$X$

$X_1$
$X_2$
$X_3$

bias

$\hat{Y}$

$\hat{Y}_K = a^{(4)}_K$

$a^{(1)}$  $a^{(2)}$  $a^{(3)}$  $a^{(4)}$

INPUT     HIDDEN      OUTPUT

layer 1   layer 2   layer 3   layer 4

$X$: input

$a$: output after activation function

$z$: output before activation function

$f$: activation function, sigmoid

$W^{(\ell)}$: weight matrix from layer $\ell$ to layer $\ell + 1$ of size $\underbrace{S_{\ell+1}}_{\substack{\text{units } m \\ \text{layer } \ell+1}} \times \underbrace{(S_\ell + 1)}_{\substack{\text{units} \\ m \\ \text{layer } \ell}} \Big\rangle$ bias

— FORWARD PASS —

$a^{(1)} = X$

$\rightarrow$ append bias: $a^{(1)} \leftarrow [a_0^{(1)} = 1, a^{(1)T}]^T = [1 \ x_1 \ x_2 \ x_3]^T$  $\underset{4}{\underbrace{(3+1) \times 1}}$

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

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

$\rightarrow$ append bias: $a^{(2)} \leftarrow [a_0^{(2)} = 1, a^{(2)T}]^T$   $\underset{5}{\underbrace{(4+1) \times 1}}$

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

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

$\rightarrow$ append bias: $a^{(3)} \leftarrow [a_0^{(3)} = 1, a^{(3)T}]^T$   $\underset{5}{\underbrace{(4+1) \times 1}}$

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

$a^{(4)} = f(z^{(4)}) = h_W(X)$ : $3 \times 1$

$\longrightarrow$ model / neural network.

Note: $W^{(\ell)} = $

bias weight

$S_{\ell+1}$

$\updownarrow S_{\ell+1}$

$W^{(\ell)}$ is the weight matrix of connections that go from units in $\ell$ to units in $\ell+1$



$\ell$   $\ell+1$

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

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

$$\delta_j^{(4)} = \overbrace{(y_j - \hat{y}_j)}^{e^{(4)}} \cdot f'(z_j^{(3)})$$

In Udacity, the es are the errors, $\delta$ is the delta

$\rightarrow$ assemble: $\underset{\sim}{\delta}^{(4)} = [\,\delta_1^{(4)}, \ldots, \delta_j^{(4)}, \ldots\,]$

$\xleftarrow{\text{num classes}}$

In Udacity it's multiplied by $f'()$; Andrew Ng didn't use $f'$ in the last layer, but I understand that depends in the activation function!

note: for the sigmoid:
$$f' = f \cdot (1-f);$$
if no activation: $f(x)=x$
$f'=1$.

- Propagate errors to all units of all layers

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

$\to$ multiply $1$ by $1$

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

I understand we could extend the bias at the end of the vector, too...

$$= [\,\underset{\text{bias}}{\delta_0^{(3)}}, \delta_1^{(3)}, \ldots\,]^T \quad : \text{ each unit in layer 3 has an error.}$$

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

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

$5\times 1$

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

initialized as $\Delta W = 0$

- Weight changes: For each sample :
  $\hat{y} = \text{forward}(x) \longrightarrow a^{(1)}, a^{(2)}, a^{(3)}, a^{(4)} = \hat{y}$
  $\text{backward}(\hat{y}) \longrightarrow \delta^{(4)}, \delta^{(3)}, \delta^{(2)}$ (no $\delta^{(1)}$!)
  $\underset{\sim}{\Delta W}^{(l)} = \underset{\sim}{\Delta W}^{(l)} + \underbrace{\delta^{(l+1)}}_{S_{l+1}\times 1} \cdot \underbrace{a^{(l)T}}_{1\times(S_l+1)}$

Note: if we are doing complete batch gradient descend, we update weights after each epoch, but if stochastic, we can update weights for every example pass

Then, after each epoch:
$$\underset{\substack{\sim \\ \text{old}}}{W}^{(l)} = \underset{\sim}{W}^{(l)} + \underbrace{\frac{\alpha}{m}\underset{\sim}{\Delta W}^{(l)}}_{\text{step}}$$

In the videos an exercises from Udacity, the weight matrix is defined the other way around, compared to how Andrew Ng does it:

- Andrew Ng: $\underset{\sim}{W} : (S_{\ell+1}, S_\ell + 1)$ : (new units, old units + 1 bias)

- Udacity: $\underset{\sim}{W} :$ (old units, new units)

Given the order how we insert the sizes the layers / weight matrices, the Udacity approach is maybe more intuitive.

Note that in that case, the multiplication order needs to be changed to match the sizes:

$$\underset{\sim}{Z} = \underset{\sim}{a} \cdot \underset{\sim}{W} \quad : \quad (batch, old) \times (old, new) = (batch, new)$$

$$\underset{\sim}{Z} \leftarrow \underset{\sim}{Z} + \underset{\sim}{b} \quad : \quad bias\ is\ added,\ not\ inserted\ in\ \underset{\sim}{W}$$

$$\underset{\sim}{b} : (batch, new)$$