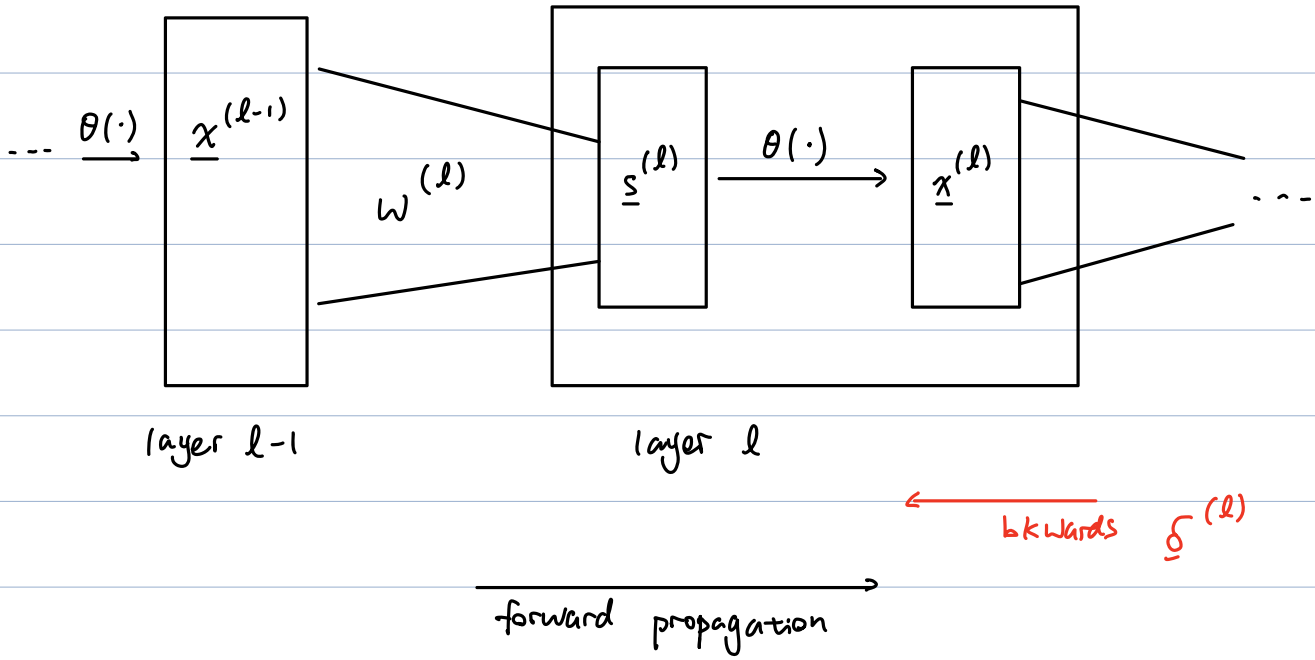


Lec 16 Backward propagation

1. Recall:



$$\hookrightarrow \Omega = \{ \omega^{(1)}, \omega^{(2)}, \dots, \omega^{(L)} \}$$

↳ Loss function

$$e_n(\Omega) \stackrel{\text{e.g.}}{=} g(x^{(L)}, y_n)$$

↳ Goal:

compute $\frac{\partial \text{en}(\Omega)}{\partial w_{ij}(\ell)} \quad \forall i, j, \ell$

But difficult \implies sol: backward propagation

2. Backward propagation Algo: (1960's)

↳ (1986, Hinton) , similar to dynamic programming

we'll use $e(\Omega)$ instead of $e_n(\Omega)$ here since same analysis for any n .

