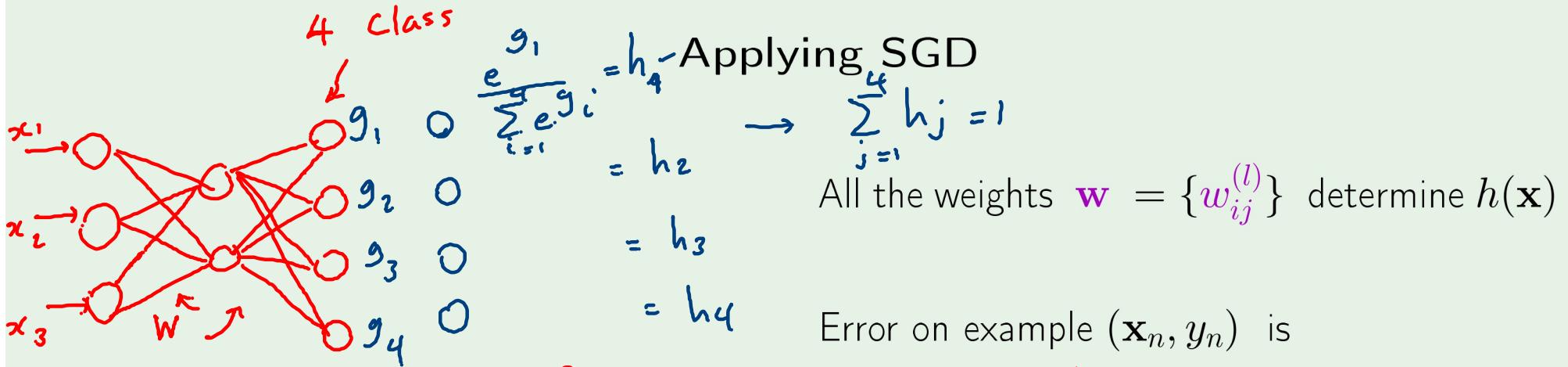


Applying SGD

All the weights $\mathbf{w} = \{\mathbf{w}_{ij}^{(l)}\}$ determine $h(\mathbf{x})$



Reg. \leftarrow Sq. Err.

$e(h(\mathbf{x}_n), y_n) = e(w)$

Cross Ent. \downarrow

predicted output \downarrow

true output

$$\min_w -\log h_3$$

$$y_1 = 3$$

$$y_2 = 1$$

$$y_5 = 4$$

Multi-Class

$$\min_w -\log h_3(x_1) - \log h_1(x_2) - \dots - \log h_4(x_5) - \log h_{y_n}(x_n)$$

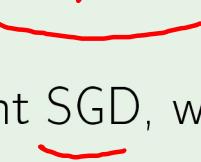
w

$l(w)$

Applying SGD

All the weights $\mathbf{w} = \{\mathbf{w}_{ij}^{(l)}\}$ determine $h(\mathbf{x})$

Error on example (\mathbf{x}_n, y_n) is

$$\mathbf{e}(h(\mathbf{x}_n), y_n) = \mathbf{e}(\mathbf{w})$$


To implement SGD, we need the gradient

$$\nabla_{\mathbf{w}} \mathbf{e}(\mathbf{w}):$$


Applying SGD

All the weights $\mathbf{w} = \{\mathbf{w}_{ij}^{(l)}\}$ determine $h(\mathbf{x})$

Error on example (\mathbf{x}_n, y_n) is

$$\mathbf{e}(h(\mathbf{x}_n), y_n) = \mathbf{e}(\mathbf{w})$$

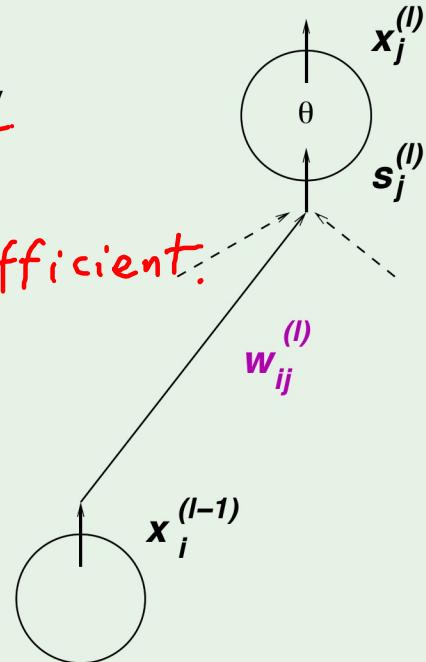
To implement SGD, we need the gradient

$$\nabla \mathbf{e}(\mathbf{w}): \frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}} \text{ for all } i, j, l$$

Computing $\frac{\partial \text{e}(\mathbf{w})}{\partial w_{ij}^{(l)}}$

We can evaluate $\frac{\partial \text{e}(\mathbf{w})}{\partial w_{ij}^{(l)}}$ one by one: analytically or numerically

not efficient.



