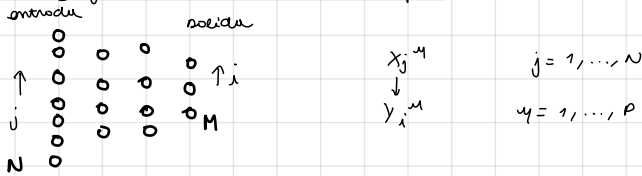


# Aprendizaje en redes multicapa



## Sistema lineal

$$g(x) = x$$

$$O_i^{(n)} = \sum_{j=1}^N w_{ij} x_j^{(n)}$$

*salida*  
*ENTRADA*

$$w_{ij} ( )$$

Poderemos encontrar  
tales que  $O_i^{(n)} = y_i^{(n)}$

*salida deseada*

Vemos que  $w_{ij}$  debe ser de la forma:

$$w_{ij} = \frac{1}{N} \sum_{n,v=1}^N y_i^{(n)} (Q^{-1})_{nv} x_j^{(n)}$$

$$Q_{nv} = \frac{1}{N} \sum_j x_j^{(n)} x_j^{(v)} \quad (n,v)$$

es el que cumple con  $O_i^{(n)} = y_i^{(n)}$

$$\begin{aligned} \sum_{j=1}^N w_{ij} x_j^{(n)} &= \frac{1}{N} \sum_{n,v,j} y_i^{(n)} (Q^{-1})_{nv} x_j^{(n)} x_j^{(v)} \\ &= \sum_{n,v} y_i^{(n)} \underbrace{(Q^{-1})_{nv} Q_{nv}}_{\delta_{nv}} = y_i^{(n)} \end{aligned}$$

## Aprendizaje por gradiente

$$E = \frac{1}{2} \sum_{n,v} (y_i^{(n)} - O_i^{(n)})^2 \quad E \geq 0$$

$$E = 0 \Leftrightarrow y_i^{(n)} = O_i^{(n)} \quad i=1, \dots, M \quad n=1, \dots, P$$

Resolver el problema de aprendizaje  $\Leftrightarrow$  mínimo absoluto de  $E$

$$\bar{w}_0$$

$$\bar{w}_0 + \Delta \bar{w}$$

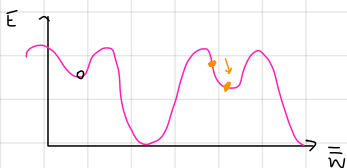
$$\Rightarrow E(\bar{w}_0 + \Delta \bar{w}) = E(\bar{w}) + \overline{\nabla} E(\bar{w}) \cdot \Delta \bar{w} + \mathcal{O}(\Delta \bar{w}^2)$$

tomamos  $\Delta \bar{w} = -\eta \overline{\nabla} E(\bar{w})$  para que (\*) sea negativo

$$= E(\bar{w}) - \eta |\overline{\nabla} E(\bar{w})|^2 + \mathcal{O}(\Delta \bar{w}^2) \leq E(\bar{w})$$

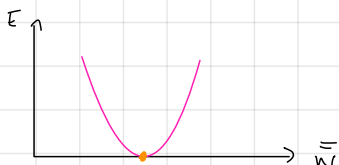
para asegurar esto pedimos  $\eta$  pequeño

grafiquemos una proyección de  $\bar{w}$



el método converge al mínimo local más cercano a los C.I

Para sistemas lineales E es cuadrática



$\Rightarrow$  el sistema es globalmente convergente

Para estos sist. tenemos:

$$E(\bar{w}) = \frac{1}{2} \sum_{i, \mu} \left( y_i^\mu - \sum_j w_{ij} x_j^\mu \right)^2$$

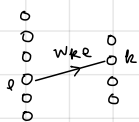
$$\frac{\partial E}{\partial w_{kl}} = - \sum_{\mu} \underbrace{\left( y_k^\mu - \sum_j w_{kj} x_j^\mu \right)}_{\delta l^\mu} x_l^\mu$$

son continuos (mo 0,1 como en el perceptron)

$$\Delta w_{kl} = \eta \sum_{\mu=1}^P \delta l^\mu x_l^\mu$$

$$\delta l^\mu = y_k^\mu - o_k^\mu$$

regla de Adeline



conexión efectiva entre neuronas  $k$  de entrada y  $e$  de salida :  $w_{ke}$

una capa, no lineal

$$O_i^u = g \left( \sum_{j=1}^N w_{ij} x_j^u \right)$$

$$O_i^u - O_i^u$$

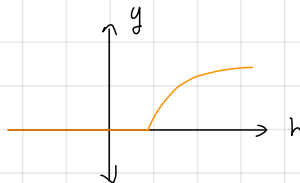
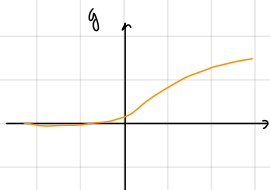
$$E = \frac{1}{2} \sum_{i,u} \left( y_i^u - g \left( \sum_{j=1}^N w_{ij} x_j^u \right) \right)^2$$

$$\frac{\partial E}{\partial w_{ke}} = \sum_u \left( y_k^u - g \left( \sum_{j=1}^N w_{kj} x_j^u \right) \right) g' \left( \sum_{j=1}^N w_{kj} x_j^u \right) x_k^u$$

