

Summary:

$$\begin{aligned}
 dz^{[2]} &= a^{[2]} - y \\
 dW^{[2]} &= dz^{[2]} a^{[1]T} \\
 db^{[2]} &= dz^{[2]} \\
 dz^{[1]} &= W^{[2]T} dz^{[2]} * g'(z^{[1]}) \\
 dW^{[1]} &= dz^{[1]} x^T \\
 db^{[1]} &= dz^{[1]}
 \end{aligned}$$

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

Vectorised:

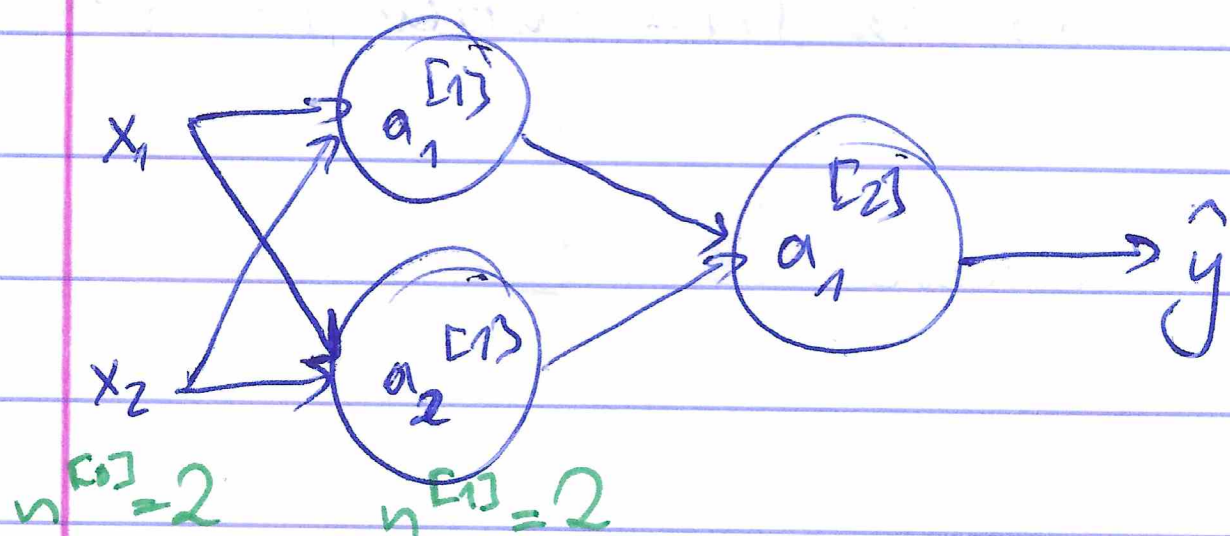
$$\begin{aligned}
 J(\dots) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) \\
 dZ^{[2]} &= A^{[2]} - y \\
 dW^{[2]} &= \frac{1}{m} dZ^{[2]} A^{[1]T} \\
 db^{[2]} &= \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims}) \\
 dz^{[1]} &= W^{[2]T} dZ^{[2]} * g'(z^{[1]}) \\
 dW^{[1]} &= \frac{1}{m} dz^{[1]} x^T \\
 db^{[1]} &= \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims})
 \end{aligned}$$

$(n^{[1]}, m)$   $(n^{[2]}, m)$   $(n^{[1]}, m)$

Full derivation is hard.

## Random Initialisation (of weights)

For log. reg., it's OK to init. weights to 0 - not for NN:



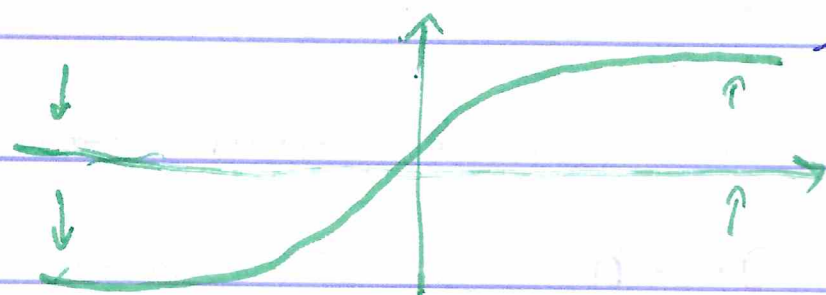
- Identical (symmetrical) hidden units - still computing the same function after iter.

$$W^{[0]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad b^{[0]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad W^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$a_1^{[0]} = a_2^{[0]}, \quad dz_1^{[0]} = dz_2^{[0]}$$

$$dW = \begin{bmatrix} u & v \\ u & v \end{bmatrix}, \quad W^{[1]} := W^{[1]} - \alpha dW, \quad W^{[1]} = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

So:



$$W^{[0]} = \text{np.random.randn}(2, 2) * 0.01$$

We like very small values  
 $z^{[0]} = W^{[0]} x + b^{[0]}$   
 $a^{[0]} = g^{[0]}(z^{[0]})$

if  $W$  too large, very large  $z$  at start of training

No symmetry breaking problem for  $b$

$$b^{[0]} = \text{np.zeros}(2, 1)$$