



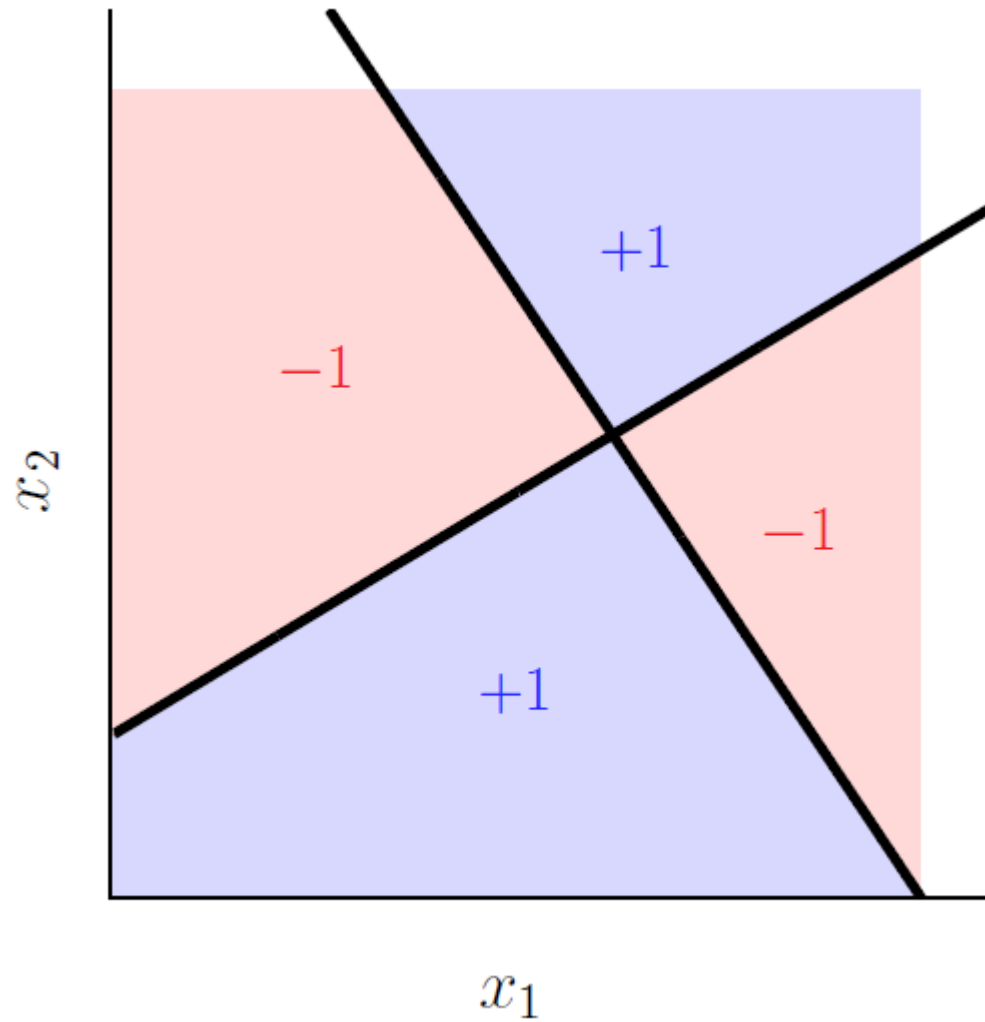
IIC 2433 Minería de Datos

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

- MLP -

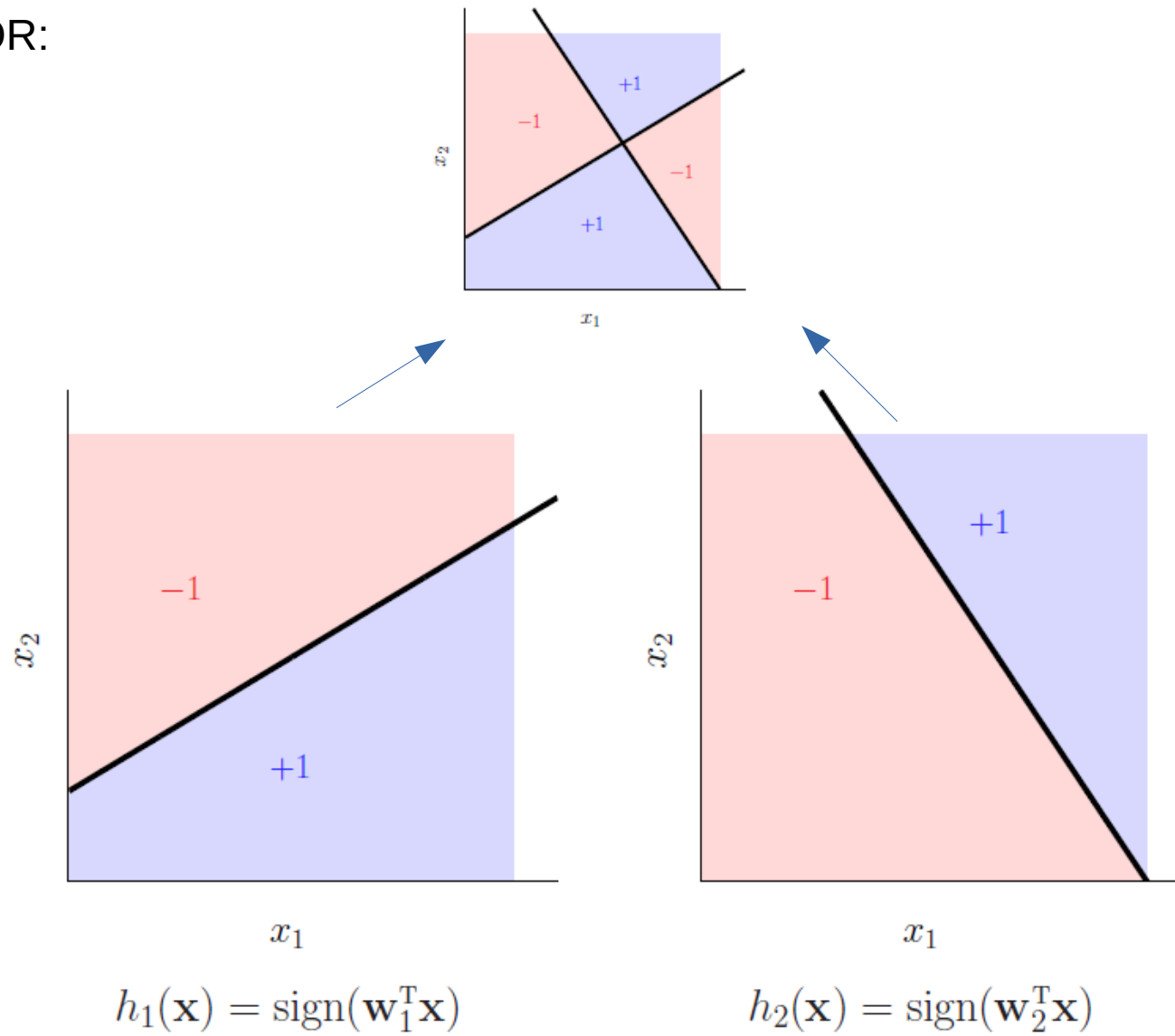
La limitación del modelo lineal

XOR:



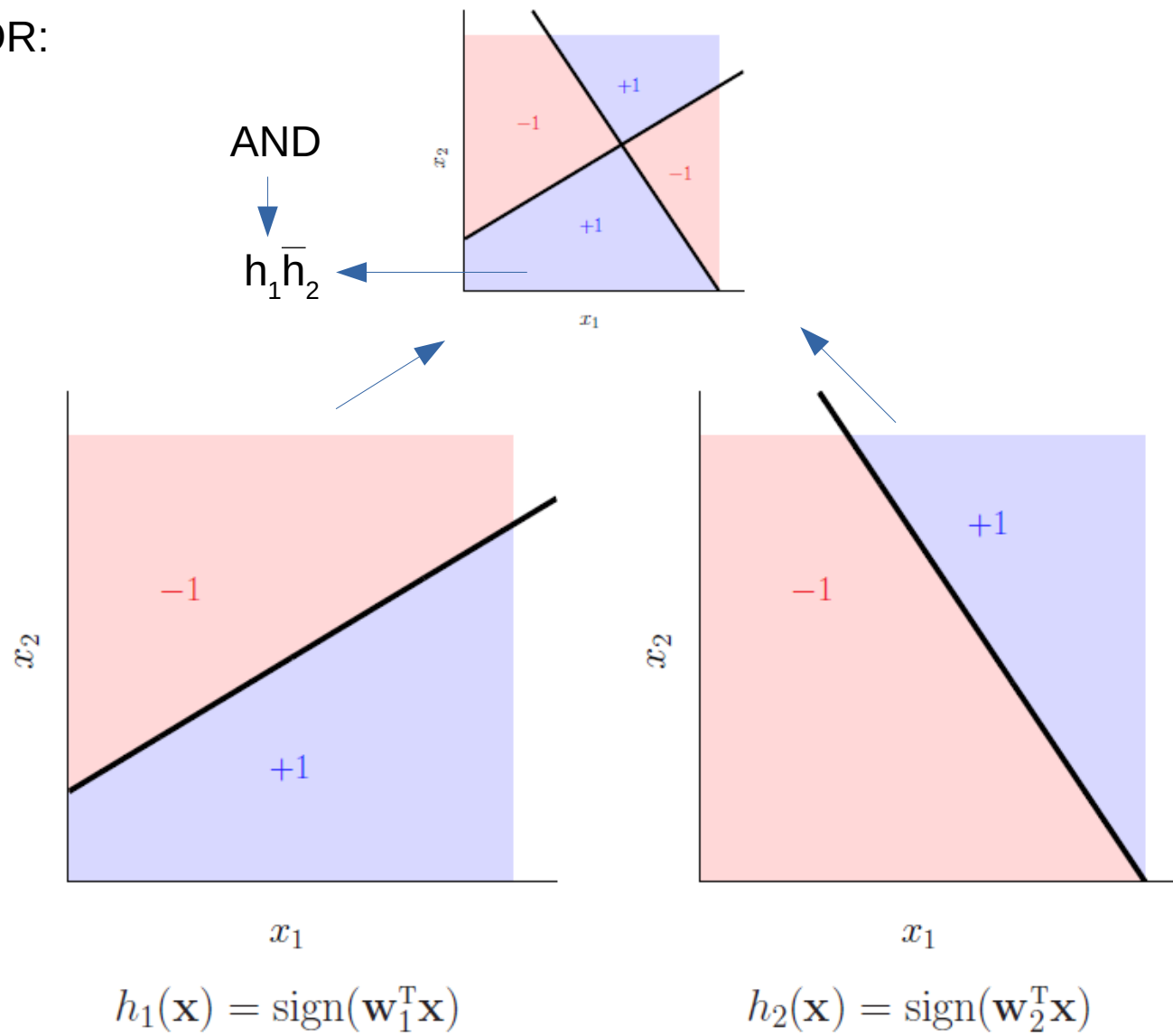
Combinación de perceptrones

XOR:



