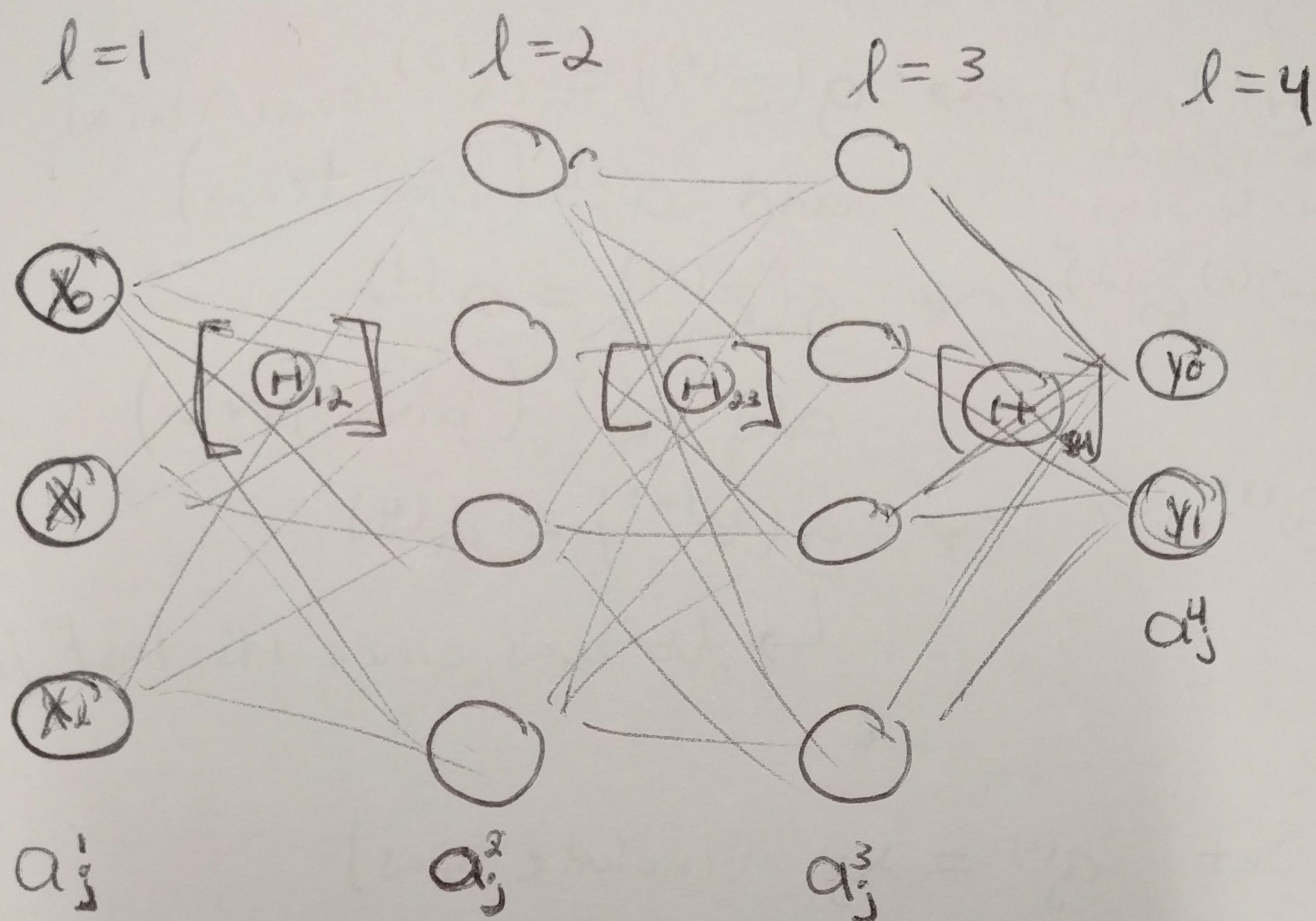


# BASIC FF-NN Implementation

①



$x_j^{(i)}$  = input from the  $i$ -th training example, describing the  $j$ -th feature.

$\textcircled{H}_{AB}^{(l)}$  = Weight matrix, mapping inputs from the  $l$  layer to the  $l+1$  layer. [Dims:  $S_{l+1} \times S_l$ ]

Note:  $S_{l+1} \times S_l + 1$ , if bias is added.

$$a_j^{(l)} := g(\textcircled{+}^{(l-1)} \cdot a_j^{(l-1)}) = g(z^{(l)}) \sim \text{activation function value}$$

$L$  := number of layers

$m$  := number of training examples.

