



# IIC 2433 Minería de Datos

<https://github.com/marcelomendoza/IIC2433>

- GRADIENTE DESCENDENTE -

## Regresión logística

Si el ajuste es adecuado:

$$P(y \mid \mathbf{x}) = \theta(y \cdot \mathbf{w}^T \mathbf{x})$$

Luego, podemos expresar la función de verosimilitud:

$$P(y_1, \dots, y_N \mid \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{n=1}^N P(y_n \mid \mathbf{x}_n).$$

# Regresión logística

Maximizamos la función de verosimilitud:

$$\max \quad \prod_{n=1}^N P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \max \quad \ln \left( \prod_{n=1}^N P(y_n \mid \mathbf{x}_n) \right)$$

$$\equiv \max \quad \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

## Regresión logística

Maximizamos la función de verosimilitud:

$$\max \quad \prod_{n=1}^N P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \max \quad \ln \left( \prod_{n=1}^N P(y_n \mid \mathbf{x}_n) \right)$$

$$\equiv \max \quad \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \text{min} \quad - \frac{1}{N} \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\equiv \min \quad \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{P(y_n \mid \mathbf{x}_n)}$$

$$\equiv \min \quad \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{\theta(y_n \cdot \mathbf{w}^T \mathbf{x}_n)}$$

# Regresión logística

Maximizamos la función de verosimilitud:

$$\max \prod_{n=1}^N P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \max \ln \left( \prod_{n=1}^N P(y_n \mid \mathbf{x}_n) \right)$$

$$\equiv \max \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \text{min} - \frac{1}{N} \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\equiv \min \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{P(y_n \mid \mathbf{x}_n)}$$

$$\equiv \min \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{\theta(y_n \cdot \mathbf{w}^T \mathbf{x}_n)}$$

$$\equiv \min \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \cdot \mathbf{w}^T \mathbf{x}_n})$$

$$\theta(s) = \frac{1}{1 + e^{-s}}.$$



# Regresión logística

Tenemos una expresión para:

Parámetros del modelo

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \cdot \mathbf{w}^T \mathbf{x}_n}) \quad \text{Cross-entropy}$$

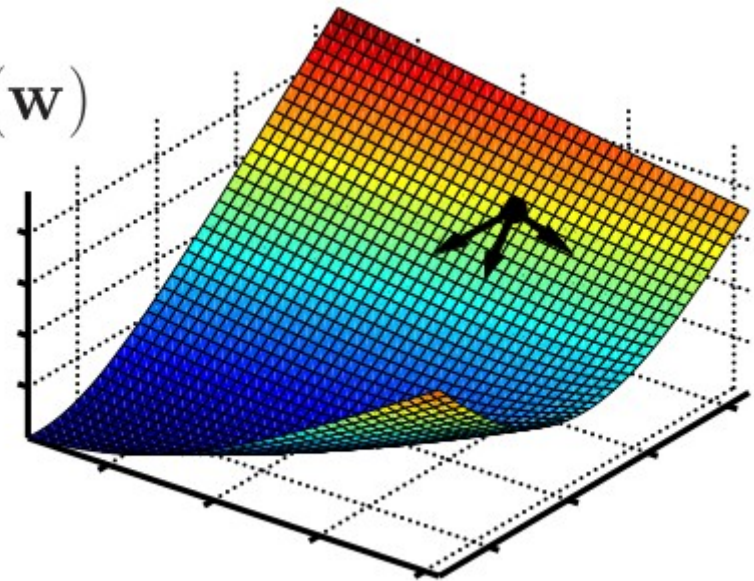
La función es convexa, por lo que podemos optimizarla de forma iterativa:

Idea del gradiente descendente:

$E_{\text{in}}(\mathbf{w})$

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \hat{\mathbf{v}}$$

¿En cuál dirección conviene moverse?



## Gradiente descendente

$$\begin{aligned}\Delta E_{\text{in}} &= E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t)) \\ &= E_{\text{in}}(\mathbf{w}(t) + \eta \hat{\mathbf{v}}) - E_{\text{in}}(\mathbf{w}(t)) \\ &= \eta \nabla E_{\text{in}}(\mathbf{w}(t))^{\text{T}} \hat{\mathbf{v}} + O(\eta^2)\end{aligned}$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(k+1)}(c)}{(k+1)!}(x - x_0)^{k+1}.$$

$$\begin{aligned}\Delta E_{\text{in}} &= E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t)) \\ &\quad \underbrace{f(x_0 + \eta \cdot x) - f(x_0)}_{\text{(expansión de Taylor de 1er orden)}} \\ &= f(x_0) + \nabla f(x_0) \cdot (x_0 + \eta \cdot x - x_0) + \dots\end{aligned}$$



## Gradiente descendente

$$\begin{aligned}
 \Delta E_{\text{in}} &= E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t)) \\
 &= E_{\text{in}}(\mathbf{w}(t) + \eta \hat{\mathbf{v}}) - E_{\text{in}}(\mathbf{w}(t)) \\
 &= \eta \nabla E_{\text{in}}(\mathbf{w}(t))^{\text{T}} \hat{\mathbf{v}} + O(\eta^2)
 \end{aligned}$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(k+1)}(c)}{(k+1)!}(x - x_0)^{k+1}.$$

$$\Delta E_{\text{in}} = E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t))$$

$$\underbrace{f(x_0 + \eta \cdot x) - f(x_0)}$$

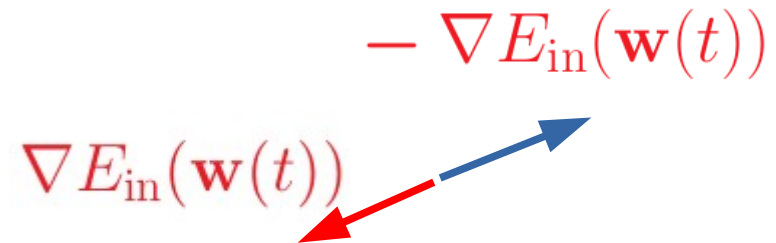
(expansión de Taylor de 1<sup>er</sup> orden)

$$\cancel{f(x_0)} + \nabla f(x_0) \cdot (\cancel{x_0} + \eta \cdot x - \cancel{x_0}) + \dots$$

## Gradiente descendente

$$= \eta \underbrace{\nabla E_{\text{in}}(\mathbf{w}(t))^T}_{\text{Minimizado en}} \hat{\mathbf{v}} + O(\eta^2)$$

$$\text{Minimizado en } \hat{\mathbf{v}} = - \frac{\nabla E_{\text{in}}(\mathbf{w}(t))}{\|\nabla E_{\text{in}}(\mathbf{w}(t))\|}$$



## Gradiente descendente

$$= \eta \underbrace{\nabla E_{\text{in}}(\mathbf{w}(t))^T \hat{\mathbf{v}}}_{\text{Minimizado en } \hat{\mathbf{v}} = -\frac{\nabla E_{\text{in}}(\mathbf{w}(t))}{\|\nabla E_{\text{in}}(\mathbf{w}(t))\|}} + O(\eta^2)$$

$$\nabla E_{\text{in}}(\mathbf{w}(t)) \quad -\nabla E_{\text{in}}(\mathbf{w}(t))$$

Gradiente descendente

1: Initialize at step  $t = 0$  to  $\mathbf{w}(0)$ .

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

3:   Compute the gradient

$$\mathbf{g}_t = \nabla E_{\text{in}}(\mathbf{w}(t)).$$

4:   Move in the direction  $\mathbf{v}_t = -\mathbf{g}_t$ .

5:   Update the weights:

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \mathbf{v}_t.$$

6:   Iterate ‘until it is time to stop’.

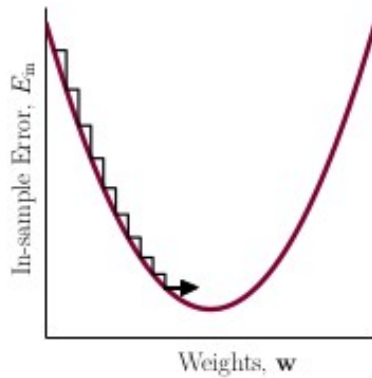
7: **end for**

8: Return the final weights.