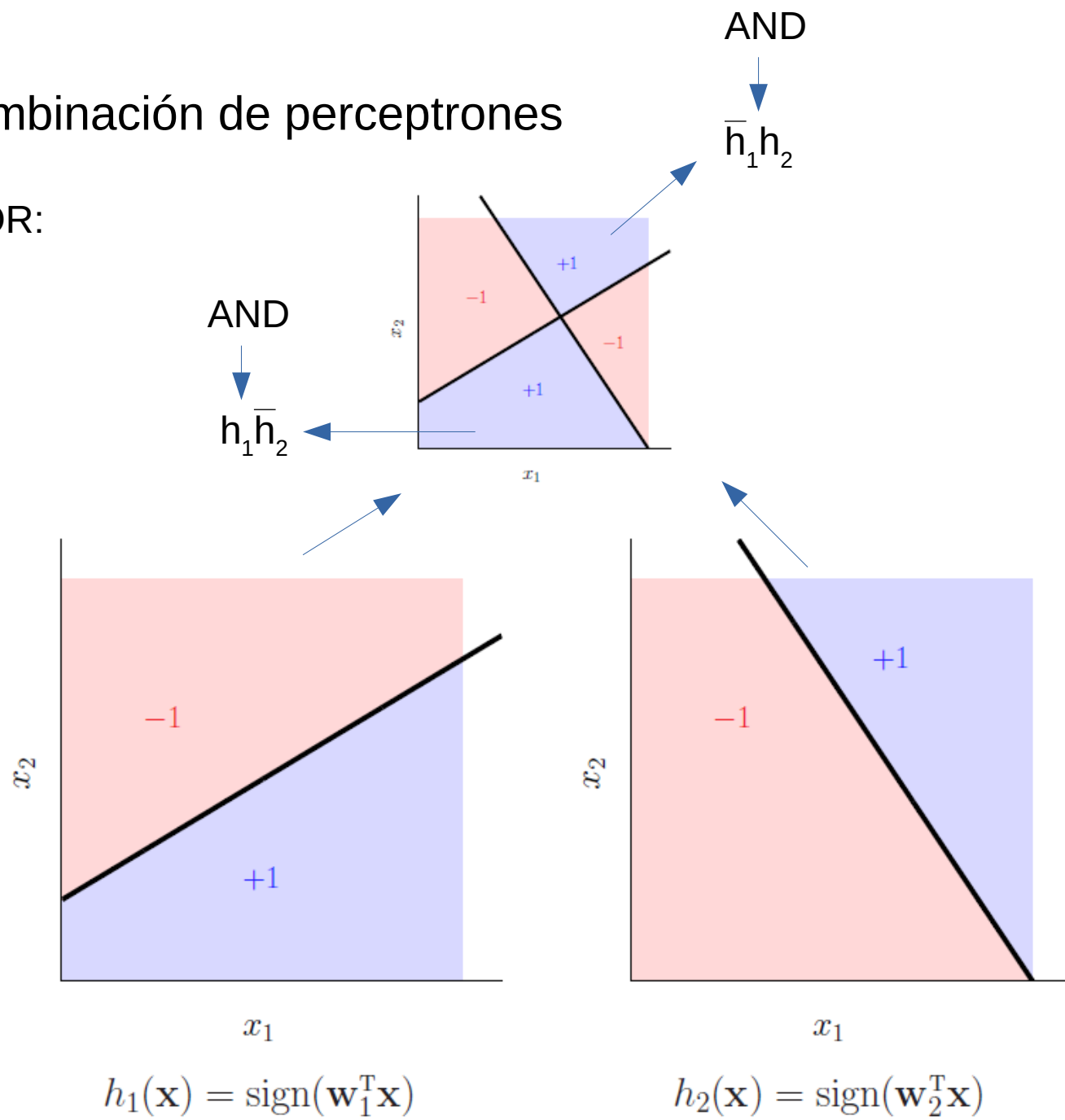
Combinación de perceptrones

XOR:



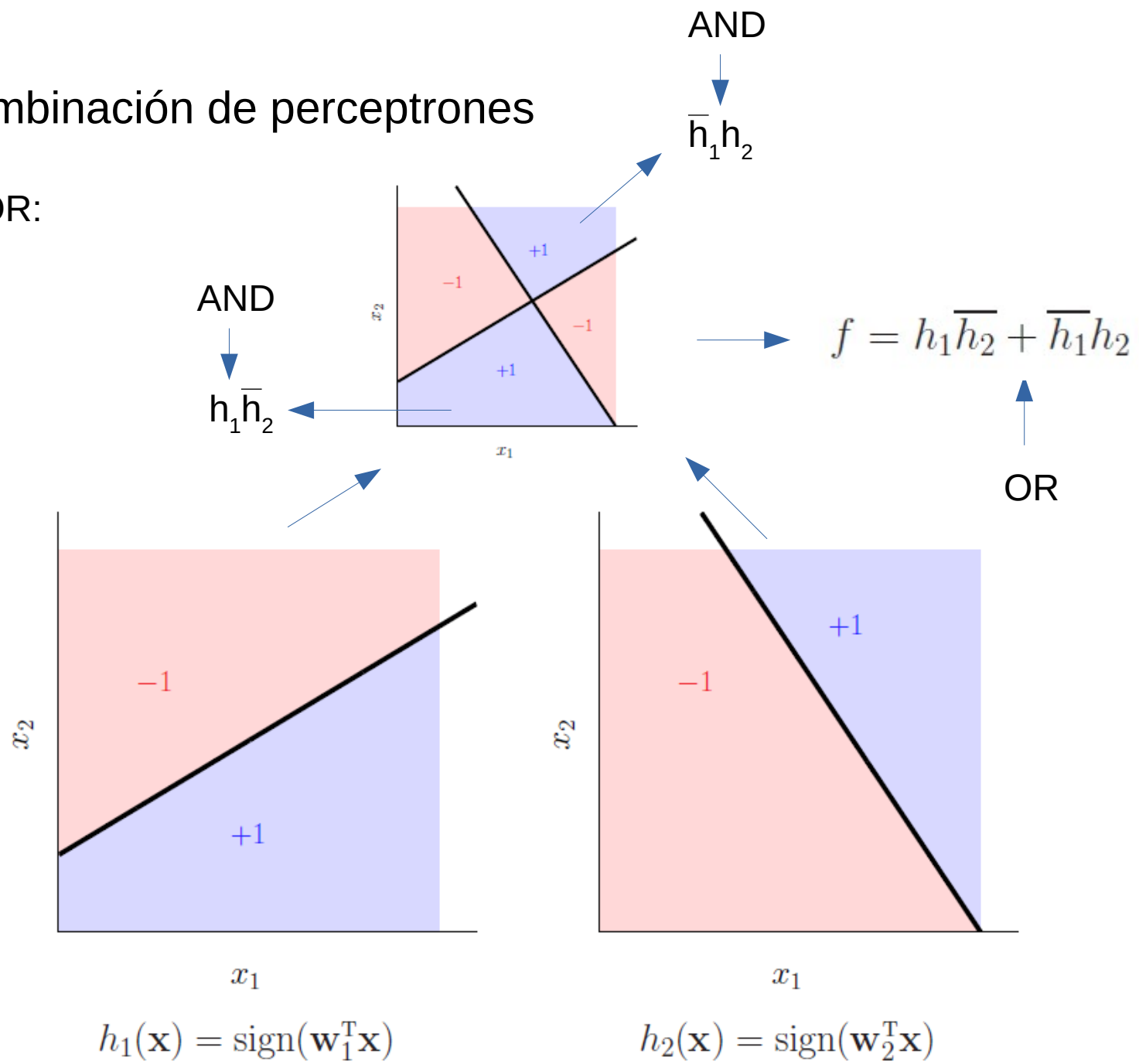
Combinación de perceptrones

XOR:



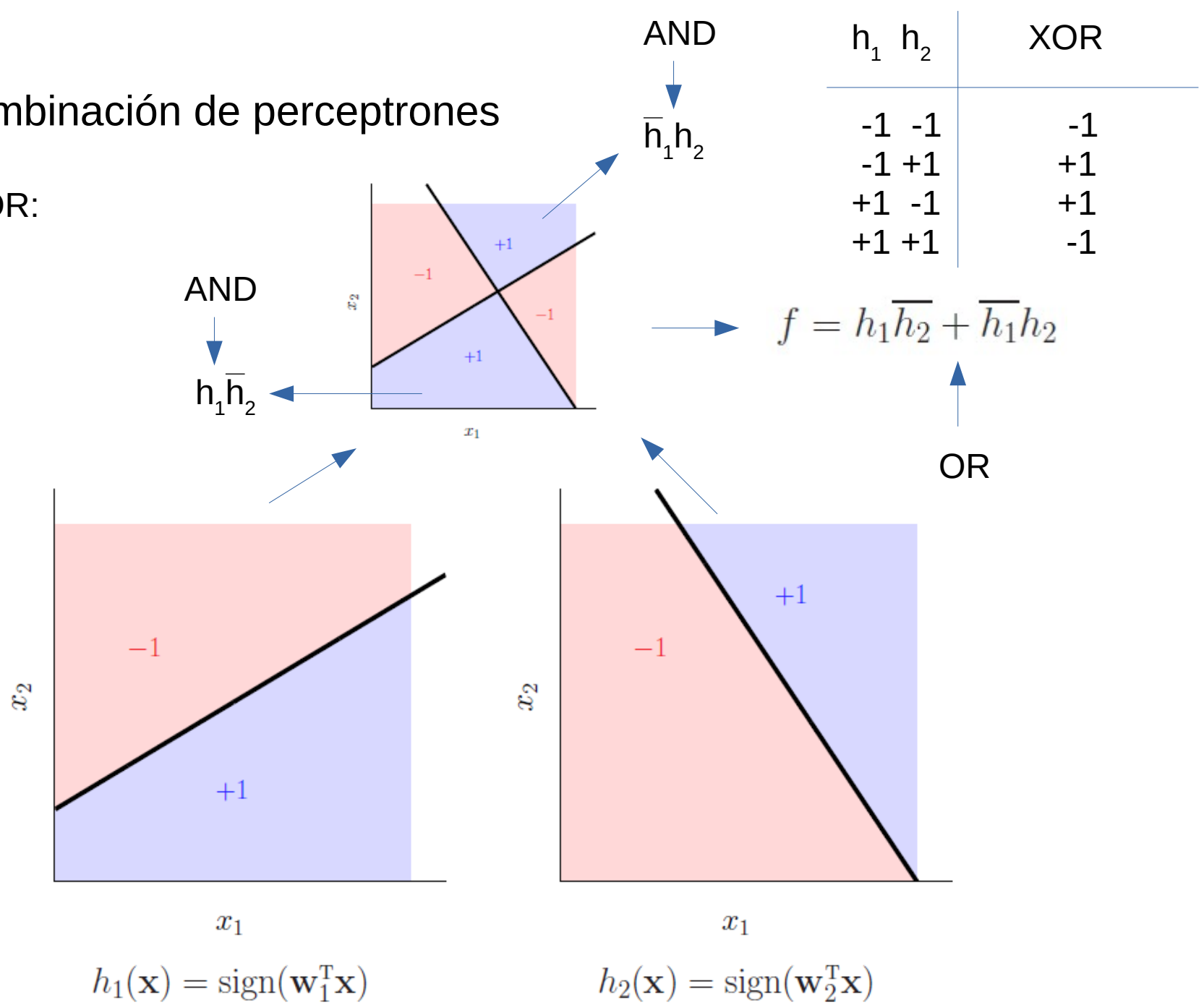
Combinación de perceptrones

XOR:



Combinación de perceptrones

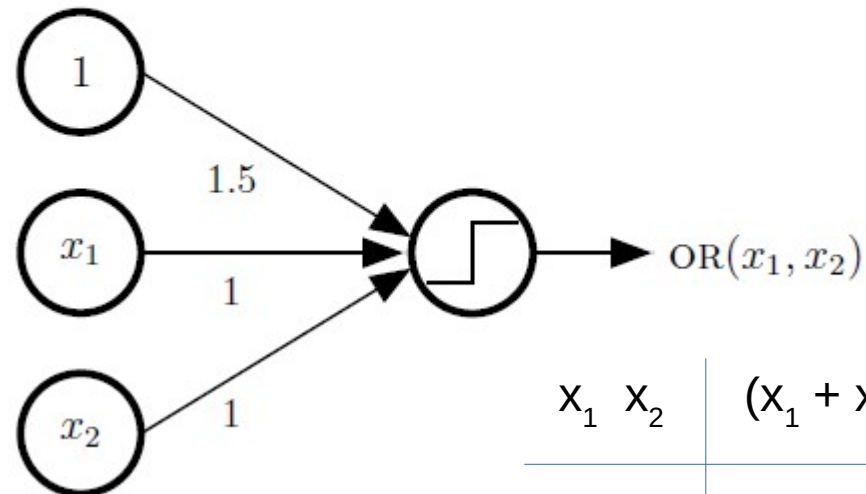
XOR:



Combinación de perceptrones

OR:

$$\text{OR}(x_1, x_2) = \text{sign}(x_1 + x_2 + 1.5)$$

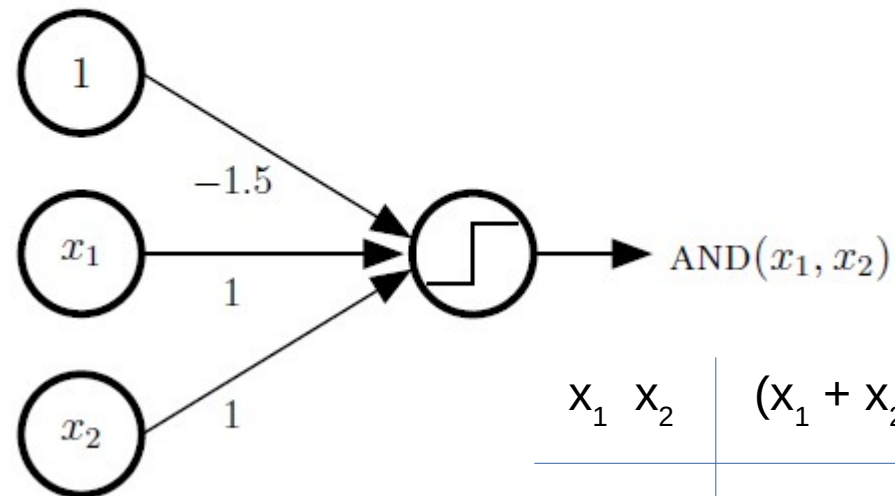


x_1	x_2	$(x_1 + x_2 + 1.5)$	signo
-1	-1	-0.5	-1
-1	+1	1.5	+1
+1	-1	1.5	+1
+1	+1	3.5	+1

Combinación de perceptrones

AND:

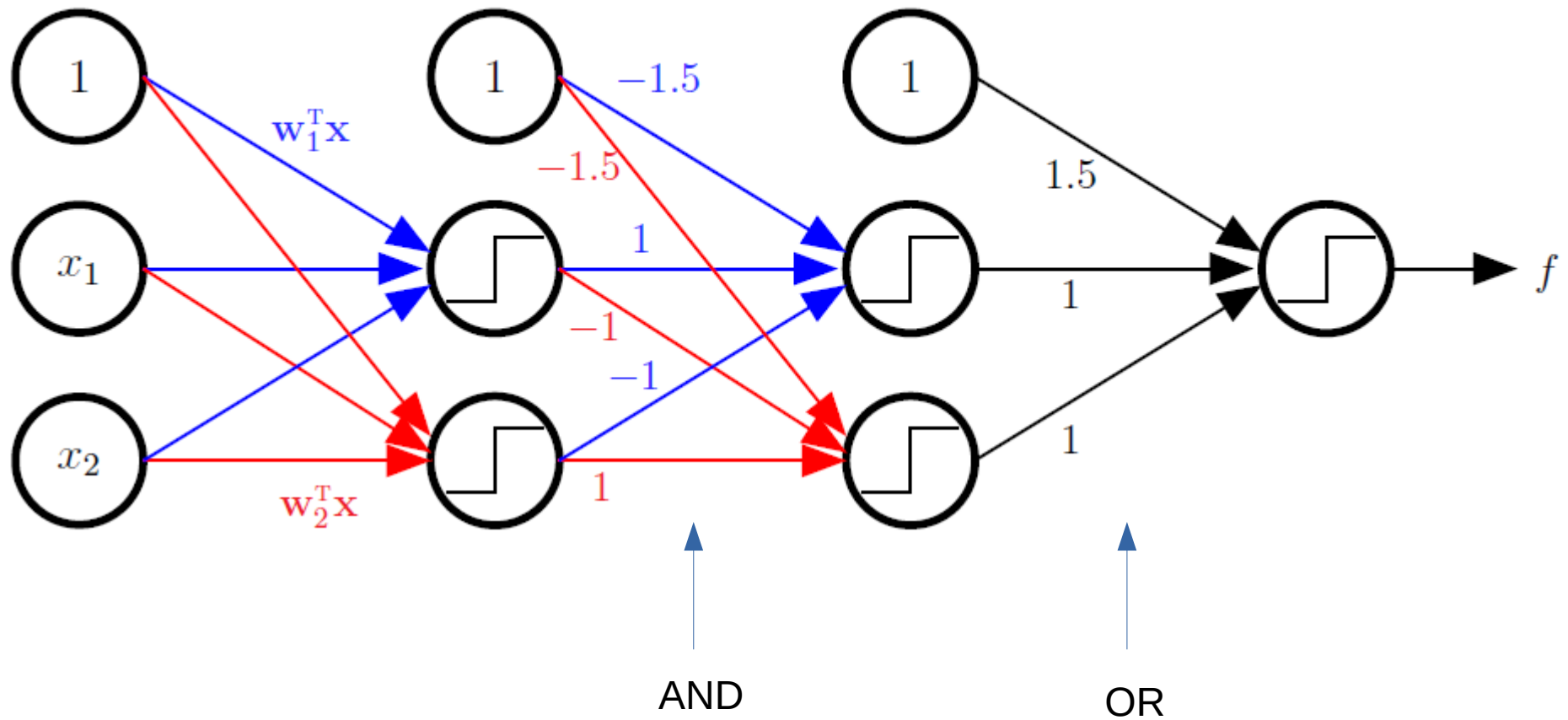
$$\text{AND}(x_1, x_2) = \text{sign}(x_1 + x_2 - 1.5)$$



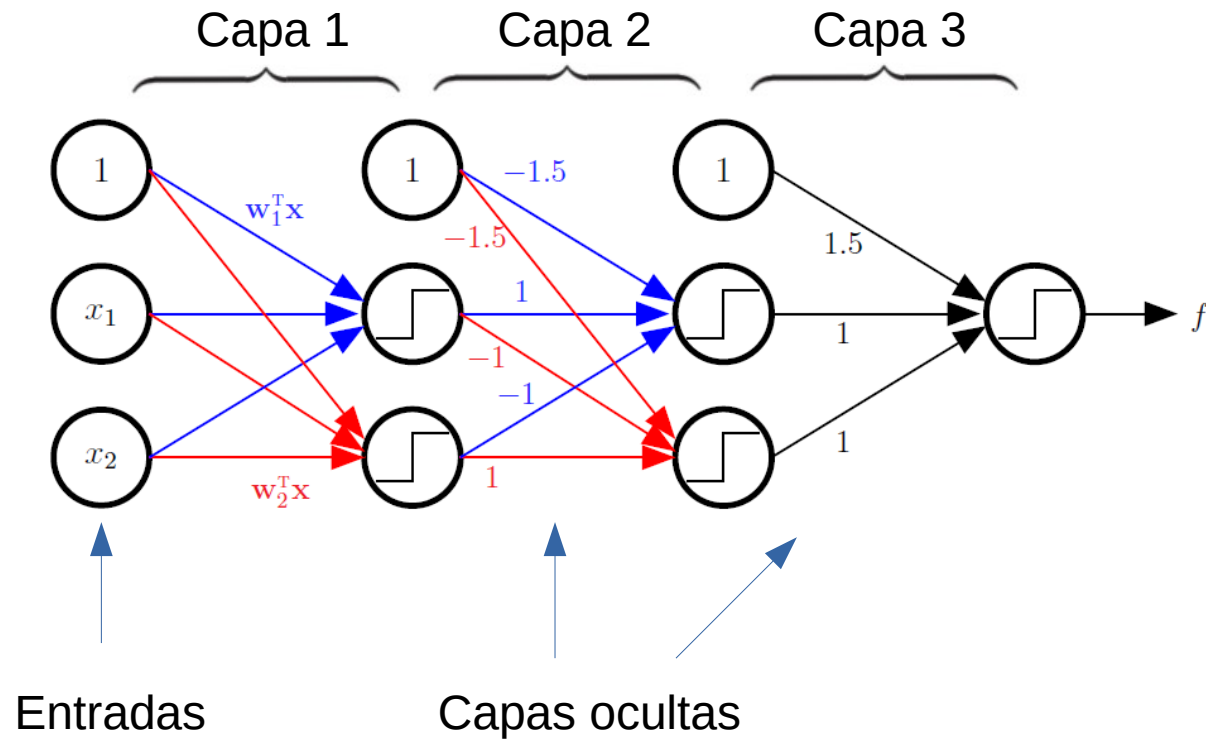
x_1	x_2	$(x_1 + x_2 - 1.5)$	signo
-1	-1	-3.5	-1
-1	+1	-1.5	-1
+1	-1	-1.5	-1
+1	+1	0.5	+1

