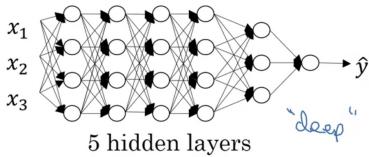
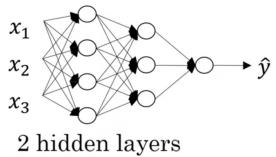
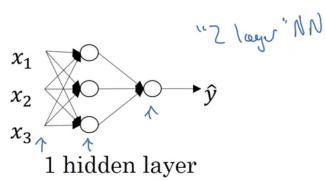
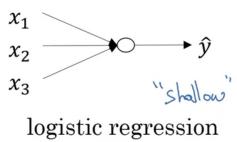
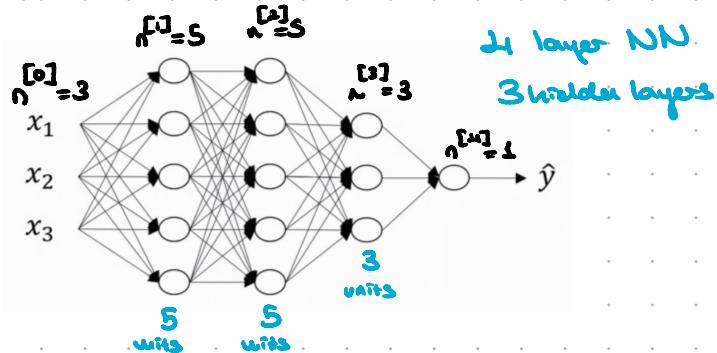


Deep 1-Layer NN



there are functions that
deep NN can learn but
shallow models cannot.

Notations



$$L = \text{number of layers}$$

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

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

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

$w^{[l]}$: weights for $\mathbf{z}^{[l]}$

$b^{[l]}$: bias for $\mathbf{z}^{[l]}$

Forward Propagation in DNN

within a for loop:

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

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

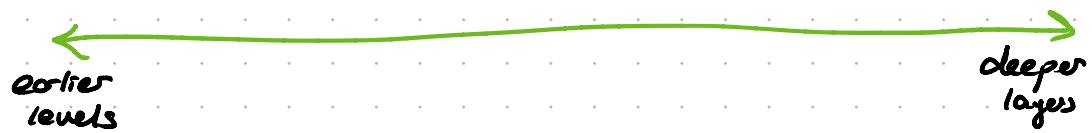
This is a location where we cannot avoid having an explicit for-loop

Why Deep Representations?

Earlier layers learn simple information, then later layers use those simpler information, and compose more complex information

great for compositional information

audio \rightarrow detect low level audio info \rightarrow detect basic units of sound \rightarrow detect words \rightarrow detect sentences



images \rightarrow detect edges \rightarrow detect eye, mouth, nose \rightarrow detect faces

Circuit Theory

informally:

there \exists functions that you can compute with "small" (low # hidden units) 1-layer deep NN

THAT

shallower NNs require exponentially more hidden units to compute

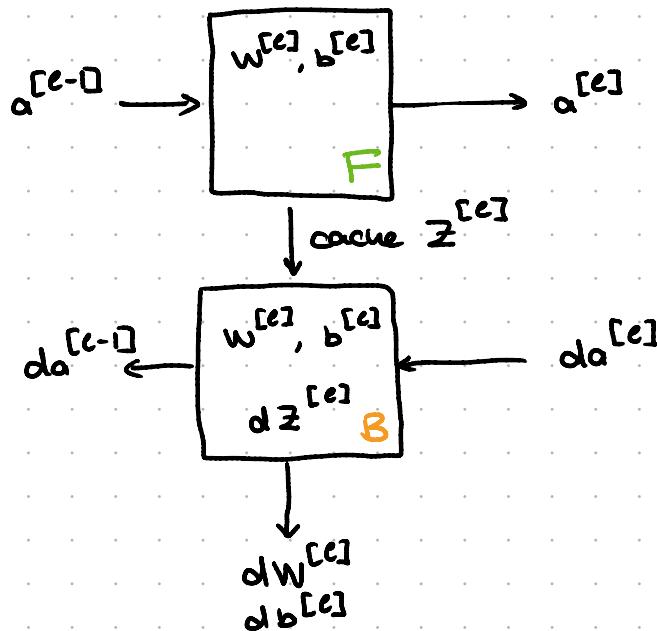
i.e. there are mathematical functions that are much easier to compute with deep networks compared to shallow networks