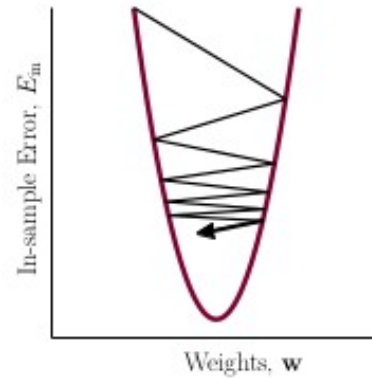
# Gradiente descendente

El efecto del **learning rate**:

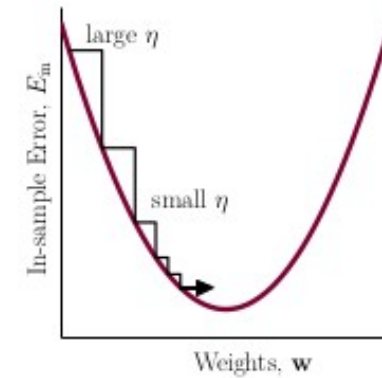
$\eta$  muy pequeño



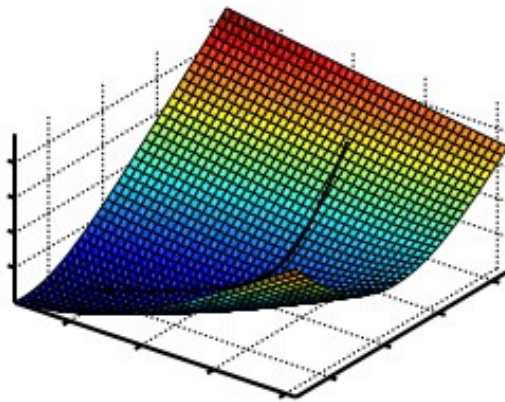
$\eta$  muy grande



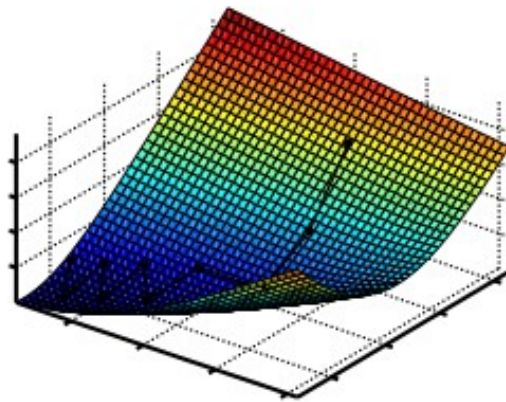
$\eta_t$



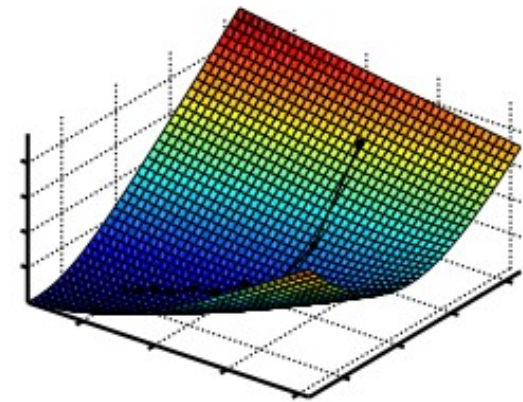
Adam



$\eta = 0.1$ ; 75 steps



$\eta = 2$ ; 10 steps



variable  $\eta_t$ ; 10 steps

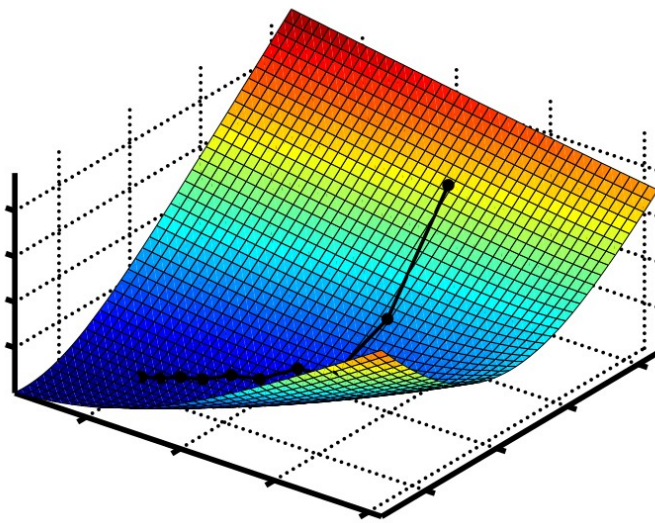
# Gradiente descendente

Una variación de gradiente descendente considera la evaluación del gradiente en una muestra del training set.

$$\mathbf{w}(t+1) \leftarrow \mathbf{w}(t) - \eta \nabla_{\mathbf{w}} e(\mathbf{w}, \mathbf{x}_*, y_*)$$

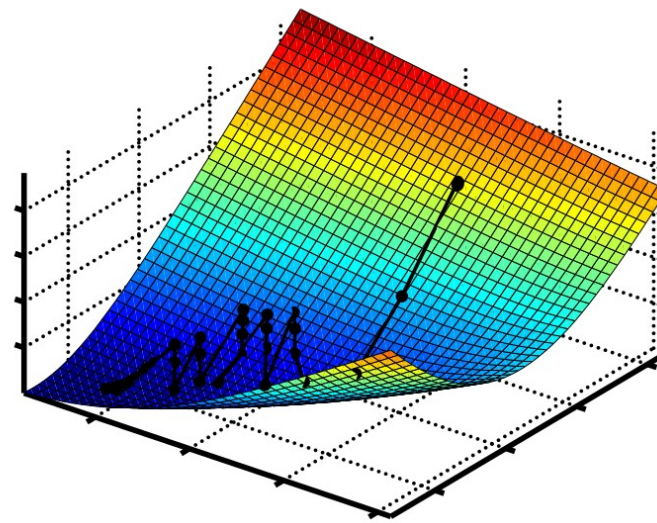
Dado que la muestra se toma al azar, se le denomina gradiente descendente estocástico.

GD



$\eta = 6$   
10 steps  
 $N = 10$

SGD



$\eta = 2$   
30 steps

Puede ayudar a escapar  
de óptimos locales

## Regresión logística + gradiente descendente

Para regresión logística encontramos que:

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \cdot \mathbf{w}^T \mathbf{x}_n})$$

Muestre que:

$$\begin{aligned} \nabla E_{\text{in}}(\mathbf{w}) &= -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T \mathbf{x}_n}} \\ &= \frac{1}{N} \sum_{n=1}^N -y_n \mathbf{x}_n \theta(-y_n \mathbf{w}^T \mathbf{x}_n). \end{aligned}$$

## Regresión logística + gradiente descendente

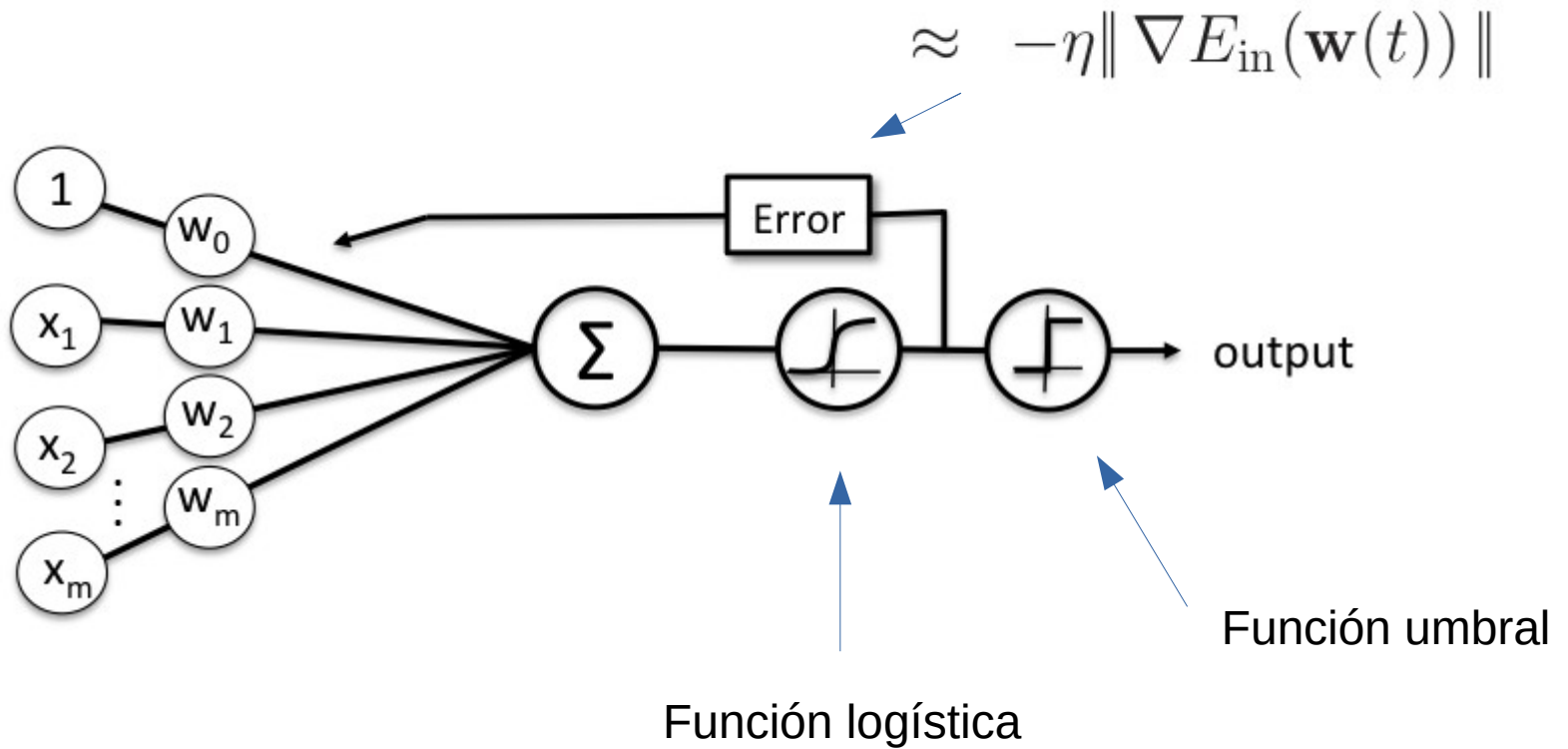
### Logistic regression algorithm:

- 1: Initialize the weights at time step  $t = 0$  to  $\mathbf{w}(0)$ .
- 2: **for**  $t = 0, 1, 2, \dots$  **do**
- 3:     Compute the gradient

$$\mathbf{g}_t = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T(t) \mathbf{x}_n}}.$$

- 4:     Set the direction to move,  $\mathbf{v}_t = -\mathbf{g}_t$ .
- 5:     Update the weights:  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \mathbf{v}_t$ .
- 6:     Iterate to the next step until it is time to stop.
- 7: Return the final weights  $\mathbf{w}$ .

## Regresión logística + gradiente descendente





- BACKPROPAGATION -

# El perceptrón multicapa (Multi-Layer Perceptron)

## Forward propagation

$$\mathbf{x} = \mathbf{x}^{(0)} \xrightarrow{W^{(1)}} \mathbf{s}^{(1)} \xrightarrow{\theta} \mathbf{x}^{(1)} \xrightarrow{W^{(2)}} \mathbf{s}^{(2)} \xrightarrow{\theta} \mathbf{x}^{(2)} \dots \xrightarrow{W^{(L)}} \mathbf{s}^{(L)} \xrightarrow{\theta} \mathbf{x}^{(L)} = h(\mathbf{x}).$$

### Forward propagation

```
1:  $\mathbf{x}^{(0)} \leftarrow \mathbf{x}$ 
2: for  $\ell = 1$  to  $L$  do
3:    $\mathbf{s}^{(\ell)} \leftarrow (W^{(\ell)})^T \mathbf{x}^{(\ell-1)}$ 
4:    $\mathbf{x}^{(\ell)} \leftarrow \begin{bmatrix} 1 \\ \theta(\mathbf{s}^{(\ell)}) \end{bmatrix}$ 
5: end for
6:  $h(\mathbf{x}) = \mathbf{x}^{(L)}$ 
```

Objetivo supervisado:

$$E_{\text{in}}(h) = E_{\text{in}}(W) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$$

Dado que  $\theta = \tanh$ ,  $E_{\text{in}}$  es diferenciable usando GD sobre

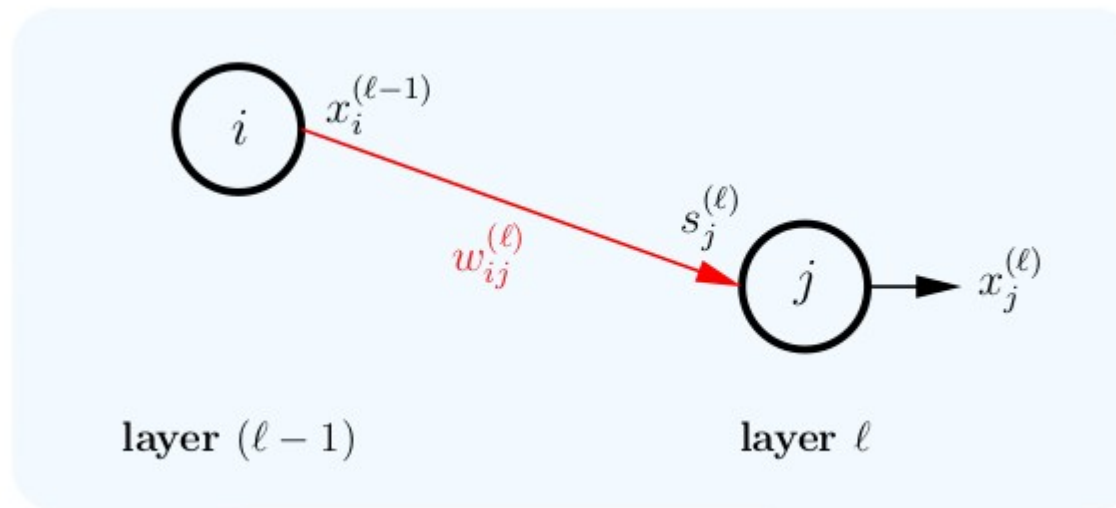
$$W = \{W^{(1)}, W^{(2)}, \dots, W^{(L)}\}$$

Parámetros del modelo

# Feed-Forward

$$\mathbf{x} = \mathbf{x}^{(0)} \xrightarrow{w^{(1)}} \mathbf{s}^{(1)} \xrightarrow{\theta} \mathbf{x}^{(1)} \xrightarrow{w^{(2)}} \mathbf{s}^{(2)} \xrightarrow{\theta} \mathbf{x}^{(2)} \dots \xrightarrow{w^{(L)}} \mathbf{s}^{(L)} \xrightarrow{\theta} \mathbf{x}^{(L)} = h(\mathbf{x}).$$

Nodo a nodo:

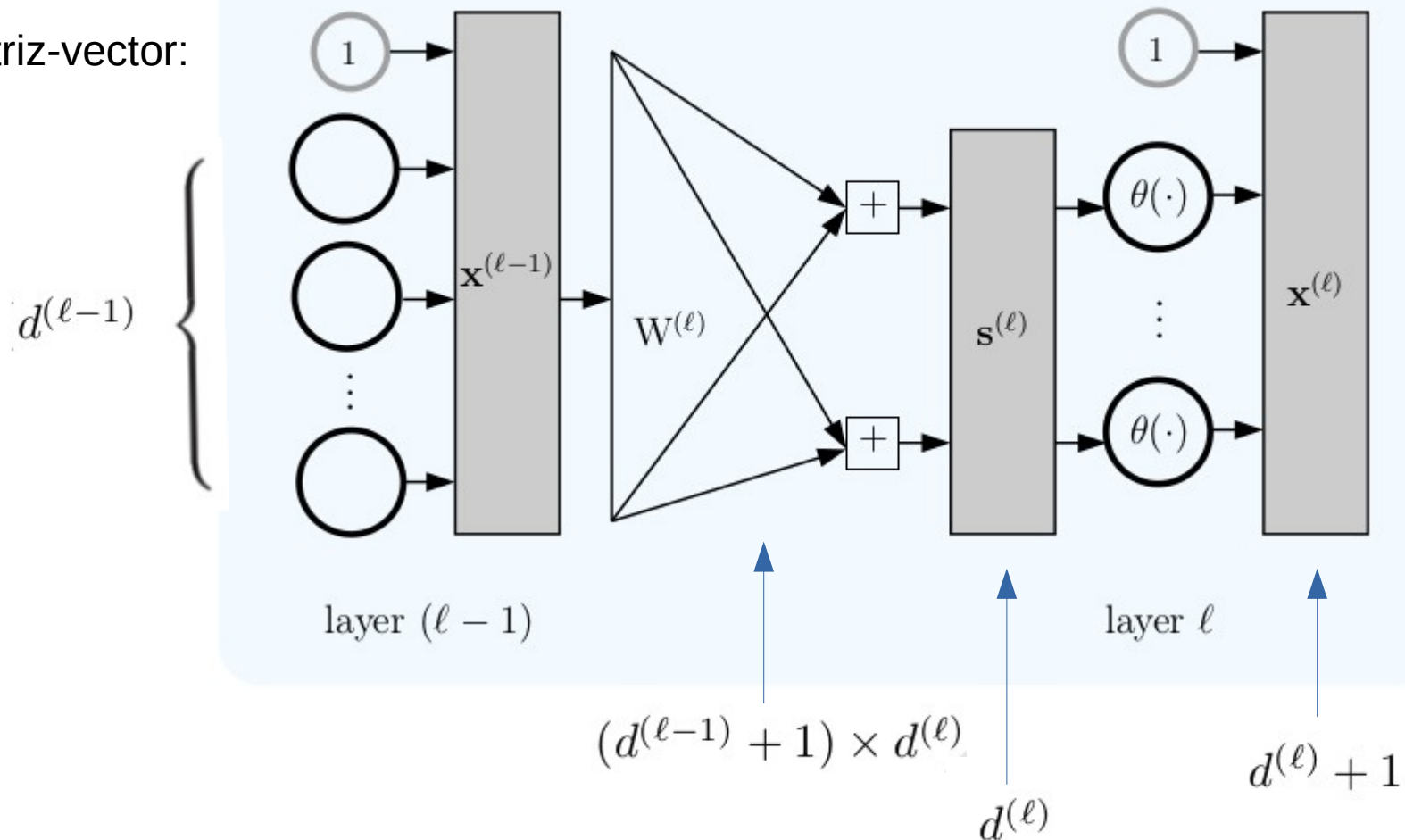


# Feed-Forward

$$\mathbf{x} = \mathbf{x}^{(0)} \xrightarrow{W^{(1)}} \mathbf{s}^{(1)} \xrightarrow{\theta} \mathbf{x}^{(1)} \xrightarrow{W^{(2)}} \mathbf{s}^{(2)} \xrightarrow{\theta} \mathbf{x}^{(2)} \dots \xrightarrow{W^{(L)}} \mathbf{s}^{(L)} \xrightarrow{\theta} \mathbf{x}^{(L)} = h(\mathbf{x}).$$

