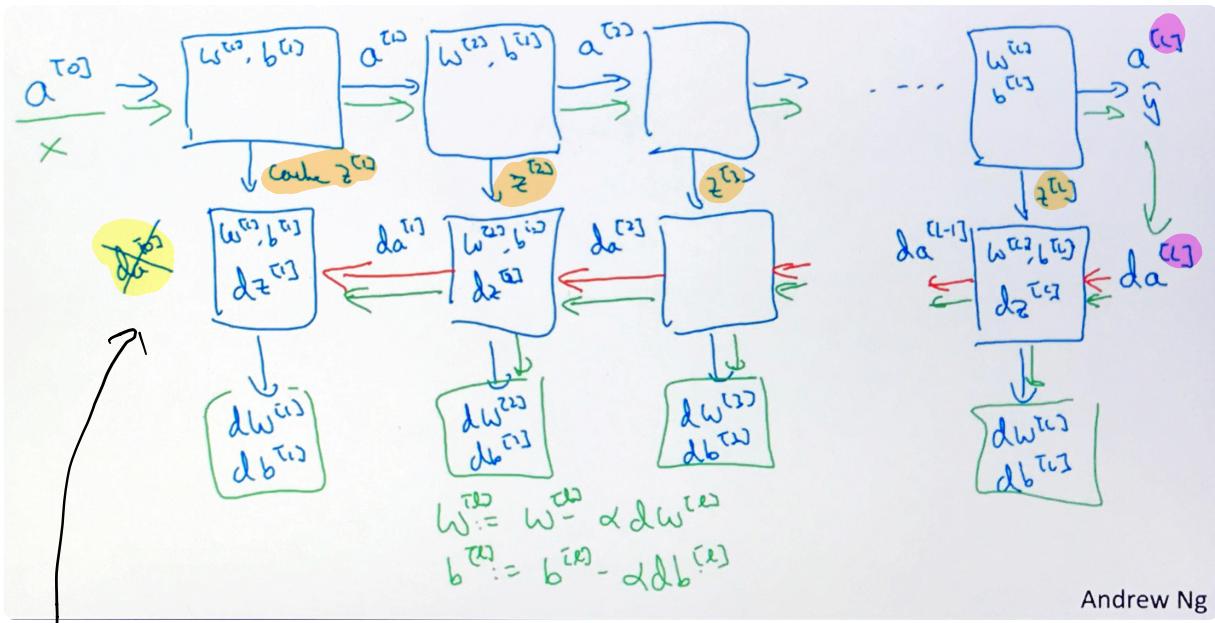
One iteration of : forward

gradient descent : back prop

Capital L

Cache data

We actually also store  $W$  and  $b$  for the back prop.



Andrew Ng

not used for training weights

### Forward and backward propagation

- Backward propagation:

$$\leftarrow \delta z^{[l]} = \delta a^{[l]} * g^{[l]'}(z^{[l]}) \quad \textcircled{1}$$

(element wise product)

$$\leftarrow \delta w^{[l]} = \delta \epsilon^{[l]} \cdot \epsilon^{[l-1]} \quad \text{this is also needed in the cache}$$

$$\leftarrow \delta b^{[l]} = \delta z^{[l]}$$

$$\leftarrow \delta a^{[l-1]} = w^{[l]T} \cdot \delta z^{[l]} \quad \textcircled{2}$$

Substituting \textcircled{2} into \textcircled{1}

$$\leftarrow \delta z^{[l]} = w^{[l+1]T} \cdot \delta z^{[l+1]} * g^{[l]'}(z^{[l]})$$

Vectorized:

$$\leftarrow \delta z^{[l]} = \delta A^{[l]} * g^{[l]'}(z^{[l]}) \quad \textcircled{1}$$

$$\leftarrow \delta w^{[l]} = \frac{1}{m} \cdot \delta z^{[l]} \cdot A^{[l-1]T}$$

$$\begin{aligned} \hookrightarrow \Delta b^{(l)} &= \frac{1}{m} \cdot \text{np.sum}(\Delta z^{(l)}, \text{axis}=1, \text{keepdims=True}) \\ \hookrightarrow \Delta A^{(l-1)} &= \omega^{(l)\top} \cdot \Delta z^{(l)} \quad (2) \end{aligned}$$

### Parameters vs Hyperparameters

• Parameters:  $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}$  (map input to output)

• Hyperparameters examples:

    └ learning rate  $\alpha$

    └ # iterations of gradient descent

    └ # hidden layers ( $L$ )

    └ # of hidden units ( $n^{(1)}, n^{(2)}\dots$ )

    └ Choice of activation functions.

\* all of these hyperparameters we have to specify to our learning algorithm in order to control the ultimate parameters ( $w$  and  $b$ )

(this is the hyperparameter that determine the final value of  $w$  and  $b$ ).

• Overall ML training today is an iterative process where we have an idea, code it, experiment and then iterate based on the results.

    └ it's difficult to know in advance the values of the hyperparameters, therefore, experimentation is required

What does Deep learning have to do with the brain?

- not a whole lot ...
  - ↳ oversimplified Seductive explanation.
- Deep learning is very good at matching complex functions, but analogies to the brain aren't as useful as they once were.