②

## \* Feed-forward pass

- Set  $a_j^{(1)} = x_j$

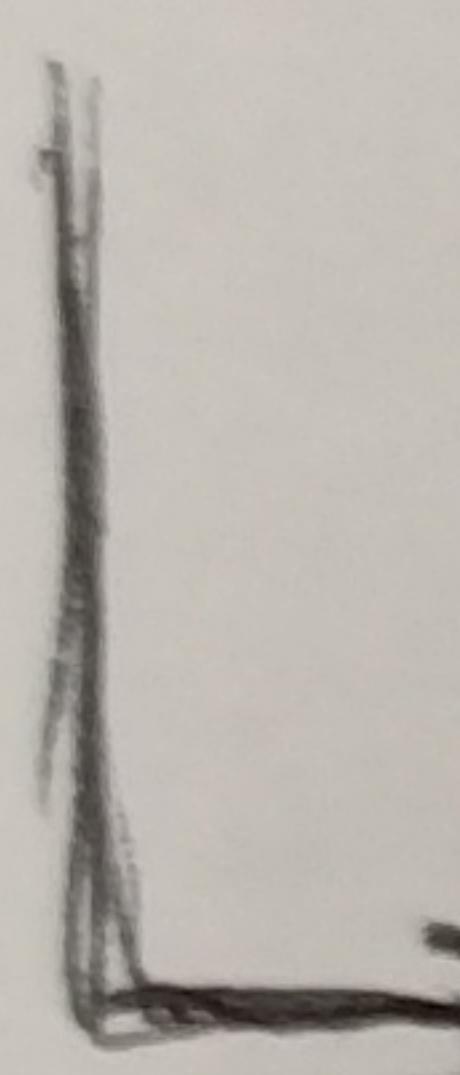
$$\rightarrow z^{(2)} = \Theta^{(1)} a^{(1)} \rightsquigarrow g(z^{(2)}) = a^{(2)} \quad \begin{matrix} s_2 \times 1 \\ s_2 \times s_1 \\ s_{l+1} \times 1 \end{matrix} \quad \begin{matrix} s_2 \times 1 / s_{l+1} \times 1 \\ \text{add } a_0^{(2)} \text{ (bias term)} \end{matrix}$$

$$\rightarrow z^{(3)} = \Theta^{(2)} a^{(2)} \rightsquigarrow g(z^{(3)}) = a^{(3)}$$

add  $a_0^{(3)}$  (bias term)

$$\rightarrow z^{(4)} = \Theta^{(3)} a^{(3)} \rightsquigarrow g(z^{(4)}) = a^{(4)}$$

↳ No bias since it's last layer.



Set  $a_j^{(l)} = x_j$  (include bias)

For  $l=1$  to  $l=L-1$ :

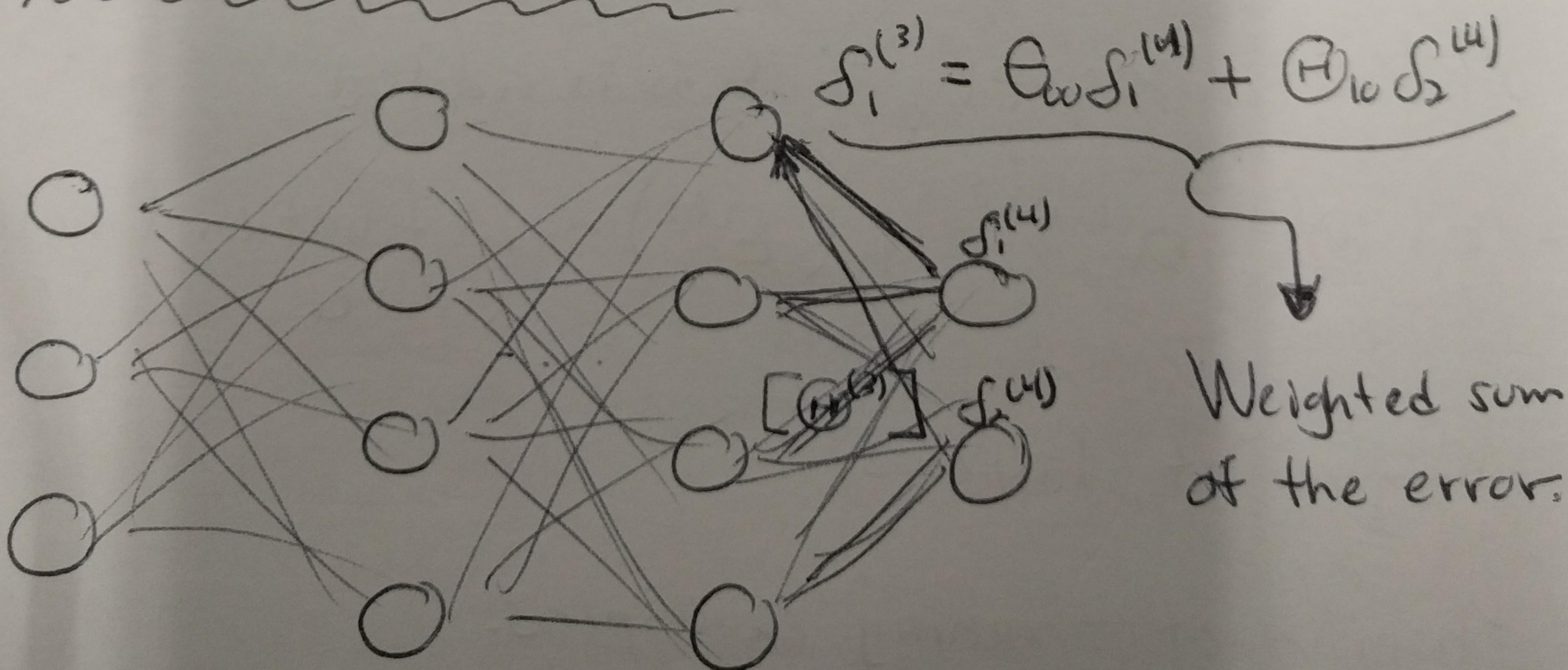
$$z^{(l+1)} = \Theta^{(l)} a^{(l)} \rightsquigarrow g(z^{(l+1)}) = a^{(l+1)}$$