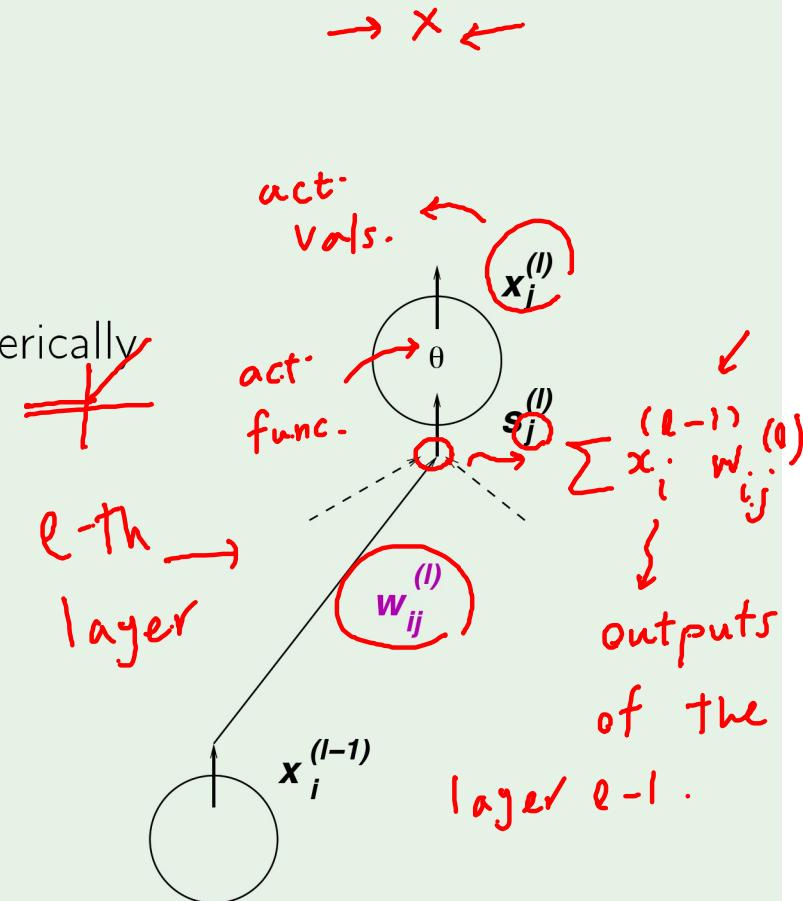
$$x_j^{(l)} = \theta(s_j^{(l)})$$

Computing $\frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}}$

We can evaluate $\frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}}$ one by one: analytically or numerically

A trick for efficient computation:

$$\frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}} = \underbrace{\frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}}}_{x_i^{(l-1)}} \times \underbrace{\frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}}}_{w_{ij}^{(l)}}$$



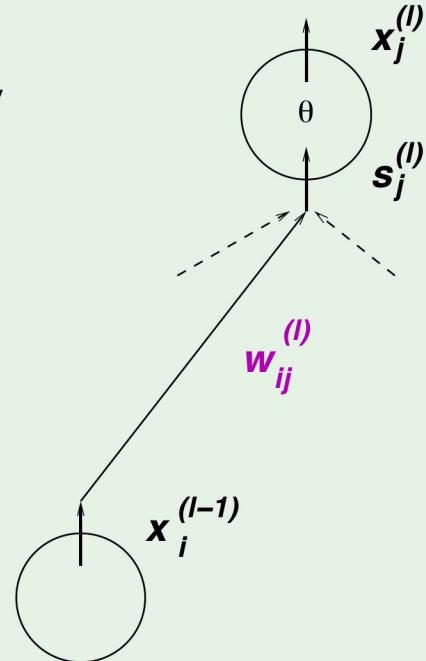
Computing $\frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}}$