$$\mathbf{s}^{(\ell)} \leftarrow (W^{(\ell)})^T \mathbf{x}^{(\ell-1)}$$

Matrix-vector:



# Backpropagation

Minimizar:  $E_{\text{in}}(h) = E_{\text{in}}(W) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$

Podemos usar la idea de gradiente descendente:

$$W(t+1) = W(t) - \eta \nabla E_{\text{in}}(W(t))$$

Dado que:  $E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \overset{e_n}{\mathbf{e}}(h(\mathbf{x}_n), y_n)$

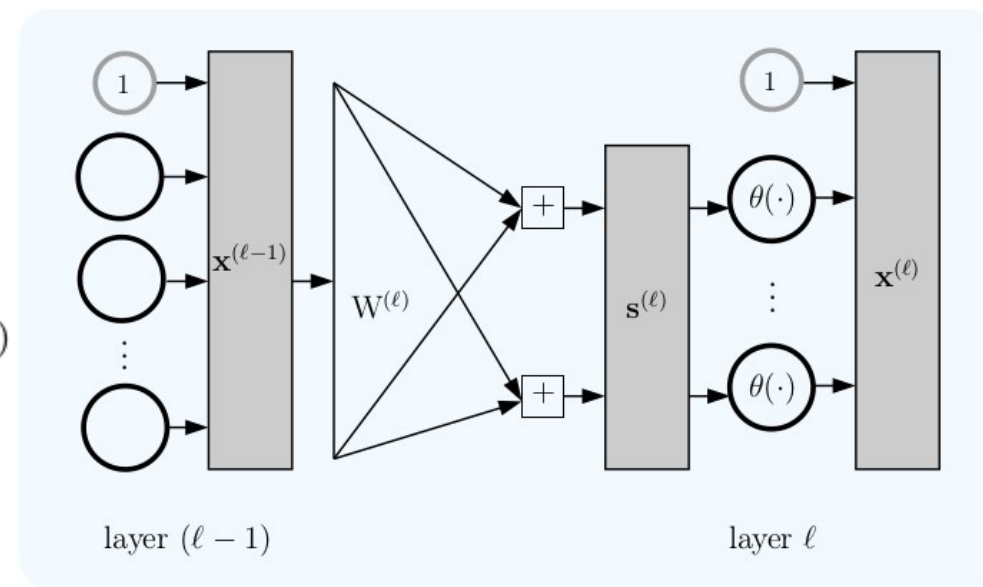
$$\frac{\partial E_{\text{in}}(\mathbf{w})}{\partial W^{(\ell)}} = \frac{1}{N} \sum_{n=1}^N \frac{\partial \mathbf{e}_n}{\partial W^{(\ell)}}$$

Necesitamos:

$$\frac{\partial \mathbf{e}(\mathbf{x})}{\partial W^{(\ell)}}$$

# Backpropagation

Vamos a usar la **regla de la cadena** para expresar las derivadas parciales de la capa  $(\ell-1)$  en función de las derivadas parciales de la capa  $\ell$ .



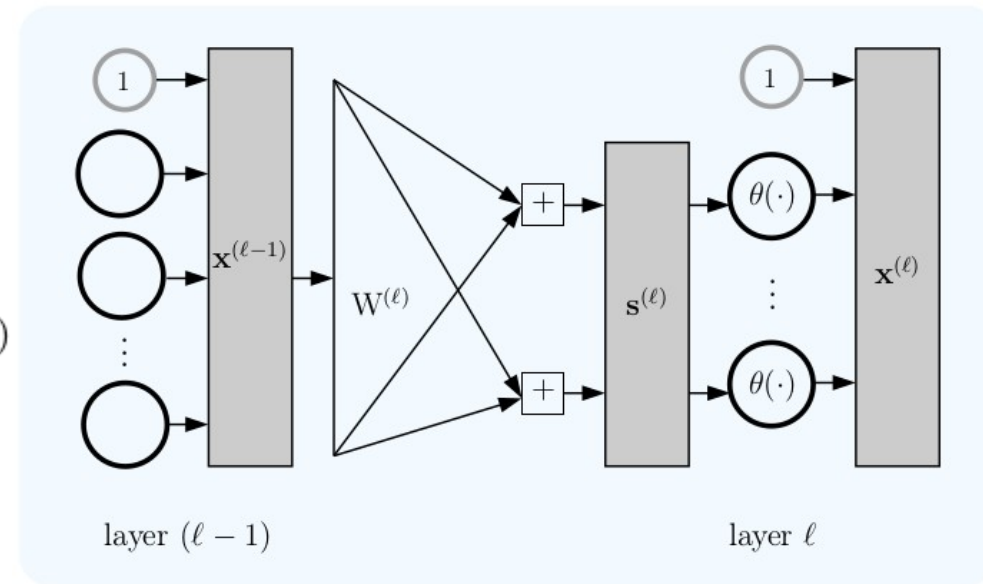
# Backpropagation

Vamos a usar la **regla de la cadena** para expresar las derivadas parciales de la capa  $(\ell-1)$  en función de las derivadas parciales de la capa  $\ell$ .

Tenemos:  $\mathbf{s}^{(\ell)} = (\mathbf{W}^{(\ell)})^T \mathbf{x}^{(\ell-1)}$

Definimos la sensibilidad de la capa  $\ell$ :

$$\boldsymbol{\delta}^{(\ell)} = \frac{\partial e}{\partial \mathbf{s}^{(\ell)}}$$



# Backpropagation

Vamos a usar la **regla de la cadena** para expresar las derivadas parciales de la capa  $(\ell-1)$  en función de las derivadas parciales de la capa  $\ell$ .

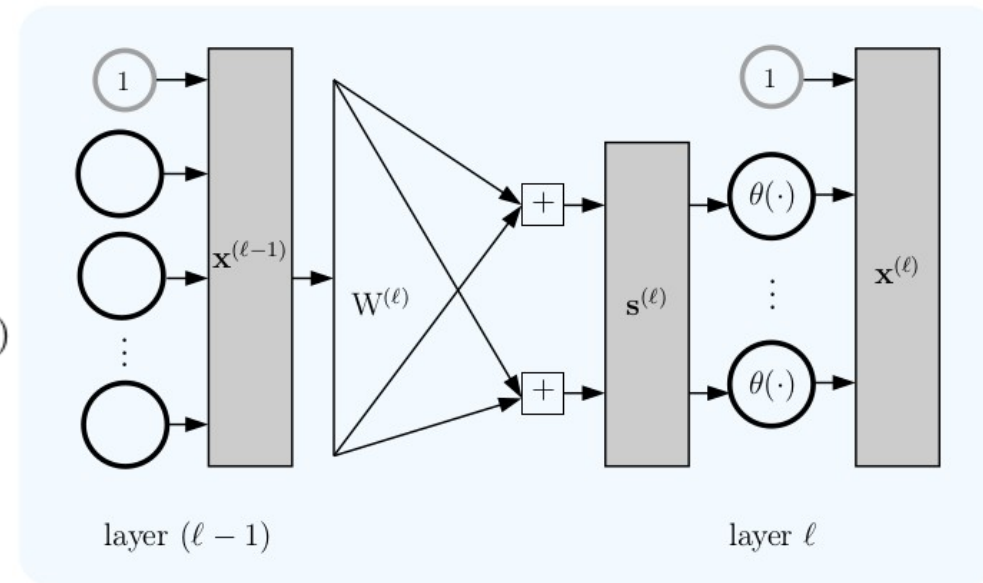
Tenemos:  $\mathbf{s}^{(\ell)} = (\mathbf{W}^{(\ell)})^T \mathbf{x}^{(\ell-1)}$