add  $a_0^{(l+1)}$  (bias)

if  $l+1 \neq L$

6

## \* BACK PROPAGATION



- ② • Roughly we can express the error and derivative as follows:

$$E = \frac{1}{2m} \sum_{j=1}^m \sum_{k=1}^{K \text{ classes}} (g(\Theta x) - y_k)^2$$

$$\frac{\partial E}{\partial \Theta_j} \approx -\frac{1}{m} (y - g(\Theta x)) g'(\Theta x) x_j$$

$$\frac{\partial E}{\partial \Theta_j} \approx (\text{error-term}) g'(\text{input})$$

activation fn. derivative.



- Computation of propagated error.

$$\delta^{(4)} = a_j^{(4)} - y_j \quad \sim \quad \delta^{(L)} = (a_j^{(L)} - y_j) g'(a_j^{(L)})$$

assumes there is no act. fn. in last layer.

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} * g'(z^{(3)})$$

$$\delta^{(l-1)} = (\Theta^{(l-1)})^T \delta^{(l)} * g'(z^{(l-1)})$$

$\underbrace{s_{l-1} \times 1}_{(1)}$      $\underbrace{s_{l-1} \times s_l}_{s_{l-1} \times 1}$      $\underbrace{s_l \times 1}_{s_{l-1} \times 1}$      $\underbrace{s_{l-1} \times 1}_{s_{l-1} \times 1}$

Mathematically it can be proved:  $\frac{\delta}{\delta \Theta_{i,j}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$

$$\delta^{(l+1)} \cdot (a^{(l)})^T = \Delta^{(l)} \rightarrow \begin{array}{l} \text{Update increment} \\ \text{for weight matrix.} \end{array}$$

$s_{l+1} \times 1$      $1 \times s_l$      $s_{l+1} \times s_l$

## \* BP-Algorithm

- Training set:  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$
- Set  $\Delta_{ij}^{(l)} = 0 \quad \forall i, j, l \rightarrow$  Accumulators for partial derivatives
- For  $i=1 \text{ to } m:$ 
  - $\rightarrow \text{Set } a^{(1)} = x^{(i)} \quad [\text{Input layer}]$
  - $\rightarrow \text{Forward Propagation} \rightarrow \text{compute all } a^{(l)}$
  - $\rightarrow \delta^{(L)} = (a^{(L)} - y^{(i)}) g'(z^{(L)}) \quad [\text{ERROR OUTPUT LAYER}]$
  - $\rightarrow \text{Compute: } \delta^l + \sum_{l < L} \delta^{l+1} \cdot \underbrace{\Theta^{(l)}}_{\text{Weight propagation}} \cdot g'(z^{(l)})$
  - $\rightarrow \text{Accumulate:}$ 

$$\begin{aligned} \Delta_{ij}^{(l)} &:= \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)} \\ &:= \Delta^{(l)} + f_i^{(l+1)} (a_j^{(l)})^T \end{aligned}$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \quad \begin{cases} \text{Final derivative.} \\ \text{Bias term} \end{cases}$$

- Update weights:

$$\Theta_{ij}^{(l)} := \Theta_{ij}^{(l)} + \alpha D_{ij}^{(l)}$$

$\hookrightarrow$  learning rate