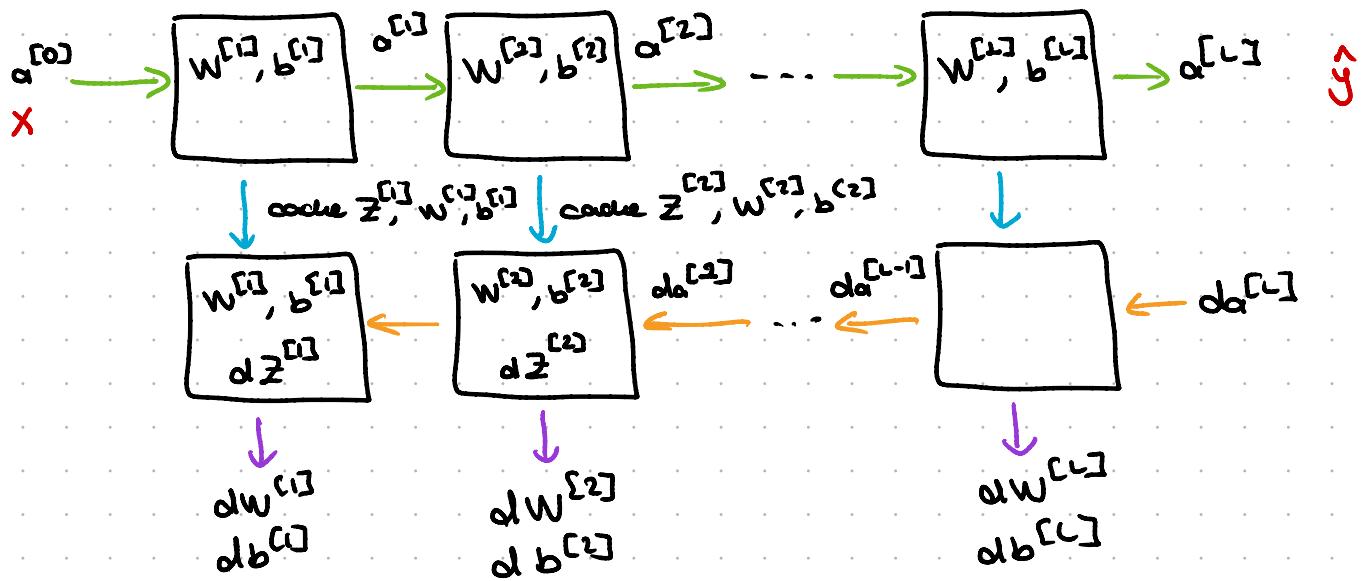
Building Blocks of DNNs

layer c : $w^{[c]}, b^{[c]}$

Forward: input $A^{[c-1]}$, output $A^{[c]}$, cache $Z^{[c]}$

Backward: input $dA^{[c]}$, output $dA^{[c-1]}$, $dW^{[c]}$, $db^{[c]}$
 $Z^{[c]}$ (cached)





BACKWARD PROPAGATION

For layer c

→ input $da^{[c]}$
 → output $da^{[c-1]}, dw^{[c]}, db^{[c]}$

Math Equations:

$$(1) \quad dz^{[c]} = da^{[c]} * g^{[c]}'(z^{[c]})$$

$$dw^{[c]} = dz^{[c]} a^{[c-1]T}$$

$$db^{[c]} = dz^{[c]}$$

$$(2) \quad da^{[c-1]} = w^{[c]T} dz^{[c]}$$

$$\text{by (1)(2)} \Rightarrow dz^{[c]} = w^{[c+1]T} dz^{[c+1]} * g^{[c]}'(z^{[c]})$$

Vectorized:

$$dz^{[c]} = dA * g^{[c]}'(z^{[c]})$$

$$dw^{[c]} = \frac{1}{m} dz^{[c]} A^{[c-1]T}$$

$$db^{[c]} = \frac{1}{m} \text{np.sum}(dz^{[c]}, \text{axis}=1, \text{keepdims=True})$$

$$dA^{[c-1]} = w^{[c]T} dz^{[c]}$$

Initialized As:

$$dA^{[L]} = -\frac{y}{A} + \frac{(1-y)}{(1-A)}$$

(we need ground before)

Parameters vs Hyperparameters

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

hyperparameters: learning rate α

iterations

hidden layers L

hidden units

choice of activation fn



these are the parameters
that control our
parameters