Definimos la sensibilidad de la capa  $\ell$ :

$$\boldsymbol{\delta}^{(\ell)} = \frac{\partial e}{\partial \mathbf{s}^{(\ell)}}$$

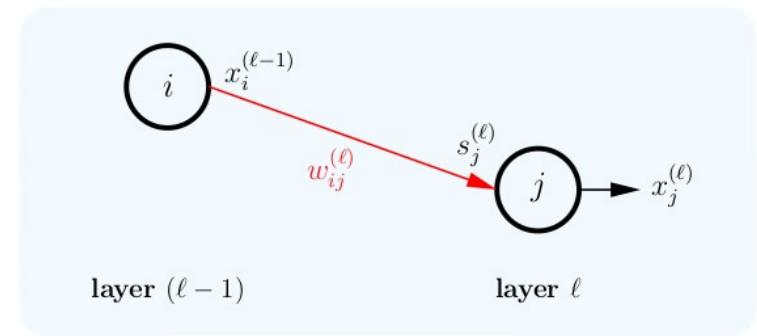
Aplicando la regla de la cadena:

$$\begin{aligned} \frac{\partial e}{\partial \mathbf{W}^{(\ell)}} &= \frac{\partial \mathbf{s}^{(\ell)}}{\partial \mathbf{W}^{(\ell)}} \cdot \left( \frac{\partial e}{\partial \mathbf{s}^{(\ell)}} \right)^T \\ &= \mathbf{x}^{(\ell-1)} (\boldsymbol{\delta}^{(\ell)})^T \end{aligned}$$





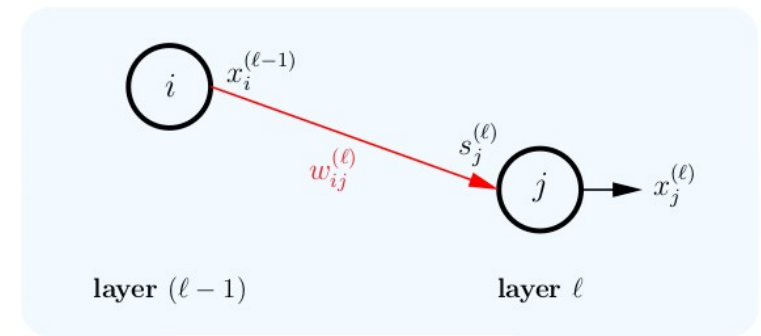
# Backpropagation



Miremos esto:

$$\frac{\partial e}{\partial W^{(\ell)}} = \frac{\partial s^{(\ell)}}{\partial W^{(\ell)}} \cdot \left( \frac{\partial e}{\partial s^{(\ell)}} \right)^T \quad \text{a nivel de un enlace.}$$
$$= \mathbf{x}^{(\ell-1)} (\boldsymbol{\delta}^{(\ell)})^T$$

# Backpropagation



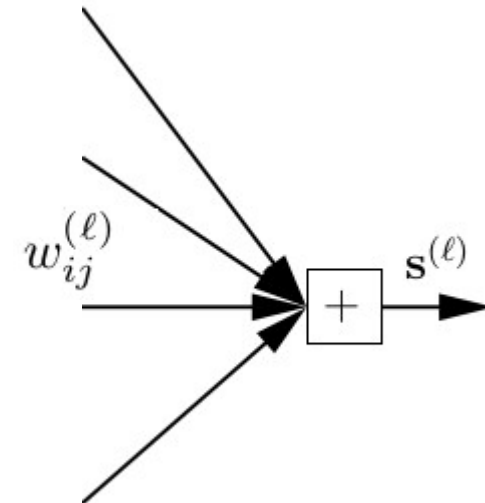
Miremos esto:  $\frac{\partial e}{\partial W^{(\ell)}} = \frac{\partial s^{(\ell)}}{\partial W^{(\ell)}} \cdot \left( \frac{\partial e}{\partial s^{(\ell)}} \right)^T$  a nivel de un enlace.

$$= \mathbf{x}^{(\ell-1)} (\boldsymbol{\delta}^{(\ell)})^T$$

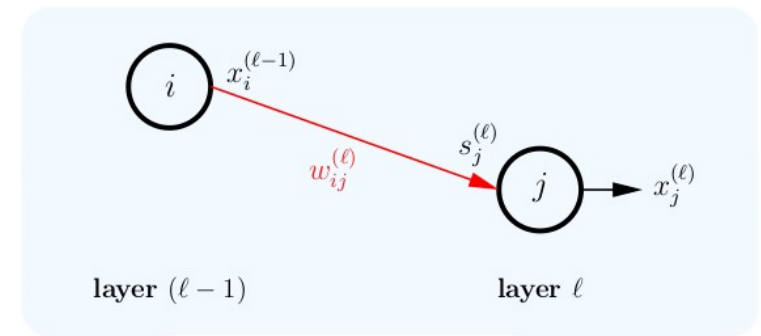
Tenemos:

$$\frac{\partial e}{\partial w_{ij}^{(\ell)}} = \frac{\partial s_j^{(\ell)}}{\partial w_{ij}^{(\ell)}} \cdot \frac{\partial e}{\partial s_j^{(\ell)}}$$

y sabemos que:  $s_j^{(\ell)} = \sum_{\alpha=0}^{d^{(\ell-1)}} w_{\alpha j}^{(\ell)} \mathbf{x}_{\alpha}^{(\ell-1)}$



# Backpropagation



Miremos esto:  $\frac{\partial e}{\partial W^{(\ell)}} = \frac{\partial s^{(\ell)}}{\partial W^{(\ell)}} \cdot \left( \frac{\partial e}{\partial s^{(\ell)}} \right)^T$  a nivel de un enlace.

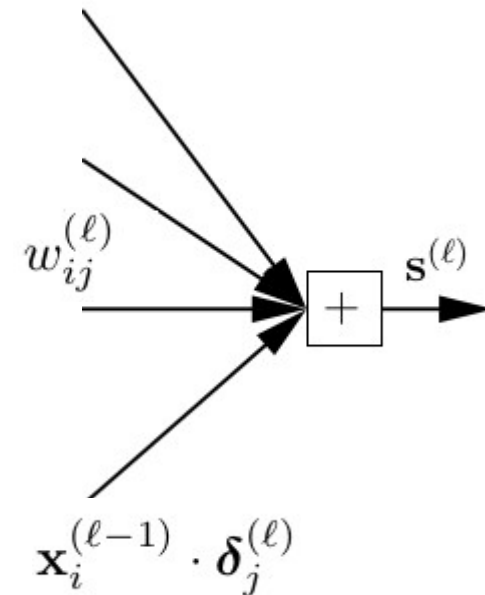
$$= \mathbf{x}^{(\ell-1)} (\boldsymbol{\delta}^{(\ell)})^T$$

Tenemos:

$$\frac{\partial e}{\partial w_{ij}^{(\ell)}} = \frac{\partial s_j^{(\ell)}}{\partial w_{ij}^{(\ell)}} \cdot \frac{\partial e}{\partial s_j^{(\ell)}}$$

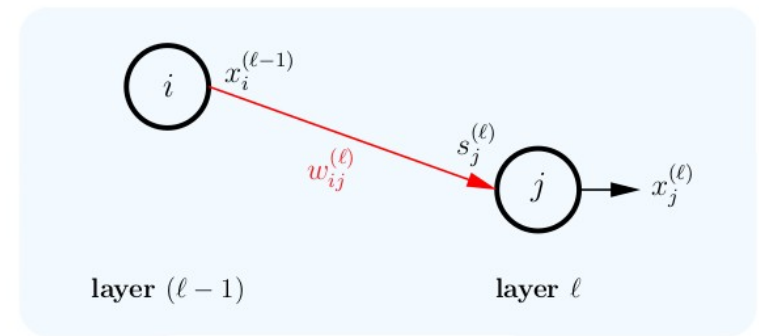
y sabemos que:  $s_j^{(\ell)} = \sum_{\alpha=0}^{d^{(\ell-1)}} w_{\alpha j}^{(\ell)} \mathbf{x}_{\alpha}^{(\ell-1)}$

por lo que al derivar con respecto a  $w_{ij}^{(\ell)}$ , queda:



$$\mathbf{x}_i^{(\ell-1)} \cdot \boldsymbol{\delta}_j^{(\ell)}$$

# Backpropagation



Miremos esto:  $\frac{\partial e}{\partial W^{(\ell)}} = \frac{\partial s^{(\ell)}}{\partial W^{(\ell)}} \cdot \left( \frac{\partial e}{\partial s^{(\ell)}} \right)^T$  a nivel de un enlace.