We can evaluate $\frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}}$ one by one: analytically or numerically

A trick for efficient computation:

$$\frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}} = \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}}$$

We have $\frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}} = x_i^{(l-1)}$



Computing $\frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}}$

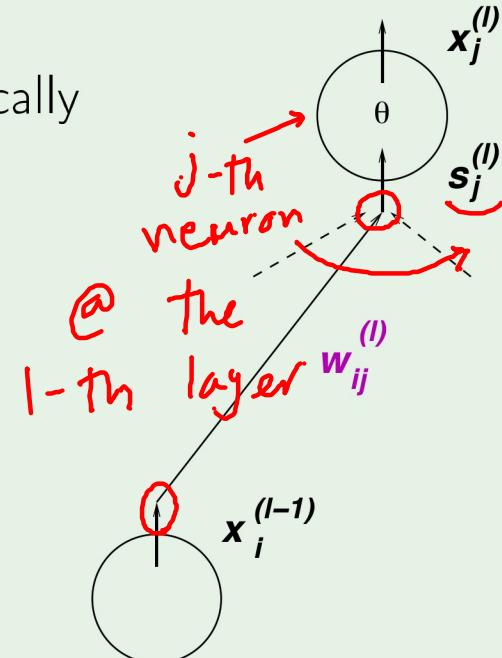
We can evaluate $\frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}}$ one by one: analytically or numerically

A trick for efficient computation:

$$\boxed{\frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}}} = \underbrace{\frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}}}_{\text{j-th neuron}} \times \underbrace{\frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}}}_{\text{@ the layer } w_{ij}^{(l)}}$$

$$\text{We have } \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}} = x_i^{(l-1)}$$

$$\text{We only need: } \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}} = \underline{\delta_j^{(l)}}$$



δ for the final layer

$$\delta_j^{(l)} = \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}}$$

For the final layer $\underline{l = L}$ and $\underline{j = 1}$:

δ for the final layer

$$\delta_j^{(l)} = \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}}$$

For the final layer $l = L$ and $j = 1$:

$$\delta_1^{(L)} = \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_1^{(L)}}$$

δ for the final layer

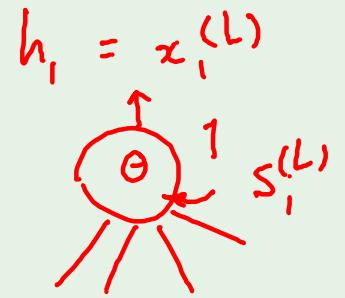
$$\delta_j^{(l)} = \frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}}$$

For the final layer $l = L$ and $j = 1$:

$$\underline{\delta_1^{(L)}} = \frac{\partial e(\mathbf{w})}{\partial s_1^{(L)}}$$

1 or -1

$$\rightarrow \underline{e(\mathbf{w})} = (\underline{x_1^{(L)}} - \widehat{y_n})^2$$
$$h_1 = \Theta(s_1^{(L)})$$



δ for the final layer

$$\delta_j^{(l)} = \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}}$$

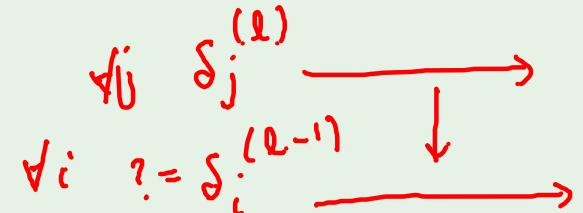
For the final layer $l = L$ and $j = 1$:

$$\delta_1^{(L)} = \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_1^{(L)}}$$

$$\mathbf{e}(\mathbf{w}) = (x_1^{(L)} - y_n)^2$$

$$x_1^{(L)} = \theta(s_1^{(L)})$$

δ for the final layer



$$\delta_j^{(l)} = \frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}}$$

For the final layer $l = L$ and $j = 1$:

$$\delta_1^{(L)} = \frac{\partial e(\mathbf{w})}{\partial s_1^{(L)}}$$

$$\frac{\partial e(\mathbf{w})}{\partial s_1^{(L)}} = 2(x_1^{(L)} - y_n) \cdot \underbrace{\theta'(s_1^{(L)})}_{1 - \theta^2(s_1^{(L)})}$$