Combinación de perceptrones

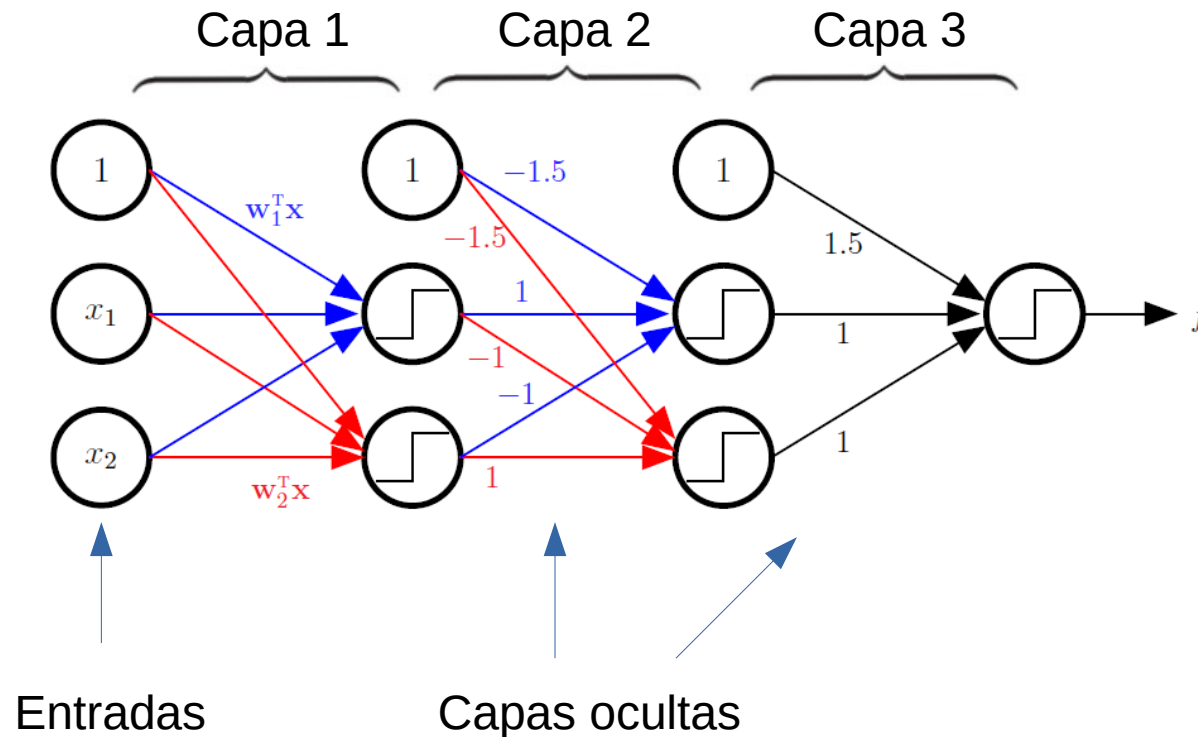
$$f = h_1 \overline{h_2} + \overline{h_1} h_2$$



El perceptrón multicapa (Multi-Layer Perceptron)



El perceptrón multicapa (Multi-Layer Perceptron)

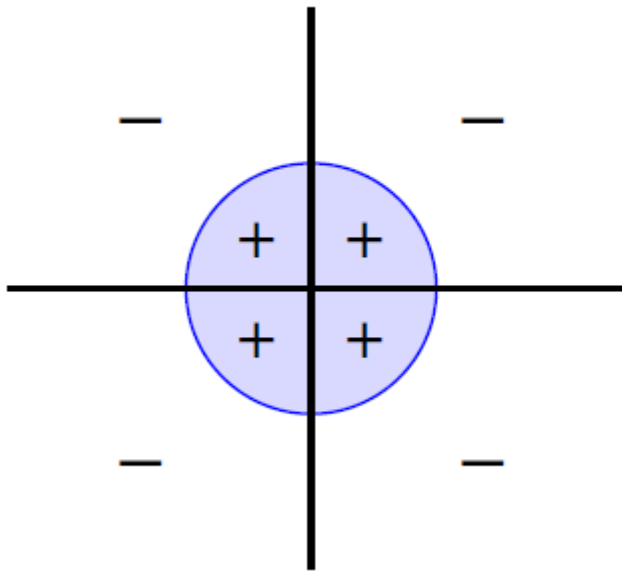


Aproximación universal

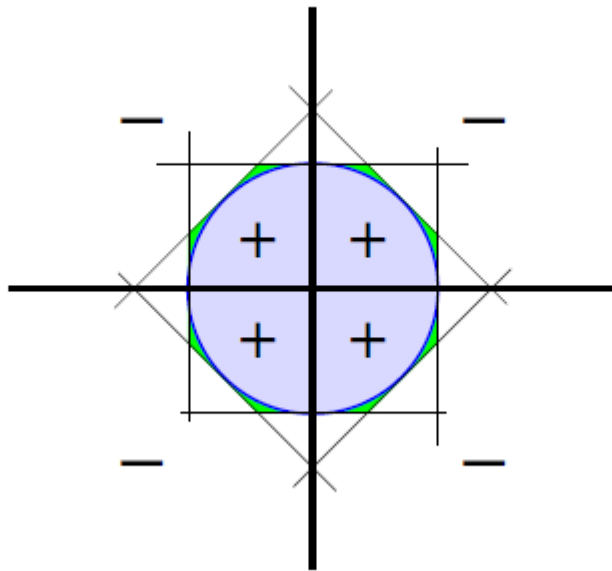
Cualquier función que se puede descomponer en separadores lineales puede ser implementada por un MLP de 3 capas

El perceptrón multicapa (Multi-Layer Perceptron)

Un separador suave puede ser aproximado por N separadores lineales.



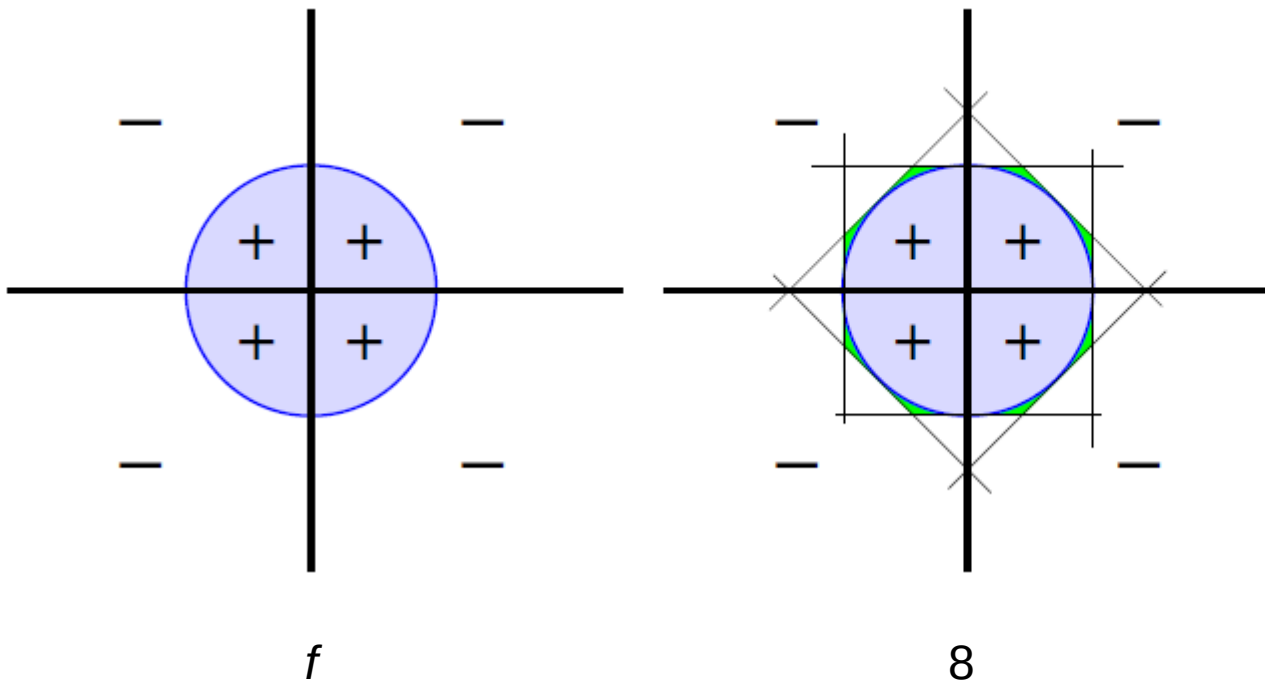
f



¿Cuántos perceptrones hay?

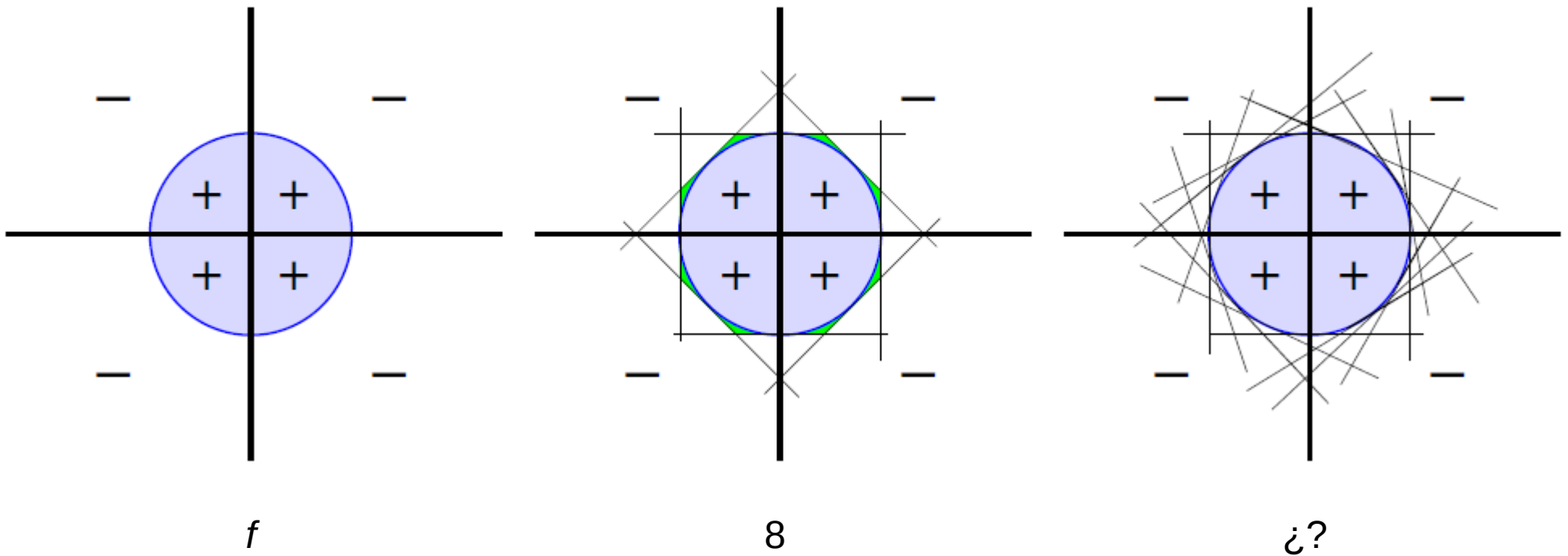
El perceptrón multicapa (Multi-Layer Perceptron)

Un separador suave puede ser aproximado por N separadores lineales.



