

Lecture - 20

Finding the  
Gradient of  $E_{in}$

Given weights  $W$ ,

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N e_n = \frac{1}{N} \sum_{n=1}^N e(h(n; w), y_n)$$

the gradient of  $E_{in}(w)$  is:

$\partial$

\* Backpropagation:-

computing the gradient: Using the chain rule.

(1) Node  $i$  only affects the output  $h(x)$  through its inputs  $s_i^{(l)}$ .

How error changes w.r. to  $w_i^{(l)}$

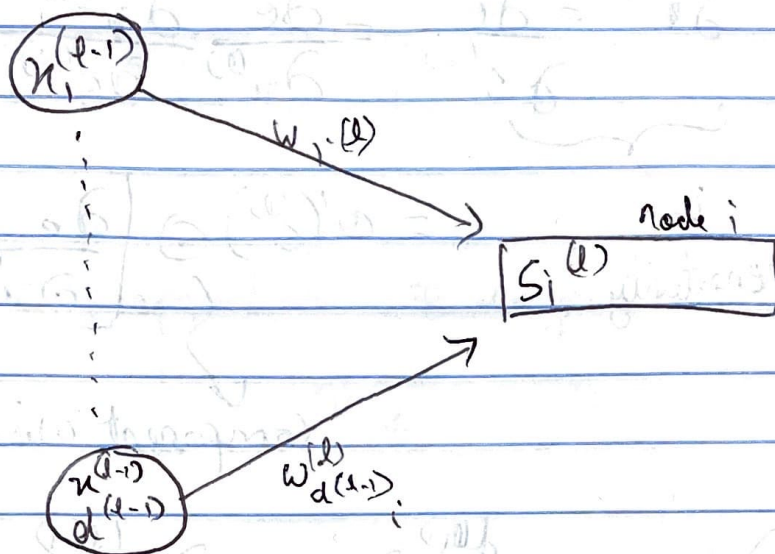
$$\frac{\partial e}{\partial w_i^{(l)}} = \frac{\partial s_i^{(l)}}{\partial w_i^{(l)}} \cdot \frac{\partial e}{\partial s_i^{(l)}}$$

How input signal changes w.r. to  $w_i^{(l)}$

(Chain rule)

Sensitivity

$$s_i^{(l)} = (w_i^{(l)})^T n^{(l-1)} \Rightarrow \frac{\partial s_i^{(l)}}{\partial w_i^{(l)}} = n^{(l-1)}$$



if  $l > 1$

$$n^{(l-1)} = \tanh \left( \begin{matrix} s^{(l-1)} \\ w^{(l-1)} n^{(l-2)} \\ d^{(l-1)} d^{(l-2)} (1 + d^{(l-3)}) \end{matrix} \right)$$

$$n^{(l-1)} = \begin{bmatrix} 0 \\ s^{(l-1)} \end{bmatrix} \rightarrow \tanh$$

$$\frac{\partial e}{\partial w_i^{(l)}} = n^{(l-1)} \cdot (s^{(l)})^T$$

Seriesing



\* Sensitivity: Using the chain rule Again:

$$\delta^{(l)} = \frac{\partial e}{\partial s^{(l)}} = \frac{\partial e}{\partial n^{(l)}} \frac{\partial n^{(l)}}{\partial s^{(l)}}$$

$\swarrow$   
 Sensitivity

$$= \sigma'(s^{(l)}) \odot \left[ \frac{\partial e}{\partial n^{(l)}} \right]^{(l)}$$

$\swarrow$   
 Component wise multiplication

$$\frac{\partial e}{\partial n^{(l)}} = \sum_{k=1}^{d^{(l+1)}} \underbrace{\frac{\partial e}{\partial s_k^{(l+1)}}}_{\delta_k^{(l+1)}} \cdot \underbrace{\frac{\partial s_k^{(l+1)}}{\partial n^{(l)}}}_{w_k^{(l+1)}}$$

\* We can now compute,

~~$\delta^{(l)} = ?$~~

$$\delta^{(l)} \leftarrow \delta^{(l+1)} \leftarrow \dots \delta^{L-1} \leftarrow \delta^L$$

$n^L$  is just the tanh applied to  $\frac{\partial n^L}{\partial s^{(L)}}$

$$\frac{\partial (u^{(L)} - y)^2}{\partial s^{(L)}}$$

$$= 2(u^{(L)} - y) \frac{\partial u^{(L)}}{\partial s^{(L)}}$$

\* Backpropagation to compute sensitivities.

$$\delta^{(1)} \leftarrow \dots \leftarrow \delta^{(L-2)} \leftarrow \delta^{(L-1)} \leftarrow \delta^{(L)}.$$

Using backpropagation we will calculate the sensitivities and then we will use these sensitivities to calculate gradient descent.

We will use Forward propagation for outputs  
Backward propagation for sensitivities