$$= \mathbf{x}^{(\ell-1)} (\boldsymbol{\delta}^{(\ell)})^T$$

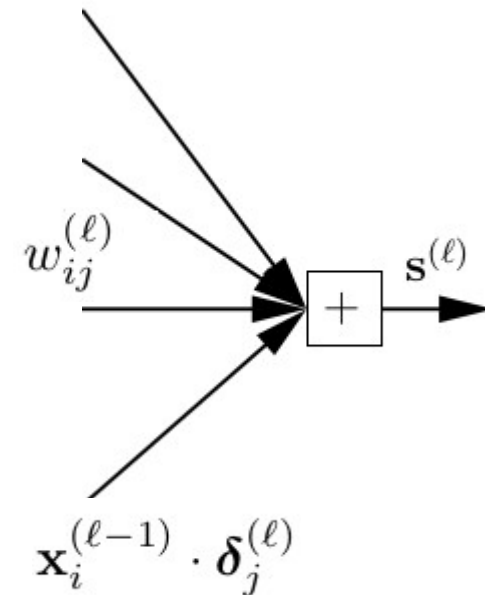
Tenemos:

$$\frac{\partial e}{\partial w_{ij}^{(\ell)}} = \frac{\partial s_j^{(\ell)}}{\partial w_{ij}^{(\ell)}} \cdot \frac{\partial e}{\partial s_j^{(\ell)}}$$

y sabemos que:  $s_j^{(\ell)} = \sum_{\alpha=0}^{d^{(\ell-1)}} w_{\alpha j}^{(\ell)} \mathbf{x}_{\alpha}^{(\ell-1)}$

por lo que al derivar con respecto a  $w_{ij}^{(\ell)}$ , queda:

Luego, haciendo lo mismo para cada parámetro, encontramos que:

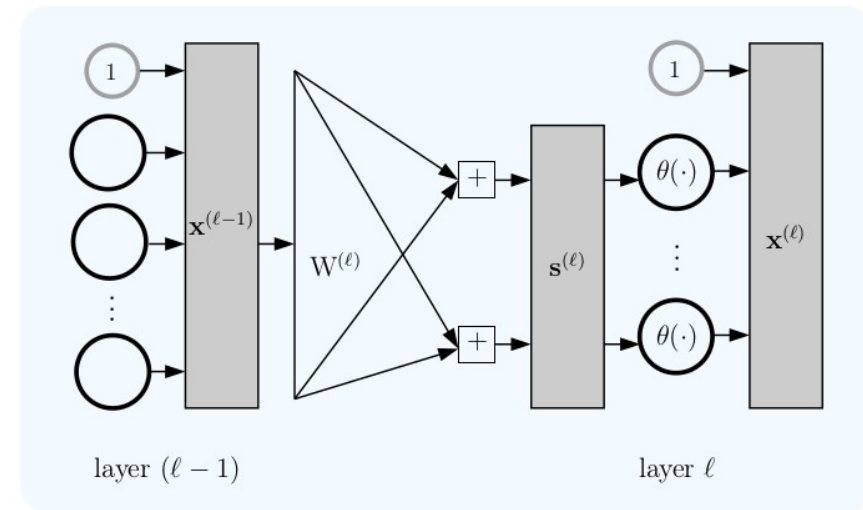


$$\frac{\partial s^{(\ell)}}{\partial W^{(\ell)}} = \mathbf{x}^{(\ell-1)} \quad \checkmark$$

# Backpropagation

Ahora trabajaremos con:

$$\delta_j^{(\ell)} = \frac{\partial e}{\partial s_j^{(\ell)}}$$



$$\frac{\partial \mathbf{s}^{(\ell)}}{\partial W^{(\ell)}} = \mathbf{x}^{(\ell-1)} \quad \checkmark$$

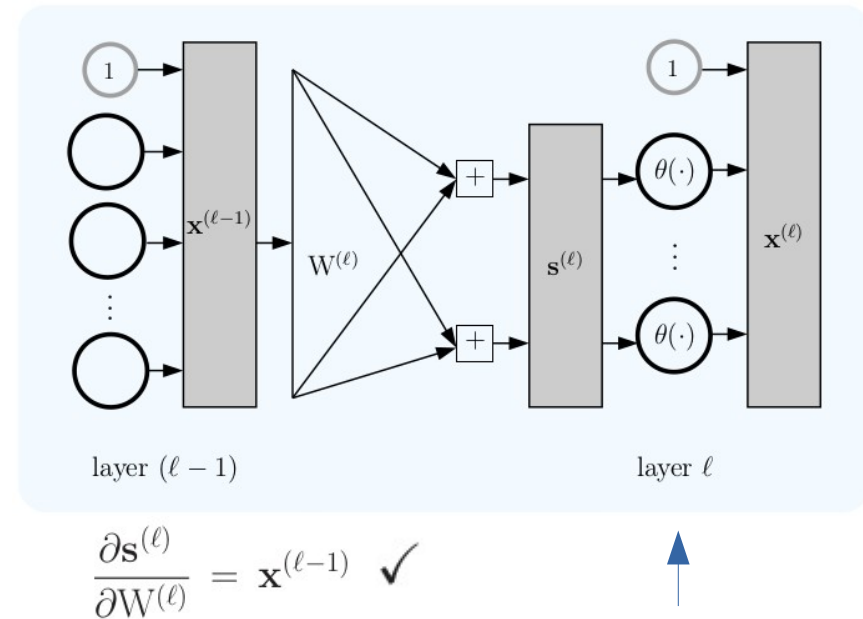


# Backpropagation

Ahora trabajaremos con:  $\delta_j^{(\ell)} = \frac{\partial e}{\partial s_j^{(\ell)}}$

Aplicamos regla de la cadena:

$$\frac{\partial e}{\partial s_j^{(\ell)}} = \frac{\partial e}{\partial x_j^{(\ell)}} \cdot \frac{\partial x_j^{(\ell)}}{\partial s_j^{(\ell)}}$$



# Backpropagation

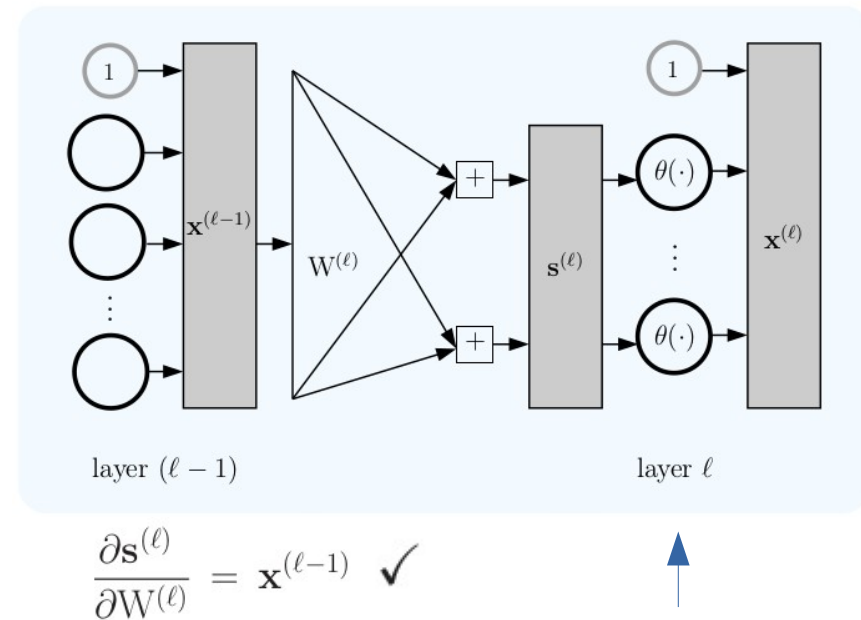
Ahora trabajaremos con:  $\delta_j^{(\ell)} = \frac{\partial e}{\partial s_j^{(\ell)}}$

Aplicamos regla de la cadena:

$$\frac{\partial e}{\partial s_j^{(\ell)}} = \frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial \mathbf{x}_j^{(\ell)}}{\partial s_j^{(\ell)}}$$

$$= \frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} \cdot \theta' \left( s_j^{(\ell)} \right)$$

Derivada de la función de activación



# Backpropagation

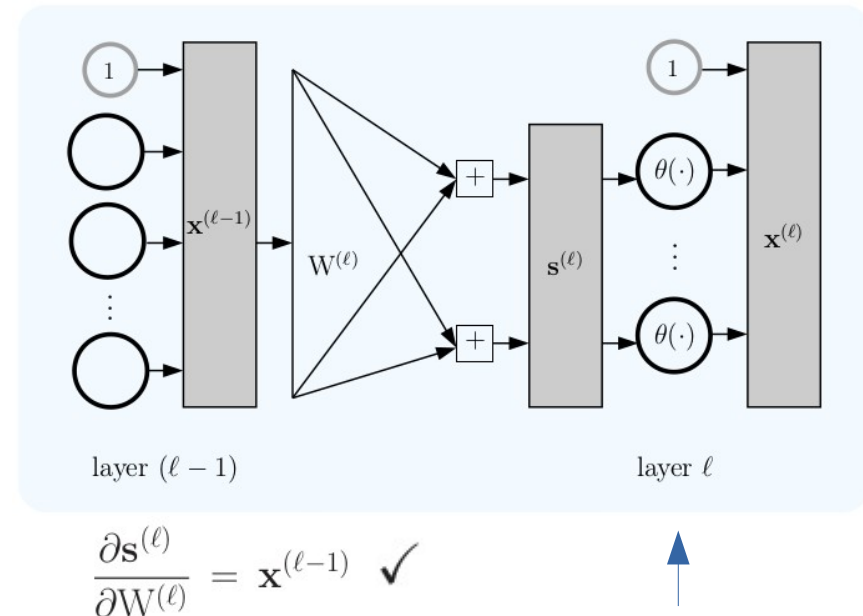
Ahora trabajaremos con:  $\delta_j^{(\ell)} = \frac{\partial e}{\partial s_j^{(\ell)}}$