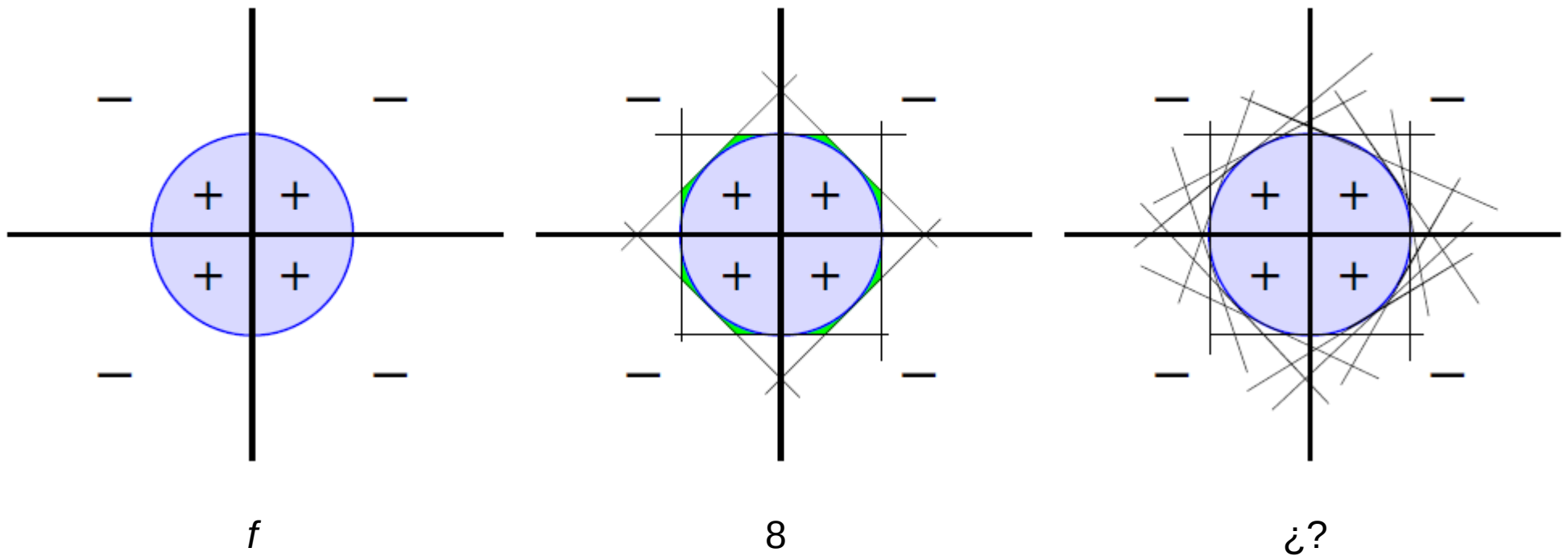
El perceptrón multicapa (Multi-Layer Perceptron)

Un separador suave puede ser aproximado por N separadores lineales.



El perceptrón multicapa (Multi-Layer Perceptron)

Un separador suave puede ser aproximado por N separadores lineales.



Tradeoff aproximación-generalización:

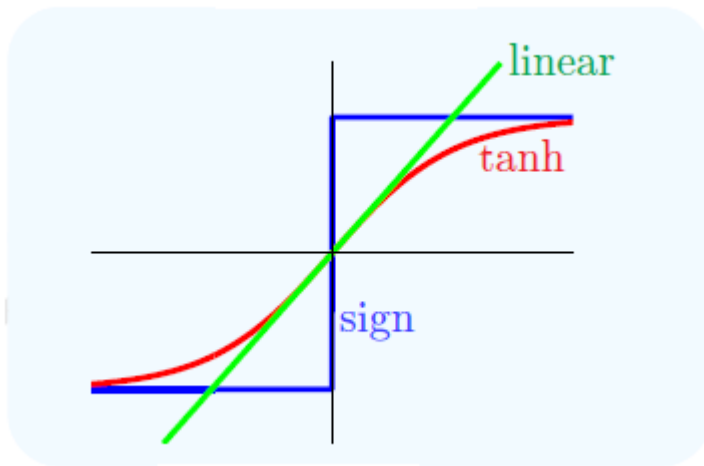
Más neuronas, mejor aproximación.

Más neuronas, más datos.

El perceptrón multicapa (Multi-Layer Perceptron)

Para entrenar un MLP necesitamos reemplazar la función signo dado que no es diferenciable.

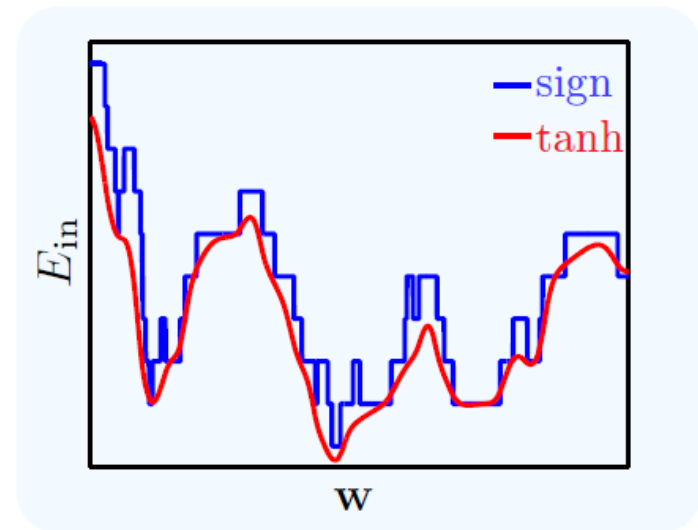
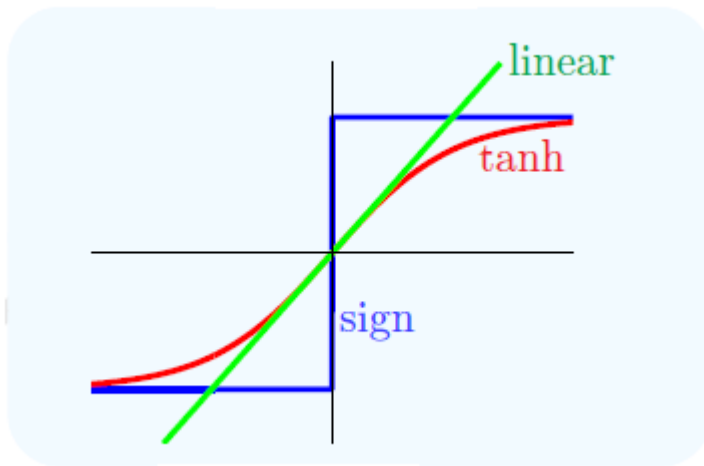
Podemos suavizar la función signo con la tangente hiperbólica, que tiene un comportamiento lineal cerca del origen y es cercana a +1 o -1 para entradas grandes.



El perceptrón multicapa (Multi-Layer Perceptron)

Para entrenar un MLP necesitamos reemplazar la función signo dado que no es diferenciable.

Podemos suavizar la función signo con la tangente hiperbólica, que tiene un comportamiento lineal cerca del origen y es cercana a +1 o -1 para entradas grandes.



- 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_{in} es diferenciable usando GD sobre

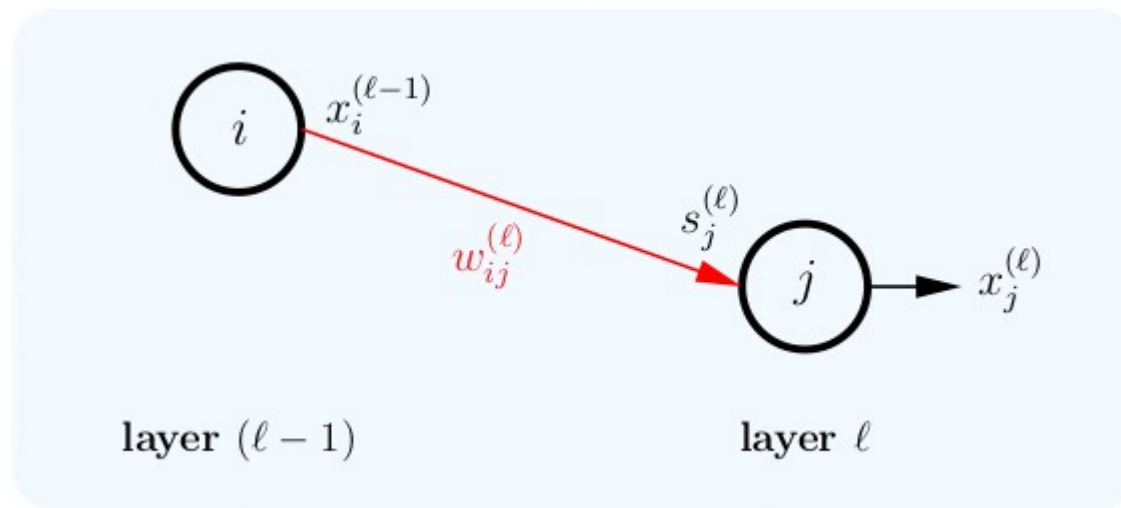
$$W = \{W^{(1)}, W^{(2)}, \dots, W^{(L)}\}$$

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:

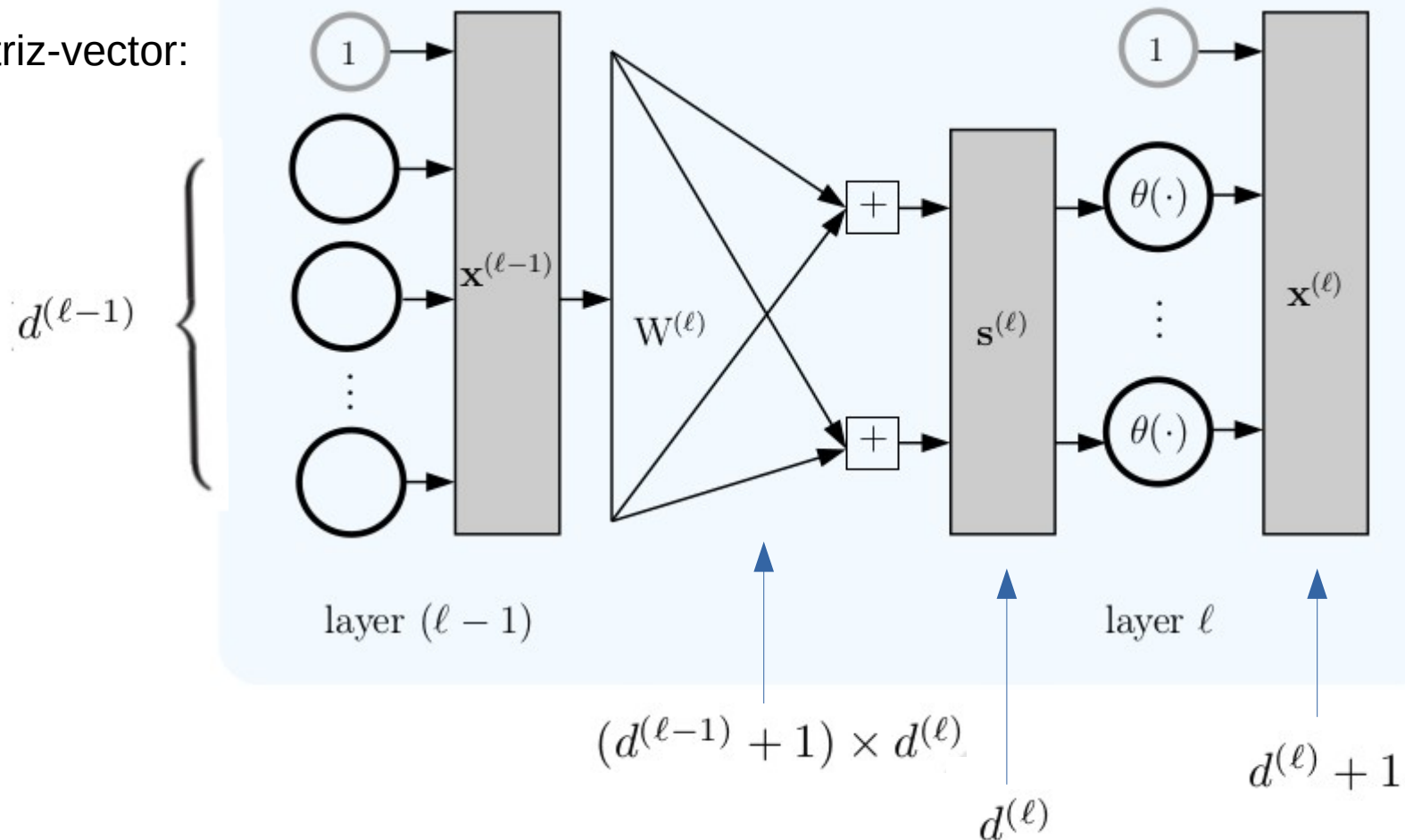


Feed-Forward

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

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

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

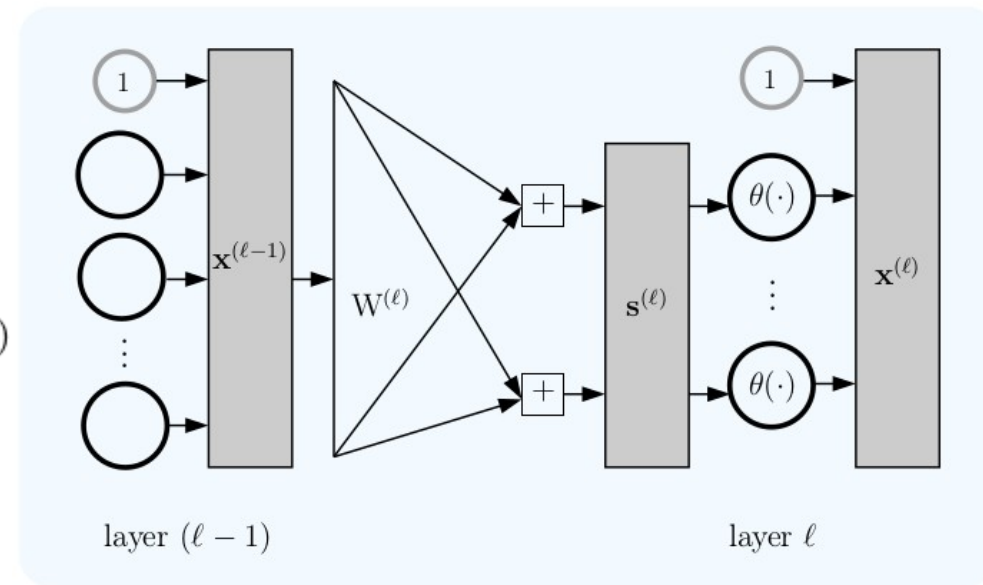
$$\frac{\partial E_{\text{in}}(\mathbf{w})}{\partial W^{(\ell)}} = \frac{1}{N} \sum_{n=1}^N \frac{\partial \mathbf{e}_n}{\partial W^{(\ell)}}$$

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 ℓ .



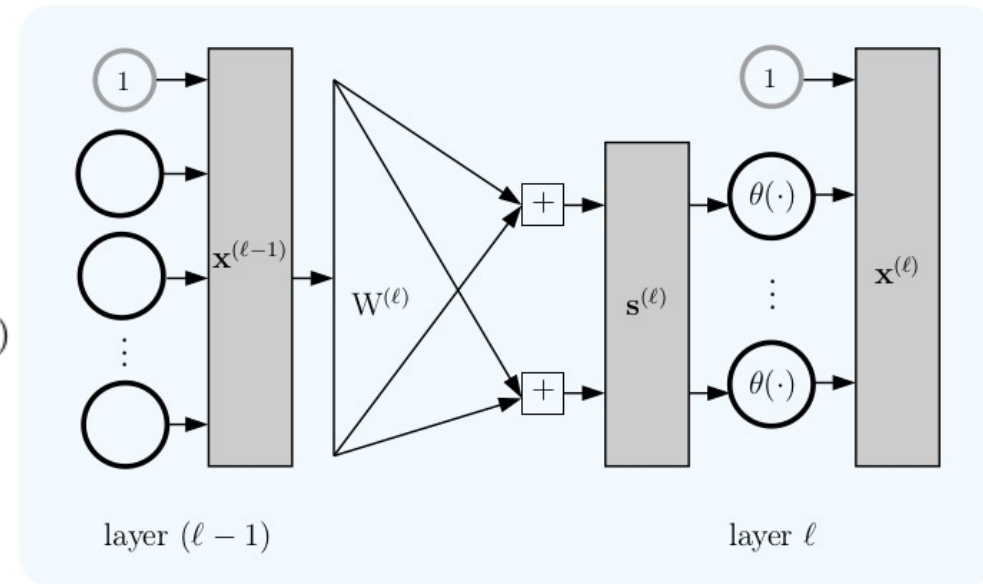
Backpropagation

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 ℓ .

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

Definimos la sensibilidad de la capa ℓ :

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



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 ℓ .

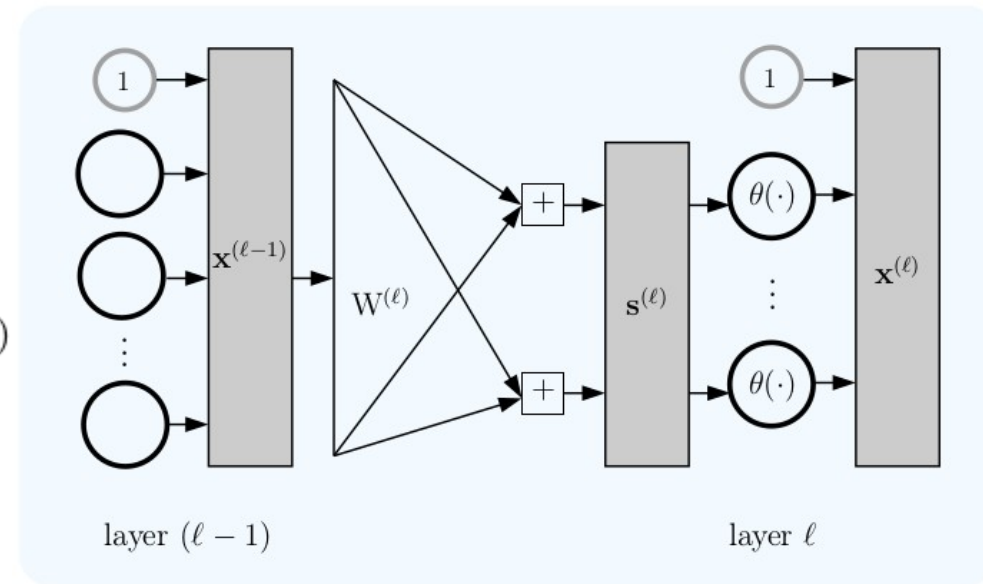
Tenemos: $\mathbf{s}^{(\ell)} = (\mathbf{W}^{(\ell)})^T \mathbf{x}^{(\ell-1)}$

Definimos la sensibilidad de la capa ℓ :

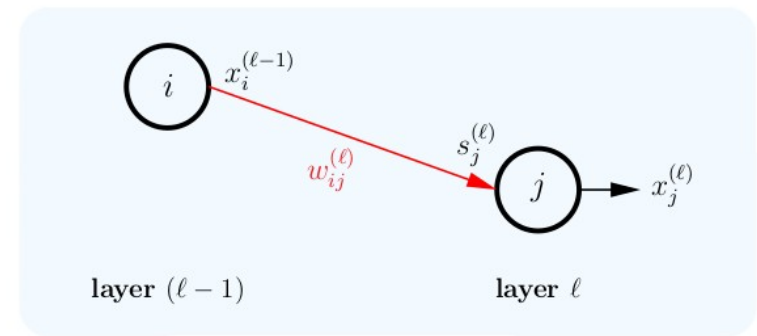
$$\boldsymbol{\delta}^{(\ell)} = \frac{\partial e}{\partial \mathbf{s}^{(\ell)}}$$

Aplicando la regla de la cadena:

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



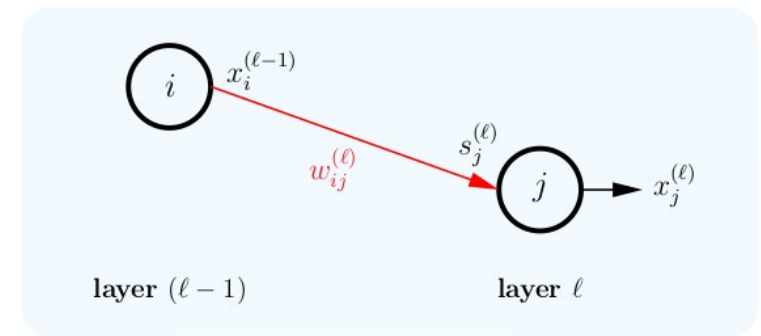
Backpropagation



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



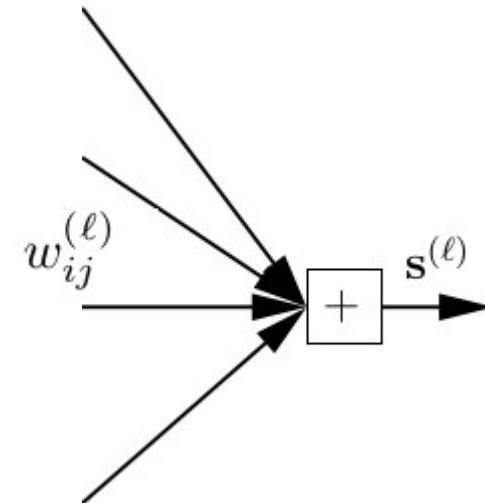
Miremos esto: $\frac{\partial e}{\partial W^{(\ell)}} = \frac{\partial s^{(\ell)}}{\partial W^{(\ell)}} \cdot \left(\frac{\partial e}{\partial s^{(\ell)}} \right)^T$ a nivel de un enlace.

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

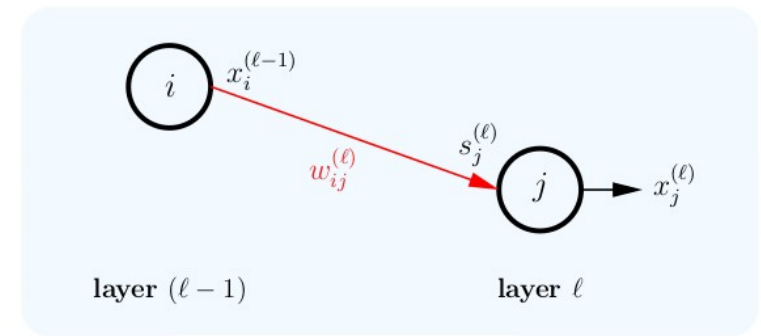
Tenemos:

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

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.