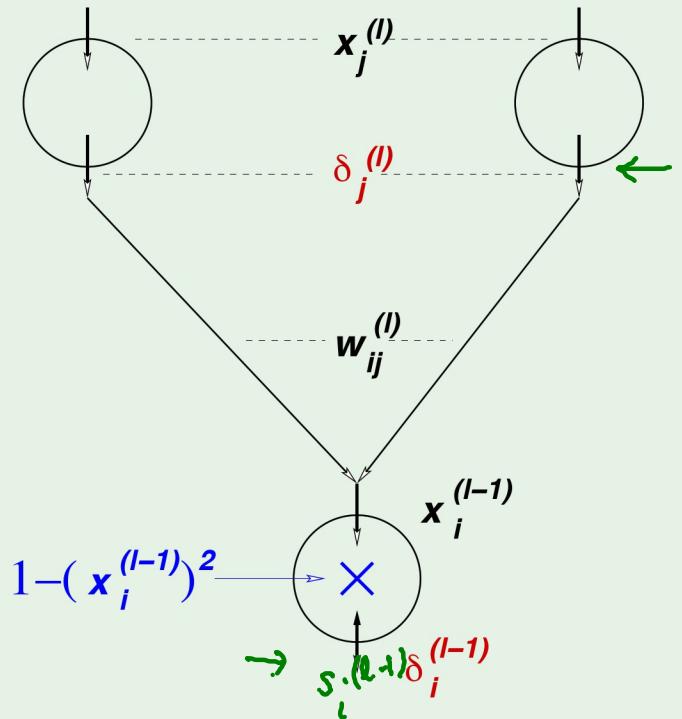
$$e(\mathbf{w}) = (\underbrace{x_1^{(L)} - y_n}_{\theta(s_1^{(L)})})^2$$

$$x_1^{(L)} = \theta(s_1^{(L)})$$

$$\theta'(s) = 1 - \theta^2(s) \quad \text{for the } \underline{\tanh}$$

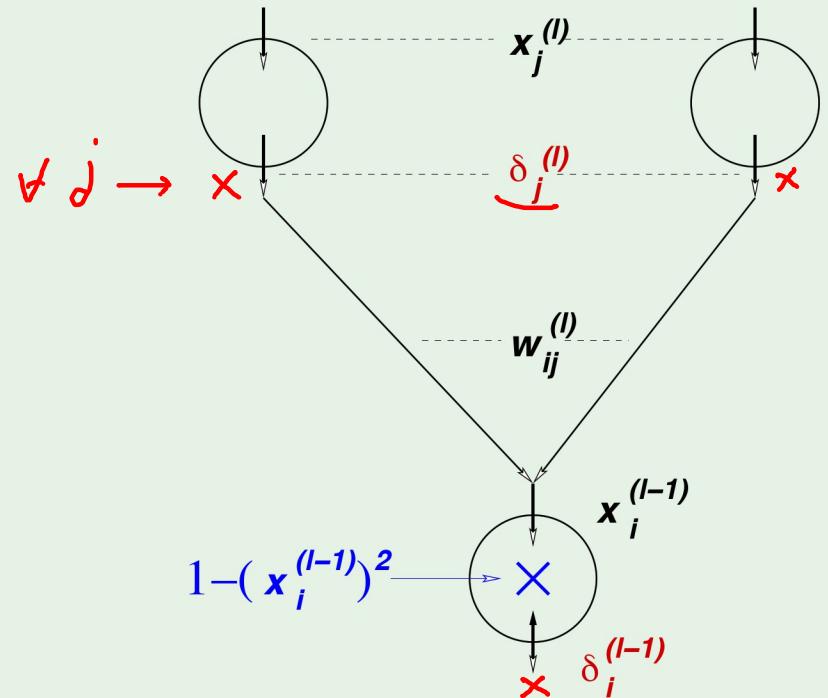
\curvearrowleft from $\delta_j^{(l)}$ Back propagation of δ

$$\underline{\delta_i^{(l-1)}} = \frac{\partial e(\mathbf{w})}{\partial s_i^{(l-1)}}$$



