

Deep neural network

- notation
- forward propagation
 - use for loop go through each layer
 - each layer compute

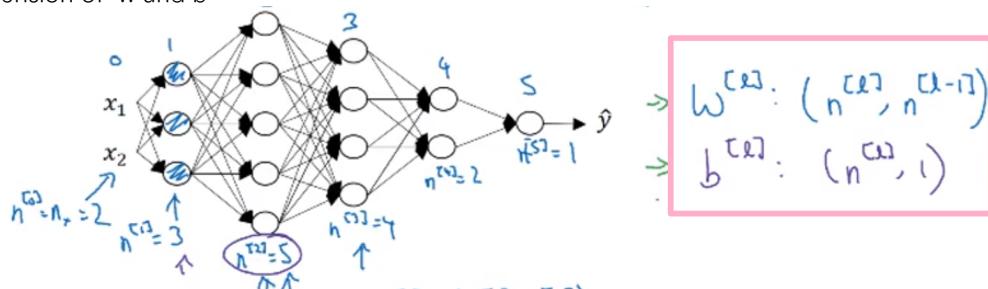
$$Z^{[n]} = W^{[n]} X A^{[n-1]} + b^{[n]}$$

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

, while $g()$ is the active function.

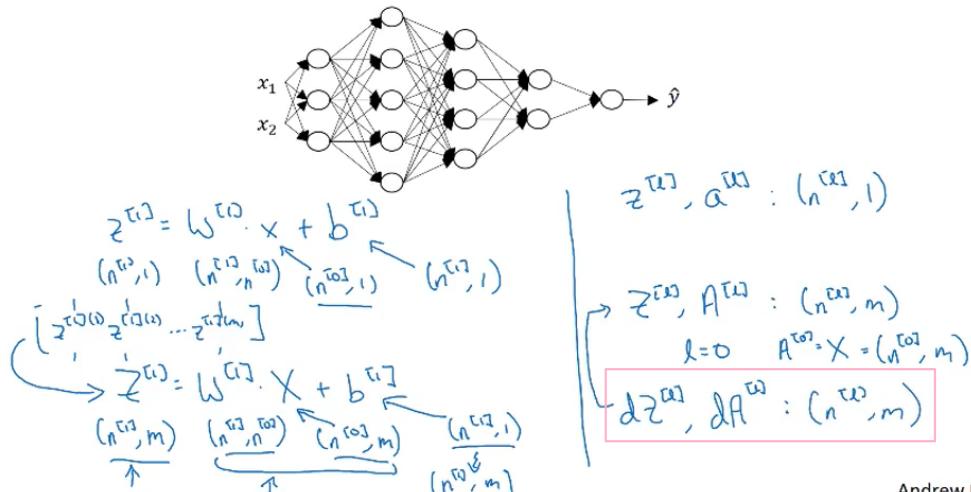
- check your matrix dimensions!

- FP: dimension of w and b



- in vectorized form (pay attention to derivative's dimension)

Vectorized implementation

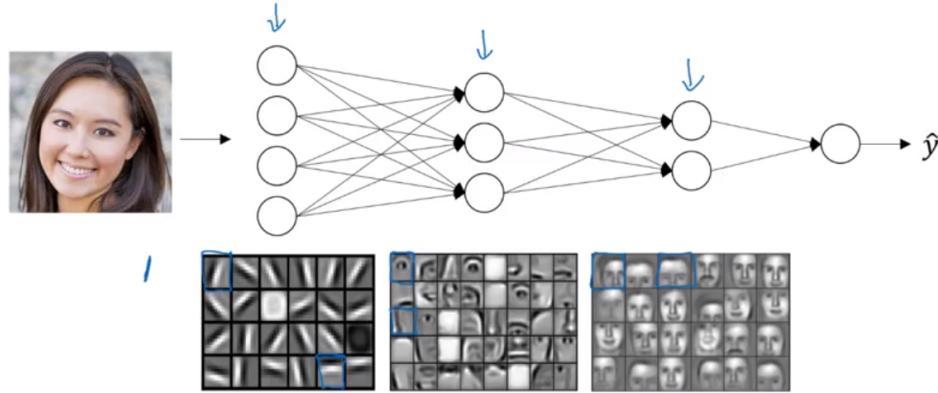


Andrew Ng

- intuition about why deep NN has such amazing power

- summary: each units(low level) \rightarrow (composing) \rightarrow each layer(high level) small to large
 - e.g. CNN

