



$\underline{x}$ : input

$\underline{a}$ : output after activation function

$\underline{z}$ : output before activation function

$f$ : activation function, sigmoid

$\underline{W}^{(l)}$ : weight matrix from layer  $l$  to layer  $l+1$  of size  $\underbrace{S_{l+1}}_{\text{units in layer } l+1} \times \underbrace{(S_l+1)}_{\text{units in layer } l + \text{bias}}$

FORWARD PASS

$$\underline{a}^{(1)} = \underline{x}$$

$$\rightarrow \text{append bias: } \underline{a}^{(1)} \leftarrow [a_0^{(1)}=1, \underline{a}^{(1)T}]^T = [1 \ x_1 \ x_2 \ x_3]^T \quad (3+1) \times 1$$

$$\underline{z}^{(2)} = \underline{W}^{(1)} \cdot \underline{a}^{(1)} \quad (4 \times 4) \times (4 \times 1) = 4 \times 1$$

$$\underline{a}^{(2)} = f(\underline{z}^{(2)})$$

$$\rightarrow \text{append bias: } \underline{a}^{(2)} \leftarrow [a_0^{(2)}=1, \underline{a}^{(2)T}]^T \quad (4+1) \times 1$$

$$\underline{z}^{(3)} = \underline{W}^{(2)} \cdot \underline{a}^{(2)} \quad (4 \times 5) \times (5 \times 1) = 4 \times 1$$

$$\underline{a}^{(3)} = f(\underline{z}^{(3)})$$

$$\rightarrow \text{append bias: } \underline{a}^{(3)} \leftarrow [a_0^{(3)}=1, \underline{a}^{(3)T}]^T \quad (4+1) \times 1$$

$$\underline{z}^{(4)} = \underline{W}^{(3)} \cdot \underline{a}^{(3)} \quad (3 \times 5) \times (5 \times 1) = 3 \times 1$$

$$\underline{a}^{(4)} = f(\underline{z}^{(4)}) = \underline{h}_{\underline{W}}(\underline{x}) : 3 \times 1$$

$\underline{W}$   $\rightarrow$  model / neural network.

## BACKWARD PASS

$\delta_j^{(l)}$ : error of node/unit  $j$  in layer  $l$

- For each output unit  $j$  in last layer  $l=4$

$$\delta_j^{(4)} = (y_j - \hat{y}_j) = (y_j - a_j^{(4)})$$

↳ assemble:  $\tilde{\delta}^{(4)} = [\delta_1^{(4)}, \dots, \delta_j^{(4)}, \dots]$

← num classes →

in Udacity it's multiplied by  $f'(z)$ :  
 $\delta_j^{(4)} = (y_j - \hat{y}_j) \cdot f'(z_j^{(3)})$

note: for the sigmoid:  
 $f' = f \cdot (1 - f)$

- Propagate errors to all units of all layers

$$\tilde{\delta}^{(3)} = \left( \underbrace{\left( \underbrace{\tilde{W}^{(3)}}_{5 \times 3} \right)^T \cdot \underbrace{\tilde{\delta}^{(4)}}_{3 \times 1}}_{5 \times 1} \right) * \left[ 1, \underbrace{f'(z^{(3)})^T}_{1 \times 4} \right]$$

↳ multiply 1 by 4

extend because of bias, so that dims match, but then it's removed!

$$= \left[ \underbrace{\delta_0^{(3)}}_{\text{bias}}, \delta_1^{(3)}, \dots \right]^T$$

: each unit in layer 3 has an error.

↳  $\delta^{(3)} \leftarrow \tilde{\delta}^{(3)} [1:]$  remove bias component:  $4 \times 1$

$$\tilde{\delta}^{(2)} = \left( \underbrace{\left( \underbrace{\tilde{W}^{(2)}}_{5 \times 4} \right)^T \cdot \underbrace{\delta^{(3)}}_{4 \times 1}}_{5 \times 1} \right) * \left[ 1, \underbrace{f'(z^{(2)})^T}_{4 \times 1} \right]$$

initialized as  $\Delta \tilde{W} = 0$

↳  $\delta^{(2)} \leftarrow \tilde{\delta}^{(2)} [1:]$  remove bias component:  $4 \times 1$

- Weight updates: For each sample:

$\hat{y} = \text{forward}(x) \rightarrow a^{(1)}, a^{(2)}, a^{(3)}, a^{(4)} = \hat{y}$

backward( $\hat{y}$ )  $\rightarrow \delta^{(4)}, \delta^{(3)}, \delta^{(2)}$

$$\Delta \tilde{W}^{(l)} = \Delta \tilde{W}^{(l)} + \underbrace{\delta^{(l+1)}}_{S_{l+1} \times 1} \cdot \underbrace{a^{(l)T}}_{1 \times (S_l + 1)}$$

Then, after each epoch

$$\tilde{W}^{(l)} = \underbrace{\tilde{W}^{(l)}}_{\text{old}} + \underbrace{\frac{\alpha}{m} \Delta \tilde{W}^{(l)}}_{\text{step}}$$