$\underbrace{\hspace{10em}}_{\delta_k^u} \quad \underbrace{\hspace{10em}}_{h_k^u}$

$$= - \sum_u \delta_k^u g'(h_k^u) x_k^u$$

posibles  $g$ :



si  $g(h) = \tanh(h) \Rightarrow g'(h) = 1 - g^2(h)$

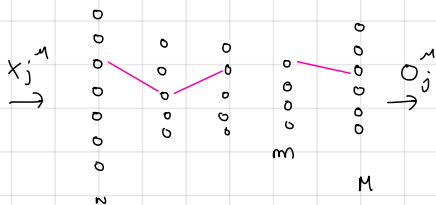
Podríamos definir la esperanza de otra manera:

$$E = \sum_{i,u} \left[ \frac{1 + y_i^u}{2} \ln \left( \frac{1 + y_i^u}{1 + O_i^u} \right) + \frac{1 - y_i^u}{2} \ln \left( \frac{1 - y_i^u}{1 - O_i^u} \right) \right]$$

$$\frac{\partial E}{\partial w_{ke}} = \sum_u \frac{g'(h_k^u)}{(1 - O_k^u)^2} [y_k^u - O_k^u] x_k^u$$

se redefinió  $\delta$  pero la forma funcional es la misma

## Caso general: redes multicapa no-lineales



$O_i^u$ : funciones de los  
entornos ( $\bar{w}_{ij} x_j^u$ )

$$E = \frac{1}{2} \sum_{i,j} (y_i^u - O_i^u (\bar{w}_{ij} x_j^u))^2$$

$$\Delta w_{kl} = -\eta \frac{\partial E}{\partial w_{kl}}$$

numéricamente:

$$\frac{\partial E}{\partial w_{kl}} \approx \frac{E(w_{kl} + \Delta) - E(w_{kl})}{\Delta}$$

tiempo de cálculo d n° de  
parámetros (P)  
x  
n° total  
de conexiones

## Back-propagation

$$(1-D) \quad f(g(x))' = f'(g(x)) g'(x)$$

$$(N-D) \quad f(g_1(x), \dots, g_n(x))' = \sum_i \frac{\partial f}{\partial g_i} (g_1(x), \dots, g_n(x)) g_i'(x)$$

$$(N-N)-D \quad \frac{\partial}{\partial x_k} f(g_1(x_1, \dots, x_n), g_2(x_1, \dots, x_n), \dots, g_n(x_1, \dots, x_n)) \\ = \sum_i \frac{\partial f}{\partial g_i} (g_1, \dots, g_n) \cdot \frac{\partial g_i}{\partial x_k} (x_1, \dots, x_n)$$

tenemos una función  $E$ , sobre la cual queremos  
calcular el gradiente  $\frac{\partial E}{\partial w_{kl}}$

$E$ : función de los neuronas de alguna capa

$$E(\{O_i^u\}) = E\left(g\left(\sum_{h_i^u} w_{ih} o_h^u\right)\right), \frac{\partial E}{\partial w_{kl}} = \frac{\partial E}{\partial h_k^u} \frac{\partial h_k^u}{\partial w_{kl}} = \frac{\partial E}{\partial h_k^u} g'(h_k^u) o_l^u$$

salida de la capa

matriz que el índice le está representado

al comenzar la salida, se comienza a introducir de las neuronas de la capa siguiente  $m$ .

$$\frac{\partial E}{\partial h_k^m} = \sum_m \frac{\partial E}{\partial h_m^m} \underbrace{\frac{\partial h_m^m}{\partial h_k^m}}_{\text{cómo varía la salida de la capa respecto a una variación en la salida de la capa anterior}} \quad h_m^m = \sum_s w_{ms} o_s^m = \sum_s w_{ms} g(h_s^m)$$

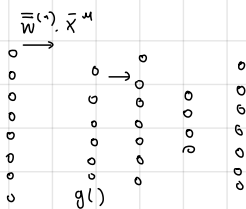
$\rightarrow \frac{\partial h_m^m}{\partial h_k^m} = w_{mk} g'(h_k^m)$

$$\frac{\partial E}{\partial h_k^m} = \sum_m \left( \frac{\partial E}{\partial h_m^m} \right) w_{mk} g'(h_k^m)$$

↓ debemos calcular esto → en columnas para la última capa y vamos para atrás

$$E = \frac{1}{2} \sum_{k,m} (y_k^m - o_k^m)^2 = \frac{1}{2} \sum (y_k^m - g(h_k^m))^2$$

$$\Rightarrow \frac{\partial E}{\partial h_m^m} = - \sum_k \underbrace{(y_k^m - g(h_k^m))}_{\text{errores}} g'(h_k^m) \rightarrow \text{para la última capa}$$



$$(\bar{y} - \bar{o}) \otimes g'(\bar{h})$$

↓  
prod.  
pto  
pto