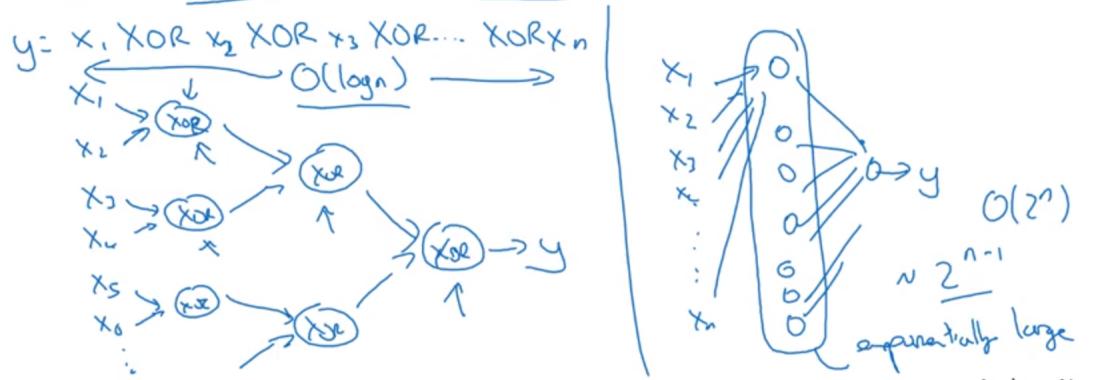
Intuition about deep representation



- e.g. take Circuit as analogy
- more layers \Rightarrow need less "gates" to complete same task

Circuit theory and 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.

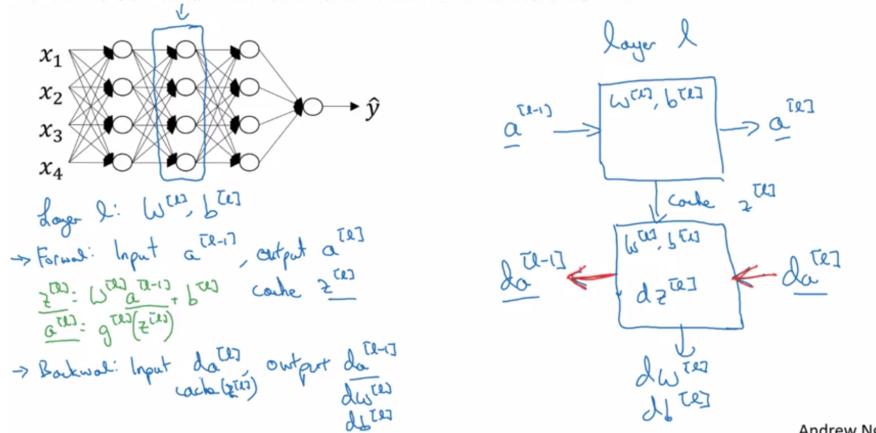


Andrew Ng

- the number of layers could be a **hyper parameters**
- implement deep neural network

- basic block
 - forward propagation
 - use a **cache** to save z
 - backward propagation

Forward and backward functions



- details

Andrew Ng

Forward propagation for layer l

$$\rightarrow \text{Input } a^{[l-1]} \leftarrow \underbrace{w^{[l]}, b^{[l]}}_{\text{Weights}}$$

$$\rightarrow \text{Output } a^{[l]}, \text{cache } (z^{[l]})$$

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

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

$x = A^{[0]} \xrightarrow{\square} \square \xrightarrow{\square} \square \xrightarrow{\square} \dots$

Vertwigl:

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

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

Backward propagation for layer l

$$\rightarrow \text{Input } da^{[l]}$$

$$\rightarrow \text{Output } da^{[l-1]}, dW^{[l]}, db^{[l]}$$

$dz^{[l]} = \cancel{da^{[l]}} \times g^{[l]}(z^{[l]})$
 $dW^{[l]} = dz^{[l]} \cdot a^{[l-1]}$
 $db^{[l]} = dz^{[l]}$
 $da^{[l-1]} = w^{[l]} \cdot dz^{[l]}$
 $dz^{[l]} = w^{[l]} \cdot dz^{[l-1]} \times g^{[l]}(z^{[l]})$

$dz^{[l]} = \cancel{dA^{[l]}} \times g^{[l]}(z^{[l]})$
 $dW^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]}$
 $db^{[l]} = \frac{1}{m} \text{np.sum}(dz^{[l]}, \text{axis}=1, \text{keepdim=True})$
 $dA^{[l-1]} = w^{[l]} \cdot dz^{[l]}$

- some hyper parameters

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

- Optimizer: Momentum, minibatch size, regularizations, ...
- just a empirical process

Forward and backward propagation

$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ Z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(Z^{[2]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) = \hat{Y} \end{aligned}$$

$$\begin{aligned} dZ^{[L]} &= A^{[L]} - Y \\ dW^{[L]} &= \frac{1}{m} dZ^{[L]} A^{[L]T} \\ db^{[L]} &= \frac{1}{m} np.\text{sum}(dZ^{[L]}, axis = 1, keepdims = True) \\ dZ^{[L-1]} &= dW^{[L]T} dZ^{[L]} g'^{[L]}(Z^{[L-1]}) \\ &\vdots \\ dZ^{[1]} &= dW^{[L]T} dZ^{[2]} g'^{[1]}(Z^{[1]}) \\ dW^{[1]} &= \frac{1}{m} dZ^{[1]} A^{[1]T} \\ db^{[1]} &= \frac{1}{m} np.\text{sum}(dZ^{[1]}, axis = 1, keepdims = True) \end{aligned}$$

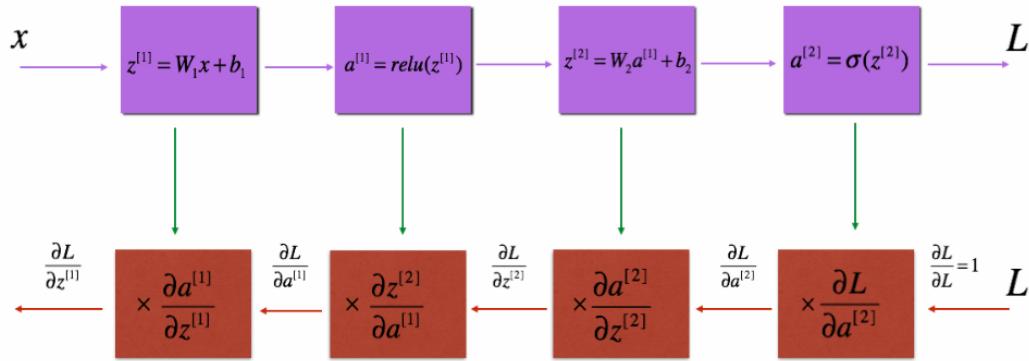


Figure 3 : Forward and Backward propagation for LINEAR->RELU->LINEAR->SIGMOID

The purple blocks represent the forward propagation, and the red blocks represent the backward propagation.

coding!!!

- use caches(datatype:dictionary)to save and pass parameters
- elegant python (programming skill)

```
# Implement [LINEAR -> RELU]*(L-1). Add "cache" to the "caches" list.
for l in range(1, L):
    A_prev = A
    ### START CODE HERE ### (≈ 2 lines of code)
    A, cache = linear_activation_forward(A_prev, parameters['W'+str(l)], parameters['b'+str(l)], activation="relu")
    caches.append(cache)
    ### END CODE HERE ###
```