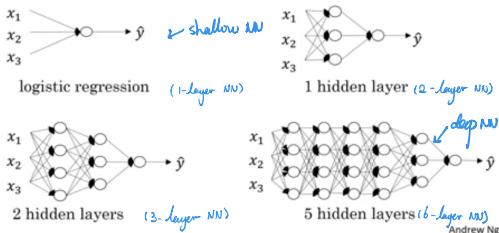
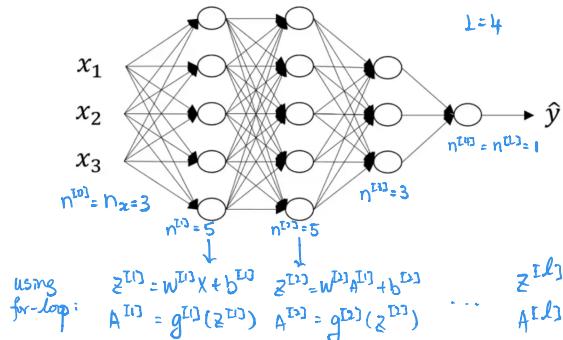


We have seen the structure and mathematical computation of shallow NN. A deep NN is quite similar to a shallow NN but just a few more hidden layers stacking up.

### Deep neural network structure



Notations:



$$\begin{array}{llll} \text{matrix size:} & & & \\ X: n_x \times m & A^{[l]}: n^{[l]} \times m & A^{[l-1]}: n^{[l-1]} \times m & A^{[L-1]}: n^{[L-1]} \times m \\ dw^{[l]}, w^{[l]}: n^{[l]} \times n_x & dw^{[l]}, W^{[l]}: n^{[l]} \times n^{[l]} & dw^{[l]}, W^{[l]}, b^{[l]}: n^{[l]} \times n^{[l-1]} & dw^{[l]}, W^{[l]}, b^{[l]}: n^{[l]} \times n^{[L-1]} \\ db^{[l]}, b^{[l]}: n^{[l]} \times 1 & db^{[l]}, b^{[l]}: n^{[l]} \times 1 & db^{[l]}, b^{[l]}: n^{[l]} \times 1 & db^{[l]}, b^{[l]}: n^{[L]} \times 1 \\ dz^{[l]}, z^{[l]}: n^{[l]} \times m & dz^{[l]}, z^{[l]}: n^{[l]} \times m & dz^{[l]}, z^{[l]}: n^{[l]} \times m & dz^{[l]}, z^{[l]}: n^{[L]} \times m \Rightarrow Y, A^{[L]}: n^{[L]} \times m \end{array}$$

### Computation

initialising: parameter ["w" + str(i)] = 0.01 \* np.random.randn(layer[i], layer[i-1])

parameter ["b" + str(i)] = 0.01 \* np.random.randn(layer[i], 1)

for layer  $l$ , given  $w^{[l]}, b^{[l]}, A^{[l-1]}$ , for  $l=1, 2, \dots, L$

$$\begin{aligned} z^{[l]} &= w^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

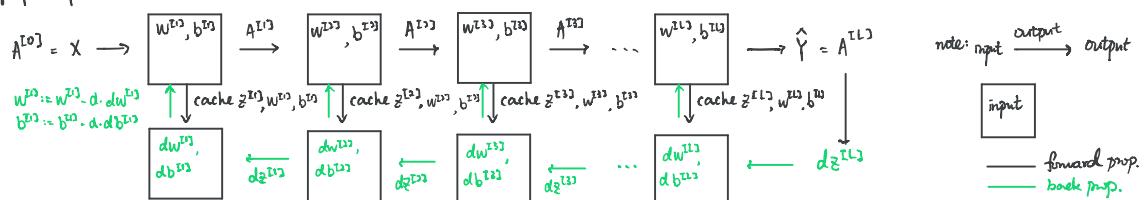
$$\text{back propagation: } dz^{[l]} = dA^{[l]} \times g'^{[l]}(z^{[l]}) \rightarrow \text{so } dz^{[l-1]} = w^{[l]T} dz^{[l]} \times g'^{[l-1]}(z^{[l-1]}) \text{ and } dz^{[L]} = A^{[L]} - Y$$

where  $dA^{[l-1]} = w^{[l]T} dz^{[l]}$  and  $dA^{[L]} = \left[ \frac{y^{(1)}}{\alpha^{(1)}} + \frac{1-y^{(1)}}{1-\alpha^{(1)}} \dots - \frac{y^{(m)}}{\alpha^{(m)}} + \frac{1-y^{(m)}}{1-\alpha^{(m)}} \right] = -\frac{Y}{A^{[L]}} + \frac{1-Y}{1-A^{[L]}}$

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

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

graph representation:



### Hyperparameter and parameter.

- parameters:  $w^{[l]}$ ,  $b^{[l]}$
- hyperparameters: learning rate  $\alpha$   
# iteration  
# layer 1  
# hidden units  $n^{[1]}, n^{[2]}, \dots$   
choice of activation functions  
momentum, minibatch size, regularization variables ...

→ control the values of output.

→ control the values of parameters,  
need tuning to find the best set.