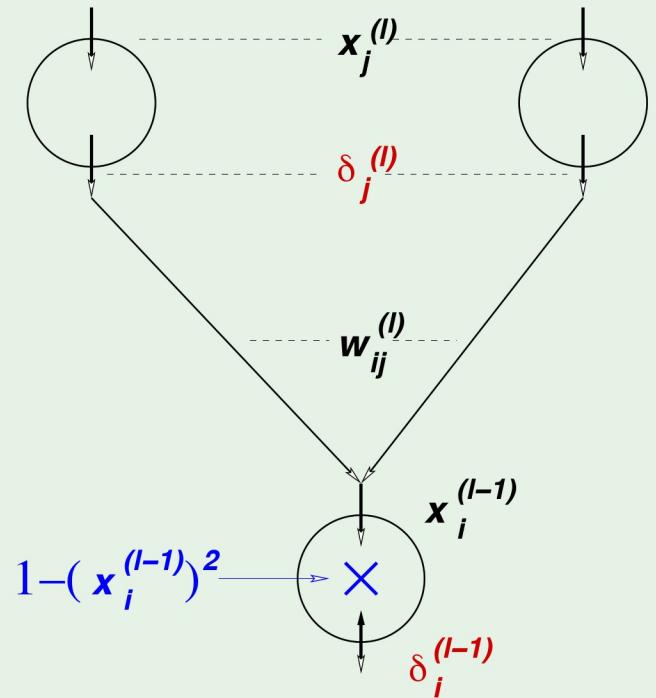
Back propagation of δ

$$\begin{aligned}\delta_i^{(l-1)} &= \frac{\partial e(\mathbf{w})}{\partial s_i^{(l-1)}} \\ &= \sum_j \frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial x_i^{(l-1)}} \times \frac{\partial x_i^{(l-1)}}{\partial s_i^{(l-1)}}\end{aligned}$$



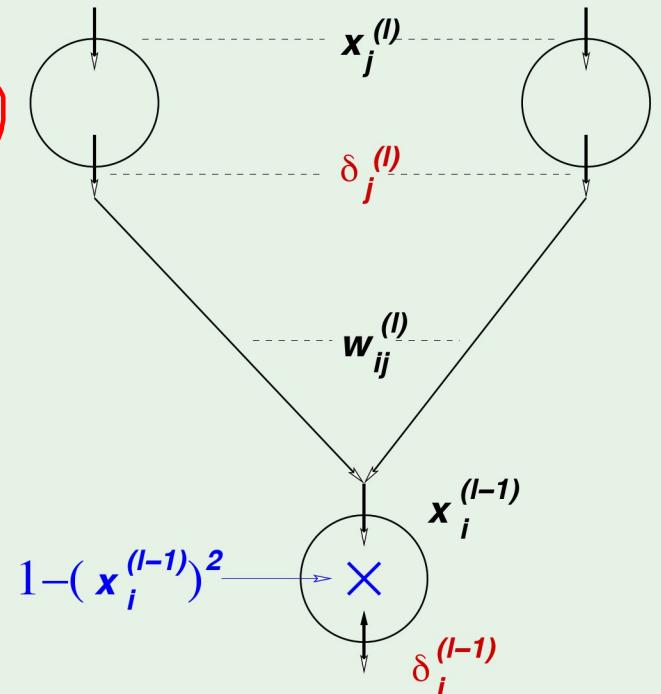
Back propagation of δ

$$\begin{aligned}\delta_i^{(l-1)} &= \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_i^{(l-1)}} \\ &= \sum_{j=1}^{d(l)} \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial x_i^{(l-1)}} \times \frac{\partial x_i^{(l-1)}}{\partial s_i^{(l-1)}}\end{aligned}$$



Back propagation of δ

$$\begin{aligned}
 \delta_i^{(l-1)} &= \frac{\partial e(\mathbf{w})}{\partial s_i^{(l-1)}} \\
 &= \sum_{j=1}^{d^{(l)}} \frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial x_i^{(l-1)}} \times \frac{\partial x_i^{(l-1)}}{\partial s_i^{(l-1)}} \\
 &= \sum_{j=1}^{d^{(l)}} \delta_j^{(l)} \times w_{ij}^{(l)} \times \theta'(s_i^{(l-1)})
 \end{aligned}$$



Vanishing Grad.

$$\delta_i^{(l-1)} = \frac{\partial e(\mathbf{w})}{\partial s_i^{(l-1)}}$$

$$= \sum_{j=1}^{d(l)} \frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial x_i^{(l-1)}} \times \frac{\partial x_i^{(l-1)}}{\partial s_i^{(l-1)}}$$

$$= \sum_{j=1}^{d(l)} \delta_j^{(l)} \times w_{ij}^{(l)} \times \theta'(s_i^{(l-1)})$$