Aplicamos regla de la cadena:

$$\frac{\partial e}{\partial s_j^{(\ell)}} = \frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial \mathbf{x}_j^{(\ell)}}{\partial s_j^{(\ell)}}$$

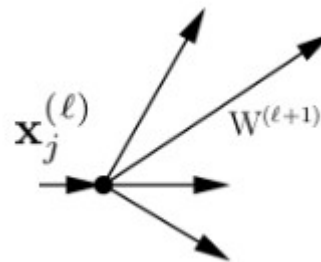
$$= \frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} \cdot \theta' \left( s_j^{(\ell)} \right)$$

Derivada de la función de activación



Ahora veamos que ocurre en:

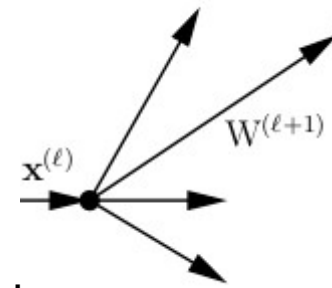
$$\frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}}$$



Hay una multiplexión



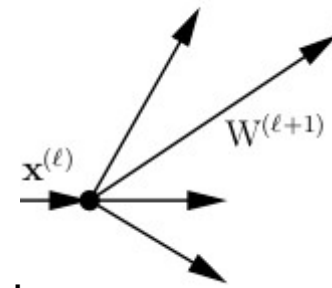
# Backpropagation



Dado que una componente de  $\mathbf{x}^{(\ell)}$  afecta a todas las componentes de  $\mathbf{s}^{(\ell+1)}$ , necesitamos sumar estas dependencias:

$$\frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} = \sum_{k=1}^{d^{(\ell+1)}} \frac{\partial \mathbf{s}_k^{(\ell+1)}}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial e}{\partial \mathbf{s}_k^{(\ell+1)}}$$

# Backpropagation

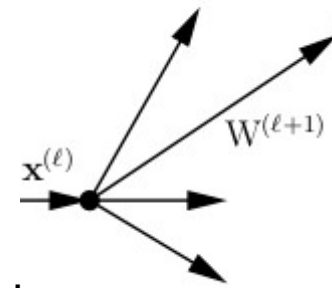


Dado que una componente de  $\mathbf{x}^{(\ell)}$  afecta a todas las componentes de  $\mathbf{s}^{(\ell+1)}$ , necesitamos sumar estas dependencias:

$$\begin{aligned}\frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} &= \sum_{k=1}^{d^{(\ell+1)}} \frac{\partial \mathbf{s}_k^{(\ell+1)}}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial e}{\partial \mathbf{s}_k^{(\ell+1)}} \\ &= \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}.\end{aligned}$$

A blue arrow points from the term  $\frac{\partial \mathbf{s}_k^{(\ell+1)}}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial e}{\partial \mathbf{s}_k^{(\ell+1)}}$  in the first equation to the term  $\frac{\partial \mathbf{x}_j^{(\ell)}}{\partial \mathbf{x}_j^{(\ell)}} \cdot w_{jk}^{(\ell+1)}$  in the second equation, indicating the simplification of the derivative of the input with respect to itself.

# Backpropagation



Dado que una componente de  $\mathbf{x}^{(\ell)}$  afecta a todas las componentes de  $\mathbf{s}^{(\ell+1)}$ , necesitamos sumar estas dependencias:

$$\begin{aligned} \frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} &= \sum_{k=1}^{d^{(\ell+1)}} \frac{\partial \mathbf{s}_k^{(\ell+1)}}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial e}{\partial \mathbf{s}_k^{(\ell+1)}} \\ &= \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)} \end{aligned}$$

The diagram shows a blue arrow pointing from the term  $\frac{\partial \mathbf{s}_k^{(\ell+1)}}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial e}{\partial \mathbf{s}_k^{(\ell+1)}}$  in the first equation to the term  $\frac{\partial \mathbf{x}_j^{(\ell)}}{\partial \mathbf{x}_j^{(\ell)}} \cdot w_{jk}^{(\ell+1)}$  in the second equation, indicating the simplification of the derivative of the input with respect to itself.

Luego:

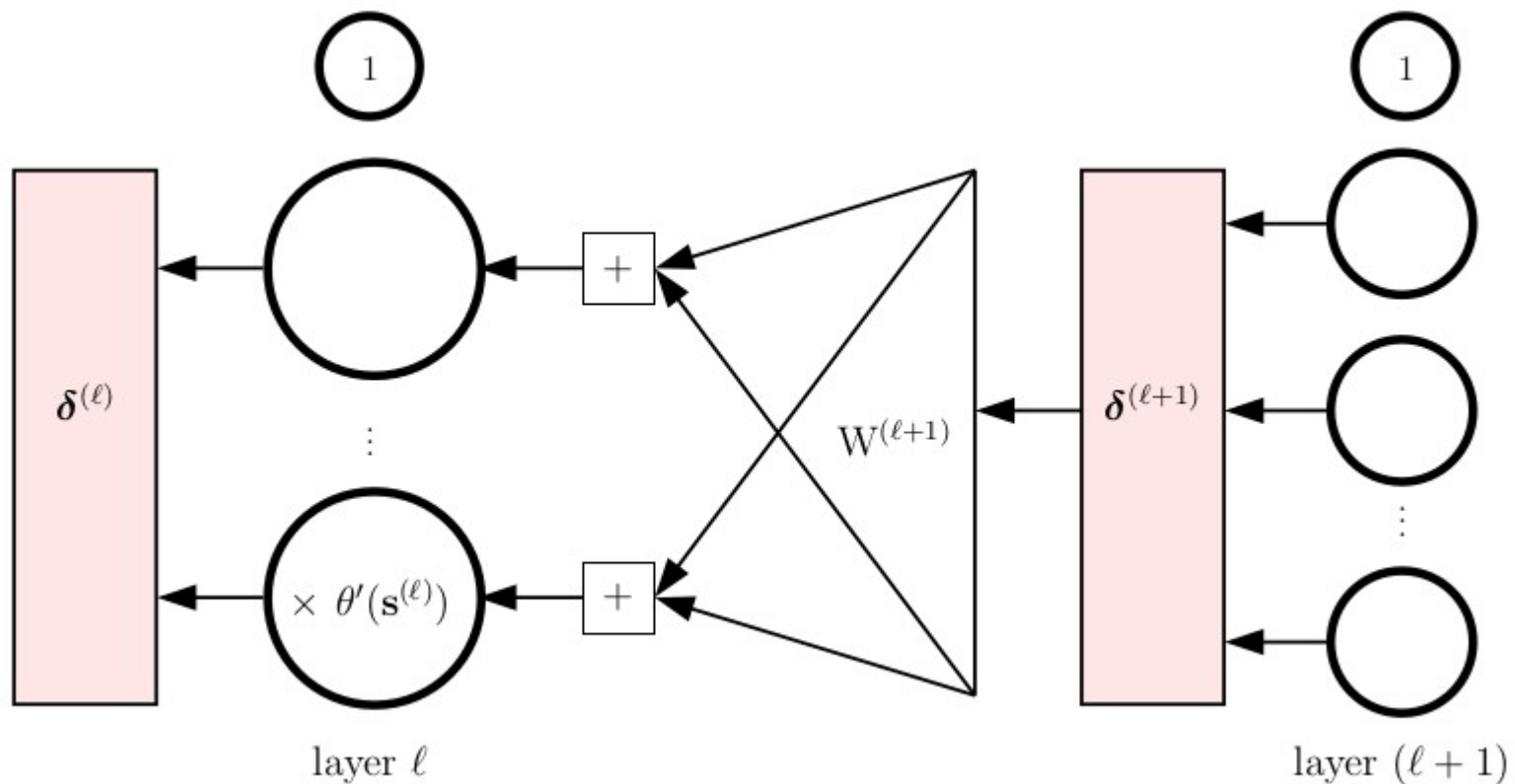
$$\begin{aligned} \frac{\partial e}{\partial \mathbf{s}_j^{(\ell)}} &= \frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial \mathbf{x}_j^{(\ell)}}{\partial \mathbf{s}_j^{(\ell)}} \\ &= \delta_j^{(\ell)} = \theta'(\mathbf{s}_j^{(\ell)}) \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)} \end{aligned}$$

The diagram shows two blue arrows pointing from the terms  $\frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}}$  and  $\frac{\partial \mathbf{x}_j^{(\ell)}}{\partial \mathbf{s}_j^{(\ell)}}$  in the first equation to the terms  $\delta_j^{(\ell)}$  and  $\theta'(\mathbf{s}_j^{(\ell)})$  in the second equation, respectively.