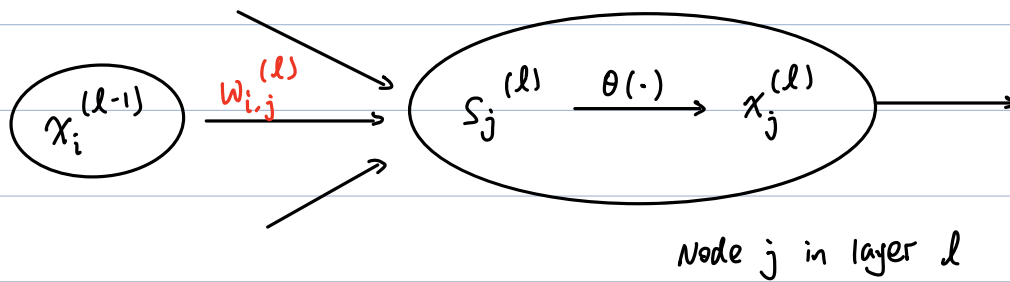
↳. Note:

① Memoryless: \rightarrow store all info up to latest t .

$$e(\mathcal{R}) = e(s^{(l)}, w^{(l+1)}, \dots, w^{(L)})$$

$$= e((s_i^{(l)}, s_j^{(l)}, \dots, s_d^{(l)}), w^{(l+1)}, \dots, w^{(L)})$$

②



• Only $S_j^{(l)}$ depends on $w_{i,j}^{(l)}$

$$\Rightarrow \frac{\partial e(\mathcal{N})}{\partial w_{i,j}^{(l)}} = \frac{\partial e(\mathcal{N})}{\partial S_j^{(l)}} \cdot \frac{\partial S_j^{(l)}}{\partial w_{i,j}^{(l)}}$$

$$\because S_j^{(l)} = w_{0,j}^{(l)} + \sum_{k=1}^{d^{(l-1)}} w_{k,j}^{(l)} x_k^{(l-1)}$$

$$\therefore \frac{\partial S_j^{(l)}}{\partial w_{i,j}^{(l)}} = x_i^{(l-1)}$$

Let $\delta_j^{(l)} \triangleq \frac{\partial e(\mathcal{N})}{\partial S_j^{(l)}}$ (sensitivity of $e(\mathcal{N})$ to layer l input)

$$\therefore \frac{\partial e(\mathcal{N})}{\partial S_j^{(l)}} = x_i^{(l-1)} \cdot \delta_j^{(l)}$$

Define $\frac{\partial e(\mathcal{N})}{\partial \mathbf{w}^{(l)}} \triangleq \left[\frac{\partial e(\mathcal{N})}{\partial w_{i,j}^{(l)}} \right]_{\substack{0 \leq i \leq d^{(l-1)} \\ 1 \leq j \leq d^{(l)}}}$

$$\therefore \underline{x}^{(l-1)} = \begin{bmatrix} x_0^{(l-1)} \\ x_1^{(l-1)} \\ \vdots \\ x_{d^{(l-1)}}^{(l-1)} \end{bmatrix}, \quad \underline{\delta}^{(l)} = \begin{bmatrix} \delta_0^{(l)} \\ \vdots \\ \delta_{d^{(l-1)}}^{(l-1)} \end{bmatrix}$$

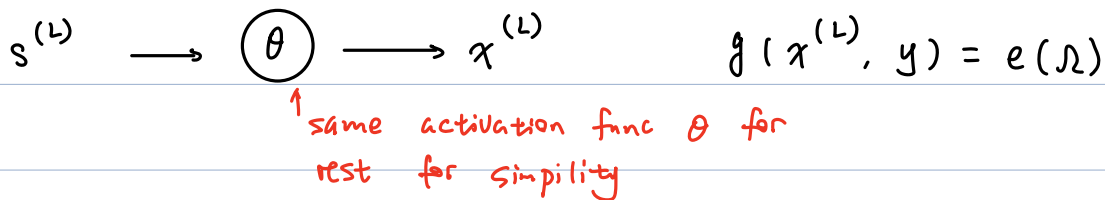
$$\therefore \frac{\partial e(\mathcal{N})}{\partial \mathbf{w}^{(l)}} = \underline{x}^{(l-1)} \cdot \underline{\delta}^{(l)T}$$

\uparrow col \uparrow row

2. Compute $\delta^{(l)}$ "backwards":

① consider $l = L$

Let $d^{(L)} = 1$ for simplicity.