① < 1

$$\delta_i^{(l-1)} = (1 - (x_i^{(l-1)})^2) \sum_{j=1}^{d(l)} w_{ij}^{(l)} \delta_j^{(l)}$$

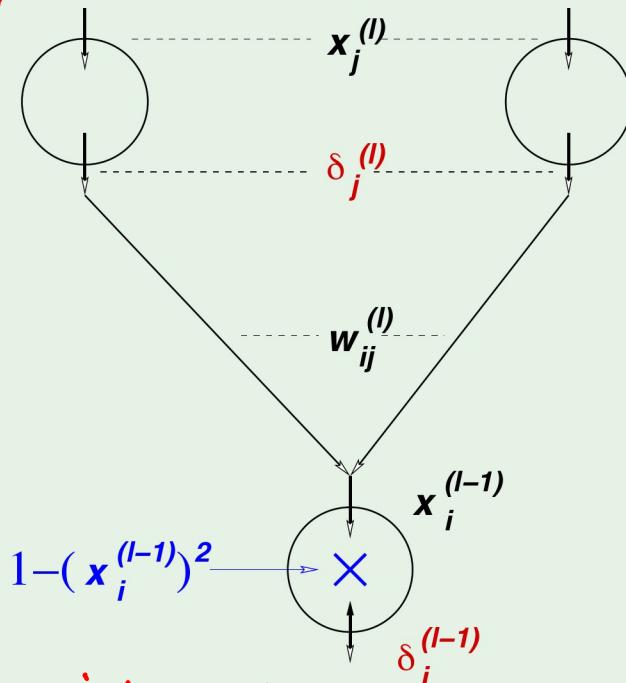
$\Theta'(s_i^{(l-1)})$ Scaling

Weighted Avg.

Back propagation of δ

$\delta_i^{(l)} \leftarrow \dots \leftarrow \delta_i^{(l-2)} \leftarrow \delta_i^{(l-1)}$

$\delta_i^{(L)}$



Weighted Avg.
Why? → Similar Forward & Backward

$$\text{Sig}(x) = \frac{1}{1 + e^{-x}}$$

$$\text{Sig}' = \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$\text{Sig}'(0) = \frac{1}{(1+1)^2} = \frac{1}{4}$$

1: Initialize all weights $w_{ij}^{(l)}$ at random

2: for $t = 0, 1, 2, \dots$ do

3: Pick $n \in \{1, 2, \dots, N\}$ \rightarrow for m in $1 \dots K$

4: Forward: Compute all $x_j^{(l)}$

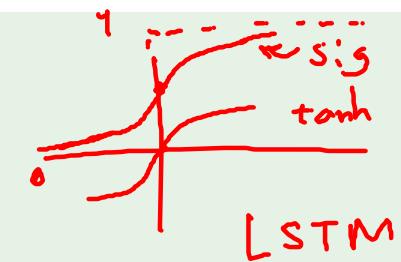
5: Backward: Compute all $\delta_j^{(l)}$

6: Update the weights: $w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta x_i^{(l-1)} \delta_j^{(l)}$

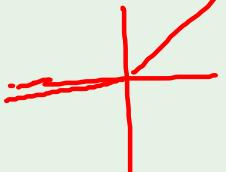
7: Iterate to the next step until it is time to stop

8: Return the final weights $w_{ij}^{(l)}$

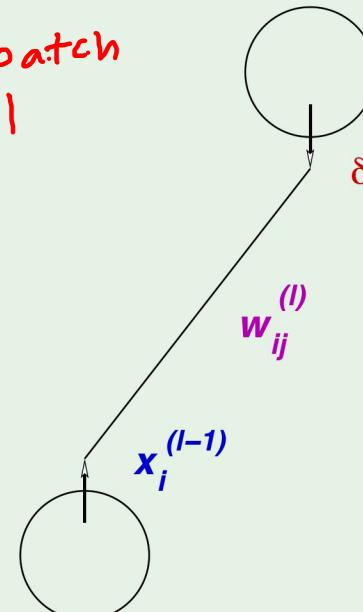
Backpropagation algorithm



ReLU



SGD
with
mini-batch
size = 1



$$\sum_{k \in M_m} x_i^{(l-1)} \delta_j^{(l)}$$

Final remark: hidden layers

learned nonlinear transform

interpretation?

