

X : input

a : output after activation function

z : output before activation function

f : activation function, sigmoid

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

FORWARD PASS

$$\tilde{a}^{(1)} = X$$

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

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

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

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

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

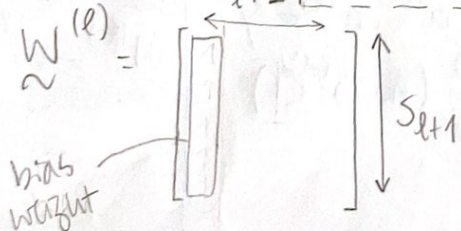
$$\tilde{a}^{(3)} = f(\tilde{z}^{(3)})$$

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

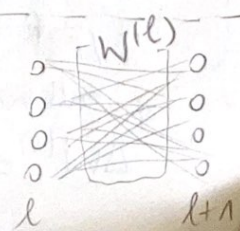
$$\tilde{z}^{(4)} = \tilde{W}^{(3)} \cdot \tilde{a}^{(3)} \quad (3 \times 5) \times (5 \times 1) = 3 \times 1$$

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

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



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



BACKWARD PASS

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

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

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

In Udaan, the e_s are the errors, δ is the delta

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

← num classes →

in Udaan it's multiplied by $f'(z)$; 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

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

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...

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

each unit in layer 3 has an error.

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

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

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

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

- Weight changes: For each sample:

$\hat{y} = \text{forward}(x) \rightarrow a^{(1)}, a^{(2)}, a^{(3)}, a^{(4)} = \hat{y}$
 $\text{backward}(\hat{y}) \rightarrow \delta^{(4)}, \delta^{(3)}, \delta^{(2)}$ (no $\delta^{(1)}$!)

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

Then, after each epoch

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

initialized as $\Delta W = 0$