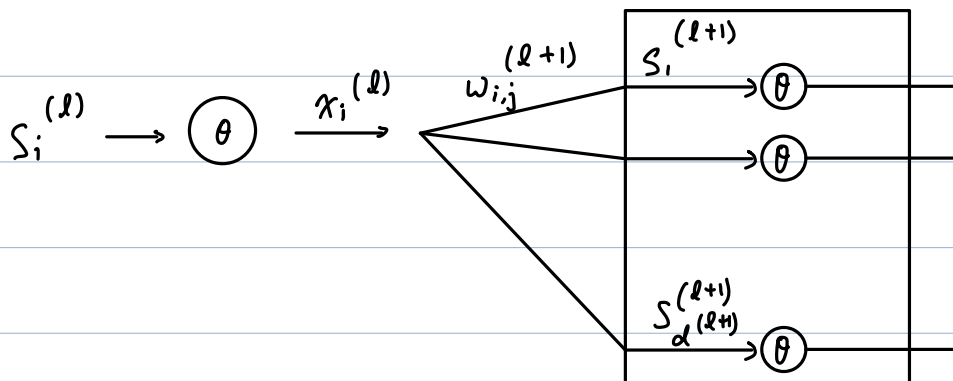
$$\begin{aligned}
 \therefore \delta^{(L)} &= \frac{\partial e(\Omega)}{\partial s^{(L)}} = \frac{\partial g(x^{(L)}, y)}{\partial s^{(L)}} \\
 &= \frac{\partial g(x^{(L)}, y)}{\partial x^{(L)}} \cdot \frac{\partial x^{(L)}}{\partial s^{(L)}} \\
 &= \frac{\partial g(x^{(L)}, y)}{\partial x^{(L)}} \cdot \theta'(s^{(L)})
 \end{aligned}$$

e.g. Regression:

$$g(x^{(L)}, y) = (x^{(L)} - y)^2$$

$$\delta^{(L)} = 2(x^{(L)} - y) \theta'(s^{(L)})$$

② Intermediate node: $\delta_i^{(l)}$



$$\delta_i^{(l)} = \frac{\partial e(\Omega)}{\partial s_i^{(l)}} = \frac{\partial e(\Omega)}{\partial x_i^{(l)}} \cdot \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}} = \frac{\partial e(\Omega)}{\partial x_i^{(l)}} \cdot \theta'(s_i^{(l)})$$

// Recall the chain rule:

$$\text{if } z = f(x, y), \text{ and } \begin{cases} x = g(s, t) \\ y = h(s, t) \end{cases}$$

$$\text{then } \frac{\partial z}{\partial s} = \frac{\partial z}{\partial x} \cdot \frac{\partial x}{\partial s} + \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial s}$$

$$\frac{\partial z}{\partial t} = \dots$$

$$\begin{aligned} \therefore \frac{\partial e(\Omega)}{\partial x_i^{(l)}} &= \sum_{j=1}^{d^{(l+1)}} \frac{\partial e(\Omega)}{\partial s_j^{(l+1)}} \cdot \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} \\ &= \sum_{j=1}^{d^{(l+1)}} \delta_j^{(l+1)} \cdot w_{i,j}^{(l+1)} \end{aligned}$$

$$\therefore \delta_i^{(l)} = \left[\sum_{j=1}^{d^{(l+1)}} \delta_j^{(l+1)} w_{i,j}^{(l+1)} \right] \cdot \theta'(s_i^{(l)})$$

Pack every i together:

$$\therefore \underline{\delta}^{(l)} = \begin{bmatrix} \delta_1^{(l)} \\ \vdots \\ \delta_{d^{(l)}}^{(l)} \end{bmatrix} \cdot \theta'(\underline{s}^{(l)}) = \begin{bmatrix} \theta'(s_1^{(l)}) \\ \vdots \\ \theta'(s_{d^{(l)}}^{(l)}) \end{bmatrix}$$

$$\hat{W}^{(l+1)} \triangleq \begin{bmatrix} w_{1,1}^{(l+1)} & \dots & w_{1,d^{(l+1)}}^{(l+1)} \\ \vdots & & \vdots \\ w_{d^{(l)},1}^{(l+1)} & \dots & w_{d^{(l)},d^{(l+1)}}^{(l+1)} \end{bmatrix}$$

↑
hat since
remove bias

$$\therefore \underline{\delta}^{(l)} = \left[\hat{W}^{(l+1)} \underline{\delta}^{(l+1)} \right] \otimes \theta'(\underline{s}^{(l)})$$



pointwise
multiplication.