

There could be better constants than 0.01 (perhaps for larger NNs?)

$$W^{[2]} = \text{np.random.randn}(1, 2) * 0.01$$

$$b^{[2]} = 0$$

### Week 3 quiz notes

-  $[l]$  → layer  
 $(n)$  → train. example  
 $i$  → node

$$X = \begin{bmatrix} \vdots & \vdots & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(n)} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

- tanh works better than  $\sigma$  for hidden units, because the mean of its output  $\approx 0$ , and so it centers the data better for the next layer

- Vectorised for prog. for layer  $l$ , where  $1 \leq l \leq L$ :

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

$$A^{[l]} = \sigma(Z^{[l]})$$

$$g(z) = \sigma(z)$$

- Cucumbers ( $y=1$ ) v. Watermelons ( $y=0$ ), which act. func. for output layer?

↳ SIGMOID.

The output value from a sigmoid can be easily understood as a probability. Sigmoid outputs  $0 \leq \text{values} \leq 1$ , so it can be classified as 0 if  $< 0.5$  and 1 if  $> 0.5$ . tanh ( $-1 \leq \text{values} \leq 1$ ) is less convenient

$$A = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}, B = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \left| \quad B = \text{np.sum}(A, \text{axis}=1, \text{keepdims=True})$$

$= 0$ : ~~vert.~~  $(y)$   
 $= 1$ : ~~hor.~~  $(x)$