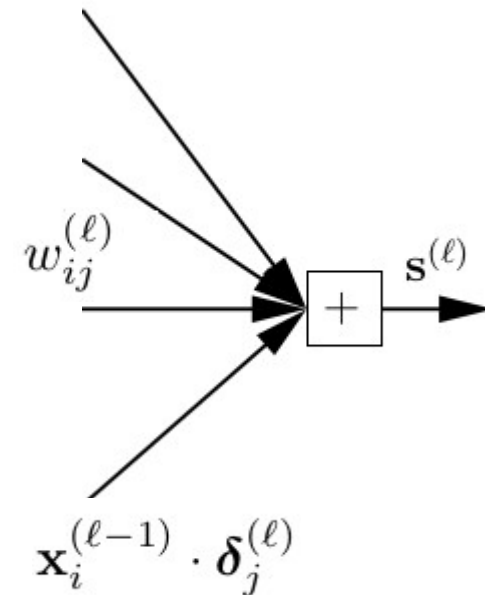
$$= \mathbf{x}^{(\ell-1)} (\boldsymbol{\delta}^{(\ell)})^T$$

Tenemos:

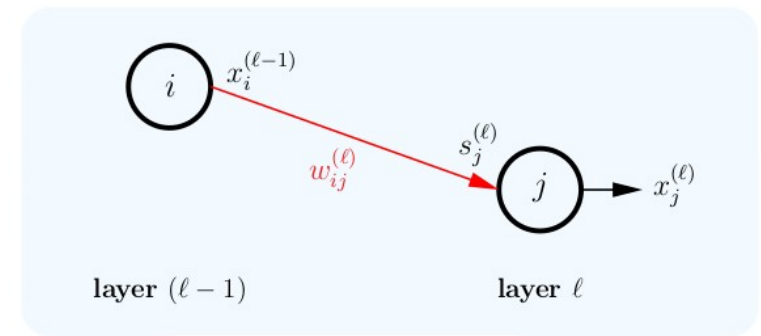
$$\frac{\partial e}{\partial w_{ij}^{(\ell)}} = \frac{\partial s_j^{(\ell)}}{\partial w_{ij}^{(\ell)}} \cdot \frac{\partial e}{\partial s_j^{(\ell)}}$$

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:



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.

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

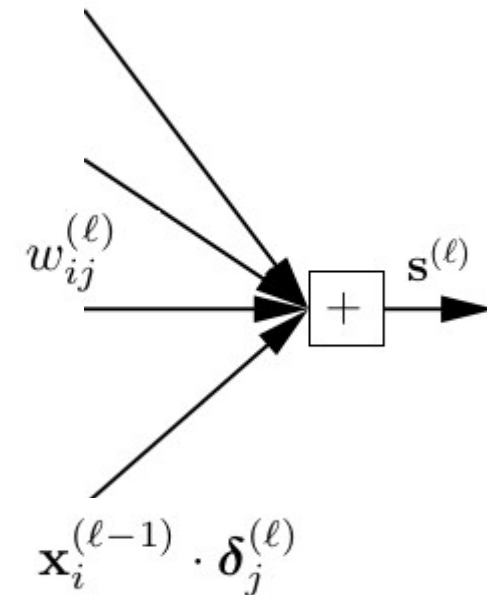
Tenemos:

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

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:

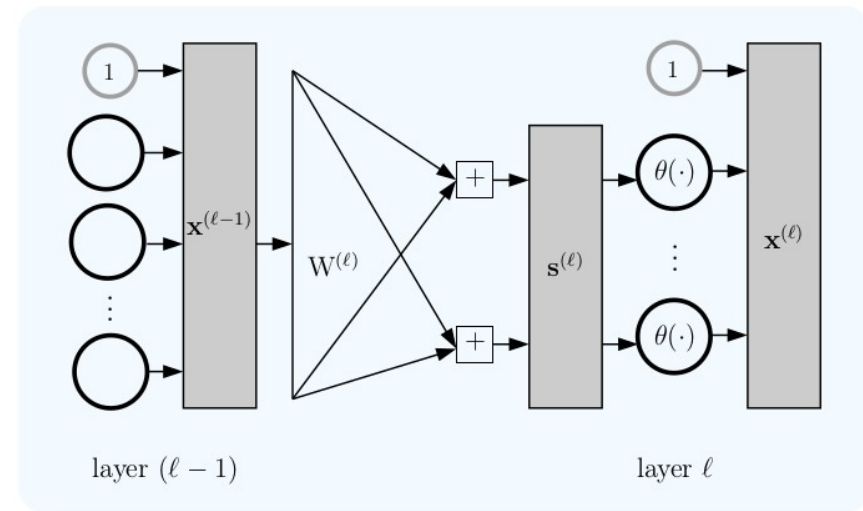


$$\frac{\partial s^{(\ell)}}{\partial W^{(\ell)}} = \mathbf{x}^{(\ell-1)} \quad \checkmark$$

Backpropagation

Ahora trabajaremos con:

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



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

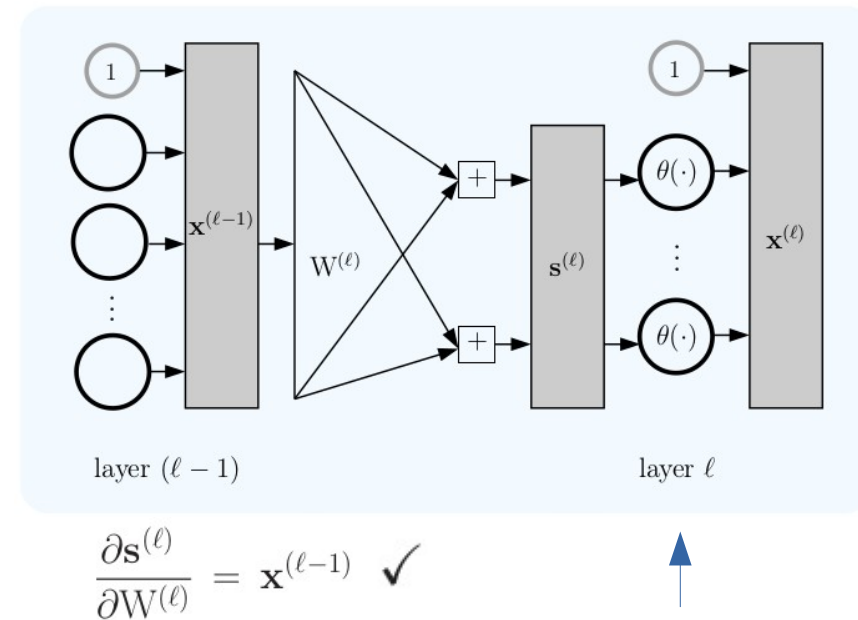


Backpropagation

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

Aplicamos regla de la cadena:

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



Backpropagation

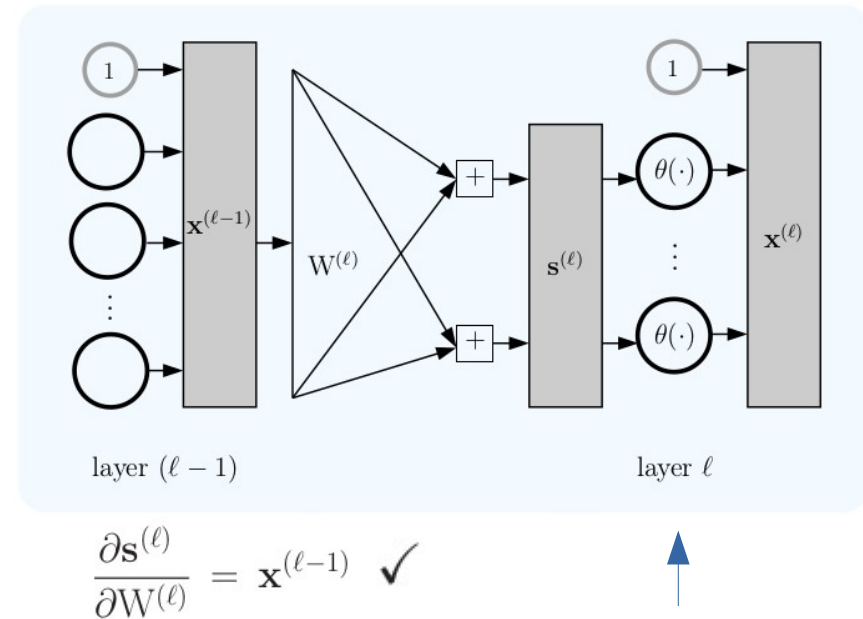
Ahora trabajaremos con: $\delta_j^{(\ell)} = \frac{\partial e}{\partial s_j^{(\ell)}}$

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

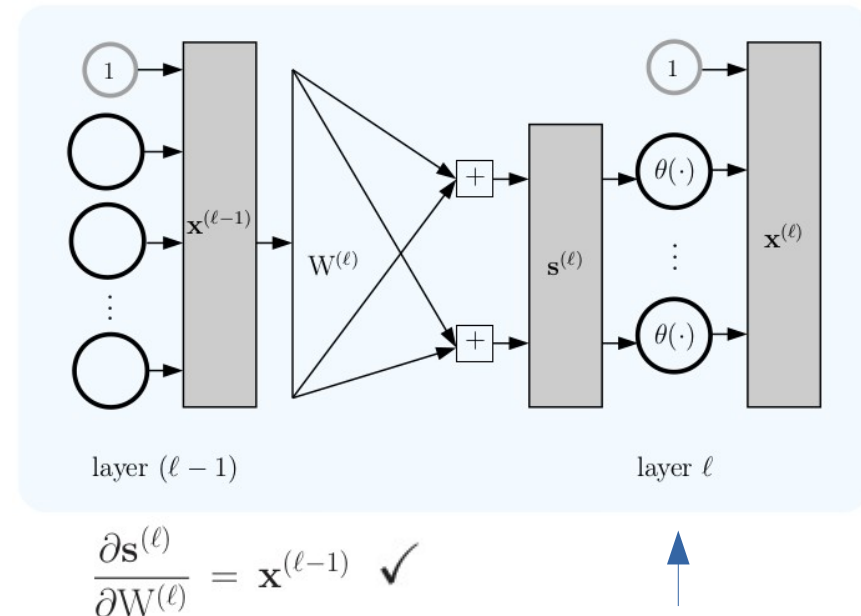
Ahora trabajaremos con: $\delta_j^{(\ell)} = \frac{\partial e}{\partial s_j^{(\ell)}}$

Aplicamos regla de la cadena:

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

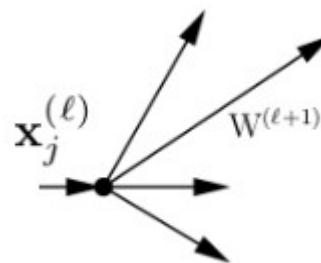
$$= \frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} \cdot \theta' \left(s_j^{(\ell)} \right)$$

Derivada de la función de activación



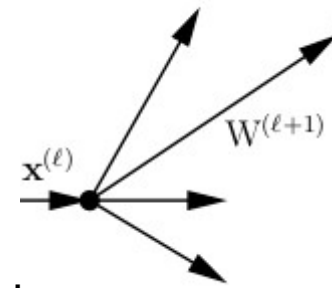
Ahora veamos que ocurre en:

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



Hay una multiplexión

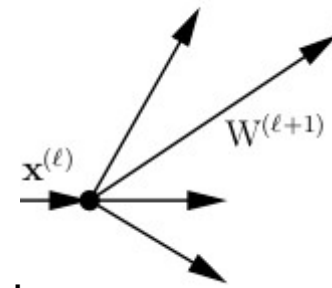
Backpropagation



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

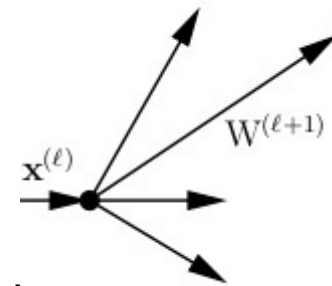


Dado que una componente de $\mathbf{x}^{(\ell)}$ afecta a todas las componentes de $\mathbf{s}^{(\ell+1)}$, necesitamos sumar estas dependencias:

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

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 to 1.