## Backpropagation

$$\delta_j^{(\ell)} = \theta'(s_j^{(\ell)}) \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}$$

$\tan h'(\mathbf{s}^{(\ell)}) = \mathbf{1} - \tan h^2(\mathbf{s}^{(\ell)})$



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

## Backpropagation

$$\delta_j^{(\ell)} = \theta'(\mathbf{s}_j^{(\ell)}) \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}$$

Backpropagation de sensibilidad:

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

Nos falta calcular  $\delta^{(L)}$  .

## Backpropagation

$$\delta_j^{(\ell)} = \theta'(\mathbf{s}_j^{(\ell)}) \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}$$

Backpropagation de sensibilidad:

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

Nos falta calcular  $\delta^{(L)}$ .

Sabemos que:  $e = (\mathbf{x}^{(L)} - y)^2 = (\theta(\mathbf{s}^{(L)}) - y)^2$

# Backpropagation

$$\delta_j^{(\ell)} = \theta'(s_j^{(\ell)}) \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}$$

Backpropagation de sensibilidad:

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

Nos falta calcular  $\delta^{(L)}$ .

Sabemos que:  $e = \underline{(\mathbf{x}^{(L)} - y)^2} = (\theta(\mathbf{s}^{(L)}) - y)^2$

Luego

$$\begin{aligned} \delta^{(L)} &= \frac{\partial e}{\partial \mathbf{s}^{(L)}} \\ &= \frac{\partial}{\partial \mathbf{s}^{(L)}} (\mathbf{x}^{(L)} - y)^2 \end{aligned}$$

# Backpropagation

$$\delta_j^{(\ell)} = \theta'(\mathbf{s}_j^{(\ell)}) \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}$$

Backpropagation de sensibilidad:

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

Nos falta calcular  $\delta^{(L)}$ .

Sabemos que:  $e = \underbrace{(\mathbf{x}^{(L)} - y)^2}_{\text{blue}} = \underbrace{(\theta(\mathbf{s}^{(L)}) - y)^2}_{\text{red}}$

Luego

$$\begin{aligned} \delta^{(L)} &= \frac{\partial e}{\partial \mathbf{s}^{(L)}} \\ &= \frac{\partial}{\partial \mathbf{s}^{(L)}} (\mathbf{x}^{(L)} - y)^2 \\ &= 2(\mathbf{x}^{(L)} - y) \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{s}^{(L)}} \\ &= 2(\mathbf{x}^{(L)} - y) \underbrace{\theta'(\mathbf{s}^{(L)})}_{\text{red}}. \end{aligned}$$

▶  $\tanh'(\mathbf{s}^{(\ell)}) = 1 - \tanh^2(\mathbf{s}^{(\ell)})$   
 $= 1 - (x^{(L)})^2$



# Backpropagation

## Backpropagation

```
1:  $\delta^{(L)} \leftarrow 2(x^{(L)} - y) \cdot \theta'(s^{(L)})$   
2: for  $\ell = L - 1$  to 1 do  
3:   Compute  $\theta'(\mathbf{s}^{(\ell)}) = \left[1 - \mathbf{x}^{(\ell)} \otimes \mathbf{x}^{(\ell)}\right]_1^{d^{(\ell)}}$   
4:    $\boldsymbol{\delta}^{(\ell)} \leftarrow \theta'(\mathbf{s}^{(\ell)}) \otimes \left[W^{(\ell+1)} \boldsymbol{\delta}^{(\ell+1)}\right]_1^{d^{(\ell)}}$   
5: end for
```

Backpropagation nos permite obtener la cadena de sensibilidades:

$$\boldsymbol{\delta}^{(1)} \longleftarrow \boldsymbol{\delta}^{(2)} \dots \longleftarrow \boldsymbol{\delta}^{(L-1)} \longleftarrow \boldsymbol{\delta}^{(L)}$$

## Backpropagation

Recordar que: 
$$\frac{\partial e}{\partial W^{(\ell)}} = \frac{\partial s^{(\ell)}}{\partial W^{(\ell)}} \cdot \left( \frac{\partial e}{\partial s^{(\ell)}} \right)^T$$

Luego, podemos calcular los gradientes para aplicar GD:

$$= \mathbf{x}^{(\ell-1)}(\boldsymbol{\delta}^{(\ell)})^T$$

**Algorithm to Compute  $E_{\text{in}}(\mathbf{w})$  and  $\mathbf{g} = \nabla E_{\text{in}}(\mathbf{w})$ :**

**Input:** weights  $\mathbf{w} = \{W^{(1)}, \dots, W^{(L)}\}$ ; data  $\mathcal{D}$ .

**Output:** error  $E_{\text{in}}(\mathbf{w})$  and gradient  $\mathbf{g} = \{G^{(1)}, \dots, G^{(L)}\}$ .

```
1: Initialize:  $E_{\text{in}} = 0$ ; for  $\ell = 1, \dots, L$ ,  $G^{(\ell)} = 0 \cdot W^{(\ell)}$  .
2: for Each data point  $\mathbf{x}_n$  ( $n = 1, \dots, N$ ) do
3:   Compute  $\mathbf{x}^{(\ell)}$  for  $\ell = 0, \dots, L$ . [forward propagation]
4:   Compute  $\boldsymbol{\delta}^{(\ell)}$  for  $\ell = 1, \dots, L$ . [backpropagation]
5:   for  $\ell = 1, \dots, L$  do
6:      $G^{(\ell)}(\mathbf{x}_n) = [\mathbf{x}^{(\ell-1)}(\boldsymbol{\delta}^{(\ell)})^T]$ 
7:      $G^{(\ell)} \leftarrow G^{(\ell)} + \frac{1}{N} G^{(\ell)}(\mathbf{x}_n)$ .
8:   end for
9: end for
```

## Backpropagation

Recordar que: 
$$\frac{\partial e}{\partial W^{(\ell)}} = \frac{\partial s^{(\ell)}}{\partial W^{(\ell)}} \cdot \left( \frac{\partial e}{\partial s^{(\ell)}} \right)^T$$

Luego, podemos calcular los gradientes para aplicar GD:

$$= \mathbf{x}^{(\ell-1)}(\boldsymbol{\delta}^{(\ell)})^T$$

**Algorithm to Compute  $E_{\text{in}}(\mathbf{w})$  and  $\mathbf{g} = \nabla E_{\text{in}}(\mathbf{w})$ :**

**Input:** weights  $\mathbf{w} = \{W^{(1)}, \dots, W^{(L)}\}$ ; data  $\mathcal{D}$ .

**Output:** error  $E_{\text{in}}(\mathbf{w})$  and gradient  $\mathbf{g} = \{G^{(1)}, \dots, G^{(L)}\}$ .

1: Initialize:  $E_{\text{in}} = 0$ ; for  $\ell = 1, \dots, L$ ,  $G^{(\ell)} = 0 \cdot W^{(\ell)}$ .

2: **for** Each data point  $\mathbf{x}_n$  ( $n = 1, \dots, N$ ) **do**

3:   Compute  $\mathbf{x}^{(\ell)}$  for  $\ell = 0, \dots, L$ . [forward propagation]

4:   Compute  $\boldsymbol{\delta}^{(\ell)}$  for  $\ell = 1, \dots, L$ . [backpropagation]

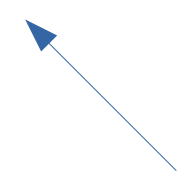
5:   **for**  $\ell = 1, \dots, L$  **do**

6:      $G^{(\ell)}(\mathbf{x}_n) = [\mathbf{x}^{(\ell-1)}(\boldsymbol{\delta}^{(\ell)})^T]$

7:      $G^{(\ell)} \leftarrow G^{(\ell)} + \frac{1}{N} G^{(\ell)}(\mathbf{x}_n)$ .

8:   **end for**

9: **end for**

$$E_{\text{in}} \leftarrow E_{\text{in}} + \frac{1}{N} (\mathbf{x}_n^{(L)} - y_n)^2.$$


# Backpropagation

GD para redes feed-forward:  $W^{(\ell)} = W^{(\ell)} - \eta G^{(\ell)}(\mathbf{x}_n)$ .