Luego:

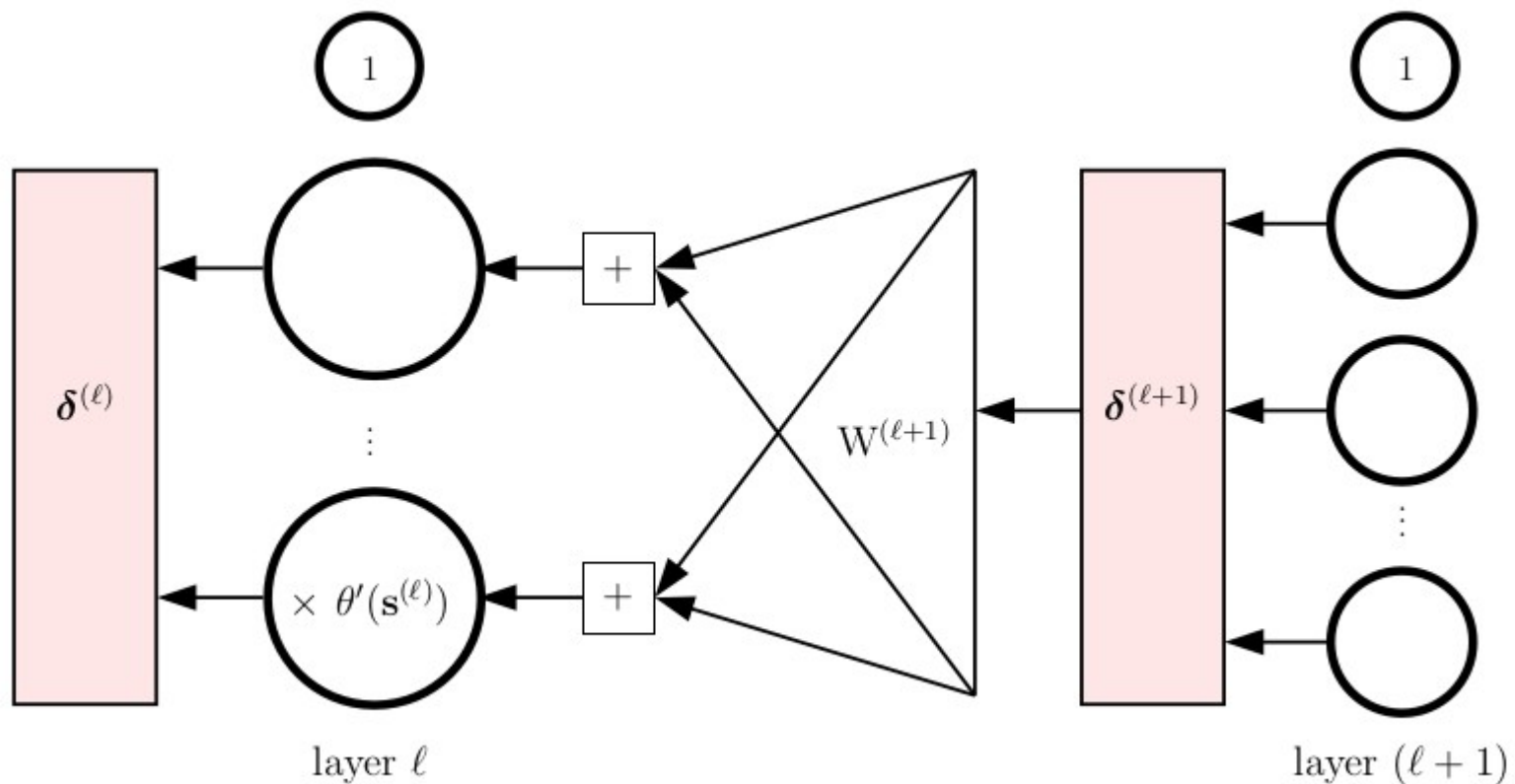
$$\begin{aligned} \frac{\partial e}{\partial \mathbf{s}_j^{(\ell)}} &= \frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial \mathbf{x}_j^{(\ell)}}{\partial \mathbf{s}_j^{(\ell)}} \\ \delta_j^{(\ell)} &= \theta'(\mathbf{s}_j^{(\ell)}) \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}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

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 $\theta'(\mathbf{s}_j^{(\ell)})$ and $\sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}$ in the second equation, indicating the simplification of the derivative of the input with respect to the output to the derivative of the activation function.

Backpropagation

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

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



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

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:

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

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

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

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

Backpropagation

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

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:

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

Esto es 0 si la red acierta

$$\begin{aligned} \tanh'(\mathbf{s}^{(\ell)}) &= \mathbf{1} - \tanh^2(\mathbf{s}^{(\ell)}) \\ &= \mathbf{1} - (x^{(L)})^2 \end{aligned}$$

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]

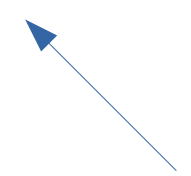
5: **for** $\ell = 1, \dots, L$ **do**

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

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

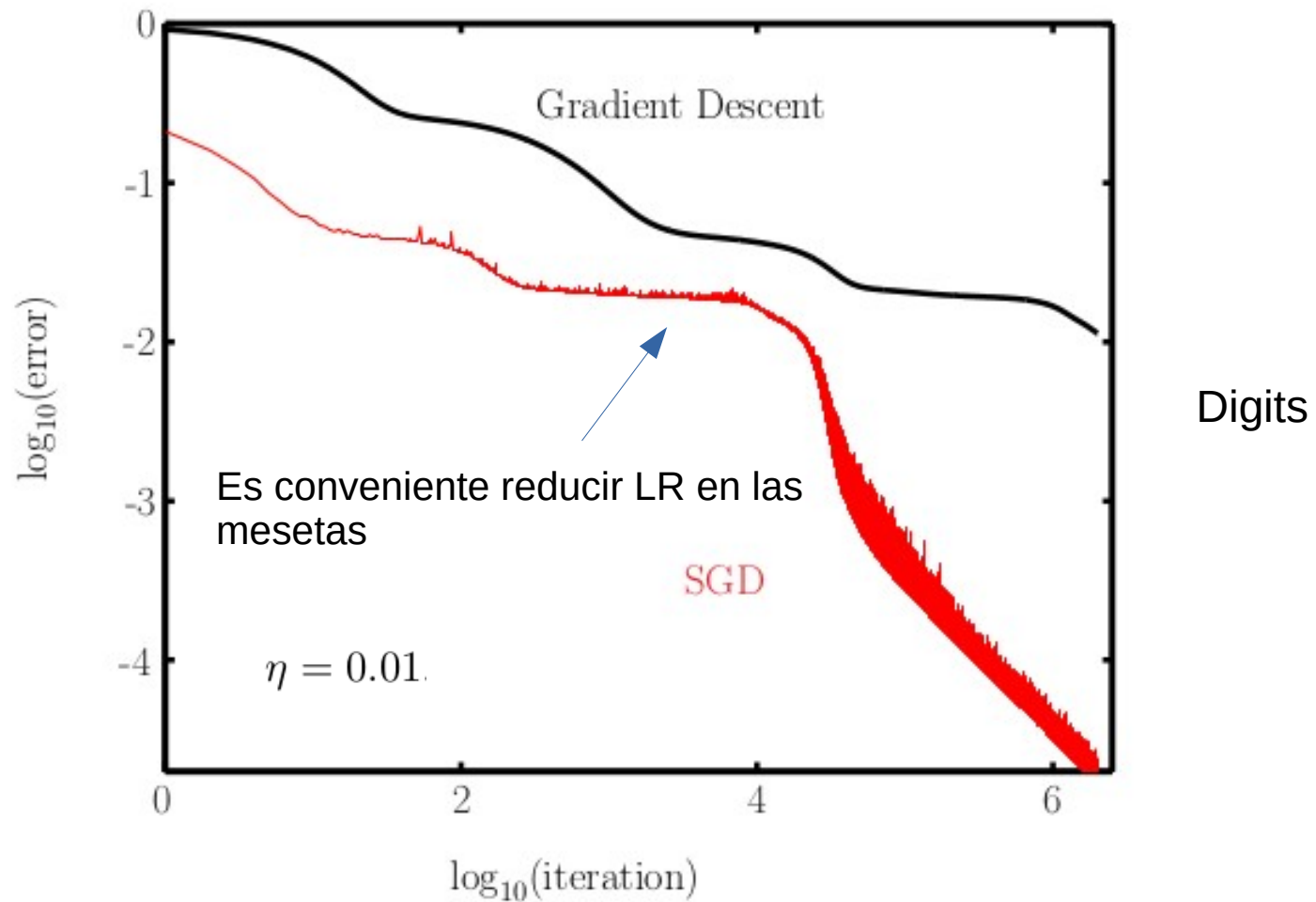
8: **end for**

9: **end for**

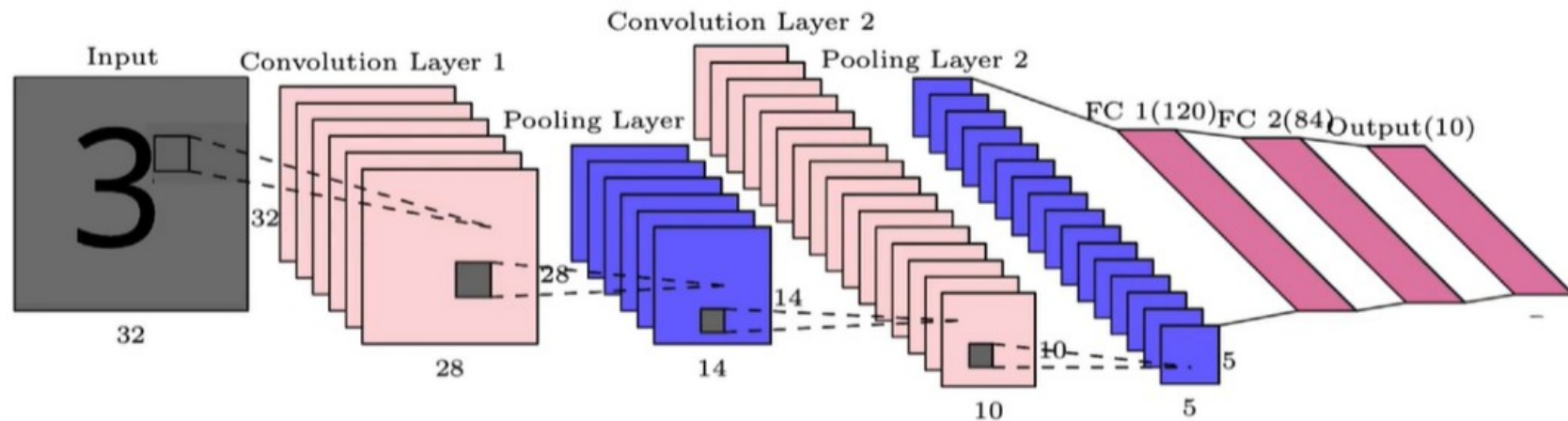
$$E_{\text{in}} \leftarrow E_{\text{in}} + \frac{1}{N} (\mathbf{x}_n^{(L)} - y_n)^2.$$


Backpropagation

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



Las redes neuronales son útiles para trabajar sobre datos raw



En este caso, la red trabaja directamente sobre la imagen. Usa dos operadores (filtro convolucional y pooling) para extraer patches desde la entrada.

Al final, usa capas densas y colapsa a una softmax.

Si la red aprende, debiera existir un gap pequeño entre accuracy de validación y training (idem para pérdida)

